



# VIOS: A Variation-Aware I/O Scheduler for Flash-Based Storage Systems

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.. VIOS: A Variation-Aware I/O Scheduler for Flash-Based Storage Systems. 13th IFIP International Conference on Network and Parallel Computing (NPC), Oct 2016, Xi'an, China. pp.3-16, 10.1007/978-3-319-47099-3\_1 . hal-01647995

**HAL Id: hal-01647995**

**<https://inria.hal.science/hal-01647995>**

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

# VIOS: A Variation-aware I/O Scheduler for Flash-based Storage Systems

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.** NAND flash memory has gained widespread acceptance in storage systems because of its superior write/read performance, shock-resistance and low-power consumption. I/O scheduling for Solid State Drives (SSDs) has received much attention in recent years for its ability to take advantage of the rich parallelism within SSDs. However, most state-of-the-art I/O scheduling algorithms are oblivious to the increasingly significant inter-block variation introduced by the advanced technology scaling. This paper proposes a variation-aware I/O scheduler by exploiting the speed variation among blocks to minimize the access conflict latency of I/O requests. The proposed VIOS schedules I/O requests into a hierarchical-batch structured queue to preferentially exploit channel-level parallelism, followed by chip-level parallelism. Moreover, conflict write requests are allocated to faster blocks to reduce access conflict of waiting requests. Experimental results shows that VIOS reduces write latency significantly compared to state-of-the-art I/O schedulers while attaining high read efficiency.

**Keywords:** process variation, solid state drive, I/O scheduling, flash memory, parallelism

## 1 Introduction

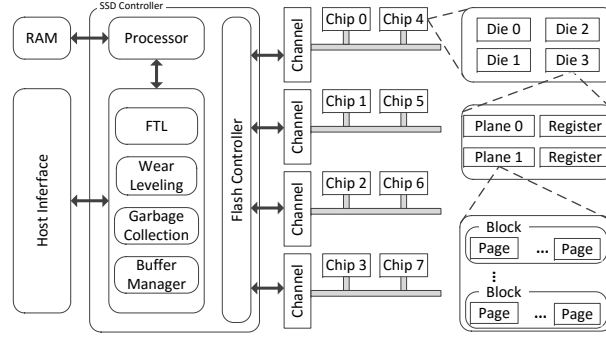
As NAND flash storage capacity becomes cheaper, NAND flash-based SSDs are being regarded as powerful alternatives to traditional Hard Disk Drives (HDDs) in a wide variety of storage systems [1]. However, SSDs introduce an out-of-place update mechanism and exhibit asymmetric I/O properties. In addition, a typical SSD usually offers rich parallelism by consisting of a number of channels with each channel connecting to multiple NAND flash chips [2, 3]. Most of conventional I/O schedulers including NOOP, CFQ and Anticipatory are designed to mitigate the high seek and rotational costs in mechanical disks, leading to many barriers to take full advantage of SSDs. Thus, I/O scheduling for SSDs has received much attention for its ability to take advantage of the unique properties within SSDs to maximize read and write performance.

Most of existing I/O scheduling algorithms for SSDs, such as PAQ [4], PIQ [5] and AOS [6], focus on avoiding resource contention resultant from shared SSD resources, while others take special consideration of Flash-Translation-Layer (FTL) [7] and garbage collection [8]. These works have demonstrated the importance of I/O scheduling for SSDs to reduce the number of read and write requests enrolled in conflict, which are the major contributors to access latency. However, little attention has been paid to dynamically optimize the data transfer latency, which could naturally reduce the access conflict latency when conflicts are unavoidable anymore.

The capacity of NAND flash memory is increasing continuously, as a result of technology scaling from 65nm to the latest 10nm technology and the bit density improvement from 1 bit per cell to the latest 6 bits per cell [9, 10, ?]. Unfortunately, for newer technology nodes, the memory block P/E cycling endurance has significantly dropped and process variation has become relatively much more significant. Recently many works have been proposed to exploit the process variation from different perspectives. Pan et al. [11] presented a dynamic BER-based greedy wear-leveling algorithm that uses BER statistics as the measurement of memory block wear-out pace by taking into account of inter-block P/E cycling endurance variation. Woo et al. [12] introduced a new measure that predicts the remaining lifetime of a flash block more accurately than the erase count based on the findings that all the flash blocks could survive much longer than the guaranteed numbers and the number of P/E cycles varies significantly among blocks. Shi et al. [13] further exploited the process variation by detecting supported write speeds during the lifetime of each flash block and allocating blocks in a way that hotter data are matched with faster blocks, but they did not reorder the requests in the I/O scheduling layer. Therefore, none of these works focuses on incorporating the awareness of inter-block variation into I/O scheduling to minimize the access conflict latency of I/O requests.

