

# Using Checkpointing to Improve Time Bounds for STM Concurrency Control in Embedded Real-Time Software

**Abstract**—We consider checkpointing with software transactional memory (STM) concurrency control for embedded multi-core real-time software, and present a modified version of FBLT contention manager called *Checkpointing FBLT* (CPFBLT). We upper bound transactional retries and task response times under CPFBLT, and identify when CPFBLT is a more appropriate alternative to FBLT without checkpointing.

## I. 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 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 [1]. 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 at least one access is a write. When that happens, a contention manager (CM) [2] 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 lock-free approach, especially for high contention and read-dominated workloads (see an example TM system's performance in [3]), and is composable [4]. TM has been proposed in hardware, 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.

Past real-time CM research proposed resolving transactional contention using dynamic and fixed priorities of parent threads. [5], [6], [7] present Earliest Deadline CM (ECM) and Rate Monotonic CM (RCM), which are used with global EDF (G-EDF) and global RMS (G-RMS) multicore real-time schedulers [8]. In particular, [6] 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 Length-based CM (LCM) in [7], increasing the coverage of STM's timeliness superiority. ECM, RCM, and LCM suffer from transitive retry. Transitive retry means one transaction aborts and retries due to another transaction with no shared objects between both transactions. Transitive retry is introduced due to access of multiple objects per transaction. Thus, ECM, RCM and LCM cannot handle multiple objects per transaction efficiently. These limitations are overcome with the Priority with Negative value and First access CM (PNF) [9]. However, PNF requires prior knowledge of all objects accessed by each transaction. This significantly limits programmability, and is incompatible with dynamic STM implementations [10]. Additionally, PNF is a centralized CM, which increases overheads and retry costs, and has a complex implementation. First Bounded, Last Timestamp CM (or FBLT) [11], in contrast to PNF, does not require prior 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. Retry cost under FBLT is close or better than retry cost under PNF. This results from prior knowledge of accessed objects per transaction under PNF. Additionally, FBLT is a decentralized CM that has a simpler implementation than PNF.

Under previous CMs, if two transactions conflict on a specific object, the aborted transaction must restart from the beginning. Even if the contended object is not initially accessed at the beginning of the aborted transaction. Thus, the time between the beginning of the aborted transaction and first access to the contended object is wasted. During this wasted time, other transactions may also conflict with the aborted transaction. Thus, increasing chances to abort it again. Checkpointing [12] can be used to solve this problem. checkpointing can reduce response time of threads with conflicting transactions. Under checkpointing, a transaction retreats to a previous control flow location upon conflict. So, an aborted transaction does not have to retreat to its beginning.

As FBLT is shown to be better or equal to ECM, RCM, LCM and PNF [11], we investigate effect of checkpointing to FBLT. We present the motivation for introducing checkpointing into FBLT in Section IV. We introduce checkpointing FBLT (CPFBLT) that combines original FBLT with checkpointing (Section V). We establish CPFBLT's retry and response time upper bounds under G-EDF and G-RMA schedulers (Section VI). We also identify the conditions under which CPFBLT is a better alternative to non-checkpointing FBLT (Section VII). We implement FBLT and CPFBLT in the Rochester STM framework [13] and conduct experimental

studies (Section VIII). Our results reveal that CPFBLT has shorter response time than non-checkpointing FBLT.

Thus, the paper's contribution is the use of checkpointing as a complementary tool to FBLT to further enhance response time. CPFBLT, thus allows programmers to reap STM's significant programmability and composability benefits for multicore embedded real-time software.

## II. 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., [14]). Despite their 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 lock-based synchronization [15]. In contrast, STM is semantically simpler [1], and is often the only viable lock-free solution for complex data structures (e.g., red/black tree) [16] and nested critical sections [3]. STM concurrency control for real-time systems has been previously studied in [17], [18], [19], [20], [16], [21], [6], [7], [9], [11].

[17] proposes a restricted version of STM for uniprocessors. [18] bounds response times in distributed systems with STM synchronization. [18] considers Pfair scheduling, limit to small atomic regions with fixed size, and limit transaction execution to span at most two quanta. [19] presents real-time scheduling of transactions and serializes transactions based on deadlines. However, the work does not bound retries and response times. [20] proposes real-time HTM. [20] assumes that the worst case conflict between atomic sections of different tasks occurs when the sections are released at the same time.

[16] upper bounds retries and response times for ECM with G-EDF, and identify the tradeoffs with locking and lock-free protocols. Similar to [20], [16] also assumes that the worst case conflict between atomic sections occurs when the sections are released simultaneously. The ideas in [16] are extended in [21], which presents three real-time CM designs.

[6] presents the ECM and RCM contention managers, and upper bounds transactional retries and task response times under them. [6] also identifies the conditions under which ECM and RCM are superior to lock-free techniques. In particular, [6] shows that, STM's superiority holds only under some ranges for the maximum atomic section length. Moreover, [6] restricts transactions to access only one object.

[7] presents the LCM contention manager, and upper bounds transactional retry cost and task response times for G-EDF and G-RMA schedulers. This work also compares (analytically and experimentally) LCM with ECM, RCM, and lock-free synchronization. However, similar to [7], [6] restricts transactions to access only one object.

[9] presents the PNF contention manager, which allows transactions to access multiple objects and avoids the consequent transitive retry effect. The work also upper bounds transactional retries and task response times under G-EDF and G-RMA. However, PNF requires a-priori knowledge of the objects accessed by each transaction, which is not always possible, limits programmability, and is incompatible with

dynamic STM implementations [10]. Additionally, PNF is a centralized CM with complex implementation.

[11] presents the FBLT contention manager. In contrast to PNF, FBLT does not require prior knowledge of required objects by each transaction. FBLT permits multiple objects per transaction. Under FBLT, each transaction can be aborted for a specific number of times. Afterwards, the transaction becomes non-preemptive. Non-preemptive transaction cannot be aborted except by another non-preemptive transaction. Non-preemptive transactions resolve conflicts based on the time they became non-preemptive.

Previous CMs try to enhance response time of real-time tasks using different policies for conflict resolution. Checkpointing does not require aborted transaction to restart from beginning. Thus, Checkpointing can be plugged into different CMs to further improve response time. [12] introduces checkpointing as an alternative to closed nesting transactions[22]. [12] uses boosted transactions [23] instead of closed nesting [24], [22], [25] to implement checkpointing. Boosted transactions are based on linearizable objects with abstract states and concrete implementation. Methods under boosted transaction have well defined semantics to transit objects from one state to another. Inverse methods are used to restore objects to previous states. Upon a conflict, a transaction does not need to revert to its beginning, but rather to a point where the conflict can be avoided. Thus, checkpointing enables partial abort. [26] applies checkpointing in distributed transactional memory using Hyflow [27].

## III. PRELIMINARY

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  $\tau_i$  is denoted  $\tau_i^k$ . Each task  $\tau_i$  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  $\tau_i^j$  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  $\tau_i$ 's (fixed) priority and it is constant for all instances of  $\tau_i$ . Under a dynamic priority scheduler such as G-EDF, a job  $\tau_i^j$ 's priority,  $p_i^j$ , differs from one instance to another.

