

Shiqiang NIE, Jie NIU, Yingzhao SHAO, Xiaobo LI, Mingming ZHANG, Weiguo WU, 2026. GC bypass: decoupling GC from the flash translation layer to eliminate GC-induced long-tail latency inside SSD. *ENGINEERING Information Technology & Electronic Engineering*, 27(2):250152.

<https://doi.org/10.1631/ENG.ITEE.2025.0152>

GC bypass: decoupling GC from the flash translation layer to eliminate GC-induced long-tail latency inside SSD

Key words: Solid-state drive (SSD); NAND flash; Garbage collection (GC); Interconnected network; Flash channel

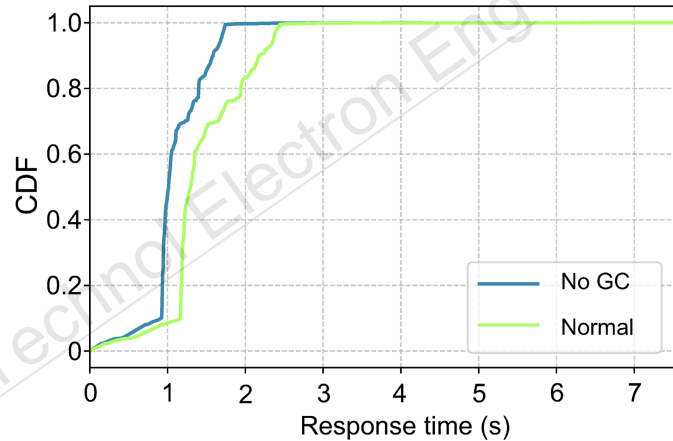
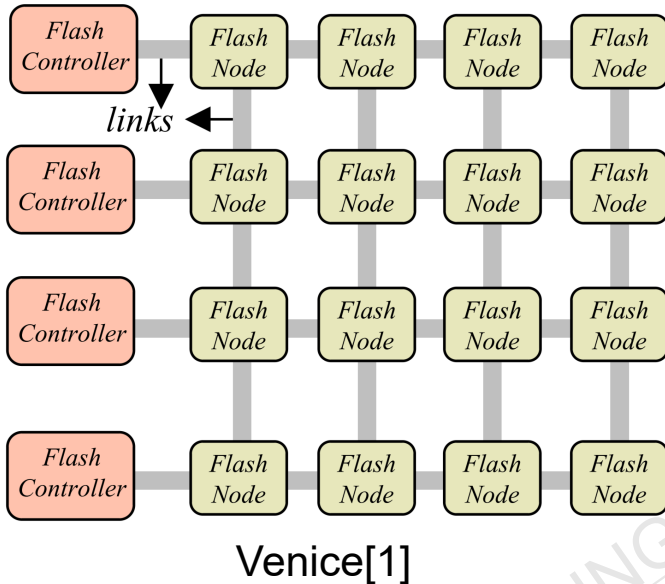
Weiguo WU

