

**On Real-Time STM Concurrency Control for Embedded
Software: Improved Time Bounds**

Journal:	<i>Transactions on Architecture and Code Optimization</i>
Manuscript ID:	TACO-2012-117
Manuscript Type:	Original Work
Date Submitted by the Author:	18-Jun-2012
Complete List of Authors:	Elshambakey, Mohammed; Virginia Tech, ECE; MUCSAT, IRI Ravindran, Binoy; Virginia Tech University, ECE
Computing Classification Systems:	Special-Purpose and Application-based Systems, Real-time and embedded systems



On Real-Time STM Concurrency Control for Embedded Software: Improved Time Bounds

We consider software transactional memory (STM) concurrency control for embedded multicore real-time software, and present a novel contention manager for resolving transactional conflicts, called FBLT. We upper bound transactional retries and task response times, and show when FBLT has better schedulability than previous contention managers and lock-free synchronization. Our implementation in the Rochester STM framework/real-time Linux reveals that FBLT yields shorter or comparable retry costs than competitors.

Categories and Subject Descriptors: C.3 [Special-Purpose and Application-based Systems]: Real-time and embedded systems

General Terms: Design, Experimentation, Measurement

Additional Key Words and Phrases: Software transactional memory (STM), real-time contention manager

ACM Reference Format:

ACM Trans. Embedd. Comput. Syst. , Article (), 25 pages.

DOI = 10.1145/0000000.0000000 <http://doi.acm.org/10.1145/0000000.0000000>

1. INTRODUCTION

Embedded systems sense physical processes and control their behavior, typically through feedback loops. Since physical processes are concurrent, computations that control them must also be concurrent, enabling them to process multiple streams of sensor input and control multiple actuators, all concurrently. Often, such computations need to concurrently read/write shared data objects. They must also process sensor input and react, while satisfying time constraints.

The de facto standard for concurrent programming is the threads abstraction, and the de facto synchronization abstraction is locks. Lock-based concurrency control has significant programmability, scalability, and composability challenges [Herlihy 2006]. Transactional memory (TM) is an alternative synchronization model for shared memory objects that promises to alleviate these difficulties. With TM, code that read/write shared objects is organized as *memory transactions*, which execute speculatively, while logging changes made to objects. Two transactions conflict if they access the same object and one access is a write. When that happens, a contention manager (CM) [Guer-raoui et al. 2005] resolves the conflict by aborting one and allowing the other to commit, yielding (the illusion of) atomicity. Aborted transactions are re-started, after rolling back the changes. In addition to a simple programming model, TM provides performance comparable to locking and lock-free approaches, especially for high contention and read-dominated workloads (see an example TM system’s performance in [Saha et al. 2006]), and is composable [Harris et al. 2008]. TM has been proposed in hard-

New Paper, Not an Extension of a Conference Paper.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org.

© ACM 1539-9087/-ART \$15.00

DOI 10.1145/0000000.0000000 <http://doi.acm.org/10.1145/0000000.0000000>

:2

ware, called HTM, and in software, called STM, with the usual tradeoffs: HTM has lesser overhead, but needs transactional support in hardware; STM is available on any hardware.

Given STM's programmability, scalability, and composability advantages, it is a compelling concurrency control technique also for multicore embedded real-time software. However, this requires bounding transactional retries, as real-time threads, which subsume transactions, must satisfy time constraints. Retry bounds under STM are dependent on the CM policy at hand (analogous to the way thread response time bounds are OS scheduler-dependent).

Past real-time CM research (Section 5) has proposed resolving transactional contention using dynamic and fixed priorities of parent threads, resulting in EDF-based CM (ECM) and RMS-based CM (RCM) [Fahmy and Ravindran 2011; El-Shambakey and Ravindran 2012b; 2012a], which are intended to be used with global EDF (G-EDF) and global RMS (G-RMS) multicore real-time schedulers [Davis and Burns 2011], respectively. In particular, [El-Shambakey and Ravindran 2012b] shows that ECM and RCM achieve higher schedulability – i.e., greater number of task sets meeting their time constraints – than lock-free synchronization only under some ranges for the maximum atomic section length. That range is significantly expanded with the LCM contention manager in [El-Shambakey and Ravindran 2012a], increasing the coverage of STM's timeliness superiority. ECM, RCM, and LCM suffer from transitive retry and cannot handle multiple objects per transaction efficiently.

These limitations are overcome with the PNF contention manager [El-Shambakey 2012]. However, PNF requires a-priori knowledge of all objects accessed by each transaction. This significantly limits programmability, and is incompatible with dynamic STM implementations [Herlihy et al. 2003]. Additionally, PNF is a centralized CM, which increases overheads and retry costs, and has a complex implementation.

We propose the First Bounded, Last Timestamp (or FBLT) contention manager (Section 5). In contrast to PNF, FBLT does not require a-priori knowledge of objects accessed by transactions. Moreover, FBLT allows each transaction to access multiple objects with shorter transitive retry cost than ECM, RCM and LCM. Additionally, FBLT is a decentralized CM and does not use locks in its implementation. Implementation of FBLT is also simpler than PNF.

We establish FBLT's retry and response time upper bounds under G-EDF and G-RMA schedulers (Section 6). We also identify the conditions under which FBLT's schedulability is better than ECM, RCM, G-EDF/LCM, G-RMA/LCM, PNF, and lock-free synchronization (Section 7).

We implement FBLT and competitor CM techniques in the Rochester STM framework [Marathe et al.] and conduct experimental studies (Section 8). Our results reveal that FBLT has shorter retry cost than ECM, RCM and LCM. FBLT has no priori knowledge of accessed objects by transactions. So, retry cost under FBLT is a little higher than PNF.

Thus, the paper's contribution is the FBLT contention manager with superior timeliness properties. FBLT, thus allows programmers to reap STM's significant programmability and composability benefits for a broader range of multicore embedded real-time software than what was previously possible.

2. RELATED WORK

Transactional-like concurrency control without using locks, for real-time systems, has been previously studied in the context of non-blocking data structures (e.g., [Anderson et al. 1995]). Despite their numerous advantages over locks (e.g., deadlock-freedom), their programmability has remained a challenge. Past studies show that they are best suited for simple data structures where their retry cost is competitive to the cost of

FBLT

:3

lock-based synchronization [Brandenburg et al. 2008]. In contrast, STM is semantically simpler [Herlihy 2006], and is often the only viable lock-free solution for complex data structures (e.g., red/black tree) [Fahmy 2010] and nested critical sections [Saha et al. 2006].

STM concurrency control for real-time systems has been previously studied in [Manson et al. 2006; Fahmy et al. 2009; Sarni et al. 2009; Schoeberl et al. 2010; Fahmy 2010; Barros and Pinho 2011; El-Shambakey and Ravindran 2012b; 2012a; El-Shambakey 2012].

[Manson et al. 2006] proposes a restricted version of STM for uniprocessors. Uniprocessors do not need contention management. [Fahmy et al. 2009] bounds response times in distributed systems with STM synchronization. They consider Pfair scheduling, limit to small atomic regions with fixed size, and limit transaction execution to span at most two quanta. In contrast, we allow transaction lengths with arbitrary duration.

[Sarni et al. 2009] presents real-time scheduling of transactions and serializes transactions based on deadlines. However, the work does not bound retries and response times. In contrast, we establish such bounds. [Schoeberl et al. 2010] proposes real-time HTM. The work does not describe how transactional conflicts are resolved. Besides, the retry bound assumes that the worst case conflict between atomic sections of different tasks occurs when the sections are released at the same time. However, we show that this is not the worst case. We develop retry and response time upper bounds based on much worse conditions.

[Fahmy 2010] upper bounds retries and response times for ECM with G-EDF, and identify the tradeoffs with locking and lock-free protocols. Similar to [Schoeberl et al. 2010], [Fahmy 2010] also assumes that the worst case conflict between atomic sections occurs when the sections are released simultaneously. The ideas in [Fahmy 2010] are extended in [Barros and Pinho 2011], which presents three real-time CM designs. But no retry bounds or schedulability analysis techniques are presented for those CMs.

[El-Shambakey and Ravindran 2012b] presents the ECM and RCM contention managers, and upper bounds transactional retries and task response times under them. The work also identifies the conditions under which ECM and RCM are superior to locking and lock-free techniques. In particular, [El-Shambakey and Ravindran 2012b] shows that, STM's superiority holds only under some ranges for the maximum atomic section length. Moreover, [El-Shambakey and Ravindran 2012b] restricts transactions to access only one object. [El-Shambakey and Ravindran 2012a] presents length-based CM (LCM), and upper bounds transactional retries and response time for G-EDF/LCM and G-RMA/LCM. [El-Shambakey and Ravindran 2012a] compares (analytically and experimentally) between ECM, RCM and LCM, as well as lock-free and LCM. [El-Shambakey and Ravindran 2012a], as [El-Shambakey and Ravindran 2012b], restricts transactions to access only one object.

[El-Shambakey 2012] presents Priority CM with Negative value and First Access (PNF). PNF was designed to avoid transitive retry effect when each transaction accesses multiple objects. PNF also optimizes processor usage by reducing priority of retrying transactions below priorities of any real-time task. PNF requires prior knowledge of accessed objects by each transaction which is not always available. Besides, PNF is a centralized CM that uses locks in its implementation. Accordingly, reduction in retry cost is wasted by high overhead. [El-Shambakey 2012] upper bounds transactional retries and response time for G-EDF and G-RMA. [El-Shambakey 2012] compares (analytically and experimentally) between PNF and ECM, RCM, LCM and lock-free.

Our work builds upon [El-Shambakey and Ravindran 2012b; 2012a; El-Shambakey 2012], allows multiple objects per transaction with no prior knowledge about these

:4

objects. We upper bound transactional retries and task response times. We identify the conditions for better schedulability for FBLT than other synchronization techniques.

3. PRELIMINARIES

We consider a multiprocessor system with m identical processors and n sporadic tasks $\tau_1, \tau_2, \dots, \tau_n$. The k^{th} instance (or job) of a task τ_i is denoted τ_i^k . Each task τ_i is specified by its worst case execution time (WCET) c_i , its minimum period T_i between any two consecutive instances, and its relative deadline D_i , where $D_i = T_i$. Job τ_i^j is released at time r_i^j and must finish no later than its absolute deadline $d_i^j = r_i^j + D_i$. Under a fixed priority scheduler such as G-RMA, p_i determines τ_i 's (fixed) priority and it is constant for all instances of τ_i . Under a dynamic priority scheduler such as G-EDF, a job τ_i^j 's priority, p_i^j , differs from one instance to another. A task τ_j may interfere with task τ_i for a number of times during an interval L , and this number is denoted as $G_{ij}(L)$.

Shared objects. A task may need to read/write shared, in-memory data objects while it is executing any of its atomic sections (transactions), which are synchronized using STM. The set of atomic sections of task τ_i is denoted s_i . s_i^k is the k^{th} atomic section of τ_i . Each object, θ , can be accessed by multiple tasks. The set of distinct objects accessed by τ_i is θ_i without repeating objects. The set of atomic sections used by τ_i to access θ is $s_i(\theta)$, and the sum of the lengths of those atomic sections is $len(s_i(\theta))$. $s_i^k(\theta)$ is the k^{th} atomic section of τ_i that accesses θ . s_i^k can access one or more objects in θ_i . So, s_i^k refers to the transaction itself, regardless of the objects accessed by the transaction. We denote the set of all accessed objects by s_i^k as Θ_i^k . While $s_i^k(\theta)$ implies that s_i^k accesses an object $\theta \in \Theta_i^k$, $s_i^k(\Theta)$ implies that s_i^k accesses a set of objects $\Theta = \{\theta \in \Theta_i^k\}$. $\bar{s}_i^k = \bar{s}_i^k(\Theta)$ refers only once to s_i^k , regardless of the number of objects in Θ . So, $|\bar{s}_i^k(\Theta)|_{\forall \theta \in \Theta} = 1$. $s_i^k(\theta)$ executes for a duration $len(s_i^k(\theta))$. $len(s_i^k) = len(s_i^k(\theta)) = len(s_i^k(\Theta)) = len(s_i^k(\Theta_i^k))$. The set of tasks sharing θ with τ_i is denoted $\gamma_i(\theta)$.

