

# Venice

## Improving Solid-State Drive Parallelism at Low Cost via Conflict-Free Accesses

**Rakesh Nadig\***, Mohammad Sadrosadati\*, Haiyu Mao,  
Nika Mansouri Ghiasi, Arash Tavakkol, Jisung Park,  
Hamid Sarbazi-Azad, Juan Gómez Luna, and Onur Mutlu



# Talk Outline

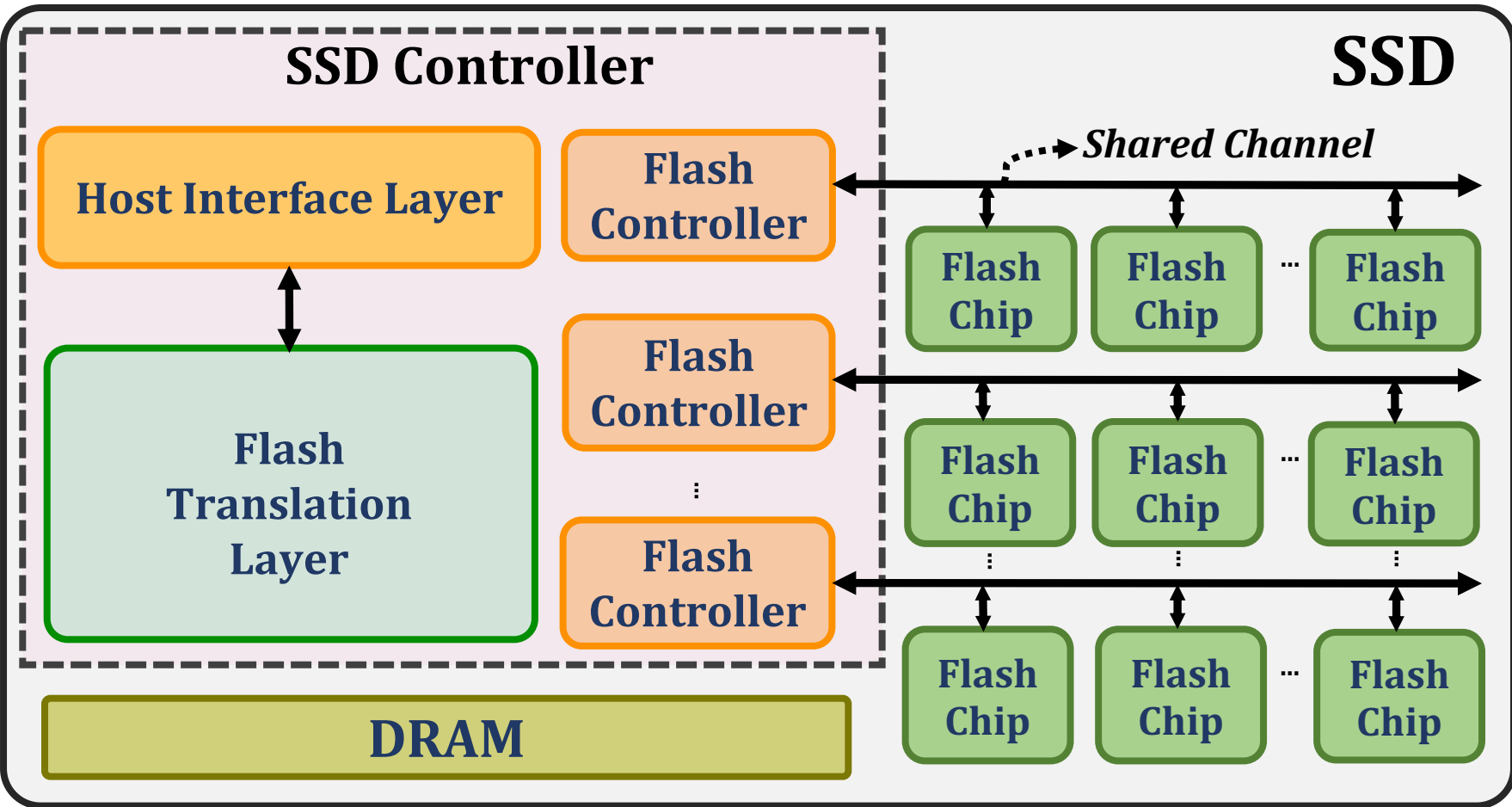
**Motivation**

Venice

**Evaluation**

**Summary**

# Overview of a Modern Solid-State Drive



# Key Problem: Path Conflicts in Modern SSDs

Multiple flash memory chips are connected to the SSD Controller using a shared channel



I/O requests attempt to simultaneously access the flash chips using a single path → *Path Conflict*



*Path Conflicts* cause I/O requests to be transferred serially on the shared channel



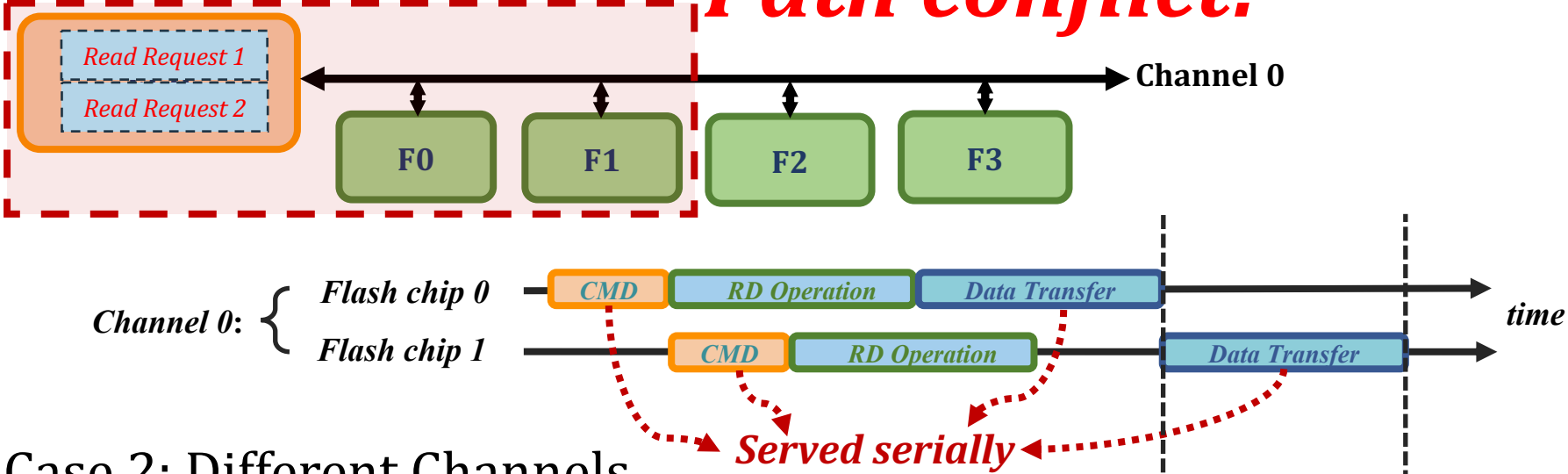
Limits SSD parallelism and reduces performance



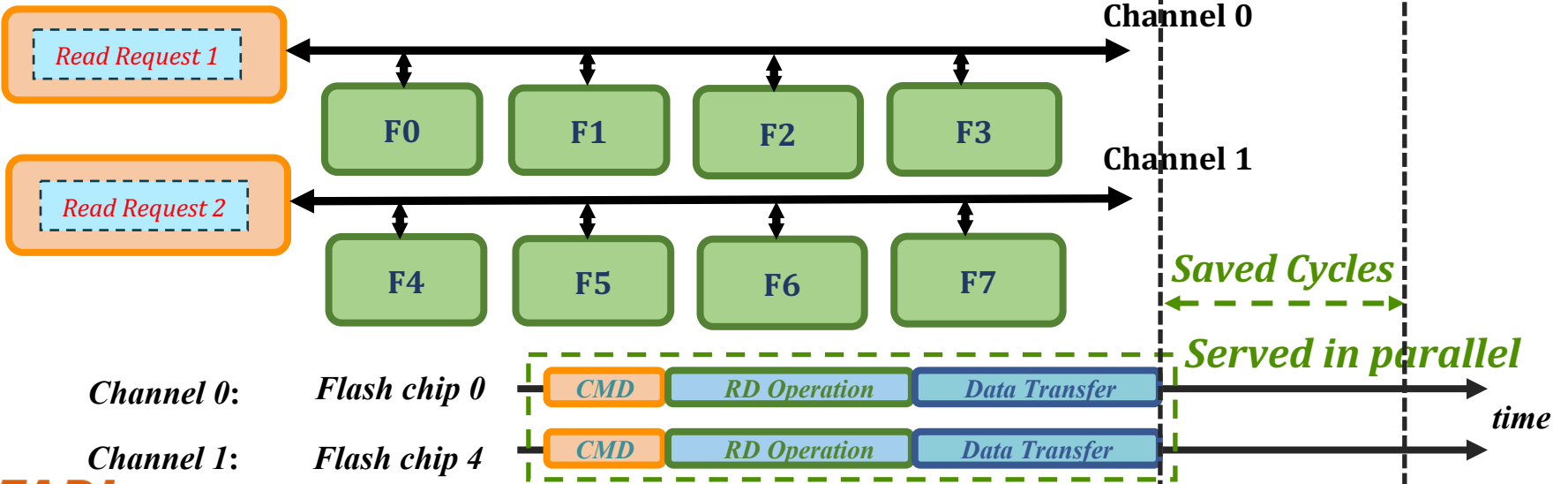
# Delay Caused by Path Conflicts

- Case 1: Same Channel

*Path conflict!*



- Case 2: Different Channels



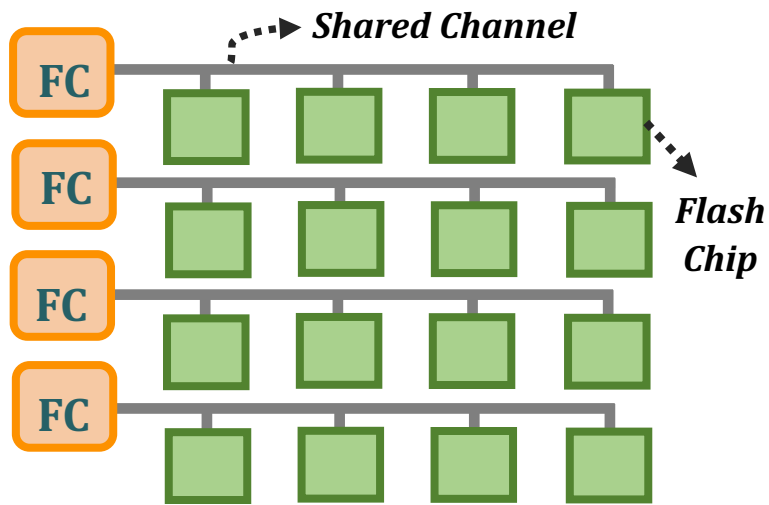
# Performance Impact of Path Conflicts

Path conflicts increase the average I/O latency by 57% in our experiments on a performance-optimized SSD

