

A Study of Leveraging Memory Level Parallelism for DRAM System on Multi-Core/Many-Core Architecture

Licheng Chen^{1,2}, Yongbing Huang^{1,2}, Yungang Bao¹, Guangming Tan¹, Zehan Cui^{1,2}, Mingyu Chen¹

¹State Key Laboratory of Computer Architecture, Institute of Computing Technology, Chinese Academy of Sciences

²University of Chinese Academy of Sciences

{chenlicheng, huangyongbing, baoyg, tgm, cuizehan, cmy}@ict.ac.cn

Abstract—DRAM system has been more and more critical on modern multi-core/many-core architecture where the Moore’s law has been made effect on increasing the number of cores integrated in a processor chip. The performance of DRAM system is usually measured in term of bandwidth and latency, which are regarded as inherently depending on Row Buffer Hit Rate (RBHR) according to previous studies. In this paper, we find that Memory Level Parallelism (MLP) exhibits a stronger correlation with the performance of DRAM system on multi-core/many-core architecture than RBHR, and promoting MLP significantly improves DRAM system performance. In order to exploit the MLP, we have evaluated various approaches including multi-bank, multi-row-buffers, multi-memory-controllers and the obsolete Virtual Channel Memory (VCM). The experimental results show that VCM is a better alternative to traditional DRAM chip on multi-core/many-core architecture than the other three approaches because VCM has almost all the advantages of the others: 1) it can improve homogeneous workloads’ IPC by 2.21X on a 16-core system with 32 virtual channels due to leveraging unexploited MLP. 2) It can also promote Quality-of-Service (QoS) of DRAM system by removing unfairness while memory controllers serve memory requests. 3) It can save energy and has low area costs. Unfortunately, VCM, which was proposed in the late 1990s, faded away before multi-core/many-core became dominated. Therefore, we suggest memory chip vendors reconsider the VCM technology for multi-core/many-core architecture.

Keywords—DRAM; Virtual Channel Memory; Memory Level Parallelism; QoS

I. INTRODUCTION

“MLP yes! ILP no!” Memory Level Parallelism (MLP) was originally proposed in term of the number of outstanding cache misses by Andrew Glew [8], in order to persuade people to do research that helps to exploit MLP. Subsequently, numerous previous studies investigated microarchitectures to enhance MLP from on-chip (processor) side, such as MLP-aware cache replacement [24], MLP-aware prefetcher [7] and runahead execution [5]. In these studies, due to limited number of cores and limited parallelism resource (such as Instruction Window, MSHR) on chip, processor was the bottleneck to exploit MLP. But this is not right for multi-core/many-core architecture any more, with the rapid increasing number of cores, the shared memory system suffers heavy pressure to service requests form all cores. Thus for multi-core/many-core architecture, the memory system has become the main bottleneck to exploit MLP due to its relative slowly increasing parallelism

resource (channel, rank, and bank). The “Memory Wall” problem under multi-core/many-core architecture becomes more and more serious.

To moderate “Memory Wall” problem, contemporary servers would adopt high memory configuration, which could provide high memory bandwidth and MLP. For example, POWER7 processor integrates 8 cores with 4 threads each and two 4-channel DDR3 memory controllers, which could provide as high as 100GBps memory bandwidth if all the channels are fully exploited. However, due to cost and power budget limitation, a large part of servers were not configured with full memory DIMMs/channels exploited. A statistical data from a server vendor company showed that, during 2011, the most popular server sold is configured with 2 sockets, each socket has 6 cores with 6 memory DIMM slots, among these, only 15% of the servers are sold with full DIMMs exploited, 50% with half DIMMs exploited, and 35% with less than half DIMMs exploited. Thus, with the budget limitation, it becomes more important to improve memory bandwidth efficiency and exploit MLP to shorten the memory wall gap.

However, from the DRAM memory system side, because each DRAM bank is integrated with only one 4KB or 8KB row-buffer (or sense-amplifier) which holds data from a 100xMB DRAM array, the Row Buffer Hit Rate (RBHR) is considered as a key factor of the performance of DRAM system. Most researchers have made significant contributions on reducing RBHR by memory access scheduling [27], address mapping [33] and so on. During the past years, the prevalence of multi-core/many-core architecture poses new challenges of performance, power and QoS to DRAM system which is shared by all cores.

Recent studies have proposed a number of solutions to address the challenges caused by multi-core system [6, 18, 20, 21, 23, 25, 26, 37]. Nevertheless, most of the studies were motivated by improving row buffer hit rate (RBHT) of DRAM system, without considering the method of enhancing MLP for DRAM system. Although previous studies show the strong correlation between RBHR and the performance of DRAM system, we find that MLP has even a stronger correlation with the performance of DRAM system than RBHR in multi-core/many-core architecture. Experimental results show that the average IPC of homogeneous multi-program workloads on a 16-core system improves by 1.95X when incrementing the DRAM banks from 8 to 32, whereas the RBHRs are almost the same for the 8-bank and 32-bank DRAM configurations, while the

bank-level MLP increases by about 2.00X (please refer to section II for detail).

On multi-core/many-core architecture, each core can generate an independent memory request stream. Memory controllers are responsible for scheduling the requests to the available DRAM banks. If there are no available banks, the requests have to queue in the request buffer of memory controller. A recent study has shown that the queuing-delay has become the dominant portion of one memory request's access latency on multicore system [31]. Given the provision of more available banks, more memory requests would be scheduled, which means there is still a large amount of unexploited MLP due to memory limited parallel resource in multi-core/many-core system.

There are several approaches to enhance MLP. As mentioned above, simply incrementing bank count within a DRAM chip is a straightforward method to enhance MLP. There are other methods to enhance MLP, such as using multiple memory controllers, splitting bank into sub-banks and incrementing the number of row buffers. For example, using multiple memory controllers is a widely used method to improve DRAM performance on multicore architecture, but the number of memory controllers is not scalable due to the limited chip pin count. Udipi et al. [31] proposed two new DRAM organizations which contain a large number of sub-banks (or sub-arrays) and show performance improvement by 54%. However, these aggressive approaches substantially change the DRAM organization, thereby cause significant re-design cost and high risks.

In this paper, we have evaluated four representative straightforward approaches to leverage MLP for DRAM system on multi-core architecture: 1) multi-bank, 2) multi-row-buffer¹, 3) multi-memory-controller and 4) Virtual Channel Memory (VCM). VCM is selected because it represents a method of providing additional cache on the DRAM chip. Furthermore, alike multi-memory-controller, VCM is a mature technology because it possessed a certain market after it was first introduced by NEC corporation in late 1990s [22], but it faded away later. (For more details, please refer to section III)

In order to exploit MLP, VCM might be an ideal alternative to traditional DRAM chip on multi-core/many-core architecture with regard to performance, QoS, power and area overheads and even design cost and risky. Unfortunately, VCM had been obsolete before we entered the multi-core/many-core era. Therefore, we suggest the multi-core/many-core vendors reconsider the VCM technology for multi-core/many-core architecture.

Overall, we have made the following contributions:

- We find that MLP has a stronger correlation with the performance of DRAM system on multi-core/many-core architecture than RBHR which is

¹ It should be noted that, since scheduling multiple outstanding memory requests to one bank requires significant changes to DRAM state transition diagram, parameters, and specification, we just simply increment row buffers but do not allow multiple outstanding memory requests. This approach is selected because it can also improve performance by increasing RBHR.

considered as the mainly inherent metric to measure the performance of traditional DRAM system.

