

Improving Parallelism in Hardware Transactional Memory

DAVE DICE, Oracle Labs

MAURICE HERLIHY, Brown University and Oracle Labs

ALEX KOGAN, Oracle Labs

Hardware transactional memory (HTM) is supported by recent processors from Intel and IBM. HTM is attractive because it can enhance concurrency while simplifying programming. Today's HTM systems rely on existing coherence protocols, which implement a requester-wins strategy. This, in turn, leads to very poor performance when transactions frequently conflict, causing them to resort to a non-speculative fallback path. Often, such a path severely limits concurrency.

In this paper, we propose very simple architectural changes to the existing requester-wins HTM implementations. The idea is to support a special mode of execution in HTM, called power mode, which can be used to enhance conflict resolution between regular and so-called power transactions. A power transaction can run concurrently with regular transactions that do not conflict with it. This permits higher levels of concurrency in cases when a (regular) transaction cannot make progress due to conflicts and would require a non-speculative fallback path otherwise.

Our idea is backward-compatible with existing HTM systems, imposing no additional cost on transactions that do not use the power mode. Furthermore, using power transactions requires no changes to target applications that employ traditional lock synchronization. Using extensive evaluation of micro- and STAMP benchmarks in a transactional memory simulator and real hardware-based emulation, we show that our technique significantly improves the performance of the baseline that does not use power mode, and performs comparably with state-of-the-art related proposals that require more substantial architectural changes.

ACM Reference format:

Dave Dice, Maurice Herlihy, and Alex Kogan. 2017. Improving Parallelism in Hardware Transactional Memory. *ACM Transactions on Architecture and Code Optimization* 1, 1, Article 1 (July 2017), 22 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Hardware transactional memory (HTM) supports a model of concurrent programming where the programmer specifies which code blocks should be atomic, but not how that atomicity is achieved. Some form of HTM is currently supported by processors from Intel [18] and IBM [6, 19, 33]. Transactional programming models are attractive because they promise simpler code structure and better concurrency compared to traditional lock-based synchronization.

An atomic code block is called a *transaction*. HTM executes transactions *speculatively*: if an attempt to execute a transaction *commits*, that code block appears to have executed instantaneously, while if it *aborts*, that code has no effect, and control passes to an abort handler; a condition code usually indicates why the transaction failed.

Permission to make digital or hard copies of all or part 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 bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2017 Association for Computing Machinery.

XXXX-XXXX/2017/7-ART1 \$15.00

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

50 As long as transactions do not conflict on the shared data they access, and as long as pathologies
51 such as capacity aborts and unsupported instructions are avoided, HTM has been shown to achieve
52 nearly linear scalability [9, 20, 37]. However, the experience shows that even with a moderate level
53 of conflicts between hardware transactions, the performance of HTM substantially deteriorates [9,
54 13, 37]. That is because most existing HTM implementations piggy-back on cache coherence
55 protocols [17], which mostly implement a *requester-wins* policy: if one transaction requests exclusive
56 access to a cache line held by another, the earlier transaction aborts and restarts. Thus, data
57 conflicts cause repetitive transactional aborts [29], which in turn force the execution to proceed
58 through a slower, non-speculative path. This path typically employs locks (e.g., in the very popular
59 transactional lock elision (TLE) method [12]), and taking this path aborts any concurrent speculative
60 transactions, even when there is no actual data conflict between the speculative and non-speculative
61 threads.

62 This paper's contribution is to propose enhancing existing requester-wins HTM systems with a
63 strikingly simple mechanism that allows threads to continue speculating on HTM despite repetitive
64 aborts. The idea is to elevate the status of a transaction that fails to commit so that any conflict
65 between this and other, non-elevated transactions can be resolved in the favor of the former.
66 Thus, the elevated transaction, which we call a *power* transaction or running in a *power mode*,
67 can execute *speculatively in parallel* with transactions it does not conflict with, while impeding
68 progress only of transactions that it does conflict with. In a nutshell, in order to support power
69 mode, each coherence request is augmented with one bit indicating whether the request is coming
70 from a thread speculating on HTM. The power transaction replies with a negative acknowledgment
71 (NACK) to coherence requests from other transactions, causing the requester to abort and allowing
72 the power transaction to proceed. At the same time, regular transactions not conflicting with power
73 transaction(s) can run in parallel with the latter.

74 As an example where this additional parallelism may be beneficial, consider a binary search tree
75 where each tree operation is run in a separate transaction. If two operations try to modify the same
76 node in the tree, they might repeatedly conflict and abort each other. Once one of those transactions
77 decides to abandon speculation and execute under lock, it will cause all other transactions, including
78 those that access a completely different set of nodes in the tree, to wait for its completion. With
79 our proposal, however, that transaction will switch into the power mode, and thus stop the other
80 operation on the same node from aborting it again; all other operations working on different nodes
81 would be able to seamlessly continue their speculation. As described later in the paper, a simple
82 software or hardware mechanism can be put in place to ensure that only one transaction enters the
83 power mode at a time. Such a mechanism is required for performance only, not correctness; that is,
84 multiple power transactions can coexist in the same system as long as they do not conflict, and
85 abort each other if they do.

86 We note that power transactions are “backward-compatible” with existing HTM systems in the
87 sense that allowing hardware transactions the ability to escalate to power mode will not break
88 existing code. Moreover, support for power transactions imposes no additional cost on transactions
89 that do not use the power mode. Furthermore, power transactions can be employed without
90 modifying target applications, e.g., in any case where TLE is applicable. In fact, when the entry to
91 the power mode is controlled by hardware, the existence of the power mode can be completely
92 hidden from the programmer.

93 We used two approaches to evaluate the utility of power transactions. First, exploiting the
94 recent support for compiling transactional programs in GCC [15], we emulated hardware power
95 transactions in software running on top of a real HTM implementation (namely, Intel Haswell
96 TSX). We conducted experiments to compare the relative performance of a number of micro-
97
98

99 and STAMP [24] benchmarks with and without power transactions. (As described below, our
100 software emulation is conservative, in the sense that it tends to understate the advantage of power
101 mode.) With one exception, every benchmark tested yielded improved performance under power
102 mode. These experiments imply that the standard dual-path code structure, in which any non-
103 speculative transaction automatically aborts all speculative transactions, fails to exploit substantial
104 opportunities for concurrent execution. We stress that the potential benefit of the power mode is in
105 unleashing the parallelism in workloads where *some* transactions conflict with each other, allowing
106 those transactions that do not conflict to keep running. In particular, the power mode is not useful
107 when all transactions conflict, nor when none of them conflict. The fact that this parallelism exists
108 in various workloads is not trivial, but our performance evaluation clearly demonstrate that it does.

109 Second, we added the power mode support to SuperTrans [30], a transactional memory simulator
110 built from SESC [32], that was recently enhanced to more accurately simulate best-effort HTM
111 similar to Intel TSX [29]. Using SuperTrans, we compared the utility of power mode to two variants
112 of PleaseTM [29], a recent related proposal for improving parallelism in HTM, as well as to the
113 baseline implementation that does not use power mode. Using STAMP benchmarks [24], we
114 show that power mode not only provides non-trivial speedup above the baseline implementation
115 (confirming our emulation-based study), but also performs similarly to (but slightly better than)
116 both variants of PleaseTM despite requiring less architectural changes.

117 118 2 RELATED WORK

119 In Intel Haswell [18] and its successors, as well as in IBM Power 8 [6], hardware transactions
120 are best-effort: no transaction is guaranteed to commit. Transactions may abort because of data
121 conflicts, cache overflow, or cache associativity issues. Transactions must not execute certain
122 instructions, such as I/O instructions and system calls.

123 In these systems, progress is usually guaranteed by combining HTM with some form of lock-
124 ing. Perhaps the simplest and most widely-used such technique is the transactional lock elision
125 (TLE) [12], where the critical section associated with a lock is first attempted speculatively, trans-
126 actionally reading but not writing the lock state. TLE is attractive, because it can be enabled, without
127 any changes to the target application, at the level of a library providing lock implementations while
128 preserving the semantics provided by the lock based synchronization [9, 12]. If the speculative lock
129 elision fails (typically, after a few retries), the thread acquires the lock and re-executes the critical
130 section non-speculatively. TLE provides the same progress guarantees as regular locking, but it has
131 a non-trivial cost: once the lock has been acquired, all concurrent speculative transactions will fail
132 and wait until the lock is released, even if there are no actual data conflicts. As a result, numerous
133 papers show that TLE is very effective when most transactions succeed, but its benefit fades once
134 the lock is acquired often [9, 13, 37]. In order to keep our usage examples of power transactions
135 concrete, we focus on the use of locks as the alternative path, effectively enhancing the standard
136 TLE technique [12]. We note, though, that power mode is equally helpful in reducing the use of any
137 fallback path, such as the one implemented using software transactional memory (STM) [7, 23],
138 lock-free techniques [21], etc.

139 The use of a special execution mode for (software or hardware) transactions has been previously
140 explored in related contexts. Blundell et al., for instance, design a system called OneTM that supports
141 unbounded hardware transactions [3]. One of the variants of OneTM, called OneTM-Concurrent,
142 supports concurrent execution of non-overflowed transactions and non-transactional code with one
143 overflowed transaction. In order to support this mode of execution, OneTM-Concurrent requires,
144 among other architectural changes, additional metadata storage and management in memory
145 controllers as well as an additional architected register, saved and restored on every context
146

148 switch. Power mode, being designed to enhance existing (bounded) HTM implementations, does
149 not require any of those complications.

