

Reducing DRAM Row Activations with Eager Read/Write Clustering

MYEONGJAE JEON, CONGLONG LI, ALAN L. COX, and SCOTT RIXNER, Rice University

This article describes and evaluates a new approach to optimizing DRAM performance and energy consumption that is based on eagerly writing dirty cache lines to DRAM. Under this approach, many dirty cache lines are written to DRAM before they are evicted. In particular, dirty cache lines that have not been recently accessed are eagerly written to DRAM when the corresponding row has been activated by an ordinary, noneager access, such as a read. This approach enables clustering of reads and writes that target the same row, resulting in a significant reduction in row activations. Specifically, for a variety of applications, it reduces the number of DRAM row activations by an average of 42% and a maximum of 82%. Moreover, the results from a full-system simulator show compelling performance improvements and energy consumption reductions. Out of 23 applications, 6 have overall performance improvements between 10% and 20%, and 3 have improvements in excess of 20%. Furthermore, 12 consume between 10% and 20% less DRAM energy, and 7 have energy consumption reductions in excess of 20%.

Categories and Subject Descriptors: B.3.2 [Memory Structures]: Design Styles; C.4 [Performance of Systems]: Design Studies

General Terms: Design, Measurement, Performance

Additional Key Words and Phrases: DRAM, performance, energy consumption, eager writeback

ACM Reference Format:

Jeon, M., Li, C., Cox, A. L., and Rixner, S. 2013. Reducing DRAM row activations with eager read/write clustering. *ACM Trans. Architect. Code Optim.* 10, 4, Article 43 (December 2013), 25 pages.
DOI: <http://dx.doi.org/10.1145/2555289.2555300>

1. INTRODUCTION

Due to the structure of modern DRAM devices, the performance and energy consumption of the main memory system is primarily determined by the locality of access within these devices. In particular, data within a device's memory cells cannot be accessed directly. Data can only be read from or written to a *row buffer*. Migrating data between the memory cells and this row buffer takes 25ns and consumes 3.9nJ in a modern DDR3 device [Micron Technology Inc. 2010]. In contrast, reading data from this row buffer takes 17.5ns and 1.44nJ and can be pipelined with other read accesses to this same row for improved performance.

Unfortunately, writes to the main memory system generally exhibit poor locality. Table II provides a breakdown of DRAM accesses for a variety of programs. Frequently, the fraction of row activations caused by writes to DRAM is greater than the fraction of accesses that are writes. In other words, writes cause a disproportionate number of row activations. The reason is that writes to DRAM are not tied to the program's memory

New Article, Not an Extension of a Conference Paper.

This work is supported by the National Science Foundation under grant CCF-1018840.

Authors' addresses: M. Jeon (corresponding author), C. Li, A. L. Cox, and S. Rixner, Department of Computer Science, Rice University, Houston, TX; email: mjjeon@rice.edu.

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

© 2013 ACM 1544-3566/2013/12-ART43 \$15.00

DOI: <http://dx.doi.org/10.1145/2555289.2555300>

access stream. They are, instead, due to the cache line evictions that are caused by the program's actual accesses. Consequently, cache line writebacks are a primary cause of performance degradation and energy consumption in modern memory systems. This article proposes a novel approach to mitigating the impact of the cache line writebacks using Eager Read/Write Clustering (ERWC).

Eager cache line writeback was originally proposed by Lee et al. [2000] to reduce the competition for bus bandwidth between cache line reads and writebacks. In particular, they performed eager writebacks when the bus was idle so that a cache line read would be less likely to be delayed by a cache line writeback.

In contrast, this article shows that eager cache line writeback can instead be used to cluster reads *and* writes that target the same row, thereby significantly reducing the number of row activations. Specifically, under ERWC, eager writebacks are triggered upon a row activation, regardless of the cause. Furthermore, eager writebacks are canceled when the rows are no longer active and prevented from occurring repeatedly for the same cache line.

In order to assess the benefits of ERWC, we run all of the programs listed in Table II on a full-system simulator. These programs include Glimpse, MineBench, SPEC CFP2006, SPEC CINT2006, and TPC-C. In a nutshell, for these programs, the geometric mean of the performance improvement and the energy savings are 7.0% and 10.6%, respectively.

In addition, this article presents a design space evaluation of different clustering and eager writeback policies for cache lines. In particular, this evaluation compares the performance of the various points in this design space, including the Virtual Write Queue (VWQ) [Stuecheli et al. 2010], DRAM-Aware Writeback [Lee et al. 2010], and ERWC. To the best of our knowledge, this article presents the most comprehensive exploration of the design space to date. In summary, ERWC achieves better overall results in terms of performance and energy consumption than the other points in this design space.

The rest of this article is organized as follows. The next section provides relevant background on the organization and operation of DRAM. Section 3 presents the motivation for our eager writeback scheme. Section 4 defines the design space of clustering and eager writeback policies, and Section 5 describes the mechanisms for realizing these policies. Section 6 describes our experimental methodology, and Section 7 presents our experimental results. Finally, Section 8 discusses related work, and Section 9 summarizes our conclusions.

2. DRAM BASICS

Modern DRAM devices are organized into banks, rows, and columns. A DRAM device has 4 to 16 banks that can operate concurrently. Memory cells within a bank cannot be directly accessed. Instead, data can only be read or written through the *row buffer*. Each bank has a dedicated row buffer that can hold one row at a time. A row is transferred from the memory cell array to the row buffer by a *row activate* command. After an arbitrary number of read and write commands, which transfer data to or from columns in the row buffer, the *precharge* command returns the row's contents to the memory cell array and prepares the bank for another row activation.

A rank ties together multiple DRAM devices to create a wider data interface. Specifically, the address lines are tied together to create the illusion of a single, wider DRAM. This leads to rows of 4 to 16 KB across the rank. A modern DIMM consists of 1 to 4 ranks of DRAM devices.

Modern DRAM timing is quite complex, and there are many constraints placed upon when DRAM commands (precharge, activate, read/write) can be issued. These

Table I. Micron DDR3-1600 Delay and Energy Costs

Operation	Latency		Energy nJ
	Cycles	ns	
Precharge	10	12.5	3.90
Activate	10	12.5	
Column read	14	17.5	1.44
Column write	12	15.0	1.44

constraints ensure that data can traverse the long internal wires of the DRAM device and that the power dissipation does not exceed the limits of the device.

The precharge and activate commands are more expensive, in terms of both energy and delay, than the read and write operations to a column within the row buffer. This is due to the fact that the column operations simply access a few bytes in a buffer, whereas the precharge and activate commands drive thousands of long wires across an entire DRAM bank.

Table I shows the timing and energy costs of the various DRAM operations for a Micron DDR3-1600 1Gb 8-bit-wide DRAM device [Micron Technology Inc. 2010]. Recall that before performing a column operation, the bank must first be precharged and then the row within that bank must be activated in order to transfer it to the row buffer. So, to read a single burst of data (64 bits in this example) would take $10 + 10 + 14 = 34$ memory cycles, or 42.5ns. Note that the column read access only occupies the data bus for the last 4 of the 14 cycles. Once a row has been activated, however, read operations can be pipelined, yielding 64 data bits every 4 memory cycles (due to the 8-bit-wide double data rate bus). As DRAM devices are aggregated into a rank, each read would occur simultaneously across eight devices yielding a 64-byte cache line for each read operation. This clearly shows that there is an enormous advantage to reading multiple cache lines per row to achieve the DRAM's peak bandwidth of one cache line every 4 cycles, instead of one cache line every 34 cycles.

Write operations are similarly expensive, with a potential 32 memory cycle delay. They too can be pipelined at a rate of one burst every 4 memory cycles. However, read and write operations cannot be freely interleaved. When transitioning from reading to writing, there is a 2-cycle delay in the Micron DRAM to allow the bus drivers to turn around. The situation is worse when transitioning from writing to reading, as the logic for dealing with the row buffer requires these operations to be isolated. This leads to a 20 memory cycle delay from the completion of a write to the completion of the read (6 empty memory cycles plus the full 14-cycle latency of a column read).

Although these latencies are quite high, they can be partially hidden in two ways. First, precharge operations can be overlapped with read operations—that is, the bank can be precharged as the last column read is proceeding. Second, while the row buffer from one bank is being read or written, precharge and activate operations can occur in parallel to other banks.

As Table I shows, the combination of a precharge and an activate command for a bank consumes 3.9nJ. Not only is this more than twice the energy of a read or write access, but these row operations can occur in parallel for multiple banks. If there were no limitations, all eight banks within a DRAM could be activated simultaneously, which would result in instantaneous power dissipation of 1.25W per DRAM device, without even transferring any data across the pins of the device. In a memory system with 8 DIMMs, each with 8 DRAM devices, that would consume 80W.

In summary, due to the latency, energy, and timing constraints of row activations, it is important to activate rows judiciously and perform as many read and write accesses per activated row as possible. In order to mitigate the cost of row activation, there have been several proposals to either reduce the row size of the DRAM [Rixner 2004; Zheng

et al. 2008; Ahn et al. 2009; Udipi et al. 2010] or improve the locality of access within a row [Sudan et al. 2010; Stuecheli et al. 2010].

3. ROW BUFFER UTILIZATION

Given the impact of row buffer utilization on performance and energy, it is important to understand how programs exercise the memory system. Across a variety of programs, there are typically only 1 to 4 column accesses per row activation. The reason for this low utilization is cache line writebacks. Writebacks interfere with the programs' read locality in the DRAM and lead to significant performance and energy problems.

This section explores this issue further and clarifies the problems with cache line writebacks. Specifically, they occur more often than might be expected, they often require row activations, and they often cause future read accesses to reactivate rows that would have already been in the row buffer if the writeback had not occurred. The data in Table II highlight these problems. The methodology for collecting this data will be described in Section 6. Briefly, we simulated these programs on a widely used memory architecture that is intended to reduce interference between reading and writing (see Section 5.4). This architecture has separate DRAM read and write queues managed using watermarks. Writes are processed either when the write queue reaches a high watermark or when there is no pending read.

3.1. Frequency of DRAM Writes

Whenever a cache line is first accessed, it must be read from the DRAM. Both read and write misses in the cache trigger DRAM reads. In systems that employ writeback caches, DRAM writes only occur when a dirty cache line is evicted. Intuitively, this should occur much less often than DRAM reads. The first column of Table II shows the number and fraction of DRAM accesses that are reads and writes for 35 programs. For 27 of these programs, writes account for more than 20% of total DRAM accesses. Moreover, writes account for more than 40% of the DRAM accesses for 7 of those 27 programs. Thus, DRAM writes are not that infrequent.

3.2. Frequency of Row Activations

Modern cache replacement policies are not concerned with the locality of data written back to the next level in the memory hierarchy. In particular, Last-Level Caches (LLCs) are not concerned with the locality of row accesses to the DRAM. Thus, DRAM writes caused by evictions frequently exhibit no locality of access with the contemporaneous reads. Moreover, they are often unrelated to other contemporaneous writes. The second column of Table II shows the number and fraction of row activations caused by reads and writes. In 21 of the programs, row activations caused by writes account for more than 30% of all row activations. In 2 programs, writes account for more than 50% of all row activations.