- In order to leverage the unexploited MLP existing in multi-core/many-core system, we have selected and evaluated four representative approaches' characteristics in term of performance, QoS, power overhead and area cost.
- According to the experimental results, we find that the obsolete VCM exhibits better than the other three approaches. We argue that memory chip vendors could reconsider the VCM technology for multi-core/many-core architecture.

The rest of the paper is organized as follows. In Section II, we introduce the observations of MLP and RBHR on multi-core/many-core architecture. In Section III, we present our evaluation scheme. We describe the experimental setup in Section IV and demonstrate experimental results and discussion in Section V. Related work and conclusion are in Section VI and Section VII respectively.

II. BACKGROUND AND MOTIVATION

A. Memory System

Figure 1 illustrates the organization of DRAM system. Contemporary multi-core/many-core processors often integrate one or more memory controllers on the chip. Each memory controller consists of 1~3 memory channels. Since adding memory channels essentially has the same effect of adding memory controllers, we use one-channel per memory controller and one-rank per channel in this paper for simplicity. So each memory controller manages multiple (usually eight) DRAM banks which can independently process multiple outstanding memory requests in parallel. Each bank is organized as a two-dimensional array of DRAM cells, consisting of multiple rows and columns. These cells are thus accessed using a DRAM address of $\langle \text{bank}, \text{row}, \text{column} \rangle$ fields, but only one row in a bank can be accessed at any given time. This row requires being stored in the row-buffer (or sense amplifier) before it could be read or written. Each bank of modern DRAM chip has only one row buffer whose size is typically 4-16KB.

If one memory request misses in the row buffer, a row buffer conflict occurs. Then the memory controller issues a *PRECHARGE* command to update the row in the row buffer back into the memory array, and then issues an *ACTIVE*

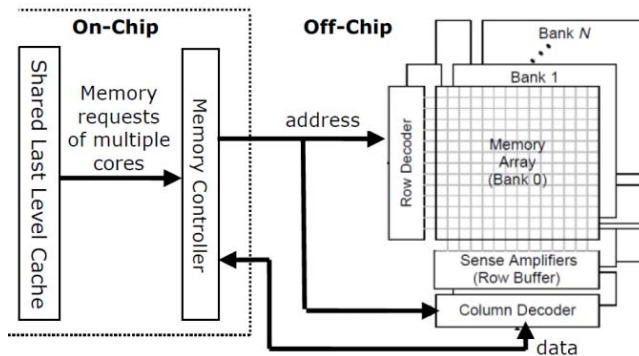


Figure 1. DRAM System Organization.

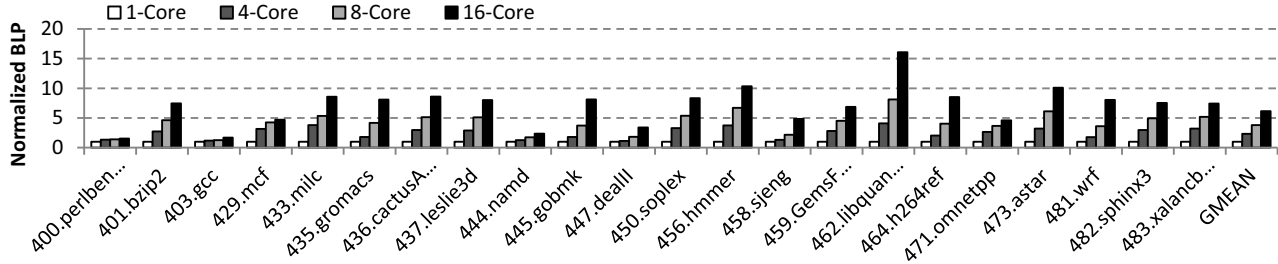


Figure 2. The Normalized BLP (Bank Level Parallelism) trend on a 32-bank memory system as the number of cores is varied from 1 to 16, where the baseline is 1-core.

command to fetch a new row into the row buffer. Therefore, the row buffer conflict causes significant memory access delay, and degrades system performance. During the past decade, numerous studies have investigated on how to improve Row Buffer Hit Rate (RBHR).

On the other hand, it is vital to keep as many banks busy as possible to improve the performance of DRAM system. This is an intuitive method for exploiting MLP. For traditional DRAM system, the maximum MLP is limited to the bank count. There is a notion called DRAM Bank-Level Parallelism (BLP) [20] which indicates the number of multiple requests being served in parallel in different DRAM banks.

B. MLP on Multicore Architecture

Previous studies have shown the strong correlation between RBHR and DRAM system performance and have proposed a number of approaches to improve RBHR during the past decade. For example, Rixner et al. [27] proposed FR-FCFS scheduling scheme which prioritized those memory requests hitting in row buffer. Recent studies investigated the RBHR on multicore architecture. Udipi et al. [31] illustrated that RBHR decreases significantly from 1 core (over 60%) to 16 cores (35%) mainly due to the row buffer conflicts caused by memory requests from different cores interfering with each other. Sudan et al. [30] also observed the same RBHR trend in their work.

In this paper, we have investigated both MLP and RBHR characterization for multi-core/many-core architecture on our simulation platform. Figure 2 shows the MLP trend in term of Normalized BLP on a 32-bank² memory system as the number of cores is varied from 1 to 16, where the baseline is 1-core system. Here we ran homogeneous multi-program workloads with each core ran the same program from SPEC CPU2006 benchmark.

We can see that except for 400.perlbench, 403.gcc, 429.mcf, 444.namd, 447.deall, 458.sjeng, and 471.omnetpp formed workloads, all the other workloads has the normalized BLP larger than 5 with 16-core, and the geometric mean of all the workloads is 6.16 with 16-core. We can also see that the geometric mean of BLP is increasing proportionally to the number of cores (or threads) increasing. The result shows that with the number of cores

increasing, the demand of MLP increasing proportionally, which would put a heavy pressure on traditional DRAM system. The least increasing of normalized BLP is 400.perlbench, it is only 1.54 with 16-core. That is because it is a memory non-intensive (the Last Level Cache MPKI is only 0.04 on 1-core) program, even with 16-core running in parallel, it still fails to exploit MLP due to its rare memory requests. But for 462.libquantum, which has the most increasing rate of normalized BLP, achieves 16.07 in 16-core, that is because it is memory-intensive and having quite good memory locality (the RBHR of it is 91.50% in 1-core). In our simulation, we adopt the bank-interleave address mapping scheme for exploiting BLP, which means we map the least bits of cache block address for bank identity. The contiguous memory accesses are mapped interleaved into multiple banks (thus exploit BLP).

