



**HAL**  
open science

## Exploiting Cross-Layer Hotness Identification to Improve Flash Memory System Performance

Jinhua Cui, Weiguo Wu, Shiqiang Nie, Jianhang Huang, Zhuang Hu, Nianjun Zou, Yinfeng Wang

► **To cite this version:**

Jinhua Cui, Weiguo Wu, Shiqiang Nie, Jianhang Huang, Zhuang Hu, et al.. Exploiting Cross-Layer Hotness Identification to Improve Flash Memory System Performance. 13th IFIP International Conference on Network and Parallel Computing (NPC), Oct 2016, Xi'an, China. pp.17-28, 10.1007/978-3-319-47099-3\_2 . hal-01647990

**HAL Id: hal-01647990**

<https://inria.hal.science/hal-01647990v1>

Submitted on 24 Nov 2017

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution 4.0 International License

# Exploiting Cross-Layer Hotness Identification to Improve Flash Memory System Performance

Jinhua Cui<sup>1</sup>, Weiguo Wu<sup>1</sup>, Shiqiang Nie<sup>1</sup>,  
Jianhang Huang<sup>1</sup>, Zhuang Hu<sup>1</sup>, Nianjun Zou<sup>1</sup>, and Yinfeng Wang<sup>2</sup>

<sup>1</sup> School of Electronic and Information Engineering, Xi'an Jiaotong University, Shaanxi 710049, China. cjhnicole@gmail.com, wgwu@xjtu.edu.cn

<sup>2</sup> Department of Software Engineering, ShenZhen Institute of Information Technology, Guangdong 518172, China

**Abstract.** Flash memory has been widely deployed in modern storage systems. However, the density improvement and technology scaling would decrease its endurance and I/O performance, which motivates the search to improve flash performance and reduce cell wearing. Wearing reduction can be achieved by lowering the threshold voltages, but at the cost of slower reads. In this paper, the access hotness characteristics are exploited for read performance and endurance improvement. First, with the understanding of the reliability characteristics of flash memory, the relationship among flash cell wearing, read latency and bit error rate is introduced. Then, based on the hotness information provided by buffer management, the threshold voltages of a cell for *write-hot* data are decreased for wearing reduction, while these for *read-hot* data are increased for read latency reduction. We demonstrate analytically through simulation that the proposed technique achieves significant endurance and read performance improvements without sacrificing the write throughput performance.

**Keywords:** NAND flash memory, Endurance, Raw bit error rate, Threshold voltage, Cross-layer

## 1 Introduction

In recent years, storage devices equipped with NAND flash memory have become widely used for a multitude of applications. Due to its high density, low power consumption, excellent IOPs performance, shock-resistance and noiselessness, NAND flash-based solid state drive (SSD) is considered as an alternative to hard disk drive (HDD) as the second storage device [1]. With semiconductor process technology scaling and cell density improvement, the capacity of NAND flash memory has been increasing continuously and the price keeps dropping. However, technology scaling inevitably brings the continuous degradation of flash memory endurance and I/O performance, which motivates the search for methods to improve flash memory performance and lifetime [2] [3].

Flash lifetime, measured as the number of erasures a block can endure, is highly correlated with the raw bit error rate (RBER), which is defined as the

number of corrupted bits per number of total bits read [4]. As the endurance of flash cells is limited, RBER is expected to grow with the number of program/erase (P/E) cycles, and a page is deemed to reach its lifetime if the combined errors are not correctable by error correction code (ECC). Many methods have been proposed to maximize the number of P/E cycles in flash memory. They include enhancing the error correction capability of ECC [5] [6] [7], distributing erasure costs evenly across the drives blocks [8] [9] [10], and reducing the threshold voltages for less wearing incurred by each P/E cycling [11] [12] [13].

Flash read latency is also highly correlated with RBER. The higher the RBER, the stronger the required ECC capability, as well as the higher the complexity of ECC scheme and the slower the read speed. Recently, several works have been proposed to reduce read latency by regulating the memory sensing precision. Zhao et al. [5] proposed the progressive soft-decision sensing strategy, which uses just-enough sensing precision for ECC decoding through a trial-and-error manner, to obviate unnecessary extra sensing latency. Cui et al. [3] sorted the read requests according to the retention age of the data, and performed fast read for data with low retention ages by decreasing the number of sensing levels.

