# BTS: Exploring Effects of Background Task-Aware Scheduling for Key-Value CSDs

**Yeohyeon Park,** Chang-Gyu Lee, Seungjin Lee,
Inhyuk Park, Soonyeal Yang, Woosuk Chung, Youngjae Kim

7[th] International Parallel Data Systems Workshop (PDSW'22)

**SOGANG UNIVERSITY**

SK hynix

# Outline

- ❏ Background
    - ❏ Computational Storage Device (CSD)
    - ❏ Intel SPDK

- ❏ Motivation

- ❏ Proposed Architecture
    - ❏ BTS : Background Task-Aware Scheduler
    - ❏ Execution Flow

- ❏ Evaluation

- ❏ Conclusion and Q&A

# Background

# Computational Storage

❑ What is computational storage device (CSD)?

# Computational Storage

❑ What is computational storage device (CSD)?
  ➢ Computational storage devices (CSD) can run computational tasks inside the storage device, reducing data transfer between the host and the device.
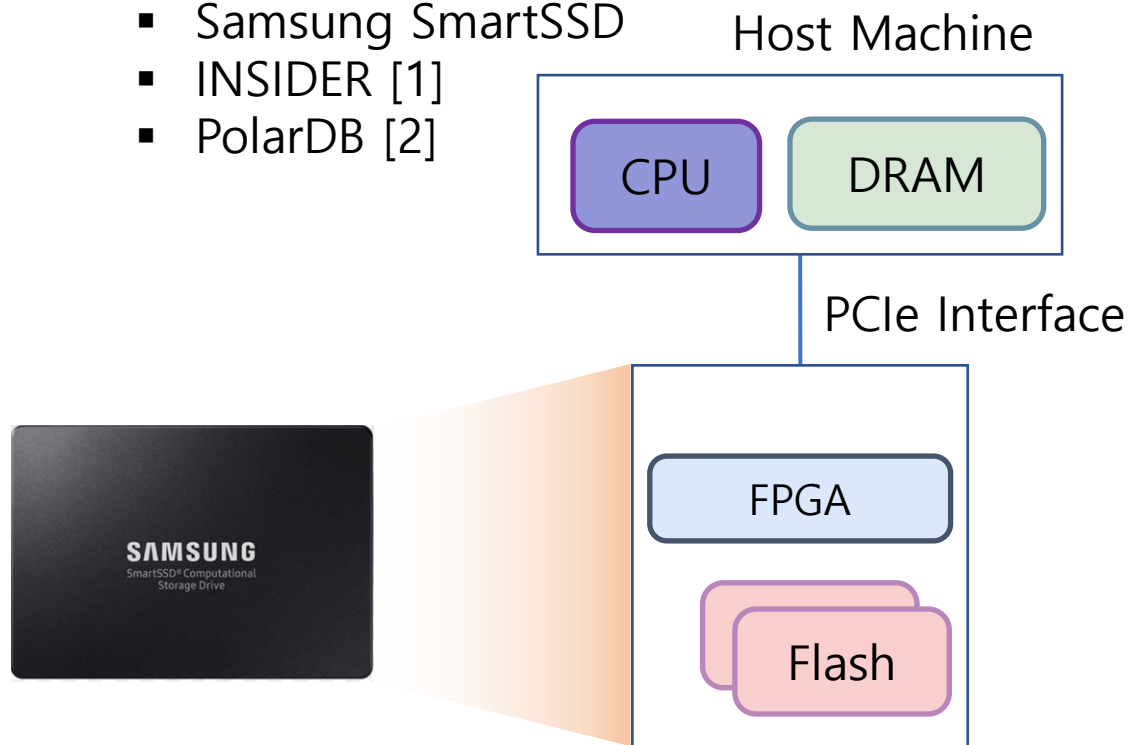
# Computational Storage

❑ What is computational storage device (CSD)?
  ➢ Computational storage devices (CSD) can run computational tasks inside the storage device, reducing data transfer between the host and the device.

❑ CSD **without OS**                    ❑ CSD **with OS**

# Computational Storage

❑ What is computational storage device (CSD)?
  ➤ Computational storage devices (CSD) can run computational tasks inside the storage device, reducing data transfer between the host and the device.

❑ CSD **without OS**

  ▪ Samsung SmartSSD
  ▪ INSIDER [1]
  ▪ PolarDB [2]

❑ CSD **with OS**

Host Machine

CPU    DRAM

PCIe Interface

FPGA

Flash

SAMSUNG
SmartSSD® Computational
Storage Drive

[1] Z. Ruan et. al., "INSIDER: Designing In-Storage Computing System for Emerging High-Performance Drive," USENIX ATC '19
[2] W. Cao et. al., "POLARDB Meets Computational Storage: Efficiently Support Analytical Workloads in Cloud-Native Relational Database," FAST '14
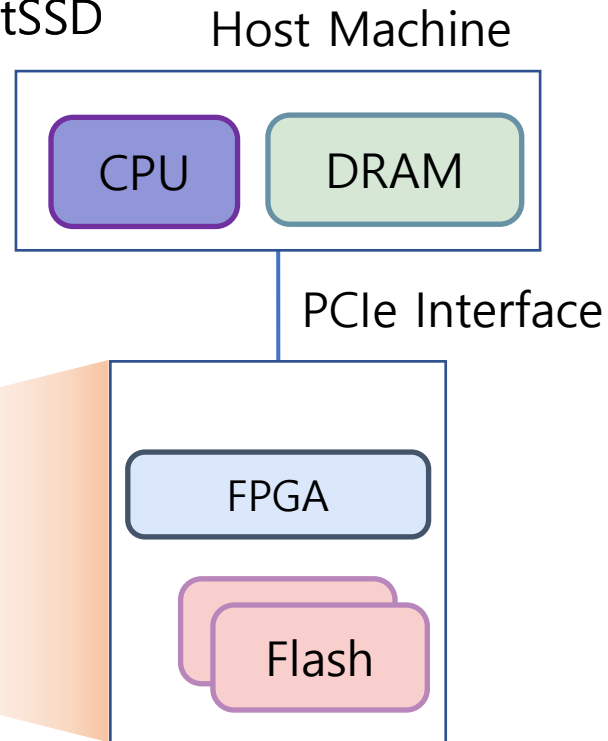
# Computational Storage

❑ What is computational storage device (CSD)?
  ➢ Computational storage devices (CSD) can run computational tasks inside the storage device, reducing data transfer between the host and the device.
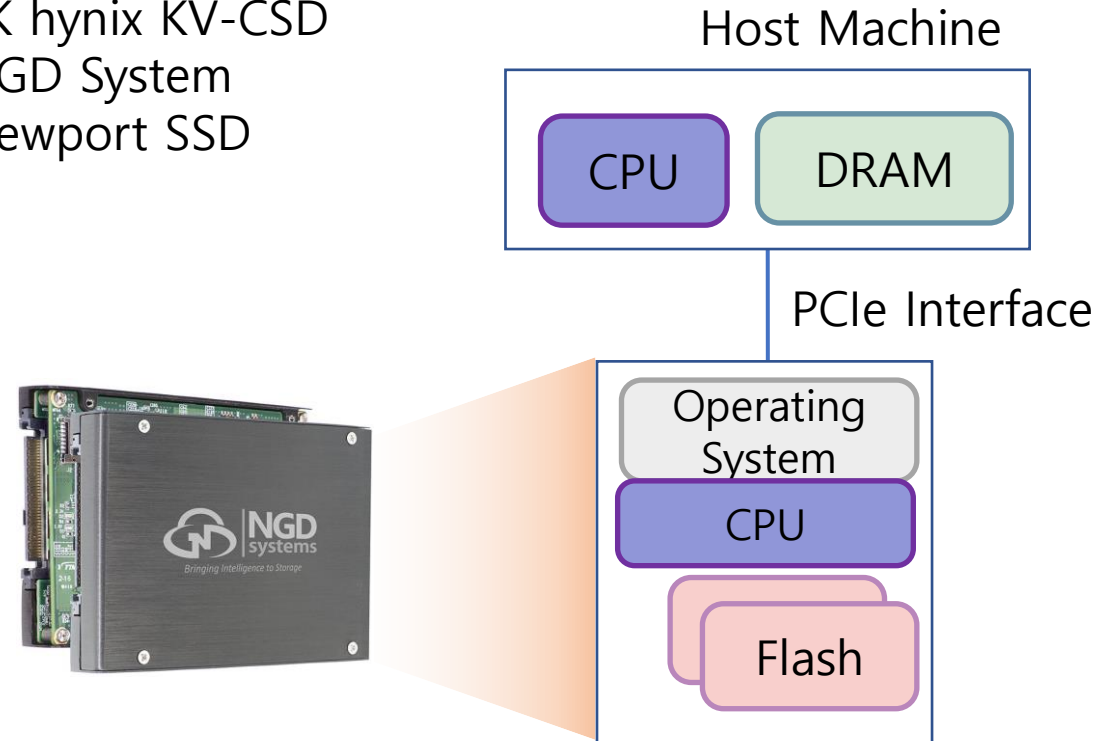
❑ CSD **without OS**

  ▪ Samsung SmartSSD
  ▪ INSIDER [1]
  ▪ PolarDB [2]

  Host Machine

  | CPU | DRAM |

  PCIe Interface

  FPGA

  Flash

❑ CSD **with OS**

  ▪ SK hynix KV-CSD
  ▪ NGD System Newport SSD

  Host Machine

  | CPU | DRAM |

  PCIe Interface

  Operating System

  CPU

  Flash

[1] Z. Ruan et. al., "INSIDER: Designing In-Storage Computing System for Emerging High-Performance Drive," USENIX ATC '19
[2] W. Cao et. al., "POLARDB Meets Computational Storage: Efficiently Support Analytical Workloads in Cloud-Native Relational Database," FAST '14

# CSD **with OS**

❑ Pros
  ➢ Programmability, and manageability

❑ Cons
  ➢ OS overhead due to frequent interrupts and context switches

[1] Intel. SPDK. https://spdk.io/.