150 In the context of software transactional memory, Ni et al. [28] describe an STM runtime library
151 that supports multiple modes of executions. One of the modes, called *obstinate*, is a software
152 equivalent of power transactions. Citing from [28], "a transaction running in obstinate mode always
153 wins all conflicts with other transactions – regular transactions are allowed to run concurrently
154 with the obstinate one, but the obstinate transaction has the highest conflict resolution priority
155 of all transactions in the system". As expected from any STM system, however, the control over
156 execution modes and conflict resolution between transactions in [28] is done entirely in software,
157 in a dedicated contention manager module of the system.

158 Since the introduction of the HTM design by Herlihy and Moss [17], numerous attempts have been
159 made to improve and extend it (e.g., [1, 22, 25, 26, 31, 36] to give just a very few examples). The scope
160 of this paper precludes elaborating on all these efforts. We note, however, that, broadly speaking,
161 the architectural changes required by most of them go far and beyond the ones needed to support
162 the power mode. Furthermore, the prime goal they pursue (e.g., supporting hardware transactions
163 with unbounded capacity [1, 22, 25, 36], running transactions across context switches [31, 36],
164 supporting nested transactions [26, 36], etc.) is typically different from the one considered in this
165 paper.

166 Perhaps the most relevant prior work to this paper is the recent publication by Park et al. on a
167 PleaseTM mechanism for requester-wins HTM systems [29]. In PleaseTM, hardware transactions
168 insert plea bit (or bits) into their responses to coherence requests. These plea bits are considered
169 by the requester and allow supporting alternative conflict resolution schemes. For instance, a
170 requester running a hardware transaction and receiving a response with the plea bit set may abort
171 its transaction, effectively achieving a responder-wins conflict resolution strategy for hardware
172 transactions. By keeping track of the number of transactionally read cache lines and encoding
173 this information in a number of plea bits, PleaseTM allows a scheme where a transaction with
174 more lines read wins a conflict. While requiring several architectural changes, PleaseTM does not
175 modify the cache coherence protocol itself, meaning that a pleading transaction releases a cache
176 line upon receiving a coherence request. Consequently, in order to ensure atomicity, the pleading
177 transaction needs to re-request the cache line and validate the line data when the line is re-acquired.
178 As a result, even when a requesting transaction decides to abort (i.e., accept the plea), it slows
179 down the responding transaction, and would do that over and over again with every retry. This
180 mechanism also puts pressure at the coherence bus, especially when the requesting transaction
181 runs on a different socket. In opposite to PleaseTM, power mode allows resolution of conflicts at
182 the time the request is received, without slowing down the responding transaction and without
183 increasing coherence traffic. Furthermore, transactions in PleaseTM can still repeatedly abort each
184 other as long as they respond with pleading bits to each other's requests. At the same time, power
185 mode provides more definitive control with respect to which transaction would receive priority
186 over others.

187 In another relevant paper, Armejach et al. consider a few hardware and hybrid (software and
188 hardware) techniques to improve performance of requester-wins HTM [2]. Perhaps the most
189 relevant technique to our work is the one called DRW (delayed requester-wins). The idea behind
190 DRW is to allow the exclusive owner of a cache line to delay response to conflicting requests, thus
191 increasing the chance for its transaction to complete. Delayed conflicting requests are queued at
192 the exclusive owners caches and are considered when the transaction ends (by commit or abort).
193 To avoid deadlocks, DRW associates timeouts with buffered requests and conservatively handles
194 a request when its timer expires. The requirement to manage the buffers of incoming conflicting
195
196

197 requests and their associated timers call for hardware changes that are much more elaborated than
198 supporting power transactions.
199

200 3 POWER TRANSACTIONS

201 We first describe the common mechanism required to support power transactions regardless of how
202 transactions enter the power mode. Next, we discuss the details of supporting a software-controlled
203 entry into the power mode, followed by details on a hardware-controlled entry. We note that those
204 two entry methods are complementary, i.e., we envision that some architectures may provide both
205 methods, alike the support for HLE and RTM in Intel TSX [18]. Note that exiting power mode
206 does not require any special treatment, i.e., power transactions commit or abort exactly as regular
207 transactions. Finally, we discuss variations of our design along with their impact on the properties
208 of power transactions.
209

210 3.1 Common Mechanism

211 Supporting the power mode requires each hardware thread to maintain a distinctive speculation
212 status that can be encoded in one bit of a thread state. That is, in addition to the status indicating
213 that a given hardware thread is speculating on HTM, we need to store a bit of information (that
214 we call the power-mode bit), which, when set, indicates that the speculating thread is running in
215 power mode. This bit will be set when a hardware thread starts a new power transaction (via one
216 of the mechanisms described in the subsequent sections) and reset when the thread completes that
217 transaction (either through commit or abort).
218

219 In addition, we require to add a speculation status bit as a simple payload to coherence request
220 messages. This bit indicates whether the request is coming from a thread speculating on HTM
221 (either in power or regular modes). It is ignored by the coherence hardware and is simply passed to
222 cache controllers, which in turn can take it into consideration when preparing the corresponding
223 coherence response.

224 Cache controllers are modified so that when the following three conditions hold, they respond
225 with a special NACK message: (1) the speculation bit in the incoming request is set, (2) the power-
226 mode bit of the target thread is set, and (3) the request is to invalidate or downgrade transactionally-
227 held data. If any of those conditions does not hold, the cache controller logic remains unchanged.
228 We note that in order to support (regular, non-power-mode) HTM, the cache controller already
229 implements logic to consider the speculation state of the target as well as whether the request is to
230 invalidate or downgrade transactionally-held data (so that the hardware transaction run by the
231 target thread can be aborted). Thus, the additional complexity of considering the speculation bit in
232 the request payload is trivial.

233 Another modification in the cache controller is related to the treatment of the NACK coherence
234 response message. Specifically, when the NACK is received and the receiving thread is speculating
235 on HTM (either in power or regular modes), the current transaction is aborted; a special abort code
236 may be used to specify that the abort occurred due to a data conflict with a power transaction.
237 Otherwise, the NACK response is ignored. This can happen only if the hardware transaction
238 that issued the coherence request, which resulted in the NACK response, has been aborted while
239 awaiting that response. Figure 1 illustrates possible interactions between a thread running a power
240 transaction (P), a thread running another (regular or power) transaction (R), and a thread running
241 non-transactional code (N).

242 The only modification to the cache coherence protocol is supporting a special NACK response
243 message (if it does not already support one). We note that numerous previous papers on computer
244 architecture considered adding a NACK message (e.g., [4, 22, 25, 36]). As opposite to most of that
245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

289

290

291

292

293

294

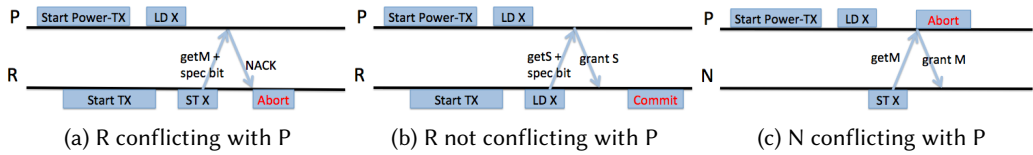


Fig. 1. Interactions between a thread running a power transaction (P) and another thread running a (regular or power) transaction (R) or non-transactional code (N).

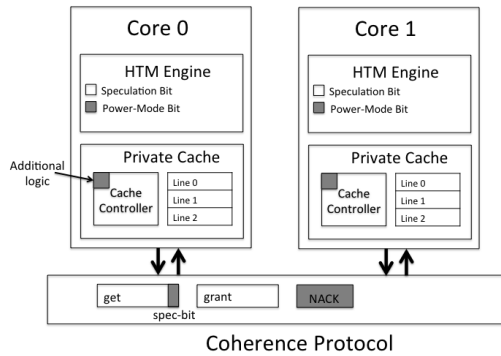


Fig. 2. Architectural changes required to support power transactions.

work, however, the very limited use of NACKs in our case does not result in any additional changes in the coherence protocol, such as managing timestamps and adding deadlock-avoidance logic, or in coherence messages, such as including timestamps in every coherence request message.

The required changes in an HTM architecture are summarized in Figure 2 with new or modified components shown in grey. We believe that all the proposed changes are simple, if not trivial. It is hard to compare in a quantitative way the amount of hardware changes required by various proposals to existing architectures, especially given that those proposals often modify different components of the architectures. Yet, adding support for power mode stands out for its simplicity, and is substantially less intrusive when compared to any related proposal known to us as described in Section 2. This makes supporting power mode more feasible in real systems.

Note that the common mechanism as described so far does not limit or control the number of power transactions that can coexist in the system. Such control is provided through entry mechanisms described in the subsequent sections.

3.2 Software-Controlled Entry

In order to allow software to control which transaction(s) would run in a power mode, one new instruction should be added to the instruction set architecture (ISA). This instruction should be virtually identical to the one used to begin a (regular) hardware transaction, but use a different opcode, which would instruct the core executing it to set the power-mode bit of the corresponding hardware thread. As mentioned above, there is no need to add a new instruction(s) for completing the power mode transaction.

295 With the software-controlled entry to power mode, it is the responsibility of the programmer to
296 ensure (or not) that only one transaction switches into the power mode at a time. As a common
297 case, at most one transaction accessing shared data would run in the power mode. However, we
298 stress that this is merely a performance consideration rather than correctness. In fact, as we discuss
299 in Section 3.4, there are cases (e.g., performance debugging) in which having multiple transactions
300 in power mode is desirable.