Multi-core/many-core architecture poses not only the negative problems (e.g., the memory contention and unfairness problem) but also exposes large amount of MLP which is the aggregation of multiple independent memory request streams generated by multiple cores. Incrementing the bank count is a straightforward approach to exploit MLP. Figure 3 illustrates that on a 16-core system, the Normalized Rate of RBHR (Row Buffer Hit Rate), BLP and IPC with 32-bank memory system, where the baseline is 8-bank setup. The 32-bank setup can exploit BLP nearly 2 times more than the 8-bank setup, thereby improve overall system performance by nearly 1.94 in term of normalized IPC. We can also see that the more BLP exploited the more IPC speedup achieved. On the other side, the Normalized RBHR of 32-bank setup is almost equal with 8-bank setup. For some workloads (such as 403.gcc, 459.GemsFDTD, 471.omnetpp), the normalized RBHR even decreased for 32-bank compared with 8-bank. The most amount of decreased workload is 471.omnetpp, the normalized RBHR is only 0.50X of the 8-bank setup. The probable reason is that with bank-interleave address mapping, the memory requests from 16-core mixed and interfered with each other, thus further decreased row buffer hits. But for 471.omnetpp, the normalized BLP speedup achieved at 2.80, which could brought the improvement of normalized IPC by 2.35 even with worse RBHR. Based on these observations, we can conclude that MLP has a stronger correlation with the performance of DRAM system than RBHR on future multi-core/many-core architecture. Leveraging MLP could effectively improve system performance.

² In DDR3, each rank can only be configured with 8 banks. Here we implement the 32-bank configuration as 4 channels, 1 rank per channel, and 8 banks per rank.

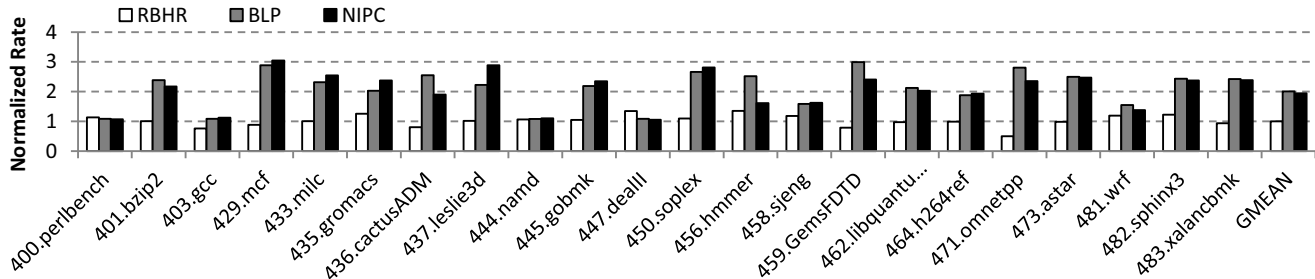


Figure 3. The Normalized Rate of RBHR (Row Buffer Hit Rate), BLP and IPC with 32-bank memory on a 16-core system, where the baseline is 8-bank memory on a 16-core system.

However, since the memory chip vendors focus on cost-per-bit and device density, they would not like to increase the number of bank because adding more banks means requiring more resources for additional sets of row decoders, sense amplifiers and column muxes etc. Actually, the number of banks integrated in a DRAM chip did not change too much in the past decade, from 4 banks in DDR SDRAM to 8 banks in DDR3 SDRAM [10]. As the number of cores increases more and more, the limited bank count leads to a large amount of unexploited MLP.

III. LEVERAGING MLP

There are two design philosophies for leveraging MLP. One is “design from scratch”, which means substantially changing DRAM organization. Several recent studies, e.g., Selective Bitline Activation (SBA) and Single Subarray Access (SSA) [31], have proposed new DRAM organizations to improve DRAM system performance by enhancing MLP. However, although this design philosophy might bring significant changes, it also might lead to unpredictable outcome and high risks. Another design philosophy is “keep it simple and stupid (KISS)”, which means looking for approaches that are either already existing or combinations of existing technologies.

In this paper, we adopt the KISS design philosophy to investigate how to leverage MLP. We select four straightforward approaches:

Multi-Row-Buffer: To keep the design simple, we only increment the row buffer count but do not change the DRAM state transition diagram, parameter and specifications. Therefore, although those row buffers hold multiple opened rows, only one row is accessible at any time according to the DRAM specification. Actually, this approach can improve RBHR other than enhance MLP, so it can be used to compare the effectiveness of improving RBHR and enhancing MLP.

Multi-Bank: Given a capacity-fixed DRAM chip, we split it into different number of banks. This approach requires additional resources for memory controller (multiple control logic modules) and DRAM chip (address decoders and row buffers etc.), but it does not need to change the DRAM specification.

Multi-Memory-Controller: We increase the number of on-chip memory controllers. This approach inherently increases the number of banks and should have the same

effect as multi-bank. However, it consumes on-chip resources, especially the pin count.

Virtual-Channel-Memory (VCM): VCM represents a method of providing additional cache on the DRAM chip. In each rank, there are 16~32 channel buffers, each holding one segment of row buffer. The DRAM specification is slightly changed to support operating channel buffers, but VCM memory controller is compatible to traditional DRAM. Furthermore, VCM is a mature technology and ever possessed a certain market around 2000.

Since VCM requires changes to DRAM specification, we would like to describe it in details. VCM was first introduced by NEC corporation in late 1990s [22]. It was intended for a wide range of applications such as multimedia and web servers. VCM puts a set of fast channel buffers within memory chips and the number of channel buffers is usually 4 or 8 times more than that of banks. Hence, VCM is expected to provide faster access as well as more concurrency.

VCM Organization: Figure 4 illustrates VCM’s conceptual organization. Channel buffers are introduced as an extra storage layer between memory controller and DRAM banks. Two commands, i.e., *PREFETCH* (reading segment data from row buffer to channel) and *RESTORE* (writing segment data from channel to row buffer), are also introduced in order to operate channels. As shown in the figure, each row buffer is divided into 4-16 segments which are transfer units between banks and channels. Memory operations (commands) are divided into **foreground** operations for channels (*READ* and *WRITE* commands) and **background** operations for DRAM banks (*ACTIVE*, *PRECHARGE*, *PREFETCH* and *RESTORE* commands). NEC’s VCM is implemented to be compatible to the industry standard SDRAM and uses the same command protocol and interface as SDRAM/DRAM. Because channels and DRAM banks are independent, **foreground** operations and **background** operations can also be executed independently. To further enhance VCM’s performance, channels and the row buffers (banks) of the original VCM are fully associative, which means that any channels can store segment data from any row buffers (banks) and can be written back to any row buffers (banks). Because channels and DRAM banks are independent, memory controller can schedule the **foreground** and **background** operations concurrently, which allows the memory controller to exploit as more MLP as possible.

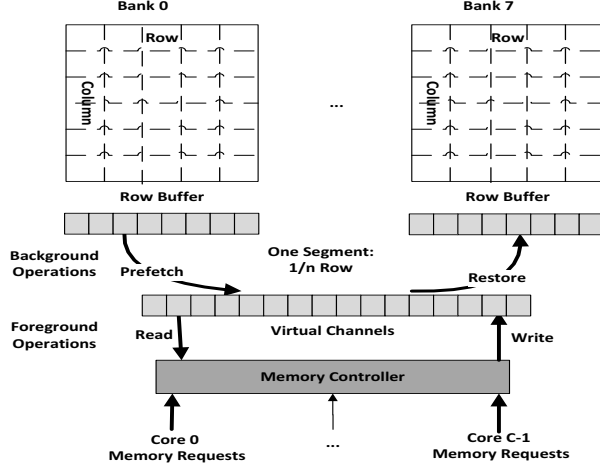


Figure 4. Conceptual Organization of VCM.