In this paper, the relationship among flash cell wearing, read latency and RBER is introduced. On one hand, flash cell wearing is reduced by lowering the threshold voltages but at the cost of less noise margins between the states of a flash cell, which in turn increase RBER and delay read operations. On the other hand, read latency can be decreased by improving the threshold voltages for lower RBER, which, however, results in more cell wearing. Based on the above relationship, we propose a comprehensive approach (HIRE) to exploit the access hotness information for improving read performance and endurance of flash memory storage systems. The basic idea is that we design a cross-layer hotness identifier in the buffer replacement management model of NAND flash memory, hence data pages in the buffer list can be classified into the following three groups, namely *read-hot*, *write-hot* and *mixed-hot*, respectively. Moreover, the fine-grained voltage controller is designed to supervise and control the appropriate threshold voltages of flash cells. In particular, the threshold voltages of a cell for *write-hot* data are decreased for wearing reduction, while these for *read-hot* data are increased for read latency reduction.

Trace-based simulations are carried out to demonstrate the effectiveness of our proposed approach. The results show that the proposed HIRE approach reduces the read response time by up to 43.49% on average and decreases the wearing of flash memory by up to 16.87% on average. In addition, HIRE does not have a negative write performance effect. Besides, the overhead of the proposed technique is negligible. In summary, this paper makes the following contributions.

- We proposed a cross-layer hotness identifier in the buffer replacement management model of NAND flash memory to guide flash read performance and endurance improvement.
- We proposed a voltage controller in the flash controller to improve flash-based system performance metrics with the guidance of the proposed hotness

identifier, which manages the appropriate threshold voltages for three types of data (*read-hot*, *write-hot*, *mixed-hot*) evicted by the buffer replacement algorithm.

- We carried out comprehensive experiments to demonstrate its effectiveness on both the wearing and read latency reduction without sacrificing write performance.

The rest of this paper is organized as follows. Section 2 presents the background and related work. Section 3 describes the design techniques and implementation issues of our hotness-guided access management for flash storage devices. Experiments and result analysis are presented in Section 4. Section 5 concludes this paper with a summary of our findings.

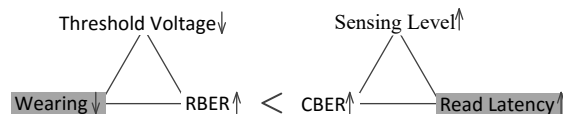
## 2 BACKGROUND AND RELATED WORK

In this section, we first present the tradeoff between flash cell wearing and read latency, which is due to two relationships. The first one is the relationship between threshold voltage, wearing and RBER. The second one is the relationship between error correction capability, read latency and the number of sensing levels when adopting LDPC as the default ECC scheme, which brings superior error correction capability as well as read response time degradation at the same time. Finally, previous studies related to this work are introduced for further work in this area.

### 2.1 Tradeoff between flash cell wearing and read latency

Firstly, the tradeoff between flash cell wearing and read latency is due to the relationship between threshold voltage, wearing and RBER. A flash chip is built from floating-gate cells whose state depends on the amount of charge they retain. Multi-level cell (MLC) flash memory uses cells with 4 or 8 states (2 or 3 bits per cell, respectively), as opposed to single-level cell (SLC) flash memory, which has 2 states (1 bit per cell). Every state is characterized by a threshold voltage ( $V_{th}$ ), which can be changed by injecting different amounts of charge onto the floating-gate. Recently, several works have showed that flash cell wearing is proportional to the threshold voltages [11] [12] [14]. The stress-induced damage in the tunnel oxide of a NAND flash memory cell can be reduced by decreasing the threshold voltages, and vice versa. Besides, the threshold voltages affect RBER significantly. When the threshold voltages are decreased, the noise margins among flash cell states are reduced, which reduces the capability for tolerating retention errors and increases RBER. Therefore, the tradeoff between wearing and RBER can be explored by controlling the threshold voltages with a wide range of settings. The less threshold voltages of a flash state, the less flash cell wearing, meanwhile, the higher RBER.