Atomic sections are non-nested (supporting nested STM is future work). The maximum-length atomic section in τ_i that accesses θ is denoted $s_{i,max}(\theta)$, while the maximum one among all tasks is $s_{max}(\theta)$, and the maximum one among tasks with priorities lower than that of τ_i is $s_{i,max}^l(\theta)$.

STM retry cost. If two or more atomic sections conflict, the CM will commit one section and abort and retry the others, increasing the time to execute the aborted sections. The increased time that an atomic section $s_i^p(\theta)$ will take to execute due to a conflict with another section $s_j^k(\theta)$, is denoted $W_i^p(s_j^k(\theta))$. If an atomic section, s_i^p , is already executing, and another atomic section s_j^k tries to access a shared object with s_i^p , then s_j^k is said to “interfere” or “conflict” with s_i^p . The job s_j^k is the “interfering job”, and the job s_i^p is the “interfered job”.

Due to *transitive retry* (introduced in Section 4.3), an atomic section $s_i^k(\Theta_i^k)$ may retry due to another atomic section $s_j^l(\Theta_j^l)$, where $\Theta_i^k \cap \Theta_j^l = \emptyset$. θ_i^* denotes the set of objects not accessed directly by atomic sections in τ_i , but can cause transactions in τ_i to retry due to transitive retry. $\theta_i^{ex}(= \theta_i + \theta_i^*)$ is the set of all objects that can cause transactions in τ_i to retry directly or through transitive retry. γ_i^* is the set of tasks that accesses objects in θ_i^* . $\gamma_i^{ex}(= \gamma_i + \gamma_i^*)$ is the set of all tasks that can directly or indirectly (through transitive retry) cause transactions in τ_i to abort and retry.

The total time that a task τ_i 's atomic sections have to retry over T_i is denoted $RC(T_i)$. The additional amount of time by which all interfering jobs of τ_j increases the response time of any job of τ_i during L , without considering retries due to atomic sections, is denoted $W_{ij}(L)$.

4. MOTIVATION

To understand the need for *First Bounded, Last Timestamp (FBLT)* contention manager, we first give a brief introduction of previous real-time CMs, ECM [El-Shambakey and Ravindran 2012b], RCM [El-Shambakey and Ravindran 2012b], LCM [El-Shambakey and Ravindran 2012a], and PNF [El-Shambakey 2012].

4.1. ECM and RCM

The *Earliest Deadline Contention Manager (ECM)* [El-Shambakey and Ravindran 2012b] is used with the G-EDF multicore real-time scheduler. ECM allows the transaction with the shortest absolute deadline to commit first. The other transactions retry and abort.

The *Rate Monotonic Contention Manager (RCM)* [El-Shambakey and Ravindran 2012b] is used with the G-RMA scheduler. RCM allows the transaction with the shortest period to commit first. ECM and RCM maintain semantic consistency with the underlying scheduler.

4.2. LCM

For both ECM and RCM, $s_i^k(\theta)$ can be totally repeated if $s_j^l(\theta)$ — which belongs to a higher priority job τ_j^b than τ_i^a — conflicts with $s_i^k(\theta)$ at the end of its execution, while $s_i^k(\theta)$ is just about to commit. The *Length-based Contention Manager (LCM)* [El-Shambakey and Ravindran 2012a], shown in Algorithm 1, uses the remaining length of $s_i^k(\theta)$ when it is interfered, as well as $len(s_j^l(\theta))$, to decide which transaction must be aborted. If p_i^k is greater than p_j^l , then $s_i^k(\theta)$ is committed, because it belongs to a higher priority job, and it started before $s_j^l(\theta)$ (step 2). Otherwise, c_{ij}^{kl} is calculated (step 4) to determine whether it is worth aborting $s_i^k(\theta)$ in favor of $s_j^l(\theta)$, because $len(s_j^l(\theta))$ is relatively small compared to the remaining execution length of $s_i^k(\theta)$. [El-Shambakey and Ravindran 2012a] assumes that:

$$c_{ij}^{kl} = len(s_j^l(\theta)) / len(s_i^k(\theta)) \quad (1)$$

where $c_{ij}^{kl} \in]0, \infty[$, to cover all possible lengths of $s_j^l(\theta)$.

Thus, LCM's key idea is to reduce the opportunity for the abort of $s_i^k(\theta)$ if it is close to committing when interfered and $len(s_j^l(\theta))$ is large. This abort opportunity is increasingly reduced as $s_i^k(\theta)$ gets closer to the end of its execution, or $len(s_j^l(\theta))$ gets larger. On the other hand, as $s_i^k(\theta)$ is interfered early, or $len(s_j^l(\theta))$ is small compared to $s_i^k(\theta)$'s remaining length, the abort opportunity is increased even if $s_i^k(\theta)$ is close to the end of its execution. To decide whether $s_i^k(\theta)$ must be aborted or not, a threshold value $\psi \in [0, 1]$ that determines α_{ij}^{kl} (step 5) is used, where α_{ij}^{kl} is the maximum percentage of $len(s_i^k(\theta))$ below which $s_j^l(\theta)$ is allowed to abort $s_i^k(\theta)$. Thus, if the already executed part of $s_i^k(\theta)$ — when $s_j^l(\theta)$ interferes with $s_i^k(\theta)$ — does not exceed $\alpha_{ij}^{kl} len(s_i^k(\theta))$, then $s_i^k(\theta)$ is aborted (step 8). Otherwise, $s_j^l(\theta)$ is aborted (step 10).

LCM reduces the retry cost of a single transaction $s_i^k(\theta)$ due to another transaction $s_j^l(\theta)$ from $2.s_{max}$ (in case of ECM and RCM) to $(1 + \alpha_{max}).s_{max}$, where s_{max} is the maximum length transaction among all tasks, and α_{max} is the maximum *alpha* for any transaction. On the other hand, LCM suffers from bounded priority inversion because a higher priority transaction can be blocked by a lower priority one. Additionally, LCM is not a centralized CM, which means that, upon a conflict, each transaction must decide whether it must commit or abort.

:6

ALGORITHM 1: The LCM Algorithm**Data:** $s_i^k(\theta) \rightarrow$ interfered atomic section. $s_j^l(\theta) \rightarrow$ interfering atomic section. $\psi \rightarrow$ predefined threshold $\in [0, 1]$. $\delta_i^k(\theta) \rightarrow$ remaining execution length of $s_i^k(\theta)$ **Result:** which atomic section of $s_i^k(\theta)$ or $s_j^l(\theta)$ aborts

```

1 if  $p_i^k > p_j^l$  then
2   |  $s_j^l(\theta)$  aborts;
3 else
4   |  $c_{ij}^{kl} = \text{len}(s_j^l(\theta)) / \text{len}(s_i^k(\theta));$ 
5   |  $\alpha_{ij}^{kl} = \ln(\psi) / (\ln(\psi) - c_{ij}^{kl});$ 
6   |  $\alpha = (\text{len}(s_i^k(\theta)) - \delta_i^k(\theta)) / \text{len}(s_i^k(\theta));$ 
7   | if  $\alpha \leq \alpha_{ij}^{kl}$  then
8     |  $s_i^k(\theta)$  aborts;
9   | else
10    |  $s_j^l(\theta)$  aborts;
11  | end
12 end

```

4.3. PNF

ECM, RCM, and LCM suffer from *transitive retry*. Transitive retry is illustrated by the following example:

Consider three atomic sections s_1^x , s_2^y , and s_3^z belonging to jobs τ_1^x , τ_2^y , and τ_3^z , with priorities $p_3^z > p_2^y > p_1^x$, respectively. Assume that s_1^x and s_2^y share objects, and s_2^y and s_3^z share objects. s_1^x and s_3^z do not share objects. Now, s_3^z can cause s_2^y to retry, which in turn will cause s_1^x to retry. This means that s_1^x will retry transitively because of s_3^z , which will increase the retry cost of s_1^x .

Now, consider another atomic section s_4^f with a priority higher than that of s_3^z . Suppose s_4^f shares objects only with s_3^z . Thus, s_4^f can cause s_3^z to retry, which in turn will cause s_2^y to retry, and finally, s_1^x to retry. Thus, transitive retry will move from s_4^f to s_1^x , increasing the retry cost of s_1^x . The situation gets worse as more higher priority tasks are added, where each task shares objects with its immediate lower priority task. τ_3^z may have atomic sections that share objects with τ_1^x , but this will not prevent the effect of transitive retry due to s_1^x .

Therefore, the analysis in [El-Shambakey and Ravindran 2012b] and [El-Shambakey and Ravindran 2012a] extend the set of objects that can cause an atomic section of a lower priority job to retry. This is done by initializing the set of conflicting objects, γ_i , to all objects accessed by all transactions of τ_i . We then cycle through all transactions belonging to all other higher priority tasks. Each transaction s_j^l that accesses at least one of the objects in γ_i adds all other objects accessed by s_j^l to γ_i . The loop over all higher priority tasks is repeated, each time with the new γ_i , until there are no more transactions accessing any object in γ_i . The final set of objects (tasks) that can cause transactions in τ_i to retry is $\theta_i^{ex}(\gamma_i^{ex})$, respectively¹.

The *Priority contention manager with Negative value and First access* (PNF) [El-Shambakey 2012] is designed to avoid transitive retry. ECM, RCM, and LCM suffer from transitive retry. PNF avoids transitive retry by concurrently executing at most m non-conflicting transactions together. These executing transactions are non-

¹However, note that, this solution may over-extend the set of conflicting objects, and may even contain all objects accessed by all tasks.

FBLT

:7

preemptive. Thus, executing transactions cannot be aborted due to direct or indirect conflict with other transactions.

However, with PNF, all objects accessed by each transaction must be known a-priori. Therefore, this is not suitable with dynamic STM implementations [Herlihy et al. 2003]. Additionally, PNF is implemented in [El-Shambakey 2012] as a centralized CM that uses locks. This increases overhead.

4.4. Case for FBLT

Thus, it is desirable to have a CM with the following goals:

- (1) reduce the retry cost of each transaction s_i^k due to another transaction s_j^l , just as LCM does compared to ECM and RCM.
- (2) avoid or bound the effect of transitive retry, similar to PNF, without prior knowledge of accessed objects by each transaction, enabling dynamic STM.
- (3) decentralized design and avoid the use of locks, thereby reducing overhead.

We propose the *First Bounded, Last Timestamp contention manager* (or (FBLT)). FBLT achieves these goals by bounding the number of times each transaction s_i^k is aborted due to other transactions to at most δ_i^k . δ_i^k includes the number of aborts due to direct conflict with other transactions, as well as transitive retry (goal 2). If a transaction s_i^k reaches its δ_i^k , it is added to an m_set in FIFO order. In the m_set , s_i^k executes non-preemptively. If transactions in the m_set conflict together, they use their FIFO order in the m_set to resolve the conflict. s_i^k can still abort after it becomes a non-preemptive transaction due to other non-preemptive transactions. The number of aborts for any non-preemptive transaction is bounded by $m-1$, where m is the number of processors, as will be shown in Section 6.

Thus, the key idea behind FBLT is to use a suitable δ_i^k for each s_i^k before it becomes a non-preemptive transaction. The choice of δ_i^k should make the total retry cost (and thus, the schedulability) of any job τ_i^x under FBLT comparable to the retry cost under ECM, RCM, LCM, and PNF. (In Section 7, we show the suitable δ_i^k for each s_i^k to have equal or better schedulability than other CMs.) Preemptive transactions resolve their conflicts using LCM. Thus, FBLT defaults to LCM if abort bounds have not been violated (goal 1). Each non-preemptive transaction s_i^k uses the time it joined the m_set to resolve conflicts with other non-preemptive transactions. Therefore, FBLT does not have to use locks and is decentralized (goal 3).

