

# T

## 1 Thread-Level Speculation

2 JOSEP TORRELLAS

3 University of Illinois at Urbana-Champaign 4231 Siebel

4 Center, M/C-258, Urbana, IL, USA

### 5 Synonyms

6 Speculative multithreading (SM); Speculative  
7 parallelization; Speculative run-time parallelization;  
8 Speculative threading; Speculative thread-level paral-  
9 lelization; Thread-level data speculation (TLDS); TLS

### 10 Definition

11 Thread-Level Speculation (TLS) refers to an environ-  
12 ment where execution threads operate speculatively,  
13 performing potentially unsafe operations, and tem-  
14 porarily buffering the state they generate in a buffer or  
15 cache. At a certain point, the operations of a thread are  
16 declared to be correct or incorrect. If they are correct,  
17 the thread commits, merging the state it generated with  
18 the correct state of the program; if they are incorrect,  
19 the thread is squashed and typically restarted from its  
20 beginning. The term TLS is most often associated to  
21 a scenario where the purpose is to execute a sequen-  
22 tial application in parallel. In this case, the compiler or  
23 the hardware breaks down the application into specu-  
24 lative threads that execute in parallel. However, strictly  
25 speaking, TLS can be applied to any environment where  
26 threads are executed speculatively and can be squashed  
27 and restarted.

### 28 Discussion

#### 29 Basic Concepts in Thread-Level

#### 30 Speculation

31 In its most common use, Thread-Level Speculation  
32 (TLS) consists of extracting units of work (i.e., tasks)  
33 from a sequential application and executing them on  
34 different threads in parallel, hoping not to violate

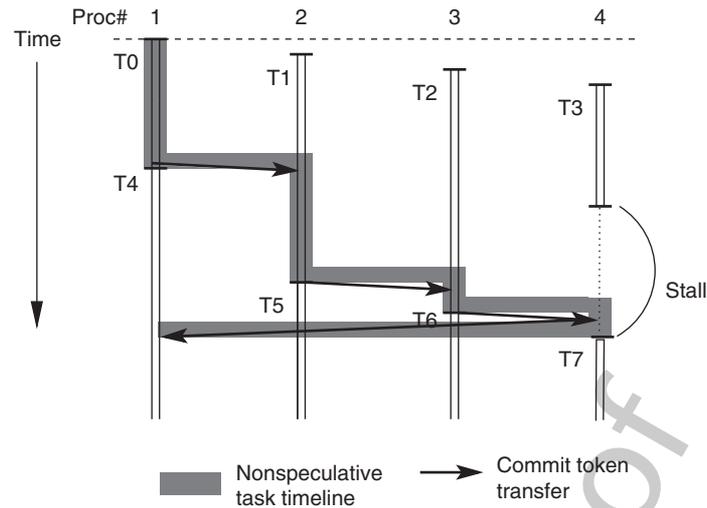
35 sequential semantics. The control flow in the sequen-  
36 tial code imposes a relative ordering between the tasks,  
37 which is expressed in terms of predecessor and suc-  
38 cessor tasks. The sequential code also induces a data  
39 dependence relation on the memory accesses issued by  
40 the different tasks that parallel execution cannot violate.

41 A task is *Speculative* when it may perform or may  
42 have performed operations that violate data or con-  
43 trol dependences with its predecessor tasks. Other-  
44 wise, the task is non-speculative. The memory accesses  
45 issued by speculative tasks are called speculative mem-  
46 ory accesses.

47 When a non-speculative task finishes execution, it is  
48 ready to *Commit*. The role of commit is to inform the  
49 rest of the system that the data generated by the task  
50 is now part of the safe, non-speculative program state.  
51 Among other operations, committing always involves  
52 passing the *Commit Token* to the immediate succes-  
53 sor task. This is because maintaining correct sequential  
54 semantics in the parallel execution requires that tasks  
55 commit in order from predecessor to successor. If a task  
56 reaches its end and is still speculative, it cannot com-  
57 mit until it acquires non-speculative status and all its  
58 predecessors have committed.