301 Following is one particular approach for ensuring that there is at most one transaction in the
302 power mode. This approach is integrated with a common implementation of the TLE mecha-
303 nism [12], and assumes that the target application uses locks to synchronize access to shared data
304 inside critical sections. As mentioned in Section 2, TLE can be enabled *without any changes* to the
305 target application, e.g., by interposing the library providing the lock implementation. The power
306 mode preserves this property – using a TLE mechanism with power transactions, such as the one
307 described below, does not require modifying target applications as well.

308 In the TLE mechanism enhanced to exploit power transactions, the entry into power mode is
309 protected by a lock. For simplicity, we use a spin lock, but other locks are possible (e.g., a queue lock
310 for fairness). Figure 3 shows pseudo-code for such an enhanced TLE mechanism. Here, a transaction
311 elevates into the power mode if it repeatedly fails to commit. The atomic compare-and-swap (CAS)
312 instruction (Line 24) ensures that only the thread that sets the powerFlag flag to its thread ID will
313 enter the power mode. We note that using thread IDs can be avoided, e.g., by using a thread-local
314 flag that tells the current thread whether it is the one that entered the power mode. Notice that
315 the code does not access the powerFlag flag inside a hardware transaction. Thus, starting (and
316 committing) a power transaction does not abort regular transactions.

317 When using power mode, regular transactions may be subject to the *lemming effect* [11] arising
318 when one transaction enters power mode and forces the rest to follow; this effect exists with
319 (regular) transactions in standard TLE as well [11]. One way to mitigate the lemming effect is to
320 give less (or even zero) weight for retries happening while the powerFlag flag is set. That is, if
321 an attempt to use a regular transaction fails and powerFlag is set, we discount this attempt by
322 decrementing the `retries` counter (Line 16). We note that the pseudo-code in Figure 3 also includes
323 a standard anti-lemming optimization in TLE, in which a transaction is retried only when the lock
324 becomes available (Line 28) [20].

3.3 Hardware-Controlled Entry

326 Along with (or instead of) a software-controlled entry into power mode, the HTM engine itself may
327 control when a regular transaction switches into the power mode. In this case, no ISA extensions
328 are required. In fact, the availability of power mode for hardware transactions may be completely
329 hidden from the programmer in this case.

330 A simple hardware-based scheme can be put in place to ensure that there is only one power trans-
331 action at a time, either system-wide or for each process. There are many ways to implement such
332 functionality, which requires the ability to arbitrate concurrent requests from multiple hardware
333 threads. One option, similar to the proposal made in [3], is to add a shared (between all hardware
334 threads) transaction status word, which resides in a fixed location in the virtual address space
335 of each process. This word acts as a mutex lock, i.e., the thread enters the power mode only if it
336 atomically sets the value of the word and exits that mode when it atomically resets it. Unlike the
337 proposal in [3], however, regular transactions do not need to monitor this word and perform any
338 special logic for conflict detection when it is set.

339 Considering the HLE mechanism in Intel TSX [18] as an example, when the hardware thread
340 encounters a lock instruction with the opcode prefix that allows speculation, it may start speculation
341
342
343

```

344 initially : (global) lock = 0;
345             (global) powerFlag = -1;
346             (thread local) tID = unique thread ID;
347
348 Lock procedure:
349
350 1  ntrials = 0;
351 2  while (true) {
352 3    // start a regular or power transaction according to the value of 'powerFlag'
353 4    if ((powerFlag != tID && begin_htm()) ||
354 5         (powerFlag == tID && begin_power_htm())) {
355 6         if (isLocked(&lock)) self-abort();
356 7         return;
357 8     }
358 9     if (!self-aborted) {
35910        if (powerFlag == tID) {
36011            // exit the power mode and fall to the lock
36112            powerFlag = -1;
36213            break;
36314        } else if (powerFlag != -1) {
36415            // avoid the lemming effect
36516            ntrials --;
36617        }
36718        // increase the counter for the number
36819        // of non-power mode trials
36920        if (++ntrials >= MAX_TLE_TRIALS) {
37021            // if we exhausted the number of non-power mode trials, check if the 'powerFlag' flag is
37122            // available and try to set it.
37223            // Note: for TLE we would simply break here and fall to the lock
37324            if (powerFlag == -1) CAS(&powerFlag, -1, tID);
37425        }
37526    }
37627    // wait for the lock to become available
37728    while (isLocked(&lock)) { pause; }
37829 }
37930 // we failed to commit a transaction, grab the lock
38031 Lock(&lock);

```

```

381
382
383
384
385 Unlock procedure:
386
387 1  if (!isLocked(&lock)) {
388 2    commit_htm();
389 3    if (powerFlag == tID) powerFlag = -1;
390 4  } else {
391 5    Unlock(&lock);
392 6  }

```

Fig. 3. TLE using power transactions

using a regular transaction. If aborted, it may try to atomically set the transaction status word, and if succeeded, it shall set its power-mode bit, and run a power transaction. Upon completion (either abort and commit), it shall reset the power-mode bit and the shared transaction status word. If the thread fails to set the transaction status word, which means that another power transaction is in progress, it may retry with a regular transaction or, if the preset retry policy instructs so, execute the lock instruction non-speculatively.

3.4 Variations

There are a few interesting extensions for the common mechanism discussed in Section 3.1. First, we may omit including the speculation status bit in the coherence request messages. A power transaction then will send NACKs in response to all invalidation and downgrading requests, not just requests from transactions. A transactional thread (regular or power) that receives a NACK simply aborts, and a non-transactional thread backs off (pauses) and resends its request. This approach has some advantages: it alleviates the need to introduce a new bit into coherence messages' payload, and it protects power transactions against conflicts with non-transactional threads (but not with other power transactions). The principal disadvantage is that care must be taken to avoid denial-of-service vulnerabilities, perhaps by limiting the duration during which a power transaction can refuse invalidations or by throttling (at the hardware level) the rate at which transactions can enter the power mode. Furthermore, the need to introduce a back-off mechanism into the cache controller logic may complicate the support for power mode.

Second, the power mode support could be easily generalized to encompass multiple levels of power transactions. Instead of a single power-mode bit, each hardware thread state may include a power-mode counter indicating the level at which the thread is running a power transaction. The payload of cache coherence messages is respectively enhanced to include this counter. Higher-priority transactions refuse invalidation and downgrading requests from lower-priority transactions, effectively providing a kind of transactional priority system, which may be a start toward adapting transactional programming to reactive systems [35]. It is straightforward to enhance both software and hardware-controlled entry mechanisms discussed above to climb through the levels of power mode before resorting to a non-speculative execution. It should be noted that in the software-controlled entry, a new ISA instruction for starting a power mode transaction should include a level argument. Along with that, no further changes are required for the hardware-controlled entry beyond increasing the number of transaction status words to match the number of power mode levels (as long as the maximum power mode level can be stored in a transaction status word that can be updated atomically).

As mentioned in Section 3.1, when a transaction receives a NACK and aborts, it may specify a special abort code providing indication to the programmer of a conflict with a power transaction. Taking this a step further, we may use a different abort code to indicate that the recipient of the NACK *was running in the power mode* as well. This abort code provides a way to detect undesired data sharing between transactions. A transactional undesired data sharing occurs when two transactions that are believed to have disjoint data sets actually have a data conflict. Such hidden conflicts can result from false sharing, from hidden data accessed by library calls, or from performance counters and related structures. Such conflicts can cause transactions to abort more often than expected, adversely affecting system performance. To test whether two transactions have disjoint data sets, run them concurrently in power mode, and if one aborts with the power-mode conflict abort code, then the transactions' data sets are not disjoint, and there is a possibly unexpected data sharing. Note that the special abort code would allow detection and, potentially, a repair of both true and false sharing issues, facilitating recent work that employs hardware performance counters for that purpose [14]. Exploring this potential benefit of supporting the power mode is left as future work.

4 EMULATION-BASED EVALUATION

We have evaluated the utility of power mode with two complementary approaches. In this section we describe our attempt to emulate power mode transactions in software (i.e., running them without HTM), while regular transactions run on top of HTM. This approach is inspired by work on hybrid transactional memory systems [8], and in particular, by the implementation of refined TLE [9].

442 We note that this is not our intent to compare power transactions to hybrid TMs (which use a
443 software-only code path), but rather to evaluate if and how the existence of power mode support
444 can increase the parallelism of hardware transactions and ultimately improve performance of
445 the existing HTM implementations. Our second evaluation approach is based on a transactional
446 memory simulator, and is described in Section 5.

447

448 4.1 Framework

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

4.1.1 *High-Level Idea.* Our experience shows that the time required for a successful execution of an atomic block using a hardware transaction is comparable to running that block in software¹. This is also echoed by results of single-thread performance in various papers [9, 27, 37]. We leverage this fact to emulate power mode transactions in software, while using an actual HTM implementation (Intel Haswell, in our case) to execute regular transactions.

In order to mimic the behavior of a hardware power-mode implementation and resolve data conflicts between power and regular transactions in favor of the former, we utilize software “metalocks”. These metalocks are implemented with the use of ownership records, or orecs, commonly used in the design of software and hybrid transactional memory systems [8, 16]. We instrument all memory accesses in transactions to read and update ownership records as appropriate, leveraging a recently introduced compiler support for a completely automatic instrumentation process. Thus, both power and regular transactions run on the instrumented path. A power transaction acquires ownership on memory words it accesses by writing into ownership records. Regular transactions check (by reading ownership records) that their memory accesses are not conflicted with those made by the concurrent power mode transaction, if such exists. The ownership records are designed in a way that regular transactions are aborted only when an actual conflict exists (as they should). In particular, a regular transaction is *not* aborted when it reads the same data as a power transaction does. Furthermore, a simple mechanism is put in place to clear all ownership records at once when the power mode transaction is completed, either by abort or commit (details are provided in Appendix A).