E-mail: wgwu@xjtu.edu.cn

 ORCID: <https://orcid.org/0009-0000-8298-0572>

Motivation

Venice: a low-cost flash chip interconnect architecture

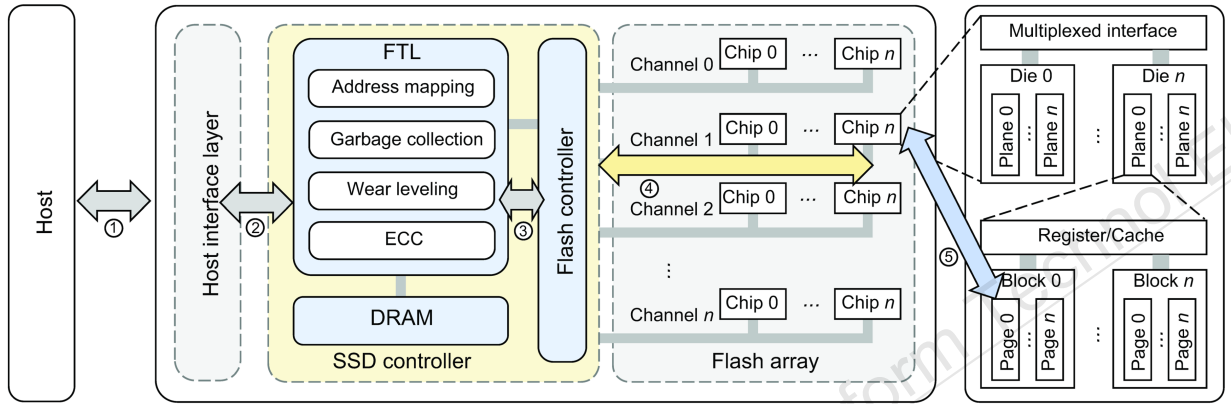


Cumulative distribution function (CDF) of response time in workload HM_1

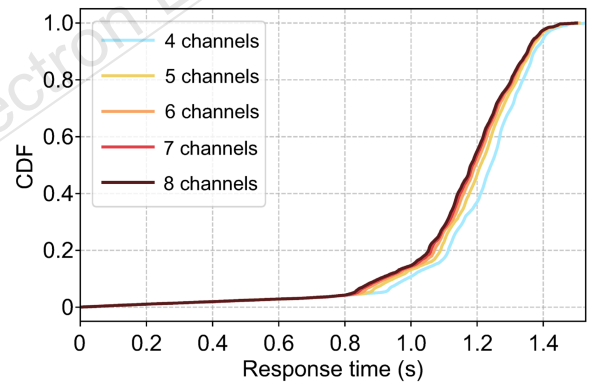
- Venice re-architects the conventional channel structure of solid-state drives (SSDs) by introducing a low-cost interconnected network between the SSD controller and flash chips to effectively mitigate path contention issues.
- While Venice demonstrates excellent performance under fair scheduling scenarios, its lack of priority awareness results in unresolved resource contention between garbage collection (GC) and regular I/O operations.

Motivation

Key question: Can Venice effectively reduce long-tail latency in SSDs?