59 Figure 1 shows an example of several tasks run-  
60 ning on four processors. In this example, when task T3  
61 executing on processor 4 finishes the execution, it can-  
62 not commit until its predecessor tasks T0, T1, and T2  
63 also finish and commit. In the meantime, depending on  
64 the hardware support, processor 4 may have to stall or  
65 may be able to start executing speculative task T7. The  
66 example also shows how the non-speculative task status  
67 changes as tasks finish and commit, and the passing of  
68 the commit token.

69 Memory accesses issued by a speculative task  
70 must be handled carefully. Stores generate *Speculative*  
71 *Versions* of data that cannot simply be merged with  
72 the non-speculative state of the program. The reason is  
73 that they may be incorrect. Consequently, these versions



**Thread-Level Speculation. Fig. 1** A set of tasks executing on four processors. The figure shows the nonspeculative task timeline and the transfer of the commit token

74 are stored in a *Speculative Buffer* local to the processor  
 75 running the task – e.g., the first-level cache. Only when  
 76 the task becomes nonspeculative are its versions safe.

77 Loads issued by a speculative task try to find the  
 78 requested datum in the local speculative buffer. If they  
 79 miss, they fetch the correct version from the memory  
 80 subsystem, i.e., the closest predecessor version from the  
 81 speculative buffers of other tasks. If no such version  
 82 exists, they fetch the datum from memory.

83 As tasks execute in parallel, the system must iden-  
 84 tify any violations of cross-task data dependences.  
 85 Typically, this is done with special hardware or soft-  
 86 ware support that tracks, for each individual task, the  
 87 data that the task wrote and the data that the task read  
 88 without first writing it. A data-dependence violation is  
 89 flagged when a task modifies a datum that has been read  
 90 earlier by a successor task. At this point, the consumer  
 91 task is *squashed* and all the data versions that it has  
 92 produced are discarded. Then, the task is re-executed.

93 Figure 2 shows an example of a data-dependence  
 94 violation. In the example, each iteration of a loop  
 95 is a task. Each iteration issues two accesses to an  
 96 array, through an un-analyzable subscripted subscript.  
 97 At run-time, iteration J writes A[5] after its succes-  
 98 sor iteration J+2 reads A[5]. This is a Read After  
 99 Write (RAW) dependence that gets violated due to  
 100 the parallel execution. Consequently, iteration J+2 is  
 101 squashed and restarted. Ordinarily, all the successor

tasks of iteration J+2 are also squashed at this time 102  
 because they may have consumed versions generated 103  
 by the squashed task. While it is possible to selectively 104  
 squash only tasks that used incorrect data, it would 105  
 involve extra complexity. Finally, as iteration J+2 re- 106  
 executes, it will re-read A[5]. However, at this time, the 107  
 value read will be the version generated by iteration J. 108

Note that WAR and WAW dependence violations do 109  
 not need to induce task squashes. The successor task has 110  
 prematurely written the datum, but the datum remains 111  
 buffered in its speculative buffer. A subsequent read 112  
 from a predecessor task (in a WAR violation) will get a 113  
 correct version, while a subsequent write from a prede- 114  
 cessor task (in a WAW violation) will generate a version 115  
 that will be merged with main memory before the one 116  
 from the successor task. 117

However, many proposed TLS schemes, to reduce 118  
 hardware complexity, induce squashes in a variety of sit- 119  
 uations. For instance, if the system has no support to 120  
 keep different versions of the same datum in different 121  
 speculative buffers in the machine, cross-task WAR and 122  
 WAW dependence violations induce squashes. More- 123  
 over, if the system only tracks accesses on a per-line 124  
 basis, it cannot disambiguate accesses to different words 125  
 in the same memory line. In this case, false sharing of a 126  
 cache line by two different processors can appear as a 127  
 data-dependence violation and also trigger a squash. 128



**Thread-Level Speculation. Fig. 2** Example of a data-dependence violation

129 Finally, while TLS can be applied to various code  
 130 structures, it is most often applied to loops. In this  
 131 case, tasks are typically formed by a set of consecutive  
 132 iterations.