3.3. Effects of Row Activations

Beyond the direct costs of row activations caused by cache writebacks, there is the collateral damage that they cause. The third column in Table II shows the number of read/write operations that occur per activated row. As the table shows, for most programs, there are only 1 to 4 read/write operations per activated row, with an average of 2.71 reads and 2.36 writes. Worse, there are fewer than 2 writes per activated row for 20 of the programs. Surprisingly, the trends for rows that are activated by reads are not much better than for rows that are activated by writes. There are fewer than 2 reads per activated row for 15 of the programs. Only 4 of the programs perform 5 or more reads per row activation.

Table II. Decomposition of DRAM Accesses, Row Activations, and Row Buffer Hit Count by Access Type

	DRAM Accesses		Row Activations		Row Buffer Hit Count		Improv. (Over RD)		
	RD	WR	RD	WR	RD	WR			
SPEC	sjeng	213,101 (0.60)	143,681 (0.40)	200,253 (0.58)	143,380 (0.42)	1.06	1.00	1.10	4%
	cactusADM	2,665,486 (0.78)	761,785 (0.22)	2,534,296 (0.82)	548,548 (0.18)	1.07	1.29	1.06	-1%
	bwaves	18,894,562 (0.89)	2,349,116 (0.11)	17,137,080 (0.97)	572,345 (0.03)	1.10	4.10	1.11	1%
	hammer	1,559,931 (0.52)	1,454,331 (0.48)	1,509,650 (0.69)	664,357 (0.31)	1.03	2.19	1.09	6%
	sphinx3	6,616,400 (0.95)	341,166 (0.05)	4,578,341 (0.96)	195,371 (0.04)	1.44	1.76	1.59	10%
	leslie3d	6,145,586 (0.78)	1,773,320 (0.22)	4,386,249 (0.83)	929,077 (0.17)	1.40	1.90	1.49	26%
	bzip2	2,131,640 (0.66)	1,098,656 (0.34)	1,174,445 (0.58)	858,813 (0.42)	1.80	1.30	1.59	52%
	zeusmp	1,578,140 (0.85)	270,828 (0.15)	946,231 (0.83)	192,603 (0.17)	1.66	1.44	1.85	11%
	calculx	20,707 (0.67)	10,236 (0.33)	10,882 (0.57)	8,168 (0.43)	1.90	1.26	1.62	96%
	astar	3,581,883 (0.65)	1,892,374 (0.35)	1,944,374 (0.58)	1,435,564 (0.42)	1.84	1.32	1.62	99%
	soplex	9,116,719 (0.79)	2,423,530 (0.21)	5,338,474 (0.76)	1,664,363 (0.24)	1.71	1.46	1.65	9%
	games	5,443 (0.87)	826 (0.13)	2,881 (0.80)	718 (0.20)	1.85	1.32	1.74	9%
	gromacs	650,233 (0.77)	198,728 (0.23)	304,434 (0.66)	158,087 (0.34)	2.13	1.26	1.84	106%
	gc	100,166 (0.85)	18,182 (0.15)	47,820 (0.75)	15,997 (0.25)	2.09	1.16	1.85	27%
	gobmk	703,772 (0.67)	339,805 (0.33)	325,251 (0.59)	227,149 (0.41)	2.13	1.55	1.89	43%
	povray	2,559 (0.78)	741 (0.22)	1,192 (0.69)	538 (0.31)	2.10	1.48	1.91	35%
	mcf	6,263,711 (0.75)	2,132,879 (0.25)	2,468,129 (0.60)	1,666,259 (0.40)	2.54	1.28	2.03	93%
	perlbench	1,418,442 (0.77)	413,012 (0.23)	568,538 (0.66)	293,042 (0.34)	2.48	1.45	2.13	71%
	h264ref	688,751 (0.69)	307,593 (0.31)	317,106 (0.71)	132,476 (0.29)	2.17	2.33	2.22	214%
	dealIII	893,867 (0.79)	232,906 (0.21)	303,577 (0.65)	163,063 (0.35)	2.93	1.45	2.41	309%
tonto	93,696 (0.58)	66,896 (0.42)	30,561 (0.56)	24,045 (0.44)	3.03	2.82	2.94	141%	
lbm	14,725,660 (0.58)	10,792,440 (0.42)	4,470,425 (0.54)	3,774,681 (0.46)	3.29	2.86	3.09	1578%	
mile	9,031,398 (0.69)	4,035,548 (0.31)	2,924,744 (0.73)	1,091,916 (0.27)	3.09	3.70	3.25	130%	
namd	87,603 (0.80)	22,344 (0.20)	20,307 (0.67)	10,017 (0.33)	4.31	2.24	3.63	60%	
omnetpp	7,059,050 (0.97)	209,709 (0.03)	1,603,319 (0.93)	129,718 (0.07)	4.40	1.65	4.19	4%	
GemsFDTD	5,914,232 (0.51)	5,785,726 (0.49)	820,227 (0.50)	835,850 (0.50)	7.21	6.92	7.06	416%	
libquantum	8,056,051 (0.52)	7,482,754 (0.48)	790,445 (0.56)	615,673 (0.44)	10.19	12.15	11.05	700%	
Utility	3,506,522 (0.88)	464,075 (0.12)	2,525,198 (0.90)	286,686 (0.10)	1.39	1.63	1.41	1109%	
ScalParC	9,809,989 (0.62)	5,966,556 (0.38)	6,559,845 (0.63)	3,776,751 (0.37)	1.49	1.59	1.53	35%	
Aprori	13,064,315 (0.51)	12,601,545 (0.49)	4,861,974 (0.34)	9,239,756 (0.66)	2.61	1.40	1.82	84%	
TPC-C	4,536,738 (0.77)	1,327,571 (0.23)	1,975,441 (0.66)	1,001,465 (0.34)	2.28	1.36	1.97	48%	
Glimpse	431,824 (0.76)	138,874 (0.24)	183,793 (0.65)	97,175 (0.35)	2.33	1.46	2.03	28%	
HOP	4,703,767 (0.66)	2,411,530 (0.34)	2,571,572 (0.84)	485,077 (0.16)	1.83	4.97	2.33	1730%	
SEMPHY	8,034,838 (0.96)	325,643 (0.04)	1,519,018 (0.88)	198,458 (0.12)	5.28	1.68	4.87	110%	
Eclat	3,481,493 (0.75)	1,174,505 (0.25)	607,372 (0.66)	307,452 (0.34)	5.73	3.82	5.09	182%	

Fractions are parenthesized under the columns DRAM Accesses and Row Activations. No Write shows the row buffer hit count when all write operations were discarded. SPEC and Other programs appear in increasing order of average row buffer hit count.

Table III. Policy Space Breakdown of Eager Writeback

Policy	Options
When	Eviction/Activation
Access Type(s)	Write-only/Read-Write
Eager Writeback Cancellation	On/Off
Repeated Eager Writeback	On/Off
Cache Line Selection	LRU/LWP
Lookup Depth	1– <i>W</i> ways
Lookup Range	1– <i>S</i> sets
Idle Prefill	On/Off

There are two reasons for the low hit rates on activated rows. First, not all programs have good spatial locality in the row buffer. Our simulated architecture maps the physical address space across ranks and banks at a row granularity, so sequential memory accesses by a program should yield numerous DRAM reads per activated row. Nonetheless, this is not happening because the cache hierarchy acts as a filter for the program’s memory access stream. Second, even when programs do have good spatial locality, cache evictions often interfere with that locality, reducing the number of reads per row activation.

The fourth column in Table II explains the damage done by the interference of cache evictions. “No Write” shows the amount of row locality in the reference stream when all write operations were discarded. This is compared to the row locality of read operations (RD) in the third column to show the improvement when there is no such interference. Although showing at least a modest increase in all programs, the number of reads per activated row increases markedly for seven of the programs. Thus, eliminating the interference of cache evictions will have a positive effect for all programs and a significant impact on some programs.

4. POLICIES

Lee et al. [2000] initially proposed using eager writeback to reduce competition for bus bandwidth between cache line reads and writebacks. In particular, they performed eager writebacks when the bus was idle so that a cache line writeback would be less likely to delay a cache line read. This article shows that the concept of eagerly writing back dirty cache lines can also be used to reduce the number of row activations in the DRAM.

The basic idea of ERWC is to write back dirty cache lines from the LLC ahead of time, when their associated DRAM row has been activated by other ordinary, noneager accesses. This strategy enables better clustering of both reads and writes that target the same DRAM row, increasing the overall number of DRAM accesses per activated row, which should improve performance and reduce energy consumption in the memory system.

To preserve normal cache behavior, lines that are eagerly written back are not actually evicted from the cache. Thus, the cache hit rate remains largely unaffected (although timing variations due to the different scheduling of DRAM accesses could result in minor differences). If those cache lines are not subsequently modified again, then their later eviction will not require a write to the DRAM.

There are many possible variations of eager writeback that might be used to reduce row activations in the DRAM. This section discusses the policy space of these variations. To organize this discussion, we present a number of orthogonal policies that collectively define a variation. Table III summarizes these policies, and Table IV explains how previous variations, including Eager [Lee et al. 2000], DAW [Lee et al. 2010], VWQ

Table IV. Policies Employed by Eager Writeback Schemes

	Eager	DAW	VWQ	LWPG	ERWC
When	Eviction	Eviction	Both	Eviction	Activation
Access	None	WR only	WR only	WR only	RD & WR
Cancel	No	No	No	No	Yes
Repeat	Yes	Yes	No	Yes	No
Selection	LRU	LRU	LRU	LWP	LRU
Depth	1-way	Dynamic	W-way	N/A	2-way
Range	None	Row	4	Row	Row
Prefill	No	No	Yes	No	No

[Stuecheli et al. 2010], and Last-Write Predictor Guided (LWPG) [Wang et al. 2012], as well as our proposed ERWC scheme, can be defined in terms of these policies.

When/Access Type(s). Dirty cache lines in the LLC are eagerly written to DRAM when they can be scheduled with other accesses to the same row. A key dimension of the policy space is with *which* types of accesses are the eager writebacks clustered.

Eviction-triggered clustering has been proposed in previous studies to cluster DRAM writes that target the same row [Lee et al. 2010; Stuecheli et al. 2010; Wang et al. 2012]. Under this policy, when a dirty cache line is evicted, other dirty cache lines in the LLC that map to the same row as the evicted line are eagerly written back. Since the eager and eviction writes are placed in the same DRAM write queue and typically are scheduled as a burst, this policy increases the row-level locality of DRAM writes.

Alternatively, we propose that dirty cache lines are eagerly written when their associated row has been activated by any type of ordinary, noneager access, including reads. This *activation-triggered clustering* should allow greater clustering of accesses, both read and write, that target the same row. By increasing the clustering of reads and writes, there will be fewer row activations, and rows should be kept open for a longer period of time to perform more accesses per activation.