This paper proposes a variation-aware I/O scheduler to exploit the speed variation among blocks for write performance improvement without degrading read performance. The key insight behind the design of VIOS is that a variation-aware scheduler can organize the blocks of each chip in a red-black tree according to their detected write speeds and maintain a global chip-state vector, where the current number of requests for each chip is recorded, so as to identify conflict requests. By scheduling arrived requests into hierarchical-batch structured queues to give channel-level parallelism a higher priority than chip-level parallelism and allocating conflict write requests to faster blocks to exploit inter-block speed variation, access conflict latency of waiting requests is reduced significantly. Trace-based simulations are carried out to demonstrate the effectiveness of the proposed variation-aware I/O scheduling algorithm.

The rest of the paper presents the background and related works in Section 2. Section 3 describes the design techniques and implementation issues of our VIOS for flash storage devices. In Section 4, experimental evaluation and comparison with several alternative I/O schedulers are illustrated. Section 5 concludes this paper.



**Fig. 1.** SSD hardware diagram

## 2 BACKGROUND AND RELATED WORK

### 2.1 SSD Organization

Flash memory chips are organized in channels and ways, as shown in Figure 1. Within each flash memory chip are one or more dies, each further consisting of multiple planes, which are the smallest unit to be accessed independently and concurrently. Each plane is composed of a number of erase units, called blocks, and a block is usually comprised of multiple pages, which are the smallest unit to read/write. There are four main levels of parallelism which can be exploited to accelerate the read/write bandwidth of SSDs. Actually, the importance of exploiting parallelism on read/write performance improvement has been testified by numerous research works from different perspectives. For example, Roh et al. [14] explored to enhance B+-trees insert and point-search performance by integrating a new I/O request concept (psync I/O) into the B+-tree which can exploit the internal parallelism of SSDs in a single process, Hu et al. [15] argued that the utilization of parallelism, primarily determined by different advanced commands, allocation schemes, and the priority order of exploiting the four levels of parallelism, will directly and significantly impact the performance and endurance of SSDs.

Since the advance command support required by the die and plane level parallelism is not widely supported by most of SSDs, the degree of parallelism is usually governed by the number of channels multiplied by the number of flash memory chips in a channel, without taking the die and plane level parallelism into consideration. In this paper, our VIOS also exploit both channel-level parallelism and chip-level parallelism by scheduling arrived requests into hierarchical-batch structured queues.

### 2.2 Process Variation of Flash Memory

Along with the bit density developments and technology scaling of NAND flash memory, the aggravating process variation (PV) among blocks has been mag-

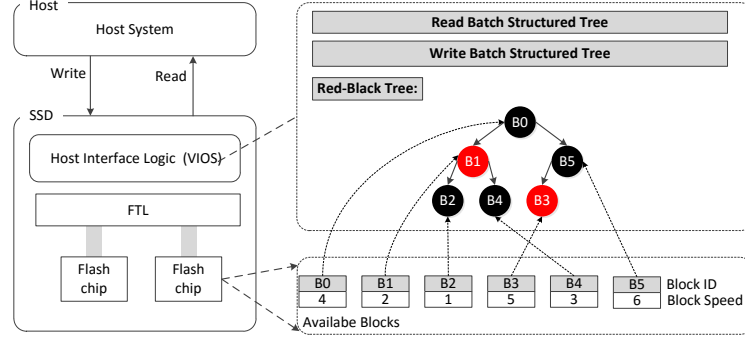
nified, which results in largely different P/E cycling endurance within different memory blocks when given the same ECC. PV is caused by the naturally occurring variation in the attributes of transistors, such as gate length, width and oxide thickness, when integrated circuits are fabricated. The distribution of the Bit Error Rates (BER) of flash blocks is characterized as the log Gaussian distribution by measuring 1000 blocks of a MLC NAND Flash memory chip manufactured in 35-nm technology at the same 15K P/E cycles [11].

Besides, there is a close relationship between BER and the speed of write operations. Typically, when program operations are carried out to write data into flash memory cells, the incremental step pulse programming (ISPP) scheme is introduced to appropriately optimize program voltage with the certain step size  $\Delta V_p$  that triggers a trade-off between write speed and BER. Using larger  $\Delta V_p$ , fewer steps are used to the desired level, thus the write latency is shorter. As the promising effect of reducing the write latency of write requests, however, the margin for tolerating retention errors is also reduced, resulting in higher BER. Therefore, with the awareness of both process variation and the BER-speed relationship, write speed for lower-BER blocks can be increased at the cost of reduced noise margins, while that for higher-BER blocks should be carefully optimized without exceeding the capability of the deployed ECC. The challenge to detect the proper write speed for each block at its current worn out state is also solved by periodically reading out the written data to find out the number of faulty bits, and the analysis indicates that overhead is negligible [13]. In this paper, the blocks of each chip are sorted according to their detected write speeds and conflict write requests are allocated to faster blocks to reduce access conflict of waiting requests.