133 The rest of this article is organized as follows:  
 134 First, the article briefly classifies TLS schemes. Then, it  
 135 describes the two major problems that any TLS scheme  
 136 has to solve, namely, buffering and managing specu-  
 137 lative state, and detecting and handling dependence  
 138 violations. Next, it describes the initial efforts in TLS,  
 139 other uses of TLS, and machines that use TLS.

**140 Classification of Thread-Level Speculation**  
**141 Schemes**

142 There have been many proposals of TLS schemes. They  
 143 can be broadly classified depending on the emphasis  
 144 on hardware versus software, and the type of target  
 145 machine.

146 The majority of the proposed schemes use hardware  
 147 support to detect cross-task dependence violations that  
 148 result in task squashes (e.g., [1, 4, 6, 8, 11, 12, 14, 16, 18, 20,  
 149 23, 27, 28, 31, 32, 36]). Typically, this is attained by using  
 150 the hardware cache coherence protocol, which sends  
 151 coherence messages between the caches when multi-  
 152 ple processors access the same memory line. Among all  
 153 these hardware-based schemes, the majority rely on a  
 154 compiler or a software layer to identify and prepare the  
 155 tasks that should be executed in parallel. Consequently,  
 156 there have been several proposals for TLS compilers  
 157 (e.g., [9, 19, 33, 34]). Very few schemes rely on the  
 158 hardware to identify the tasks (e.g., [1]).

159 Several schemes, especially in the early stages of TLS  
 160 research, proposed software-only approaches to TLS  
 161 (e.g., [7, 13, 25, 26]). In this case, the compiler typically  
 162 generates code that causes each task to keep shadow  
 163 locations and, after the parallel execution, checks if mul-  
 164 tiple tasks have updated a common location. If they  
 165 have, the original state is restored.

Most proposed TLS schemes target small shared- 166  
 memory machines of about two to eight processors 167  
 (e.g., [14, 18, 27, 29]). It is in this range of paral- 168  
 lelism that TLS is most cost effective. Some TLS pro- 169  
 posals have focused on smaller machines and have 170  
 extended a superscalar core with some hardware units 171  
 that execute threads speculatively [1, 20]. Finally, some 172  
 TLS proposals have targeted scalable multiprocessors 173  
 [4, 23, 28]. This is a more challenging environment, 174  
 given the longer communication latencies involved. It 175  
 requires applications that have significant parallelism 176  
 that cannot be analyzed statically by the compiler. 177

**178 Buffering and Managing Speculative State**

179 The state produced by speculative tasks is unsafe, since  
 180 such tasks may be squashed. Therefore, any TLS scheme  
 181 must be able to identify such state and, when neces-  
 182 sary, separate it from the rest of the memory state.  
 183 For this, TLS systems use structures, such as caches  
 184 [4, 6, 12, 18, 28], and special buffers [8, 14, 23, 32], or  
 185 undo logs [7, 11, 36]. This section outlines the chal-  
 186 lenges in buffering and managing speculative state. A  
 187 more detailed analysis and a taxonomy is presented by  
 188 Garzaran et al. [10].

**189 Multiple Versions of the Same Variable**  
**190 in the System**

191 Every time that a task writes for the first time to a  
 192 variable, a new version of the variable appears in the  
 193 system. Thus, two speculative tasks running on different  
 194 processors may create two different versions of the same  
 195 variable [4, 12]. These versions need to be buffered sep-  
 196 arately, and special actions may need to be taken so that  
 197 a reader task can find the correct version out of the sev-  
 198 eral coexisting in the system. Such a version will be the  
 199 version created by the producer task that is the closest  
 200 predecessor of the reader task.

A task has at most a single version of any given 201  
 variable, even if it writes to the variable multiple times. 202

203 The reason is that, on a dependence violation, the whole  
 204 task is undone. Therefore, there is no need to keep  
 205 intermediate values of the variable.

## 206 Multiple Speculative Tasks per Processor

207 When a processor finishes executing a task, the task  
 208 may still be speculative. If the TLS buffering support is  
 209 such that the processor can only hold state from a single  
 210 speculative task, the processor stalls until the task com-  
 211 mits. However, to better tolerate task load imbalance,  
 212 the local buffer may have been designed to buffer state  
 213 from several speculative tasks, enabling the processor to  
 214 execute another speculative task. In this case, the state  
 215 of each task must be tagged with the ID of the task.