# CSD **with OS**

❑ Pros
  ➢ Programmability, and manageability

❑ Cons
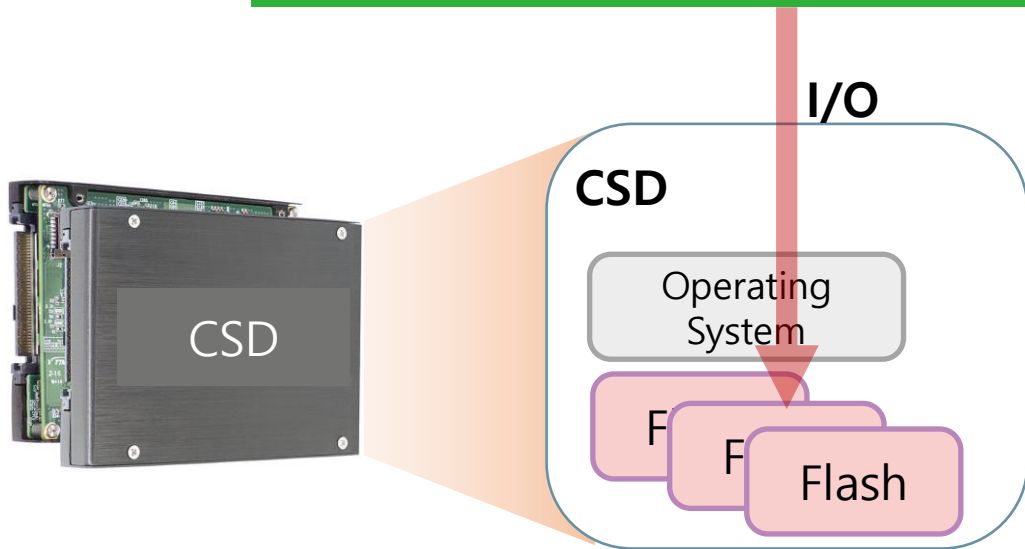  ➢ OS overhead due to frequent interrupts and context switches
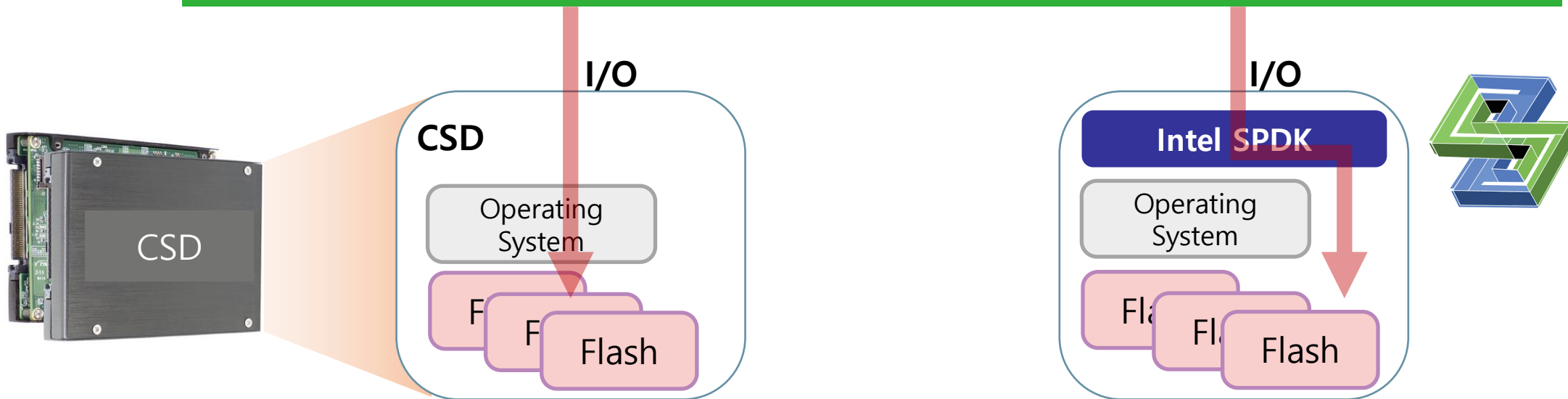
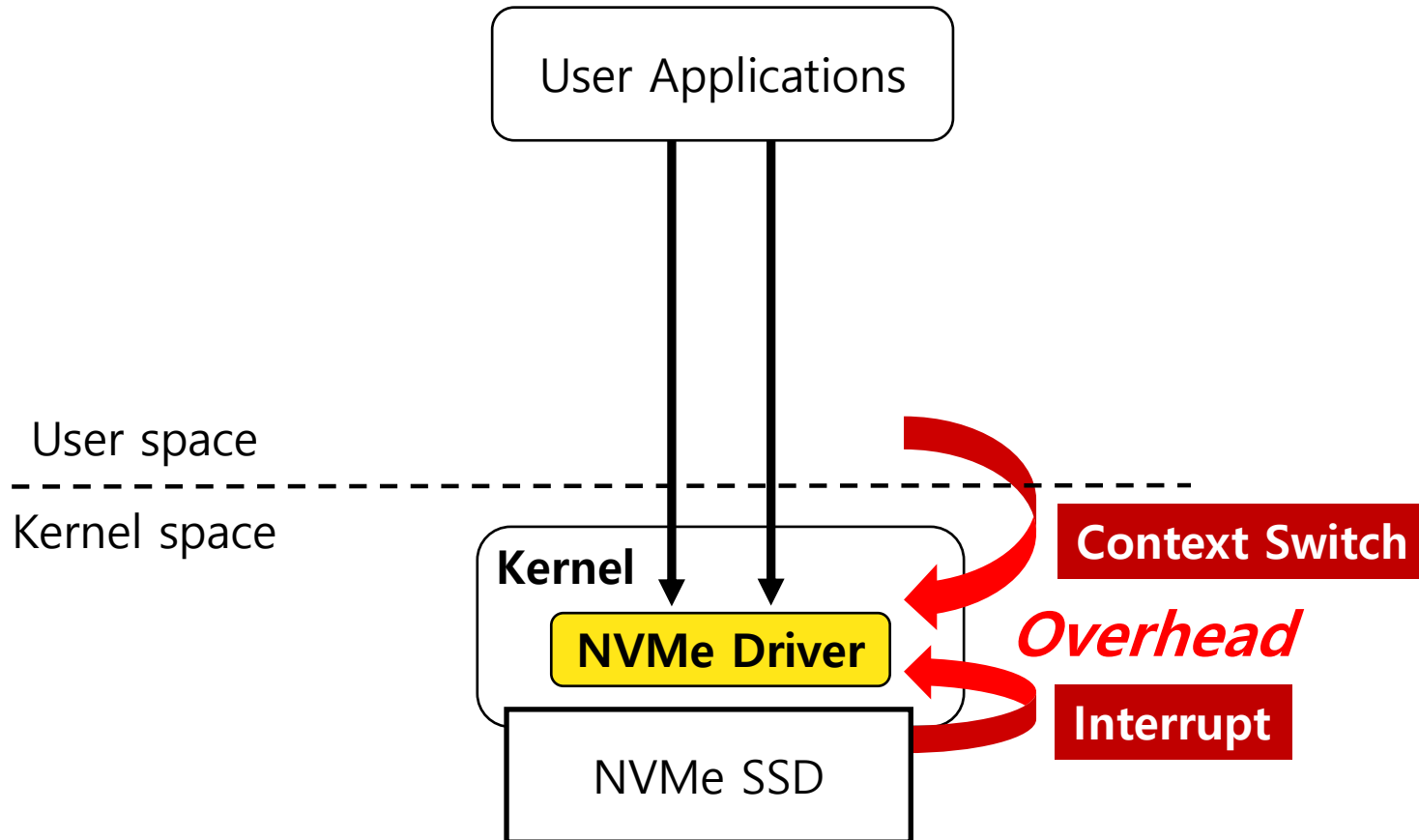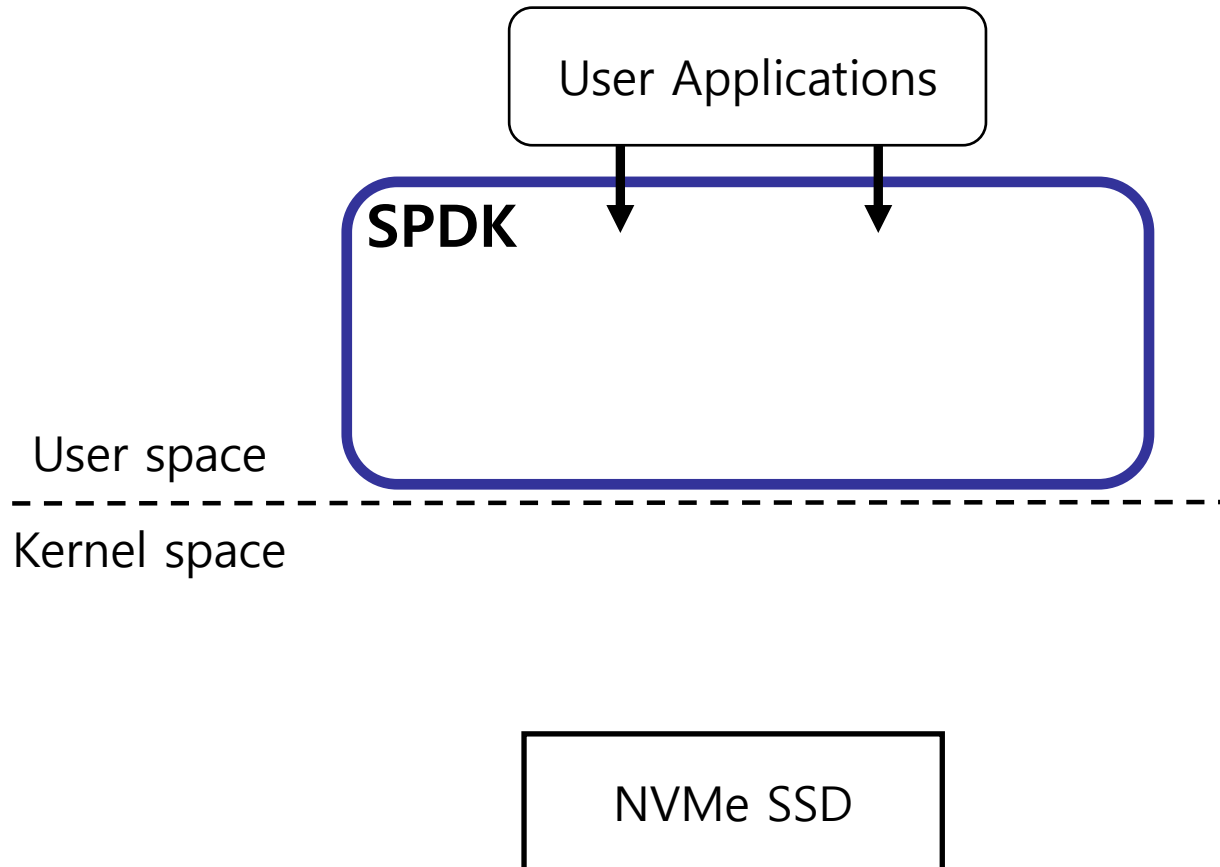> To reduce OS overhead, CSD can adopt Intel SPDK[1].

[1] Intel. SPDK. https://spdk.io/.

# CSD **with OS**

❑ Pros
  ➢ Programmability, and manageability

❑ Cons
  ➢ OS overhead due to frequent interrupts and context switches



To reduce OS overhead, CSD can adopt Intel SPDK[1].

**I/O**

**CSD**

Operating System

Flash
Flash
Flash

[1] Intel. SPDK. https://spdk.io/.

# CSD **with OS**

❑ Pros
  ➢ Programmability, and manageability

❑ Cons
  ➢ OS overhead due to frequent interrupts and context switches

To reduce OS overhead, CSD can adopt Intel SPDK[1].



**CSD**

I/O

Operating System

F

F

Flash

**Intel SPDK**

I/O

Operating System

Fla

Fl

Flash

[1] Intel. SPDK. https://spdk.io/.

# Traditional I/O Stack

❑ Kernel-level NVMe Driver

# Intel SPDK[1]

❏ User-level NVMe Driver

User Applications

**SPDK**

User space

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Kernel space

NVMe SSD

[1] Intel. SPDK. https://spdk.io/.

# Intel SPDK[1]

❑ User-level NVMe Driver

User Applications

**SPDK**

User space

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Kernel space

NVMe SSD

**1. User level NVMe device driver**

[1] Intel. SPDK. https://spdk.io/.

# Intel SPDK[1]

❑ User-level NVMe Driver



1. **User level NVMe device driver**

2. **Binds I/O at the core**

[1] Intel. SPDK. https://spdk.io/.