Workflow of read/write operations in SSDs

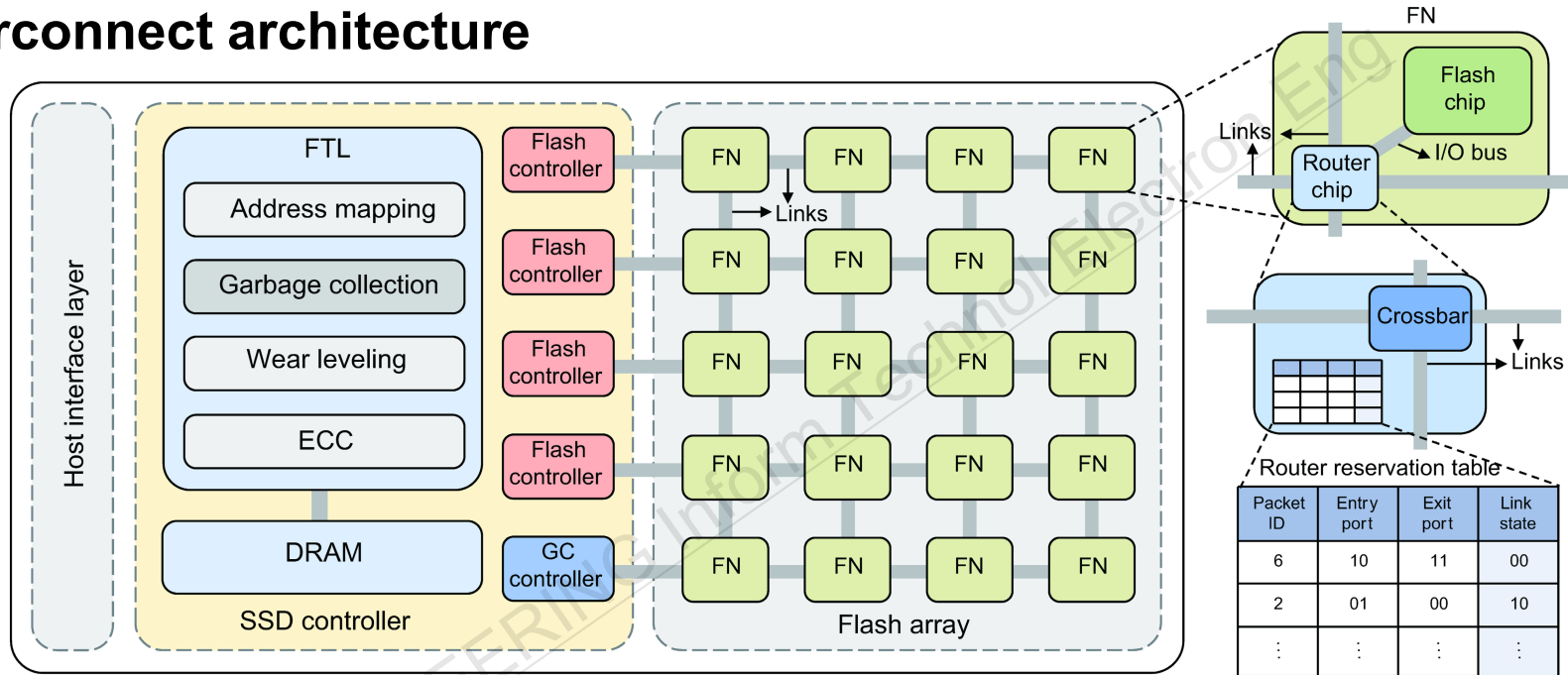


Cumulative distribution function (CDF) of response time under more channels

- The operational latency inside NAND flash chips, which depends on physical processes, is significantly higher than the electrical signal transmission delay over the channels. The primary source of read/write latency originates from inside the flash chip.
- The long-tail latency remains at comparable levels despite the increased number of network entry points. This phenomenon stems from Venice's inability to specifically address GC-I/O contention.

Design

GC bypass: a GC-induced long-tail latency suppressive flash chip interconnect architecture



Overview of GC bypass

- **Dedicated GC controller**

- Isolate entry points of GC vs. normal I/O requests in flash chip interconnect
- Integrate key functions: effective page caching, path management, etc.

- **Priority-aware scheduling & preemption**

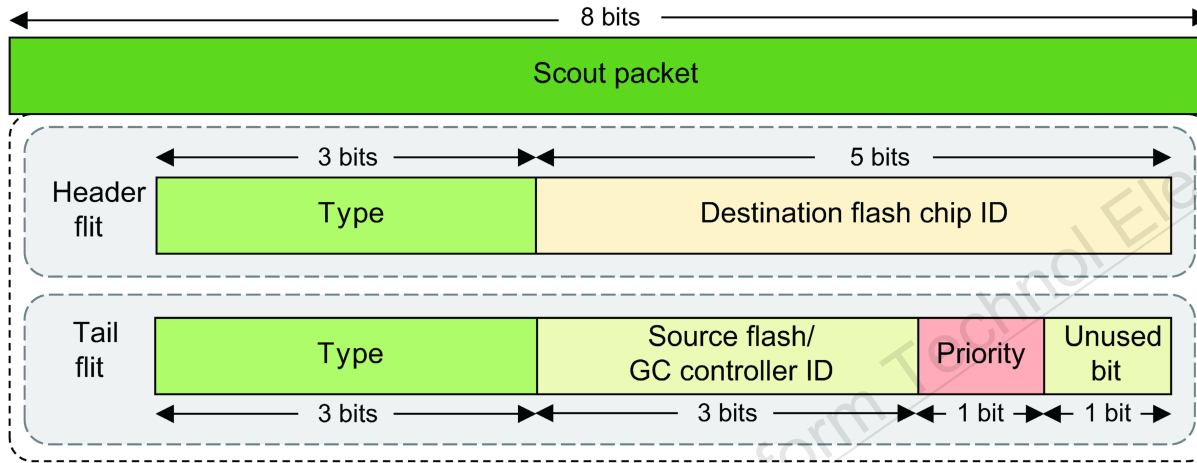
- High priority: normal I/O requests + critical GC sub-operations (valid page reads, victim block erases)
- Low priority: GC valid page write-back/programming
- High-priority requests preempt low-priority ones by taking over pre-reserved paths