5. THE FBLT CONTENTION MANAGER

Algorithm 2 illustrates FBLT. Each transaction s_i^k can be aborted during T_i for at most δ_i^k times. η_i^k records the number of times s_i^k has already been aborted up to now. If s_i^k and s_j^l have not joined the m_set yet, then they are preemptive transactions. Preemptive transactions resolve conflicts using Algorithm 1 (step 2). Thus, FBLT defaults to LCM when no transaction reaches its δ . If only one of the transactions is in the m_set , then the non-preemptive transaction (the one in m_set) aborts the other one (steps 15 to 26). η_i^k is incremented each time s_i^k is aborted as long as $\eta_i^k < \delta_i^k$ (steps 5 and 18). Otherwise, s_i^k is added to the m_set and its priority is increased to m_prio (steps 7 to 9 and 20 to 22). When the priority of s_i^k is increased to m_prio , s_i^k becomes a non-preemptive transaction. Non-preemptive transactions cannot be aborted by other preemptive transactions, nor by any other real-time job. The m_set can hold at most m concurrent transactions because there are m processors in the system. $r(s_i^k)$ records the time s_i^k joined the m_set (steps 8 and 21). When non-preemptive transactions conflict together (step 27), the transaction with the smaller $r()$ commits first (steps 29 and 31). Thus, non-preemptive transactions are executed in FIFO order of the m_set .

:8

ALGORITHM 2: The FBLT Algorithm**Data:** s_i^k : interfered transaction; s_j^l : interfering transactions; δ_i^k : the maximum number of times s_i^k can be aborted during T_i ; η_i^k : number of times s_i^k has already been aborted up to now; m_set : contains at most m non-preemptive transactions. m is number of processors; m_prio : priority of any transaction in m_set . m_prio is higher than any priority of any real-time task; $r(s_i^k)$: time point at which s_i^k joined m_set ;**Result:** atomic sections that will abort

```

1 if  $s_i^k, s_j^l \notin m\_set$  then
2   Apply Algorithm 1 (default to LCM);
3   if  $s_i^k$  is aborted then
4     if  $\eta_i^k < \delta_i^k$  then
5       Increment  $\eta_i^k$  by 1;
6     else
7       Add  $s_i^k$  to  $m\_set$ ;
8       Record  $r(s_i^k)$ ;
9       Increase priority of  $s_i^k$  to  $m\_prio$ ;
10    end
11  else
12    Swap  $s_i^k$  and  $s_j^l$ ;
13    Go to Step 3;
14  end
15 else if  $s_j^l \in m\_set, s_i^k \notin m\_set$  then
16   Abort  $s_i^k$ ;
17   if  $\eta_i^k < \delta_i^k$  then
18     Increment  $\eta_i^k$  by 1;
19   else
20     Add  $s_i^k$  to  $m\_set$ ;
21     Record  $r(s_i^k)$ ;
22     Increase priority of  $s_i^k$  to  $m\_prio$ ;
23   end
24 else if  $s_i^k \in m\_set, s_j^l \notin m\_set$  then
25   Swap  $s_i^k$  and  $s_j^l$ ;
26   Go to Step 15;
27 else
28   if  $r(s_i^k) < r(s_j^l)$  then
29     Abort  $s_j^l$ ;
30   else
31     Abort  $s_i^k$ ;
32   end
33 end

```

5.1. Illustrative Example

We now illustrate FBLT's behavior with the following example:

(1) Transaction $s_i^k(\theta_1, \theta_2)$ is released while $m_set = \emptyset$. $\eta_i^k = 0$ and $\delta_i^k = 3$.

(2) Transaction $s_a^b(\theta_2)$ is released while $s_i^k(\theta_1, \theta_2)$ is running. $p_a^b > p_i^k$ and $\eta_i^k < \delta_i^k$. Applying LCM, $s_i^k(\theta_1, \theta_2)$ is aborted in favor of s_a^b and η_i^k is incremented to 1.

(3) $s_a^b(\theta_2)$ commits. $s_i^k(\theta_1, \theta_2)$ runs again. Transaction $s_c^d(\theta_2)$ is released while $s_i^k(\theta_1, \theta_2)$ is running. $p_c^d > p_i^k$. Applying LCM, $s_i^k(\theta_1, \theta_2)$ is aborted again in favor of $s_c^d(\theta_2)$. η_i^k is incremented to 2.

FBLT

:9

- (4) $s_e^d(\theta_2)$ commits. $s_e^f(\theta_2, \theta_3)$ is released. $p_e^f > p_i^k$ and $\eta_e^f = 2$. $s_i^k(\theta_1, \theta_2)$ is aborted in favor of $s_e^f(\theta_2, \theta_3)$ and η_i^k is incremented to 3.
- (5) $s_j^l(\theta_3)$ is released. $p_j^l > p_e^f$. $s_e^f(\theta_2, \theta_3)$ is aborted in favor of $s_j^l(\theta_3)$ and η_e^f is incremented to 1.
- (6) $s_i^k(\theta_1, \theta_2)$ and $s_e^f(\theta_2, \theta_3)$ are compared again. $\eta_i^k = \delta_i^k$, $s_i^k(\theta_1, \theta_2)$ is added to m_set . $m_set = \{s_i^k(\theta_1, \theta_2)\}$. $s_i^k(\theta_1, \theta_2)$ becomes a non-preemptive transaction. As $s_e^f(\theta_2, \theta_3)$ is a preemptive transaction, $s_e^f(\theta_2, \theta_3)$ is aborted in favor of $s_i^k(\theta_1, \theta_2)$, despite p_e^f being greater than the original priority of $s_i^k(\theta_1, \theta_2)$. η_e^f is incremented to 2.
- (7) $s_j^l(\theta_3)$ commits but $s_g^h(\theta_3)$ is released. $p_g^h > p_e^f$ but $\eta_e^f = \delta_e^f$. So, $s_e^f(\theta_2, \theta_3)$ becomes a non-preemptive transaction. $m_set = \{s_i^k(\theta_1, \theta_2), s_g^h(\theta_2, \theta_3)\}$.
- (8) $s_i^k(\theta_1, \theta_2)$ and $s_g^h(\theta_2, \theta_3)$ are now non-preemptive transactions. $s_i^k(\theta_1, \theta_2)$ and $s_g^h(\theta_2, \theta_3)$ still conflict together. So, they are executed according to their addition order to the m_set . So, $s_i^k(\theta_1, \theta_2)$ commits first, followed $s_g^h(\theta_2, \theta_3)$.
- (9) $s_g^h(\theta_3)$ will continue to abort and retry in favor of $s_e^f(\theta_2, \theta_3)$ until $s_e^f(\theta_2, \theta_3)$ commits or $\eta_g^h = \delta_g^h$. Even if $s_g^h(\theta_3)$ joined the m_set , $s_g^h(\theta_3)$ will still abort and retry in favor of $s_e^f(\theta_2, \theta_3)$, because $s_e^f(\theta_2, \theta_3)$ joined the m_set earlier than $s_g^h(\theta_3)$.

It is seen from steps 2 to 6 that $s_i^k(\theta_1, \theta_2)$ can be aborted due to direct conflict with other transactions, or due to transitive retry. Irrespective of the reason for the conflict, once a transaction has reached its maximum allowed δ , the transaction becomes a non-preemptive one (steps 6 and 7). Non-preemptive transactions have higher priority than other preemptive transactions (steps 6 and 7). Non-preemptive transactions execute in their arrival order to the m_set .

6. RETRY COST AND RESPONSE TIME BOUNDS

We now derive an upper bound on the retry cost of any job τ_i^x under FBLT during an interval $L \leq T_i$. Since all tasks are sporadic (i.e., each task τ_i has a minimum period T_i), T_i is the maximum study interval for each task τ_i .

CLAIM 1. *The total retry cost for any job τ_i^x under FBLT due to 1) conflicts between its transactions and transactions of other jobs during an interval $L \leq T_i$ and 2) release of higher priority jobs is upper bounded by:*

$$RC_{to}(L) \leq \sum_{\forall s_i^k \in s_i} \left(\delta_i^k \text{len}(s_i^k) + \sum_{\forall s_{iz}^k \in \chi_i^k} \text{len}(s_{iz}^k) \right) + RC_{re}(L) \quad (2)$$

where χ_i^k is the set of at most $m - 1$ maximum length transactions conflicting directly or indirectly (through transitive retry) with s_i^k . Each transaction $s_{iz}^k \in \chi_i^k$ belongs to a distinct task τ_j . $RC_{re}(L)$ is the retry cost resulting from the release of higher priority jobs which preempt τ_i^x . $RC_{re}(L)$ is calculated by (6.8) in [El-Shambakey 2012] for G-EDF, and (6.10) in [El-Shambakey 2012] for G-RMA schedulers.

PROOF. By the definition of FBLT, $s_i^k \in \tau_i^x$ can be aborted a maximum of δ_i^k times before s_i^k joins the m_set . Before joining the m_set , s_i^k can be aborted due to higher priority transactions, or transactions in the m_set . The original priority of transactions in the m_set can be higher or lower than p_i^x . Thus, the maximum time s_i^k is aborted before joining the m_set occurs if s_i^k is aborted for δ_i^k times.

Transactions preceding s_i^k in the m_set can conflict directly with s_i^k , or indirectly through transitive retry. The worst case scenario for s_i^k after joining the m_set occurs if s_i^k is preceded by $m - 1$ maximum length conflicting transactions. Hence, in the worst

:10

case, s_i^k has to wait for the previous $m - 1$ transactions to commit first. The priority of s_i^k after joining the m_set is higher than any real-time job. Therefore, s_i^k is not aborted by any job. If s_i^k has not joined the m_set yet, and a higher priority job τ_j^y is released while s_i^k is running, then s_i^k may be aborted if τ_j^y has conflicting transactions with s_i^k . τ_j^y causes only one abort in τ_i^x because τ_j^y preempts τ_i^x only once. If s_i^k has already joined the m_set , then s_i^k cannot be aborted by the release of higher priority jobs. Thus, the maximum number of times transactions in τ_i^x can be aborted due to the release of higher priority jobs is less than or equal to the number of interfering higher priority jobs to τ_i^x . Claim follows. \square

CLAIM 2. Under FBLT, the blocking time of a job τ_i^x due to lower priority jobs is upper bounded by:

$$D(\tau_i^x) = \min \left(\max_1^m (s_{j_{max}}, \forall \tau_j^l, p_j^l < p_i^x) \right) \quad (3)$$

where $s_{j_{max}}$ is the maximum length transaction in any job τ_j^l with original priority lower than p_i^x . The right hand side of (6) is the minimum of the m maximum transactional lengths in all jobs with lower priority than τ_i^x .

PROOF. τ_i^x is blocked when it is initially released and all processors are busy with lower priority jobs with non-preemptive transactions. Although τ_i^x can be preempted by higher priority jobs, τ_i^x cannot be blocked after it is released. If τ_i^x is preempted by a higher priority job τ_j^y , then, when τ_j^y finishes execution, the underlying scheduler will not choose a lower priority job than τ_i^x before τ_i^x . So, after τ_i^x is released, there is no chance for any transaction s_u^v belonging to a lower priority job than τ_i^x to run before τ_i^x . Thus, s_u^v cannot join the m_set before τ_i^x finishes. Consequently, the worst case blocking time for τ_i^x occurs when the maximum length m transactions in lower priority jobs than τ_i^x are executing non-preemptively. After the minimum length transaction in the m_set finishes, the underlying scheduler will choose τ_i^x or a higher priority job to run. Claim follows. \square

CLAIM 3. The response time of any job τ_i^x during an interval $L \leq T_i$ under FBLT is upper bounded by:

$$R_i^{up} = c_i + RC_{to}(L) + D(\tau_i^x) + \left\lceil \frac{1}{m} \sum_{\forall j \neq i} W_{ij}(R_i^{up}) \right\rceil \quad (4)$$

where $RC_{to}(L)$ is calculated by (5), $D(\tau_i^x)$ is calculated by (6), and $W_{ij}(R_i^{up})$ is calculated by (11) in [El-Shambakey and Ravindran 2012b] for G-EDF, and (17) in [El-Shambakey and Ravindran 2012b] for G-RMA schedulers. (11) and (17) in [El-Shambakey and Ravindran 2012b] inflates c_j of any job $\tau_j^y \neq \tau_i^x$, $p_j^y > p_i^x$ by the retry cost of transactions in τ_j^y .