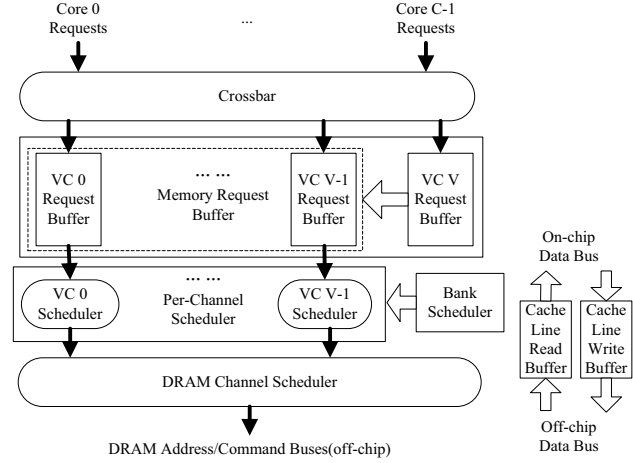


Figure 5. VCM Memory Controller.

VCM for Multicore Architecture: We have implemented VCM on our simulation platform. Figure 5 illustrates the design of VCM memory controller which is responsible for scheduling memory requests from different cores. There are three main distinctions between a VCM memory controller and a traditional DRAM controller:

- 1) The channels are decoupled with bank, which means that the channels can hold any segment data from any bank/row. Thus each channel needs a tag register to record memory address of the segment it holds (i.e., $\langle \text{bank}, \text{row}, \text{segment} \rangle$).
- 2) There are channel-level schedulers to schedule VCM commands to channels. Those schedulers need to keep track of their corresponding channels' state to determine which command should be issued.
- 3) There are bank-level schedulers for keeping track of the state of the banks. The bank state information and the channel state information are used to select requests in request buffer.

NEC's VCM usually has 8 to 32 channels. In fact, channel buffers operate like a fully associative cache with write through policy. Take reading data from a memory bank as an example. Memory controller first looks up registers to determine whether the required data is already in a channel; if the read request misses in all channels, the bank-level controller will issue an *Active* command to the corresponding bank to fetch a row into the row buffer; after 2 (tPAD) cycles, a *Prefetch* command issued by a channel-level scheduler follows to transfer one segment of the row buffer to a channel which memory controller selects for replacing; then 2 (tPCD) cycles later, the *Read* command can be issued to the new refilled channel and after another 2 cycles (read latency) the data can be present in the data bus; meanwhile, if using a close-page policy, a *Precharge* command is issued to the bank to close the row; for an 8 burst length, the data transfer requires 4 cycles with double data rate (DDR). As for writing data into a channel buffer, *Write*, *Restore* and *Active* commands are issued consecutively, and finally a *Precharge* command is used to

close the corresponding row. According to the datasheet [22], we conclude that the original NEC's VCMs adopt write-through policy for channel buffers and close-page policy for row buffer.

For more details of VCM, please refer to [26]³.

IV. EXPERIMENTAL SETUP

A. Evaluation Tools

We implement and evaluate the four approaches using an in house cycle-accurate x86 CMP simulator. The functional front-end is based on Pin [15] and iDNA [4]. We model the memory system in detail, faithfully capturing bandwidth limitations, contention, and enforcing bank/channel/bus conflicts. Table I shows the major DRAM and processor parameters. We model a modest multi-core (16 core) system with one channel as baseline to produce heavy access pressure on memory system to simulate the environment in future multi-core/many-core system, meanwhile limiting simulation time. The VCM was implemented in the memory system, the baseline configuration is 32 Virtual Channels within each DRAM chip. We use CACTI 6.5[1] to evaluate area and power parameters.

B. Workloads

We use the SPEC CPU2006 benchmarks for evaluation. We compile each benchmark using gcc 4.1.2 with -O3 optimizations and choose a representative simulation phase using PinPoints [15]. We select memory intensive benchmarks and memory non-intensive benchmarks from SPEC CPU2006. Table II and Table III list their characteristics (including IPC, MPKI, RBHR and BLP). We run multiple programs on multicore system where each core is dedicated to one program. We use homogeneous workloads (multiple instances of the same program) to evaluate performance and heterogeneous workloads (the combinations of different programs) to evaluate QoS. All programs are run with their reference (maximum size) input

³ Since the original VCM datasheet of NEC [22] has been outdated on the Internet.

TABLE I. SIMULATED TARGET SYSTEM CONFIGURATION

Processor Pipeline	3 GHz processor, 128-entry instruction window (64-entry issue queue, 64-entry store queue), 12-stage pipeline
Fetch/Exec/Commit width	3 instructions per cycle in each core; only 1 can be a memory operation
L1 Caches	32 K-byte per-core, 4-way set associative, 64-byte block size, 2-cycle latency
L2 Caches	512 K-byte per core, 8-way set associative, 64-byte block size, 12-cycle latency, 32 MSHRs
Baseline DRAM controller (on-chip)	FR-FCFS, close-page row buffer policy for VCM, open-page for other approaches; 128-entry request buffer, 64-entry write data buffer, reads prioritized over writes, XOR-based address-to-bank mapping [33].
DRAM chip parameters	Micron DDR3-1600 timing parameters , tCL=15ns, tRCD=15ns, tRP =15ns; 8 banks, 8K-byte row-buffer per bank
DIMM configuration	Single-rank, 8 DRAM chips put together on a DIMM (dual in-line memory module) to provide a 64-bit wide channel to DRAM
Round-trip L2 miss latency with VCM	For a 64-byte cache line, row-buffer hit: 200 cycles, closed: 300 cycles, conflict: 400 cycles. VCM additional Latency: Active-to-Prefetch Latency (100 cycles), Prefetch-to-Read/Write Latency (80 cycles), Read Latency (200 cycles)
VCM parameters	1K-byte per segment, 8 segments per row-buffer, 8 requests per VC request buffer and 128 requests for the Recycle Request Buffer

TABLE II. MEMORY INTENSIVE BENCHMARK CHARACTERISTICS

Benchmark	Type	IPC	MPKI	RBHR	BLP
401.bzip2	INT	1.07	3.42	80.93	1.58
429.mcf	INT	0.18	71.86	10.82	5.05
433.milc	FP	0.27	20.59	86.54	1.29
436.cactusADM	FP	0.45	6.31	19.70	1.32
437.leslie3d	FP	0.35	17.10	71.04	1.71
450.soplex	FP	0.19	38.99	88.27	1.71
456.hmmer	INT	0.77	5.19	49.76	1.29
459.GemsFDTD	FP	2.41	3.78	52.70	2.87
462.libquantum	INT	1.25	6.90	91.50	1.03
464.h264ref	INT	3.28	2.39	85.20	1.12
470.lbm	FP	0.56	35.26	85.85	3.31
471.omnetpp	INT	2.10	5.85	62.70	3.56
473.astar	INT	1.18	6.39	51.80	1.47
481.wrf	FP	3.07	1.98	71.39	1.03
482.sphinx3	FP	2.98	2.47	81.97	1.86
483.xalancbmk	INT	2.22	3.70	68.68	2.05

sets. For multi-program workloads: we fast forward 20 million instructions for each process to warm up the simulator, and then execute another 100 million instructions for each core, and then collect simulation data such as IPC, memory access latency and power consumption.