Router reservation table

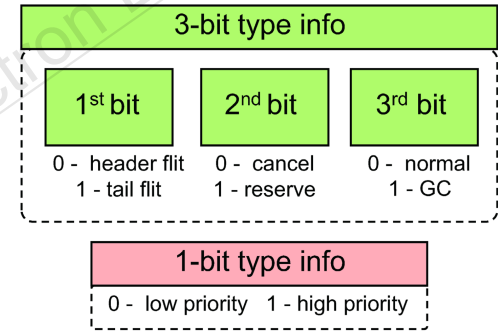
Packet ID	Entry port	Exit port	Link state
6	10	11	00
2	01	00	10
⋮	⋮	⋮	⋮

Design

Adaptive routing algorithm with priority



Structure of the scout packet

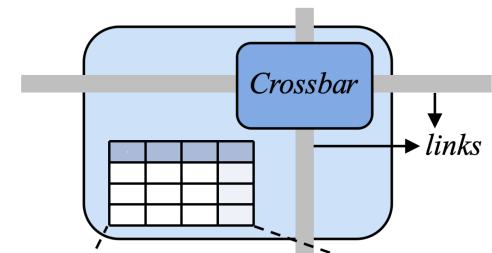


- Probe packet structure

Composed of header flit + tail flit

- Key modifications

- Reuse Venice's 1-bit reserved field → as priority flag
- 3-bit field for packet type encoding
- Include path reservation table (inherited from Venice)
- Newly added 2-bit link status field
- Indicate current link state
- Embed priority information for scheduling/preemption decisions