PROOF. The response time of a job is calculated directly from FBLT's behavior. The response time of any job τ_i^x is the sum of its worst case execution time c_i , plus the retry cost of transactions in τ_i^x ($RC_{to}(L)$), plus the blocking time of τ_i^x ($D(\tau_i^x)$), and the workload interference of higher priority jobs. The workload interference of higher priority jobs scheduled by G-EDF is calculated by (11) in [El-Shambakey and Ravindran 2012b], and by (17) in [El-Shambakey and Ravindran 2012b] for G-RMA. Claim follows. \square

FBLT

:11

We now derive an upper bound on the retry cost of any job τ_i^x under FBLT during an interval $L \leq T_i$. Since all tasks are sporadic (i.e., each task τ_i has a minimum period T_i), T_i is the maximum study interval for each task τ_i .

CLAIM 4.

The total retry cost for any job τ_i^x under FBLT due to: 1) conflicts between its transactions and transactions of other jobs during an interval $L \leq T_i$. 2) release of higher priority jobs, is upper bounded by:

$$RC_{to}(L) \leq \sum_{\forall s_i^k \in s_i} \left(\delta_i^k \text{len}(s_i^k) + \sum_{\forall s_{iz}^k \in \chi_i^k} \text{len}(s_{iz}^k) \right) + RC_{re}(L) \quad (5)$$

where χ_i^k is the set of at most $m - 1$ maximum length transactions conflicting directly or indirectly (through transitive retry) with s_i^k . Each transaction $s_{iz}^k \in \chi_i^k$ belongs to a distinct task τ_j . $RC_{re}(L)$ is the retry cost resulting from release of higher priority jobs which preempt τ_i^x . $RC_{re}(L)$ is calculated by (6.8) in [El-Shambakey 2012] for G-EDF, and (6.10) in [El-Shambakey 2012] for G-RMA.

PROOF.

By definition of FBLT, $s_i^k \in \tau_i^x$ can be aborted at maximum δ_i^k times before s_i^k joins m_set . Before joining m_set , s_i^k can be aborted due to higher priority transactions, or transactions in the m_set . Original priority of transactions in m_set can be of higher or lower priority than p_i^x . Thus, the maximum time s_i^k is aborted before joining m_set occurs if s_i^k is aborted for δ_i^k times. Transactions preceding s_i^k in m_set can conflict directly with s_i^k , or indirectly through transitive retry. The worst case scenario for s_i^k after joining m_set occurs if s_i^k is preceded by $m - 1$ maximum length conflicting transactions. Hence, in worst case, s_i^k has to wait for the previous $m - 1$ transactions to commit first. Priority of s_i^k after joining m_set is higher than any real-time job. So, s_i^k is not aborted by any job. If s_i^k has not joined m_set yet, and a higher priority job τ_j^y is released while s_i^k is running, then s_i^k may be aborted if τ_j^y has conflicting transactions with s_i^k . τ_j^y causes only one abort in τ_i^x because τ_j^y preempts τ_i^x only once. If s_i^k has already joined m_set , then s_i^k cannot be aborted by release of higher priority jobs. So, the maximum number of abort times to transactions in τ_i^x due to release of higher priority jobs is less or equal to number of interfering higher priority jobs to τ_i^x . Claim follows.

□

CLAIM 5.

The blocking time for a job τ_i^x due to lower priority jobs is upper bounded by:

$$D(\tau_i^x) = \min \left(\max_1^m (s_{j_{max}, \forall \tau_j^l, p_j^l < p_i^x}) \right) \quad (6)$$

where $s_{j_{max}}$ is the maximum length transaction in any job τ_j^l with original priority lower than p_i^x . The right hand side of (6) is the minimum of the m maximum transactional lengths in all jobs with lower priority than τ_i^x .

PROOF.

τ_i^x is blocked when it is initially released and all processors are busy with lower priority jobs with non-preemptive transactions. Although τ_i^x can be preempted by higher priority jobs, τ_i^x cannot be blocked after it is released. If τ_i^x is preempted by a higher

:12

priority job τ_j^y , then τ_j^y finishes execution, the underlying scheduler will not choose a lower priority job than τ_i^x before τ_i^x . So, after τ_i^x is released, there is no chance for any transaction s_u^v belonging to a lower priority job than τ_i^x to run before τ_i^x . Thus, s_u^v cannot join m_set before τ_i^x finishes. Consequently, the worst case blocking time for τ_i^x occurs when the maximum length m transactions in lower priority jobs than τ_i^x are executing non-preemptively. After the minimum length transaction in the m_set finishes, the underlying scheduler will choose τ_i^x or a higher priority job to run. Claim follows.

□

CLAIM 6.

Response time of any job τ_i^x during an interval $L \leq T_i$ under FBLT is upper bounded by

$$R_i^{up} = c_i + RC_{to}(L) + D(\tau_i^x) + \left\lceil \frac{1}{m} \sum_{\forall j \neq i} W_{ij}(R_i^{up}) \right\rceil \quad (7)$$

where $RC_{to}(L)$ is calculated by (5), $D(\tau_i^x)$ is calculated by (6), and $W_{ij}(R_i^{up})$ is calculated by (11) in [El-Shambakey and Ravindran 2012b] for G-EDF, and (17) in [El-Shambakey and Ravindran 2012b] for G-RMA. (11) and (17) in [El-Shambakey and Ravindran 2012b] inflates c_j of any job $\tau_j^y \neq \tau_i^x$, $p_j^y > p_i^x$ by retry cost of transactions in τ_j^y .

PROOF.

Response time of any job τ_i^x is calculated directly from FBLT's behaviour. Response time of any job τ_i^x is the sum of its worst case execution time c_i , plus retry cost of transactions in τ_i^x ($RC_{to}(L)$), plus blocking time of τ_i^x ($D(\tau_i^x)$), and the workload interference of higher priority jobs. Workload interference of higher priority jobs scheduled by G-EDF is calculated by (11) in [El-Shambakey and Ravindran 2012b], and by (17) in [El-Shambakey and Ravindran 2012b] for G-RMA. Claim follows.

□

7. SCHEDULABILITY COMPARISON

We now (formally) compare the schedulability of G-EDF (G-RMA) with FBLT against ECM, RCM, LCM, PNF, and lock-free synchronization [El-Shambakey and Ravindran 2012b; 2012a; Devi et al. 2006; El-Shambakey 2012]. Such a comparison will reveal when FBLT outperforms the others. Toward this, we compare the total utilization under G-EDF (G-RMA)/FBLT with that under the other synchronization methods. In this comparison, we use the inflated execution time of the task, which is the sum of the worst-case execution time of the task and its retry cost, in the utilization calculation of the task.

Note that, for a job τ_i^x , no processor is available during its blocking time. Since each processor is busy with some job other than τ_i^x , $D(\tau_i^x)$ is not added to the inflated execution time of τ_i^x . Hence, $D(\tau_i^x)$ is not added to the utilization calculation of τ_i^x .

Let $RC_A(T_i)$ and $RC_B(T_i)$ denote the retry cost of a job τ_i^x during T_i using the synchronization method A and synchronization method B , respectively. Now, schedulability of A is comparable to B if:

$$\sum_{\forall \tau_i} \frac{c_i + RC_A(T_i)}{T_i} \leq \sum_{\forall \tau_i} \frac{c_i + RC_B(T_i)}{T_i}$$

FBLT

:13

$$\sum_{\forall \tau_i} \frac{RC_A(T_i)}{T_i} \leq \sum_{\forall \tau_i} \frac{RC_B(T_i)}{T_i} \quad (8)$$

7.1. FBLT vs. ECM

CLAIM 7. *The schedulability of FBLT is equal to or better than ECM's when the maximum abort number of any preemptive transaction s_i^k is less than or equal to the number of transactions directly conflicting with s_i^k in all other jobs with higher priority than τ_i 's current job.*

PROOF.

By substituting $RC_A(T_i)$ and $RC_B(T_i)$ in (8) with (5) and (6.7) in [El-Shambakey 2012], respectively, we get:

$$\begin{aligned} & \sum_{\forall \tau_i} \frac{\sum_{s_i^k \in s_i} \left(\delta_i^k \text{len}(s_i^k) + \sum_{s_{iz}^k \in \chi_i^k} \text{len}(s_{iz}^k) \right) + RC_{re}(T_i)}{T_i} \\ & \leq \sum_{\forall \tau_i} \frac{\left(\sum_{\tau_j \in \gamma_i^{ex}} \sum_{\theta \in \theta_i^{ex}} \left(\left\lceil \frac{T_i}{T_j} \right\rceil \sum_{s_j^h \in \bar{s}_j^h(\theta)} \text{len}(s_j^h(\theta) + s_{max}^j(\theta)) \right) \right) + RC_{re}(T_i)}{T_i} \end{aligned} \quad (9)$$

Let $\theta_i^{ex} = \theta_i + \theta_i^*$, where θ_i^* is the set of objects not accessed directly by τ_i but can cause transactions in τ_i to retry due to transitive retry. Let $\gamma_i^{ex} = \gamma_i + \gamma_i^*$, where γ_i^* is the set of tasks that access objects in θ_i^* . $\bar{s}_j^h(\theta)$ can access multiple objects, so $s_{max}^j(\theta)$ is the maximum length transaction conflicting with $\bar{s}_j^h(\theta)$. $s_j^h(\theta)$ is included only once for all $\theta \in \Theta_j^h$. Each $\theta \in \theta_i^{ex}$ has its own $s_{max}^j(\theta)$. But s_i^h can access multiple objects, denoted as Θ_j^h . So, $s_{max}^j(\theta)$ is replaced by $s_{max}^j(\Theta_j^h)$, where $s_{max}^j(\Theta_j^h) = \max\{s_{max}^j(\theta), \forall \theta \in \Theta_j^h\}$. $s_{max}^j(\Theta_j^h)$ is included once for each $\theta \in \theta_i$.

Each job τ_i^x has the same interference pattern from higher priority jobs, τ_j^h , under FBLT and ECM. Hence, $RC_{re}(T_i)$ for τ_i^x is the same under FBLT and ECM. Consequently, (9) becomes:

$$\begin{aligned} & \sum_{\forall \tau_i} \frac{\sum_{s_i^k \in s_i} \left(\delta_i^k \text{len}(s_i^k) + \sum_{s_{iz}^k \in \chi_i^k} \text{len}(s_{iz}^k) \right)}{T_i} \\ & \leq \sum_{\forall \tau_i} \frac{\left(\sum_{\tau_j \in \gamma_i^{ex}} \sum_{s_j^h \in \bar{s}_j^h(\Theta_j^h), \Theta_j^h \in \theta_i^{ex}} \left(\left\lceil \frac{T_i}{T_j} \right\rceil \text{len}(s_j^h(\Theta_j^h) + s_{max}^j(\Theta_j^h)) \right) \right)}{T_i} \end{aligned} \quad (10)$$

Although different s_i^k s can have common conflicting transactions \bar{s}_j^h , no more than one s_i^k can be preceded by the same \bar{s}_j^h in the m_set . This happens because transactions in the m_set are non-preemptive. The original priority of transactions preceding s_i^k in the m_set can be lower or higher than the original priority of s_i^k . Since under G-EDF, τ_j can have at least one job of higher priority than τ_i^x , $\left\lceil \frac{T_i}{T_j} \right\rceil \geq 1$. Thus, each one of the s_{iz}^k term in the left hand side of (10) is included in one of the $\bar{s}_j^h(\theta)$ term in the right hand side of (10). Now, (10) holds if:

$$\begin{aligned} & \sum_{\forall \tau_i} \frac{\sum_{s_i^k \in s_i} \delta_i^k \text{len}(s_i^k)}{T_i} \\ & \leq \sum_{\forall \tau_i} \frac{\sum_{\tau_j \in \gamma_i^{ex}} \sum_{s_j^h \in \bar{s}_j^h(\Theta_j^h), \Theta_j^h \in \theta_i^{ex}} \left(\left\lceil \frac{T_i}{T_j} \right\rceil \text{len}(s_{max}^j(\Theta_j^h)) \right)}{T_i} \end{aligned} \quad (11)$$

:14

Since FBLT is required to bound the effect of transitive retry, only θ_i (not the whole θ_i^{ex}) will be considered in (11). Thus, ECM acts as if there were no transitive retry. Consequently, (11) holds if:

$$\begin{aligned} & \sum_{\forall \tau_i} \frac{\sum_{s_i^k \in s_i} \delta_i^k \text{len}(s_i^k)}{T_i} \\ & \leq \sum_{\forall \tau_i} \frac{\sum_{\tau_j \in \gamma_i} \sum_{s_j^h(\Theta), \Theta \in (\theta_i \cap \theta_j^h)} \left(\left\lceil \frac{T_i}{T_j} \right\rceil \text{len}(s_{max}^j(\Theta)) \right)}{T_i} \end{aligned} \quad (12)$$

where $s_{max}^j(\Theta) \leq s_{max}^j(\Theta_j^h)$.