Our framework effectively adds the power mode to an *existing* HTM implementation, leveraging all of its properties for running and managing (regular) hardware transactions. In the emulated system, all instructions but loads and stores have absolutely the same latency as provided by the native platform. Load and store instructions are slowed down due to the use of instrumentation. We note that both power and regular transactions are slowed down, so the *relative* performance of these transaction types using the software metalocks is a way to estimate their relative performance in a hardware implementation. Indeed, because power transactions, unlike regular transactions, might be required to write each time they read (to acquire corresponding metalocks), our estimation is conservative, favoring the relative performance of regular transactions. Despite the impacts of instrumentation, which depend on the number of loads and stores in a critical section [9], we believe that the ability to exploit an actual HTM implementation as well as the ability to use arbitrary benchmarks make our framework an interesting tool able to provide important insights on performance of power mode. In the following subsection, we expand on implementation aspects of our framework.

4.1.2 *Implementation Details.* The GCC compiler [15] (starting from version 4.8) provides the libitm interface for transactional programs. The compiler translates critical sections implemented

¹ This is true at least as long as transactions are not tiny, so the overhead of starting and committing a hardware transaction is negligible.

491 as atomic transactions into two distinct code paths: instrumented and uninstrumented. The instru-
492 mented code path includes calls to *instrumentation barriers*, functions invoked on each transactional
493 memory access. The libitm library provides instrumentation barriers for a few standard synchron-
494 ization mechanisms, such as TLE, STM, or lock synchronization, as well as the opportunity to
495 provide customized instrumentation barriers and functions to be called when transactions commit
496 or abort. For our framework, however, we used our own custom implementation of the libitm
497 interface to reduce instrumentation overheads².

498 We associate two metalocks with each cache line accessed by a transaction, one for read access
499 and another for write access. A power transaction (run without HTM) acquires metalocks for the
500 cache lines it accesses, according to the access mode desired (read or write) by writing a value into
501 the corresponding metalock. A regular transaction (run with HTM) reads the metalocks associated
502 with the cache lines it accesses, and aborts if it finds a metalock held in a conflicting mode. If
503 the metalock does not conflict, the transaction proceeds to access the intended data. This scheme
504 emulates power mode semantics, ensuring that any conflicting (and only conflicting) request by a
505 regular transaction for data accessed by a power mode transaction is refused, causing that regular
506 transaction to abort.

507 The instrumentation increases every transaction’s data footprint, doubling the number of accessed
508 cache lines. However, regular transactions access the additional (metalock) cache lines only for
509 reading, which does not stress Haswell HTM, whose read set capacity is relatively large [27].
510 Although power transactions access metalocks for writing, they do not use HTM and thus are not
511 limited by its write capacity.

512 Power transactions sustain conflicts with regular transactions, but they can abort for other
513 reason, and in particular, due to capacity limitations. A direct way to emulate capacity aborts for
514 a power transaction is to detect when a capacity limit is reached, roll back that transaction, and
515 restart it using locks. This direct approach, however, requires logging each transaction’s write set
516 and reverting its memory updates on abort, further increasing instrumentation overhead. Instead,
517 we opted for the following less intrusive emulation. First, each power mode transaction takes a
518 timestamp when it begins its execution. Second, we track the number of cache lines accessed by a
519 power mode transaction. Once this number goes beyond a preset limit, we calculate the time period
520 δ elapsed since the transaction started, switch to the locking mode, and spin for another time period
521 δ , effectively charging twice for the transaction so far. Once a power mode transaction switches
522 to the locking mode (simply by setting a Boolean flag), all regular transactions are aborted (as in
523 standard TLE) and wait for the lock to become available again. By spinning after lock acquisition (for
524 the time period δ), we “charge” for the time required to re-execute the same atomic block without
525 actually rolling back the changes made by the power mode transaction and without reapplying
526 them under lock. A reader interested in further implementation details, including the metalock
527 structure and the pseudocode of instrumentation barriers, is referred to Appendix A.

528 Our experiments were run on an Intel Haswell (Core i7-4770) 4-core hyper-threaded machine (8
529 hardware threads in total). Before starting measurements, all threads were set to spin for a few
530 seconds to allow the system to warm up. Our goal was to compare standard TLE [12] with one
531 that makes use of power transactions (henceforth *PowerTLE*). The pseudo-code for PowerTLE is
532 provided in Figure 3. To evaluate the benefit of the additional concurrency provided by power
533 mode, and to reduce the impact of other unrelated factors, such as the cost of instrumentation or
534 transactions’ increased memory footprints, we used exactly the same instrumentation barriers for
535

536 ²The library implementing the libitm interface in GCC is dynamically linked to an executable, resulting in an expensive
537 function call for every memory access on the instrumented path. Our custom implementation of the libitm interface
538 supports static linkage with the target executable.

540 TLE as well. We emphasize that the prime difference between TLE and PowerTLE in our framework
541 is the ability of the latter to use a power transaction that runs concurrently with regular transactions
542 as long as those two kinds of transactions do not actually conflict on shared data.

543 Each critical section was attempted ten times using regular transactions before reverting to lock
544 (in TLE) or power mode (in PowerTLE). As demonstrated in Figure 3, a power mode transaction is
545 tried only once (or more, in case of self-abort indicating that the lock is taken). We note, however,
546 that other, more sophisticated retry policies that use power mode can be put into place. Although
547 finding an optimal lock elision retry policy is an interesting question by itself [10, 13], it falls out of
548 scope of this paper.

551 4.2 Skip List-based Priority Queues

552 Figure 4 shows throughput results of a priority queue microbenchmark that uses a standard
553 skip list implementation as an underlying data structure. The results shown are the average of
554 ten runs performed in the same configuration. The breakdown of operations between different
555 modes of executions, e.g., regular transactions, power transactions, etc., is presented in Figure 5.
556 For PowerTLE, we report separately regular transactions completed without any power mode
557 transaction running concurrently with them (denoted as NonC TXs) and those completed while
558 some power mode transaction was running (denoted as C TXs).

559 For the experiment reported in Figure 4 (a), the queue is initialized with 100K elements, and all
560 threads run a total number of 100K RemoveMin operations, divided equally among the participating
561 threads. We measure the time from the start till the last thread is done with its operations, and
562 calculate throughput by dividing the total number of performed operations (100K) by this time.
563 In this particular workload, all threads compete with each other over the minimal element in the
564 queue. Not surprisingly, except for two threads, power mode does not increase throughput, since
565 a power mode transaction conflicts with every other regular transaction and thus aborts them.
566 This is echoed by results in Figure 5 (a) showing that only very few regular transactions manage
567 to complete, while the majority of operations is executed using a lock (in TLE) or power mode
568 transactions (in PowerTLE). The case of two threads is slightly different, and shows that substantial
569 portion of regular transactions completes concurrently with a power-mode transaction. Indeed,
570 this is the only point where PowerTLE beats TLE by a large gap (cf. Figure 4 (a)). We believe this
571 happens because a regular transaction (running a very short RemoveMin operation) manages to
572 “sneak in” without any contention while another thread transitions into power mode and before the
573 actual data conflict occurs. In TLE, all transactions are aborted at the moment the lock is acquired,
574 and thus transactions do not have enough time to complete when another thread switches to lock.
575 When we increase the number of threads, this benefit of PowerTLE fades as regular transactions
576 conflict with each other.

577 In the experiment reported in Figure 4 (b), the queue is initialized with 100K elements, and each
578 thread runs loop iterations for 5 seconds, where in each iteration it chooses randomly to remove
579 a minimal element or insert a random element into the queue. Here the increased concurrency
580 provided by power mode starts to take effect as the number of threads increases. This is because
581 when a thread runs, e.g., an Insert operation in power mode, other threads can proceed concurrently
582 to apply their non-conflicting operations. As a result, at 8 threads, PowerTLE achieves almost 2x
583 more throughput than TLE. Figure 5 (b) shows that, indeed, some portion of regular transactions
584 manages to complete concurrently with a power transaction, and this portion grows with the
585 number of threads. Interestingly, the portion of regular transactions completing non-concurrently
586 with a power transaction is also larger for PowerTLE than the portion of transactions in TLE. We
587
588

589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637

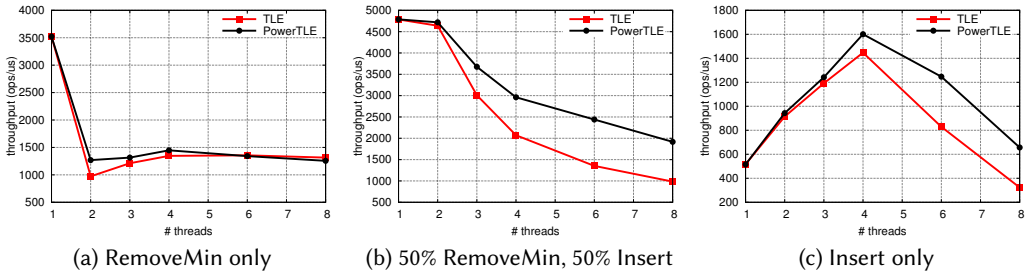


Fig. 4. Skip list-based priority queue throughput. Higher is better.