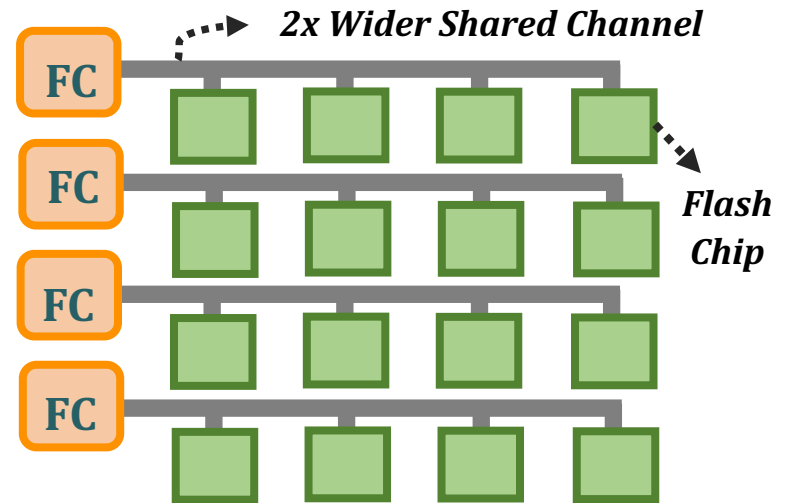
The performance overhead of path conflicts increases by 1.6x in our experiments for high-I/O-intensity workloads

# Prior Approaches

*Baseline SSD*



*Packetized SSD (pSSD) [1]*



# Prior Approaches

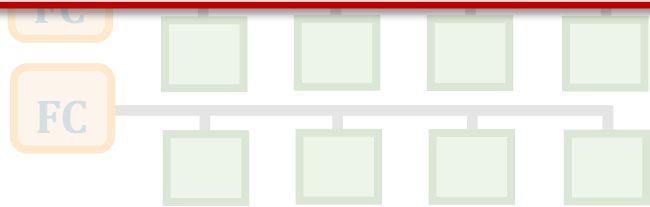
*Baseline SSD*

...► *Shared Channel*

*Packetized SSD (pSSD) [1]*

...► *2x Wider Shared Channel*

**Baseline SSD and Packetized SSD  
do *not* provide path diversity to flash chips**



# Prior Approaches

Baseline SSD

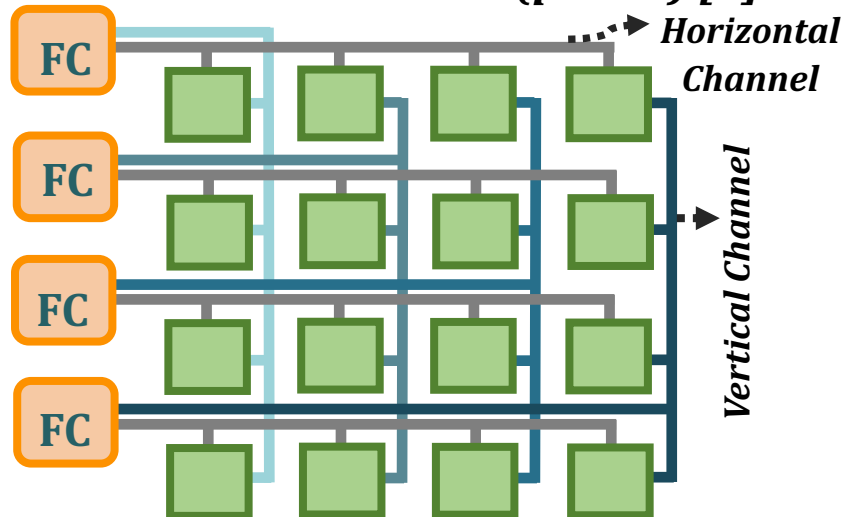
Shared Channel

Packetized SSD (pSSD) [1]

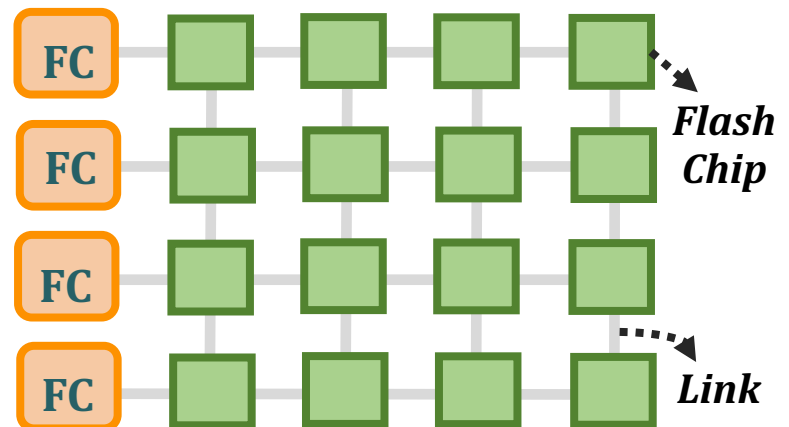
2x Wider Shared Channel

Baseline SSD and Packetized SSD  
do *not* provide path diversity to flash chips

Packetized Network SSD (pnSSD) [1]



Network-on-SSD (NoSSD) [2]



[1] Kim+, "Networked SSD: Flash Memory Interconnection Network for High-Bandwidth SSD", MICRO 2022

[2] Tavakkol+, "Network-on-SSD: A Scalable and High-Performance Communication Design Paradigm for SSDs", IEEE CAL 2012

# Prior Approaches

Baseline SSD

Shared Channel

Packetized SSD (pSSD) [1]

2x Wider Shared Channel

Baseline SSD and Packetized SSD  
do *not* provide path diversity to flash chips

Packetized Network SSD (pnSSD) [1]

Network-on-SSD (NoSSD) [2]

Packetized Network SSD and Network-on-SSD

1. do *not* effectively utilize the path diversity
2. incur large area & cost overheads

[1] Kim+, "Networked SSD: Flash Memory Interconnection Network for High-Bandwidth SSD", MICRO 2022

[2] Tavakkol+, "Network-on-SSD: A Scalable and High-Performance Communication Design Paradigm for SSDs", IEEE CAL 2012

# Our Goal

To fundamentally address the path conflict problem in SSDs by

1. increasing the number of paths to each flash chip (i.e., path diversity) at low cost
2. effectively utilizing the increased path diversity for communication between the SSD controller and flash chips

# Talk Outline

Motivation

Venue

Evaluation

Summary



# Our Proposal



## Venice



A low-cost interconnection network of flash chips in the SSD



Conflict-free path reservation for each I/O request



A non-minimal fully-adaptive routing algorithm for path identification

*Named after the network of canals in the city of Venice*  
<https://en.wikipedia.org/wiki/Venice>

# Our Proposal



## Venice



A low-cost interconnection network of flash chips in the SSD



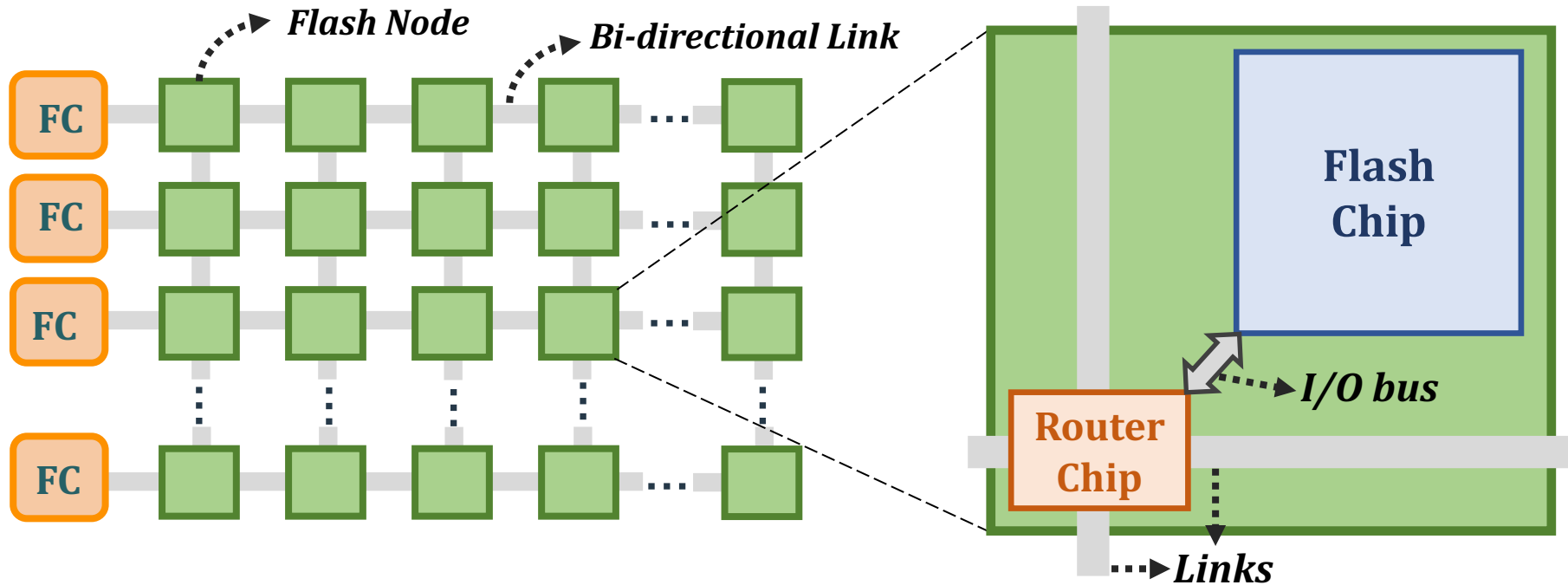
Conflict-free path reservation for each I/O request