### 2.3 I/O Scheduler for Flash-Based SSDs

An increasing number of I/O scheduling algorithms have been proposed to improve flash memory storage system performance from different perspectives. The first two algorithms for Flash-Based SSDs called IRBW-FIFO and IRBW-FIFO-RP were designed by Kim et al. [16]. The basic idea is to arrange write requests into bundles of an appropriate size while read requests are independently scheduled by taking the read/write interference into consideration. Numerous research works enhanced the IRBW-FIFO and IRBW-FIFO-RP by exploiting rich parallelism in I/O scheduling, such as PAQ [4] and PIQ [5]. In addition, there is also recognition on the importance of fairness in multi-programmed computer systems and multi-tenant cloud systems, such as FIOS [17] and FlashFQ [18].

However, little attention has been paid to optimize the data transfer latency dynamically, which could naturally reduce the access conflict latency when conflicts can not be avoided anymore. Our VIOS focuses on incorporating the awareness of inter-block variation into I/O scheduling to minimize the access conflict latency of I/O requests. Fortunately, all these existing algorithms are somewhat orthogonal to our work, and can be used concurrently with the proposed VIOS to optimize the efficiency of flash-based I/O scheduling.



**Fig. 2.** Organizational view of the VIOS scheduler

### 3 Details of VIOS

Motivated by the increasingly significant inter-block variation introduced by the advanced technology scaling, a variation-aware I/O scheduler is proposed to minimize the access conflict latency of I/O requests. The proposed VIOS scheduler is implemented in the host interface logic (HIL), where the knowledge of both the interface protocols and SSD specific characteristics contributes to better device-specific scheduling decisions for I/O requests. The main idea behind the design of VIOS is to allocate conflict write requests to faster blocks so as to reduce access conflict latency of waiting requests. To achieve this, our scheduler organizes the blocks of each chip in a red-black tree according to their detected write speeds and maintains a global chip-state vector, where the current number of requests for each chip is recorded, so as to identify conflict requests.

#### 3.1 Block Management

After detecting the proper write speed for each block at its current worn out state by periodically reading out the written data to find out the number of faulty bits, it is important to manage blocks with different write speeds for VIOS to easily schedule requests to appropriate blocks. In VIOS, the red-black tree structure is adopted as the main data structure to sort its blocks in detected speed order. Since the advance command support required by the die and plane level parallelism is not widely supported by most of SSDs and the relatively higher chip-level conflicts are the major contributors to access latency, the blocks of each chip are associated with a red-black tree respectively. Figure 2 shows the main data structures associated with VIOS. Once all the pages of a block have been programmed, it will be set as an invalid node in the red-black tree. When the prepared empty blocks are used up, a time-consuming task called Garbage Collection (GC) is triggered to reclaim stale pages for free write space and then the states of erased blocks become valid again. The blocks are evicted and inserted into another place of the red-black tree only when the write speed

detection process is triggered and its  $\Delta V_p$  is decreased, corresponding to reduced write speed.

### 3.2 Global Chip-state Vector

To build a scheduling algorithm that can identify conflicts and exploit the speed variation based on the degree of conflicts, we propose a global chip-state vector to track the current number of requests for each chip. The global chip-state vector mainly depends on the location of data, which is determined by different data allocation schemes and the logical sector number (LSN) of I/O requests. Lets take an SSD where the number of chips is  $ChipNum$  and the size of each page is  $PageSize$  as an example. For a given static allocation scheme where the priority order of parallelism is channel-level parallelism first, followed by the chip-level parallelism [19], the assemblage of chips accessed by request  $r$  can be defined as:

$$A_r = \left\{ \mu | \mu = \left( \frac{(lsn(r)+\lambda)*SectorSize}{PageSize} \right) \% ChipNum \right\} \\ \text{where } 0 \leq \lambda < len(r) \quad (1)$$