# Intel SPDK[1]

❑ User-level NVMe Driver



1. **User level NVMe device driver**

2. **Binds I/O at the core**

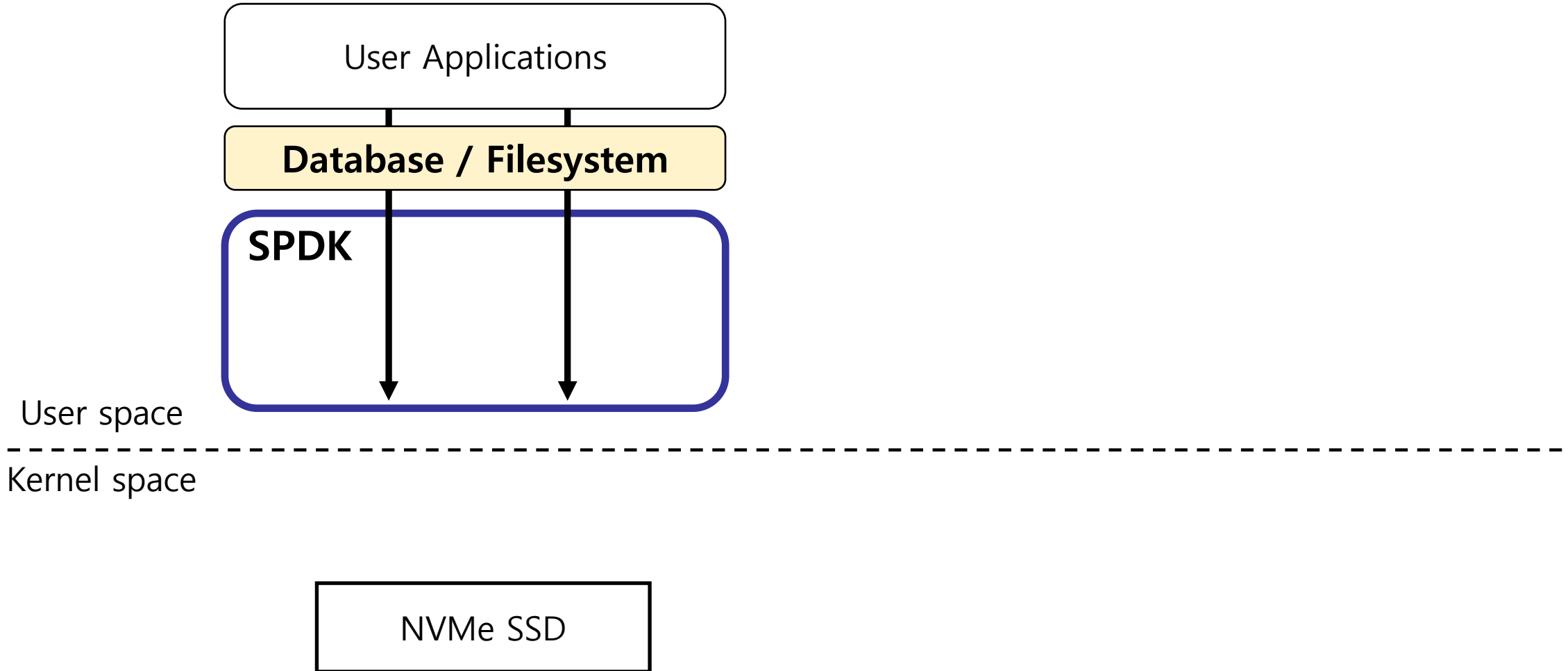3. **Uses polled mode**

[1] Intel. SPDK. https://spdk.io/.

# Use Case

❑ Storage Applications

# Use Case

❑ Storage Applications

```
┌─────────────────────────────────────┐
│         User Applications            │
└─────────────────────────────────────┘
        │              │
┌─────────────────────────────────────┐
│      Database / Filesystem          │
└─────────────────────────────────────┘
        │              │
┌─────────────────────────────────────┐
│  SPDK                                │
│        │              │              │
│        ▼              ▼              │
└─────────────────────────────────────┘
```

User space

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Kernel space

```
┌─────────────────────────────────────┐
│              NVMe SSD                │
└─────────────────────────────────────┘
```

# Use Case

☐ Storage Applications

☐ Storage Services (Deduplication)

# Research Problem

However, SPDK has a problem in that foreground I/O and background service tasks compete for CPU cores.
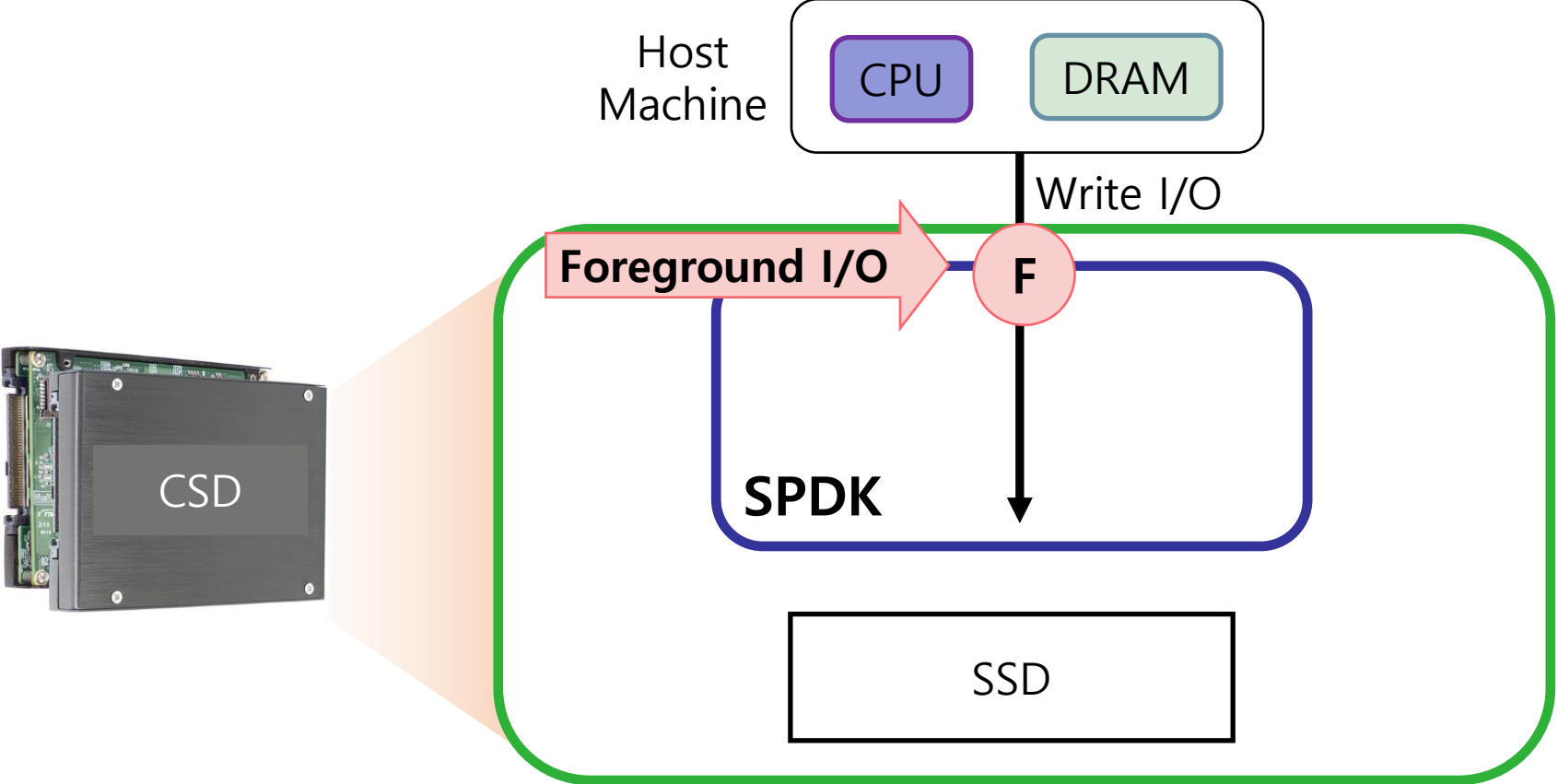
# Research Problem

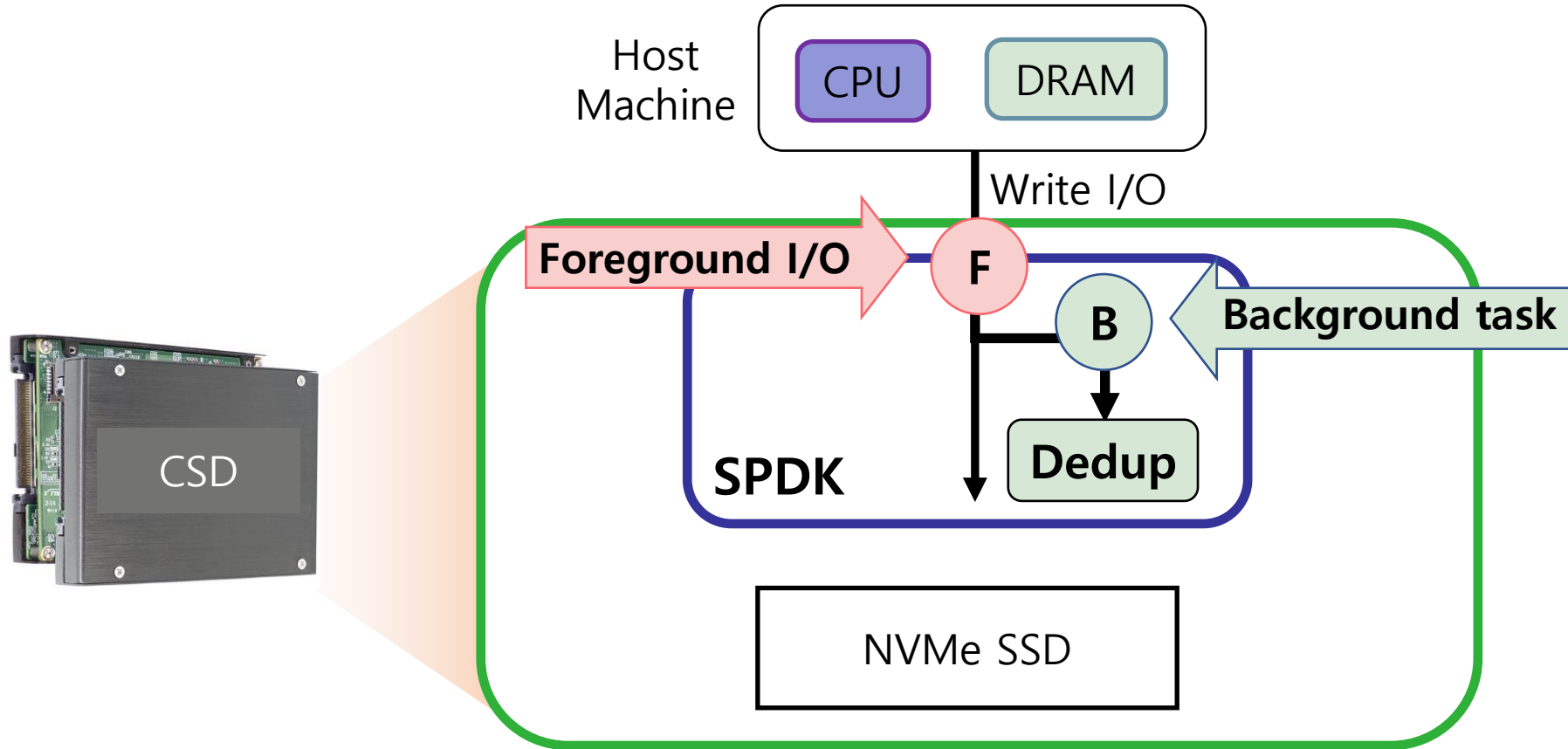However, SPDK has a problem in that foreground I/O and background service tasks compete for CPU cores.