A non-minimal fully-adaptive routing algorithm for path identification

*Named after the network of canals in the city of Venice*  
<https://en.wikipedia.org/wiki/Venice>

# Venice: Architecture



Venice provides increased path diversity at low cost

No modifications to existing flash chips in Venice

# Our Proposal



## Venice



A low-cost interconnection network of flash chips in the SSD



Conflict-free path reservation for each I/O request



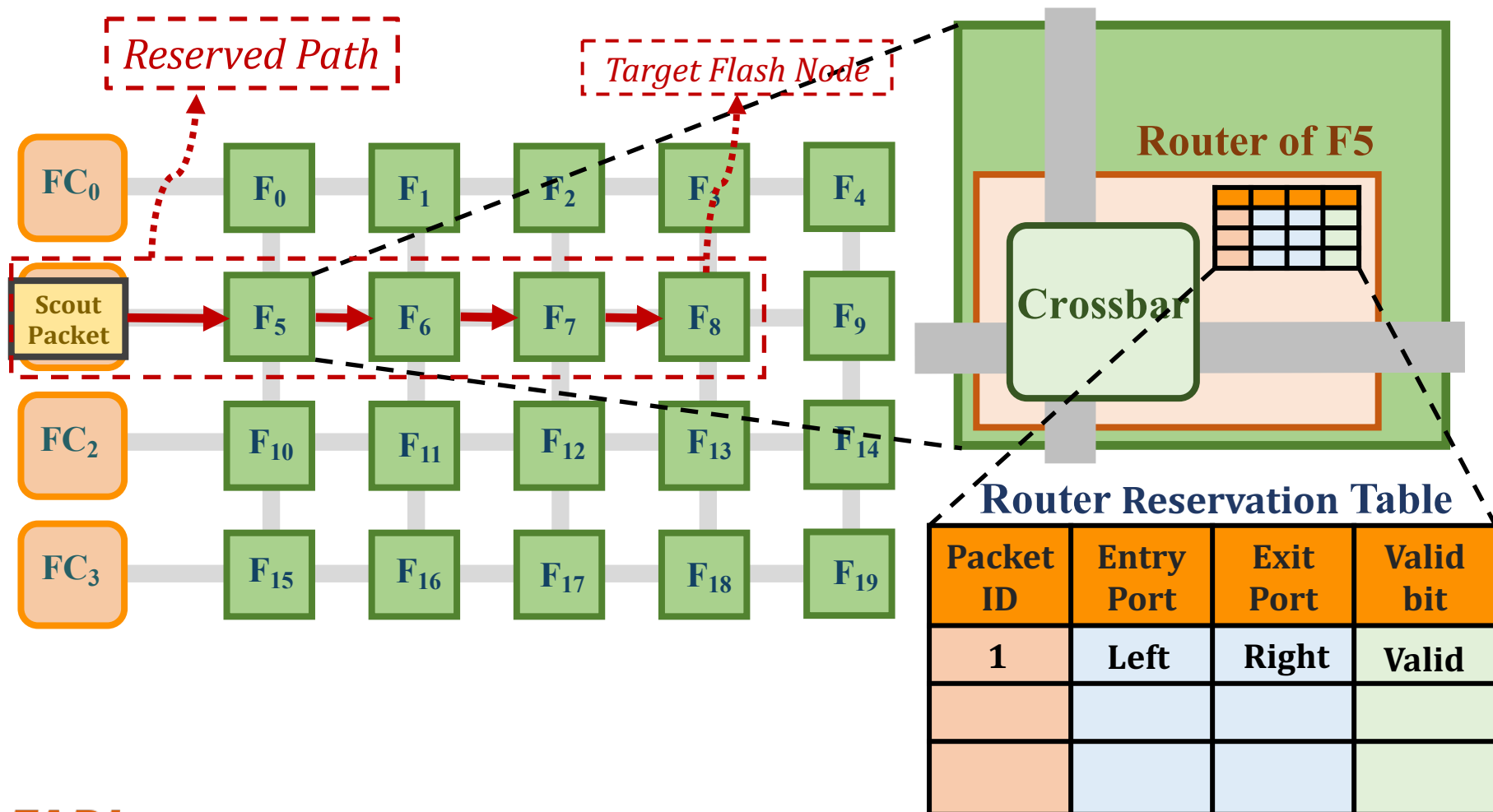
A non-minimal fully-adaptive routing algorithm for path identification

*Named after the network of canals in the city of Venice*  
<https://en.wikipedia.org/wiki/Venice>



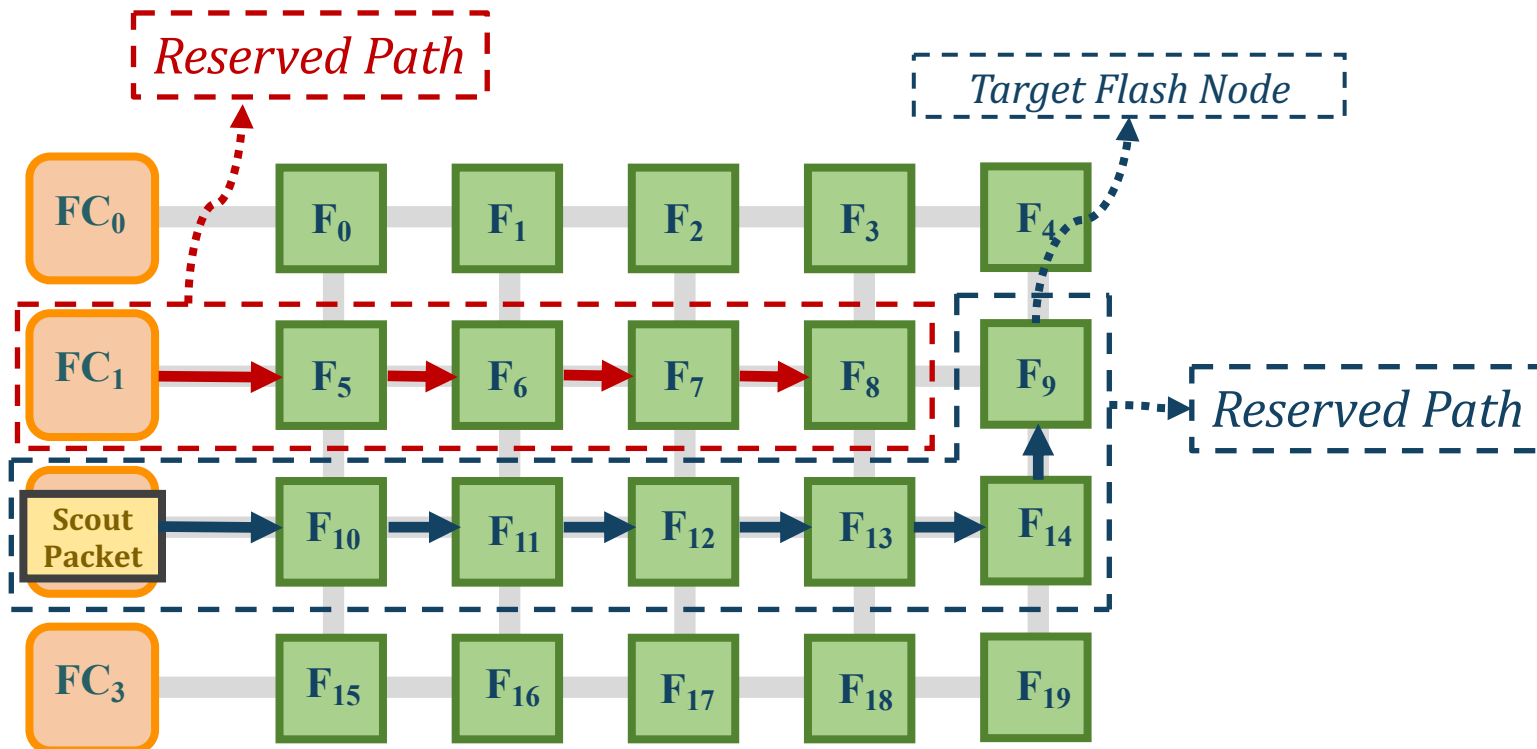
# Venice: Path Reservation (I)

- Venice uses a small *scout packet* to reserve a conflict-free path for each I/O request



# Venice: Path Reservation (II)

- Venice uses a small *scout packet* to reserve a conflict-free path for each I/O request