Secondly, the tradeoff is due to the significant relationship between error correction capability, read latency and the number of sensing levels when adopting



**Fig. 1.** Relationship between flash cell wearing and read latency.

LDPC scheme. The flash controller reads data from each cell by recursively applying several read reference voltages to the cell in a level-by-level manner to identify its threshold voltage. Therefore, sensing latency is linearly proportional to the number of sensing levels. In addition, when  $N$  sensing levels quantize the threshold voltage of each memory cell into  $N+1$  regions, a unique  $\lceil \log_2(N+1) \rceil$ -bit number is used to represent each region, indicating that transferring latency is proportional to the logarithm of  $N+1$ . Although read requests are delayed by slower sensing and transferring when using more sensing levels, more accurate input probability information of each bit for LDPC code decoding can be obtained, thus improving error correction capability (CBER).

Based on the precondition that RBER should be within CBER of the deployed LDPC code, the tradeoff between flash cell wearing and read latency can be concluded from above two relationships. As shown in Figure 1, flash cell wearing can be reduced by lowering the threshold voltages but at the cost of less noise margins between the states of a flash cell, which in turn increase RBER and delay read operations.

## 2.2 Related Work

Several methods for improving NAND flash memory I/O performance and endurance have been suggested by exploiting either of the two described relationships. By tasking advantage of the relationship between threshold voltage, wearing and RBER, Peleato et al. [12] proposed to optimize the target voltage levels to achieve a trade-off between lifetime and reliability, which tried to maximize lifetime subject to reliability constraints, and vice versa. Jeong et al. [11] presented a new system-level approach called dynamic program and erase scaling to improve the NAND endurance, by exploiting idle times between consecutive write requests to shorten the width of threshold voltage distributions so that blocks can be slowly erased with a lower erase voltage. However, both of them would induce the negative effects such as increased error and decreased write throughput.

Another set of approaches takes the relationship between error correction capability, read latency and the number of sensing levels into account. For example, Zhao et al. [5] proposed the progressive sensing level strategy to achieve latency reduction, which uses soft-decision sensing only triggered after the hard-decision decoding failure. Cai et al. [4] presented a retention optimized reading (ROR) method that periodically learns a tight upper bound and applies the optimal read reference voltage for each flash memory block online. Cui et al. [3] sorted

the read requests according to the retention age of the data, and performed fast read for data with low retention ages by decreasing the number of sensing levels. These state-of-the-art retention-aware methods improve the read performance significantly, and fortunately they are orthogonal to our work.

These studies demonstrate that wearing reduction by adjusting voltage and read performance improvement by soft-decision memory sensing are useful. However, none of these works consider the tradeoff between read latency and wearing when exploiting both of the two described relationships. In this paper, we focus on reducing both the read latency and wearing by controlling the threshold voltages based on the hotness information of each request, which can be easily acquired from the buffer manager.

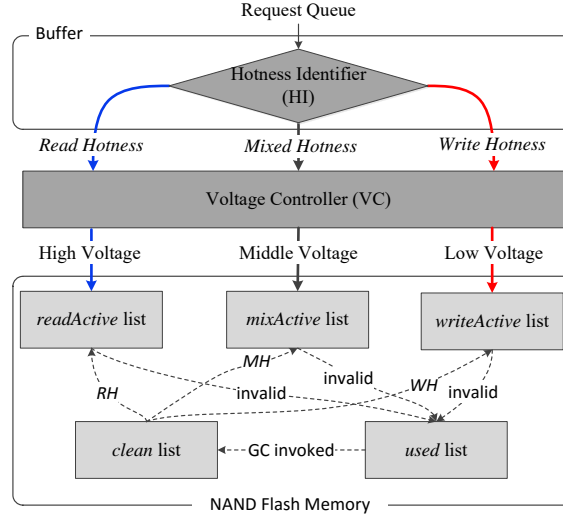
### 3 Exploiting Cross-Layer Hotness Identification to Improve Read and Endurance Performance (HIRE)

In this section, we propose HIRE, a wearing and read reduction approach, which includes two new management strategies: *Hotness Identifier (HI)* captures the access hotness characteristics at the flash buffer management level and *Voltage Controller (VC)* modulates the appropriate threshold voltages of a flash cell. We first present the cross-layer study on the access hotness characteristics of each data page in several workloads based on the buffer replacement policy. Then, on the basis of the observations of this cross-layer study, we propose a hotness-guided voltage controller to reduce the wearing and read latency. Finally, we present the overhead analysis.