## 216 Multiple Versions of the Same Variable 217 in a Single Processor

218 When a processor buffers state from multiple specu-  
 219 lative tasks, it is possible that two such tasks create  
 220 two versions of the same variable. This occurs in load-  
 221 imbalanced applications that exhibit private data pat-  
 222 terns (i.e., WAW dependences between tasks). In this  
 223 case, the buffer will have to hold multiple versions of  
 224 the same variable. Each version will be tagged with a  
 225 different task ID. This support introduces complication  
 226 to the buffer or cache. Indeed, on an external request,  
 227 extra comparisons will need to be done if the cache has  
 228 two versions of the same variable.

## 229 Merging of Task State

230 The state produced by speculative tasks is typically  
 231 merged with main memory at task commit time; how-  
 232 ever, it can instead be merged as it is being generated.  
 233 The first approach is called *Architectural Main Memory*  
 234 (*AMM*) or *Lazy Version Management*; the second one  
 235 is called *Future Main Memory (FMM)* or *Eager Version*  
 236 *Management*. These schemes differ on whether the main  
 237 memory contains only safe data (*AMM*) or it can also  
 238 contain speculative data (*FMM*).

239 In *AMM* systems, all speculative versions remain  
 240 in caches or buffers that are kept separate from the  
 241 coherent memory state. Only when a task becomes non-  
 242 speculative can its buffered state be merged with main  
 243 memory. In a straightforward implementation, when  
 244 a task commits, all the buffered dirty cache lines are  
 245 merged with main memory, either by writing back the

lines to memory [4] or by requesting ownership for  
 them to obtain coherence with main memory [28].

In *FMM* systems, versions from speculative tasks are  
 merged with the coherent memory when they are gen-  
 erated. However, to enable recovery from task squashes,  
 when a task generates a speculative version of a variable,  
 the previous version of the variable is saved in a log.  
 Note that, in both approaches, the coherent memory  
 state can temporarily reside in caches, which function  
 in their traditional role of extensions of main memory.

## Detecting and Handling Dependence Violations

### Basic Concepts

The second aspect of TLS involves detecting and han-  
 dling dependence violations. Most TLS proposals focus  
 on data dependences, rather than control dependences.  
 To detect (cross-task) data-dependence violations, most  
 TLS schemes use the same approach. Specifically, when  
 a speculative task writes a datum, the hardware sets a  
 Speculative Write bit associated with the datum in the  
 cache; when a speculative task reads a datum before it  
 writes to it (an event called *Exposed Read*), the hard-  
 ware sets an Exposed Read bit. Depending on the TLS  
 scheme supported, these accesses also cause a tag asso-  
 ciated with the datum to be set to the ID of the task.

In addition, when a task writes a datum, the cache  
 coherence protocol transaction that sends invalidations  
 to other caches checks these bits. If a successor task has  
 its Exposed Read bit set for the datum, the successor  
 task has prematurely read the datum (i.e., this is a RAW  
 dependence violation), and is squashed [18].

If the Speculative Write and Exposed Read bits are  
 kept on a per-word basis, only dependences on the same  
 word can cause squashes. However, keeping and main-  
 taining such bits on a per-word basis in caches, network  
 messages, and perhaps directory modules is costly in  
 hardware. Moreover, it does not come naturally to the  
 coherence protocol of multiprocessors, which operate  
 at the granularity of memory lines.

Keeping these bits on a per-line basis is cheaper and  
 compatible with mainstream cache coherence proto-  
 cols. However, the hardware cannot then disambiguate  
 accesses at word level. Furthermore, it cannot combine  
 different versions of a line that have been updated in dif-  
 ferent words. Consequently, cross-task RAW and WAW

291 violations, on both the same word and different words  
 292 of a line (i.e., false sharing), cause squashes.