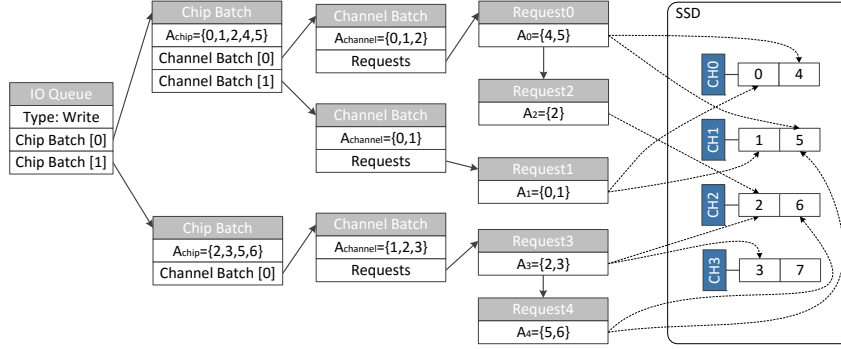
Where  $lsn(r)$  and  $len(r)$  are the accessed logical sector number and data size in sectors of request  $r$  respectively, while  $SectorSize$  is the sector size in bytes. For a global chip-state vector defined as  $(NR_0, NR_1, \dots, NR_i)$  where  $NR_i$  is the current number of requests for chip  $i$ , when pushing arrived requests into the queue or issuing chosen requests to SSDs, the  $NR$  of each chip accessed by requests is updated as follows:

$$NR_i = \begin{cases} NR_i + 1 & \text{arriving} \\ NR_i - 1 & \text{issued} \end{cases}, i \in A_r \quad (2)$$

### 3.3 Conflict Optimized Scheduling Mechanism

Next, we propose the conflict optimized scheduling technique, which aims to reduce access conflict latency by exploiting the rich parallelism within SSDs and the speed variation among blocks. It consists of two components: 1) a hierarchical-batch structured queue to avoid conflicts from chip to channel. 2) a variation-aware block allocation technique that assigns conflict write requests to faster blocks to reduce access conflict of waiting requests.

**Hierarchical-batch Structured Queue** Since there are four main levels of parallelism which can be exploited to accelerate the read/write bandwidth of SSDs, the conflicts can also be classified into four types based on the physical resources contended by arrived requests. Among them, channel conflicts and chip conflicts are taken into account for the reason that the advance command support required by the die and plane level parallelism is not widely supported by most of SSDs. In the hierarchical-batch structured queue, the requests are grouped



**Fig. 3.** An example that shows VIOS conflict reduction with the hierarchical-batch structured queue

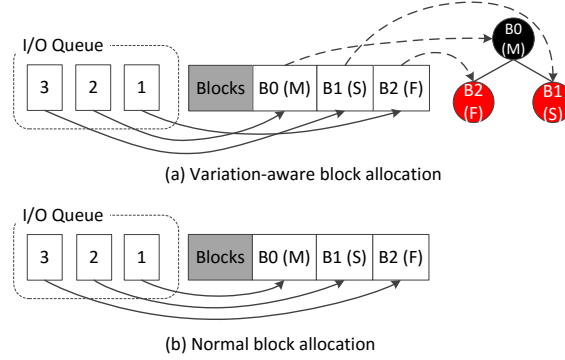
into batches from bottom up based on the chip-level and channel-level conflict detection respectively. Requests in the same chip batch can be serviced in chip-level parallel, while requests belong to the same channel batch can be executed in completely independent channel-level parallel. Each chip batch and channel batch use  $A_{chip}$  and  $A_{channel}$  respectively to track the assemblage of chips and channels accessed by all its requests. Each time when a new request  $r$  arrives, the detection of chip-level conflicts is enabled by repeated intersection operations of  $A_{chip}$  with  $A_r$ , whose resulting assemblage of conflict chips ( $A_{conflict}$ ) can be modeled as:

$$A_{conflict} = A_{chip} \cap A_r \quad (3)$$

Once a chip batch that has no conflict with  $r$  is found, which means  $A_{conflict}$  is empty, the request is added to the chip batch and  $A_{chip}$  is updated. Otherwise, a new chip batch is created for the new request. After that, the detection of channel-level conflicts is performed in the same manner to further exploit channel-level parallelism of requests within the same chip batch. For example, after scheduling the first five requests into the hierarchical-batch structured queue shown in Figure 3, if one more request accessing chips  $\{6, 7\}$  arrives, it is found that there is no chip-level conflict between the new request and the first chip batch, where the new request is thus added and  $A_{chip}$  is updated to  $\{0, 1, 2, 4, 5, 6, 7\}$ . Then each channel batch of the first chip batch is checked and the second channel batch where the new request has no channel-level conflict is chosen with  $A_{channel}$  being updated to  $\{0, 1, 2, 3\}$ .