For each $s_i^k \in s_i$, there are a set of zero or more $s_j^h(\Theta) \in \tau_j, \forall \tau_j \neq \tau_i$ that are conflicting with s_i^k . Assuming this set of transactions conflicting with s_i^k is denoted as $\nu_i^k = \{s_j^h(\Theta) \in \tau_j : (\Theta \in \theta_i \cap \theta_j^h) \wedge (\forall \tau_j \neq \tau_i) \wedge (s_j^h(\Theta) \notin \nu_i^l, l \neq k)\}$.

The last condition $s_j^h(\Theta) \notin \nu_i^l, l \neq k$ in the definition of ν_i^k ensures that common transactions s_j^h that can conflict with more than one transaction $s_i^k \in \tau_i$ are split among different $\nu_i^k, k = 1, \dots, |s_i|$. This condition is necessary, because in ECM, no two or more transactions of τ_i^x can be aborted by the same transaction of τ_j^h , where $p_j^h > p_i^x$. By substitution of ν_i^k in (12), we get:

$$\begin{aligned} & \sum_{\forall \tau_i} \frac{\sum_{s_i^k \in s_i} \delta_i^k \text{len}(s_i^k)}{T_i} \\ & \leq \sum_{\forall \tau_i} \frac{\left(\sum_{k=1}^{|s_i|} \sum_{s_j^h(\Theta) \in \nu_i^k} \left(\left\lceil \frac{T_i}{T_j} \right\rceil \text{len}(s_{max}^j(\Theta)) \right) \right)}{T_i} \end{aligned} \quad (13)$$

(13) holds if for each $s_i^k \in \tau_i$:

$$\delta_i^k \leq \frac{\sum_{s_j^h(\Theta) \in \nu_i^k} \left(\left\lceil \frac{T_i}{T_j} \right\rceil \text{len}(s_{max}^j(\Theta)) \right)}{\text{len}(s_i^k)} \quad (14)$$

Since $\text{len}(s_{max}^j(\Theta)) \geq \text{len}(s_i^k)$, (14) holds if $\delta_i^k \leq \sum_{s_j^h(\Theta) \in \nu_i^k} \left\lceil \frac{T_i}{T_j} \right\rceil \cdot \sum_{s_j^h(\Theta) \in \nu_i^k} \left\lceil \frac{T_i}{T_j} \right\rceil$ is the maximum number of transactions directly conflicting with s_i^k in all jobs with higher priority than τ_i . Claim follows. \square

7.2. FBLT vs. RCM

CLAIM 8. *The schedulability of FBLT is equal to or better than RCM's if*

$$\delta_i^k \leq \left(\sum_{s_j^h(\Theta) \in \nu_i^k} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \right) - \sum_{u=1, s_{u_{max}} \in \epsilon}^{\min(n,m)-1} s_{u_{max}}$$

$\sum_{s_j^h(\Theta) \in \nu_i^k} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right)$ is number of transactions directly conflicting with s_i^k in all jobs with higher priority than τ_i . $\sum_{u=1, s_{u_{max}} \in \epsilon}^{\min(n,m)-1} s_{u_{max}}$ is the sum of the maximum $m - 1$ transactional lengths in all tasks

PROOF. By substituting $RC_A(T_i)$ and $RC_B(T_i)$ in (8) with (5) and (6.9) in [El-Shambakey 2012], respectively, we get:

$$\sum_{\forall \tau_i} \frac{\sum_{s_i^k \in s_i} \left(\delta_i^k \text{len}(s_i^k) + \sum_{s_i^k \in \chi_i^k} \text{len}(s_{iz}^k) \right) + RC_{re}(T_i)}{T_i} \quad (15)$$

FBLT

:15

$$\leq \sum_{\forall \tau_i} \frac{\left(\sum_{\forall \tau_j^* \in \gamma_i^{ex}} \sum_{\forall \theta \in \theta_i^{ex}} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \sum_{\forall s_j^h(\theta)} \text{len}(s_j^h(\theta) + s_{max}^j(\theta)) \right) + RC_{re}(T_i)}{T_i}$$

where $\tau_j^* = \{\tau_j : (\tau_j \neq \tau_i) \wedge (p_j > p_i)\}$.

Let $\theta_i^{ex} = \theta_i + \theta_i^*$, where θ_i^* is the set of objects not directly accessed by any job of τ_i , but can cause transactions in τ_i to retry due to transitive retry. Let $\gamma_i^{ex} = \gamma_i + \gamma_i^*$, where γ_i^* is the set of tasks that access objects in θ_i^* . $s_j^h(\theta)$ can access multiple objects, so $s_{max}^j(\theta)$ is the maximum length transaction conflicting with $s_j^h(\theta)$. $s_j^h(\theta)$ is included only once for all $\theta \in \Theta_j^h$. Each $\theta \in \Theta_j^h$ has its own $s_{max}^j(\theta)$. But s_i^h can access multiple objects, denoted as Θ_j^h . So, $s_{max}^j(\theta)$ is replaced by $s_{max}^j(\Theta_j^h)$, where $s_{max}^j(\Theta_j^h) = \max\{s_{max}^j(\theta), \forall \theta \in \Theta_j^h\}$. $s_{max}^j(\Theta_j^h)$ is included once for each $\theta \in \theta_i$.

Each τ_i^x has the same interference pattern from higher priority jobs, τ_j^h , under FBLT and RCM. Hence, $RC_{re}(T_i)$ for τ_i^x is the same under FBLT and RCM. Consequently, (15) becomes:

$$\begin{aligned} & \sum_{\forall \tau_i} \frac{\sum_{s_i^k \in s_i} \left(\delta_i^k \text{len}(s_i^k) + \sum_{s_{iz}^k \in \chi_i^k} \text{len}(s_{iz}^k) \right)}{T_i} \\ & \leq \sum_{\forall \tau_i} \frac{\sum_{\tau_j^* \in \gamma_i^{ex}} \sum_{s_j^h(\Theta_j^h), \Theta_j^h \in \theta_i^{ex}} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \text{len}(s_j^h(\Theta_j^h) + s_{max}^j(\Theta_j^h))}{T_i} \end{aligned} \quad (16)$$

Although different s_i^k s can have common conflicting transactions s_j^h , no more than one s_i^k can be preceded by the same s_j^h in the m_set . This happens because transactions in the m_set are non-preemptive. The original priority of transactions preceding s_i^k in the m_set can be of lower or higher priority than the original priority of s_i^k . Under G-RMA, $p_j > p_i$, which means that $T_j \leq T_i$. Therefore, $\left\lceil \frac{T_i}{T_j} \right\rceil \geq 1$. For each $s_i^k \in s_i$, there are a set of zero or more $s_j^h(\Theta_j^h) \in \tau_j^*$ that are conflicting with s_i^k . Assuming this set of transactions conflicting with s_i^k is denoted as $\nu_i^k = \{s_j^h(\Theta_j^h) \in \tau_j^* : (\Theta_j^h \in \theta_i^{ex}) \wedge (s_j^h(\Theta_j^h) \notin \nu_i^l, l \neq k)\}$.

The last condition $s_j^h(\theta) \notin \nu_i^l, l \neq k$ in the definition of ν_i^k ensures that common transactions s_j^h that can conflict with more than one transaction $s_i^k \in \tau_i$ are split among different $\nu_i^k, k = 1, \dots, |s_i|$. This condition is necessary, because in RCM, no two or more transactions of τ_i^x can be aborted by the same transaction of τ_j^h , where $p_j^h > p_i^x$. By substitution of ν_i^k in (16), we get:

$$\begin{aligned} & \sum_{\forall \tau_i} \frac{\sum_{s_i^k \in s_i} \left(\delta_i^k \text{len}(s_i^k) + \sum_{s_{iz}^k \in \chi_i^k} \text{len}(s_{iz}^k) \right)}{T_i} \\ & \leq \sum_{\forall \tau_i} \frac{\left(\sum_{k=1}^{|s_i|} \sum_{s_j^h(\Theta_j^h) \in \nu_i^k} \left(\left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \text{len}(s_j^h(\Theta_j^h) + s_{max}^j(\Theta_j^h)) \right) \right)}{T_i} \end{aligned} \quad (17)$$

s_j^h belongs to higher priority jobs than τ_i . s_{max}^j belongs to higher priority jobs than τ_i or τ_i itself. s_{max}^j has a lower priority than τ_j . Transactions in the m_set can belong to jobs with original priority higher or lower than τ_i . Thus, (17) holds if for each $s_i^k \in \tau_i$:

$$\delta_i^k \text{len}(s_i^k) \leq \left(\sum_{s_j^h(\Theta_j^h) \in \nu_i^k} \left(\left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \text{len}(s_j^h(\Theta_j^h) + s_{max}^j(\Theta_j^h)) \right) \right) - \sum_{s_{iz}^k \in \chi_i^k} \text{len}(s_{iz}^k) \quad (18)$$

:16

Then,

$$\delta_i^k \leq \left(\sum_{\bar{s}_j^h(\Theta_j^h) \in \nu_i^k} \left(\left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \text{len} \left(\frac{\bar{s}_j^h(\Theta_j^h) + s_{max}^j(\Theta_j^h)}{s_i^k} \right) \right) \right) - \sum_{s_{iz}^k \in \chi_i^k} \text{len} \left(\frac{s_{iz}^k}{s_i^k} \right) \quad (19)$$

Let $\epsilon = \{s_{u_{max}} : (1 \leq u \leq n) \wedge (s_{u1_{max}} \geq s_{u2_{max}}, u1 < u2)\}$, where n is the number of tasks, and $s_{u_{max}}$ is the maximum transactional length in any job of τ_u . Thus, ϵ is the set of maximum transactional lengths of all tasks in non-increasing order. Each $s_{u_{max}} \in \epsilon$ belongs to a distinct task. Thus, $\sum_{s_{iz}^k \in \chi_i^k} \text{len} \left(\frac{s_{iz}^k}{s_i^k} \right) \leq \sum_{u=1, s_{u_{max}} \in \epsilon}^{min(n,m)-1} s_{u_{max}}$. $\sum_{u=1, s_{u_{max}} \in \epsilon}^{min(n,m)-1} s_{u_{max}}$ is the sum of at most maximum $m - 1$ transactional lengths of all tasks. $|\chi_i^k| \leq m - 1$ and $\text{len}(s_{max}^j(\Theta_j^h)) \geq \text{len}(s_i^k)$. So, (19) holds if:

$$\delta_i^k \leq \left(\sum_{\bar{s}_j^h(\Theta_j^h) \in \nu_i^k} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \right) - \sum_{u=1, s_{u_{max}} \in \epsilon}^{min(n,m)-1} s_{u_{max}} \quad (20)$$