Contention of these tasks for CPU cores increases the response time of foreground I/O.
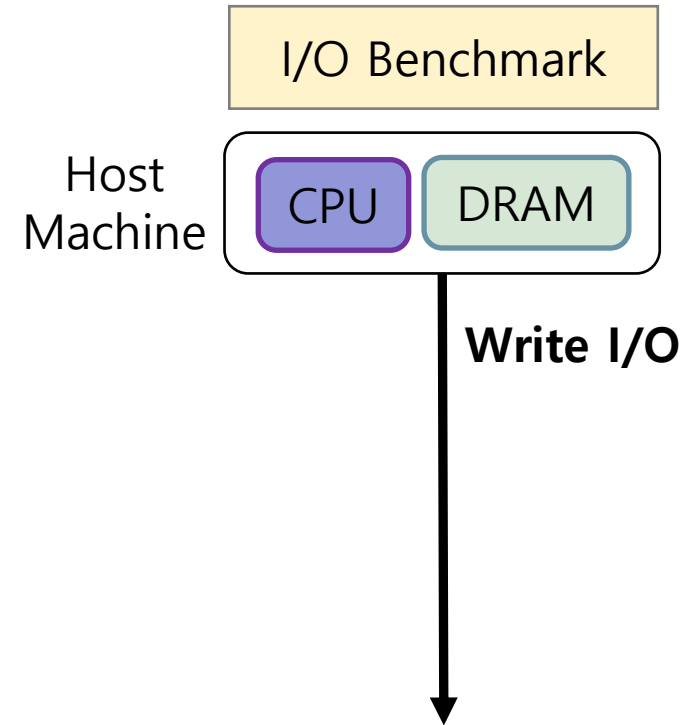
# Motivation

# Executing Background Tasks in SPDK

# Executing Background Tasks in SPDK



In SPDK, background tasks are derived from foreground I/O.
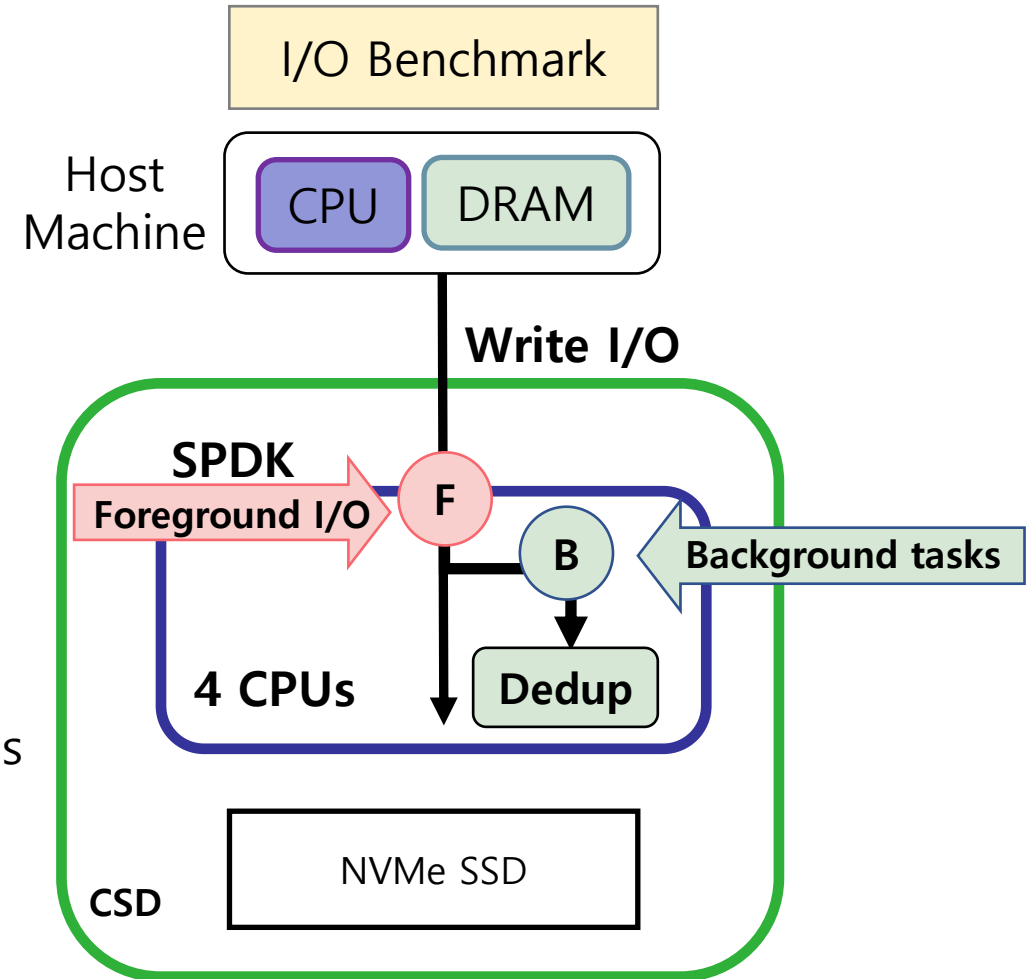
# Experimental Setup

- ❏ Host machine
  - ❏ Running a db_bench
    (Two I/O threads issue write I/Os)
  - ❏ I/O request size = 16KB

- ❏ CSD
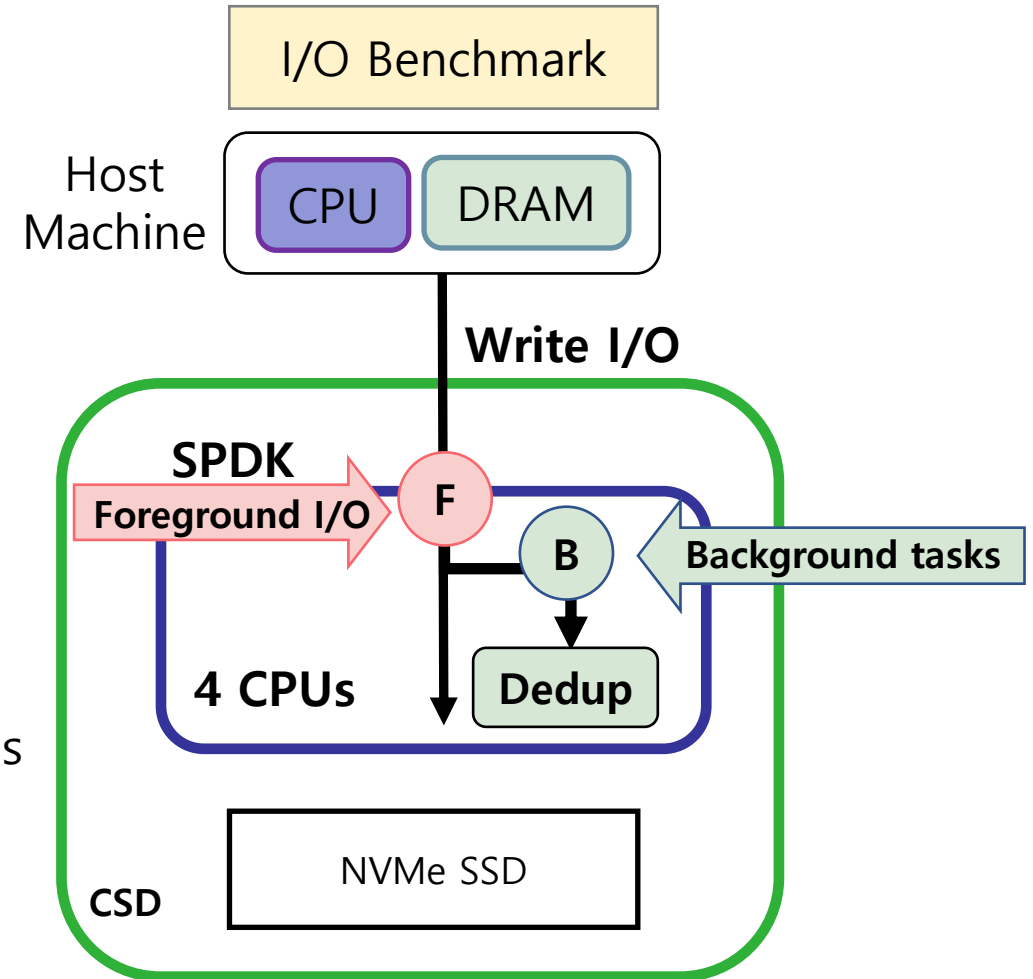  - ❏ 4 Core device
  - ❏ Running a Linux OS using Intel SPDK

- ❏ Background task
  - ❏ Offline deduplication
  - ❏ Fingerprinting using a SHA-1 hash algorithm
  - ❏ Light deduplication : 1KB chunk size, SHA1 16 times
  - ❏ Heavy dedulication : 0.5KB chunk size, SHA1 32 times

- ❏ Comparisons
  - ❏ Only foreground I/O
  - ❏ Foreground I/O + Background task (light)
  - ❏ Foreground I/O + Background task (heavy)