**Variation-aware Block Allocation** Motivated by the findings that speed variation among blocks can be easily detected, we propose a variation-aware block allocation algorithm to optimize the data transfer latency dynamically, which could naturally reduce the access conflict latency when conflicts can not be avoided anymore. Each time when a request is issued, the NR of each chip accessed by the request is checked and then updated. Since a request processing





**Fig. 4.** Block allocation process. Program speed represented by Slow (S), Medium (M) and Fast (F).

may access multiple chips in NAND flash memory, we scatter the request into separate sub-requests. Each sub-request is only able to gain access to one chip. The sub-request accessing the chip whose NR is more than 1 will be allocated to a faster block from the red-black tree of the chip. Otherwise, a slower block is chosen for the sub-request. Figure 4 shows the process of scheduling three conflict sub-requests when only three pages in three blocks with different speeds are empty. As the Figure 4(a) shows, the variation-aware block allocation algorithm assigns the first two sub-requests to currently faster blocks for the reason that one or more sub-requests are still waiting in the queue. The last sub-request is allocated to a slower block because no conflict sub-request is waiting at this moment. Assuming that the write latency of fast, medium and slow blocks are  $150\mu s$ ,  $180\mu s$ , and  $210\mu s$  respectively, the average request response time of the proposed algorithm is  $(150+330+540)/3 = 340 \mu s$ , while that of the normal algorithm (Figure 4(b)) which distributes blocks in order is  $(180+390+540)/3 = 370 \mu s$ . Therefore, by incorporating the awareness of inter-block variation into I/O scheduling, the access conflict latency of I/O requests is reduced significantly.

**Overhead Analysis** The overheads of the proposed I/O scheduler are analyzed as follows. According to the detailed descriptions of components above, the implementation of VIOS needs to maintain hierarchical-batch structured queues in the I/O queue. Since the number of channels and chips in NAND flash memory is limited, all sets can be stored as binary words and the set-intersection/set-union operation can be performed as an  $O(1)$ -time bitwise-AND/OR operation. This storage overhead is negligible for an I/O queue. Furthermore, the complexity of adding an incoming I/O request into the hierarchical-batch structured queues is proportional to the sum of the number of chip batches and the number of channel batches, and it is less than the queue length of I/O scheduler, which also has negligible cost.

## 4 Experimental Results

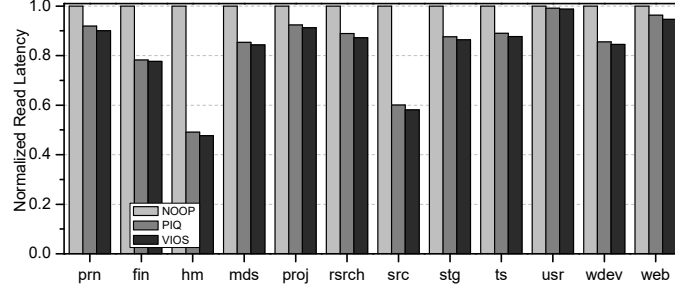
To evaluate our proposed variation-aware I/O scheduler, we perform a series of trace driven simulations and analysis. We implement VIOS as well as baseline NOOP scheduling and state-of-the-art PIQ scheduling within an event-driven simulator named as SSDSim [19], which provides the detailed and accurate simulation of each level of parallelism. Note that write speed detection technique is implemented with all of the schedulers. We simulate a 128GB SSD with 8 channels, each of which is connected to 8 chips. For the flash micro-architecture configuration, each flash chip employs 2 dies, where each die contains 4 planes and each plane consists of 2048 blocks. Each flash block contains 64 pages with a page size of 2KB. All these settings are consistent with previous works [5]. Page mapping FTL is configured to maintain a full map of logical pages to physical ones and greedy garbage collection scheme is implemented.

The BER growth rate that follows Bounded Gaussian distribution is used to simulate the process variation of flash memory, where the mean  $\mu$  and the standard deviation  $\sigma$  are set as  $3.7 \times 10^{-4}$  and  $9 \times 10^{-5}$  respectively [11]. The maximal possible write step size is set to 0.6 and the step of decreasing  $\Delta V_p$  is set to 0.03. We use  $600\mu s$  as the 2bit/cell NAND flash memory program latency when  $\Delta V_p$  is 0.3,  $20\mu s$  as memory sensing latency and 1.5ms as erase time. Four different wear-out stages corresponding to 15K, 12K, 9K and 6K P/E cycles are evaluated. We evaluate our design using real world workloads from MSR Cambridge traces [20] and the write-dominated Financial1 trace [21], where 500000 I/Os of each trace are used in accordance with previous work.