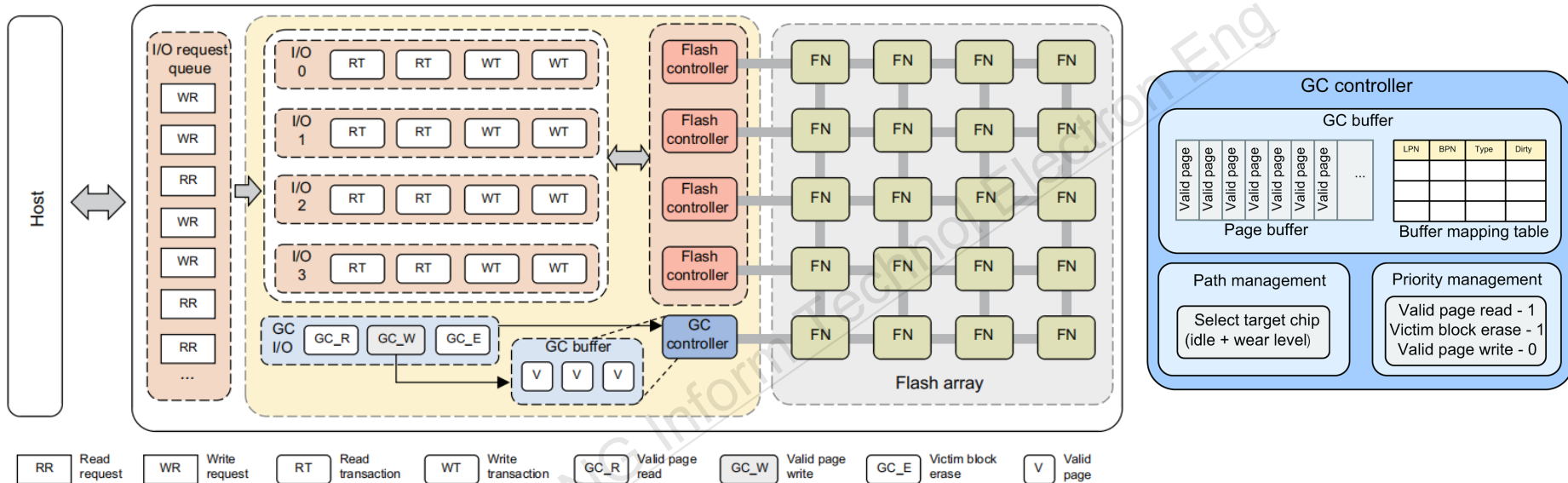


Router Reservation Table

Packet ID	Entry Port	Exit Port	Link State
6	10	11	00
2	01	00	10
⋮	⋮	⋮	⋮

Design

Design of GC bypass



Structure of the GC controller

GC controller architecture

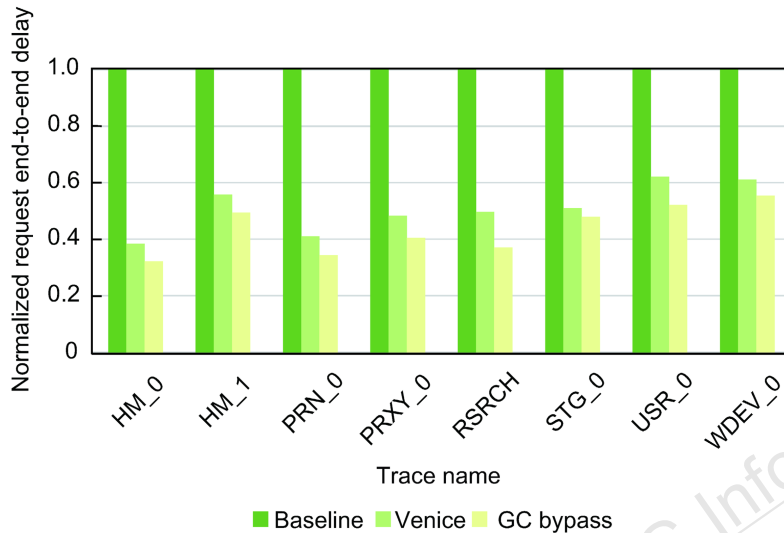
- Cache management module
- Priority management module
- Path management module

Preemption policy

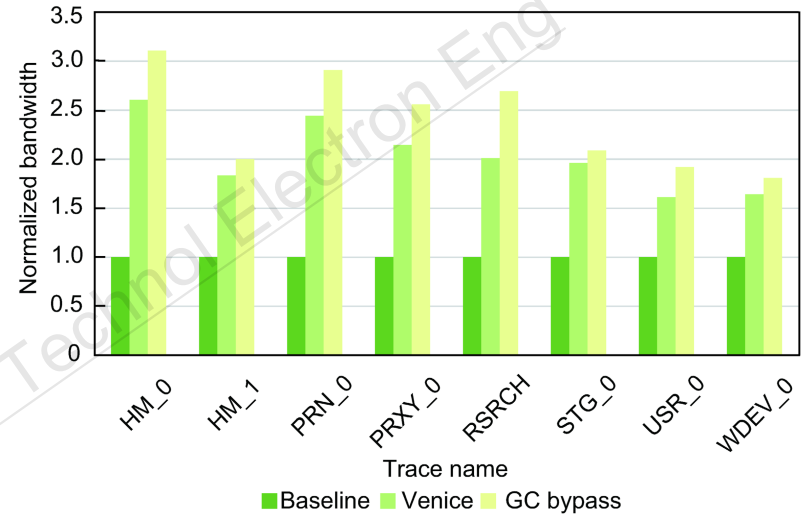
- High-priority: I/O requests, GC reads, GC erases
- Can preempt low-priority GC write reserved paths