Both ERWC and VWQ perform activation-triggered clustering. The distinction between them is that activation-triggered clustering in VWQ only applies to activations caused by write bursts to the DRAM.

Eager Writeback Cancellation. If the memory controller activates a new row while eager writebacks remain for the old row, those eager writebacks would later trigger another row activation. In such circumstances, it may be better to *cancel* all eager writebacks to a row once that bank is precharged. That way, they will either be reissued speculatively if eager writeback with the original row is ever attempted again, or they will simply occur when the cache lines are finally evicted, as if no eager writeback had been attempted.

Repeated Eager Writeback. Since we are speculatively writing dirty cache lines from the LLC, a cache line may be modified again after it is eagerly written back. In that case, it might be eagerly written back multiple times before it is evicted, causing an increase in overall writeback traffic. After a cache line has been eagerly written back, further eager writebacks could be prohibited for that cache line, minimizing the excess writeback traffic. This choice can impact programs where data writes are more likely to hit in the cache, cache lines are less likely to be evicted, and DRAM writebacks are relatively infrequent.

Cache Line Selection. The goal of the cache line selection policy is to identify dirty cache lines that are unlikely to be written to again before they are evicted. Most of the eager writeback schemes, including ERWC, use the LRU information in the LLC for

cache line selection. Lee et al. [2000] analyzed the frequency of writes to dirty cache lines as a function of their LRU position in the cache set. Their results showed that dirty cache lines that reach the LRU position are rarely written to again. In contrast, cache lines in the MRU position are highly likely to be written to again.

More recently, Wang et al. [2012] have proposed the LWPG approach that does not depend on the use of LRU replacement by the cache. Their approach is based upon the following observation: if an instruction has previously been the source of the last write before eviction for one cache line, then that instruction is likely to be the source of last writes to other cache lines. Thus, they propose a Last-Write Predictor (LWP) that selects cache lines based upon the past behavior of the instruction at a given address.

Lookup Depth. Given a W -way, set-associative LLC, using LRU replacement, an eager writeback scheme can be configured to check a certain number, *depth*, of the LRU ways. In effect, the depth determines the aggressiveness of the scheme. Varying the depth results in an interesting trade-off—less aggressive configurations that consider fewer LRU ways result in fewer cache line writebacks, and more aggressive configurations that consider more LRU ways result in more writebacks. However, the latter may write back more cache lines that will be written to again.

Lookup Range. The eager writeback schemes evaluated in this article require that the cache be probed for dirty lines that match the given row address. If the row size is 8KB and the cache line size is 64 bytes, then 128 cache lines can map to a row and these lines will be placed in 128 different cache sets. Thus, to achieve the full benefits of the eager writeback schemes, 128 sets may have to be probed. Although this full search may increase the likelihood of finding dirty cache lines, it can also consume energy. Another question that our evaluation will address is whether or not a comparable reduction in the number of row activations can still be achieved while probing a smaller number of cache sets.

Idle Prefill. When no cache line transfer is in progress, the memory controller can preemptively ask the LLC for dirty cache lines belonging to ranks with no pending writes. Ideally, the LLC will provide dirty cache lines belonging to the same row. In any case, otherwise idle DRAM cycles will be utilized by these writes. This *idle prefill* policy can be combined with any eager writeback scheme. Its effectiveness depends on many factors, including how many and how often idle DRAM cycles occur, and how DRAM reacts to incoming reads during a burst of writes.

5. MECHANISMS

Eager writebacks may be initiated as soon as the memory controller begins to activate a row or detects an evicted cache line. These are *triggering events*, which have a *trigger address* that is the physical address accessed by the event.

The memory controller and LLC must be modified to (1) detect these triggering events, (2) transmit the trigger address from the memory controller to the LLC, if necessary, (3) find dirty cache lines to eagerly write back, and (4) properly manage the eager writebacks. Although the first two steps are obvious extensions to existing hardware, there are several design options for the final two steps, which will be discussed in this section.

5.1. Cache Line Lookup Mechanism

To support cache line lookup, there needs to be enough lookup bandwidth in the LLC. We believe that sufficient bandwidth exists because Lee et al. [2010] have shown that there is substantial idle time in the LLC. Specifically, in a single-core system, for 15

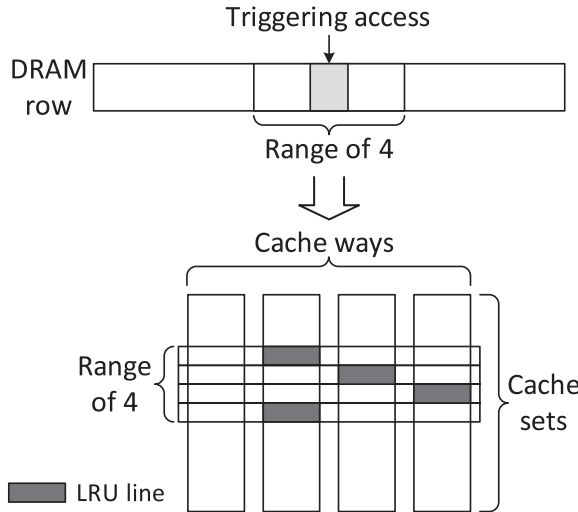


Fig. 1. An example of finding dirty cache lines in 4-set range and 1-way depth of lookup.

out of 16 SPEC2000/2006 programs, they found that the LLC is idle 95% of the time. Moreover, even in 4-core and 8-core systems, the idle time is around 90%.

An 8KB DRAM row is composed of 128 64B cache lines. Therefore, the LLC would potentially have to be probed 127 times to find all dirty cache lines for eager writeback after every triggering event. The lookup range decreases the number of sets that should be searched. To facilitate a more efficient range-based search, the cache sets mapping to a row in the DRAM are partitioned into equal range-sized groups, as shown in Figure 1. Only the partition corresponding to the trigger address is searched.

When the LRU information is used for cache line selection, the lookup depth limits the selection of dirty cache lines. Within each set, a dirty cache line will be selected for eager writeback only if its LRU position is within the lookup depth. Note that only one cache line per set can correspond to a particular DRAM row, so all of these conditions can easily be detected using the normal cache line eviction and LRU logic.

With a lookup range of k , the LLC controller must access the cache $k - 1$ times using cache line addresses neighboring the trigger address. Any matching dirty lines within the lookup depth can be forwarded to the memory controller back-to-back for eager writeback as they are found in the cache. Considering that a modern memory system has multiple channels and ranks, the LLC controller may need to interleave accesses for several active banks.

To reduce lookup overhead, Stuecheli et al. [2010] proposed the Set State Vector (SSV). The SSV consists of one bit per cache set that indicates whether the designated LRU ways contain any dirty cache lines. By first checking the SSV bit for a cache set before performing the cache line lookup, the LLC controller can avoid pointless lookups that cannot possibly yield a dirty cache line.

5.2. Eager Cancellation Support

If a speculative writeback is guaranteed to be performed by the memory controller, then the dirty bit in the cache can be cleared as soon as the eager writeback is initiated. However, there can be an advantage to allowing the memory controller to *cancel* speculative writebacks in order to activate other rows to perform pending read accesses. To support this optimization, both the memory and cache controllers need to be modified.

Table V. Processor and DRAM System Parameters

Processor	
Processor ISA	AMD Family 10h
Frequency	4.0GHz
Re-Order-Buffer	256 entries
Pipeline	out-of-order, 8-wide issue/decode
L1 Cache	16KB Inst/16KB Data, 2-way, 1-cycle 64-byte line, write-through, no write-allocate
L2 Cache	512 KB, 4-way, 10-cycle 64-byte line, write-back, write-allocate
DRAM	
Memory Configuration	2 DIMMs/channel, 1 rank/DIMM 8 devices/rank, 8-bit output/device, 64-bit channel
DRAM Device Parameters	Micron MT41J128M8 DDR3-1600 [Micron Technology Inc. 2010] Timing $t_{RCD}-t_{RP}-t_{CL} = 10-10-10$ (12.5ns) 8 banks/device, 16,384 rows/bank, 1,024 columns/row 8KB row buffer
Total DRAM Capacity	1Gbit/device \times 8 devices/rank \times 2 ranks = 2GB
DRAM Queue Capacity	Total 96 entries per channel
High/Low Watermarks	32/16 (in 48 write queue entries)

Eager Write Command. The memory controller requires a new command type for eager write in order to differentiate between a write caused by a normal cache eviction and a write caused by a speculative writeback. This way, whenever a bank is precharged, the memory controller can discard all speculative writebacks to the row that used to be active in that bank. Furthermore, after the eager writeback has occurred to the DRAM, the memory controller must notify the LLC that the cache line is no longer dirty. In all other respects, the new command behaves identically to a conventional writeback.

Maintaining Correct Cache State. With cancellation, the dirty bit for a cache line cannot be cleared until the memory controller confirms that the eager write has been performed. However, just because the eager writeback has been committed to the DRAM does not mean that the dirty bit can be cleared, as the cache line may have been updated after the eager writeback was initiated.

An additional “eager” bit is necessary for each cache line to support cancellation. When an eager writeback is initiated, the eager bit is set. If the cache line is modified while the eager bit is set, the bit should be cleared. If the memory controller notifies the cache that the eager writeback completed, then the dirty bit will be cleared if and only if the eager bit is still set. The eager bit can always be safely cleared at that point.

A dirty cache line can be evicted while its eager bit is set. In that case, there may still be an eager writeback pending in the memory controller (if it has not been canceled). The memory controller should detect this situation and remove the eager writeback in favor of the actual eviction. Of course, this is a performance optimization that will not affect correctness.

Finally, when the memory controller precharges a bank, it can notify the LLC that all eager writebacks to that row have been canceled. The LLC can clear the associated eager bits and discontinue further lookups on the associated trigger. Again, these are performance optimizations that do not affect correctness.

Table V describes the baseline system that we simulate in our evaluation. Since the 512KB LLC contains 8,192 cache lines, the cancellation mechanism requires an extra 8,192 bits inside the LLC.

5.3. Eager Writeback Management

The Eager Writeback Manager (EWM) in the LLC is responsible for managing all eager writeback operations. Upon a triggering event, the memory controller sends the trigger address to the EWM. The EWM then creates an entry in an Eager Writeback Queue (EWQ) corresponding to that trigger.

For ERWC, there may be at most one active trigger per DRAM bank. Thus, for an 8-bank, 8-rank DRAM system, the EWQ would need to hold 64 entries. Each entry in the EWQ must store the next cache line address to search and the number of remaining accesses, which is initialized to $k - 1$ for a lookup range of k . The EWM services the EWQ in a round-robin fashion to ensure fairness in eager writeback traffic among the triggers. Whenever the LLC is otherwise idle, the EWM services the next EWQ entry by accessing the stored cache line address and initiating an eager writeback, if appropriate. After the access, the cache line address is updated subject to the DRAM address mapping, and the number of remaining accesses is decremented until it reaches zero. Upon reaching zero, or if the EWM is notified of a cancellation, the entry is discarded.