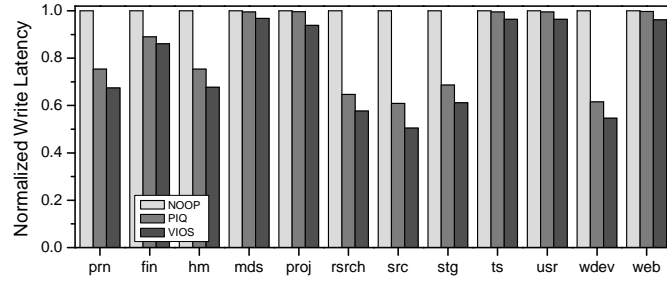
### 4.1 Performance Analysis of VIOS

Our experiments evaluate scheduling performance with read and write latency. Figure 5 shows the average read latency for NOOP, PIQ and VIOS tested under the P/E cycling of 12K. As can be observed, VIOS improves the average read latency by about 17.66% compared to NOOP, indicating that the hierarchical-batch structured read queue helps VIOS exploit multilevel parallelism inside SSDs by resolving resource conflicts. However, the improvements in average read latency brought by VIOS are not significantly higher than those obtained when using PIQ. This is because the variation-aware block allocation technique of VIOS mainly serves write request, and read requests are always preferentially scheduled in both PIQ and VIOS without being affected by write performance improvement.

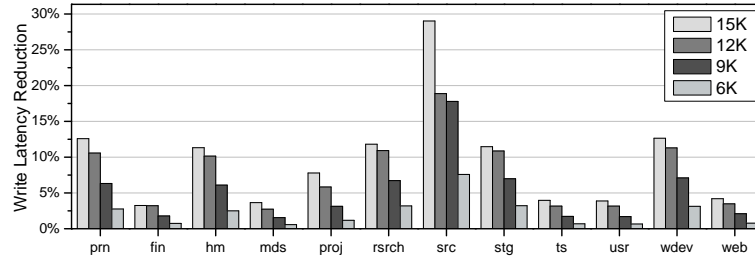
Figure 6 plots the simulation results on average write latency when different scheduling algorithms are used in the variation-induced SSDs. To facilitate the comparison, the average write latency is normalized against the case of using NOOP algorithm. The first thing to observe is that VIOS outperforms NOOP and PIQ with write latency reduction by 22.93% and by 7.71% on average, respectively. This is because both the hierarchical-batch structured write queue and the variation-aware block allocation algorithm reduce access conflict of write requests. However, the write performance improvements under different traces



**Fig. 5.** Average read latencies for three different types of schedulers (normalized to the NOOP scheduler)



**Fig. 6.** Average write latencies for three different types of schedulers (normalized to the NOOP scheduler)



**Fig. 7.** A comparison of the write latency reduction relative to PIQ among four different wear-out stages

vary greatly. For example, compared to PIQ, the greatest improvement made in the *src* trace is 17.17%, but the slightest improvement made in the *mds* trace is only 2.73%. This is due to the different percentages of requests enrolled in conflict – VIOS works for I/O intensive applications where more requests can be processed in parallel and optimized. Table I shows the percentages of conflicts collected under the P/E cycling of 12K with NOOP scheduler. One can also observe that the percentage of conflicts in *src* is 77.74%, which has an impact

**Table 1.** Characteristics of used workloads

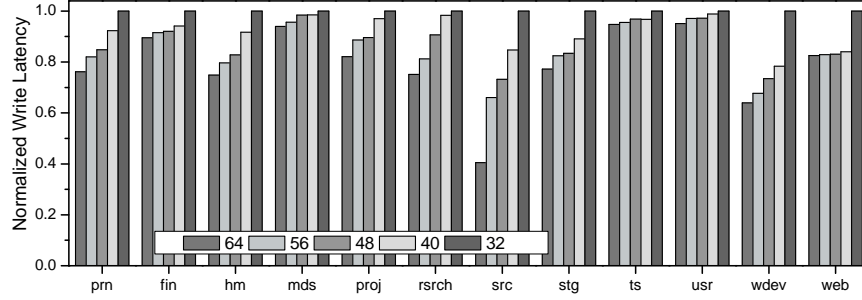
<i>Traces</i>	<i>ReadI/O</i>	<i>WriteI/O</i>	<i>ReadConflicts</i>	<i>WriteConflicts</i>	<i>ConflictsRatio</i>
<i>prn</i>	61422	438578	37678	261787	59.89
<i>fin</i>	129342	370658	12135	129379	28.30
<i>hm</i>	139527	360473	54387	168864	44.65
<i>mds</i>	26301	473699	10947	12570	4.70
<i>proj</i>	151549	348451	38615	32852	14.29
<i>rsrch</i>	43803	456197	19413	217828	47.45
<i>src</i>	81377	418623	28146	360549	77.74
<i>stg</i>	43818	456182	21657	210355	46.40
<i>ts</i>	76687	423313	27566	7161	6.95
<i>usr</i>	206983	293017	91391	9068	20.09
<i>wdev</i>	102663	397337	48986	170202	43.84
<i>web</i>	239458	260542	89753	19184	21.79

on improving efficiency of VIOS. In contrast, that of *mds* is 4.70%, indicating that fewer conflicts lead to slighter write performance improvement.