#### 3.1 Cross-Layer Study for Hotness Identifier

In order to guide wearing and read latency reduction, the access hotness characteristics of each data page are needed. The hotness in this work means the frequency of read or write operations on each data page for a given period of time. We find that the hotness characteristics can be archived by the buffer manager. Buffer replacement policy can optimize the I/O sequence and reduce storage accesses, thus improving the overall efficiency of the storage system. In this study, we use the simple and efficient Least Recently Used (LRU) policy to acquire these information. Note that other buffer replacement algorithms for NAND flash-based storage systems, e.g., PT-LRU [15] or HA-LRU [16], are completely orthogonal to our work, and can be also used with the proposed HIRE approach to improve flash-based system performance metrics.

We implement the *Hotness Identifier (HI)* strategy in the buffer manager. LRU uses a linked list to manage all cached data page, and when data page is evicted out of the buffer, the historical statistical information about the read/write hit can be used to identify the access hotness characteristics because the read/write hit statistical information reflect the access history. In order to collect the access hotness characteristics, each data page in the buffer list adds two attributes: buffer read hit count  $C_r$  and buffer write hit count  $C_w$ . If one



**Fig. 2.** Flow in the wearing and read latency reduction approach.

data page is first referenced, it will be added into the MRU position of the buffer linked list, besides, its corresponding buffer read/write hit count value will plus one according to its read/write operation, respectively. When the data page in the buffer is referenced again, LRU adjusts the position of the data page to the MRU position of the linked list, meanwhile, its corresponding buffer read/write hit count value will also plus one according to its read/write operation.

During the eviction procedure, we classify the data access characteristics in the buffer into three types, as shown in Figure 2. When the buffer does not have a free page slot for the new access data, LRU preferentially evicts the data page in the LRU position of the linked list, which is the least recently accessed data page. At this time, when the read hit ratio grows to a high watermark  $\frac{C_r}{C_r+C_w} > w$  (in this work  $w = 95\%$ ), the most hit statistical information of a data page are read, and we determine it as *read-hot* (*RH*). For instance, the request data pages in the *WebSearch* trace from a popular search engine [17] are all read only, therefore, all data pages in this trace will be identified as *RH*. When the write hit ratio grows to a high watermark  $\frac{C_w}{C_r+C_w} > w$ , the most hit statistical information of a data page are write, and we determine it as *write-hot* (*WH*). Hence, other pages mixed with read and write hit are identified as *mixed-hot* (*MH*). As a result, data page are classified into three groups according to their access hotness characteristics.

### 3.2 Voltage Controller in HIRE

Based on the access hotness characteristics of each data page, the *Voltage Controller* (*VC*) strategy, aiming to the wearing and read latency reduction, is proposed, as shown in Figure 2. Furthermore, blocks transition between five states

are used in this work to cooperate with the fine-grained voltage controller strategy, namely *clean*, *readActive*, *writeActive*, *mixActive* and *used*. *Clean* state is an initial state of a block which receives none program operation, or is erased.

To improve the read performance, VC implemented in the flash controller boosts the threshold voltages of flash cells that store *RH* data, and hence the noise margins between the states of a flash cell increase, increasing the wearing of flash memory, but it leads to lower raw bit error rates so that its read performance is further improved. Note that *RH* data defines the resources stored in *RH* data page. Moreover, VC performs this *RH* data on a block with *readActive* state. If there is no *readActive* block, a block with *clean* state will be chosen as the current *readActive* block. When all pages on the *readActive* block have been written, it moves to the *used* state, and a new *readActive* block will be produced by the above method. Note that although the wearing is estimated to somewhat higher than that of the horizontal voltage line, fewer P/E (Program/Erase) operations in this *read-hot* data compensates for the wearing increment. Thus, this relatively small increment of wearing is acceptable to achieve the read performance improvement in the *read-hot* data.