*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.  $s_i^k$  is the  $k^{th}$  atomic section of  $\tau_i$ . Each object,  $\theta$ , can be accessed by multiple tasks. The set of distinct objects accessed by  $s_i^k$  is  $\Theta_i^k$ .  $s_i^k$  executes for a duration  $len(s_i^k)$ . Let  $s_i^k$  accesses objects  $\theta_1, \theta_2, \dots, \theta_g, \dots, \theta_z$  in that order. If all objects before  $\theta_g$  are not shared between  $s_i^k$  and any other transaction, then  $\nabla_{i*}^k$  is the time interval between start of  $s_i^k$  and the first access to  $\theta_g$  by  $s_i^k$ . So,  $\nabla_{i*}^k$  is the time interval between start of  $s_i^k$  and the first access to any shared object between  $s_i^k$  and any other transaction.  $rt(\nabla_{i*}^k) = \frac{\nabla_{i*}^k}{len(s_i^k)}$  is the maximum ratio of length of  $s_i^k$  where objects can be shared between  $s_i^k$  and other transactions.  $rt(\nabla_{i*}^k) = 1 - rt(\nabla_{i*}^k)$  is the minimum ratio of length of  $s_i^k$  where no objects can be shared between

$s_i^k$  and other transactions.  $rt(\nabla) = \max\left(rt(\nabla_{i*}^k)\right), \forall s_i^k$  is the maximum  $rt(\nabla_{i*}^k)$  between all transactions.

**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 transactions. 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”. The total time that a task  $\tau_i$ ’s atomic sections have to retry over  $T_i$  is denoted  $RC(T_i)$ .

Due to *transitive retry* [9], [11], an atomic section  $s_i^k$  may retry due to another atomic section  $s_j^l$ , where  $s_i^k$  and  $s_j^l$  have no shared objects. Transitive retry happens when  $s_j^l$  conflicts with another transaction  $s_h^u$ , and  $s_h^u$  conflicts with  $s_i^k$ .  $s_i^k$  is said to transitively, or indirectly, retry due to  $s_j^l$ .

#### IV. MOTIVATION

Under checkpointing, a transaction  $s_i^k \in \tau_i$  does not need to restart from the beginning upon a conflict on object  $\theta$ .  $s_i^k$  just needs to return to the first point it accessed  $\theta$ . If  $s_i^k$  needs to restart from its beginning, then the time between the beginning of  $s_i^k$  and the first access to  $\theta$  is wasted. Besides, restarting  $s_i^k$  from its beginning increases the chances of aborting  $s_i^k$  by other transactions. Thus, response time of  $\tau_i$  can be improved by checkpointing unless  $s_i^k$  acquires all its objects at its beginning. While ECM, RCM, LCM, PNF and FBLT without checkpointing try to resolve conflicts using proper strategies, checkpointing enhances performance by reducing aborted part of each transaction. Thus, checkpointing acts as a complementary component to different CMs to further enhance response time.

Behaviour of some CMs, like PNF [9], can make checkpointing useless. PNF requires a priori knowledge of accessed objects within transactions. Only the first  $m$  non-conflicting transactions are allowed to execute concurrently and non-preemptively. Thus, PNF makes no use of checkpointing because there is no conflict between non-preemptive transactions.

Other CMs (e. g., FBLT[11]) allow conflicting transaction to run concurrently. So, FBLT can benefit from checkpointing. FBLT, by definition, depends on LCM. LCM, in turn, depends on ECM for G-EDF and RCM for G-RMA. Experimental results show superiority of FBLT over LCM, ECM and RCM[11]. Due to prior knowledge of accessed objects per transaction in PNF, retry cost under FBLT is close or better than retry cost under PNF [11]. So, we extend FBLT to checkpointing FBLT (CPFBLT) to improve response time than the non-checkpointing FBLT (NCPFBLT).

#### V. CHECKPOINTING FBLT (CPFBLT)

CPFBLT depends on LCM [7] with checkpointing. So, we initially illustrate LCM integrated with checkpointing (Section V-A). Afterwards, we illustrate FBLT with checkpointing in (Section V-B).

#### A. Checkpointing LCM (CPLCM)

Algorithm 1 presents LCM [7] integrated with checkpointing to give CPLCM. A new checkpoint is recorded for

---

#### Algorithm 1: CPLCM

---

**Data:**

$s_i^k \rightarrow$  interfered transaction.

$s_j^l \rightarrow$  interfering transaction with  $s_i^k$  on object  $\theta_{ij}^{kl}$ .

$\Psi \rightarrow$  user predefined threshold  $\in [0, 1]$ .

$\epsilon_i^k \rightarrow$  remaining execution length of  $s_i^k$  when interfered by  $s_j^l$ .

$cp_h^u(\theta) \rightarrow$  recorded checkpoint in transaction  $s_h^u$  for newly accessed object  $\theta$

**Result:** which transaction of  $s_i^k$  or  $s_j^l$  aborts

1 **foreach** newly accessed  $\theta$  requested by any transaction  $s_h^u$  **do**

2 | Add a checkpoint  $cp_h^u(\theta)$

3 **end**

4 **if**  $p_i^k > p_j^l$  **then**

5 |  $s_j^l$  aborts and retreats to  $cp_j^l(\theta_{ij}^{kl})$ ;

6 | Remove all checkpoints in  $s_j^l$  recorded after  $cp_j^l(\theta_{ij}^{kl})$

7 **else**

8 |  $c_{ij}^{kl} = \text{len}(s_j^l) / \text{len}(s_i^k)$ ;

9 |  $\alpha_{ij}^{kl} = \ln(\Psi) / (\ln(\Psi) - c_{ij}^{kl})$ ;

10 |  $\alpha_i^k = (\text{len}(s_i^k) - \epsilon_i^k) / \text{len}(s_i^k)$ ;

11 **if**  $\alpha_i^k \leq \alpha_{ij}^{kl}$  **then**

12 |  $s_i^k$  aborts and retreats to  $cp_i^k(\theta_{ij}^{kl})$ ;

13 | Remove all checkpoints in  $s_i^k$  recorded after  $cp_i^k(\theta_{ij}^{kl})$

14 **else**

15 |  $s_j^l$  aborts and retreats to  $cp_j^l(\theta_{ij}^{kl})$ ;

16 | Remove all checkpoints in  $s_j^l$  recorded after  $cp_j^l(\theta_{ij}^{kl})$