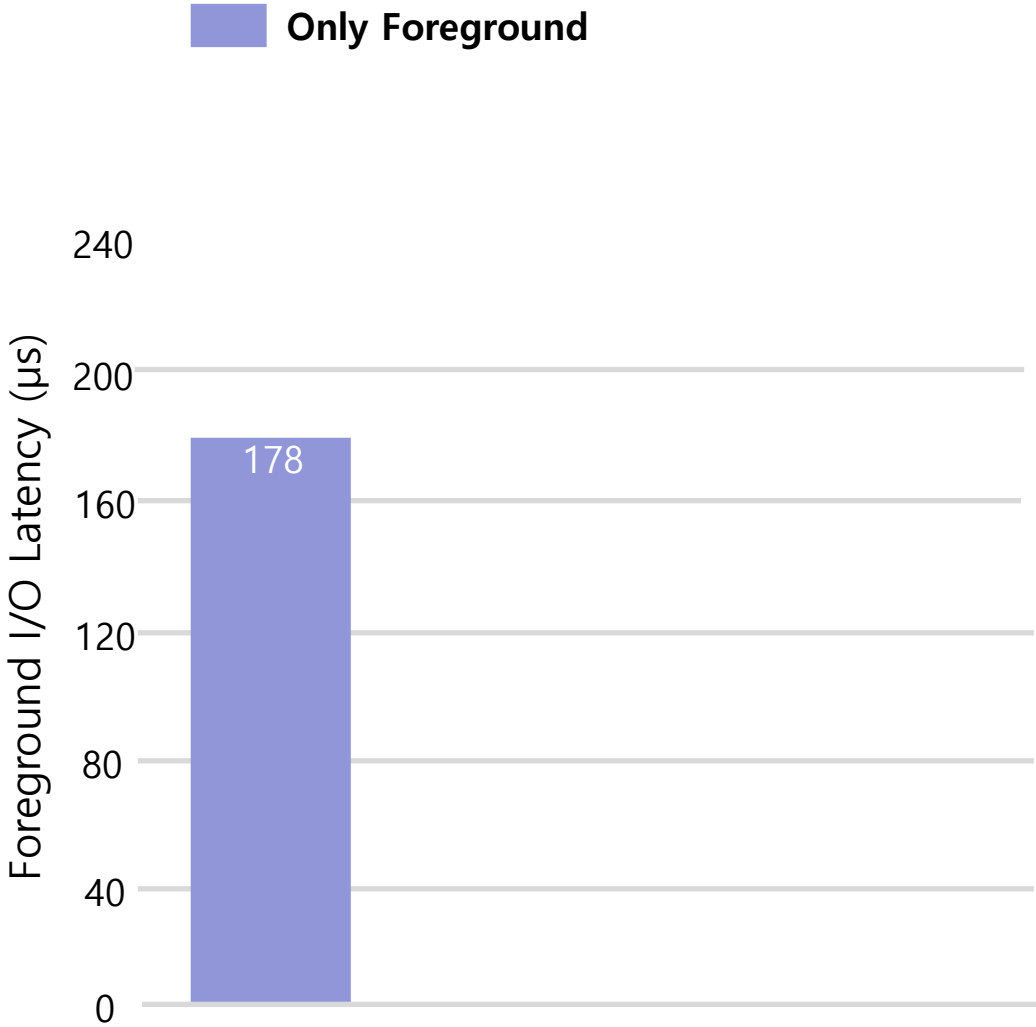I/O Benchmark

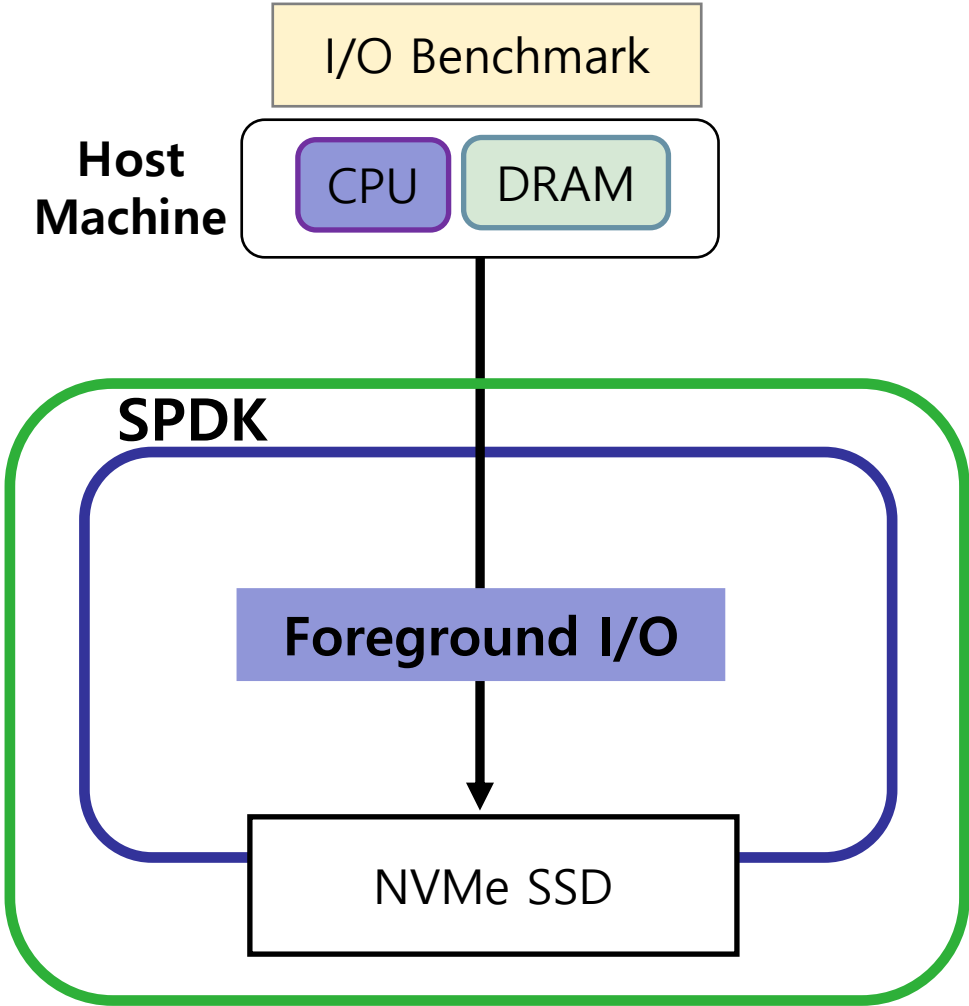Host Machine | CPU | DRAM

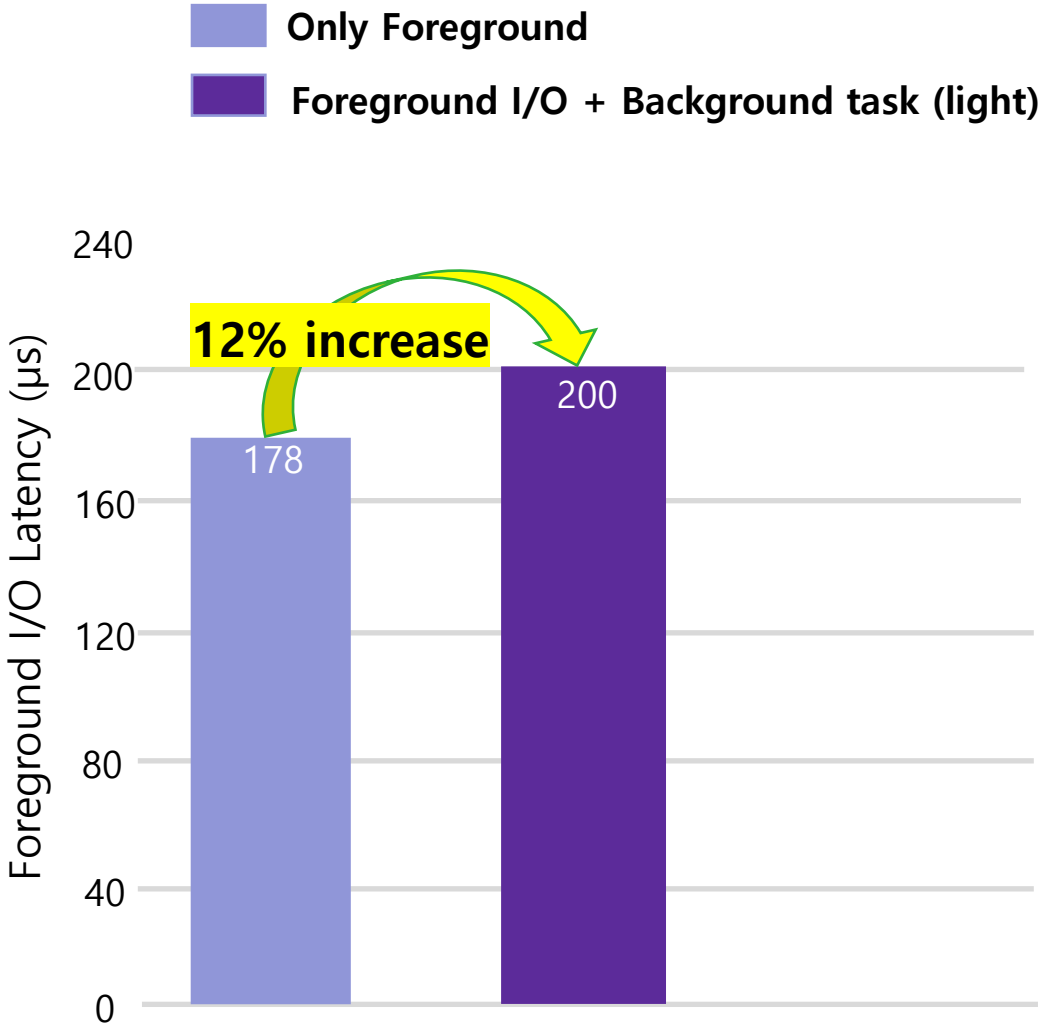**Write I/O**
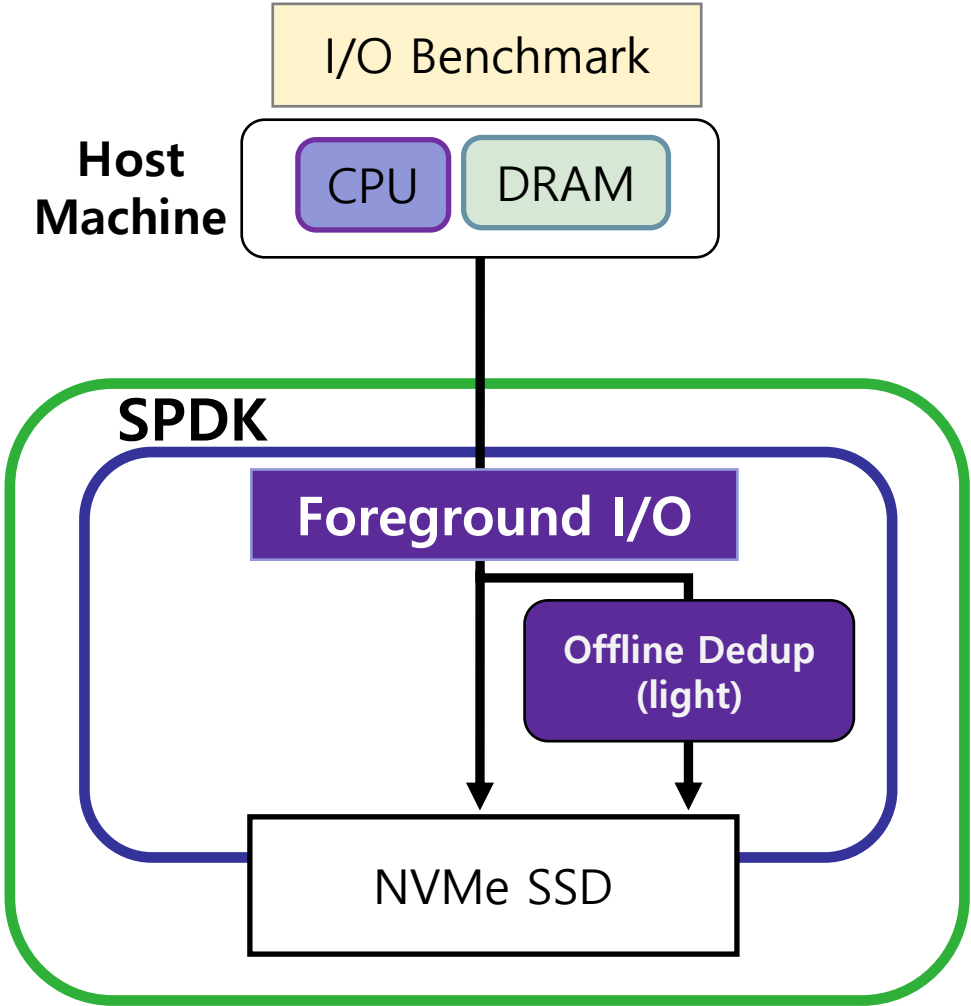
# Experimental Setup