Our baseline system, shown in Table V, has an 8-bank, 2-rank DRAM system. Thus, the EWQ needs to hold 16 entries. Further, each EWQ entry needs to store one cache line address (58 bits, since the last 6 bits of the cache line address are always zero) and the number of remaining accesses (7 bits, since there are 128 cache lines in a row buffer). Therefore, ERWC's EWM requires an extra $(58 + 7) * 16 = 1,040$ bits inside the LLC.

For eviction-triggered clustering, including VWQ, the EWQ could be much larger, as there is no bound on the number of writeback operations that would cause triggering events. To avoid a large EWQ, VWQ uses an SSV (as described in Section 5.1). Under VWQ, the EWM is called the *Cache Cleaner*. When the Cache Cleaner receives the trigger address, it will read the corresponding section of the SSV. Then, the Cache Cleaner will lookup the cache sets corresponding to the SSV bits that show the existence of dirty cache lines and send eager writeback requests to the cache controller. In other words, the Virtual Write Queue, as its name implies, uses the SSV to maintain a virtual EWQ, thereby avoiding a large source of overhead in the LLC.

Our baseline system, shown in Table V, has a 512KB, 4-way LLC with a 64-byte line size. Thus, there are $512 * 1,024/4/64 = 2,048$ cache sets in the LLC. Since the SSV requires 1 bit per cache set, VWQ requires an extra 2,048 bits inside the LLC.

As described in Section 4, an LWP is used as the EWM under the LWPG scheme. An LWP is composed of a lightweight LLC write simulator and a prediction table. The address of the instruction accessing the cache line is stored in the prediction table and used as an input to the simulation. The simulator updates the prediction table based on the simulated write behavior of the LLC. When a dirty cache line is evicted from the L1 cache to the LLC, the LWP consults the prediction table to make a prediction of whether this cache line will be written to again. When the LWP selects a dirty cache line for eager writeback, it will put a corresponding write request into the last-write buffer.

Our baseline system, shown in Table V, has a 16KB, L1 cache with a 64-byte size line. Since the LWP keeps a 16-bit partial instruction pointer for each L1 cache line, it requires an extra 4,096 bits inside the L1 cache. For a 512KB, 4-way LLC with a 64-byte line size, the LLC write simulator requires at least 128 entries. Since each entry requires 36 bits, the simulator requires an extra 4,608 bits in the LLC. Therefore, the storage overhead of LWPG is 8,704 bits.

5.4. Memory Controller Queue Management

Separate read and write queues are used in the memory controller to provide a static partitioning of the available entries, isolating read and write accesses. High/low

watermarks are managed on the write queue to schedule DRAM operations in bursts [Stuecheli et al. 2010; Chatterjee et al. 2012]. The basic mechanism is to process a burst of reads until the write queue reaches a high watermark (or there are no pending reads). Once reached, writes begin to be scheduled over reads and drained from the write queue until a low watermark is reached. Even when a read/write clustering policy is used, this reduces switching between reads and writes at the DRAM, amortizing the delay incurred due to the write turnaround.

5.5. Design Complexity

Based on this discussion of mechanisms, we can quantitatively compare the different eager writeback schemes in Table IV in terms of their design complexity. The approaches described in Table IV that require the most complex designs are VWQ, LWPG, and ERWC.

Prior to this work, VWQ and LWPG both represented the state-of-the-art in eager writeback schemes. Previous work had shown that under LRU replacement LWPG has very similar performance to VWQ [Wang et al. 2012]. Moreover, as discussed in Section 5.3, LWPG has similar design complexity to VWQ in terms of storage overhead. Therefore, we include only one of these schemes, VWQ, in our evaluation. We do, however, include both of the less complex schemes, Eager [Lee et al. 2000] and DAW [Lee et al. 2010], in our evaluation.

ERWC has similar design complexity to VWQ. Both ERWC and VWQ require communication between the memory controller and the LLC for selecting dirty cache lines for eager writeback. ERWC's EWM and VWQ's Cache Cleaner require similar mechanisms for cache line lookup and eager writeback management. ERWC's primary storage overhead, the EWQ, depends on the number of DRAM banks. In contrast, VWQ's primary storage overhead, the SSV, depends on the number of LLC cache sets. To support the cancellation technique, ERWC requires an extra storage overhead as described in Section 5.2.

6. METHODOLOGY

We use a full-system simulator to evaluate the policy space. This simulator is based on HP Labs' COTSon [Argollo et al. 2009], which we have extended to simulate a DRAM system using DRAMSim2 [Rosenfeld et al. 2011]. DRAMSim2 accurately models all characteristics of the main memory system, including the state of all channels, ranks, banks, and rows, in a cycle-accurate manner. COTSon uses AMD's SimNow [Bedichek 2004] to emulate an x86-64 processor and to capture all instructions that are executed for the operating system and user-level processes. Table V shows the important parameters of the simulated system.

In all of the experiments, the baseline memory access scheduler is based on a column-first scheduling policy [Rixner et al. 2000] and gives reads higher priority than writes. The ERWC memory access scheduler is extended from the baseline and gives row-hit reads higher priority than row-hit eager writes. As a result, even when a few eager writebacks are issued in each bank, they can be gathered across the banks to form a long burst without any row activations.

The address space is mapped to a single channel using page interleaving. Detailed timing and power parameters for the memory devices are based upon the Micron MT41J128M8 DDR3-1600 datasheet [Micron Technology Inc. 2010]. Energy and power are calculated using Micron's power calculation methodology, assuming x8 DRAM devices [Micron Technology Inc. 2007]. We adopt the same DRAM queue parameters as in Chatterjee et al. [2012] for managing high/low watermarks.

We use a wide variety of programs, including Glimpse, TPC-C, and a large number of programs from the MineBench, SPEC CFP2006, and SPEC CINT2006 suites. Glimpse

is a text indexing and search program. We use a collection of program source and text files, whose total size is 2.3GB, as the input to Glimpse. TPC-C is a distributed, Online Transaction Processing (OLTP) benchmark specification. We use an open-source implementation of TPC-C, TPCC-UVa [Llanos and Palop 2006]. MineBench is a suite of data mining programs. We use *HOP*, *Utility*, *Apriori*, *Eclat*, *SEMPHY*, and *ScalParC* on large datasets. For each of these programs, we simulate 0.5 billion instructions. To eliminate the initialization phase from our simulations, we apply appropriate fast forwarding to each program. We have run some simulations for a larger number of instructions and obtained similar results.

7. EXPERIMENTAL RESULTS

This section evaluates various points in the eager writeback policy space, including *ERWC*, *VWQ*, and *DAW* as described in Table III. Unless stated otherwise, *ERWC* uses 2-way depth and 128-set range of lookup, and *VWQ* uses 2-way depth. For *DAW*, repeated writebacks are disabled because that has shown the best results. Throughout this section, unless stated otherwise, all results are normalized to the baseline system, *No Eager WB*, which does not perform eager writeback.

7.1. Comparison of *ERWC*, *DAW*, *VWQ* to *No Eager WB*

Row Activations and Bus Turnaround. Table VI presents the number of row activations (Row Act) and bus turnarounds (Turnaround) under *No Eager WB*, *ERWC*, *DAW*, and *VWQ* for each of the programs. The fractions in the Row Act and Turnaround columns for *ERWC*, *DAW*, and *VWQ* are relative to the total number of row activations and bus turnarounds, respectively, under *No Eager WB*.

Overall, all three schemes significantly reduce the number of row activations. Moreover, the number of row activations never increases for any of the programs. *ERWC* shows the greatest reduction in the number of row activations. With *ERWC*, the number of row activations is reduced by an average of 42% and a maximum of 82%. Eager writeback eliminates two kinds of row activations: row activations due to a DRAM write and row activations due to interference from a cache eviction (RD Interfered). RD Interfered represents the situation where a row is activated for a writeback operation that closes a row to which a subsequent read arrives. Were the writeback not to occur, that subsequent read would not require an additional row activation. Out of the total reduction in row activations that *ERWC* achieves, 38.4% comes from eliminating RD Interfered activations.

For most programs, *ERWC* reduces the number of bus turnarounds. *ERWC* may cause more bus turnarounds than *No Eager WB* when read and write accesses are highly interleaved. However, we only see this phenomenon in two programs (*hammer* and *cactusADM*) out of the 35 evaluated.

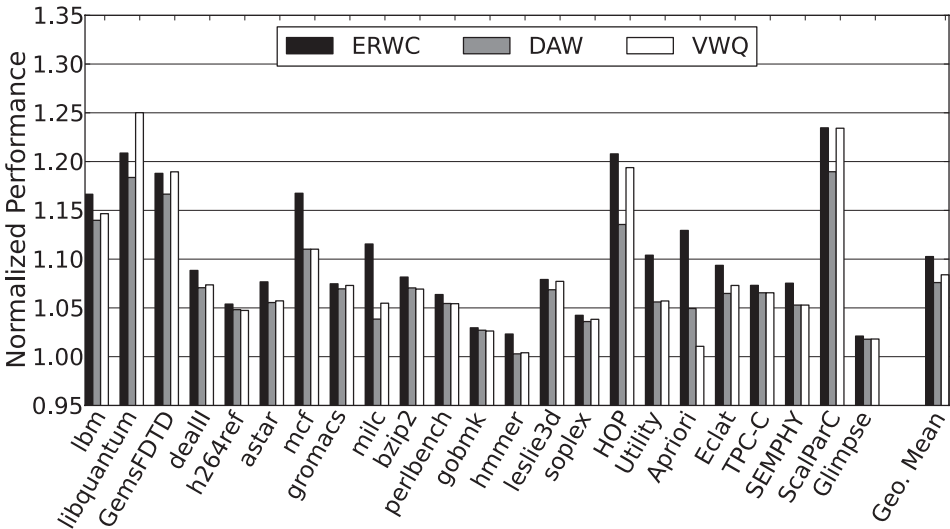
As *DAW* and *VWQ* both cluster only the writes, they achieve on average similar reductions in the number of row activations. *VWQ*, however, shows an overall greater reduction in bus turnarounds than *DAW* because *activation-triggered clustering* is able to form a longer write burst than can be stored in the physical write queue. Similarly, as *ERWC* is based on activation-triggered clustering, it often achieves a greater decrease in bus turnarounds than *DAW*.

A decrease in row activations and bus turnarounds can lead to improvements in performance and energy consumption. However, not all programs will experience substantial improvements. We found that programs that perform infrequent row activations or are highly read intensive get very minor benefits (<1% performance improvement) from eager writeback. From now on, we will exclude these programs from the evaluation and focus on the remaining 23 programs for the remainder of this section.