293 Task squash is a very costly operation. The cost  
 294 is threefold: overhead of the squash operation itself,  
 295 loss of whatever correct work has already been per-  
 296 formed by the offending task and its successors, and  
 297 cache misses in the offending task and its successors  
 298 needed to reload state when restarting. The latter over-  
 299 head appears because, as part of the squash opera-  
 300 tion, the speculative state in the cache is invalidated.  
 301 Figure 3a shows an example of a RAW violation across  
 302 tasks  $i$  and  $i+j+1$ . The consumer task and its successors  
 303 are squashed.

**304 Techniques to Avoid Squashes**

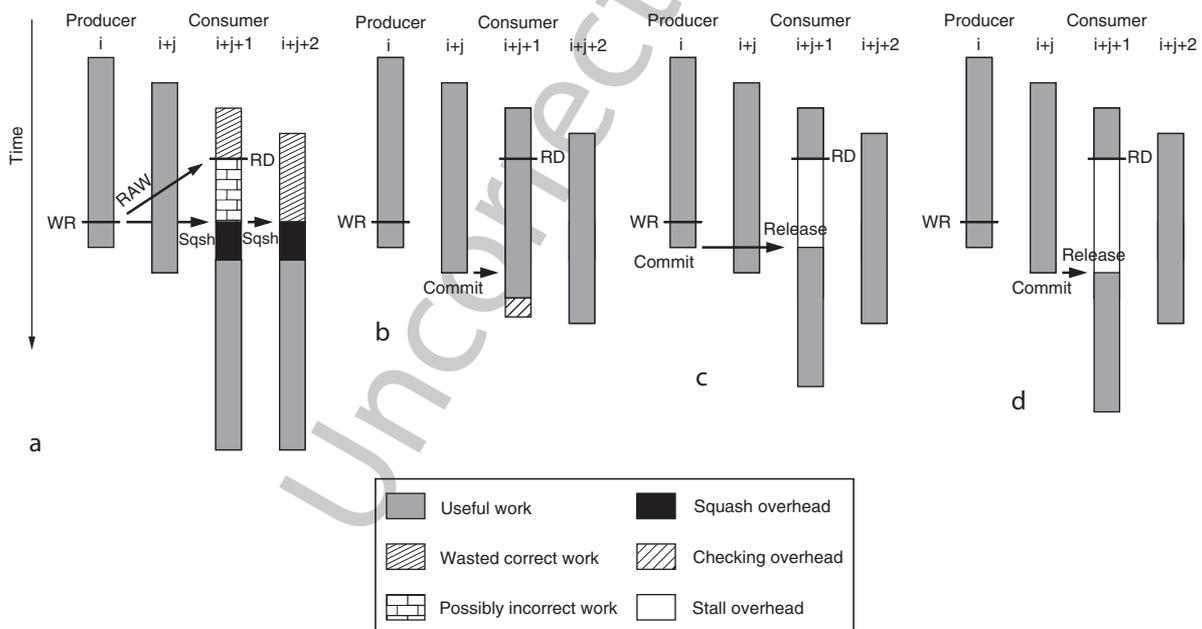
305 Since squashes are so expensive, there are techniques to  
 306 avoid them. If the compiler can conclude that a certain  
 307 pair of accesses will frequently cause a data-dependence  
 308 violation, it can statically insert a synchronization opera-  
 309 tion that forces the correct task ordering at runtime.

310 Alternatively, the machine can have hardware sup-  
 311 port that records, at runtime, where dependence viola-  
 312 tions occur. Such hardware may record the program  
 313 counter of the read or writes involved, or the address  
 314 of the memory location being accessed. Based on this

315 information, when these program counters are reached  
 316 or the memory location is accessed, the hardware can  
 317 try one of several techniques to avoid the violation. This  
 318 section outlines some of the techniques that can be used.  
 319 A more complete description of the choices is presented  
 320 by Cintra and Torrellas [5]. Without loss of generality, a  
 321 RAW violation is assumed.

322 Based on past history, the predictor may predict  
 323 that the pair of conflicting accesses are engaged in false  
 324 sharing. In this case, it can simply allow the read to pro-  
 325 ceed and then the subsequent write to execute silently,  
 326 without sending invalidations. Later, before the con-  
 327 sumer task is allowed to commit, it is necessary to  
 328 check whether the sections of the line read by the con-  
 329 sumer overlap with the sections of the line written by  
 330 the producer. This can be easily done if the caches  
 331 have per-word access bits. If there is no overlap, it was  
 332 false sharing and the squash is avoided. Figure 3b shows  
 333 the resulting time line.