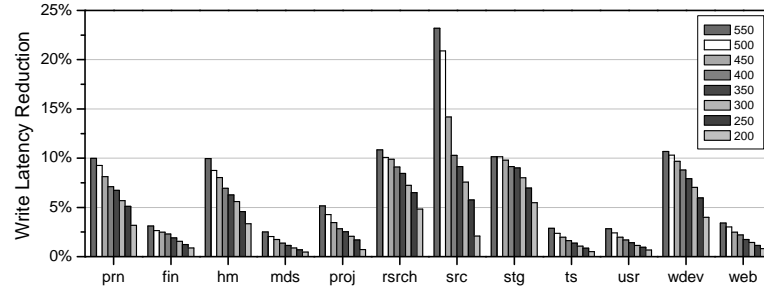
Figure 7 gives a comparison of the write latency reduction relative to PIQ among four different wear-out stages. One can observe from these results that the write latency reduction is positive all the time, which means VIOS always outperforms PIQ under four different wear-out stages. Furthermore, VIOS improves write performance by 2.25%, 5.25%, 7.71% and 9.63% on average under the P/E cycling of 6K, 9K, 12K and 15K respectively. As can be observed, with the increased number of P/E cycling, the write performance improvement brought by VIOS gets greater. This is a very reasonable result since BER spread grows as flash memory cells gradually wear out with the P/E cycling, corresponding to more significant variation among blocks, which improves the efficiency of the variation-aware block allocation strategy in VIOS. Overall, these results clearly demonstrate the effectiveness of VIOS in reducing the write latency during the entire flash memory lifetime.

## 4.2 Sensitivity Analysis of VIOS

To measure the sensitivity of VIOS to the I/O intensity, we repeated our experiments by varying the number of chips and the baseline program latency. Either fewer chips or slower program speeds increase the probability of access conflict. Figure 8 plots the normalized average write latency for each trace under 64, 56, 48, 40 and 32 chips when using VIOS. From the plots, it can be seen that the write latency increases as the number of chips decreases. For most traces, varying the number of chips from 64 to 32 increases the write latency by less than 25%. However, for traces *src* and *wdev*, the increase in write latency is 59.48% and 36.09% respectively. By comparing the results with the percentages of write conflicts shown in Table I, it can be observed that the increment in average write latency is greater when the number of write conflicts is larger. For



**Fig. 8.** Average write latencies for five different numbers of chips when using VIOS (normalized to 32 chips)



**Fig. 9.** A comparison of the write latency reduction relative to PIQ among eight different baseline program latencies

example, the write latency of *src* that has most write conflicts (360549) is increased with maximum rate (59.48%), while these of *mds*, *ts* and *usr* which have fewer write conflicts (12570, 7161 and 9068) are increased with minimum rate (6.05%, 5.29% and 4.91%). On one hand, the number of conflicts is proportional to the quotient of access density and the number of chips, which means that more conflicts occur when reducing chips for traces with more intensive I/O. On the other hand, the average write latency is proportional to the square of the number of write conflicts, amplifying the effect of each new conflict.

Figure 9 plots the impact of the baseline program latency on the write latency reduction relative to PIQ. The x-axis is the baseline program latency for  $\Delta V_p=0.3$ , varying from  $200\mu s$  to  $550\mu s$ . The number of conflict requests increases as the program latency is delayed, thus improving the benefit from hierarchical-batch structured queues and variation-aware block allocation technique. However, the effect of program latency delay is greater than that of reduction in the number of chips. For example, from  $200\mu s$  to  $550\mu s$ , the write latency reduction for *rsrch* varies from 4.83% to 10.85%, compared to a slighter variation from 10.92% to 13.63% as the number of chips varies from 64 to 32. The major reason is that delaying program latency not only increases the number of conflict

requests, but also amplifies the access conflict latency, which is the dominant factor for slow write operations.

## 5 Conclusion

In this paper, we propose a variation-aware I/O scheduler (VIOS) for NAND flash-based storage systems. The process variation is exploited to reduce the access conflict latency of SSDs when conflicts are unavoidable anymore. VIOS organizes the blocks of each chip in a red-black tree according to their detected write speeds and allocate conflict write requests to faster blocks to exploit inter-block speed variation. In addition, the hierarchical-batch structured queue that focuses on the exploration of the parallelism of SSDs is presented. Furthermore, with diverse system configurations such as wear-out stages, the number of chips and the baseline program latency, VIOS reduces write latency significantly compared to the state-of-the-art NOOP and PIQ while attaining high read efficiency.