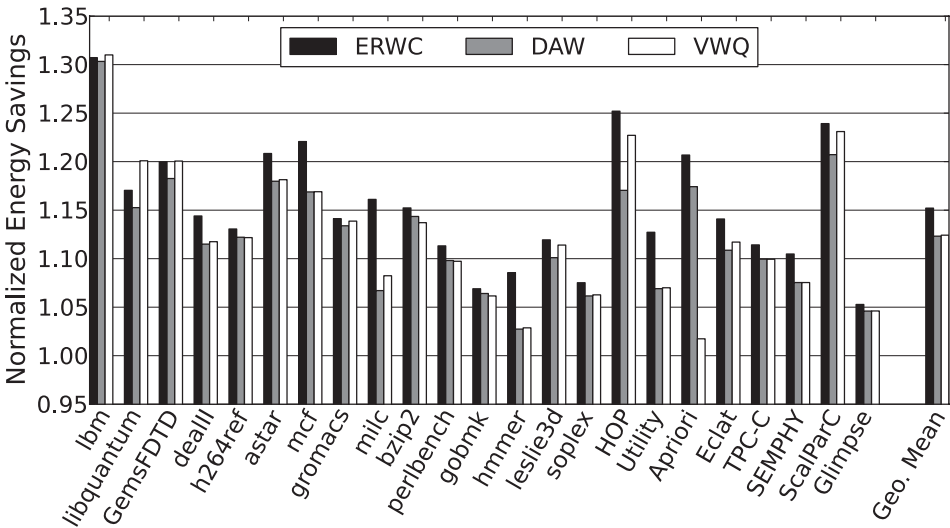
Table VI. Decomposition of Row Activations and Bus Turnarounds

	No Eager WB		ERWC		DAW		VWQ		
	Row Act	Turnaround	Row Act	Turnaround	Row Act	Turnaround	Row Act	Turnaround	
SPEC	lbm	8,245,106	795,196	1,481,934 (0.18)	367,303 (0.46)	1,336,779 (0.16)	338,519 (0.43)	1,224,276 (0.15)	315,585 (0.40)
	GemsFDTD	1,656,077	337,381	412,536 (0.25)	109,395 (0.32)	491,636 (0.30)	139,123 (0.41)	423,415 (0.26)	83,841 (0.25)
	libquantum	1,406,118	395,974	452,685 (0.32)	92,081 (0.23)	522,085 (0.37)	123,980 (0.31)	285,924 (0.20)	42,719 (0.11)
	dealll	466,640	164,958	177,018 (0.38)	47,462 (0.29)	234,356 (0.50)	72,750 (0.44)	229,301 (0.49)	70,963 (0.43)
	h264ref	449,582	71,979	199,772 (0.44)	41,679 (0.58)	211,565 (0.47)	52,082 (0.72)	211,461 (0.47)	51,906 (0.72)
	gromacs	462,521	136,919	202,790 (0.44)	30,248 (0.52)	214,454 (0.46)	32,127 (0.23)	207,552 (0.45)	30,608 (0.22)
	astar	3,379,938	511,793	1,693,132 (0.50)	428,359 (0.84)	1,882,007 (0.56)	426,641 (0.83)	1,873,619 (0.55)	421,266 (0.82)
	calculus	19,050	6,610	9,570 (0.50)	3,132 (0.47)	9,059 (0.48)	2,941 (0.44)	9,068 (0.48)	2,937 (0.44)
	milc	4,016,660	1,227,723	2,134,263 (0.53)	413,677 (0.34)	3,129,644 (0.78)	902,097 (0.73)	3,009,261 (0.75)	809,741 (0.66)
	mcf	4,134,388	767,829	2,219,199 (0.54)	450,326 (0.59)	2,563,544 (0.62)	522,441 (0.68)	2,561,289 (0.62)	521,775 (0.68)
	namd	30,324	8,254	17,672 (0.58)	3,052 (0.37)	18,505 (0.61)	3,534 (0.43)	19,988 (0.66)	3,253 (0.39)
	bzip2	2,033,258	441,150	1,226,393 (0.60)	300,438 (0.68)	1,248,487 (0.61)	266,384 (0.60)	1,287,795 (0.63)	262,724 (0.60)
	perlbench	861,580	247,099	513,567 (0.60)	125,926 (0.51)	557,722 (0.65)	149,095 (0.60)	559,655 (0.65)	147,685 (0.60)
	tonto	54,606	17,739	35,064 (0.64)	11,656 (0.66)	36,244 (0.66)	12,930 (0.73)	51,853 (0.95)	16,690 (0.94)
gobmk	552,400	165,748	366,036 (0.66)	103,577 (0.62)	376,576 (0.68)	96,953 (0.58)	382,818 (0.69)	93,348 (0.56)	
gcc	63,817	13,899	43,474 (0.68)	6,822 (0.49)	44,045 (0.68)	5,095 (0.37)	42,549 (0.67)	2,279 (0.16)	
povray	1,730	493	1,239 (0.72)	363 (0.74)	1,321 (0.76)	367 (0.74)	1,578 (0.91)	402 (0.82)	
hammer	2,174,007	505,813	1,564,088 (0.72)	897,927 (1.78)	1,950,614 (0.90)	452,936 (0.90)	1,952,980 (0.90)	447,160 (0.88)	
leslie3d	5,315,326	896,517	3,845,258 (0.72)	235,045 (0.26)	4,074,521 (0.77)	269,202 (0.30)	3,923,055 (0.74)	243,450 (0.27)	
zeusmp	1,138,834	168,491	903,849 (0.79)	99,248 (0.59)	916,398 (0.80)	62,358 (0.37)	957,032 (0.84)	65,626 (0.39)	
soplex	7,002,837	999,078	5,703,923 (0.81)	760,790 (0.76)	5,955,756 (0.85)	795,694 (0.80)	5,953,533 (0.85)	792,006 (0.79)	
cactusADM	3,082,844	477,943	2,580,802 (0.84)	514,823 (1.08)	2,661,445 (0.86)	228,540 (0.48)	2,660,496 (0.86)	227,027 (0.48)	
sjeng	343,633	89,265	293,603 (0.85)	84,054 (0.94)	308,894 (0.90)	77,675 (0.87)	319,945 (0.93)	54,699 (0.61)	
sphinx3	4,773,712	200,468	4,353,758 (0.91)	65,411 (0.33)	4,490,430 (0.94)	103,720 (0.52)	4,489,672 (0.94)	103,589 (0.52)	
garness	3,599	700	3,336 (0.93)	615 (0.88)	3,364 (0.93)	585 (0.84)	3,492 (0.97)	595 (0.85)	
omnetpp	1,733,037	66,706	1,628,320 (0.94)	44,421 (0.67)	1,648,553 (0.95)	45,200 (0.65)	1,647,949 (0.95)	44,923 (0.67)	
bwaves	17,709,425	468,428	17,229,483 (0.97)	431,716 (0.92)	17,296,536 (0.98)	306,401 (0.68)	17,277,111 (0.98)	299,882 (0.64)	
Other	HOP	3,056,649	584,078	1,213,863 (0.40)	136,143 (0.23)	1,749,612 (0.57)	199,583 (0.34)	1,357,534 (0.44)	216,965 (0.37)
	Eclat	914,824	226,619	413,102 (0.45)	76,354 (0.34)	506,599 (0.55)	116,504 (0.51)	483,626 (0.53)	104,519 (0.46)
	SEMPHY	1,717,476	246,476	843,845 (0.49)	62,198 (0.25)	1,074,476 (0.63)	109,817 (0.45)	1,034,523 (0.60)	109,172 (0.44)
	Apriori	14,101,730	779,383	7,558,263 (0.54)	713,859 (0.92)	7,715,157 (0.55)	738,206 (0.95)	12,484,292 (0.89)	506,190 (0.65)
	ScalParC	10,336,596	1,174,103	6,414,952 (0.62)	362,996 (0.31)	6,803,180 (0.66)	385,918 (0.33)	6,625,802 (0.64)	298,794 (0.25)
	TPC-C	2,976,906	901,631	2,007,013 (0.67)	504,069 (0.56)	2,145,059 (0.72)	648,006 (0.72)	2,145,475 (0.72)	647,763 (0.72)
	Glimpse	280,968	75,884	191,841 (0.68)	44,293 (0.58)	202,687 (0.72)	47,240 (0.62)	202,670 (0.72)	46,314 (0.61)
	Utility	2,811,884	283,106	2,050,784 (0.73)	145,587 (0.51)	2,394,654 (0.85)	219,635 (0.78)	2,388,839 (0.85)	220,138 (0.78)
	Geo. Mean			(0.58)	(0.52)	(0.64)	(0.55)	(0.63)	(0.50)

The SPEC and Other programs are presented in increasing order of the fraction of row activations in ERWC.



(a) Performance improvement



(b) Energy savings

Fig. 2. ERWC/DAW/VWQ performance improvement and energy savings for 23 programs (normalized to the results for No Eager WB). Larger numbers are better.

Performance and Energy. Figure 2 shows the improvements in performance and energy consumption for the three competing schemes. The results show that ERWC outperforms VWQ for 21 out of 23 programs and always outperforms DAW. Specifically, ERWC achieves an average of 10.3% performance improvement compared to No Eager WB, which is 3.8% and 2.0% larger than the improvements achieved by DAW and VWQ, respectively. Moreover, ERWC improves energy consumption by an average of 15.2%, which is almost 3.0% greater than that achieved by the other schemes. It is interesting to note that VWQ outperforms ERWC in *libquantum*.

Table VII. Improvements (%) in Different Configurations for ERWC/DAW/VWQ for 23 Programs

Configuration	Perf.	Energy	Row Act.	Turn.
ERWC	10.3	15.2	49.5	54.2
ERWC, 8 burst	13.1	18.3	61.9	76.4
ERWC, 8 burst, Idle Prefill	13.1	18.3	61.9	76.5
DAW	6.5	12.3	42.7	46.1
DAW, 8 burst	10.4	16.3	58.6	79.3
VWQ	8.3	12.4	44.7	52.1
VWQ, 8 burst	11.6	16.8	58.8	80.8

performs a considerable number of read accesses to each activated row (see Table II). VWQ outperforms ERWC because it significantly reduces the number of bus turnarounds.

With ERWC, out of the 23 total programs, 6 have overall performance improvements between 10% and 20%, and 3 have improvements in excess of 20%. Moreover, 12 consume between 10% and 20% less DRAM energy, and 7 have energy savings in excess of 20%. The magnitude of the improvements corresponds closely to the reduction in row activations from Table VI.

In fact, the improvements from using eager writeback not only depend on the clustering policy but also depend on other features of the system configuration. Table VII shows a summary of the improvements for a variety of configurations with selected features.