- ❑ Host machine
  - ❑ Running a db_bench
    (Two I/O threads issue write I/Os)
  - ❑ I/O request size = 16KB

- ❑ CSD
  - ❑ 4 Core device
  - ❑ Running a Linux OS using Intel SPDK

I/O Benchmark

Host Machine

CPU    DRAM

**Write I/O**

**SPDK**

F

**4 CPUs**

NVMe SSD

CSD

# Experimental Setup

- Host machine
  - Running a db_bench
    (Two I/O threads issue write I/Os)
  - I/O request size = 16KB

- CSD
  - 4 Core device
  - Running a Linux OS using Intel SPDK

- Background task
  - Offline deduplication
  - Fingerprinting using a SHA-1 hash algorithm
  - Light deduplication : 1KB chunk size, SHA1 16 times
  - Heavy dedulication : 0.5KB chunk size, SHA1 32 times

I/O Benchmark

Host Machine

CPU    DRAM

**Write I/O**

**SPDK**
**Foreground I/O**    F

B    **Background tasks**

**4 CPUs**    **Dedup**

**CSD**

NVMe SSD

# Experimental Setup

- ❑ Host machine
  - ❑ Running a db_bench
    (Two I/O threads issue write I/Os)
  - ❑ I/O request size = 16KB

- ❑ CSD
  - ❑ 4 Core device
  - ❑ Running a Linux OS using Intel SPDK

- ❑ Background task
  - ❑ Offline deduplication
  - ❑ Fingerprinting using a SHA-1 hash algorithm
  - ❑ Light deduplication : 1KB chunk size, SHA1 16 times
  - ❑ Heavy dedulication : 0.5KB chunk size, SHA1 32 times

- ❑ Comparisons
  - ❑ Only foreground I/O
  - ❑ Foreground I/O + Background task (light)
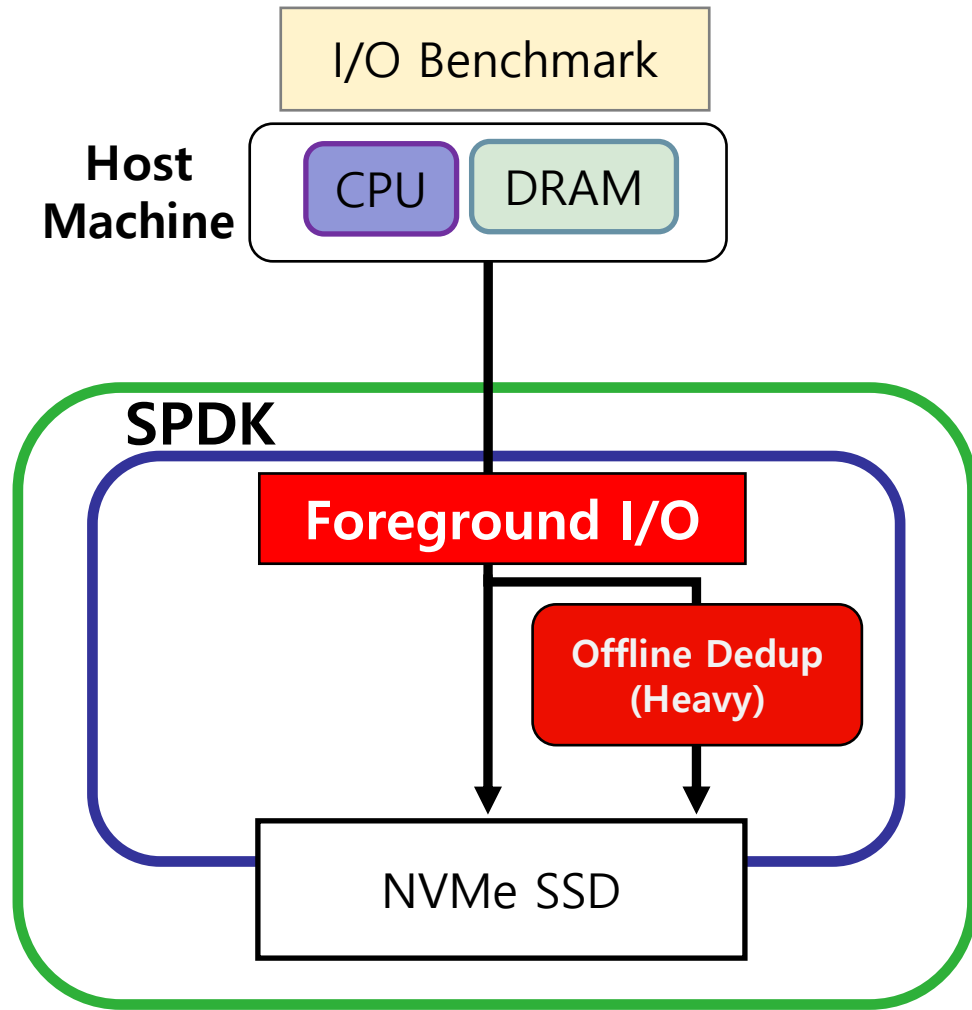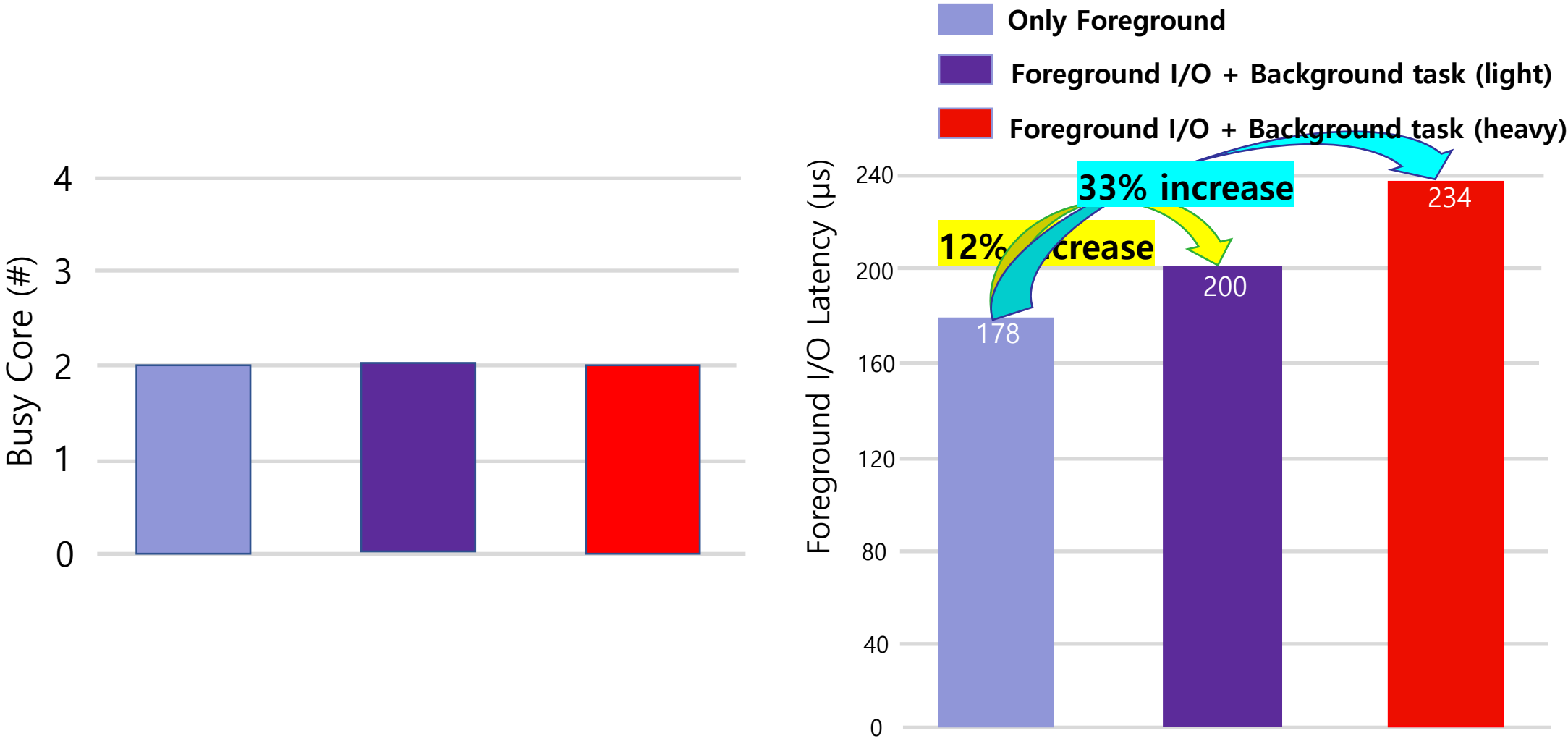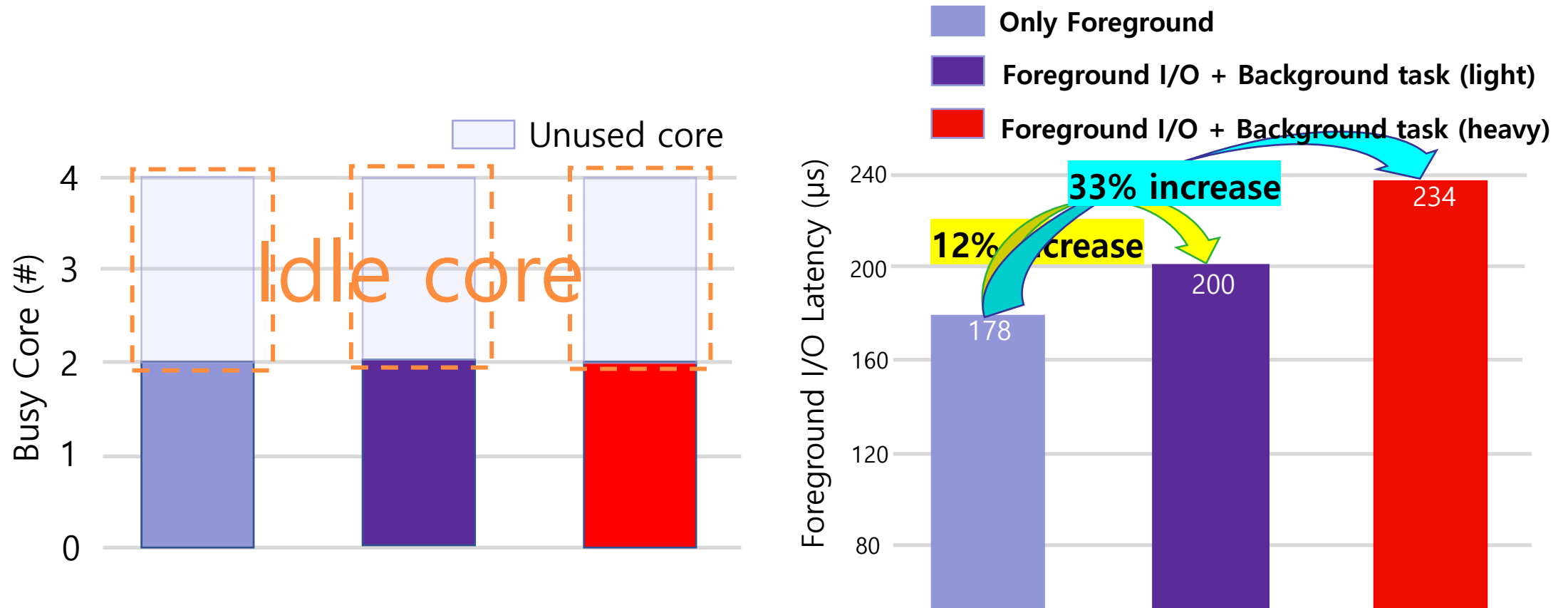  - ❑ Foreground I/O + Background task (heavy)