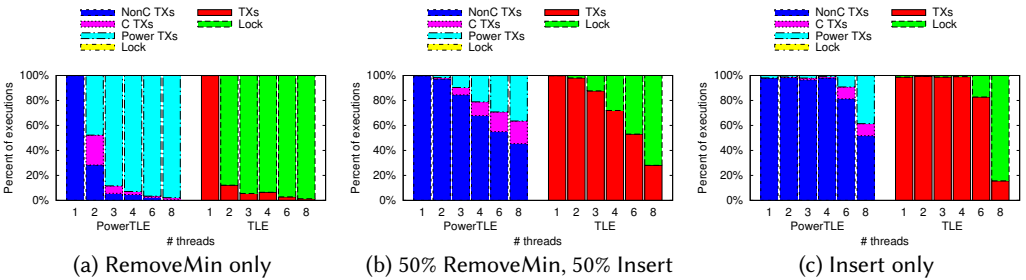


Fig. 5. Breakdown of execution modes for operations in skip list-based priority queue benchmark. “C TXs” (“NonC TXs”) stand for transactions executed concurrently (non concurrently, respectively) with a power transaction.

attribute that to the decreased lemming effect [11] that power transactions have comparing to lock, as the former do not abort all transactions but only those conflicting with them.

The benefit of PowerTLE over TLE increases even further when we consider only Insert operations that are less likely to conflict with each other compared to RemoveMin operations. Figure 4 (c) shows the results of the experiment where the queue is initially empty and all threads perform a total number of 100K insert operations, divided equally among threads. At 8 threads, PowerTLE achieves more than 2x throughput of TLE. When the number of threads grows, the improved concurrency of PowerTLE becomes evident with the increase in the portion of regular transactions executed while a power mode transaction was running (cf. Figure 5 (c)).

4.3 AVL Tree-based Sets

In this section, we discuss results of a set microbenchmark implemented on top of AVL trees. The AVL tree implementation is similar to the one found in OpenSolaris OS. In all experiments, each thread runs iterations for 5 seconds, and in each iteration it chooses an operation and a key. The operations are randomly selected from a given workload distribution, while the key is randomly selected from a given range from 0 to 511. The set is initialized to contain half of the given key range (256 keys).

Figure 6 (a) shows results for the read only workload where all threads perform only Find operations. Here, the vast majority of operations succeed without any retries, and thus power mode is not used. Thus, both TLE and PowerTLE yield identical performance. The breakdown of

638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686

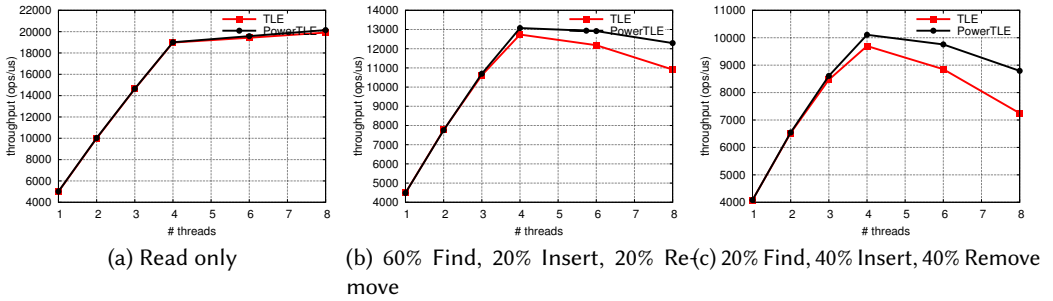


Fig. 6. AVL tree-based set throughput. Higher is better.

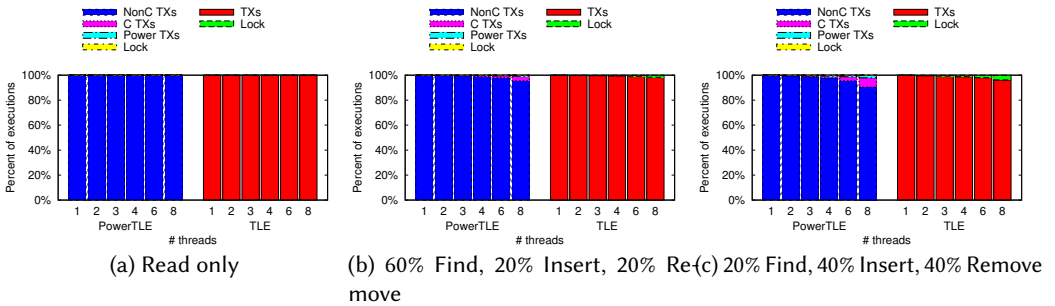


Fig. 7. Breakdown of execution modes for operations in AVL tree-based set benchmark.

execution modes shows that, indeed, virtually all operations succeed using regular transactions (cf. Figure 7 (a)). This is not surprising, as the operations do not conflict with each other.

The workloads in Figure 6 (b) and (c) include update operations. Specifically, in the former threads perform 60% Find operations, while in the latter threads perform 20% Find operations; the rest is divided equally between Insert and Remove. Here, as the number of threads grows, some transactions fall back to the lock (in TLE) as they experience conflicts on data they access. As a result, the benefit of increased concurrency provided by PowerTLE becomes more significant as the number of threads and/or the portion of update operations increases. The breakdown of execution modes for these workloads (Figures 7 (b) and (c), respectively) confirms that as the number of threads increases, more regular transactions manage to complete concurrently with a power transaction in PowerTLE, rather than falling to the lock as they would with TLE.

4.4 STAMP

This section presents results measured with the STAMP benchmarking suite [24], which is used extensively in transactional memory research³. For each benchmark, we used a standard ('native') set of command line parameters. Figure 8 shows running time reported by each benchmark, averaged over ten runs. We omit the results for one of the STAMP benchmarks (namely, bayes) due to extremely high variance (which was also observed by others [29, 37]).

³ We used a version of STAMP available at <https://github.com/mfs409/stamp>.

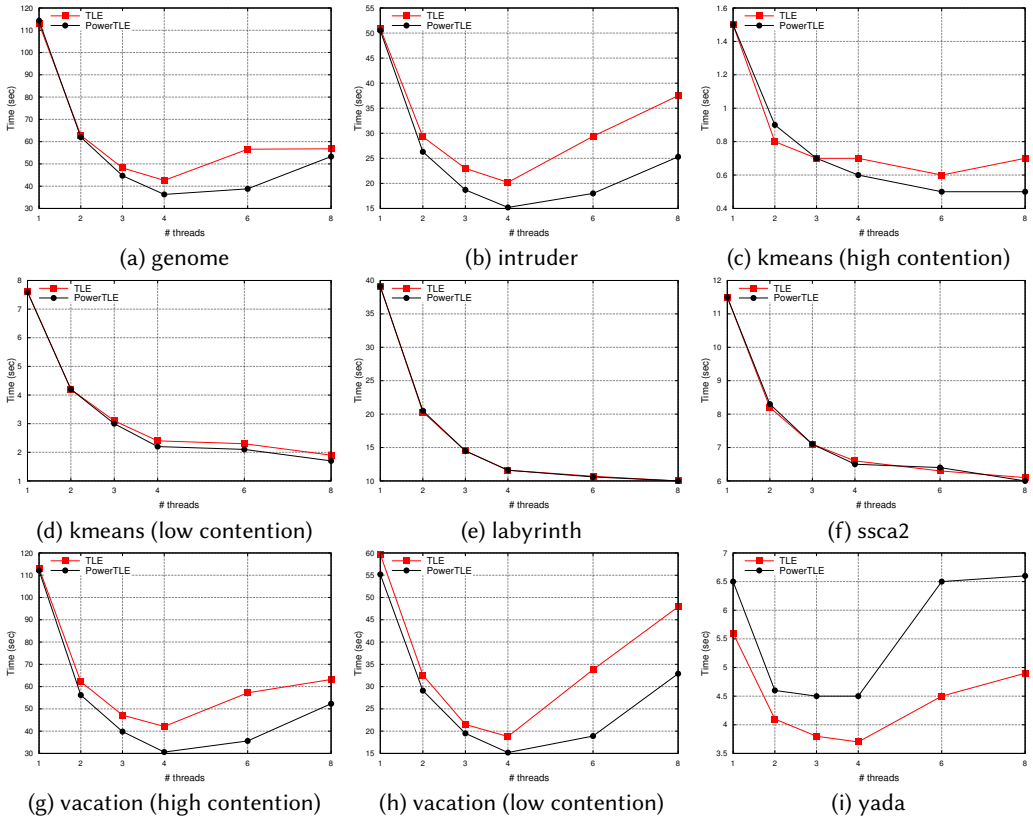


Fig. 8. STAMP running time measurements. Lower is better.

The results in Figure 8 show that power mode can be very helpful in certain cases, while it is harmful in one particular case. Specifically, in five cases (genome, intruder, kmeans-high, vacation-high and vacation-low), PowerTLE beats TLE by substantial margin, while it harms the performance of yada. Notably, in all three cases where PowerTLE performs on par with or only slightly improves over TLE (kmeans-low, labyrinth and ssca2), the TLE variant exhibits scalability up to 8 threads, thus limiting the benefits of power mode.

The breakdown of execution modes for critical sections of various STAMP benchmarks is presented in Figure 9, and sheds some light on the performance of PowerTLE compared to TLE. First, just like in the case of microbenchmarks reported in Sections 4.2 and 4.3, power mode appears to be helpful when substantial amount of transactions fail to lock (in TLE) and these transactions manage to commit using power mode. This happens in all five cases where PowerTLE beats TLE.