334 When there is a true data dependence between tasks,  
 335 a squash can be avoided with effective use of value pre-  
 336 diction. Specifically, the predictor can predict the value  
 337 that the producer will produce, speculatively provide it  
 338 to the consumer’s read, and let the consumer proceed.



**Thread-Level Speculation. Fig. 3** RAW data-dependence violation that results in a squash (a) or that does not cause a squash due to false sharing or value prediction (b), or consumer stall (c and d)

339 Again, before the consumer is allowed to commit, it is  
 340 necessary to check that the value provided was correct.  
 341 The timeline is also shown in Fig. 3b.

342 In cases where the predictor is unable to predict the  
 343 value, it can avoid the squash by stalling the consumer  
 344 task at the time of the read. This case can use two pos-  
 345 sible approaches. An aggressive approach is to release  
 346 the consumer task and let it read the current value as  
 347 soon as the predicted producer task commits. The time  
 348 line is shown in Fig. 3c. In this case, if an intervening  
 349 task between the first producer and the consumer later  
 350 writes the line, the consumer will be squashed. A more  
 351 conservative approach is not to release the consumer  
 352 task until it becomes nonspeculative. In this case, the  
 353 presence of multiple predecessor writers will not squash  
 354 the consumer. The time line is shown in Fig. 3d.

### 355 Initial Efforts in Thread-Level Speculation

356 An early proposal for hardware support for a form of  
 357 speculative parallelization was made by Knight [16] in  
 358 the context of functional languages. Later, the Multi-  
 359 scalar processor [27] was the first proposal to use a form  
 360 of TLS within a single-chip multithreaded architec-  
 361 ture. A software-only form of TLS was proposed in the  
 362 LRPD test [25]. Early proposals of hardware-based TLS  
 363 include the work of several authors [14, 17, 21, 29, 35].

### 364 Other Uses of Thread-Level Speculation

365 TLS concepts have been used in environments that  
 366 have goals other than trying to parallelize sequen-  
 367 tial programs. For example, they have been used to  
 368 speed up explicitly parallel programs through Spec-  
 369 ulative Synchronization [22], or for parallel program  
 370 debugging [24] or program monitoring [37]. Similar  
 371 concepts to TLS have been used in systems supporting  
 372 hardware transactional memory [15] and continuous  
 373 atomic-block operation [30].

### 374 Machines that Use Thread-Level 375 Speculation

376 Several machines built by computer manufacturers have  
 377 hardware support for some form of TLS – although  
 378 the specific implementation details are typically not dis-  
 379 closed. Such machines include systems designed for  
 380 Java applications such as Sun Microsystems’ MAJC  
 381 chip [31] and Azul Systems’ Vega processor [2].  
 382 The most high-profile system with hardware support

for speculative threads is Sun Microsystems’ ROCK 383  
 processor [3]. Other manufacturers are rumored to be 384  
 developing prototypes with similar hardware. 385

### Related Entries 386

- ▶Instruction-Level Speculation 387
- ▶Speculative Synchronization 388
- ▶Transactional Memory 389