# Motivation



Host Machine

I/O Benchmark

CPU   DRAM

SPDK

Foreground I/O

NVMe SSD

Only Foreground

Foreground I/O Latency (µs)

178

# Motivation

# Motivation



**I/O Benchmark**

**Host Machine**
CPU | DRAM

**SPDK**

**Foreground I/O**

**Offline Dedup (Heavy)**

NVMe SSD

Legend:
- Only Foreground
- Foreground I/O + Background task (light)
- Foreground I/O + Background task (heavy)

Foreground I/O Latency (μs)

178 — 200 — 234

**12% increase**

**33% increase**

# Motivation



Busy Core (#)

Foreground I/O Latency (μs)

**Only Foreground**

**Foreground I/O + Background task (light)**

**Foreground I/O + Background task (heavy)**

**33% increase**

**12% increase**

178

200

234

# Motivation



**Only two cores are fully-utilized.**
**That is, the remaining cores are under-utilized.**

# Problem Definition

① **SPDK is not aware of background tasks.**

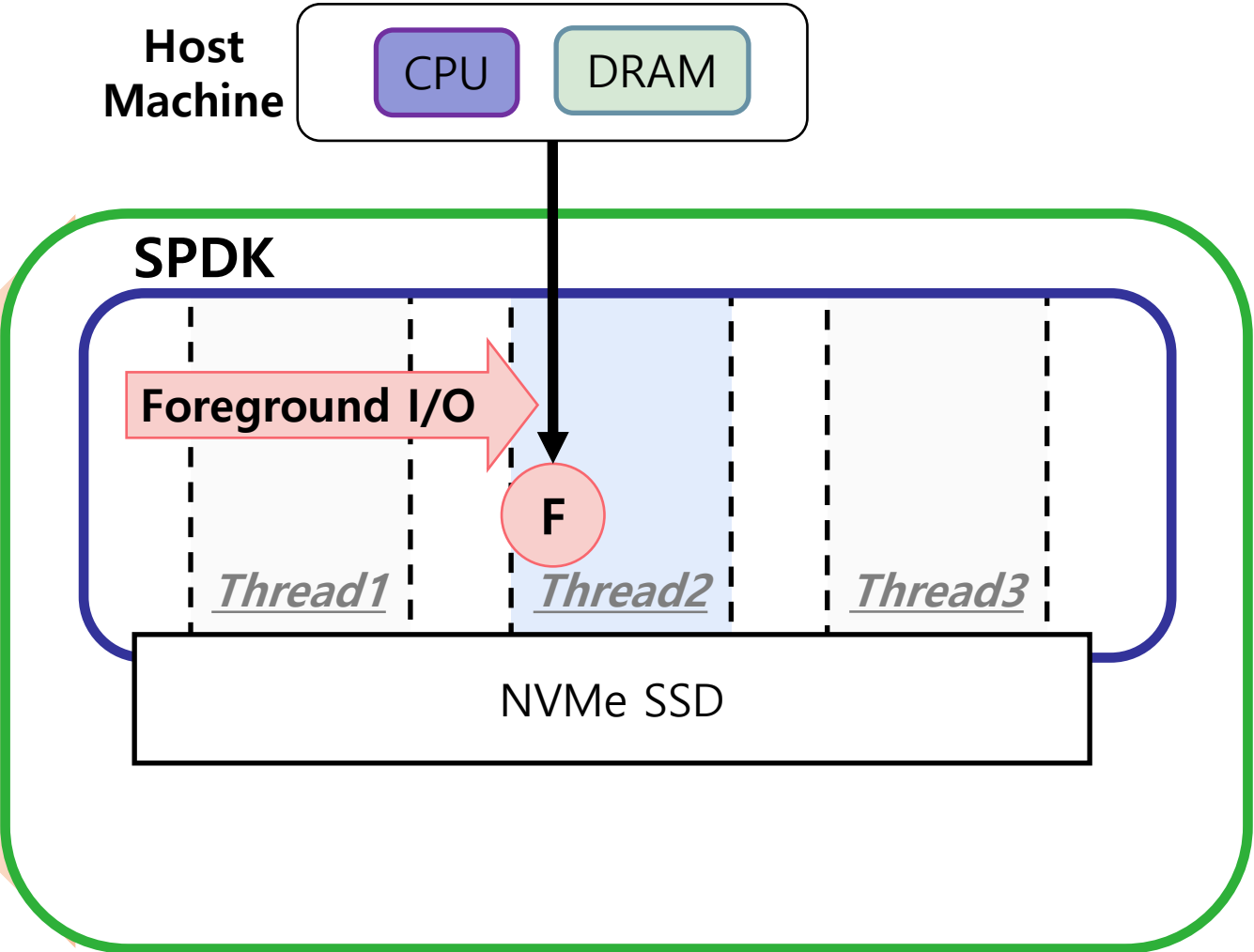② **SPDK cannot perform dynamic task scheduling considering the load of each CPU core.**

→ **BTS : Background Task-aware Scheduler**
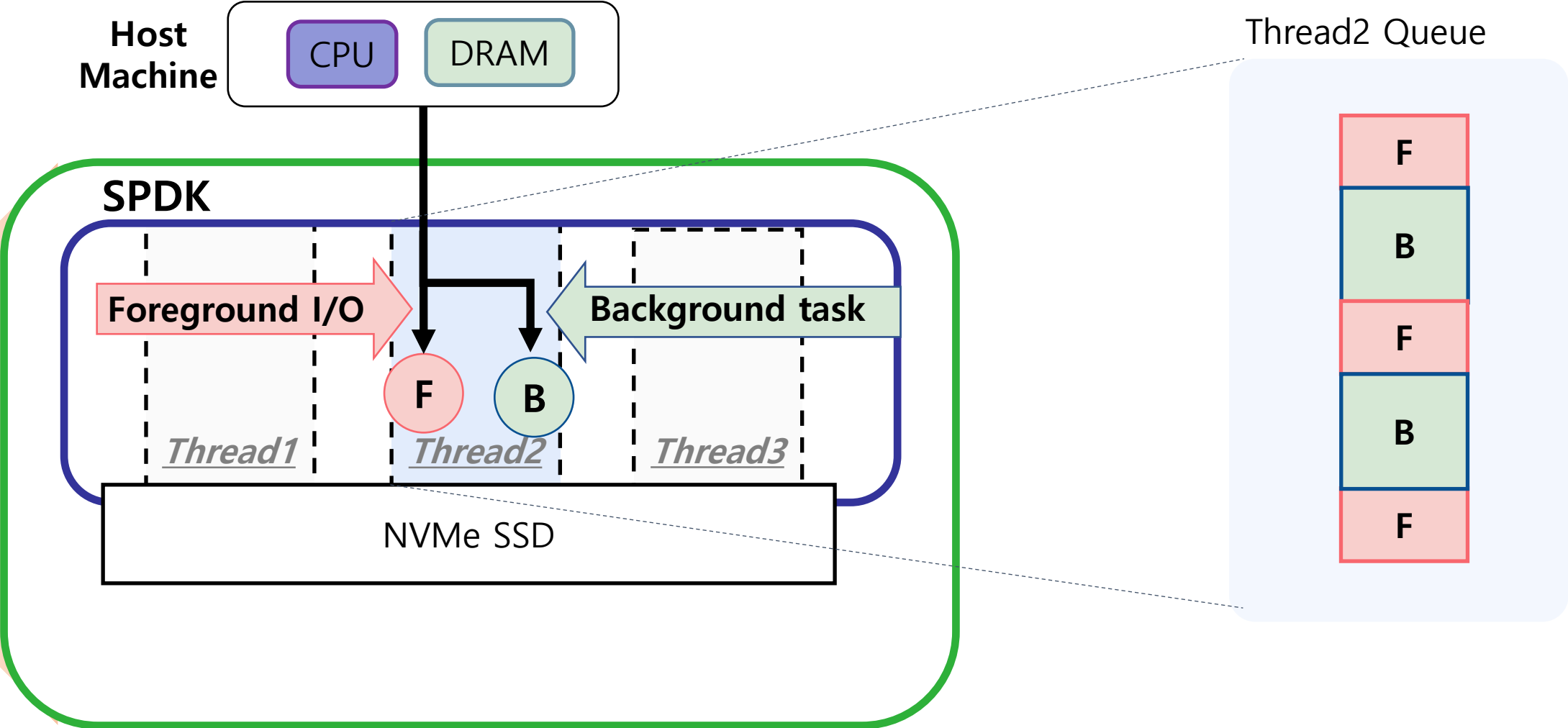
# BTS : Background Task-Aware Scheduler