Second, in two of the three cases where PowerTLE and TLE perform almost the same (kmeans-low and ssca2), the vast majority of critical sections execute using regular transactions only. In fact, the only place where PowerTLE improves slightly over TLE in kmeans-low is when a small fraction of critical sections fail to the lock (in TLE) or revert to power mode (in PowerTLE) as the number of threads grows. Along with that, the case of labyrinth shows a different picture (cf. Figure 9(e)). Despite almost half of critical sections being executed using locks (in TLE), only a small portion of

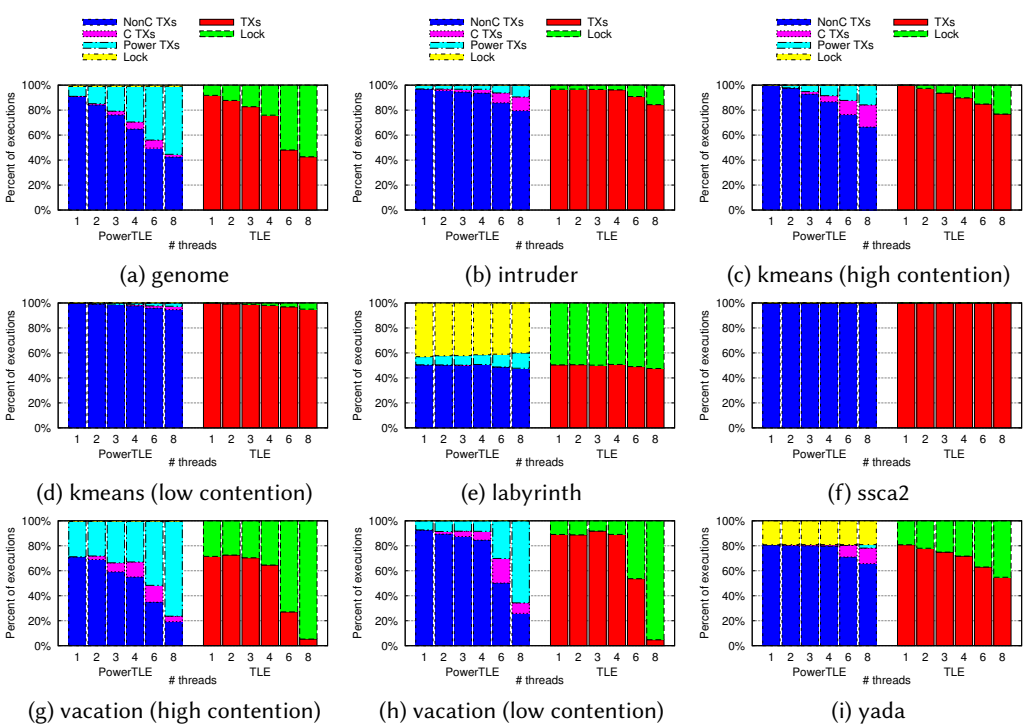


Fig. 9. Breakdown of execution modes for critical sections in STAMP benchmarks.

them is executed using power mode transactions (in PowerTLE), suggesting that the majority of those transactions fail due to capacity reasons. These results suggest that most of the time in this particular benchmark is spent outside of critical sections, explaining why despite the overhead of failed power mode transactions, PowerTLE achieves essentially the same results as TLE for this benchmark.

Finally, while yada shows a similar pattern to labyrinth (i.e., executions fail to commit using power mode transactions and therefore switch to lock), its running time is more sensitive to the performance of its critical sections. Here, the cost of failed power mode transactions is detrimental to the performance of PowerTLE. This benchmark shows that power transactions, like any kind of speculative execution, are effective only when speculation is mostly successful. We note, though, that a relatively straightforward optimization in PowerTLE that might eliminate performance degradation in yada, is to avoid using the power mode if a regular transaction fails due to capacity. This optimization should be used with care, as at times transactions that fail due to capacity do manage to commit if retried [5]. Exploring the impact of this optimization is in our future work.

5 SIMULATOR-BASED EVALUATION

In addition to the framework discussed above, we added support for power mode into SuperTrans [30], a transactional memory simulator built on top of SESC [32]. As reported in [29], SuperTrans was enhanced with a best-effort HTM support similar to Intel TSX.

Benchmark	Baseline	MoreReadsWins	ResponderWins	PowerTLE
genome	1.000	0.912	0.859	0.928
intruder	1.000	0.847	1.180	1.100
kmeans (low)	1.000	0.999	0.999	1.003
kmeans (high)	1.000	0.527	0.527	0.667
labyrinth	1.000	0.950	1.126	0.782
ssca2	1.000	1.002	1.003	1.014
vacation (low)	1.000	0.994	1.011	1.047
vacation (high)	1.000	0.948	0.971	1.035
yada	1.000	0.952	0.978	0.521
mean	1.000	0.903	0.962	0.900

Table 1. Relative performance of STAMP (lower is better).

Benchmark	Baseline	MoreReadsWins	ResponderWins	PowerTLE
genome	1.3	1.1	0.9	0.0
intruder	15.2	9.3	20.3	0.1
kmeans (low)	0.0	0.0	0.0	0.0
kmeans (high)	5.6	0.0	0.0	0.0
labyrinth	25.4	26.3	26.8	0.0
ssca2	0.0	0.0	0.0	0.0
vacation (low)	0.0	0.0	0.0	0.0
vacation (high)	0.4	0.1	0.3	0.0
yada	28.1	17.1	16.8	0.0

Table 2. Percent of transactions that fall to the lock.

In our evaluation, we use the default configuration file provided with the simulator, with minor configuration modifications for more realistic cache structure⁴. Specifically, we model a CMP machine with 64 cores connected through a 8 by 8 mesh network. Each core has private 8-way associative 64KB L1 instruction and data caches, a private 16-way associative 256KB L2 cache and a shared L3 cache with 8MB capacity. The L1 caches have hit latency of 3 cycles, the L2 caches have hit latency of 18 cycles, and the L3 cache has hit latency of 34 cycles.

We simulate a system with 16 threads running STAMP [24] with recommended inputs for a simulator environment. Note that these input sets are different from the native ones used in Section 4.4, as they are intended to produce shorter workloads that can be simulated in a reasonable time. Thus, some benchmarks might exhibit different contention patterns.

We modify the existing TLE implementation to use power transactions following the pseudocode in Figure 3. In line with the previous section, we refer to this modified implementation as PowerTLE. We compare PowerTLE to TLE running on top of the baseline (requester-wins, best-effort) HTM as well as on top of HTM modified according to the PleaseTM proposal [29]. For the latter, we use two variations called ResponderWins and MoreReadsWins. In the former, the requester running a hardware transaction and receiving a line with the plea bit set, aborts its transaction. In the latter, each core tracks the number of cache lines read transactionally and includes this counter along with the plea bit. The requester compares this counter to its own, and aborts its transaction if it read less lines than the responder. Both these variations are discussed in detail in [29] and implemented in SuperTrans by their authors.

Table 1 summarizes the relative performance of all variants compared to the baseline HTM implementation. That is, for each variant, we divide the simulated running time of the benchmark (as reported in the "Time=" output line by each benchmark) to the one measured with the baseline. The simulator results show that PowerTLE outperforms TLE, and performs on par with (but marginally better, on average, than) both PleaseTM variants. In general, the average gains of PowerTLE over

⁴ We verified that our configuration modifications did not have any impact on the results.

834 TLE are more modest compared to those measured with our emulation framework, in part due
835 to the different workload settings. Yet, PowerTLE is able to significantly (by 22–48%) improve
836 performance of 3 benchmarks, while not harming the performance of any other benchmark by
837 more than 10%.

838 Table 2 provides details on the number of transactions that end up falling to the lock for each of
839 the variants. Two observations can be drawn from these data. First, the percentage of transactions
840 falling to the lock in the baseline TLE is different, across all STAMP benchmarks, from the percentage
841 presented in Figure 9 for the emulation-based evaluation. This suggests that the contention patterns
842 are indeed different, and explain the difference in the overall performance results. Second, PowerTLE
843 eliminates virtually all failures to the lock. This property is important for several benchmarks, such
844 as yada and labyrinth, helping PowerTLE there to achieve better parallelism between hardware
845 transactions that leads to impressive gains over TLE. We note, though, that the lack of failures to
846 the lock does not translate to performance advantage for all benchmarks, as same transactions that
847 fall to the lock in TLE might be unable to make progress in PowerTLE due to conflicts with a power
848 mode transaction.
849

850 6 CONCLUSION

851 HTM is a promising tool to ease the development and accelerate the performance of concurrent
852 code. Most existing HTM implementations rely on requester-wins cache coherence protocols and
853 provide best-effort guarantees to concurrent transactions. The first property means that concurrent
854 transactions abort frequently when data conflicts are common, as demonstrated by multitude
855 of previous work [9, 13, 37]. The second property means that in order to guarantee progress,
856 concurrent programs must include a non-speculative fallback path [18]. This path is typically
857 implemented using a lock [12]; once a thread switches to this path, all other transactions have to
858 wait even if they do not conflict with the holder of the lock.