# Our Proposal



## Venice



A low-cost interconnection network of flash chips in the SSD



Conflict-free path reservation for each I/O request



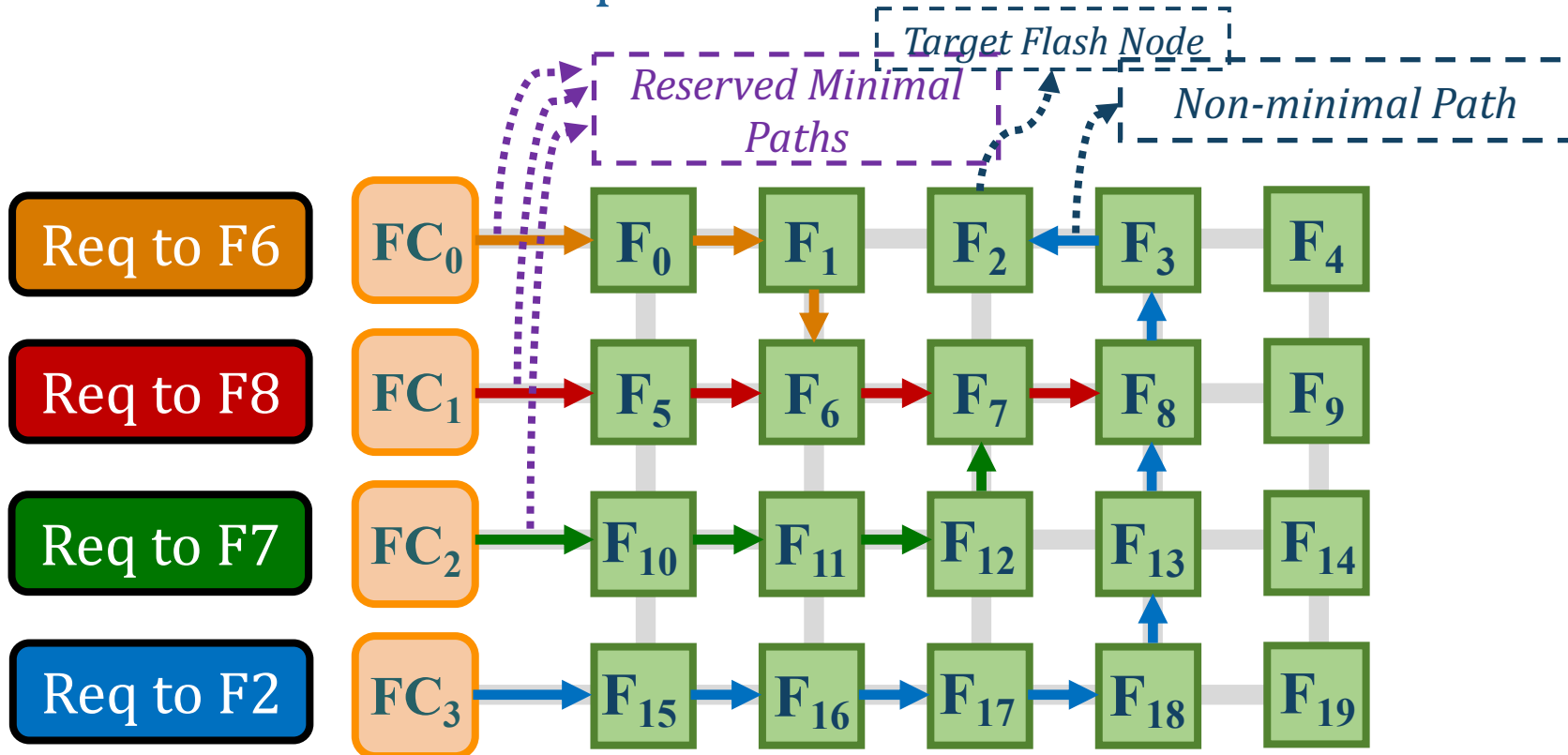
A non-minimal fully-adaptive routing algorithm for path identification

*Named after the network of canals in the city of Venice*  
<https://en.wikipedia.org/wiki/Venice>



# Venice: Non-Minimal Fully Adaptive Routing

- Venice uses a **non-minimal fully-adaptive routing algorithm** to route *scout packets* when a **minimal path** is unavailable
- **Effectively utilizes the idle links** in the interconnection network to find a **conflict-free path**



# More in the Paper

- Venice's **non-minimal fully-adaptive routing algorithm**
- **Handling deadlock and livelock** scenarios
- **Overhead of exercising a non-minimal path**
- **Analysis of prior architectures** proposed to mitigate the **path conflict problem**
- Detailed background on modern SSD architecture

# Talk Outline

Motivation

Venice

**Evaluation**

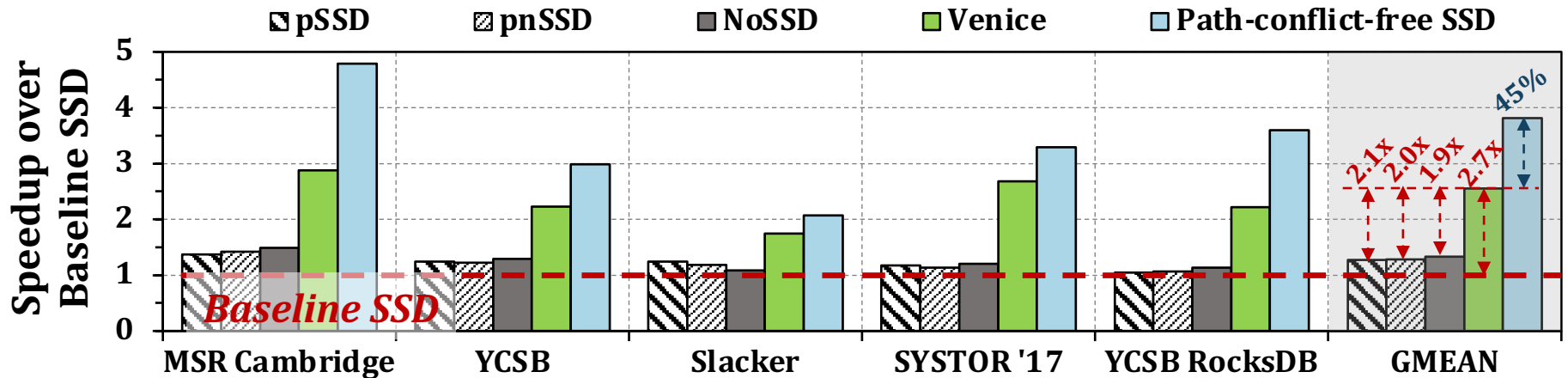
Summary

# Evaluation Methodology

- **Using MQSim [Tavakkol+, FAST'18]**, a state-of-the-art SSD simulator
- **Two SSD configurations**
  - **Performance-Optimized** (Samsung Z-NAND SSD)
  - **Cost-Optimized** (Samsung PM9A3)
- **Nineteen data-intensive workloads from**
  - MSR Cambridge, YCSB, Slacker, SYSTOR '17 and RocksDB
- **Prior Approaches**
  - **Baseline SSD**: A typical multi-channel shared bus SSD
  - **Packetized SSD (pSSD)** [Kim+, MICRO'22]: Uses packetization to double the flash channel bandwidth
  - **Packetized Network SSD (pnSSD)** [Kim+, MICRO'22]: Increases path diversity by introducing vertical channels
  - **Network-on-SSD (NoSSD)** [Tavakkol+, CAL 2012]: Proposes an interconnection network of flash chips with simple deterministic routing
  - **Path-conflict-free SSD**: An *ideal SSD* with no path conflicts

# Results: Performance Analysis (I)

- Performance-Optimized SSD

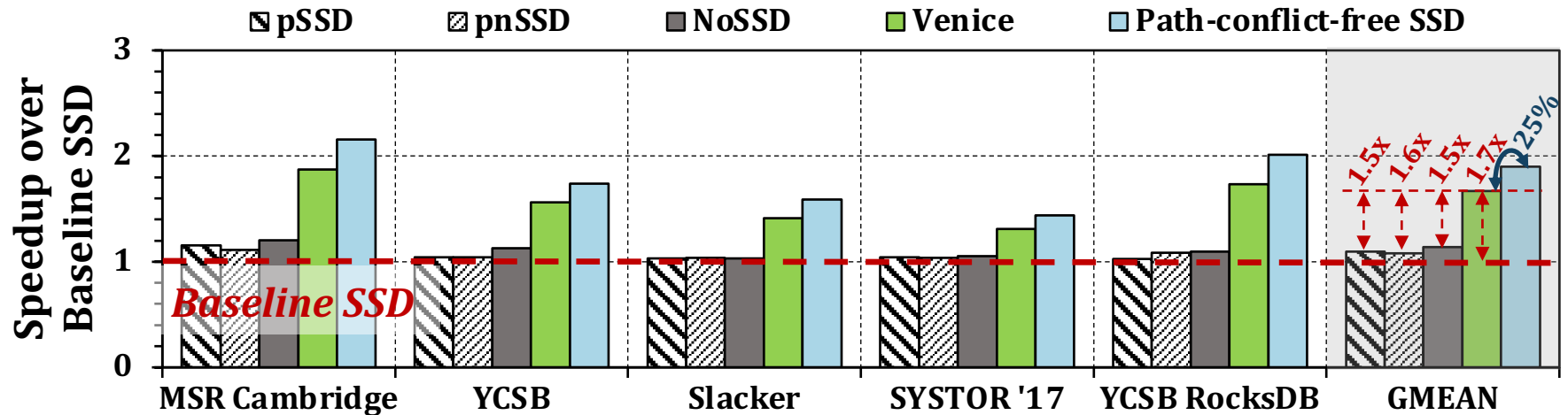


Venice improves SSD performance by 1.9x on average over the best-performing prior work

Venice's performance is within 45% of the performance of a Path-conflict-free SSD