17 **end**

18 **end**

---

each newly accessed object  $\theta$  by any transaction  $s_h^u$  (step 2). Checkpoint is recorded when  $\theta$  is accessed for the first time because any further changes to  $\theta$  will be discarded upon conflict. CPLCM uses priorities of  $s_i^k$  and  $s_j^l$ , the remaining length of  $s_i^k$  when it is interfered, as well as  $\text{len}(s_j^l)$ , to decide which transaction must be aborted. If  $p_i^k > p_j^l$ , then  $s_j^l$  would be the transaction to abort because of lower priority of  $s_j^l$ , and the start of  $s_i^k$  before  $s_j^l$  (step 5). Otherwise,  $c_{ij}^{kl}$ ,  $\alpha_{ij}^{kl}$  and  $\alpha_i^k$  are calculated (steps 8, 9 and 10) to determine whether it is worth aborting  $s_i^k$  in favour of  $s_j^l$ . If  $\text{len}(s_j^l)$  is relatively small compared to  $\text{len}(s_i^k)$ , then  $c_{ij}^{kl}$  approaches its minimum value (i.e., 0), and  $\alpha_{ij}^{kl}$  approaches its maximum value (i.e., 1) (steps 8, 9). Otherwise,  $c_{ij}^{kl}$  approaches its maximum value (i.e.,  $\infty$ ), and  $\alpha_{ij}^{kl}$  approaches its minimum value (i.e., 0).  $\Psi$  is a predefined threshold that lies in  $[0, 1]$  as defined in [7]. The remaining execution length of  $s_i^k$  (i.e.,  $\epsilon_i^k$ ) is used to calculate  $\alpha_i^k$  (step 10). If  $s_i^k$  still has much work to do, then  $\epsilon_i^k$  approaches  $\text{len}(s_i^k)$  and  $\alpha_i^k$  approaches 0. Otherwise,  $\alpha_i^k$  approaches 1. If  $\text{len}(s_i^k)$  is much longer than  $\text{len}(s_j^l)$ , or  $s_i^k$  still has much work to do when interfered by  $s_j^l$ , then  $\alpha_i^k$  tends to be smaller than  $\alpha_{ij}^{kl}$ . Consequently,  $s_i^k$  aborts in favour of  $s_j^l$ . When  $s_i^k$  aborts upon a conflict with  $s_j^l$  on object  $\theta_{ij}^{kl}$ , then checkpoints in  $s_i^k$  recorded after  $cp_i^k(\theta_{ij}^{kl})$  are removed (step 13). Prior checkpoints to  $cp_i^k(\theta_{ij}^{kl})$  remain the same. Also, if  $s_j^l$  aborts in favour of  $s_i^k$ , then all checkpoints in  $s_j^l$  recorded after  $cp_j^l(\theta_{ij}^{kl})$  are removed (steps 6, 16).

## B. Design of CPFBLT

Algorithm 2 illustrates FBLT [11] integrated with checkpointing to give CPFBLT. A new checkpoint is recorded for

---

### Algorithm 2: The CPFBLT Algorithm

---

**Data:**  
 $s_i^k$ : interfered transaction.  
 $s_j^l$ : interfering transaction.  
 $\delta_i^k$ : maximum number of times  $s_i^k$  can be aborted during  $T_i$ .  
 $\eta_i^k$ : 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  $m\_set$ .  
 $cp_h^u(\theta) \rightarrow$  recorded checkpoint in transaction  $s_h^u$  for newly accessed object  $\theta$   
**Result:** which transaction,  $s_i^k$  or  $s_j^l$ , aborts

```

1 foreach newly accessed  $\theta$  requested by any transaction  $s_h^u$  do
2   | Add a checkpoint  $cp_h^u(\theta)$ 
3 end
4 if  $s_i^k, s_j^l \notin m\_set$  then
5   | Apply CPLCM (Algorithm 1);
6   | if  $s_i^k$  is aborted then
7     | if  $\eta_i^k < \delta_i^k$  then
8       | Increment  $\eta_i^k$  by 1;
9     | else
10      | Add  $s_i^k$  to  $m\_set$ ;
11      | Record  $r(s_i^k)$ ;
12      | Increase priority of  $s_i^k$  to  $m\_prio$ ;
13    | end
14   | else
15     | Swap  $s_i^k$  and  $s_j^l$ ;
16     | Go to Step 6;
17   | end
18 else if  $s_j^l \in m\_set, s_i^k \notin m\_set$  then
19   |  $s_i^k$  aborts and retreats to  $cp_i^k(\theta_{ij}^{kl})$ ;
20   | Remove all checkpoints in  $s_i^k$  recorded after  $cp_i^k(\theta_{ij}^{kl})$ ;
21   | if  $\eta_i^k < \delta_i^k$  then
22     | Increment  $\eta_i^k$  by 1;
23   | else
24     | Add  $s_i^k$  to  $m\_set$ ;
25     | Record  $r(s_i^k)$ ;
26     | Increase priority of  $s_i^k$  to  $m\_prio$ ;
27   | end
28 else if  $s_i^k \in m\_set, s_j^l \notin m\_set$  then
29   | Swap  $s_i^k$  and  $s_j^l$ ;
30   | Go to Step 18;
31 else
32   | if  $r(s_i^k) < r(s_j^l)$  then
33     |  $s_j^l$  aborts and retreats to  $cp_j^l(\theta_{ij}^{kl})$ ;
34     | Remove all checkpoints in  $s_j^l$  recorded after  $cp_j^l(\theta_{ij}^{kl})$ ;
35   | else
36     |  $s_i^k$  aborts and retreats to  $cp_i^k(\theta_{ij}^{kl})$ ;
37     | Remove all checkpoints in  $s_i^k$  recorded after  $cp_i^k(\theta_{ij}^{kl})$ ;
38   | end
39 end

```

---

each newly accessed object  $\theta$  by any transaction  $s_h^u$  (step 2). Checkpoint is recorded when  $\theta$  is accessed for the first time because any further changes to  $\theta$  will be discarded upon conflict. Each transaction  $s_i^k$  can be aborted during  $T_i$  for at most  $\delta_i^k$  times.  $\eta_i^k$  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 CPLCM (step 5). Thus, CPFBLT defaults to CPLCM when the conflicting transactions

( $s_i^k$  and  $s_j^l$ ) have not reached their  $\delta$ s ( $\delta_i^k$  and  $\delta_j^l$ ).  $\eta_i^k$  is incremented each time  $s_i^k$  is aborted as long as  $\eta_i^k < \delta_i^k$  (steps 8 and 22). Otherwise,  $s_i^k$  is added to the  $m\_set$  and priority of  $s_i^k$  is increased to  $m\_prio$  (steps 10 to 12 and 24 to 26). 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 (steps 18 to 30). On the other hand, non-preemptive transactions can abort each other. 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 11 and 25). When non-preemptive transactions conflict together (step 31), the transaction that joined  $m\_set$  first becomes the transaction that commits first (steps 33 and 36). When  $s_i^k$  aborts due to a conflict on  $\theta_{ij}^{kl}$  with  $s_j^l$ , then  $s_i^k$  retreats to  $cp_i^k(\theta_{ij}^{kl})$ . All checkpoints recorded after  $cp_i^k(\theta_{ij}^{kl})$  are removed (steps 20, and 37). Similarly,  $s_j^l$  removes all checkpoints recorded after  $cp_j^l(\theta_{ij}^{kl})$  if aborted by  $s_i^k$  (step 34).