C. Metrics

We evaluate the benchmark by four main metrics, i.e., performance, QoS, power consumption and area cost. For individual benchmark, we evaluate performance in term of IPC. For whole system, we use the Unfairness metric to estimate QoS and evaluate performance in term of System_Throughput [12][29]:

$$\text{Slowdown}_i = \frac{IPC_i^{\text{shared}}}{IPC_i^{\text{alone}}}, \text{ Unfairness} = \frac{\text{Max}\{\text{Slowdown}_i\}}{\text{Min}\{\text{Slowdown}_i\}}$$

$$\text{System_Throughput} = \sum_i \{\text{Slowdown}_i\}$$

D. Experimental schemes

We adopt three experimental schemes: 1) we use homogeneous workloads to evaluate system performance in term of average IPC. 2) We use heterogeneous workloads to evaluate QoS effect of the approaches in term of Unfairness and System_Throughput. 3) We demonstrate the area cost and power consumption by the CACTI tool.

TABLE III. MEMORY NON-INTENSIVE BENCHMARK CHARACTERISTICS

Benchmark	Type	IPC	MPKI	RBHR	BLP
400.perlbench	INT	2.05	0.04	68.54	1.33
403.gcc	INT	1.73	0.22	62.12	1.65
435.gromacs	FP	1.55	0.95	80.85	1.43
444.namd	FP	2.44	0.05	91.51	1.05
445.gobmk	INT	1.82	0.60	58.92	1.35
447.dealII	FP	1.85	0.08	83.85	1.17
453.povray	FP	1.88	0.00	88.54	1.58
454.calculix	FP	1.73	0.01	84.70	1.17
458.sjeng	INT	1.94	0.43	24.20	1.54
465.tonto	FP	2.10	0.16	10.85	1.73

Here, **IPC**: Instruction per Cycle, **MPKI**: L2 Cache Misses per 1000 Instructions, **RBHR**: Row-Buffer Hit Rate, **BLP**: Bank Level Parallelism.

V. EXPERIMENTAL RESULTS

A. Performance

Figure 6 shows the normalized IPC speedup of memory-intensive homogenous workloads running on a baseline multicore system which has 16-core and 8 memory banks. We apply the four approaches to the baseline system and measure their IPC speedups.

For the multi-row-buffer approach, even using eight 8-KB row buffers in one bank, the GMEAN performance improvement is only 1.22X. As mentioned above, this approach can improve RBHR but is still limited to bank-level parallelism. Figure 7 further illustrates the speedups on systems with different number of cores. For the 1-core system, incrementing row-buffer can achieve almost the same improvement as the other approaches. However, as the number of cores increases, its performance scalability is very poor, compared with other approaches. The multi-bank and the multi-memory-controller approaches exhibit good IPC speedups. According to Figure 6, 4 memory controllers exhibit the best improvement, by 2.44X while 2 memory controllers setup improves performance by 1.63X. The 32-bank scheme and the 16-bank scheme exhibit performance improvement by 2.44X and 1.63X respectively, and 32-bank can achieve nearly the same improvement as four memory controllers.

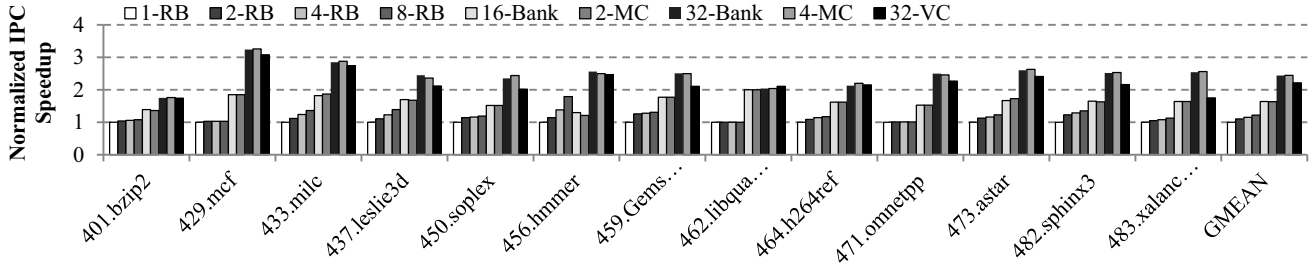


Figure 6. The Normalized IPC Speedup on 16-core system, the baseline is 8-bank with 1-row buffer. Here, RB: Row Buffer, MC: Memory Controller.

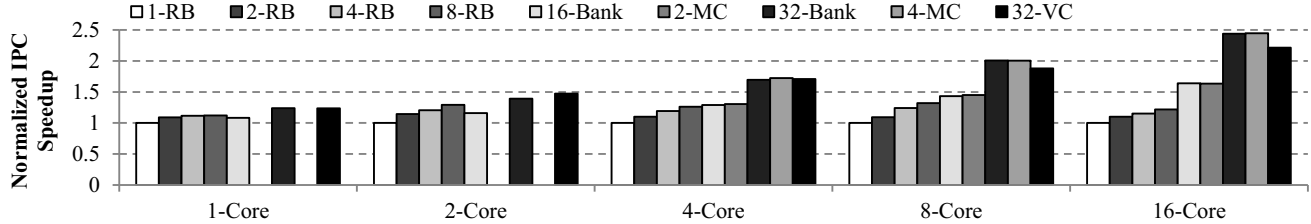


Figure 7. Performance scalability against the increasing core numbers, in terms of Normalized IPC.

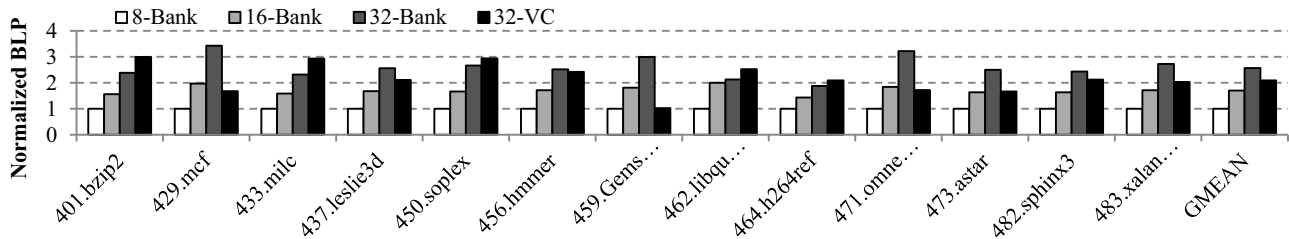


Figure 8. The Normalized BLP on a 16-core system, the baseline is 8-Bank.