# Results: Performance Analysis (II)

- Cost-Optimized SSD

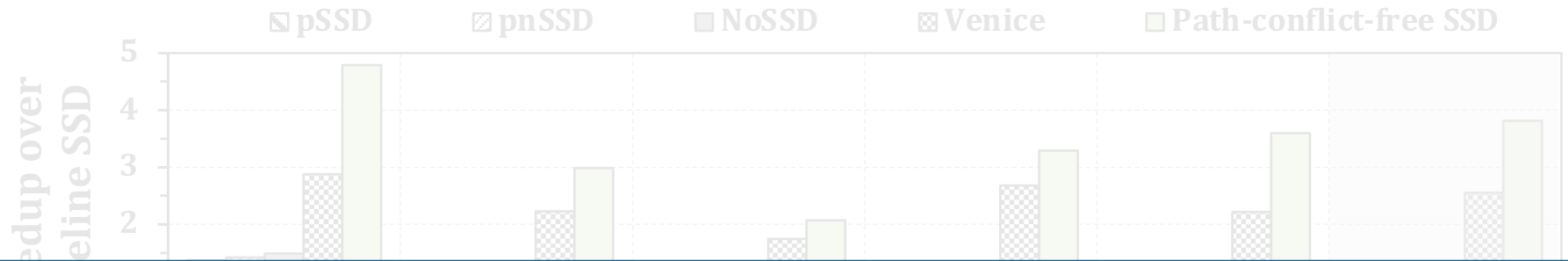


Venice improves SSD performance by 1.5x on average over the best-performing prior work

Venice's performance is within 25% of the performance of a Path-conflict-free SSD

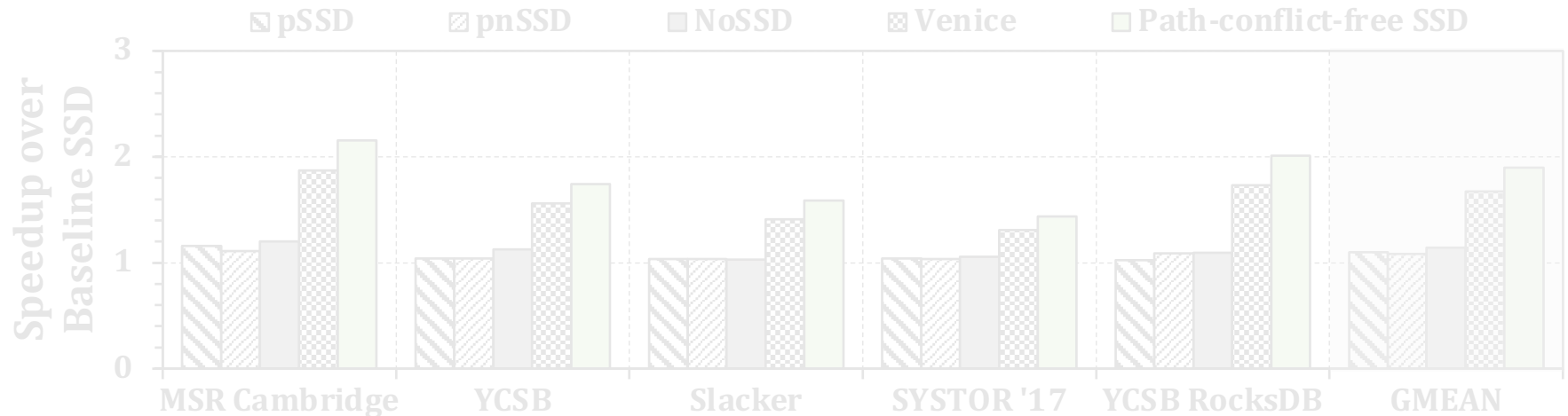
# Results: Performance Analysis (III)

- Performance-Optimized SSD

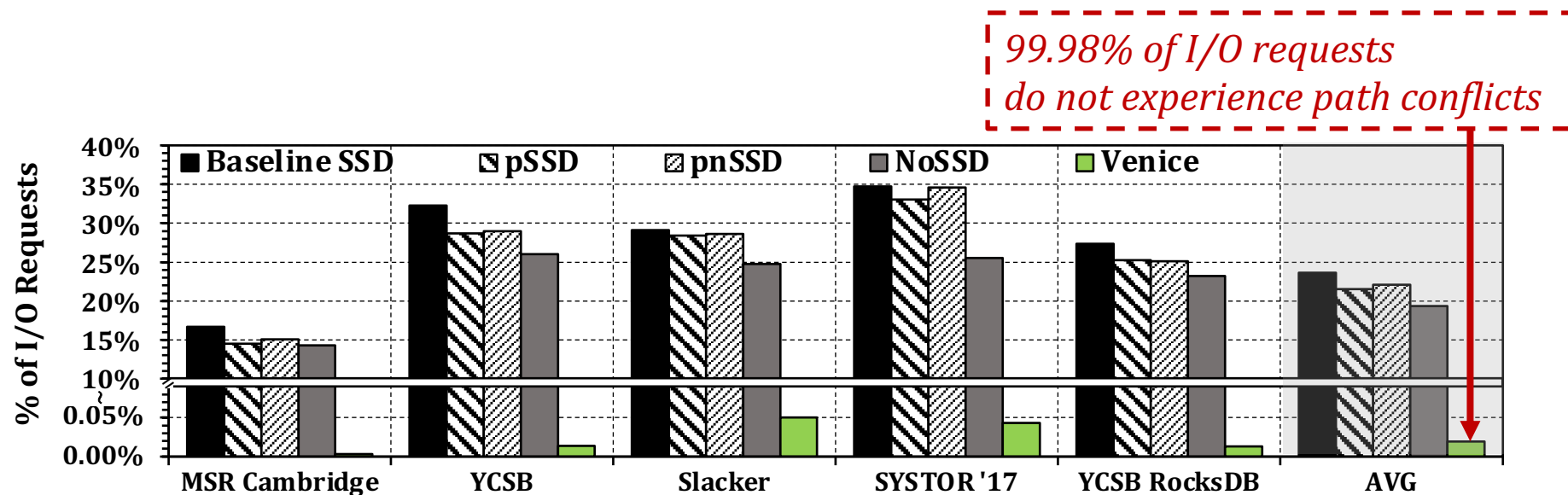


Venice provides significant improvement in performance over all prior approaches

- Cost-Optimized SSD



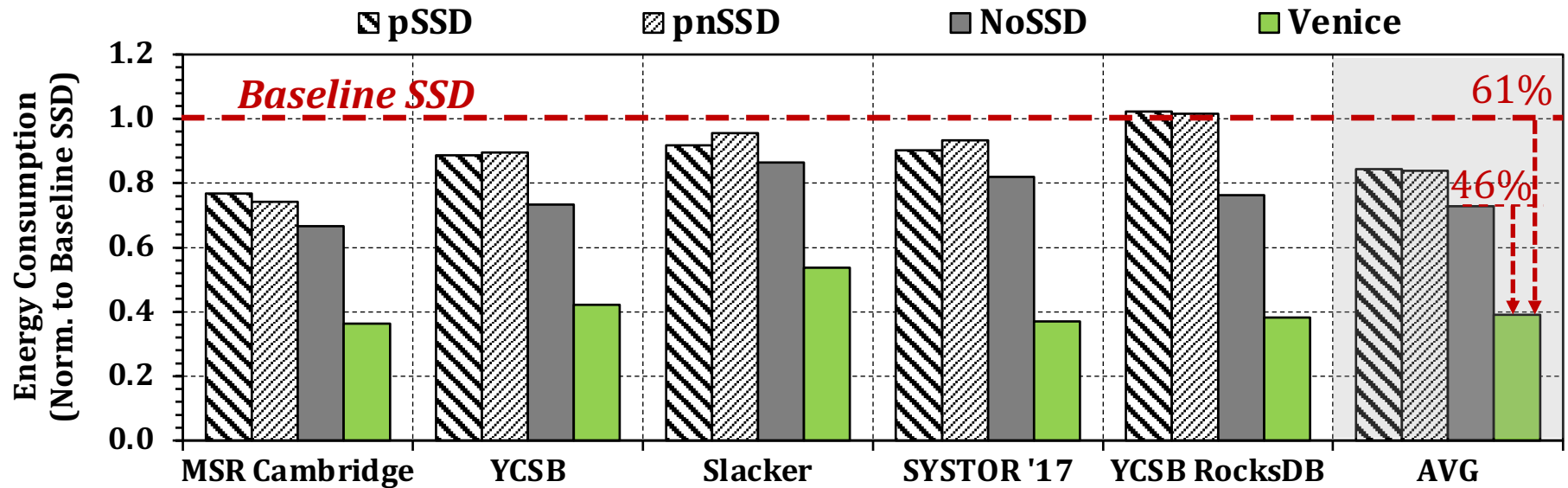
# Results: Reduction in Path Conflicts



Venice mitigates path conflicts by using path reservation and effective utilization of path diversity



# Results: SSD Energy Consumption



Venice reduces the SSD energy consumption by 46% on average over the most efficient prior work

# More in the Paper

- Power and area overhead analysis
- Tail latency analysis
- Sensitivity to interconnection network configurations
- Performance on mixed workloads
- Detailed evaluation methodology

# More in the Paper

## Venice: Improving Solid-State Drive Parallelism at Low Cost via Conflict-Free Accesses

\*Rakesh Nadig<sup>§</sup>   \*Mohammad Sadrosadati<sup>§</sup>   Haiyu Mao<sup>§</sup>   Nika Mansouri Ghiasi<sup>§</sup>  
Arash Tavakkol<sup>§</sup>   Jisung Park<sup>§∇</sup>   Hamid Sarbazi-Azad<sup>†‡</sup>   Juan Gómez Luna<sup>§</sup>   Onur Mutlu<sup>§</sup>  
*§ETH Zürich*   *∇POSTECH*   *†Sharif University of Technology*   *‡IPM*



<https://arxiv.org/abs/2305.07768>

# Talk Outline

Motivation

Venue

Evaluation

Summary

# Venice: Summary



Mitigates path conflicts by efficiently utilizing the path diversity of the SSD interconnection network