## 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. Min, C., Kim, K., Cho, H., Lee, S.W. and Eom, Y.I.: SFS: random write considered harmful in solid state drives. in *USENIX Conference on File and Storage Technologies*, pp. 12, Feb (2012)
2. 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)
3. Chen, F., Lee, R. and Zhang, X.: Essential roles of exploiting internal parallelism of flash memory based solid state drives in high-speed data processing. in *IEEE 17th International Symposium on High Performance Computer Architecture*, pp. 266–277, Feb (2011)
4. Jung, M., Wilson III, E.H. and Kandemir, M.: Physically addressed queueing (PAQ): improving parallelism in solid state disks. *ACM SIGARCH Computer Architecture News*, Vol. 40, No. 3, pp. 404–415, June (2012)
5. Gao, C., Shi, L., Zhao, M., Xue, C.J., Wu, K. and Sha, E.H.: Exploiting parallelism in i/o scheduling for access conflict minimization in flash-based solid state drives. in: *30th Symposium on Mass Storage Systems and Technologies* pp. 1–11, June (2014)

6. Li, P., Wu, F., Zhou, Y., Xie, C. and Yu, J.: AOS: Adaptive Out-of-order Scheduling for Write-caused Interference Reduction in Solid State Disks. in *Proceedings of the International MultiConference of Engineers and Computer Scientists*, Vol. 1 (2015)
7. Wang, M. and Hu, Y.: An I/O scheduler based on fine-grained access patterns to improve SSD performance and lifespan. in *Proceedings of the 29th Annual ACM Symposium on Applied Computing*, pp. 1511–1516, March (2014)
8. Jung, M., Choi, W., Srikantaiah, S., Yoo, J. and Kandemir, M.T.: HIOS: a host interface I/O scheduler for solid state disks. in *ACM SIGARCH Computer Architecture News*, pp. 289–300, June (2014)
9. Ho, K.C., Fang, P.C., Li, H.P., Wang, C.Y.M. and Chang, H.C.: A 45nm 6b/cell charge-trapping flash memory using LDPC-based ECC and drift-immune soft-sensing engine. in *IEEE International Solid-State Circuits Conference Digest of Technical Papers*, pp. 222–223, Feb (2013)
10. Zuloaga, S., Liu, R., Chen, P.Y. and Yu, S.: Scaling 2-layer RRAM cross-point array towards 10 nm node: A device-circuit co-design. in *IEEE International Symposium on Circuits and Systems*, pp. 193–196, May (2015)
11. 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)
12. Woo, Y.J. and Kim, J.S.: Diversifying wear index for MLC NAND flash memory to extend the lifetime of SSDs. in *Proceedings of the Eleventh ACM International Conference on Embedded Software*, pp. 6, Sep (2013)
13. Shi, L.H., Di, Y., Zhao, M., Xue, C.J., Wu, K. and Sha, E.H.M.: Exploiting Process Variation for Write Performance Improvement on NAND Flash Memory Storage Systems. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 99, pp.1–4 (2015)
14. Roh, H., Park, S., Kim, S., Shin, M. and Lee, S.W.: B+-tree index optimization by exploiting internal parallelism of flash-based solid state drives. in *Proceedings of the VLDB Endowment*, pp. 286–297 (2011)
15. 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)
16. Kim, J., Oh, Y., Kim, E., Choi, J., Lee, D. and Noh, S.H.: Disk schedulers for solid state drivers. In *Proceedings of the seventh ACM international conference on Embedded software*, pp. 295–304, Oct (2009)
17. Park, S. and Shen, K.: FIOS: a fair, efficient flash I/O scheduler. in *USENIX Conference on File and Storage Technologies*, pp. 13, Feb (2012)
18. Shen, K. and Park, S.: FlashFQ: A Fair Queueing I/O Scheduler for Flash-Based SSDs. In *USENIX Annual Technical Conference*, pp. 67–78, Jun (2013)
19. Hu, Y., Jiang, H., Feng, D., Tian, L., Luo, H., Zhang, S.: Performance impact and interplay of SSD parallelism through advanced commands, allocation strategy and data granularity. in *proceedings of the international conference on Supercomputing*, pp. 96–107, May (2011)
20. Narayanan, D., Thereska, E., Donnelly, A., Elnikety, S., Rowstron, A.: Migrating server storage to SSDs: analysis of tradeoffs. in *the European Conference on Computer Systems*, pp. 145–158, April (2009)
21. UMass Trace Repository, <http://traces.cs.umass.edu>
22. Cui, J., Wu, W., Zhang, X., Huang, J., Wang, Y.: Exploiting process variation for read and write performance improvement of flash memory. in: *32th International Conference on Massive Storage Systems and Technology*, May (2016)