Experiments

Experimental results analysis



Normalized request end-to-end delay

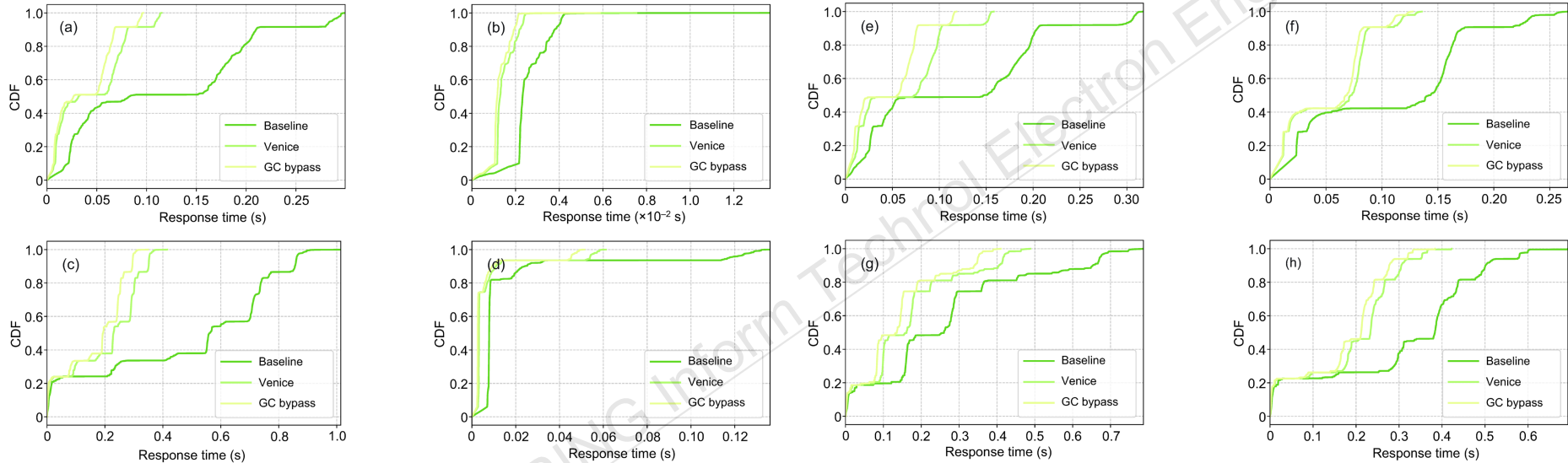


Normalized bandwidth

- GC bypass achieves up to 62% lower average end-to-end request latency compared to Baseline. GC bypass achieves up to 26% lower average end-to-end request latency compared to Venice.
- GC bypass achieves the highest bandwidth across all workloads, showing improvements of up to 3.1 and 1.3 times compared to Baseline and Venice, respectively.

Experiments

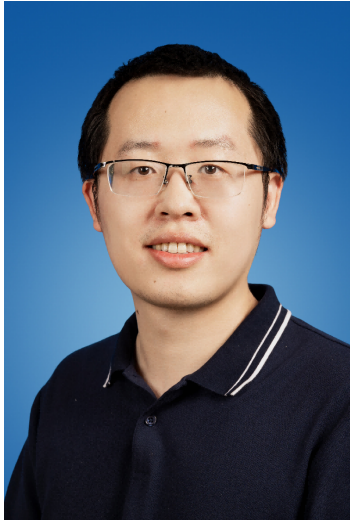
Experimental results analysis



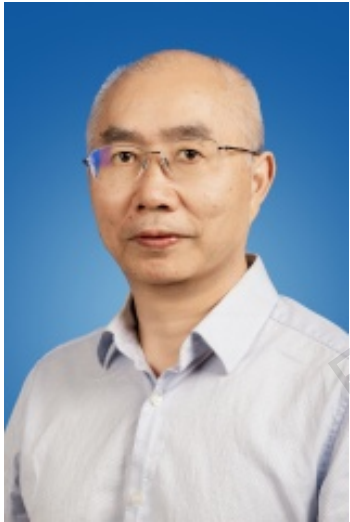
GC bypass maintains the lowest latency distribution across all workloads. Under the RSRCH workload, GC bypass reduces the tail latency of 99.99% I/O requests by 62.7% and 25.2% compared to Baseline and Venice, respectively.

Conclusions

We propose GC bypass to optimize the long-tail latency caused by page migration during GC in SSDs. By introducing a dedicated GC controller, GC requests are transmitted through independent interfaces, and the process of writing valid pages is assigned low priority, allowing high-priority requests to preempt path reservations from low-priority ones. Experimental results indicate that GC bypass reduces the 99.99th percentile long-tail latency by up to 25% compared to Venice.



Shiqiang NIE is a member of the China Computer Federation. He received his Ph.D. degree in computer science and technology from Xi'an Jiaotong University, China, in 2021. Currently, he is an Associate Researcher with the Department of Computer Science and Technology, Xi'an Jiaotong University. His research is supported by the National Natural Science Foundation of China. His research interests include optimizations for non-volatile memory and architecture, distributed storage systems, and cloud computing systems.



Weiguo WU received the B.S., M.S., and Ph.D. degrees in computer science and technology from Xi'an Jiaotong University, China, in 1986, 1993, and 2006, respectively. He is currently a Professor with the School of Computer Science and Technology, Xi'an Jiaotong University. His research interests include high performance computer architecture, storage system, cloud computing, and embedded system.