If one data is classified as *WH*, the threshold voltages of corresponding flash cells are reduced, which reduces the wearing of flash memory. Simultaneously, VC performs this data on a block with *writeActive* state. When all pages on the *writeActive* block have been written, this block also moves to the *used* state, and a new *clean* block will be chosen as the current *writeActive* block. Although RBER is consequently higher in *WH*, the number of errors within a code word is not beyond the superior error correction capability of LDPC code. Moreover, the corresponding block is more likely to trigger garbage collection (GC) because of the *write-hot* access characteristic.

Traditional voltage control without the fluctuation of threshold voltages, is applied in *MH* data. These *MH* data will be performed on a block with *mixActive* state. When all pages on the *mixActive* block have been written, it also moves to the *used* state, and VC will chose a *clean* block as the current *mixActive* block. And when a *used* block is chosen by GC, it will be erased and move back to the *clean* state. Similar to hot/cold data separation policy [22] [23], separating *readActive*, *writeActive* and *mixActive* blocks improves the efficiency of garbage collection.

### 3.3 Overhead Analysis

The implementation of the proposed HIRE approach includes method implemented in the flash buffer manager and information recorded in the flash controller. In the flash buffer manager, we need to maintain the access hotness information. In order to figure out these information, the buffer list adds two attributes, including buffer read and write hit counts. We assume that the most buffer hit read/write counts is  $N_h$ , and the buffer size is  $N_c$  which specifies the maximum number of pages cached by buffer pool in this work, then the maximal size of storage required by the access hotness information is  $\lceil \log_2 N_h \rceil \times 2 \times N_c$  bits. This storage overhead is negligible for a state-of-the-art SSD. In addition,



we also extend each mapping entry in the flash translation layer of flash controller with a hotness field, using 2 bits to record three types of access hotness (*RH*, *WH* and *MH*), which is also negligible for a state-of-the-art SSD. Thus, it can be seen that the storage overhead is negligible.

## 4 Performance Evaluation

In this section, we first present the experimental methodology. Then, the I/O performance and endurance improvement of the proposed voltage optimization are presented. For comparison purpose, we have implemented several works which are closely related to our proposed HIRE approach.

### 4.1 Methodology

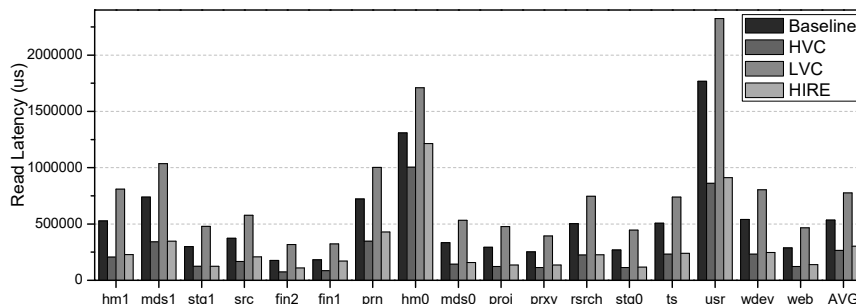
In this paper, we use an event-driven simulator to further demonstrate the effectiveness of the proposed HIRE. We simulate a 128GB SSD with 8 channels, each of which is connected to 8 flash memory chips. We implement dynamic page mapping scheme between logical and physical locations as the FTL mapping scheme. Greedy garbage collection and dynamic wear-leveling scheme are also implemented in the FTL of the simulator. All these settings are consistent with previous works [3].