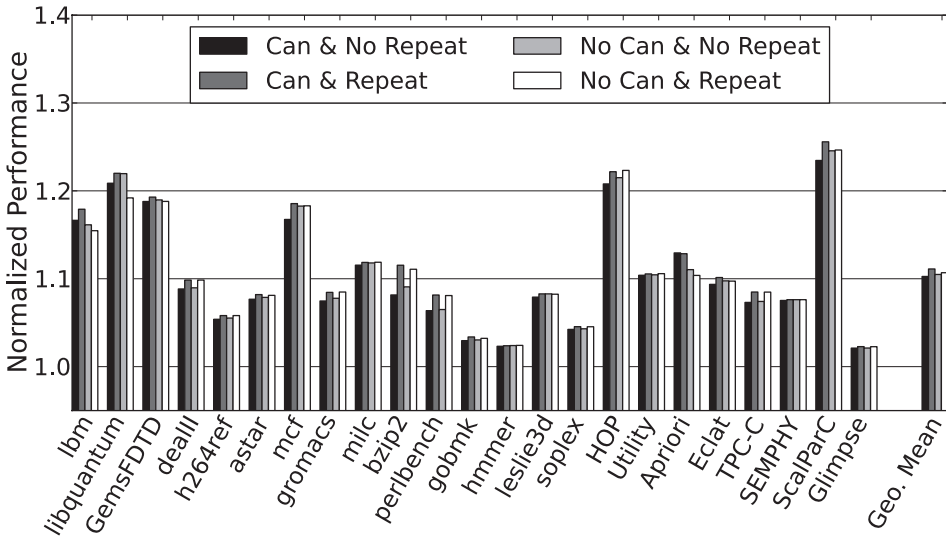
An interesting finding is that having a minimum burst size for processing writes has a positive effect on all of the schemes. With this feature, if a burst of writes is in progress to the DRAM in response to idle DRAM cycles, the memory controller attempts to process a minimum number of writes without being preempted by reads. Typically, under any of the eager writeback schemes, a burst of writes can be effectively pipelined to activated rows once the initial cost has been paid. Therefore, the delay to reads can be amortized by processing multiple row-hit writes in a fairly short time. However, the benefits are not proportional to the length of the bursts. Based on our observations, 8-write bursts are sufficient and longer bursts do not offer any significant additional benefit.

With *Idle Prefill*, the memory controller can preemptively ask the LLC for dirty cache lines in order to promptly react to rank idleness with writes. We find that Idle Prefill does not beneficially combine with the eager writeback schemes. If a minimum burst is used, Idle Prefill adds little to performance. Even worse, without the burst for writes, Idle Prefill just keeps precharging rows that have been activated by a read and will be utilized by subsequent accesses.

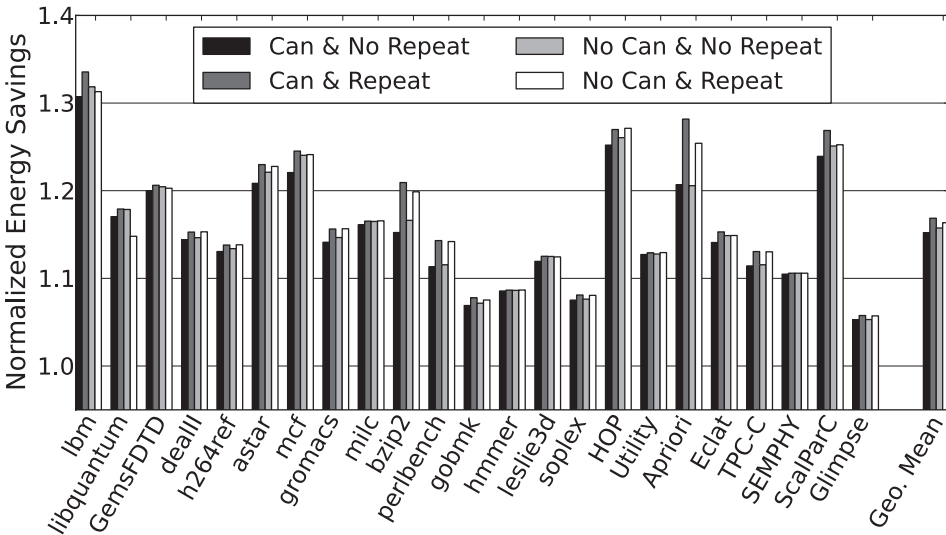
Finally, we observe that *eviction-triggered clustering* (i.e., DAW) does not perform as well overall as *activation-triggered clustering* (i.e., ERWC or VWQ). By clustering accesses earlier than is actually required, eviction-triggered clustering quickly fills up the physical write queue and subsequently delays the processing of reads more often. In contrast, with activation-triggered clustering, the dirty cache lines to cluster are brought into the DRAM at the time the physical write queue drains. So, the high watermark is less often reached.

7.2. Eager WB Cancellation, Repetition, and Clustering

Figure 3 shows the trade-offs involved in the four possible combinations of cancellation (*Can*) and repeated eager writebacks (*Repeat*) when read/write clustering is performed. These combinations are evaluated under 2-way depth and 128-set range of lookup. As the figure shows, all variations provide improvements over No Eager WB for all



(a) Performance improvement



(b) Energy savings

Fig. 3. Performance improvement and energy savings with variations in Eager WB Cancellation (Can) and Repeated Eager WB (Repeat).

programs. There is a trade-off between having less repeated write traffic that competes for memory bandwidth with cache line reads and serving more writes on activated rows. For most programs, the use of Repeat is beneficial since LRU lines are rarely written again [Lee et al. 2000]. The figure also shows that when repeated eager writebacks are allowed, cancellation is critical to avoid triggering excess row activations to perform these writebacks. Moreover, cancellation helps to avoid the interference between writes and incoming read requests. However, cancellation, by itself, seems to generate many short write bursts, as observed in Can & No Repeat.

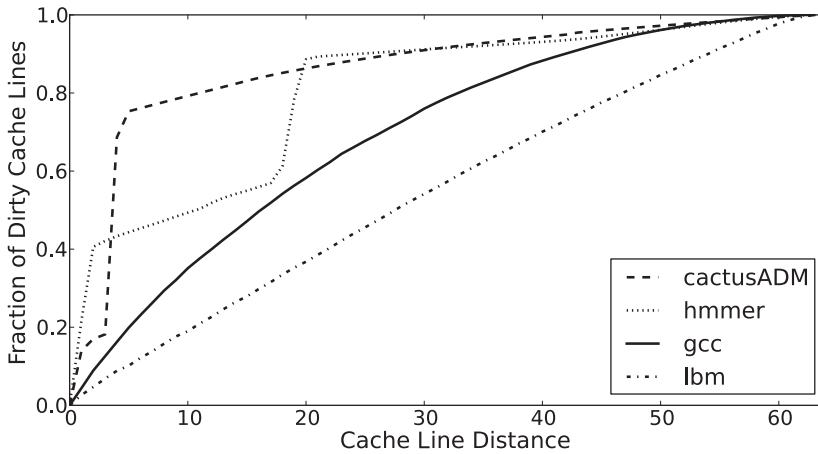


Fig. 4. CDF of distance from trigger to dirty cache lines.

7.3. Lookup Depth/Range Analysis

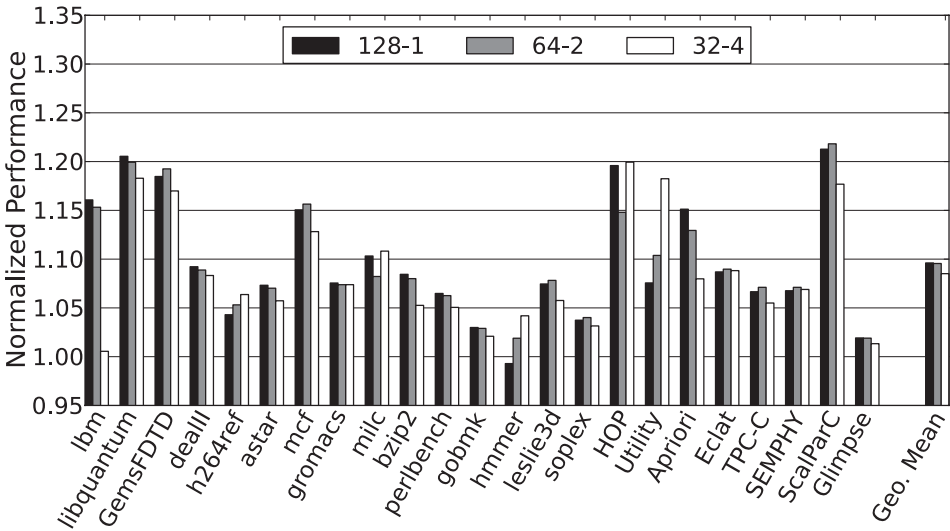
Figure 4 shows the Cumulative Distribution Function (CDF) of the distance (in cache lines) from an eager writeback triggering address to the other dirty cache lines in that row. The figure shows four representative programs. For all but nine of the programs that we evaluated, the CDFs are similar to *gcc*. As the figure shows, dirty cache lines are not often near the triggering access. In fact, to find 50% of the dirty cache lines, 15 cache lines must be checked (in both directions) around the triggering access. This suggests that larger lookup ranges are going to perform significantly better.

Figure 5 shows the effect of the lookup depth and range on performance and energy consumption. Each benchmark has a set of three bars, which are the results of a 128-line lookup. The 128-line lookups are organized as the following (range, depth) combinations: (128, 1), (64, 2), (32, 4).

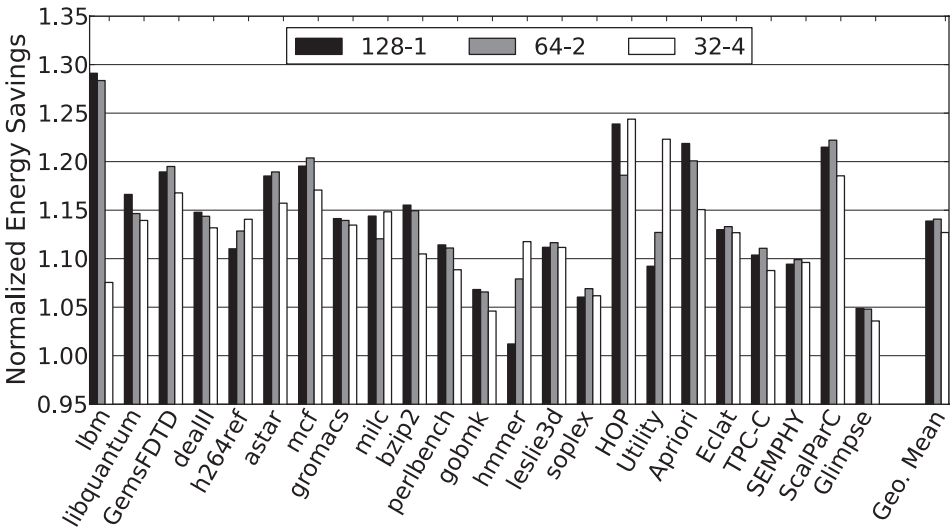
Using a depth of four results in many eager writebacks of cache lines from more recently used ways that will later be written again before they are evicted. So, these useless eager writebacks consume bandwidth in the cache and memory system for no benefit. Although for some programs the 4-way depth configurations achieve better results, for most programs the useless eager writes lead to worse performance than the 2-way depth configurations. It is also important to note that the 1-way depth configurations achieve similar results to the 2-way depth configurations. The trends are similar in 64-line lookups, (64, 1), (32, 2), (16, 4), and 32-line lookups (32, 1), (16, 2), (8, 4), but overall improvements decrease as the lookup size decreases.