## VI. CPFBLT RETRY COST

**Claim 1.** Assume only two transaction  $s_i^k$  and  $s_j^l$  conflicting together. Let  $s_i^k$  begins at time  $S(s_i^k)$  and  $s_j^l$  begins at time  $S(s_j^l)$ . Let  $\Delta = S(s_j^l) - S(s_i^k)$ . In the absence of checkpointing, retry cost of  $s_i^k$  due to  $s_j^l$  is given by

$$BASE\_RC_{ij}^{kl} \leq \begin{cases} len(s_j^l) + \Delta & , -len(s_j^l) \leq \Delta \leq len(s_i^k) \\ 0 & , \text{Otherwise} \end{cases} \quad (1)$$

$BASE\_RC_{ij}^{kl}$  is upper bounded by

$$len(s_j^l + s_i^k) \quad (2)$$

which is the same upper bound given by Proofs of Claims 1 and 3 in [6]

*Proof:* Due to absence of checkpointing,  $s_i^k$  aborts and retries from its beginning due to  $s_j^l$ . So,  $s_i^k$  retries for the period starting from  $S(s_i^k)$  to the end of execution of  $s_j^l$ .  $s_j^l$  ends execution at  $S(s_j^l) + len(s_j^l)$ . If  $S(s_j^l) < S(s_i^k) - len(s_j^l)$ , then  $s_j^l$  finishes execution before start of  $s_i^k$  and there will be no conflict. Also, if  $S(s_j^l) > S(s_i^k) + len(s_i^k)$ , then  $s_j^l$  starts execution after  $s_i^k$  finishes execution and there will be no conflict. Thus, (1) follows. Equation (2) is derived by substitution of  $\Delta$  by its maximum value (i.e.,  $(s_i^k)$ ). Claim follows. ■

**Claim 2.** Assume only two transactions  $s_i^k$  and  $s_j^l$  conflicting on one object  $\theta$ . Let  $\nabla_j^l$  be the time interval between the start of  $s_j^l$  and the first access to  $\theta$ . Similarly, let  $\nabla_i^k$  be the time interval between the start of  $s_i^k$  and the first access to  $\theta$ . Let  $\Delta$  be the time difference between start of  $s_j^l$  to the start of  $s_i^k$ . So,  $\Delta < 0$  if  $s_j^l$  starts before  $s_i^k$ . Under checkpointing,  $s_i^k$  aborts

and retries due to  $s_j^l$  for

$$RC0_{ij}^{kl} \leq \begin{cases} \text{len}(s_j^l) - \nabla_i^k + \Delta, & \text{if } \Delta \geq \nabla_i^k - \text{len}(s_j^l) \\ 0, & \text{Otherwise} \end{cases} \quad (3)$$

$RC0_{ij}^{kl}$  is upper bounded by

$$\text{len}(s_j^l + s_i^k) - \nabla_j^l - \nabla_i^k \quad (4)$$

*Proof:* As  $s_i^k$  and  $s_j^l$  conflict only on one object  $\theta$ , there will be no conflict before both  $s_i^k$  and  $s_j^l$  access  $\theta$ . Retry cost of  $s_i^k$  due to  $s_j^l$  is derived by Claim 1 excluding parts of  $s_i^k$  and  $s_j^l$  before both transactions access  $\theta$ . Thus, excluding the parts of  $s_i^k$  and  $s_j^l$  that do not cause conflict. So,  $\text{len}(s_i^k)$  in Claim 1 is substituted by  $\text{len}(s_i^k) - \nabla_i^k$ .  $\text{len}(s_j^l)$  is substituted by  $\text{len}(s_j^l) - \nabla_j^l$ .  $\Delta$  in Claim 1 is substituted by  $\Delta + \nabla_j^l - \nabla_i^k$ . Claim follows. ■

**Claim 3.** Assume only two transactions  $s_i^k$  and  $s_j^l$  conflicting on a number of objects  $\theta_1, \theta_2 \dots \theta_z$ . Let  $\nabla_{i*}^k$  be the time interval between start of  $s_i^k$  and the first access to the first object accessed by  $s_i^k$  and shared with  $s_j^l$  (e.g.,  $\theta_i$ ). Let  $\nabla_{j*}^l$  be the time interval between start of  $s_j^l$  and the first access to the first object accessed by  $s_j^l$  and shared with  $s_i^k$  (e.g.,  $\theta_j$ ).  $\theta_i$  and  $\theta_j$  may not be the same. With checkpointing, retry cost of  $s_i^k$  due to  $s_j^l$  is upper bounded by

$$RC1_{ij}^{kl} \leq \text{len}(s_i^k + s_j^l) - \nabla_{i*}^k - \nabla_{j*}^l \quad (5)$$

*Proof:* Proof follows directly from Claim 2 by maximizing (4).  $\text{len}(s_i^k)$ , as well as,  $\text{len}(s_j^l)$  in (4) cannot be changed. Thus, by choosing minimum values of  $\nabla_i^k$  and  $\nabla_j^l$  that correspond to shared objects between  $s_i^k$  and  $s_j^l$ , (4) is maximized. Claim follows. ■

**Claim 4.** If  $s_j^l$  is conflicting indirectly (through transitive retry) with  $s_i^k$ , then it is safe to ignore  $\nabla_{i*}^k$  in calculating the upper bound of retry cost of  $s_i^k$  due to  $s_j^l$ .

*Proof:* If  $s_j^l$  is conflicting indirectly with  $s_i^k$ , then  $s_j^l$  is accessing an object  $\theta$  that does not belong to  $\Theta_i^k$ . In this case, to get an upper bound for the retry cost of  $s_i^k$  due to  $s_j^l$ ,  $\nabla_{i*}^k$  assumes its minimum value in (5). Thus,  $\nabla_{i*}^k = 0$ . Claim follows. ■

**Claim 5.** Assume only two non-preemptive transactions  $s_i^k$  and  $s_j^l$  under CPFBLT. With checkpointing, retry cost of  $s_i^k$  due to direct or indirect conflict with  $s_j^l$  is upper bounded by

$$RC2_{ij}^{kl} \leq \text{len}(s_j^l) - \nabla_{i*}^k \quad (6)$$

where  $\nabla_{i*}^k = 0$  in case of indirect conflict.

*Proof:* Proof follows directly from Claims 2, 3 and 4 except that  $s_j^l$  must have become non-preemptive before  $s_i^k$ .