859 In this paper, we introduced special power transactions with the aim of alleviating these issues.
860 These transactions, running in a so-called power mode, receive priority in conflict resolution with
861 other, regular transactions. We show that supporting power transactions requires very simple,
862 almost trivial changes to existing best-effort requester-wins HTM implementations. Our extensive
863 experimental evidence using micro- and STAMP benchmarks, collected with emulation on top
864 of a real HTM implementation as well as with a transactional memory simulator, demonstrates
865 that power mode can improve parallelism between hardware transactions, leading to significant
866 benefits for HTM that supports power transactions over the one that does not.
867

868 REFERENCES

- 869
- 870 [1] C. S. Ananian, K. Asanovic, B. C. Kuszmaul, C. E. Leiserson, and S. Lie. Unbounded transactional memory. In
871 *Proceedings of the International Symposium on High-Performance Computer Architecture (HPCA)*, pages 316–327, 2005.
 - 872 [2] Adrià Armejach, Ruben Titos-Gil, Anurag Negi, Osman S. Unsal, and Adrián Cristal. Techniques to improve per-
873 formance in requester-wins hardware transactional memory. *ACM Trans. Archit. Code Optim.*, 10(4):42:1–42:25,
874 2013.
 - 875 [3] Colin Blundell, Joe Devietti, E. Christopher Lewis, and Milo M. K. Martin. Making the fast case common and the
876 uncommon case simple in unbounded transactional memory. In *Proceedings of the International Symposium on*
877 *Computer Architecture (ISCA)*, pages 24–34, 2007.
 - 878 [4] Jayaram Bobba, Kevin E. Moore, Haris Volos, Luke Yen, Mark D. Hill, Michael M. Swift, and David A. Wood. Performance
879 pathologies in hardware transactional memory. In *Proceedings of the International Symposium on Computer Architecture*
880 *(ISCA)*, pages 81–91, 2007.
 - 881 [5] Trevor Brown, Alex Kogan, Yossi Lev, and Victor Luchangco. Investigating the performance of hardware transactions
882 on a multi-socket machine. In *Proceedings of the ACM Symposium on Parallelism in Algorithms and Architectures*
(SPAA), pages 121–132, 2016.

- 883 [6] Harold W. Cain, Maged M. Michael, Brad Frey, Cathy May, Derek Williams, and Hung Le. Robust architectural support
884 for transactional memory in the power architecture. In *Proceedings of the International Symposium on Computer*
885 *Architecture (ISCA)*, pages 225–236, 2013.
- 886 [7] Luke Dalessandro, François Carouge, Sean White, Yossi Lev, Mark Moir, Michael L. Scott, and Michael F. Spear. Hybrid
887 NOrec: a case study in the effectiveness of best effort hardware transactional memory. In *Proceedings of the 16th*
888 *International Conference on Architectural Support for Programming Languages and Operating Systems ASPLOS*, pages
39–52, 2011.
- 889 [8] Peter Damron, Alexandra Fedorova, Yossi Lev, Victor Luchangco, Mark Moir, and Daniel Nussbaum. Hybrid transac-
890 tional memory. In *Proceedings of the International Conference on Architectural Support for Programming Languages and*
891 *Operating Systems (ASPLOS)*, pages 336–346, 2006.
- 892 [9] Dave Dice, Alex Kogan, and Yossi Lev. Refined transactional lock elision. In *Proceedings of the Symposium on Principles*
893 *and Practice of Parallel Programming (PPOPP)*, 2016.
- 894 [10] Dave Dice, Alex Kogan, Yossi Lev, Timothy Merrifield, and Mark Moir. Adaptive integration of hardware and software
895 lock elision techniques. In *ACM Symposium on Parallelism in Algorithms and Architectures (SPAA)*, pages 188–197,
896 2014.
- 897 [11] Dave Dice, Yossi Lev, Mark Moir, Daniel Nussbaum, and Marek Olszewski. Early experience with a commercial
898 hardware transactional memory implementation. Technical report, Sun Labs, 2009.
- 899 [12] David Dice, Yossi Lev, Mark Moir, and Daniel Nussbaum. Early experience with a commercial hardware transactional
900 memory implementation. In *Proceedings of the International Conference on Architectural Support for Programming*
901 *Languages and Operating Systems (ASPLOS)*, pages 157–168, 2009.
- 902 [13] Nuno Diegues, Paolo Romano, and Luís Rodrigues. Virtues and limitations of commodity hardware transactional
903 memory. In *Proceedings of the 23rd International Conference on Parallel Architectures and Compilation (PACT)*, pages
904 3–14, 2014.
- 905 [14] Ariel Eizenberg, Shiliang Hu, Gilles Pokam, and Joseph Devietti. Remix: Online detection and repair of cache contention
906 for the JVM. In *Proceedings of the ACM Conference on Programming Language Design and Implementation (PLDI)*, pages
907 251–265, 2016.
- 908 [15] GNU. Transactional memory in GCC. Retrieved from <https://gcc.gnu.org/wiki/TransactionalMemory> on 3 August
909 2015.
- 910 [16] Tim Harris and Keir Fraser. Language support for lightweight transactions. In *Proceedings of the Conference on*
911 *Object-oriented Programing, Systems, Languages, and Applications (OOPSLA)*, pages 388–402, 2003.
- 912 [17] Maurice Herlihy and J. Eliot B. Moss. Transactional memory: architectural support for lock-free data structures. In
913 *Proceedings of the 20th Annual International Symposium on Computer Architecture (ISCA'93)*, pages 289–300. ACM
914 Press, 1993.
- 915 [18] Intel Corporation. Transactional Synchronization in Haswell. Retrieved from [http://software.intel.com/en-us/blogs/](http://software.intel.com/en-us/blogs/2012/02/07/transactional-synchronization-in-haswell/)
916 [2012/02/07/transactional-synchronization-in-haswell/](http://software.intel.com/en-us/blogs/2012/02/07/transactional-synchronization-in-haswell/), 8 September 2012.
- 917 [19] Christian Jacobi, Timothy Slegel, and Dan Greiner. Transactional memory architecture and implementation for IBM
918 System Z. In *Proceedings of the IEEE/ACM International Symposium on Microarchitecture (MICRO)*, pages 25–36, 2012.
- 919 [20] Andi Kleen. Scaling existing lock-based applications with lock elision. *ACM Queue*, 12(1), 2014.
- 920 [21] Yujie Liu, Tingzhe Zhou, and Michael F. Spear. Transactional acceleration of concurrent data structures. In *Proceedings*
921 *of Symposium on Parallelism in Algorithms and Architectures (SPAA)*, pages 244–253, 2015.
- 922 [22] M. Lupon, G. Magklis, and A. Gonzalez. Fastm: A log-based hardware transactional memory with fast abort recovery.
923 In *Proceedings of the International Conference on Parallel Architectures and Compilation Techniques (PACT)*, pages
293–302, 2009.
- 924 [23] Alexander Matveev and Nir Shavit. Reduced hardware norec: A safe and scalable hybrid transactional memory. In
925 *Proceedings of the International Conference on Architectural Support for Programming Languages and Operating Systems*
926 *(ASPLOS)*, pages 59–71, 2015.
- 927 [24] Chi Cao Minh, JaeWoong Chung, Christos Kozyrakis, and Kunle Olukotun. STAMP: stanford transactional applications
928 for multi-processing. In *Proceedings of the International Symposium on Workload Characterization (IISWC)*, pages 35–46,
929 2008.
- 930 [25] Kevin E. Moore, Jayaram Bobba, Michelle J. Moravan, and Mark D. Hill and David A. Wood. Logtm: Log-based
931 transactional memory. In *Proceedings of the IEEE Symposium on High-Performance Computer Architecture (HPCA)*,
pages 258–269, 2006.
- [26] Michelle J. Moravan, Jayaram Bobba, Kevin E. Moore, Luke Yen, Mark D. Hill, Ben Liblit, Michael M. Swift, and
David A. Wood. Supporting nested transactional memory in logTM. In *Proceedings of the International Conference on*
Architectural Support for Programming Languages and Operating Systems (ASPLOS), pages 359–370, 2006.

- [27] Takuya Nakaike, Rei Odaira, Matthew Gaudet, Maged Michael, and Hisanobu Tomari. Quantitative Comparison of Hardware Transactional Memory for Blue Gene/Q, zEnterprise EC12, Intel Core, and POWER8. In *Proceedings of the International Symposium on Computer Architecture (ISCA)*, 2015.
- [28] Yang Ni, Adam Welc, Ali-Reza Adl-Tabatabai, Moshe Bach, Sion Berkowits, James Cownie, Robert Geva, Sergey Kozhukow, Ravi Narayanaswamy, Jeffrey Olivier, Serguei Preis, Bratin Saha, Ady Tal, and Xinmin Tian. Design and implementation of transactional constructs for C/C++. In *Proceedings of the ACM SIGPLAN Conference on Object-Oriented Programming, Systems, Languages, and Applications (OOPSLA)*, pages 195–212, 2008.
- [29] S. Park, M. Prvulovic, and C. J. Hughes. PleaseTM: Enabling transaction conflict management in requester-wins hardware transactional memory. In *Proceedings of the IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pages 285–296, 2016.
- [30] J. Poe, C. B. Cho, and T. Li. Using analytical models to efficiently explore hardware transactional memory and multi-core co-design. In *Proceedings of the International Symposium on Computer Architecture and High Performance Computing*, pages 159–166, 2008.
- [31] Ravi Rajwar, Maurice Herlihy, and Konrad Lai. Virtualizing transactional memory. In *Proceedings of the International Symposium on Computer Architecture (ISCA)*, pages 494–505, 2005.
- [32] Jose Renau, Basilio Fraguera, James Tuck, Wei Liu, Milos Prvulovic, Luis Ceze, Smruti Sarangi, Paul Sack, Karin Strauss, and Pablo Montesinos. SESC simulator, January 2005. <http://sesc.sourceforge.net>.
- [33] Amy Wang, Matthew Gaudet, Peng Wu, José Nelson Amaral, Martin Ohmacht, Christopher Barton, Raul Silvera, and Maged Michael. Evaluation of Blue Gene/Q hardware support for transactional memories. In *Proceedings of the International Conference on Parallel Architectures and Compilation Techniques (PACT)*, pages 127–136, 2012.
- [34] Thomas Wang. Integer hash function, 2007. Retrieved from <http://web.archive.org/web/20071223173210/http://www.concentric.net/~Ttwang/tech/inthash.htm>, 3 August 2015.
- [35] Roel Wieringa. *Design Methods for Reactive Systems: Yourdan, StateMATE, and the UML*. Morgan Kaufmann Publishers, Boston, 2003.
- [36] L. Yen, J. Bobba, M. R. Marty, K. E. Moore, H. Volos, M. D. Hill, M. M. Swift, and D. A. Wood. LogTM-SE: Decoupling hardware transactional memory from caches. In *Proceedings of the International Symposium on High Performance Computer Architecture (HPCA)*, pages 261–272, 2007.
- [37] Richard M. Yoo, Christopher J. Hughes, Konrad Lai, and Ravi Rajwar. Performance evaluation of Intel® transactional synchronization extensions for high-performance computing. In *International Conference for High Performance Computing, Networking, Storage and Analysis (SC)*, 2013.