On average, VCM with 32 1KB-channel buffers achieves 2.21X performance speedup and also exhibits good scalability, from 1.24X for single-core to 2.21X for 16-core. In fact, an interesting phenomenon is that its RBHR⁴ is much less than the multi-row-buffer approach. However, according to Figure 8, the MLP (in term of BLP here) improves significantly. Even on a 16-core system with only 8 banks, the VCM with 32-channel can exploit normalized MLP nearly 2.09x compared with 8-Bank. The VCM’s capability of exploiting MLP benefits from its organization that channels and DRAM banks are independent so that memory controller can schedule the bank operations and the channel operations concurrently. This phenomenon is strong evidence that the performance on multicore architecture is more dependent of MLP than RBHR. We further investigate how the latter three approaches exploit more MLP. Figure 9 illustrates the average access latency of memory requests. The latency is divided into two parts: DRAM accessing latency and queuing latency. In fact, the DRAM accessing latencies are almost the same for all approaches, however, the queuing delay at memory controller side reduces significantly for the VCM approach and the multi-bank (the same as multi-MC) approach. Take VCM as an example, for

each individual request, although the latency of background operations (i.e., accessing DRAM latency) increases slightly by 7.3%, the queuing latency substantially decreases, from 8.75x (of baseline) to 3.50x for 32-channel VCM.

B. QoS

The memory QoS problem means that unfair servicing of different cores’ requests by the memory controller can lead to application/core starvation [21][23] and even denial of memory service for some cores [19]. For example, Mutlu et al. [21] showed that the slowdowns for some memory non-intensive applications can increase from 7.74X for 4-core system to 11.35X for 8-core system whereas the memory-intensive application experienced the slowdowns of only 1.04X and 1.09X respectively. Some other works [18-20, 25] also demonstrated the similar unfair phenomenon. Recent studies on memory controller optimization, such as PAR-BS [20] and ATLAS [13], have shown their effectiveness in improving QoS.

According to our above analysis, VCM is better than the other three approaches for performance. Thus in this section, we mainly evaluate QoS for VCM. Moreover, we have integrated the PAR-BS and ATLAS scheme into VCM. We use the following metrics to evaluate QoS: *Unfairness* and *System_Throughput*.

⁴For VCM, RBHR means the virtual channel hit rate rather than the traditional row buffer hit rate for banks.

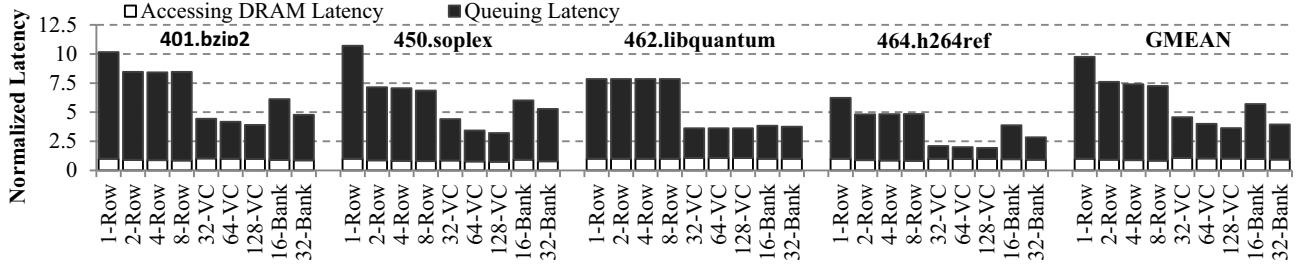


Figure 9. Breakdown of memory access latency on a 16-core system, where the baseline is Accessing DRAM Latency of 1-Row.

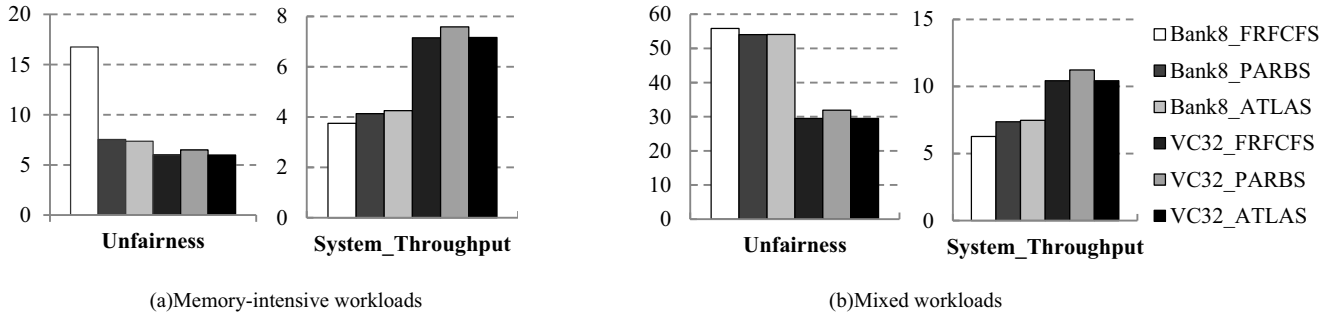


Figure 10. The Unfairness and System Throughput on 16-core system with and without VCM: (a) memory-intensive workloads; (b) mixed memory-intensive and memory-non-intensive workloads.

Figure 10(a) shows the average *Unfairness* and *System_Throughput* of 16 memory-intensive workloads on a 16-core 8-bank system with and without VCM. Each workload includes 16 memory-intensive programs which are generated by randomly selected from Table II. We can see that without VCM, the two state of the art scheduling algorithm could efficiently improve system QoS by reducing unfairness from 16.74 to 7.53 (reduce 55.02%) for PARBS and 7.36 (reduce 56.03%) for ATLAS. Meanwhile they can improve *System_Throughput* by 10.13% for PARBS (4.13 vs. 3.75) and 13.33% for ATLAS (4.25 vs. 3.75). On the other hand, VCM exhibits good performance both in QoS and *System_Throughput* even with simple memory schedule algorithm (FR-FCFS). It can reduce unfairness by 64.16% (6.00 vs. 16.74) and significantly improve system throughput by 90.67% (7.15 vs. 3.75), both of these are even better than the two state of the art memory scheduling algorithm without VCM. There are two reasons: 1) the fewer memory requests stay in queue, the less is the probability of unfairness. As shown in Figure 9, VCM can significantly reduce queuing latency, so it can also substantially eliminates unfairness; 2) Channels are independent of banks and they can operate simultaneously, hence more memory requests can be issued without the limitation of bank count. Therefore, those requests with low priority have more chances to be scheduled. Although PAR-BS and ATLAS algorithms can largely reduce unfairness by ranking threads' requests, they still subject to the bank-level parallelism. When applying PAR-BS and ATLAS to VCM, the Unfairness metrics change slightly to 6.49 (+8.17%) and 5.99 (-0.17%) respectively, and that is nearly the same for *System_Throughput* by 6.01% and 0.14%. Thus we can conclude that VCM could perform well with simple

FRFCFS memory scheduling algorithm for memory-intensive workloads, further, figure 10(b) could also prove this view for mixed memory-intensive and memory-non-intensive workloads.

Figure 10(b) shows that for mixed memory-intensive and memory-non-intensive workloads, without VCM the PARBS and ATLAS reduced unfairness quite slightly (3.24% and 3.13%), but the VCM with FRFCFS could also effectively reduce unfairness by 47.18% (29.49 vs. 55.83) and improve system throughput by 66.19% (6.27 vs. 10.42). These results show that VCM could be suitable for various workloads in multi-core/many-core architecture.