So,  $s_j^l$  starts execution non-preemptively before  $s_i^k$ . Otherwise, by definition of CPFBLT,  $s_j^l$  will not be able to abort  $s_i^k$ . Thus,  $\Delta$  must not exceed 0. Claim follows. ■

**Claim 6.** Let  $s_i^k$  be a non-preemptive transaction under CPFBLT. Let  $\chi_i^k$  be the set of transactions conflicting (directly or indirectly) with  $s_i^k$ . Each transaction  $s_j^l \in \chi_i^k$  belongs to a distinct task. Transactions in  $\chi_i^k$  are organized in non-increasing order of  $RC2_{ij}^{kl}$  for each  $s_j^l \in \chi_i^k$ . Total retry cost of non-preemptive transaction  $s_i^k$  due to other non-preemptive transactions is upper bounded by

$$RC3_i^k \leq \sum_{a=1}^{a=\min(|\chi_i^k|, m-1)} RC2_i^k(\chi_i^k(a)) \quad (7)$$

where  $\chi_i^k(a)$  is the  $a^{\text{th}}$  transaction in  $\chi_i^k$ . If  $\chi_i^k(a) = s_j^l$ , then  $RC2_i^k(\chi_i^k(a)) = RC2_{ij}^{kl}$ .

*Proof:* By definition of CPFBLT, a transaction  $s_i^k$  can be preceded by at most  $m-1$  non-preemptive transactions. As non-preemptive transactions are organized in the order they become non-preemptive, no two non-preemptive transactions can belong to the same task. Maximum retry cost of non-preemptive  $s_i^k$  occurs when: 1)  $s_i^k$  is preceded by at most  $m-1$  transactions conflicting with  $s_i^k$ . 2) Each conflicting transaction  $s_j^l$  to  $s_i^k$  must have one of the highest  $m-1$  values for  $RC2_{ij}^{kl}$ . 3) Non-preemptive transactions preceding  $s_i^k$  are executing sequentially. Thus, retry cost of non-preemptive  $s_i^k$  can be upper bounded by Claim 5 for at most the first  $m-1$  transactions in  $\chi_i^k$ . If the third condition is not satisfied, then (7) still gives a correct, but not tight, upper bound. Claim follows. ■

**Claim 7.** Under CPFBLT, a preemptive transaction  $s_i^k$  aborts and retries for at most

$$RC4_i^k \leq \delta_i^k (\text{len}(s_i^k) - \min(\nabla_{i*}^k)) \quad (8)$$

where  $\min(\nabla_{i*}^k)$  is the minimum  $\nabla_{i*}^k$  for  $s_i^k$  and any other conflicting transaction  $s_j^l$ . If there are indirectly conflicting transactions with  $s_i^k$ , then  $\min(\nabla_{i*}^k) = 0$ .

*Proof:* No transaction will make preemptive  $s_i^k$  aborts and retries before  $\min(\nabla_{i*}^k)$ . By checkpointing,  $s_i^k$  will not retreat earlier than  $\min(\nabla_{i*}^k)$ . By definition of CPFBLT, preemptive  $s_i^k$  aborts for at most  $\delta_i^k$  times before it becomes non-preemptive. Claim follows. ■

**Claim 8.** The total retry cost of any job  $\tau_i^x$  under CPFBLT due to 1) conflicts with other transactions during an interval  $L$ . 2) release of higher priority jobs during execution of preemptive transactions is upper bounded by

$$RC(L)_{i0}^i = \sum_{s_i^k \in s_i} (RC4_i^k + RC3_i^k) + RC_{re}(L) \quad (9)$$

where  $RC_{re}(L)$  is the retry cost resulting from the release of higher priority jobs during execution of preemptive transactions.  $RC_{re}(L)$  is calculated by Claim 1 in [11].

*Proof:* Following Claims 4, 6, 7 and Claim 1 in [11], Claim follows. ■

Any newly released task  $\tau_i^x$  can be blocked by  $m$  lower priority non-preemptive transactions. Blocking time  $D_i$  of any job  $\tau_i^x$  is independent of checkpointing. Thus,  $D_i$  is calculated by Claim 3 in [11]. Claim 2 in [11] is used to calculate response time under CPFBLT where  $RC_{io}(T_i)$  is calculated by (9).

## VII. CPFBLT vs. NCPFBLT

**Claim 9.** *Schedulability of CPFBLT is better or equal to schedulability of NCPFBLT if shared objects within each transaction  $s_i^k$  are accessed close to the end of execution  $s_i^k$ .*

*Proof:* Let upper bound on retry cost of any task  $\tau_i^x$  during  $T_i$  under NCPFBLT be denoted as  $RC_i^{ncp}$ .  $RC_i^{ncp}$  is calculated by Claim 1 in [11]. Let upper bound on retry cost of any task  $\tau_i^x$  during  $T_i$  under CPFBLT be denoted as  $RC_i^{cp}$ .  $RC_i^{cp}$  is calculated by (9). Let  $D_i$  be the upper bound on blocking time of any newly released task  $\tau_i^x$  during  $T_i$  due to lower priority jobs.  $D_i$  is the same for both CPFBLT and NCPFBLT.  $D_i$  is calculated by Claim 2 in [11]. For CPFBLT schedulability to be better than schedulability of NCPFBLT:

$$\sum_{\forall \tau_i} \frac{c_i + RC_i^{cp} + D_i}{T_i} \leq \sum_{\forall \tau_i} \frac{c_i + RC_i^{ncp} + D_i}{T_i} \quad (10)$$

$\therefore D_i$  and  $c_i$  are the same for each  $\tau_i$  under CPFBLT and NCPFBLT, then (10) holds if:

$$\forall \tau_i, RC_i^{cp} \leq RC_i^{ncp}$$

$$\begin{aligned} & \delta_i^k (len(s_i^k) - \min(\nabla_{i*}^k)) + \sum_{a=1}^{\min(|\chi_i^k|, m-1)} (len(\chi_i^k(a)) - \nabla_{i*}^k) \\ & \leq \delta_i^k len(s_i^k) + \sum_{a=1}^{\min(|\gamma_i^k|, m-1)} (len(\gamma_i^k(a))) \end{aligned} \quad (11)$$

where  $\gamma_i^k$  is the set of at most  $m-1$  longest transactions conflicting directly or indirectly with  $s_i^k$ . Thus,  $\gamma_i^k(a) \geq \chi_i^k(a), \forall a$ . Thus, by increasing  $\nabla_{i*}^k$ , (11) holds. Claim follows. ■

## VIII. EXPERIMENTAL EVALUATION

We now would like to understand how CPFBLT's retry cost and response time compare with NCPFBLT in practice. Since this can only be understood experimentally, we implement CPFBLT and NCPFBLT and conduct experiments.