**Table 1.** Parameters in this work

<i>Threshold Voltages</i>	<i>Sensing</i> ( $\mu$ s)	<i>Transfer</i> ( $\mu$ s)	<i>Program</i> ( $\mu$ s)
(1.40, 2.85, 3.55, 4.25)	30	40	600
(1.17, 2.37, 2.96, 3.54)	90	80	600
(0.93, 1.90, 2.37, 2.83)	210	100	600

As shown in Table 1, under the traditional voltage control, (1.17, 2.37, 2.96, 3.54) represents the threshold voltage of the four states in the simulated 2 bit/cell MLC NAND flash memory chip, and we use 600  $\mu$ s as the program latency when  $\Delta V_p$  is 0.25, 90  $\mu$ s as memory sensing latency and 80  $\mu$ s as data transfer latency when using LDPC with seven reference voltages. For the boosted threshold voltages, (1.40, 2.85, 3.55, 4.25) represents the corresponding threshold voltage of the four flash cell states when  $\Delta V_p$  is 0.3, 30  $\mu$ s as memory sensing latency and 40  $\mu$ s as data transfer latency. When the threshold voltages reduce, the threshold voltage of four flash cell states is (0.93, 1.90, 2.37, 2.83) and  $\Delta V_p$  is 0.2, the sensing latency is 210  $\mu$ s and the data transfer latency is 100  $\mu$ s.

For validation, we implement HIRE as well as baseline, HVC and LVC. We treat traditional voltage control strategy without any further targeted optimization as the baseline case in our contrastive experiments. HVC (High Voltage



**Fig. 3.** The read latency under different voltage control approaches.

Controller) approach increases threshold voltages of all flash cells to maximize read latency reduction, while LVC (Low Voltage Controller) approach reduces threshold voltages of all the flash cells to maximize wearing reduction. We evaluate our design using real world workloads from the MSR Cambridge traces [18] and two financial workloads [17], which are widely used in previous works to study SSD system performance metrics [19] [20] [21].

## 4.2 Experiment Results

In this section, the experimental results are presented and analyzed. Seventeen datasets were used in the experiments including a variety of application scenarios. To test the performance of our approach, the I/O response time and wearing are presented below.

Figure 3 shows the read latency of the proposed HIRE under different datasets, as compared to baseline, HVC and LVC approaches. It can be seen that HIRE outperforms baseline under all datasets, reducing read latency by 43.49% on average. From Figure 3, we can also see that the read performance improvement of HIRE under different datasets is distinctly different. For the *stg1* trace, the percent read latency reduction between baseline and HIRE is 58.59%. For the *fin1* trace, the percent read latency reduction between baseline and HIRE is 6.89%. HIRE increases voltages for the *read-hot* data, thus more *read-hot* significantly reduces the read latency. Figure 3 also shows that HVC approach archives the maximize read performance improvement, which is even better than HIRE. This significant improvement achieved for HVC approach is attributed to the fact that persistently increasing voltages leads to the minimize read response time, while HIRE only increases voltages of the *read-hot* data. However, HVC will also get the worst flash wearing at the same time, which can be seen in Figure 4.

Figure 4 shows the wearing weight for the seventeen traces under four approaches. Wearing weight is invoked as the metric to show the wearing degree of the proposed approach. The lower wearing weight, the longer endurance of storage system. It can be seen that compared with the traditional baseline approach, HIRE achieves significant wearing reduction. HIRE outperforms baseline

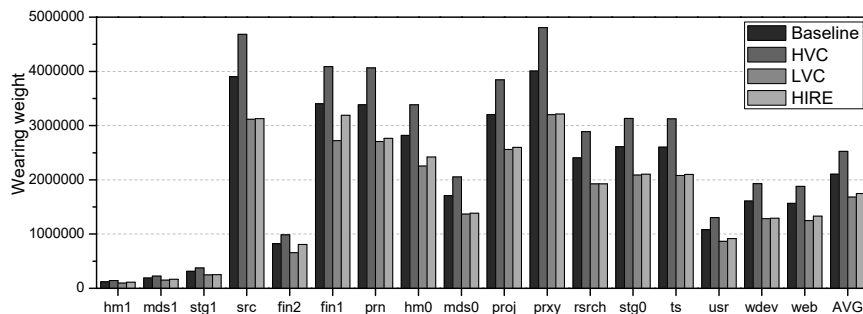


Fig. 4. The wearing weight under four approaches.