### Bibliography 390

1. Akkary H, Driscoll M (1998) A dynamic multithreading proces- 391  
 sor. In: International symposium on microarchitecture, Dallas, 392  
 November 1998 393
2. Azul Systems. Vega 3 Processor. [http://www.azulsystems.com/](http://www.azulsystems.com/products/vega/processor) 394  
[products/vega/processor](http://www.azulsystems.com/products/vega/processor) 395
3. Chaudhry S, Cypher R, Ekman M, Karlsson M, Landin A, Yip S, 396  
 Zeffer H, Tremblay M (2009) Simultaneous speculative threading: 397  
 a novel pipeline architecture implemented in Sun’s ROCK Pro- 398  
 cessor. In: International symposium on computer architecture, 399  
 Austin, June 2009 400
4. Cintra M, Martínez JF, Torrellas J (2000) Architectural support 401  
 for scalable speculative parallelization in shared-memory multi- 402  
 processors. In: International symposium on computer architec- 403  
 ture, Vancouver, June 2000, pp 13–24 404
5. Cintra M, Torrellas J (2002) Eliminating squashes through 405  
 learning cross-thread violations in speculative parallelization for 406  
 multiprocessors. In: Proceedings of the 8th High-Performance 407  
 computer architecture conference, Boston, Feb 2002 408
6. Figueiredo R, Fortes J (2001) Hardware support for extract- 409  
 ing coarse-grain speculative parallelism in distributed shared- 410  
 memory multiprocessors. In: Proceedings of the international 411  
 conference on parallel processing, Valencia, Spain, September 412  
 2001 413
7. Frank M, Lee W, Amarasinghe S (2001) A software framework 414  
 for supporting general purpose applications on raw computation 415  
 fabrics. Technical report, MIT/LCS Technical Memo MIT-LCS- 416  
 TM-619, July 2001 417
8. Franklin M, Sohi G (1996) ARB: a hardware mechanism for 418  
 dynamic reordering of memory references. *IEEE Trans Comput* 419  
 45(5):552–571 420
9. Garcia C, Madriles C, Sanchez J, Marcuello P, Gonzalez A, 421  
 Tullsen D (2005) Mitosis compiler: An infrastructure for specu- 422  
 lative threading based on pre-computation slices. In: Conference 423  
 on programming language design and implementation, Chicago, 424  
 Illinois, June 2005 425
10. Garzarán M, Prvulovic M, Llabería J, Viñals V, Rauchwerger L, 426  
 Torrellas J (2005) Tradeoffs in buffering speculative memory 427  
 state for thread-level speculation in multiprocessors. *ACM Trans* 428  
*Archit Code Optim* 429
11. Garzarán MJ, Prvulovic M, Llabería JM, Viñals V, Rauchwerger L, 430  
 Torrellas J (2003) Using software logging to support multi- 431  
 version buffering in thread-level speculation. In: International 432