We used the ChronOS real-time Linux kernel [28] and the Rochester STM (RSTM) library [13] in our implementation. We implemented G-EDF and G-RMA schedulers in ChronOS, and modified RSTM to include implementations of CPFBLT and NCPFBLT. Checkpointing is implemented using *setjmp/longjmp* instructions. When an object  $\theta$  is accessed for the first time, *setjmp* is used to record a checkpoint. Additionally, a copy of the object is recorded to be restored in case of partial abort. A transaction partially aborts using *longjmp*. We used an 8 core, 2GHz AMD Opteron platform. We used three task sets consisting of 4, 8, and 20 periodic tasks. Each task runs in its own thread and has a set of atomic sections. Atomic section properties are controlled using three parameters: the maximum (*max*) and minimum (*min*) lengths of any atomic section within a task, and the total length (*total*) of atomic sections within any task. Each one of *min*, *max*, *total* lies in  $\{0.2, 0.5, 0.8\}$  provided that  $min \leq max \leq total$ . For each run,

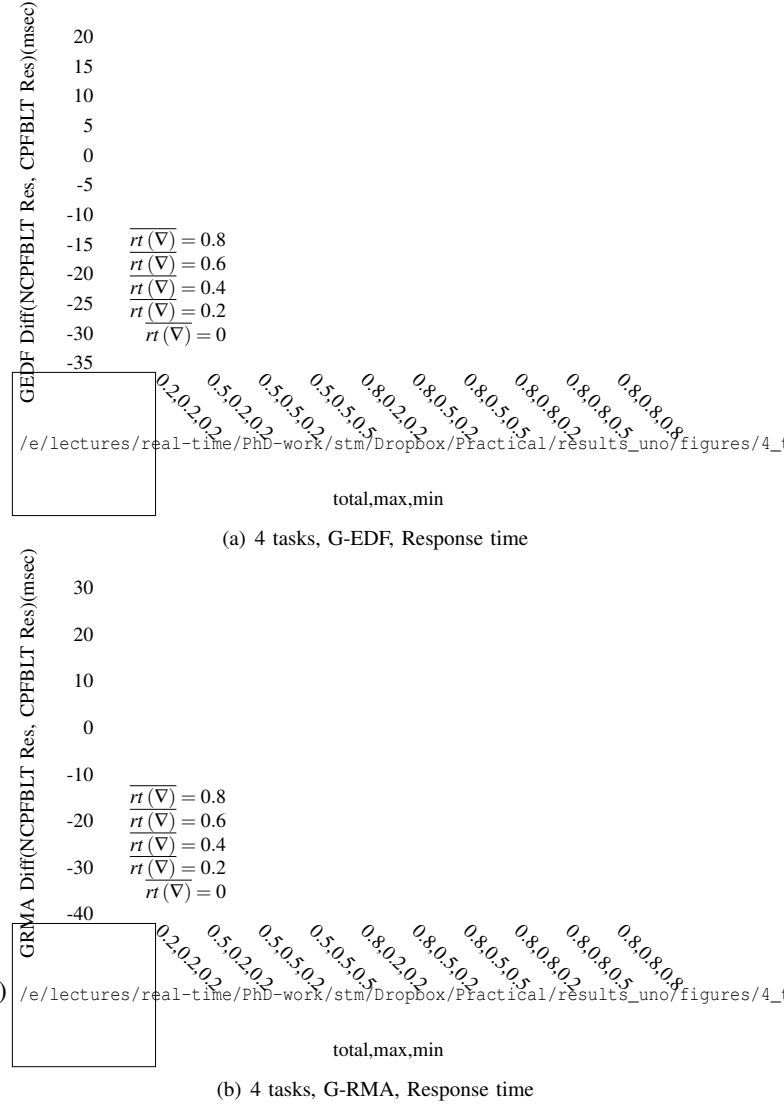
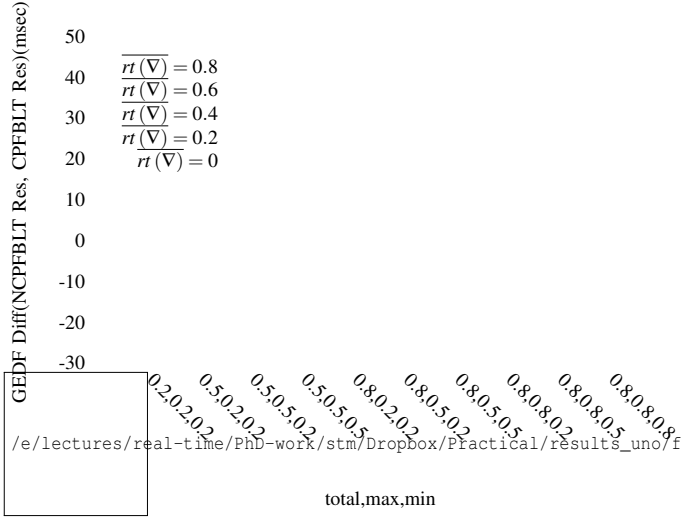


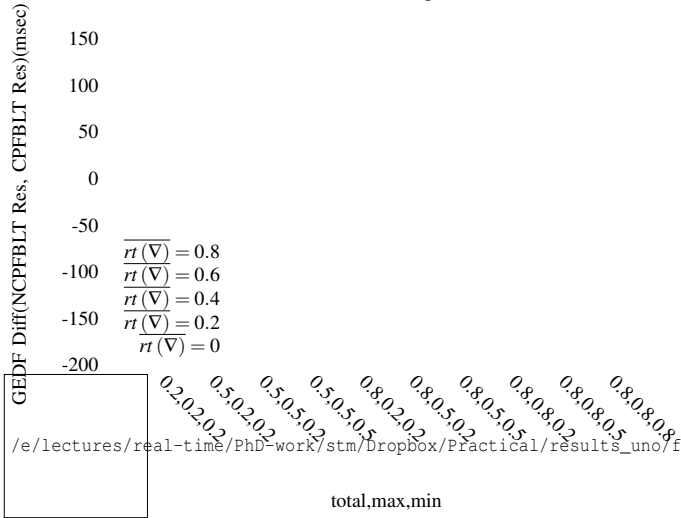
Fig. 1. Response time difference between NCPFBLT and CPFBLT for 4 tasks.

$\overline{rt}(\nabla) = \max(1 - rt(\nabla_{i*}^k))_{\forall s_i^k}$  lies within  $\{0, 0.2, 0.4, 0.6, 0.8\}$ .  $\overline{rt}(\nabla)$  represents the maximum ratio of  $len(s_i^k)$ ,  $\forall s_i^k$  after which accessed objects by  $s_i^k$  can be shared with other transactions. Response time difference between NCPFBLT and CPFBLT for the 4, 8 and 20 task sets are shown in figures 1 and 2. Response time for 8 and 20 tasks under G-RMA show the same trend as response time under G-EDF for 8 and 20 tasks. Response time for 8 and 20 tasks under G-RMA are not shown here due to space limitation. Results show benefits of checkpointing when combined with FBLT. As  $\nabla = 0$ , then all objects accessed by  $s_i^k$  can be shared with other transactions. Consequently,  $s_i^k$  retreats to its beginning, under both CPFBLT and NCPFBLT, if contended objects are accessed at the beginning of  $s_i^k$ . Thus, NCPFBLT can show better response time than CPFBLT as shown in the figures.