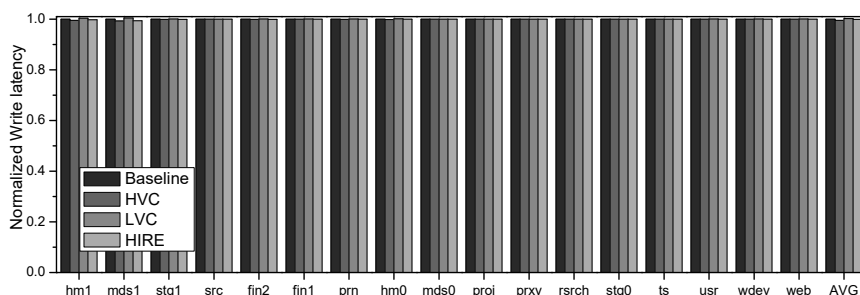


Fig. 5. The normalized write latency under different voltage control strategies.

with wearing reduction by 16.87% on average. On the other hand, compared to baseline, the greatest wearing reduction observed in *stg1* is 19.99%, while the smallest reduction observed in *fin2* is only 1.66%. This is because HIRE drops voltages for the *write-hot* data, thus more *write-hot* data significantly reduces the wearing. Moreover, HVC gets the maximum wearing weight although it realizes the maximum read latency reduction, which has been discussed above.

To further demonstrate that our HIRE approach does not sacrifice the write performance, the normalized write latency in all experiments are measured and presented in Figure 5. It can be seen that, very intuitively, the write latency of HIRE is comparable to those of the three competitor approaches. We can thus get the conclusion that the voltage control strategy in the proposed HIRE approach does not affect the write latency of the storage system.

## 5 Conclusion

In this paper, we proposed HIRE, a wearing and read reduction approach, that applies two new management strategies in NAND flash memory storage system: *Hotness Identifier (HI)* captures the hotness characteristics at the flash buffer management level and *Voltage Controller (VC)* manages the appropriate threshold voltages of a flash cell in the flash controller. The key insight behind

the design of HIRE is that based on the access hotness characteristics provided by *HL*, *VR* decreases the threshold voltages of a flash cell for the *write-hot* data for wearing reduction and increases these for the *read-hot* data for read latency reduction. Extensive experimental results and detailed comparisons show that the proposed approach is effective in various types of workloads for NAND flash memory storage system. On average, HIRE improves read performance by up to 43.49% and decreases the wearing of flash memory by up to 16.87% over previous voltage management approaches. In addition, HIRE does not have a negative write performance effect, and the overhead of the proposed approach is negligible.

## Acknowledgment

The authors would like to thank the anonymous reviewers for their detailed and thoughtful feedback which improved the quality of this paper significantly. This work was supported in part by the National Natural Science Foundation of China under grant NO.91330117, the National High-tech R&D Program of China (863 Program) under grant No.2014AA01A302, the Shaanxi Social Development of Science and Technology Research Project under grant No. 2016SF-428, the Shenzhen Scientific Plan under grant No.JCYJ20130401095947230 and No.JSGG20140519141854753.

## References

1. Margaglia, F., Yadgar, G., Yaakobi, E., Li, Y., Schuster, A. and Brinkmann, A.: The Devil is in the Details: Implementing Flash Page Reuse with WOM Codes. In Proceedings of Conference on File and Storage Technologies. February (2016)
2. Zhang, X., Li, J., Wang, H., Zhao, K. and Zhang, T.: Reducing Solid-State Storage Device Write Stress Through Opportunistic In-Place Delta Compression. In Proceedings of Conference on File and Storage Technologies. pp. 111–124, February (2016)
3. Cui, J., Wu, W., Zhang, X., Huang, J. and Wang, Y.: Exploiting Latency Variation for Access Conflict Reduction of NAND Flash Memory. In 32nd International Conference on Massive Storage Systems and Technology. May (2016)
4. Schroeder, B., Lagisetty, R. and Merchant, A.: Flash Reliability in Production: The Expected and the Unexpected. In Proceedings of Conference on File and Storage Technologies. pp. 67–80, February (2016)
5. Zhao, K., Zhao, W., Sun, H., Zhang, X., Zheng, N. and Zhang, T.: LDPC-in-SSD: making advanced error correction codes work effectively in solid state drives. In Proceedings of Conference on File and Storage Technologies. pp. 243–256 (2013)
6. Dong, G., Xie, N. and Zhang, T.: Enabling nand flash memory use soft-decision error correction codes at minimal read latency overhead. IEEE Transactions on Circuits and Systems I: Regular Papers. vol. 60, no. 9, pp. 2412–2421 (2013)
7. Wu, G., He, X., Xie, N. and Zhang, T.: Exploiting workload dynamics to improve SSD read latency via differentiated error correction codes. ACM Transactions on Design Automation of Electronic Systems. vol. 18, no. 4, pp. 55, (2013)