To bound the effect of transitive retry, only objects that belong to θ_i (not whole θ_i^{ex}) will be considered. Thus, RCM acts as if there were no transitive retry. Thus, ν_i^k is modified to $\bar{\nu}_i^k = \{\bar{s}_j^h(\Theta) \in \tau_j^* : (\Theta \in \Theta_j^h \cap \theta_i) \wedge (\bar{s}_j^h(\Theta) \notin \nu_i^l, l \neq k)\}$. Since $\bar{\nu}_i^k \subseteq \nu_i^k$, (20) still holds if ν_i^k is replaced with $\bar{\nu}_i^k$. Consequently, (20) holds if:

$$\delta_i^k \leq \left(\sum_{\bar{s}_j^h(\Theta) \in \bar{\nu}_i^k} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \right) - \sum_{u=1, s_{u_{max}} \in \epsilon}^{min(n,m)-1} s_{u_{max}} \quad (21)$$

$\sum_{\bar{s}_j^h(\Theta) \in \bar{\nu}_i^k} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right)$ represents the number of transactions directly conflicting with s_i^k in all jobs with higher priority than τ_i . Claim follows. \square

7.3. FBLT vs. G-EDF/LCM

CLAIM 9. *The schedulability of FBLT is equal to or better than G-EDF/LCM's when*

$$\delta_i^k \leq \left(\sum_{\bar{s}_j^h(\Theta) \in \nu_i^k} \left(\left\lceil \frac{T_i}{T_j} \right\rceil \alpha_{max}^{jh} \right) \right)$$

α_{max}^{jh} is the maximum α with which s_j^h can conflict with the maximum length transaction sharing objects with s_i^k and s_j^h

PROOF. 7. The proof is similar to that of Claim 7 and is therefore skipped for brevity. It can be found in [El-Shambakey and Ravindran 2011]. \square

7.4. FBLT vs. G-RMA/LCM

CLAIM 10. *The schedulability of FBLT is equal to or better than G-RMA/LCM's when*

$$\delta_i^k \leq \left(\sum_{\bar{s}_j^h(\Theta) \in \bar{\nu}_i^k} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \alpha_{max}^{jh} \right) - \sum_{u=1, s_{u_{max}} \in \epsilon}^{min(n,m)-1} s_{u_{max}}$$

FBLT

:17

$\left(\sum_{s_j^h(\Theta) \in \nu_i^k} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1\right)\right)$ is the sum of the total number each transaction s_j^h can directly conflict with s_i^k . $\alpha_{max}^{j\bar{h}}$ is the maximum α with which s_j^h can conflict with the maximum length transaction sharing objects with s_i^k and s_j^h .

PROOF. The proof is similar to that of Claim 8 and is therefore skipped for brevity. It can be found in [El-Shambakey and Ravindran 2011].

□

7.5. FBLT vs. PNF

CLAIM 11. Let $\rho_i^j(k) = \left(\sum_{s_j^h(\Theta) \in \nu_i^k(j)} \text{len}(s_j^h(\Theta))\right) - s_{i_{max}}$, $\tau_j \in \gamma_i^k$. $\rho_i^j(k)$ is the difference between the sum of transactional lengths of all transactions in τ_j conflicting with s_i^k , and the maximum length transaction in τ_i . Let $\sum_{u=1, s_{u_{max}} \in \epsilon}^{min(n,m)-1} s_{u_{max}}$ be sum of at most maximum $m-1$ transactional lengths in all tasks. Schedulability of FBLT is better or equal to PNF's when

$$\delta_i^k \leq \left(\sum_{\forall \tau_j \in \gamma_i^k} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1\right) \text{len} \left(\frac{\left(\sum_{s_j^h(\Theta) \in \nu_i^k(j)} \text{len}(s_j^h(\Theta))\right) - s_{i_{max}}}{s_i^k} \right)\right) - \sum_{u=1, s_{u_{max}} \in \epsilon}^{min(n,m)-1} s_{u_{max}}$$

PROOF. By substituting $RC_A(T_i)$ and $RC_B(T_i)$ in (8) with (5) and (6.1) in [El-Shambakey 2012], respectively, we get:

$$\begin{aligned} & \sum_{\forall \tau_i} \frac{\sum_{s_i^k \in s_i} \left(\delta_i^k \text{len}(s_i^k) + \sum_{s_{iz}^k \in \chi_i^k} \text{len}(s_{iz}^k)\right) + RC_{re}(T_i)}{T_i} \\ & \leq \sum_{\forall \tau_i} \frac{\sum_{\tau_j \in \gamma_i} \sum_{\theta \in \theta_i} \left(\left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1\right) \sum_{s_j^h(\Theta) \in \nu_i^k} \text{len}(s_j^h(\Theta))\right)}{T_i} \end{aligned} \quad (22)$$

$s_j^h(\theta)$ can access multiple objects. $s_j^h(\theta)$ is included only once for all objects accessed by it. $RC_{re}(T_i)$ is given by (6.8) in [El-Shambakey 2012] in case of G-EDF, and (6.10) in [El-Shambakey 2012] in case of G-RMA. Substituting $RC_{re}(T_i)$ given by (6.8) and (6.10) in [El-Shambakey 2012] with $RC_{re}(T_i) = \sum_{\tau_j \in \gamma_i} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1\right) s_{i_{max}}$, we ensure correctness of (22). If τ_j has no shared objects with τ_i , then the release of any higher priority job τ_j^y will not abort any transaction in any job of τ_i . Thus, (22) holds if:

$$\begin{aligned} & \sum_{\forall \tau_i} \frac{\sum_{s_i^k \in s_i} \left(\delta_i^k \text{len}(s_i^k) + \sum_{s_{iz}^k \in \chi_i^k} \text{len}(s_{iz}^k)\right)}{T_i} \\ & \leq \sum_{\forall \tau_i} \frac{\sum_{\tau_j \in \gamma_i} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1\right) \left(\left(\sum_{s_j^h(\Theta) \in \nu_i^k} \text{len}(s_j^h(\Theta))\right) - s_{i_{max}}\right)}{T_i} \end{aligned} \quad (23)$$

For each $s_i^k \in s_i$, there are a set of zero or more $s_j^h(\Theta) \in \tau_j, \forall \tau_j \neq \tau_i$ that are conflicting with s_i^k . Assuming this set of transactions conflicting with s_i^k is denoted as $\nu_i^k(j) = \left\{s_j^h(\Theta) \in \tau_j : (\Theta \in \theta_i) \wedge (\tau_j \neq \tau_i) \wedge (s_j^h(\Theta) \notin \nu_i^l, l \neq k)\right\}$. The last condition $s_j^h(\Theta) \notin \nu_i^l, l \neq k$ in the definition of ν_i^k ensures that common transactions s_j^h that can conflict with more than one transaction $s_i^k \in \tau_i$ are split among different $\nu_i^k(j), k = 1, \dots, |s_i|$. This condition is necessary, because in PNF, no two or more transactions of τ_i^x can be aborted by the same transaction of τ_j^h .

:18

Let γ_i^k be the subset of γ_i that contains tasks with transactions conflicting directly with s_i^k . By substitution of ν_i^k and γ_i in (23) by $\nu_i^k(j)$ and γ_i^k , (23) holds if for each s_i^k :

$$\begin{aligned} & \delta_i^k + \sum_{s_{iz}^k \in \chi_i^k} \text{len}\left(\frac{s_{iz}^k}{s_i^k}\right) \\ & \leq \sum_{\tau_j \in \gamma_i^k} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \left(\left(\sum_{\forall s_j^h(\Theta) \in \nu_i^k(j)} \text{len}\left(\frac{s_j^h(\Theta)}{s_i^k}\right) \right) - \text{len}\left(\frac{s_{imax}^k}{s_i^k}\right) \right) \end{aligned} \quad (24)$$

Non-preemptive transactions preceding s_i^k in the m -set can directly or indirectly conflict with s_i^k . Under PNF, transactions can only directly conflict with s_i^k . Thus, s_{iz}^k on the left hand side of (24) is not necessarily included in $\bar{s}_j^h(\Theta)$ on the right hand side of (24). Let $\epsilon = \{s_{u_{max}} : (1 \leq u \leq n) \wedge (s_{u1_{max}} \geq s_{u2_{max}}, u1 < u2)\}$, where n is the number of tasks, and $s_{u_{max}}$ is the maximum transactional length in any job of τ_u . Thus, ϵ is the set of maximum transactional lengths of all tasks in non-increasing order. Each $s_{u_{max}} \in \epsilon$ belongs to a distinct task. Thus, $\sum_{s_{iz}^k \in \chi_i^k} \text{len}\left(\frac{s_{iz}^k}{s_i^k}\right) \leq \sum_{u=1, s_{u_{max}} \in \epsilon}^{min(n,m)-1} s_{u_{max}}$. $\sum_{u=1, s_{u_{max}} \in \epsilon}^{min(n,m)-1} s_{u_{max}}$ is the sum of at most maximum $m-1$ transactional lengths of all tasks. $|\chi_i^k| \leq m-1$. Then (24) holds if:

$$\begin{aligned} \delta_i^k & \leq \left(\sum_{\tau_j \in \gamma_i^k} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \text{len}\left(\frac{\left(\sum_{\forall s_j^h(\Theta) \in \nu_i^k(j)} \text{len}\left(\bar{s}_j^h(\Theta)\right) \right) - s_{imax}^k}{s_i^k}\right) \right) \\ & \quad - \sum_{u=1, s_{u_{max}} \in \epsilon}^{min(n,m)-1} s_{u_{max}} \end{aligned}$$

Since $\rho_i^j(k) = \left(\sum_{\forall s_j^h(\Theta) \in \nu_i^k(j)} \text{len}\left(\bar{s}_j^h(\Theta)\right) \right) - s_{imax}^k$, $\tau_j \in \gamma_i^k$, and $\left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right)$ is the total number of jobs of τ_j interfering with τ_i , claim follows. \square

7.6. FBLT vs. Lock-free

CLAIM 12. *Under G-EDF and G-RMA, the schedulability of FBLT is equal or better than that under lock-free synchronization if $s_{max} \leq r_{max}$. If transactions execute in FIFO order (i.e., $\delta_i^k = 0, \forall s_i^k$) and contention is high, s_{max} can be much larger than r_{max} .*

PROOF. Lock-free synchronization [Devi et al. 2006; Herlihy 2006] allows only one object to be synchronized at a given time (e.g., a lock-free stack). Thus, for comparing FBLT's schedulability with lock-free, we limit the number of accessed objects per transaction under FBLT to one.

By substituting $RC_A(T_i)$ and $RC_B(T_i)$ in (8) with (5) and (6.17) in [El-Shambakey 2012], respectively, we get:

$$\begin{aligned} & \sum_{\tau_i} \frac{\sum_{s_i^k \in s_i} \left(\delta_i^k \text{len}(s_i^k) + \sum_{s_{iz}^k \in \chi_i^k} \text{len}(s_{iz}^k) \right) + RC_{re}(T_i)}{T_i} \\ & \leq \sum_{\tau_i} \frac{\left(\sum_{\tau_j \in \gamma_i} \left(\left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \beta_{i,j} r_{max} \right) \right) + RC_{re}(T_i)}{T_i} \end{aligned} \quad (25)$$

where $\beta_{i,j}$ is the number of retry loops of τ_j that access the same objects as accessed by any retry loop of τ_i [Devi et al. 2006] and r_{max} is the maximum execution cost of a single iteration of any retry loop of any task [Devi et al. 2006].