Figures 3 and 4 show retry cost difference between NCPFBLT and CPFBLT. Retry cost for 8 and 20 tasks under G-RMA show the same trend as retry cost under G-EDF for 8 and 20 tasks. Retry cost for 8 and 20 tasks under G-RMA



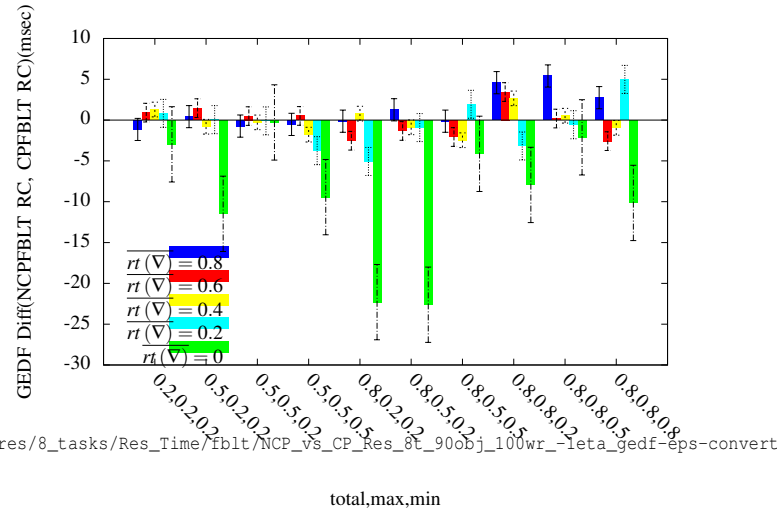
(a) 8 tasks, G-EDF, Response time



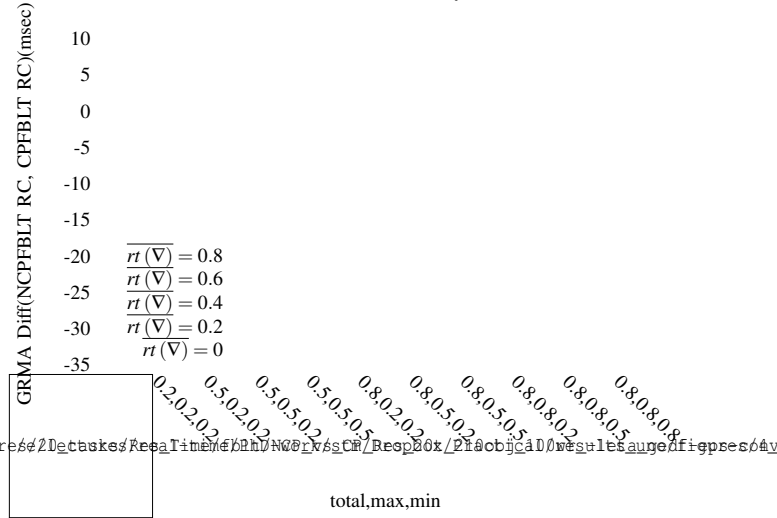
(b) 20 tasks, G-EDF, Response time

Fig. 2. Response time difference between NCPFBLT and CPFBLT under G-EDF for a) 8 tasks. b) 20 tasks.

are not shown here due to space limitation. Retry cost results can be misleading because of the negative difference between retry cost of NCPFBLT and CPFBLT. But this is natural due to behaviour of CPFBLT. Under NCPFBLT, a transaction  $s_i^k$  returns to its beginning upon a conflict with  $s_j^l$  on object  $\theta$ . Whereas, under CPFBLT,  $s_i^k$  returns to the first point it accessed  $\theta$ . Thus, under CPFBLT,  $s_i^k$  tries to access  $\theta$  directly after returning to the proper checkpoint. But  $s_j^l$  is still holding  $\theta$ . Accordingly,  $s_i^k$  will abort and retry again. This retrial (denoted as  $RC_{cp}s_i^k$ ) is added to the accumulated retry cost of  $s_i^k$  under CPFBLT. Under NCPFBLT,  $s_i^k$  returns to its beginning when it conflicts with  $s_j^l$ . Thus, when  $s_i^k$  reaches  $\theta$  once again,  $s_j^l$  may have finished execution. Thus,  $s_i^k$  will not have to abort again due to a conflict with  $s_j^l$  upon  $\theta$ . Let the time between start of  $s_i^k$  and first access to  $\theta$  be  $\nabla_i^k(\theta)$ . Accordingly, retry cost of  $s_i^k$  under NCPFBLT can be less than retry cost of  $s_i^k$  under CPFBLT. Whereas,  $RC_{cp}s_i^k$  can be much less than  $\nabla_i^k(\theta)$ . Thus,  $RC_{cp}s_i^k$  contributes by a smaller value to the response time of  $\tau_i$ , in contrast to  $\nabla_i^k(\theta)$ . This why response time for CPFBLT



(a) 4 tasks, G-EDF, Retry cost



(b) 4 tasks, G-RMA, Retry cost

Fig. 3. Retry cost difference between NCPFBLT and CPFBLT for 4 tasks.

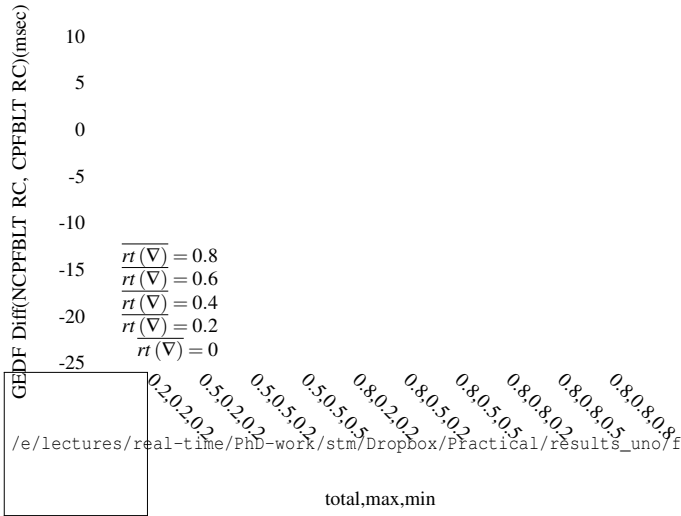
is better than NCPFBLT, whereas, retry cost of NCPFBLT is less than CPFBLT.

## IX. CONCLUSION

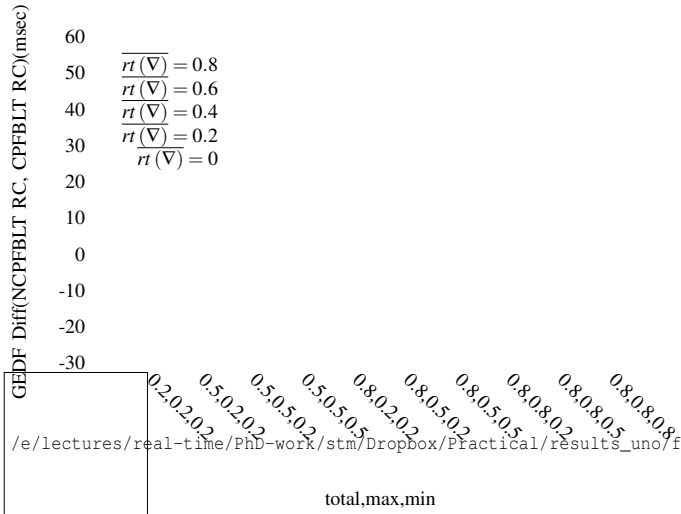
Past research on real-time CMs focused on developing different conflict resolution strategies for transactions. Except for LCM [7], no policy was made to reduce the length of conflicting transactions. In this paper, we analysed effect of checkpointing over FBLT CM. Analysis shows that response time of CPFBLT can be reduced than NCPFBLT if shared objects are accessed close to the end of execution of each transaction. Experimental evaluation reveals better response time for CPFBLT than NCPFBLT. Despite retry cost of NCPFBLT is lower than retry cost of CPFBLT, but this is natural as explained previously. Some CMs make no use of checkpointing due to behaviour of that CM (e.g, under PNF, all non-preemptive transactions are non-conflicting).