A ADDITIONAL DETAILS ON EMULATION-BASED EVALUATION

Figures 10 and 11 provide additional implementation details of our emulation of the power mode, including the definition of metalocks and other auxiliary data structures (Figures 10), and the pseudo-code of read and write instrumentation barriers (Figure 11). Although we target the Intel Haswell architecture, the evaluation framework design is architecture-independent, and can be used with other HTM systems.

The mapping between an address (or more precisely, a cache line) and its corresponding metalock uses a fast pseudo-uniform hash function described in [34]. In our framework, we used very large arrays of 4M words representing metalocks (Lines 27 and 28) to reduce the chance that two cache lines will be mapped to the same metalock. Moreover, large arrays and a pseudo-uniform hash function mean that the chance that two cache lines accessed in the same transaction are mapped into adjacent metalock words is negligible, mitigating the chance for false sharing on the accessed metadata. As a result, it was not necessary to pad metalock words to avoid false sharing. Other fields in the State structure (cf. Figures 10), however, are properly padded (not shown for clarity).

Entering power mode is protected by a simple test-test-set lock (similar to the one shown in Figure 3) augmented with a sequence number (cf. Line 24 and 25). The latter is incremented after every lock acquisition (that is, right after a transaction enters the power mode) and before lock release (that is, right before a power mode transaction commits). The sequence number serves the purpose of efficient release of all acquired metalocks. Specifically, an execution that uses a regular transaction stores the current sequence number in a thread-local variable (`localSeqNumber` in the `ThreadInfo` structure, Line 37) *before* starting on HTM (and thus any change to this number by a

```

981 1 #define NUM_META_LOCKS (4 * 1024 * 1024)
982 2 #define CACHE_LINE_SIZE (64)
983
984 4 #define READ_CAPACITY (256)
985 5 #define WRITE_CAPACITY (64)
986
987 7 // fast pseudo-uniform hash function that maps a given key
988 8 // into a number between 0 and mask
989 9 uint64_t fast_hash(uintptr_t key, uint64_t mask) { ... }
990
991 11 // These macros translate from an address to a
992 12 // read/write meta lock protecting the cache line
993 13 // where the address belongs to.
994 14 #define ADDR_TO_READ_LOCK(addr) \
995 15     (&rMetadata[fast_hash( \
996 16         addr&~(CACHE_LINE_SIZE-1), \
997 17         NUM_META_LOCKS-1)])
998 18 #define ADDR_TO_WRITE_LOCK(addr) \
999 19     (&wMetadata[fast_hash( \
1000 20         addr&~(CACHE_LINE_SIZE-1), \
1001 21         NUM_META_LOCKS-1)])
1002
1003 23 struct State {
1004 24     uint64_t powerFlag;
1005 25     uint64_t seqNumber;
1006 26     bool isLocked;
1007 27     uint64_t rMetadata[NUM_META_LOCKS];
1008 28     uint64_t wMetadata[NUM_META_LOCKS];
1009 29     uint32_t uniqRCacheLines;
1010 30     uint32_t uniqWCacheLines;
1011 31     uint64_t lastPowerModeStartTime;
1012 32     ...
1013 33 } g_State;
1014
1015 35 struct ThreadInfo {
1016 36     bool myPowerFlag;
1017 37     uint64_t localSeqNumber;
1018 38     ...
1019 39 }

```

Fig. 10. Implementation details for power mode support

```

41 T read_barrier(void *addr) {
42     ThreadInfo *tx = getThreadInfo ();
43     if (!tx->myPowerFlag) {
44         uint64_t seqNumber = tx->localSeqNumber;
45         if (*ADDR_TO_WRITE_LOCK(addr) >= seqNumber)
46             htm_abort();
47     } else {
48         if (*ADDR_TO_READ_LOCK(addr) < g_State.seqNumber) {
49             *ADDR_TO_READ_LOCK(addr) = g_State.seqNumber;
50             membarstoreload();
51             if (!g_State.isLocked &&
52                 ++g_State.uniqRCacheLines > READ_CAPACITY) {
53                 // switch to the lock-based execution ;
54                 // this aborts all regular transactions ,
55                 // and forces them to wait for the lock
56                 // to become available again
57                 g_State.isLocked = true;
58                 membarstoreload();
59                 // calculate how much time we wasted
60                 // so far on this power mode transaction
61                 uint64_t timer = read_hw_clock();
62                 uint64_t delta = timer -
63                     g_State.lastPowerModeStartTime;
64                 // charge this amount of time for the
65                 // re-execution under lock
66                 while (read_hw_clock() - timer < delta);
67             }
68         }
69     }
70     return *addr;
71 }
72
73 void write_barrier(void *addr, T val) {
74     ThreadInfo *tx = getThreadInfo ();
75     if (!tx->myPowerFlag) {
76         uint64_t seqNumber = tx->localSeqNumber;
77         if (*ADDR_TO_READ_LOCK(addr) >= seqNumber)
78             htm_abort();
79         if (*ADDR_TO_WRITE_LOCK(addr) >= seqNumber)
80             htm_abort();
81     } else {
82         int seqNumber = *ADDR_TO_WRITE_LOCK(addr);
83         if (seqNumber < g_State.seqNumber) {
84             *ADDR_TO_WRITE_LOCK(addr) = g_State.seqNumber;
85             if (!g_State.isLocked &&
86                 ++g_State.uniqWCacheLines > WRITE_CAPACITY) {
87                 /* same as Lines 53-66 */
88             }
89         }
90     }
91     *addr = val;
92 }

```

Fig. 11. Instrumentation barriers used to implement the libitm interface of GCC

power mode transaction does not abort running regular transactions). Regular transactions use this number to check whether the metalock is “locked” by a power mode transaction (see Lines 44 and 76). Thus, once the sequence number is incremented at the end of the power transaction, any

1030 regular transaction starting and reading this number afterwards can deduce that all metalocks have
1031 been released.

1032 The power transaction (whose `myPowerFlag` is set) stores the current sequence number into
1033 the corresponding metalock word (Lines 49 and 84). We use an if-statement (Lines 48 and 83) to
1034 check whether the store is actually required to avoid writing the same value when the same cache
1035 line is accessed multiple times by a power transaction. (This if-statement also helps to keep track
1036 of the number of unique cache lines accesses for read and for write; the concrete use of these
1037 numbers is described later.) This optimization is more important for the read barrier, which requires
1038 a store-load memory fence (Line 50) to ensure that the metalock update becomes visible to regular
1039 transactions before the power transaction performs its read; otherwise, a power transaction may
1040 read inconsistent data. We note that in TSO architectures, such as Intel Haswell, the store-load
1041 memory fence is not required in the write barrier due to the total order on memory writes.

1042 Notice that in the read instrumentation barrier, a regular transaction accesses a write metalock
1043 only (Line 45), while in the write barrier, it accesses both read and write metalocks (Lines 77 and 79).
1044 Thus, a regular transaction is able to share cache lines accessed by a power transaction for read,
1045 but it cannot acquire ownership of (i.e., write to) cache lines accessed by a power transaction for
1046 read or for write, as required.

1047 The `uniqRCacheLines` and `uniqWCacheLines` fields of the `State` structure are used to keep
1048 track of the number of unique cache lines accessed by a power transaction for read and for write,
1049 respectively. As described above, we use these numbers to emulate capacity aborts by power
1050 transactions and re-execution under lock. Based on data in [27], the read capacity of HTM in
1051 Intel Core i7-4770 machine (which is the machine we used for our evaluation) is several tens of
1052 thousands of cache lines, while the write capacity is a few hundreds of cache lines. Factors like
1053 cache associativity and hyper threading limit the effective capacity of hardware transactions. In
1054 fact, our experiments show that in some cases, transactions experience capacity aborts when they
1055 access only a few hundreds of cache lines for read and even less than that for write. As a result, we
1056 chose very conservative capacity limits for our evaluation (cf. Lines 4 and 5).

1057 Taking the read barrier as an example, once the number of unique read cache lines goes beyond a
1058 threshold (Line 52), we switch to the lock-based execution by turning the `iSLocked` flag on (Line 57).
1059 After that, we calculate how much time has passed since we started the power transaction, and spin
1060 for that amount of time, charging the lock-based execution for running the prefix of the (effectively,
1061 aborted) power transaction (Lines 61–66). Note that once the power transaction transitions into
1062 the lock-based execution, other, regular transactions are aborted and wait for the lock to become
1063 available again. Thus, the lock-based execution in a real system would take the same path as the
1064 power-mode transaction, as it would access the same memory locations and read same values. As a
1065 result, the emulation of time cost required to abort a power transaction and re-execute it under
1066 lock is realistic. Note that once the execution of the power transaction continues under lock, it goes
1067 through same barriers, to keep the cost of memory access comparable across all execution modes.

1068
1069
1070
1071
1072
1073
1074
1075
1076
1077
1078