- 433 conference on parallel architectures and compilation techniques,  
434 New Orleans, Sept 2003
- 435 12. Gopal S, Vijaykumar T, Smith J, Sohi G (1998) Speculative ver-  
436 sioning cache. In: International symposium on high-performance  
437 computer architecture, Las Vegas, Feb 1998
- 438 13. Gupta M, Nim R (1998) Techniques for speculative run-time par-  
439 allelization of loops. In: Proceedings of supercomputing 1998,  
440 ACM Press, Melbourne, Australia, Nov 1998
- 441 14. Hammond L, Willey M, Olukotun K (1998) Data speculation sup-  
442 port for a chip multiprocessor. In: International conference on  
443 architectural support for programming languages and operating  
444 systems, San Jose, California, Oct 1998, pp 58–69
- 445 15. Herlihy M, Moss E (1993) Transactional memory: architectural  
446 support for lock-free data structures. In: International sympo-  
447 sium on computer architecture, IEEE Computer Society Press,  
448 San Diego, May 1993
- 449 16. Knight T (1986) An architecture for mostly functional languages.  
450 In: ACM lisp and functional programming conference, ACM  
451 Press, New York, Aug 1986, pp 500–519
- 452 17. Krishnan V, Torrellas J (1998) Hardware and software sup-  
453 port for speculative execution of sequential binaries on a chip-  
454 multiprocessor. In: International conference on supercomputing,  
455 Melbourne, Australia, July 1998
- 456 18. Krishnan V, Torrellas J (1999) A chip-multiprocessor archi-  
457 tecture with speculative multithreading. *IEEE Trans Comput*  
458 48(9):866–880
- 459 19. Liu W, Tuck J, Ceze L, Ahn W, Strauss K, Renau J, Torrellas J  
460 (2006) POSH: A TLS compiler that exploits program structure.  
461 In: International symposium on principles and practice of parallel  
462 programming, San Diego, Mar 2006
- 463 20. Marcuello P, Gonzalez A (1999) Clustered speculative multi-  
464 threaded processors. In: International conference on supercom-  
465 puting, Rhodes, Island, June 1999, pp 365–372
- 466 21. Marcuello P, Gonzalez A, Tubella J (1998) Speculative multi-  
467 threaded processors. In: International conference on supercom-  
468 puting, ACM, Melbourne, Australia, July 1998
- 469 22. Martinez J, Torrellas J (2002) Speculative synchronization: apply-  
470 ing thread-level speculation to explicitly parallel applications. In:  
471 International conference on architectural support for program-  
472 ming languages and operating systems, San Jose, Oct 2002
- 473 23. Prvulovic M, Garzaran MJ, Rauchwerger L, Torrellas J (2001)  
474 Removing architectural bottlenecks to the scalability of specu-  
475 lative parallelization. In: Proceedings of the 28th international  
476 symposium on computer architecture (ISCA'01), New York, June  
477 2001, pp 204–215
- 478 24. Prvulovic M, Torrellas J (2003) ReEnact: using thread-level spec-  
479 ulation to debug data races in multithreaded codes. In: Inter-  
480 national symposium on computer architecture, San Diego, June  
481 2003
- 482 25. Rauchwerger L, Padua D (1995) The LRPD test: speculative run-  
483 time parallelization of loops with privatization and reduction  
484 parallelization. In: Conference on programming language design  
485 and implementation, La Jolla, California, June 1995
26. Rundberg P, Stenstrom P (2000) Low-cost thread-level data 486  
dependence speculation on multiprocessors. In: Fourth work- 487  
shop on multithreaded execution, architecture and compilation, 488  
Monterrey, Dec 2000 489
27. Sohi G, Breach S, Vijaykumar T (1995) Multiscalar processors. In: 490  
International Symposium on computer architecture, ACM Press, 491  
New York, June 1995 492
28. Steffan G, Colohan C, Zhai A, Mowry T (2000) A scalable 493  
approach to thread-level speculation. In: Proceedings of the 27th 494  
Annual International symposium on computer architecture, Van- 495  
couver, June 2000, pp 1–12 496
29. Steffan G, Mowry TC (1998) The potential for using thread- 497  
level data speculation to facilitate automatic parallelization. In: 498  
International symposium on high-performance computer archi- 499  
tecture, Las Vegas, Feb 1998 500
30. Torrellas J, Ceze L, Tuck J, Cascaval C, Montesinos P, Ahn W, 501  
Prvulovic M (2009) The bulk multicore architecture for improved 502  
programmability. *Communications of the ACM*, New York 503
31. Tremblay M (1999) MAJC: microprocessor architecture for java 504  
computing. *Hot Chips*, Palo Alto, Aug 1999 505
32. Tsai J, Huang J, Amlo C, Lilja D, Yew P (1999) The superthreaded 506  
processor architecture. *IEEE Trans Comput* 48(9):881–902 507
33. Vijaykumar T, Sohi G (1998) Task selection for a multiscalar pro- 508  
cessor. In: International symposium on microarchitecture, Dallas, 509  
Nov 1998, pp 81–92 510
34. Zhai A, Colohan C, Steffan G, Mowry T (2002) Compiler opti- 511  
mization of scalar value communication between speculative 512  
threads. In: International conference on architectural support for 513  
programming languages and operating systems, San Jose, Oct 514  
2002 515
35. Zhang Y, Rauchwerger L, Torrellas J (1998) Hardware for specula- 516  
tive run-time parallelization in distributed shared-memory multi- 517  
processors. In: Proceedings of the 4th International symposium 518  
on high-performance computer architecture (HPCA), Phoenix, 519  
Feb 1998, pp 162–174 520
36. Zhang Y, Rauchwerger L, Torrellas J (1999) Hardware for specu- 521  
lative parallelization of partially-parallel loops in DSM multi- 522  
processors. In: Proceedings of the 5th international symposium 523  
on high-performance computer architecture, Orlando, Jan 1999, 524  
pp 135–139 525
37. Zhou P, Qin F, Liu W, Zhou Y, Torrellas (2004) iWatcher: efficient 526  
architectural support for software debugging. In: International 527  
symposium on computer architecture, IEEE Computer society, 528  
München, June 2004 529