For most of the programs, given the same lookup depth, configurations with larger lookup ranges produce more improvements. For example, results of *milc* show that (128, 1) improves the performance 3.1% and 6.8% more than (64, 1) and (32, 1), respectively. These results intuitively make sense (especially with smaller lookup depth) because a larger lookup range provides a higher chance to find dirty cache lines that would be eagerly written back.

The trade-offs on the lookup depth and range that we described earlier for read/write clustering also occur for write clustering. Our analysis shows that for 128-line lookup, ERWC produces better performance improvements of 3.3%, 3.0%, and 3.1% for depths of 1, 2, and 4, respectively, than DAW. Moreover, for the same lookup, ERWC shows better energy savings of 3.6%, 3.3%, and 2.7% for depths of 1, 2, and 4, respectively.



(a) Performance improvement



(b) Energy savings

Fig. 5. Performance improvement and energy savings with a variety of lookup range/depth combinations. Legends denote range-depth combinations.

Optimization Opportunities. Stuecheli et al. [2010] proposed the SSV under VWQ to reduce the lookup overhead for dirty cache lines by keeping track of their existence in the designated LRU ways of a cache set. We can similarly use the SSV in ERWC to reduce the lookup overhead. If the SSV’s bit for a cache set is zero, then we do not need to examine that set because there are no dirty cache lines in its LRU ways. For the programs that we evaluated, we found that by using the SSV, the number of lookups that must be performed for eager writebacks is reduced by an average of 57.7% and a maximum of 98%.

Cache Lookup Energy Overhead. In terms of energy consumption, there are two kinds of cache lookup overhead for any eager writeback scheme. The first of these overheads is the tag array access energy, which is consumed when the EWM accesses the tag bits to find dirty cache lines. The second overhead is the data array access energy, which is consumed when the EWM retrieves a dirty cache line's data. To quantify both of these overheads, we used CACTI 6.5 [Muralimanohar et al. 2009] to model a 32nm, 512KB, 4-way L2 cache with a 64-byte line size. These results show that the tag array access energy is 0.007nJ/access and the data array access energy is 0.127nj/access.

To quantify the overall cache lookup energy, we calculated the total energy consumption by tag accesses and data accesses for the 0.5 billion instructions that we simulated for each program. For ERWC+SSV, the lookup energy overhead is 0.738mJ on average (ranging from 0.024mJ to 5.223mJ) across the programs, which is on average 1% (ranging from 0.05% to 2.5%) of the energy reduction in the DRAM. For VWQ, the lookup energy overhead is 0.742mJ on average, ranging from 0.019mJ to 10.823mJ. Breaking down the overall lookup energy overhead for ERWC+SSV, the tag access energy overhead accounts for the largest portion (0.695mJ). Only 0.29% of the extra tag accesses result in extra data accesses. Thus, the data access energy overhead (0.043mJ) is comparatively small. As a result, the cache lookup energy overhead for ERWC+SSV is small compared to the energy reduction in DRAM.

7.4. Sensitivity Studies

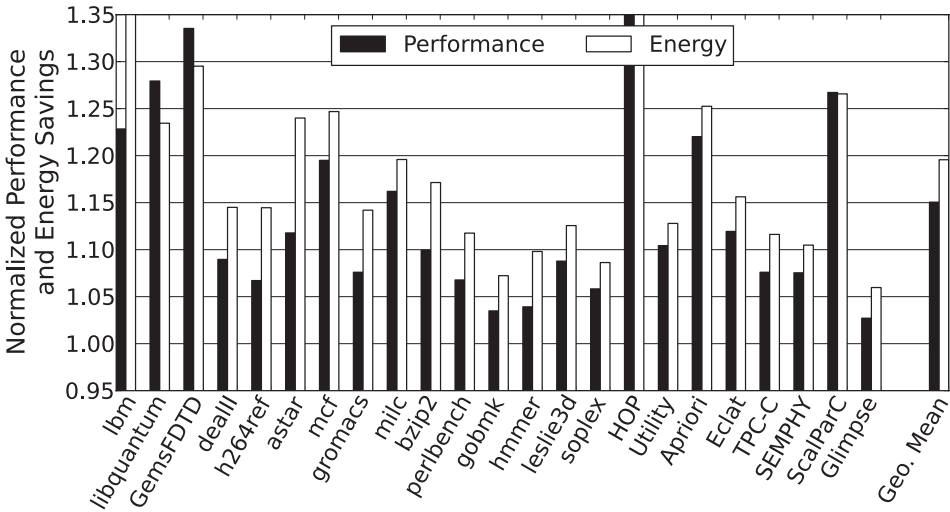
This section evaluates the sensitivity of the results for ERWC to a variety of different system design parameters.

LLC Size. When the LLC size is increased from 512KB to 2MB, the benefits of ERWC go from 10.3% and 15.2% down to 8.0% and 11.9% for performance and energy consumption, respectively. These slight reductions are caused by the larger LLC filtering DRAM accesses in such a way that there are fewer overall row activations. In particular, the average number of row activations caused by reads and writes are reduced by 26% and 14%, respectively. This reduces the opportunities for ERWC to eliminate row activations.

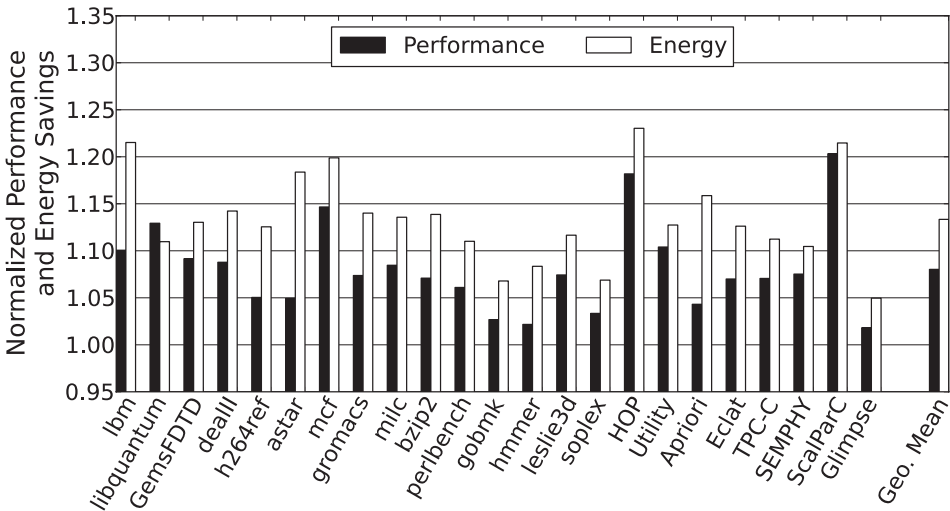
Cache Latency. With a faster cache access, an instruction is more likely delayed by DRAM latency than cache latency. ERWC in fact achieves similar improvements in both performance and energy consumption even when a cache hit at any level of the memory hierarchy adds only one-cycle latency. Specifically, we see the improvements of 10.0% and 14.9% for performance and energy consumption, respectively. Improvements with ERWC are closely related to access patterns and pressure on DRAM, and cache latency seems less relevant. A program's DRAM access pattern may not be largely changed with cache latency. Moreover, there are several factors, such as CPU pipeline bandwidth and dependency between instructions, that affect the amount of outstanding DRAM accesses.

DRAM Queue Parameters. The use of high and low watermarks affects the duration and frequency of write bursts. Since reads must be stalled while draining writes, they directly affect system performance. Longer write bursts induce longer read stalls, but they are less frequent. To see how these stalls interact with ERWC, we perform simulations for two different configurations with high watermarks of 16 and 64 and low watermarks of 8 and 32, respectively.

Figure 6 shows the effects of changing the high and low watermarks. The figure shows the performance of ERWC with these watermarks normalized to No Eager WB with the same watermarks. These figures can be compared to Figure 2, which shows the same comparison with the default high/low watermarks of 32/16. As the figure



(a) 16/8 high/low watermarks



(b) 64/32 high/low watermarks

Fig. 6. Performance improvement and energy savings with different write queue sizes. In (a), *libm* energy (1.39), *HOP* performance (1.86), and *HOP* energy (1.94) are cut off due to the space.

shows, the ERWC improvements are higher with smaller watermarks. The geometric mean of performance improvement and energy savings are 15.1% and 19.6% for the small case (Figure 6a), 10.3% and 15.2% for the default case (Figure 2), and 8.0% and 13.3% for the large case (Figure 6b). This results from ERWC being very effective at clustering reads and writes in the same row. With 16/8 watermarks, the write queue must be serviced frequently, causing many disruptive row activations without ERWC. In summary, although eager writeback provides larger improvements with smaller watermarks, it is still effective at improving performance and reducing energy with much larger watermarks.

Bank Parallelism. The number of banks in the main memory system affects how frequently eager writes are able to be scheduled on the activated row. Each bank attempts to produce a DRAM command that is issuable every cycle. With more banks, reads tend to be spread across banks, so there are fewer row conflicts. This provides less opportunity to schedule eager writebacks, as the reads are effectively using the available DRAM bandwidth. To test this effect, we ran the programs on a DRAM of 64 banks (8 ranks of 8 banks each). Even with larger bank parallelism, ERWC still provides performance and energy improvements of 4.5% and 6.0%, respectively.

8. RELATED WORK

Lee et al. [2000] first proposed *eager writeback* to optimize systems that experience competition for memory bandwidth between cache writebacks and data fetches. Eager writeback can effectively redistribute and balance the memory traffic for those systems and thereby improve system performance. VWQ [Stuecheli et al. 2010] exploited eager writeback to increase row-level locality of DRAM writes. This technique coordinates DRAM access scheduling with cache writeback to increase the opportunities to find dirty cache lines in the LLC that can hit on the activated row. LWPG [Wang et al. 2012] writeback extended this idea to take the similar advantage in any cache replacement policies (including NRU and random). Our work also focuses on increasing row buffer locality but is not limited to DRAM writes. Moreover, none of these works explored the policy space.

Clustering of accesses can be realized across row boundaries. Sudan et al. [2010] recently observed that a large number of accesses within heavily accessed OS pages happen to small, contiguous chunks of cache lines and proposed the co-location of the clusters from different OS pages in a row buffer to improve its utilization.

Some efforts have been made to try to avoid access interference within banks. The XOR-based address mapping [Zhang et al. 2000] scheme avoids row buffer conflicts between reads and cache writebacks under page interleaving by generating bank addresses pseudorandomly. This technique effectively redirects writes that would conflict in a bank to another bank. A write buffer with read bypass [Chen and Baer 1992] could also alleviate row buffer conflicts by postponing writebacks. This potentially would allow consecutive reads to be grouped together. However, these schemes do not eliminate write interference but rather attempt to reduce it, and neither scheme reduces the number of row activations caused by writes.