## REFERENCES

- [1] M. Herlihy, "The art of multiprocessor programming," in *PODC*, 2006, pp. 1–2.



(a) 8 tasks, G-EDF, Retry cost



(b) 20 tasks, G-EDF, Retry cost

Fig. 4. Retry cost difference between NCPFBT and CPFBLT under G-EDF for a) 8 tasks. b) 20 tasks.

[2] R. Guerraoui, M. Herlihy, and B. Pochon, "Toward a theory of transactional contention managers," in *PODC*, 2005, pp. 258–264.

[3] B. Saha, A.-R. Adl-Tabatabai *et al.*, "McRT-STM: a high performance software transactional memory system for a multi-core runtime," in *PPoPP*, 2006, pp. 187–197.

[4] T. Harris, S. Marlow, S. P. Jones, and M. Herlihy, "Composable memory transactions," *Commun. ACM*, vol. 51, pp. 91–100, Aug 2008.

[5] S. Fahmy and B. Ravindran, "On STM concurrency control for multicore embedded real-time software," in *International Conference on Embedded Computer Systems*, July 2011, pp. 1–8.

[6] M. Elshambakey and B. Ravindran, "STM concurrency control for multicore embedded real-time software: time bounds and tradeoffs," in *Proceedings of the 27th SAC*. ACM, 2012, pp. 1602–1609.

[7] —, "STM concurrency control for embedded real-time software with tighter time bounds," in *Proceedings of the 49th DAC*. ACM, 2012, pp. 437–446.

[8] R. I. Davis and A. Burns, "A survey of hard real-time scheduling for multiprocessor systems," *ACM Comput. Surv.*, vol. 43, no. 4, pp. 35:1–35:44, Oct. 2011.

[9] M. Elshambakey and B. Ravindran, "On real-time STM concurrency control for embedded software with improved schedulability," *18th ASP-DAC*, 2013, available as [http://www.ssrgece.vt.edu/papers/dac\\_](http://www.ssrgece.vt.edu/papers/dac_)

[asp\\_final.pdf](#).

[10] M. Herlihy *et al.*, "Software transactional memory for dynamic-sized data structures," in *Proceedings of the 22nd PODC*. ACM, 2003, pp. 92–101.

[11] M. Elshambakey and B. Ravindran, "FBLT: A real-time contention manager with improved schedulability," *DATE*, 2013, available as [http://www.ssrgece.vt.edu/papers/DATE13\\_SH.pdf](http://www.ssrgece.vt.edu/papers/DATE13_SH.pdf).

[12] E. Koskinen and M. Herlihy, "Checkpoints and continuations instead of nested transactions," in *Proceedings of the twentieth annual symposium on Parallelism in algorithms and architectures*, ser. SPAA '08. New York, NY, USA: ACM, 2008, pp. 160–168.

[13] Marathe *et al.*, "Lowering the overhead of nonblocking software transactional memory," in *Workshop on Languages, Compilers, and Hardware Support for Transactional Computing*, 2006.

[14] J. Anderson, S. Ramamurthy, and K. Jeffay, "Real-time computing with lock-free shared objects," in *RTSS*, 1995, pp. 28–37.

[15] B. Brandenburg, J. Calandrino, A. Block, H. Leontyev, and J. Anderson, "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*, april 2008, pp. 342–353.

[16] S. F. Fahmy, "Collaborative scheduling and synchronization of distributable real-time threads," Ph.D. dissertation, Virginia Tech, 2010.

[17] J. Manson, J. Baker *et al.*, "Preemptible atomic regions for real-time Java," in *RTSS*, 2006, pp. 10–71.

[18] S. Fahmy, B. Ravindran, and E. D. Jensen, "On bounding response times under software transactional memory in distributed multiprocessor real-time systems," in *DATE*, 2009, pp. 688–693.

[19] T. Sarni, A. Queudet, and P. Valduriez, "Real-time support for software transactional memory," in *RTCSA*, 2009, pp. 477–485.

[20] M. Schoeberl, F. Brandner, and J. Vitek, "RTTM: Real-time transactional memory," in *ACM SAC*, 2010, pp. 326–333.

[21] A. Barros and L. Pinho, "Managing contention of software transactional memory in real-time systems," in *IEEE RTSS, Work-In-Progress*, 2011.

[22] A. Turcu, B. Ravindran, and M. Saad, "On closed nesting in distributed transactional memory," in *Seventh ACM SIGPLAN workshop on Transactional Computing*, 2012.

[23] M. Herlihy and E. Koskinen, "Transactional boosting: a methodology for highly-concurrent transactional objects," in *Proceedings of the 13th ACM SIGPLAN Symposium on Principles and practice of parallel programming*, ser. PPoPP '08. New York, NY, USA: ACM, 2008, pp. 207–216.

[24] S. Peri and K. Vidyasankar, "Correctness of concurrent executions of closed nested transactions in transactional memory systems," in *Proceedings of the 12th international conference on Distributed computing and networking*. Springer-Verlag, 2011, pp. 95–106. [Online]. Available: <http://dl.acm.org/citation.cfm?id=1946143.1946152>

[25] J. Kim and B. Ravindran, "Scheduling closed-nested transactions in distributed transactional memory," in *IEEE 26th International Parallel Distributed Processing Symposium (IPDPS)*, may 2012, pp. 179–188.

[26] A. Turcu, "On improving distributed transactional memory through nesting and data partitioning," PhD Proposal, Virginia Tech, 2012, available as [http://www.ssrgece.vt.edu/theses/PhdProposal\\_Turcu.pdf](http://www.ssrgece.vt.edu/theses/PhdProposal_Turcu.pdf).

[27] M. M. Saad and B. Ravindran, "Hyflow: a high performance distributed software transactional memory framework," in *Proceedings of the 20th international symposium on High performance distributed computing*, ser. HPDC '11. New York, NY, USA: ACM, 2011, pp. 265–266.

[28] M. Dellinger, P. Garyali, and B. Ravindran, "ChronOS Linux: a best-effort real-time multiprocessor linux kernel," in *Proceedings of the 48th DAC*. ACM, 2011, pp. 474–479.