FBLT

:19

For G-EDF (G-RMA), any job τ_i^x under FBLT has the same pattern of interference from higher priority jobs as ECM (RCM), respectively. $RC_{re}(T_i)$ for ECM, RCM, and lock-free are given by Claims 25, 26, and 27 in [El-Shambakey 2012], respectively. $RC_{re}(T_i) = \left\lceil \frac{T_i}{T_j} \right\rceil s_{i_{max}}, \forall \tau_j \neq \tau_i$ for G-EDF/FBLT and G-RMA/FBLT. $RC_{re}(T_i) = \left\lceil \frac{T_i}{T_j} \right\rceil r_{i_{max}}, \forall \tau_j \neq \tau_i$ for G-EDF/lock-free and G-RMA/lock-free. (31) becomes:

$$\begin{aligned} & \sum_{\forall \tau_i} \frac{\sum_{\forall s_i^k \in s_i} \left(\delta_i^k \text{len}(s_i^k) + \sum_{s_i^k \in \chi_i^k} \text{len}(s_i^k) \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil s_{i_{max}}}{T_i} \\ & \leq \sum_{\forall \tau_i} \frac{\left(\sum_{\forall \tau_j \in \gamma_i} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \beta_{i,j} r_{max} \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil r_{i_{max}}}{T_i} \end{aligned} \quad (26)$$

Since $s_{max} \geq s_{i_{max}}, \text{len}(s_i^k), \text{len}(s_{iz}^k), \forall i, z, k$ and $r_{max} \geq r_{i_{max}}$ (32) holds if:

$$\begin{aligned} & \sum_{\forall \tau_i} \frac{\left(\left(\sum_{\forall s_i^k \in s_i} (\delta_i^k + |\chi_i^k|) \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil \right) s_{max}}{T_i} \\ & \leq \sum_{\forall \tau_i} \frac{\left(\left(\sum_{\forall \tau_j \in \gamma_i} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \beta_{i,j} \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil \right) r_{max}}{T_i} \end{aligned} \quad (27)$$

(33) holds if for each τ_i :

$$\begin{aligned} & \left(\left(\sum_{\forall s_i^k \in s_i} (\delta_i^k + |\chi_i^k|) \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil \right) s_{max} \\ & \leq \left(\left(\sum_{\forall \tau_j \in \gamma_i} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \beta_{i,j} \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil \right) r_{max} \end{aligned} \quad (28)$$

$$\frac{s_{max}}{r_{max}} \leq \frac{\left(\sum_{\forall \tau_j \in \gamma_i} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \beta_{i,j} \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil}{\left(\sum_{\forall s_i^k \in s_i} (\delta_i^k + |\chi_i^k|) \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil} \quad (29)$$

It appears from (35) that as δ_i^k and $|\chi_i^k|$ increases, s_{max}/r_{max} decreases. So, to get the lower bound on s_{max}/r_{max} , let $\sum_{\forall s_i^k \in s_i} (\delta_i^k + |\chi_i^k|)$ reach its maximum value. This maximum value is the total number of interfering transactions belonging to any job $\tau_j^l, j \neq i$. The priority of τ_j^l can be higher or lower than the current instance of τ_i . Beyond this maximum value, there will be no more transactions that conflict with s_i^k . Thus, higher values for any δ_i^k beyond the maximum value will be ineffective. $\sum_{\forall s_i^k \in s_i} (\delta_i^k + |\chi_i^k|) \leq \sum_{\forall \tau_j \in \gamma_i} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right)$. Consequently, (35) will be:

$$\frac{s_{max}}{r_{max}} \leq \frac{\left(\sum_{\forall \tau_j \in \gamma_i} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \beta_{i,j} \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil}{\left(\sum_{\forall \tau_j \in \gamma_i} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil} \quad (30)$$

Since we are seeking the lower bound on $\frac{s_{max}}{r_{max}}$, let $\beta_{i,j}$ assume its minimum value. Thus, $\beta_{i,j} = 1$. (36) holds if $\frac{s_{max}}{r_{max}} \leq 1$.

Let $\delta_i^k(T_i) \rightarrow 0$ in (35). This means that transactions approximately execute according to their arrival order. Let $\beta_{i,j} \rightarrow \infty, \left\lceil \frac{T_i}{T_j} \right\rceil \rightarrow \infty$ in (35). This implies high contention. Consequently, $\frac{s_{max}}{r_{max}} \rightarrow \infty$. Hence, if transactions execute in FIFO order and contention is high, s_{max} can be much larger than r_{max} . Claim follows. \square

:20

CLAIM 13.

Under G-EDF and G-RMA, schedulability of FBLT is equal or better than lock-free's if $s_{max} \leq r_{max}$. If transactions execute in FIFO order (i. e., $\delta_i^k = 0, \forall s_i^k$) and contention is high, s_{max} can be much larger than r_{max} .

PROOF.

Lock-free synchronization [Devi et al. 2006; El-Shambakey and Ravindran 2012b] accesses only one object. Thus, the number of accessed objects per transaction in FBLT is limited to one. This allows us to compare the schedulability of FBLT with the lock-free algorithm.

By substituting $RC_A(T_i)$ and $RC_B(T_i)$ in (8) with (5) and (6.17) in [El-Shambakey 2012] respectively.

$$\begin{aligned} & \sum_{\forall \tau_i} \frac{\sum_{s_i^k \in s_i} (\delta_i^k \text{len}(s_i^k) + \sum_{s_{iz}^k \in \chi_i^k} \text{len}(s_{iz}^k)) + RC_{re}(T_i)}{T_i} \\ & \leq \sum_{\forall \tau_i} \frac{\left(\sum_{\forall \tau_j \in \gamma_i} \left(\left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \beta_{i,j} r_{max} \right) \right) + RC_{re}(T_i)}{T_i} \end{aligned} \quad (31)$$

where $\beta_{i,j}$ is the number of retry loops of τ_j that access the same objects as accessed by any retry loop of τ_i [Devi et al. 2006]. r_{max} is the maximum execution cost of a single iteration of any retry loop of any task [Devi et al. 2006]. For G-EDF(G-RMA), any job τ_i^x under FBLT has the same pattern of interference from higher priority jobs as ECM(RCM) respectively. $RC_{re}(T_i)$ for ECM, RCM and lock-free are given by Claims 25, 26 and 27 in [El-Shambakey 2012] respectively. $RC_{re}(T_i) = \left\lceil \frac{T_i}{T_j} \right\rceil s_{imax}, \forall \tau_j \neq \tau_i$ covers $RC_{re}(T_i)$ for G-EDF/FBLT and G-RMA/FBLT. $RC_{re}(T_i) = \left\lceil \frac{T_i}{T_j} \right\rceil r_{imax}, \forall \tau_j \neq \tau_i$ covers retry cost for G-EDF/lock-free and G-RMA/lock-free. (31) becomes

$$\begin{aligned} & \sum_{\forall \tau_i} \frac{\sum_{s_i^k \in s_i} (\delta_i^k \text{len}(s_i^k) + \sum_{s_{iz}^k \in \chi_i^k} \text{len}(s_{iz}^k)) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil s_{imax}}{T_i} \\ & \leq \sum_{\forall \tau_i} \frac{\left(\sum_{\forall \tau_j \in \gamma_i} \left(\left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \beta_{i,j} r_{max} \right) \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil r_{imax}}{T_i} \end{aligned} \quad (32)$$

Since $s_{max} \geq s_{imax}$, $\text{len}(s_i^k), \text{len}(s_{iz}^k), \forall i, z, k$ and $r_{max} \geq r_{imax}$ (32) holds if

$$\begin{aligned} & \sum_{\forall \tau_i} \frac{\left(\left(\sum_{s_i^k \in s_i} (\delta_i^k + |\chi_i^k|) \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil \right) s_{max}}{T_i} \\ & \leq \sum_{\forall \tau_i} \frac{\left(\left(\sum_{\forall \tau_j \in \gamma_i} \left(\left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \beta_{i,j} \right) \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil \right) r_{max}}{T_i} \end{aligned} \quad (33)$$

(33) holds if for each τ_i

$$\begin{aligned} & \left(\left(\sum_{s_i^k \in s_i} (\delta_i^k + |\chi_i^k|) \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil \right) s_{max} \\ & \leq \left(\left(\sum_{\forall \tau_j \in \gamma_i} \left(\left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \beta_{i,j} \right) \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil \right) r_{max} \end{aligned} \quad (34)$$

$$\frac{s_{max}}{r_{max}} \leq \frac{\left(\sum_{\forall \tau_j \in \gamma_i} \left(\left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \beta_{i,j} \right) \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil}{\left(\sum_{s_i^k \in s_i} (\delta_i^k + |\chi_i^k|) \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil} \quad (35)$$

FBLT

:21

It appears from (35) that as δ_i^k , as well as $|\chi_i^k|$, increases, then s_{max}/r_{max} decreases. So, to get the lower bound on s_{max}/r_{max} , let $\sum_{\forall s_i^k \in s_i} (\delta_i^k + |\chi_i^k|)$ reaches its maximum value. This maximum value is the total number of interfering transactions belonging to any job τ_j^l , $j \neq i$. Priority of τ_j^l can be higher or lower than current instance of τ_i . Beyond this maximum value, there will be no more transactions to conflict with s_i^k . So, higher values for any δ_i^k beyond maximum value will be ineffective. $\sum_{\forall s_i^k \in s_i} (\delta_i^k + |\chi_i^k|) \leq \sum_{\forall \tau_j \in \gamma_i} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right)$. Consequently, (35) will be

$$\frac{s_{max}}{r_{max}} \leq \frac{\left(\sum_{\forall \tau_j \in \gamma_i} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \beta_{i,j} \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil}{\left(\sum_{\forall \tau_j \in \gamma_i} \left(\left\lceil \frac{T_i}{T_j} \right\rceil + 1 \right) \right) + \sum_{\forall \tau_j \neq \tau_i} \left\lceil \frac{T_i}{T_j} \right\rceil} \quad (36)$$

As we look for the lower bound on $\frac{s_{max}}{r_{max}}$, let $\beta_{i,j}$ assumes its minimum value. So, $\beta_{i,j} = 1$. (36) holds if $\frac{s_{max}}{r_{max}} \leq 1$.

Let $\delta_i^k(T_i) \rightarrow 0$ in (35). This means transactions approximately execute in their arrival order. Let $\beta_{i,j} \rightarrow \infty$, $\left\lceil \frac{T_i}{T_j} \right\rceil \rightarrow \infty$ in (35). This means contention is high. Consequently, $\frac{s_{max}}{r_{max}} \rightarrow \infty$. So, if transactions execute in FIFO order and contention is high, s_{max} can be much larger than r_{max} . Claim follows.

□

8. EXPERIMENTAL EVALUATION

We now would like to understand how FBLT's retry cost compares with competitors in practice (i.e., on average). Since this can only be understood experimentally, we implement FBLT and the competitors and conduct experiments.

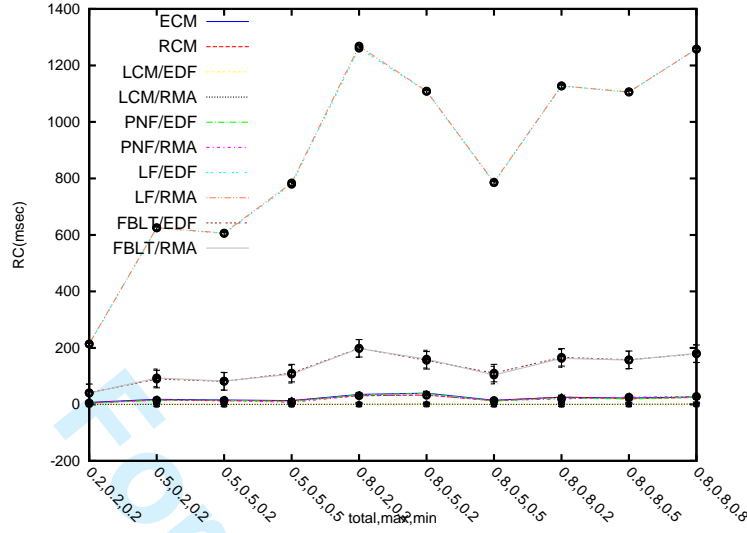
We used the ChronOS real-time Linux kernel [Dellinger et al. 2011] and the RSTM library [Marathe et al.] in our implementation. We implemented G-EDF and G-RMA schedulers in ChronOS, and modified RSTM to include implementations of ECM, RCM, LCM, PNF and FBLT. For the retry-loop lock-free synchronization, we used a loop that reads an object and attempts to write to it using a CAS instruction. The task retries until the CAS succeeds. We used an 8 core, 2GHz AMD Opteron platform. The average time taken for one write operation by RSTM on any core is $0.0129653375 \mu s$, and the average time taken by one CAS-loop operation on any core is $0.0292546250 \mu s$.