Improves performance  
by 1.9x/1.5x over the best-performing prior work  
on performance-optimized/cost-optimized SSD



Reduces energy consumption  
by 46% on average over the most efficient prior work



Low-cost and requires  
no changes to commodity flash chips

# Venice



## Improving Solid-State Drive Parallelism at Low Cost via Conflict-Free Accesses

**Rakesh Nadig\***, Mohammad Sadrosadati\*, Haiyu Mao,  
Nika Mansouri Ghiasi, Arash Tavakkol, Jisung Park,  
Hamid Sarbazi-Azad, Juan Gómez Luna, and Onur Mutlu

<https://arxiv.org/abs/2305.07768>



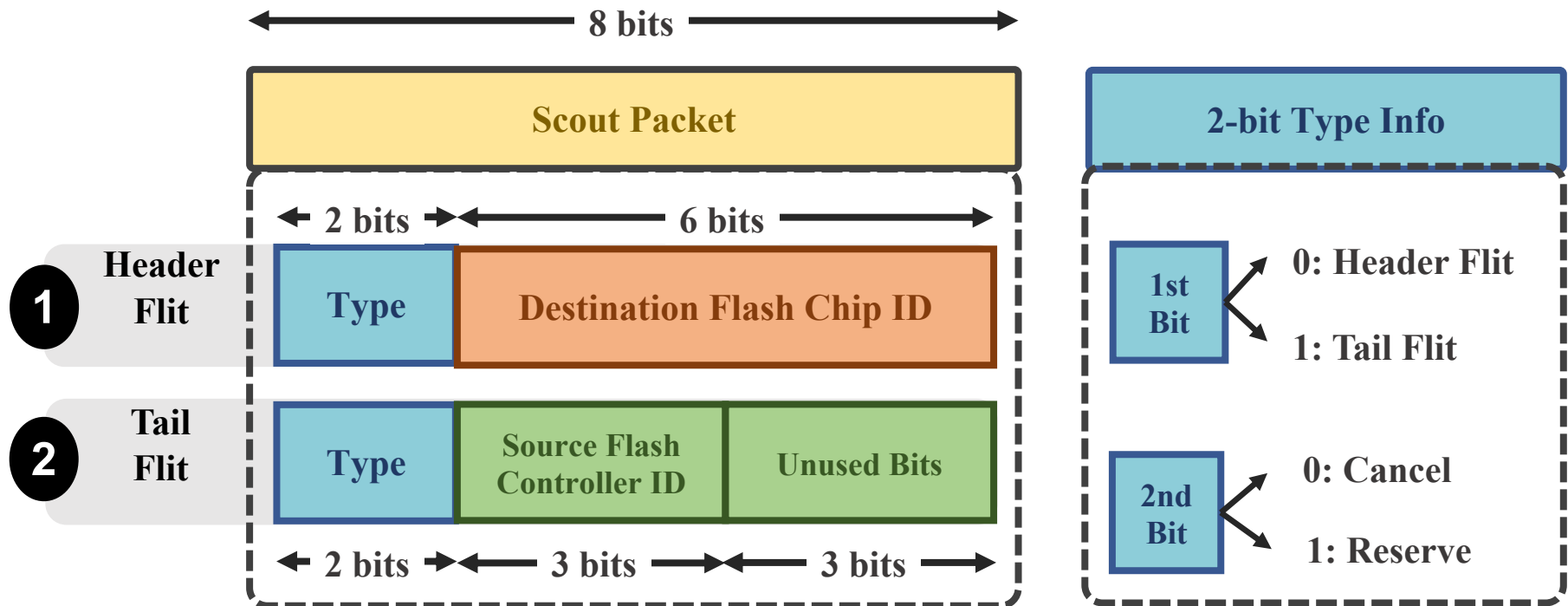
# FAQ

- Scout Packet Structure
- Power Overhead
- Area Overhead
- Evaluated Configurations
- Workload Characteristics
- Mixed Workloads
- Results
  - Throughput Analysis
  - Tail Latency
  - Power Consumption
  - Performance on Mixed Workloads
  - Sensitivity to Interconnection Network Configuration

# Structure of a Scout Packet



- A *scout packet* consists of **two 8 bit flits**, a **header flit** and a **tail flit**
- The **flash controller** sends a **scout packet** to **identify** a **conflict-free path** for the I/O request







- **Router**

- We implement the HDL and synthesize it using UMC 65nm technology node
- Router consumes 0.241mW for a 4KB page transfer

- **Network Link**

- ORION 3.0 power model tool
- Each network link consumes about 1.08mW for a 4KB page transfer
- **Link capacitance is lower than bus capacitance** -> 90% less power than that of the shared channel bus
  - Links are shorter and thinner than a shared bus
  - Two drivers in links compared to several drivers in a bus

| Component | # of Instances         | Avg. Power [mW]<br>for 4KB page transfer | Area                     |
|-----------|------------------------|--|--------------------------|
| Router    | 1 per flash node       | 0.241                                    | 8% of flash chip area    |
| Link      | Up to 4 per flash node | 1.08                                     | 0.04× flash channel area |



- Router
  - Area overhead estimated using router's HDL model
  - Each router has
    - an area of  $614 \mu\text{m}^2 + 40 \text{ I/O}$
    - A total area of  $8\text{mm}^2 \rightarrow 8\%$  of a typical  $100\text{mm}^2$  flash chip
- Network Link
  - ORION 3.0 model for area analysis of network links
  - 112 network links for a  $8 \times 8$  flash array configuration
  - 44% lower area than a baseline multi-channel shared bus architecture
  - **Links are thinner and require lower pitch sizes**

# Evaluated Configurations



|   |   |
|---|---|
| <b>Performance-optimized<br/>SSD [31, 99]</b> | 240GB, Z-NAND [31, 99, 119],<br>8-GB/s External I/O bandwidth (4-lane PCIe Gen4);<br>1.2-GB/s Flash Channel I/O rate  |
|   | <b>NAND Config:</b> 8 channels, 8 chips/channel,<br>1 die/chip, 2 planes/die, 128Gb die capacity,<br>1024 blocks/plane, 768 pages/block, 4KB page   |
|   | <b>Latencies:</b> Read( $t_R$ ): $3\mu s$ ; Erase ( $t_{BERS}$ ): $1ms$<br>Program ( $t_{PROG}$ ): $100\mu s$   |
| <b>Cost-optimized<br/>SSD [55]</b>            | 1TB, 3D TLC NAND Flash,<br>8-GB/s External I/O bandwidth (4-lane PCIe Gen4);<br>1.2-GB/s Flash Channel I/O rate   |
|   | <b>NAND Config:</b> 8 channels, 8 chips/channel,<br>1 die/chip, 2 planes/die, 1024 blocks/die, 16KB page  |
|   | <b>Latencies:</b> Read ( $t_R$ ): $45\mu s$ ; Erase ( $t_{BERS}$ ): $3.5ms$<br>Program ( $t_{PROG}$ ): $650\mu s$   |
| <b>Venice Design Parameters</b>               | <b>Topology.</b> $8 \times 8$ 2D mesh topology, 8-bit 1 GHz links,<br>One router next to each flash chip<br><b>Router Architecture.</b> Two 8-bit buffers per port,<br>1 GHz frequency<br><b>Routing Algorithm.</b> Non-minimal fully-adaptive<br><b>Switching.</b> Circuit switching [102] |

# Workload Characteristics



|                     | Traces   | Read % | Avg. Request Size (KB) | Avg. Inter-request Arrival Time ( $\mu$ s) |
|---------------------|----------|--------|------------------------|--|
| MSR Cambridge [122] | hm_0     | 36     | 8.8                    | 58   |
|                     | mds_0    | 12     | 9.6                    | 268  |
|                     | proj_3   | 95     | 9.6                    | 19   |
|                     | prxy_0   | 3      | 7.2                    | 242  |
|                     | rsrch_0  | 9      | 9.6                    | 129  |
|                     | src1_0   | 56     | 43.2                   | 49   |
|                     | src2_1   | 98     | 59.2                   | 50   |
|                     | usr_0    | 40     | 22.8                   | 98   |
|                     | wdev_0   | 20     | 9.2                    | 162  |
|                     | web_1    | 54     | 29.6                   | 67   |
| YCSB [123]          | YCSB_B   | 99     | 65.7                   | 13   |
|                     | YCSB_D   | 99     | 62                     | 14   |
| Slacker [124]       | jenkins  | 94     | 33.4                   | 615  |
|                     | postgres | 82     | 13.3                   | 382  |
| SYSTOR '17 [125]    | LUN0     | 76     | 20.4                   | 218  |
|                     | LUN2     | 73     | 16                     | 320  |
|                     | LUN3     | 7      | 7.7                    | 3127                                       |
| YCSB RocksDB [126]  | ssd-00   | 91     | 90                     | 5  |
|                     | ssd-10   | 99     | 11.5                   | 2  |