C. Area and Power Cost

We use CACTI to estimate area and power cost for the four approaches. Table IV illustrates the area and energy parameters for 2GB DRAM with 32nm technology. Given a fixed DRAM size, the area is dependent on the layout of those components. For multi-bank, the area increases significantly mainly because of the area of large row buffers. It should be noted that when the DRAM chip is divided into more banks, the DRAM arrays become smaller so that the CACTI can figure out a smaller area (e.g., area of 128-bank is smaller than area of 64-bank). The multi-row-buffer approach and the multi-memory-controller approach also suffer from the area penalty. In particular, multi-memory-controller has the biggest area cost because adding one memory controller means adding a set of banks as well as consuming chip pins which have been scarce resources in multicore systems. However, VCM almost has no area cost, increasing area by only 0.5% for 32 1KB-channels and 7.7% for 128 channels.

TABLE IV. AREA AND ENERGY PARAMETERS

con-fig	Area (mm ²)	High Performance Mode			Low Power Standby Mode			Con-fig	Area (mm ²)	High Performance Mode			Low Power Standby Mode		
		DRPA	LP	AE	DRPA	LP	AE			DRPA	LP	AE	DRPA	LP	AE
Multiple Banks						Multiple Row Buffers									
8	617.4	2.13	3660.1	0.13	3.39	30.87	0.20	1	617.4	2.13	3660.1	0.13	3.39	30.87	0.20
16	928.4	2.43	4832.3	0.15	3.54	31.35	0.22	2	662.2	2.19	3798.5	0.14	3.49	31.52	0.21
32	718.4	2.27	4330.6	0.14	3.62	34.86	0.22	4	736.6	2.29	3812.5	0.14	3.64	31.60	0.22
64	1093.6	2.59	6271.7	0.16	3.70	29.49	0.22	8	886.7	2.48	3841.3	0.15	3.96	31.77	0.23
128	895.4	2.50	6363.4	0.15	4.00	47.00	0.23	16	1189.7	2.86	3931.0	0.17	4.43	21.36	0.26
Multiple Memory Controllers						Virtual Channel Memory									
1	617.4	2.13	3660.1	0.13	3.39	30.87	0.20	16	618.2	2.19	3660.9	0.17	3.45	30.88	0.25
2	1234.7	4.26	7320.1	0.26	6.78	61.74	0.40	32	620.2	2.21	3661.2	0.18	3.47	30.88	0.26
4	2469.5	8.52	14640.	0.52	13.56	123.4	0.80	64	629.2	2.27	3663.2	0.21	3.50	30.89	0.28

Here, DRPA: Dynamic Read per Access (nJ); LP: Leakage Power (mW); AE: Active Energy (nJ).

We evaluate power consumption in both high performance (HP) mode and low power standby (LPS) mode. For individual DRAM operations, experimental results show that there are no significant differences among different configurations. For 32-channel VCM in HP mode, each read operation consumes about 2.19 nJ, increasing slightly by about 3.6% compared to the baseline 8-bank DRAM.

VI. RELATED WORKS

Performance Issue: Rixner et al. proposed FR-FCFS [27] scheduling algorithm which could improve memory bandwidth by 40%~93% for streaming applications. McKee et al. [16] showed that dynamically reordering memory requests can increase the row buffer hit rate for scientific and multimedia applications. Hur et al. [9] proposed the adaptive history-based (AHB) scheduler which improves performance by 7.6%~15.6%. There is still a lot of work [28, 37] focusing on reordering memory requests. However, most of these studies are aim to take full advantage of the row buffer within each bank, still being limited to bank-level parallelism. Agrawal et al. [2] proposed virtually pipelined memory (VPM) which provided a deeper pipeline for handling memory requests. Although VPM might increase access individual request latency, it is able to improve effective memory bandwidth. This is because that pipelining allows more on-the-fly memory requests. Nevertheless, these schemes do not consider enhancing MLP.

VCM represents a number of approaches which add additional buffer into DRAM chip to exploit MLP. The similar idea was first proposed by Alexander and Kedem [3]. They proposed to integrate some small buffers in DRAM chip as prediction table for a DRAM prefetching scheme. Zhao et al. [34] implemented an additional off chip SRAM Cache to exploit locality for larger workload with higher bandwidth. Jiang et al.[11] further improved the DRAM Cache with some effective filters to cache only hot pages.

Several studies focus on changing DRAM's organization in order to improve performance as well as reduce power. Zheng et al. proposed Mini-rank [35] and decoupled-DIMM [36] to address DRAM power issue. Yamauchi et al. [32] present hierarchical multi-bank DRAM composed of 8 sub-

banks, improving performance about 65%. More recently, Udipi et al. [31] argued that traditional DRAM system has already not suitable for multicore system because of overfetch problem which wastes significant energy. They proposed two main aggressive schemes including Selective Bit-line Activation (SBA) and Single Subarray Access (SSA) which are able to reduce energy by 5X~6X and performance by 54% due to reduced queuing delays.

QoS Issue: Both of Rafique et al. [25] and Nesbit et al. [23] studied QoS based on fair queuing mechanism. Mutlu et al. proposed STFM [21] and PAR-BS [20] to solve unfairness problems for multicore system. Furthermore, Kim et al. [13] proposed ATLAS algorithm for multiple memory controllers system which could improve system throughput by 8.4%~10.8%. Eiman et al. proposed FST algorithm [6] which throttles down cores causing unfairness by limiting the number of their available MSHRs.

VII. CONCLUSIONS

In this paper, we have found that MLP has a stronger correlation with the performance of DRAM system on multi-core/many-core architecture than RBHR which is considered as the only inherent metric to measure the performance of DRAM system in the past decade. In order to leverage the unexploited MLP existing in multi-core/many-core system, we have selected and evaluated four representative approaches by measuring performance, QoS, power overhead and area cost. According to the experimental results, we have found that the obsolete VCM exhibits better than the other three approaches. We argue that memory chip vendors could reconsider the VCM technology for multi-core/many-core architecture.

ACKNOWLEDGMENT

The authors thank the anonymous reviewers for their constructive suggestions. This research is supported by the National Natural Science Foundation of China (NSFC) under grant numbers 60925009, 60921002, 60903046, 61272134, 61033009 and the National Basic Research Program of China (973 Program) under a grant number 2011CB302502.