We used four task sets consisting of 4, 5, 8, and 20 periodic tasks. Each task runs in its own thread and has a set of atomic sections. Atomic section properties are probabilistically controlled using three parameters: the maximum and minimum lengths of any atomic section within a task, and the total length of atomic sections within any task. Since lock-free synchronization cannot handle more than one object per atomic section, we first compare FBLT's retry cost with that of lock-free (and other CMs) for one object per transaction. We then compare FBLT's retry cost with that of other CMs for multiple objects per transaction.

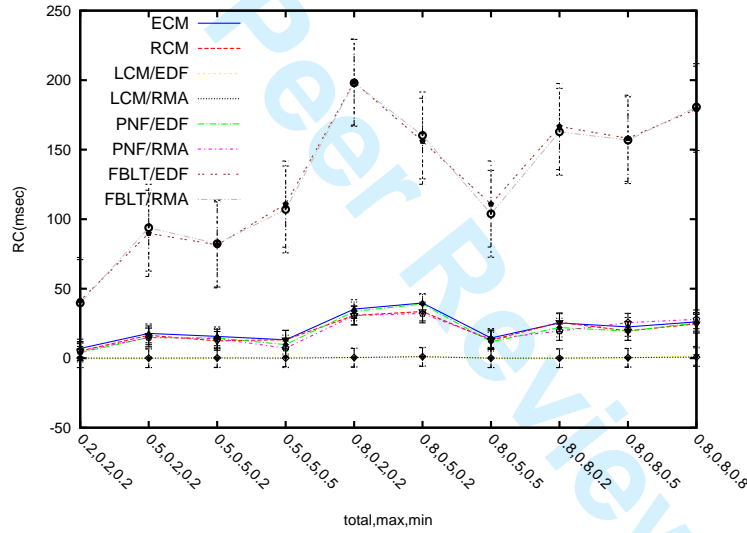
Figures 1 show the average retry cost for the 5 task case sharing 1 object. On the x-axis of the figures, we record 3 parameters x , y , and z . x is the ratio of the total length of all atomic sections of a task to the task WCET. y is the ratio of the maximum length of any atomic section of a task to the task WCET. z is the ratio of the minimum length of any atomic section of a task to the task WCET. Confidence level of all data points is 0.95.

While Figure 1(a) includes all methods, Figure 1(b) excludes lock-free. From these figures, we observe that lock-free has the largest retry cost, as it provides no conflict resolution. FBLT has the longest retry cost among CMs because transactions share

:22



(a) ECM, RCM, LCM, PNF, Lock-Free



(b) ECM, RCM, LCM, PNF

Fig. 1. Avg. retry cost (one object/transaction).

only 1 object. For multi-objects per transaction, FBLT provides equal or shorter retry cost than LCM as shown in Figures 2 and 3. PNF has an advantage over FBLT. PNF knows a priori all objects accessed by each transaction, but FBLT does not. Consequently, retry cost under PNF is shorter than FBLT. Experiments show that by proper adjustment of maximum abort number of each transaction, retry cost can be shorter than ECM, RCM and LCM. Retry cost can be comparable to PNF.

Similar trends were observed for the other task sets those are omitted here for brevity and due to space limitations. Other results are shown in [El-Shambakey and Ravindran 2011].

FBLT

:23

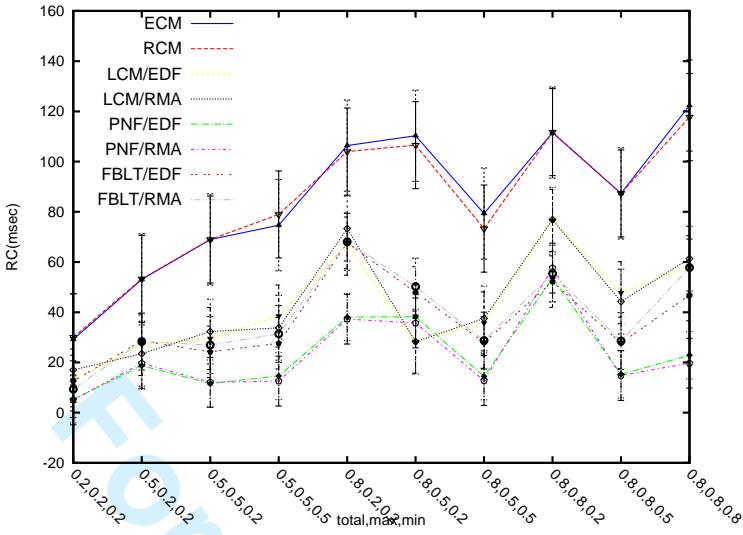


Fig. 2. Avg. retry cost (20 shared objects, 4 tasks).

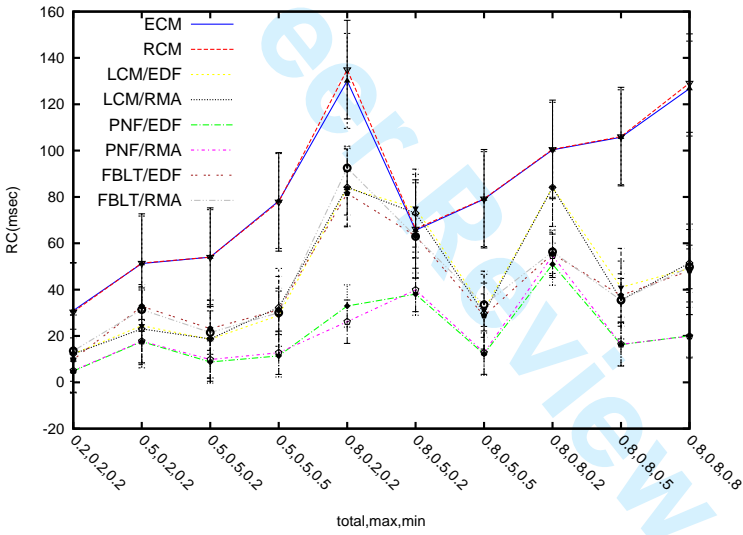


Fig. 3. Avg. retry cost (40 shared objects, 4 tasks).

9. CONCLUSIONS

Transitive retry increases transactional retry cost under ECM, RCM, and LCM. PNF avoids transitive retry by avoiding transactional preemptions. PNF avoids retry cost by concurrently executing non-conflicting. Non-conflicting objects are non-preemptive. PNF requires priori knowledge about objects accessed by each transaction. This is not suitable with dynamic STM. So, we introduce the First Bounded, Last Timestamp (FBLT) CM. Each transaction is allowed to abort for a specified number. Afterwards, transaction becomes non-preemptive. Non-preemptive transactions have higher prior-

:24

ities than other preemptive transactions and real-time jobs. Non-preemptive transactions resolve their conflicts using FIFO order.

By proper adjustment of maximum abort number of each transaction, FBLT schedulability is equal or better than other synchronization techniques. For FBLT's schedulability to be equal or better than lock-free, the upper bound on s_{max}/r_{max} is 1. Upper bound on s_{max}/r_{max} can be higher than 1 if transactions execute in their arrival order and contention is high.

Experimental results show that FBLT has equal or shorter retry cost than ECM, RCM and LCM. PNF knows a priori all objects accessed by each transaction. This is an advantage for PNF over FBLT. Consequently, retry cost under PNF is shorter than retry cost under FBLT.

Future work includes choosing another criterion to resolve conflicts of non-preemptive transactions. Also, using feedback from the system to adjust maximum abort number of each transaction. Consequently, retry cost can be reduced over time.

REFERENCES

- ANDERSON, J., RAMAMURTHY, S., AND JEFFAY, K. 1995. Real-time computing with lock-free shared objects. In *RTSS*. 28–37.
- BARROS, A. AND PINHO, L. 2011. Managing contention of software transactional memory in real-time systems. In *IEEE RTSS, Work-In-Progress*.
- BRANDENBURG, B., CALANDRINO, J., BLOCK, A., LEONTYEV, H., AND ANDERSON, J. 2008. Real-time synchronization on multiprocessors: To block or not to block, to suspend or spin? In *Real-Time and Embedded Technology and Applications Symposium, 2008. RTAS '08. IEEE*. 342–353.
- DAVIS, R. I. AND BURNS, A. 2011. A survey of hard real-time scheduling for multiprocessor systems. *ACM Comput. Surv.* 43, 4, 35:1–35:44.
- DELLINGER, M., GARYALI, P., AND RAVINDRAN, B. 2011. Chronos linux: a best-effort real-time multiprocessor linux kernel. In *Proceedings of the 48th Design Automation Conference*. DAC '11. ACM, New York, NY, USA, 474–479.
- DEVI, U., LEONTYEV, H., AND ANDERSON, J. 2006. Efficient synchronization under global edf scheduling on multiprocessors. In *18th Euromicro Conference on Real-Time Systems*. 10 pp. –84.
- EL-SHAMBAKEY, M. 2012. Phd proposal. Ph.D. thesis, Virginia Tech.
- EL-SHAMBAKEY, M. AND RAVINDRAN, B. 2011. On the design of real-time stm contention managers. Tech. rep., ECE Department, Virginia Tech. Available as <http://www.real-time.ece.vt.edu/tech-report-rt-stm-cm11.pdf>.
- EL-SHAMBAKEY, M. AND RAVINDRAN, B. 2012a. Stm concurrency control for embedded real-time software with tighter time bounds. In *Proceedings of the 49th Annual Design Automation Conference*. DAC '12. ACM, New York, NY, USA, 437–446.
- EL-SHAMBAKEY, M. AND RAVINDRAN, B. 2012b. Stm concurrency control for multicore embedded real-time software: time bounds and tradeoffs. In *Proceedings of the 27th Annual ACM Symposium on Applied Computing*. SAC '12. ACM, New York, NY, USA, 1602–1609.
- FAHMY, S. AND RAVINDRAN, B. 2011. On stm concurrency control for multicore embedded real-time software. In *International Conference on Embedded Computer Systems*. SAMOS. 1–8.
- FAHMY, S., RAVINDRAN, B., AND JENSEN, E. D. 2009. On bounding response times under software transactional memory in distributed multiprocessor real-time systems. In *DATE*. 688–693.
- FAHMY, S. F. 2010. Collaborative scheduling and synchronization of distributable real-time threads. Ph.D. thesis, Virginia Tech.
- GUERRAoui, R., HERLIHY, M., AND POCHON, B. 2005. Toward a theory of transactional contention managers. In *PODC*. 258–264.
- HARRIS, T., MARLOW, S., JONES, S. P., AND HERLIHY, M. 2008. Composable memory transactions. *Commun. ACM* 51, 91–100.
- HERLIHY, M. 2006. The art of multiprocessor programming. In *PODC*. 1–2.
- HERLIHY, M., LUCHANGCO, V., MOIR, M., AND SCHERER, III, W. N. 2003. Software transactional memory for dynamic-sized data structures. In *Proceedings of the twenty-second annual symposium on Principles of distributed computing*. PODC '03. ACM, New York, NY, USA, 92–101.
- MANSON, J., BAKER, J., ET AL. 2006. Preemptible atomic regions for real-time Java. In *RTSS*. 10–71.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

FBLT :25

MARATHE, V., SPEAR, M., HERIOT, C., ACHARYA, A., EISENSTAT, D., SCHERER III, W., AND SCOTT, M. Lowering the overhead of nonblocking software transactional memory. In *Workshop on Languages, Compilers, and Hardware Support for Transactional Computing*. TRANSACT.

SAHA, B., ADL-TABATABAI, A.-R., ET AL. 2006. McRT-STM: a high performance software transactional memory system for a multi-core runtime. In *PPoPP*. 187–197.

SARNI, T., QUEUDET, A., AND VALDURIEZ, P. 2009. Real-time support for software transactional memory. In *RTCSA*. 477–485.

SCHOEBERL, M., BRANDNER, F., AND VITEK, J. 2010. RTTM: Real-time transactional memory. In *ACM SAC*. 326–333.

Received February 2007; revised March 2009; accepted June 2009