# Mixed Workloads

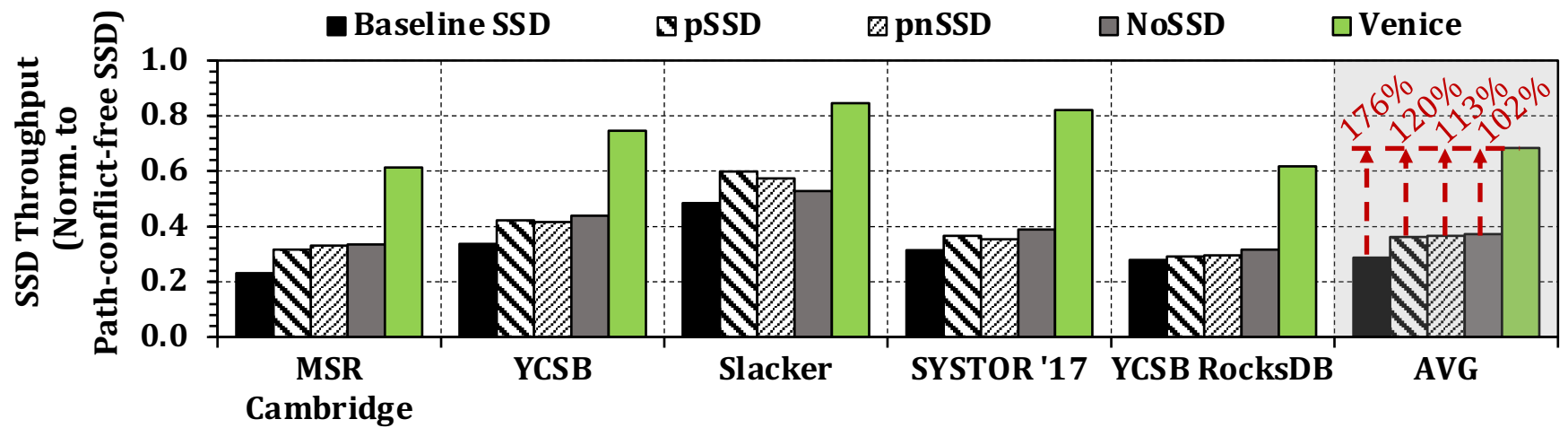


| <b>Mix</b>  | <b>Constituent Workloads [122, 123]</b> | <b>Description</b>   | <b>Avg. Inter-request Arrival Time (<math>\mu</math>s)</b> |
|-------------|---|--|--|
| <b>mix1</b> | src2_1 and proj_3                       | Both workloads are read-intensive  | 5.8  |
| <b>mix2</b> | src2_1, proj_3 and YCSB_D               | All three workloads are read-intensive   | 8.4  |
| <b>mix3</b> | prxy_0 and rsrch_0                      | Both workloads are write-intensive   | 93   |
| <b>mix4</b> | prxy_0, rsrch_0 and mds_0               | All three workloads are write-intensive  | 56   |
| <b>mix5</b> | prxy_0 and src2_1                       | prxy_0 is write-intensive and src2_1 is read-intensive                                     | 5  |
| <b>mix6</b> | prxy_0, src2_1 and usr_0                | prxy_0 is write-intensive, src2_1 is read-intensive and usr_0 has 60% writes and 40% reads | 3  |

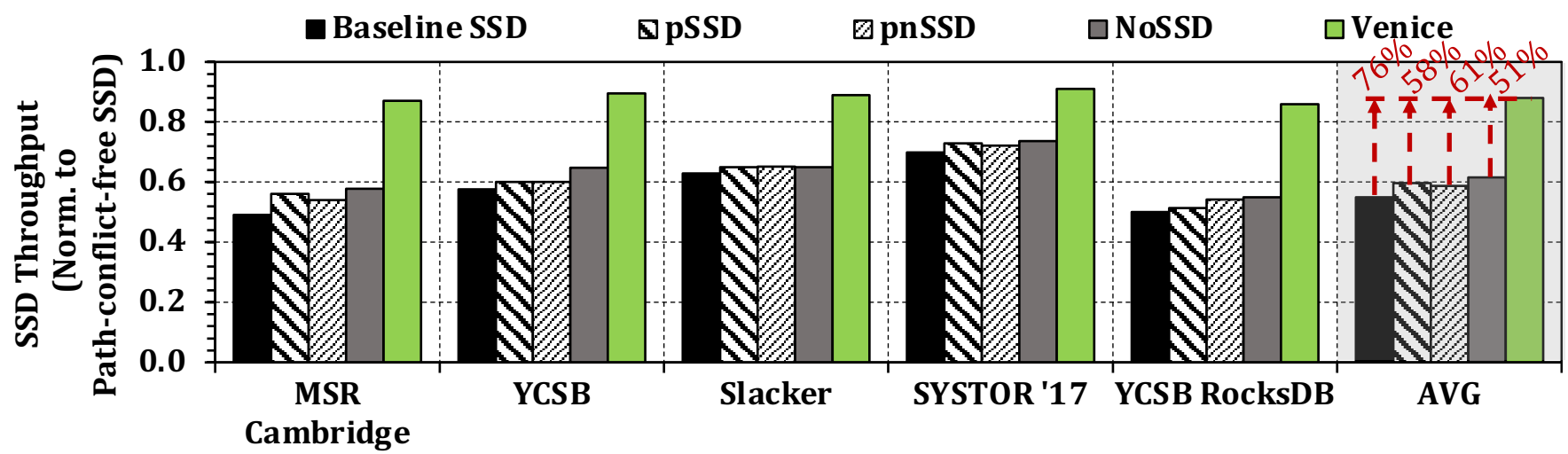
# SSD Throughput Analysis



- Performance-Optimized SSD



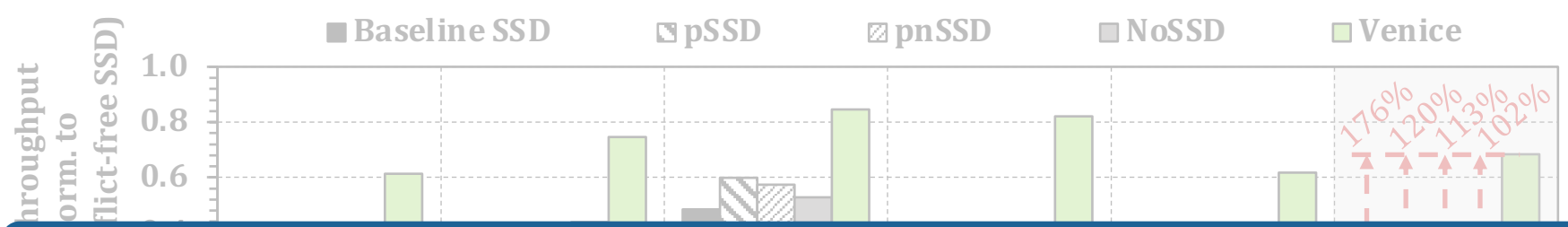
- Cost-Optimized SSD



# SSD Throughput Analysis



- Performance-Optimized SSD



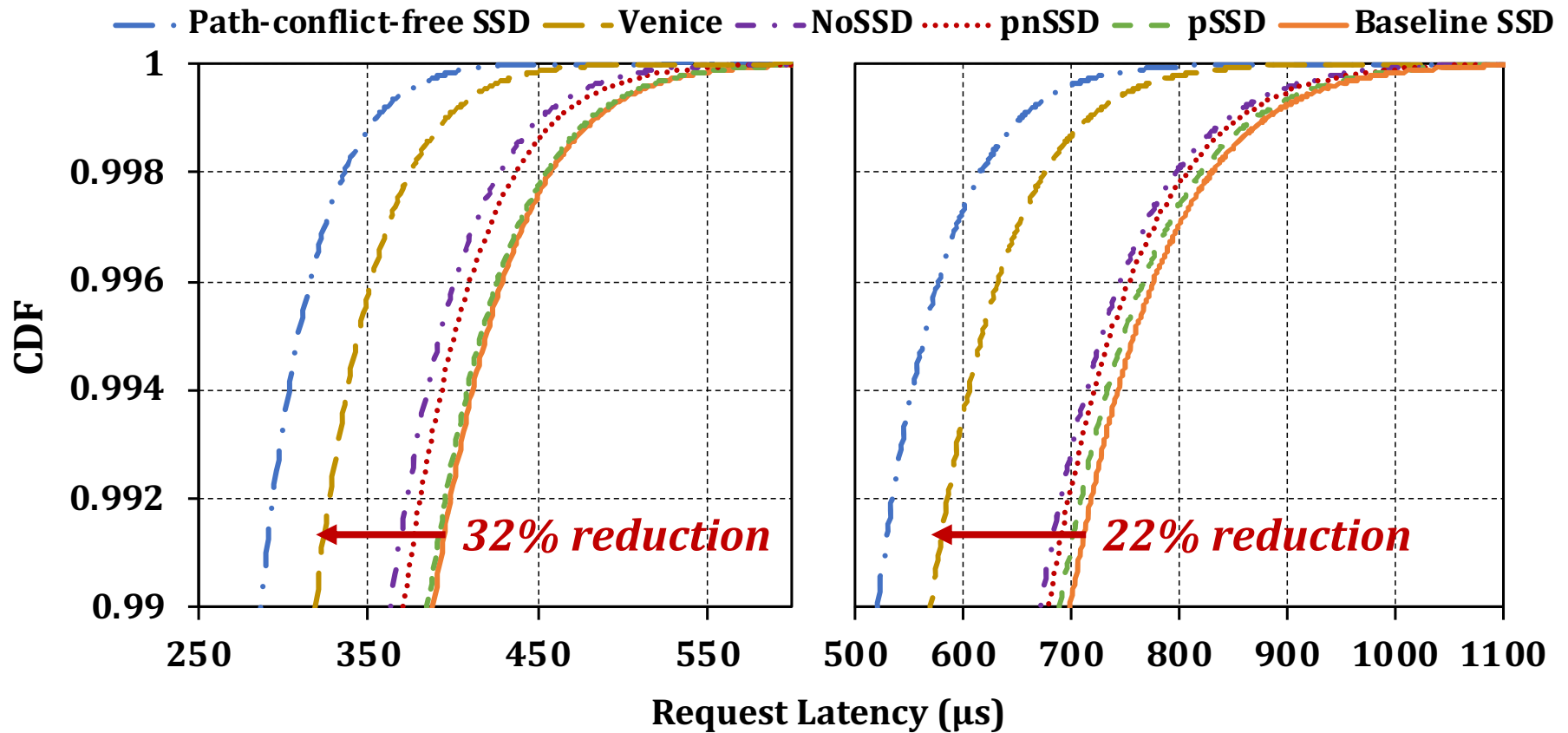
Venice improves SSD throughput over prior approaches by effectively mitigating path conflicts