8. Jimenez, X., Novo, D. and Ienne, P.: Wear unleveling: improving NAND flash lifetime by balancing page endurance. In *Proceedings of Conference on File and Storage Technologies*. pp. 47–59 (2014)
9. Pan, Y., Dong, G. and Zhang, T.: Error rate-based wear-leveling for NAND flash memory at highly scaled technology nodes. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*. vol. 21, no. 7, pp. 1350–1354 (2013)
10. Agrawal, N., Prabhakaran, V., Wobber, T., Davis, J.D., Manasse, M.S. and Panigrahy, R.: Design Tradeoffs for SSD Performance. In *USENIX Annual Technical Conference*. pp. 57–70 (2008)
11. Jeong, J., Hahn, S.S., Lee, S. and Kim, J.: Lifetime improvement of NAND flash-based storage systems using dynamic program and erase scaling. In *Proceedings of Conference on File and Storage Technologies*. pp. 61–74 (2014)
12. Peleato, B. and Agarwal, R.: Maximizing MLC NAND lifetime and reliability in the presence of write noise. In *Proceedings of IEEE International Conference on Communications*. pp. 3752–3756, June (2012)
13. Jeong, J., Hahn, S.S., Lee, S. and Kim, J.: Improving NAND Endurance by Dynamic Program and Erase Scaling. In *USENIX Workshop on Hot Topics in Storage and File Systems*. June (2013)
14. Pan, Y., Dong, G. and Zhang, T.: Exploiting Memory Device Wear-Out Dynamics to Improve NAND Flash Memory System Performance. In *Proceedings of Conference on File and Storage Technologies*. pp. 18 (2011)
15. Cui, J., Wu, W., Wang, Y. and Duan, Z.: PT-LRU: a probabilistic page replacement algorithm for NAND flash-based consumer electronics. *IEEE Transactions on Consumer Electronics*. vol. 60, no. 4, pp. 614–622 (2014)
16. Lin, M., Yao, Z. and Xiong, J.: History-aware page replacement algorithm for NAND flash-based consumer electronics. *IEEE Transactions on Consumer Electronics*, vol. 62, no. 1, pp. 23–29 (2016)
17. Storage Performance Council traces. <http://traces.cs.umass.edu/storage/>
18. Narayanan, D., Thereska, E., Donnelly, A., Elnikety, S., Rowstron, A.: Migrating server storage to SSDs: Analysis of tradeoffs. In *Proceedings of the 4th ACM European conference on Computer systems*. Nuremberg, Germany, pp. 145–158 (2009)
19. Hu, Y., Jiang, H., Feng, D., Tian, L., Luo, H. and Zhang, S.: Performance impact and interplay of SSD parallelism through advanced commands, allocation strategy and data granularity. In *Proceedings of Proceedings of the international conference on Supercomputing*, pp. 96–107, May (2011)
20. Hu, Y., Jiang, H., Feng, D., Tian, L., Luo, H. and Ren, C.: Exploring and exploiting the multilevel parallelism inside ssds for improved performance and endurance. *IEEE Transactions on Computers*. vol. 62, no. 6, pp. 1141–1155 (2013)
21. Jung, M., Kandemir, M.: An evaluation of different page allocation strategies on high-speed ssds. In *Proceedings of USENIX Conference on File and Storage Technologies*, pp. 9 (2012)
22. Jung, S., Lee, Y., Song, Y.: A process-aware hot/cold identification scheme for flash memory storage systems. *IEEE Transactions on Consumer Electronics*. vol. 56, no. 2, pp.339-347 (2010)
23. Park, D. and Du, D: hot data identification for flash memory using multiple bloom filters, In *Proc. of USENIX Conference on File and Storage Technologies*, October (2011)