There are several works, such as McKee et al. [2000], Rixner et al. [2000], and Hur and Lin [2004, 2007], that present techniques to reschedule memory accesses to improve overall system performance. These works introduce a variety of scheduling mechanisms and algorithms to increase performance. However, they do not consider energy consumption and are orthogonal to our eager writeback scheme, which introduces more opportunity for read and write clustering in the accesses to be scheduled.

Energy and latency could further be reduced by utilizing smaller rows within the DRAM. Rixner [2004] and Cooper-Balis and Jacob [2010] presented memory controller policies that can make effective use of commercial DRAM architectures that support the use of subsets of rows to further reduce the average latency of the DRAM. These types of proposals are being revisited in the context of multicore systems where access streams are mixed creating more interference in the row buffers. Udipi et al. [2010] proposed a redesign of DRAM to activate subsets of a row and to keep inactive subarrays in low-power sleep modes. Others have proposed a *rank subsetting* that enables the smaller number of devices to be involved for a memory access [Zheng et al. 2008; Ahn et al. 2009]. All of these techniques are orthogonal to our read and write clustering strategy, and we would expect the ideas to be complementary.

Taking advantage of low-power modes has also been focused in the memory controller. Hur and Lin [2008] extended the Adaptive History-Based Scheduler (AHB) with power consumption to increase the average idle duration of each rank, thus increasing the utility of the power-down unit. Huang et al. [2005] identified memory access traffic often random when the OS arbitrarily maps virtual pages to physical pages. They thus proposed a way to reshape such random memory traffic to produce longer idle periods by which a more aggressive power-saving mode is actively utilized.

9. CONCLUSIONS

This article has shown that cache line writebacks to the DRAM have a disproportionate impact on performance and energy consumption in the main memory system. In particular, the fraction of row activations due to writes is higher than the fraction of writes in the DRAM access stream.

To address this problem, this article has introduced ERWC to optimize DRAM performance and energy efficiency. Whenever a row is activated, either by a read or write access, eager writeback operations to that row are triggered. Clustering both read and write operations in this manner, instead of just write operations, leads to superior improvements in both performance and energy consumption. Read/write clustering increases the row-level locality of access and reduces the number of row activations, increasing bandwidth and reducing energy. In fact, ERWC reduces row activations by an average of 42% for 35 programs. For the programs with more than 15% DRAM writes, row activations are reduced by 43%, and for the programs with more than 30% DRAM writes, row activations are reduced by 48%. This reduction results in improvements in both performance and energy efficiency. Specifically, 9 programs have performance improvements in excess of 10%, and 19 programs have energy consumption reductions in excess of 10%.

The article has also explored the design space of eager writeback schemes. This has yielded three key insights. First, dirty cache lines are not usually near the cache line that triggers the eager writeback. This necessitates larger lookup ranges in order to provide the best improvements. However, smaller lookup depths usually perform best. This validates the use of the SSV to mitigate the cost of these lookups [Stuecheli et al. 2010]. Second, the improvements of read/write clustering increase as the write queue watermarks are decreased. Therefore, eager writeback enables the use of smaller write buffers and shorter write bursts. Third, canceling eager writebacks in read/write clustering is important, whereas it is detrimental with eviction-triggered write clustering. With write clustering, eager writebacks are always triggered by an actual writeback operation. Therefore, it is guaranteed that row will be activated, so it makes sense to always perform the pending eager writebacks to that row in a burst. In contrast, with read/write clustering, a read operation may have activated the row, so the bus turnaround penalty may make it less desirable to try to perform all of the eager writeback operations.

This article used one out-of-order single-core processor in all of its experiments. Running multiple threads or programs on a multicore system is known to have two effects. First, the reduced locality of reference in the shared LLC results in additional cache misses and DRAM accesses. However, our experiments with a shorter cache latency, which also results in more outstanding DRAM accesses at a time, still showed substantial improvements in both performance and energy consumption with ERWC. Specifically, with single-cycle cache latency, the benefits of ERWC only went down slightly, from 10.3% and 15.2% to 10.0% and 14.9% for performance and energy consumption, respectively. Moreover, Lee et al. [2010] have reported that in 4- and 8-core systems, the shared LLC still has around 90% idle cycles, and so we do not foresee the additional load on the LLC to be a problem. The second effect is

that additional DRAM accesses may introduce more read interference. However, several techniques, including XOR-based address mapping [Zhang et al. 2000] and a write buffer with read bypass [Chen and Baer 1992], can alleviate row buffer conflicts between reads and cache writebacks. Therefore, we expect that eager writeback schemes, like ERWC, will still achieve substantial improvements in performance and energy consumption on multicore systems.

Eager writeback is a promising approach to reducing row activations in DRAM. Our article has shown the trade-offs along several dimensions to motivate the proposed eager writeback policies. This has shown that many policy decisions improve performance and energy consumption and that the proposed ERWC scheme outperforms previous strategies.

REFERENCES

- AHN, J. H., LEVERICH, J., SCHREIBER, R., AND JOUPPI, N. P. 2009. Multicore DIMM: An energy efficient memory module with independently controlled DRAMs. *IEEE Comput. Archit. Lett.* 8, 5–8.
- ARGOLLO, E., FALCÓN, A., FARABOSCHI, P., MONCHIERO, M., AND ORTEGA, D. 2009. COTSon: Infrastructure for full system simulation. *SIGOPS Oper. Syst. Rev.* 43, 1, 52–61.
- BEDICHEK, R. 2004. SimNow: Fast platform simulation purely in software. In *Proceedings of the Symposium on High Performance Chips (Hot Chips'04)*.
- CHATTERJEE, N., MURALIMANOHAR, N., BALASUBRAMONIAN, R., DAVIS, A., AND JOUPPI, N. P. 2012. Staged reads: Mitigating the impact of DRAM writes on DRAM reads. In *Proceedings of the 18th International Symposium on High-Performance Computer Architecture (HPCA'12)*. 1–12.
- CHEN, T.-F. AND BAER, J.-L. 1992. Reducing memory latency via non-blocking and prefetching caches. In *Proceedings of the 5th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'92)*. 51–61.
- COOPER-BALIS, E. AND JACOB, B. 2010. Fine-grained activation for power reduction in dram. *IEEE Micro* 30, 34–47.
- HUANG, H., SHIN, K. G., LEFURGY, C., AND KELLER, T. 2005. Improving energy efficiency by making DRAM less randomly accessed. In *Proceedings of the International Symposium on Low Power Electronics and Design (ISLPED'05)*. 393–398.
- HUR, I. AND LIN, C. 2004. Adaptive history-based memory schedulers. In *Proceedings of the 37th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO'04)*. 343–354.
- HUR, I. AND LIN, C. 2007. Memory scheduling for modern microprocessors. *ACM Transactions on Computer Systems* 25, 4.
- HUR, I. AND LIN, C. 2008. A comprehensive approach to DRAM power management. In *Proceedings of the 14th International Conference on High-Performance Computer Architecture (HPCA'08)*. 305–316.
- LEE, C. J., NARASIMAN, V., EBRAHIMI, E., MUTLU, O., AND PATT, Y. N. 2010. *DRAM-aware Last-Level Cache Writeback: Reducing Write-caused Interference in Memory Systems*. Technical Report TR-HPS-2010-002. The University of Texas at Austin.
- LEE, H.-H. S., TYSON, G. S., AND FARRENS, M. K. 2000. Eager writeback—a technique for improving bandwidth utilization. In *Proceedings of the 33rd International Symposium on Microarchitecture (MICRO'00)*. 11–21.
- LLANOS, D. R. AND PALOP, B. 2006. TPCC-UVa: an open-source TPC-C implementation for parallel and distributed systems. In *Proceedings of the 20th International Parallel and Distributed Processing Symposium*.
- McKEE, S., WULF, W., AYLOR, J., KLENKE, R., SALINAS, M., HONG, S., AND WEIKLE, D. 2000. Dynamic Access Ordering for Streamed Computations. *IEEE Transactions on Computers* 49, 11, 1255–1271.
- MICRON TECHNOLOGY, INC. 2007. Calculating Memory System Power for DDR3. Technical Report TN-41-01.
- MICRON TECHNOLOGY INC. 2010. Micron MT41J128M8 DDR3-1600. Retrieved from http://download.micron.com/pdf/datasheets/dram/ddr3/1Gb_DDR3_SDRAM.pdf.
- MURALIMANOHAR, N., BALASUBRAMONIAN, R., AND JOUPPI, N. P. 2009. *Cacti 6.0: A Tool to Model Large Caches*. HP Laboratories.
- RIXNER, S. 2004. Memory controller optimizations for web servers. In *Proceedings of the 37th International Symposium on Microarchitecture (MICRO'04)*. 355–366.
- RIXNER, S., DALLY, W. J., KAPASI, U. J., MATTSON, P., AND OWENS, J. D. 2000. Memory access scheduling. In *Proceedings of the 27th Annual International Symposium on Computer Architecture (ISCA'00)*. 128–138.

- ROSENFELD, P., COOPER-BALIS, E., AND JACOB, B. 2011. DramSim2: A cycle accurate memory system simulator. *Computer Architecture Letters* 10, 1, 16–19.
- STUECHELI, J., KASERIDIS, D., DALY, D., HUNTER, H. C., AND JOHN, L. K. 2010. The virtual write queue: coordinating DRAM and last-level cache policies. In *Proceedings of the 37th International Symposium on Computer Architecture (ISCA'10)*. 72–82.
- SUDAN, K., CHATTERJEE, N., NELLANS, D., AWASTHI, M., BALASUBRAMONIAN, R., AND DAVIS, A. 2010. Micro-pages: Increasing DRAM efficiency with locality-aware data placement. In *Proceedings of the 15th International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS'10)*. 219–230.
- UDIPI, A. N., MURALIMANO HAR, N., CHATTERJEE, N., BALASUBRAMONIAN, R., DAVIS, A., AND JOUPPI, N. P. 2010. Rethinking DRAM design and organization for energy-constrained multi-cores. In *Proceedings of the 37th International Symposium on Computer Architecture (ISCA'10)*. 175–186.
- WANG, Z., KHAN, S. M., AND JIMÉNEZ, D. A. 2012. Improving writeback efficiency with decoupled last-write prediction. In *Proceedings of the 39th International Symposium on Computer Architecture (ISCA'12)*. 309–320.
- ZHANG, Z., ZHU, Z., AND ZHANG, X. 2000. A permutation-based page interleaving scheme to reduce row-buffer conflicts and exploit data locality. In *Proceedings of the 33rd International Symposium on Microarchitecture (MICRO'00)*. 32–41.
- ZHENG, H., LIN, J., ZHANG, Z., GORBATOV, E., DAVID, H., AND ZHU, Z. 2008. Mini-rank: Adaptive DRAM architecture for improving memory power efficiency. In *Proceedings of the 41st International Symposium on Microarchitecture (MICRO'08)*. 210–221.

Received May 2013; revised September 2013; accepted November 2013