# Tail Latency



- Comparison of tail latencies in the 99th percentile of I/O requests



(a) src1\_0

(b) hm\_0



# Tail Latency

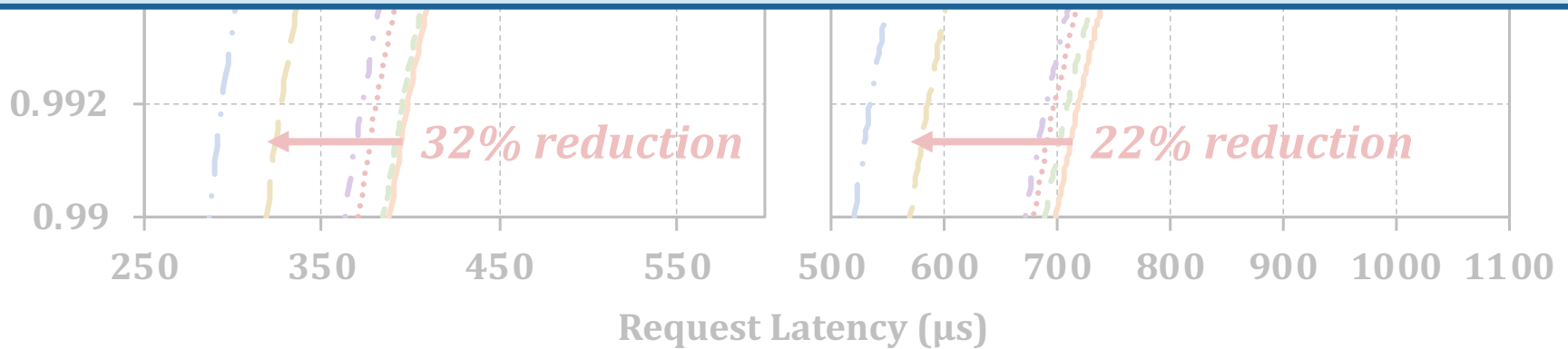


- Comparison of tail latencies in the 99th percentile of I/O requests

— Path-conflict-free SSD — Venice — NoSSD ..... pnSSD — pSSD — Baseline SSD

1

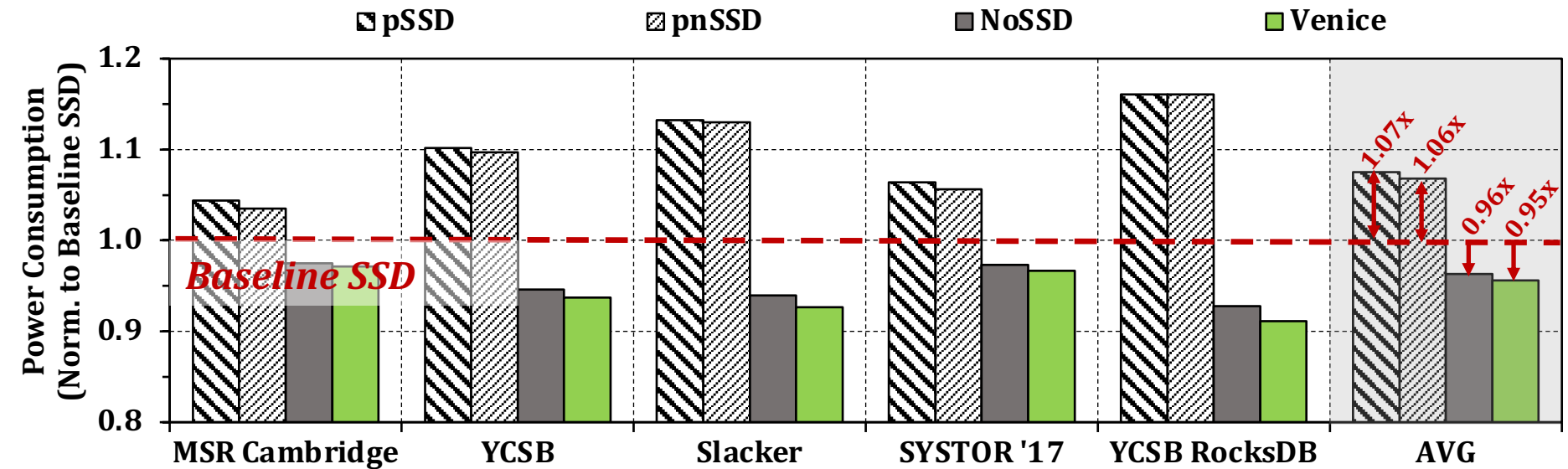
Venice reduces tail latencies  
by effectively mitigating path conflicts



(a) src1\_0

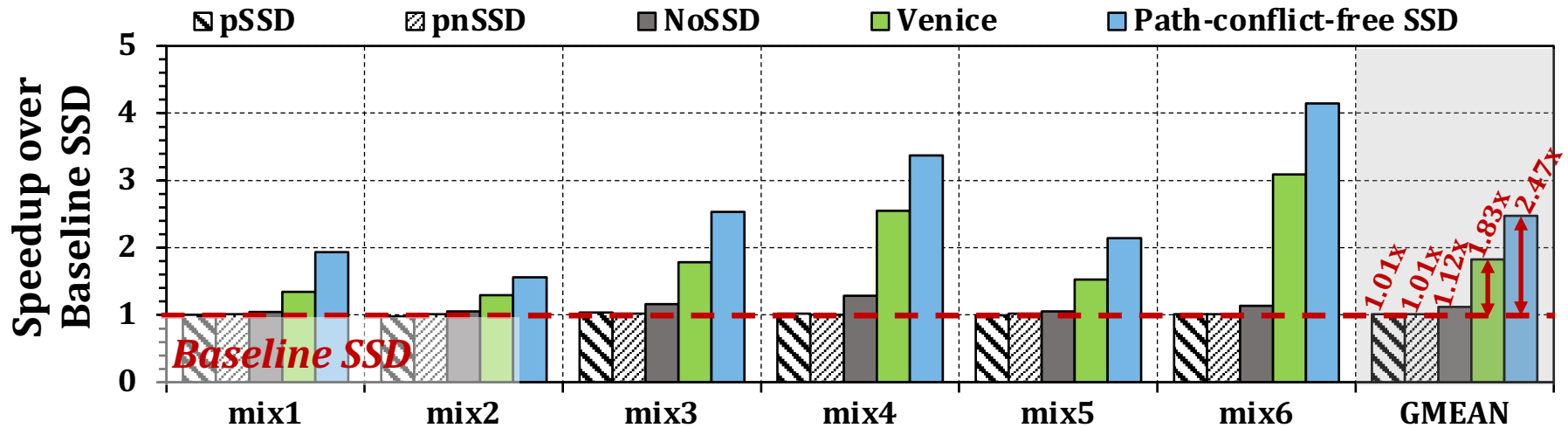
(b) hm\_0

# Power Consumption (II)



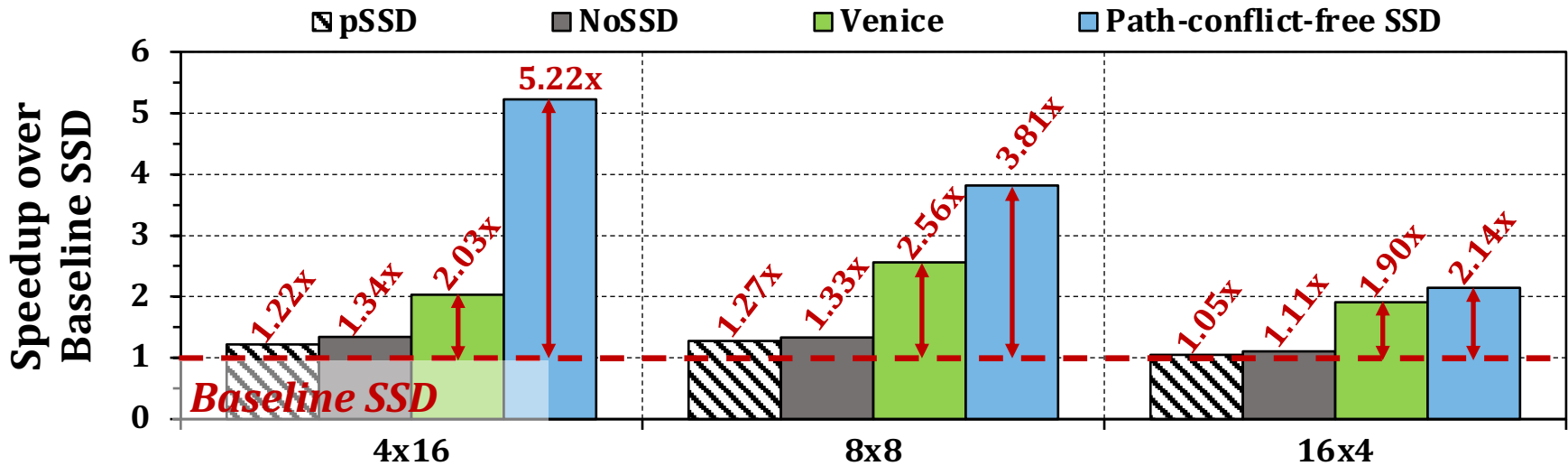
Venice reduces the average power consumption by 4% over Baseline SSD

# Performance on Mixed Workloads



Venice outperforms prior approaches on high-intensity mixed workloads by effectively mitigating path conflicts

# Sensitivity to Network Configurations



Venice provides higher performance improvement for 8x8 compared to 4x16 and 16x4

# Prior Approaches to Address Path Conflicts

- **Network-On-SSD [2]**

- Replaces a multi-channel shared bus architecture with an interconnection network of flash chips
- Significantly increases path diversity than a typical SSD