REFERENCES

- [1] CACTI: An Integrated Cache and Memory Access Time, Cycle Time, Area, Leakage, and Dynamic Power Model. <http://www.hpl.hp.com/research/cacti/>.
- [2] B. Agrawal and T. Sherwood, Virtually Pipelined Network Memory, in Proceedings of the 39th Annual IEEE/ACM International Symposium on Microarchitecture, 2006.
- [3] T. Alexander and G. Kedem, Distributed Prefetch-buffer/Cache Design for High Performance Memory Systems, in Proceedings of the 2nd IEEE Symposium on High-Performance Computer Architecture, 1996.
- [4] S. Bhansali, W.-K. Chen, S. D. Jong, A. Edwards, R. Murray, M. Drinic, D. Mihocka, and J. Chau, Framework for instruction-level tracing and analysis of program executions, in Proceedings of the 2nd international conference on Virtual execution environments, 2006.
- [5] Y. Chou, B. Fahs, and S. Abraham, Microarchitecture Optimizations for Exploiting Memory-Level Parallelism, in Proceedings of the 31st annual international symposium on Computer architecture, 2004.
- [6] E. Ebrahimi, L. Chang Joo, M. Onur, and N. P. Yale, Fairness via source throttling: a configurable and high-performance fairness substrate for multi-core memory systems, in Proceedings of the fifteenth Architectural support for programming languages and operating systems (ASPLOS), 2010.
- [7] S. Everman and L. Eeckhout, A Memory-Level Parallelism Aware Fetch Policy for SMT Processors, in Proceedings of the 2007 IEEE 13th International Symposium on High Performance Computer Architecture, 2007.
- [8] A. Glew. "MLP yes! ILP no!" in ASPLOS Wild and Crazy Idea Session. 1998.
- [9] I. Hur and C. Lin, Memory scheduling for modern microprocessors. *ACM Trans. Comput. Syst.*, 2007. 25(4): p. 10.
- [10] Jeduc, DDR3 SDRAM STANDARD, 2010.
- [11] X. Jiang, N. Madan, L. Zhao, M. Upton, R. Iyer, S. Makineni, D. Newell, Y. Solihin, and R. Balasubramonian, CHOP: Integrating DRAM Caches for CMP Server Platforms. *IEEE Micro*, 2010. 31(1): p. 99-108.
- [12] K. Luo, J. Gummaraju, and M. Franklin. Balancing throughput and fairness in SMT processors. in *IEEE International Symposium on Performance Analysis of Systems and Software*. 2001.
- [13] Y. Kim, D. Han, O. Mutlu, and M. Harchol-Balter. ATLAS: A Scalable and High-Performance Scheduling Algorithm for Multiple Memory Controllers. in Proceedings of the 16th International Symposium on High-Performance Computer Architecture (HPCA). 2010.
- [14] C. J. Lee, V. Narasiman, O. Mutlu, and Y. N. Patt, Improving memory bank-level parallelism in the presence of prefetching, in Proceedings of the 42nd Annual IEEE/ACM International Symposium on Microarchitecture, 2009.
- [15] C.-K. Luk, R. Cohn, R. Muth, H. Patil, A. Klauser, G. Lowney, S. Wallace, V. J. Reddi, and K. Hazelwood, Pin: building customized program analysis tools with dynamic instrumentation, in Proceedings of the 2005 ACM SIGPLAN conference on Programming language design and implementation, 2005.
- [16] S. A. Mcke, W. A. Wulf, J. H. Aylor, M. H. Salinas, R. H. Klenke, S. I. Hong, and D. A. B. Weikle, Dynamic Access Ordering for Streamed Computations. *IEEE Trans. Comput.*, 2000. 49(11): p. 1255-1271.
- [17] Micron, 1Gb DDR2 SDRAM Component: MT47H128M8HQ-25. <http://download.micron.com/pdf/datasheets/dram/ddr2/1GbDDR2.pdf>, May 2007.
- [18] T. Moscibroda and O. Mutlu, Distributed order scheduling and its application to multi-core dram controllers, in Proceedings of the twenty-seventh ACM symposium on Principles of distributed computing, 2008.
- [19] T. Moscibroda and O. Mutlu, Memory performance attacks: denial of memory service in multi-core systems, in Proceedings of 16th USENIX Security Symposium on USENIX Security Symposium, 2007.
- [20] O. Mutlu and T. Moscibroda, Parallelism-Aware Batch Scheduling: Enhancing both Performance and Fairness of Shared DRAM Systems, in Proceedings of the 35th Annual International Symposium on Computer Architecture, 2008.
- [21] O. Mutlu and T. Moscibroda, Stall-Time Fair Memory Access Scheduling for Chip Multiprocessors, in Proceedings of the 40th Annual IEEE/ACM International Symposium on Microarchitecture, 2007.
- [22] Nec, 64M-bit Virtual Channel SDRAM data sheet, 1998.
- [23] K. J. Nesbit, N. Aggarwal, J. Laudon, and J. E. Smith, Fair Queuing Memory Systems, in Proceedings of the 39th Annual IEEE/ACM International Symposium on Microarchitecture, 2006.
- [24] M. Qureshi, D. Lynch, O. Mutlu, and Y. Patt, A Case for MLP-Aware Cache Replacement, in Proceedings of the 33rd annual international symposium on Computer Architecture, 2006.
- [25] N. Rafique, W.-T. Lim, and M. Thottethodi, Effective Management of DRAM Bandwidth in Multicore Processors, in Proceedings of the 16th International Conference on Parallel Architecture and Compilation Techniques (PACT), 2007.
- [26] S. Rixner, Memory Controller Optimizations for Web Servers, in Proceedings of the 37th annual IEEE/ACM International Symposium on Microarchitecture, 2004.
- [27] S. Rixner, W. J. Dally, U. J. Kapasi, P. Mattson, and J. D. Owens. Memory Access Scheduling. in Proceedings of the 27th annual international symposium on Computer architecture. 2000.
- [28] J. Shao and B. T. Davis, A Burst Scheduling Access Reordering Mechanism, in Proceedings of the 2007 IEEE 13th International Symposium on High Performance Computer Architecture, 2007.
- [29] A. Snaveley and D. M. Tullsen, Symbiotic jobscheduling for a simultaneous multithreaded processor, in Proceedings of the ninth international conference on Architectural support for programming languages and operating systems 2000.
- [30] K. Sudan, N. Chatterjee, D. Nellans, M. Awasthi, R. Balasubramonian, and A. Davis, Micro-pages: increasing DRAM efficiency with locality-aware data placement, in Proceedings of the fifteenth Architectural support for programming languages and operating systems (ASPLOS), 2010.
- [31] A. N. Udipi, N. Muralimanohar, N. Chatterjee, R. Balasubramonian, A. Davis, and N. P. Jouppi, Rethinking DRAM design and organization for energy-constrained multi-cores, in Proceedings of the 37th annual international symposium on Computer architecture, 2010.
- [32] T. Yamauchi, L. Hammond, and K. Olukotun. The Hierarchical Multi-Bank DRAM: A High-Performance Architecture for Memory Integrated with Processors. in Proceedings of the 19th Conference on Advanced Research in VLSI. 1997.
- [33] Z. Zhang, Z. Zhu, and X. Zhang, A permutation-based page interleaving scheme to reduce row-buffer conflicts and exploit data locality, in Proceedings of the 33rd annual ACM/IEEE international symposium on Microarchitecture, 2000.
- [34] L. Zhao, R. Iyer, R. Illikkal, and D. Newell. Exploring DRAM cache architectures for CMP server platforms, *Computer Design*. in 25th International Conference on Computer Design (ICCD). 2007.
- [35] H. Zheng, J. Lin, Z. Zhang, E. Gorbato, H. David, and Z. Zhu, Mini-rank: Adaptive DRAM architecture for improving memory power efficiency, in Proceedings of the 41st annual IEEE/ACM International Symposium on Microarchitecture, 2008.
- [36] H. Zheng, J. Lin, Z. Zhang, and Z. Zhu, Decoupled DIMM: building high-bandwidth memory system using low-speed DRAM devices, in Proceedings of the 36th annual international symposium on Computer architecture, 2009.
- [37] H. Zheng, J. Lin, Z. Zhang, and Z. Zhu, Memory Access Scheduling Schemes for Systems with Multi-Core Processors, in 37th International Conference on Parallel Processing (ICPP), 2008.