# SPDK Problem

# SPDK Problem

# SPDK Problem

# SPDK Problem

# SPDK Problem



**Host Machine**

CPU  DRAM

**SPDK**

Foreground I/O    Background task

F    B

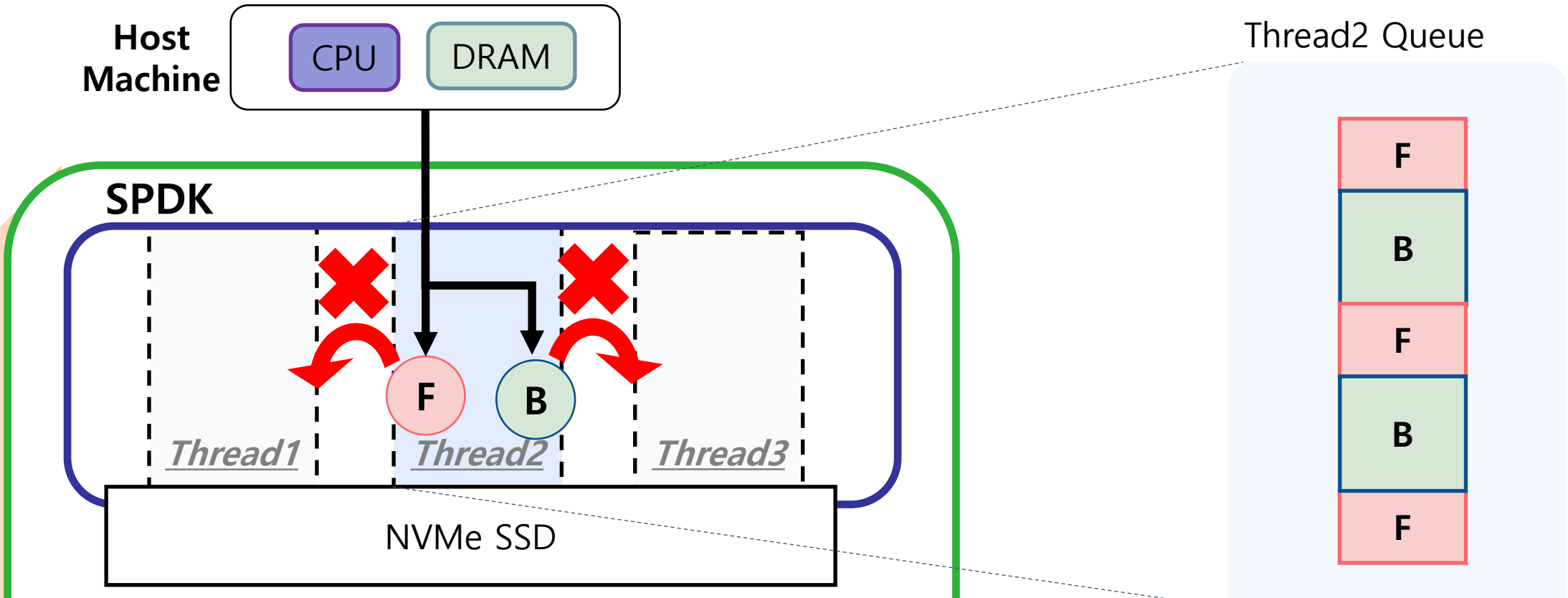*Thread1*    *Thread2*    *Thread3*

NVMe SSD

Thread2 Queue

F
B
F
B
F

The SPDK places a background task on the core where the foreground I/O is being processed.
Therefore, the two tasks compete for the same core.
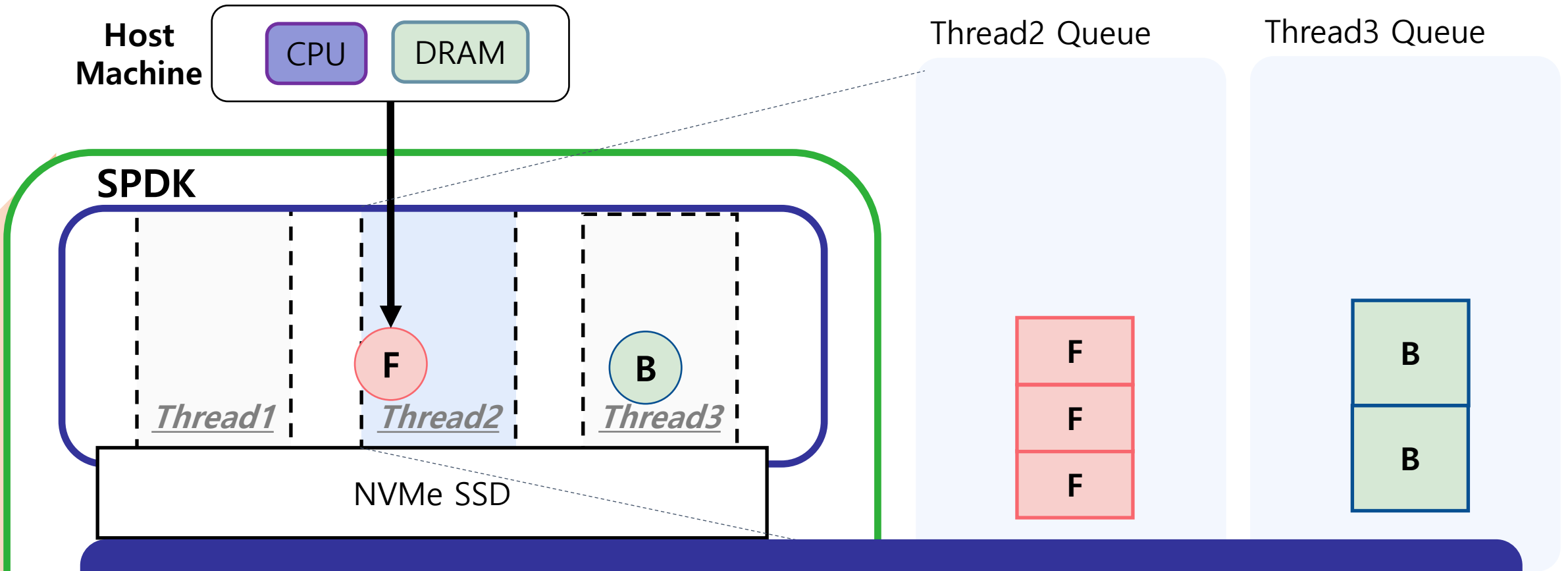
# SPDK Problem



Host Machine

CPU   DRAM

SPDK

Thread1   Thread2   Thread3

F   B

NVMe SSD

Thread2 Queue

F
B
F
B
F

SPDK does not schedule flexible relocation of tasks to idle CPU core. Therefore, idle CPU cores cannot be utilized.

# SPDK Problem

# SPDK Problem



BTS dynamically schedules background tasks to relocate to idle CPU cores. This allows the SPDK to actively utilize idle CPU cores.

# BTS Scheduler

**Background Task-aware Scheduler**

# BTS Scheduler

**Background Task-aware Scheduler**
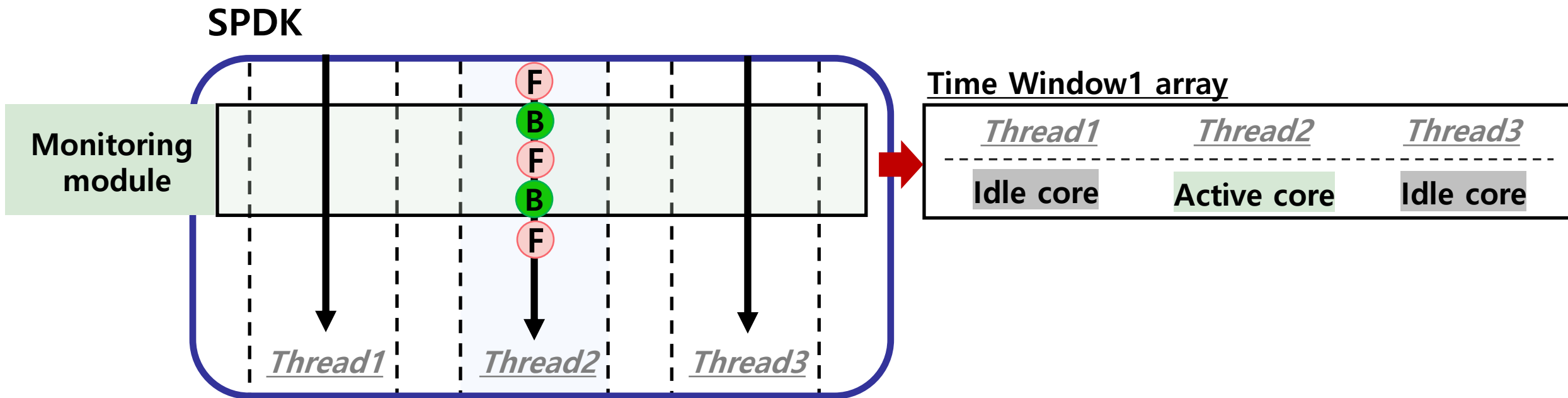
**(1) Monitoring module**

# BTS Scheduler

**Background Task-aware Scheduler**

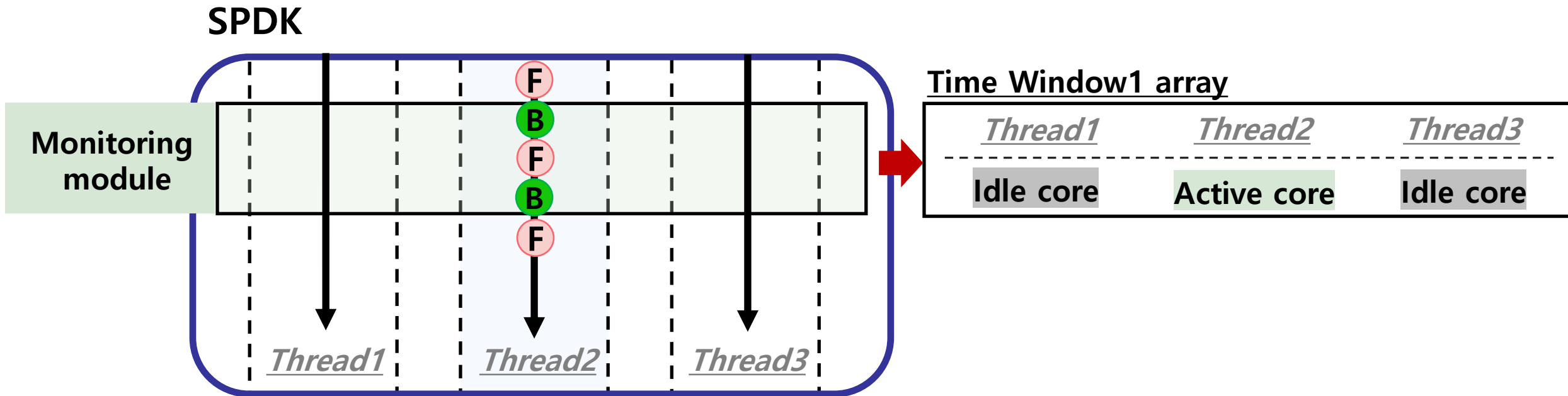**(1) Monitoring module**

**(2) Core selection module**

# 1) Monitoring Module

❑ Monitors the utilization of each core

    ❑ Because SPDK randomly changes the core that handles foreground I/O
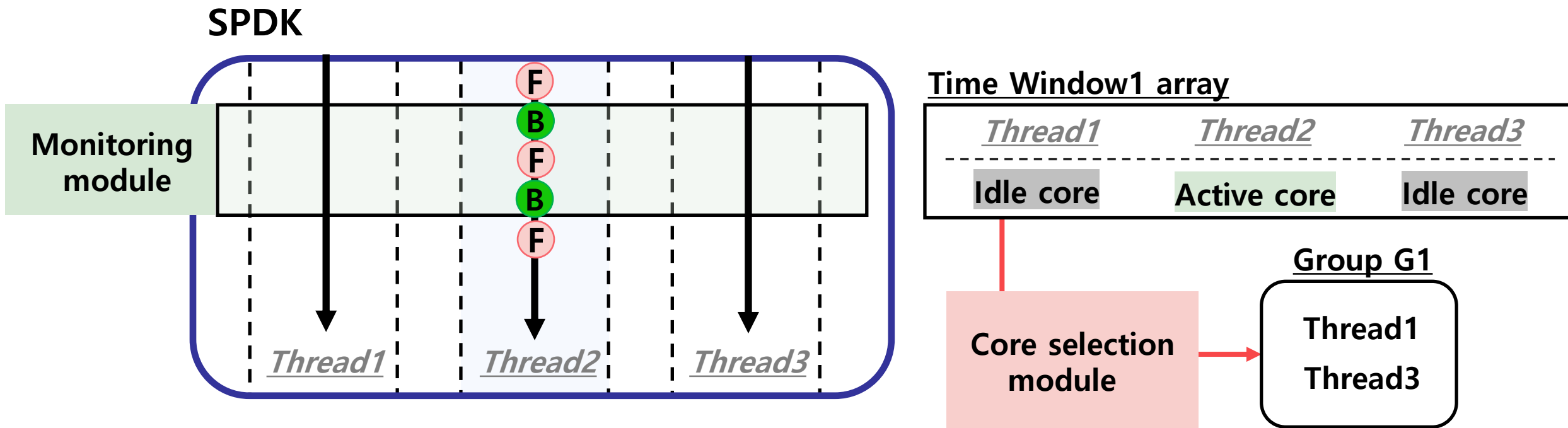
# 1) Monitoring Module

- ❑ Monitors the utilization of each core
    - ❑ Because SPDK randomly changes the core that handles foreground I/O
- ❑ The monitoring module periodically tracks the utilization of active cores
    - ❑ Active core is CPU core that processes at least one foreground I/O
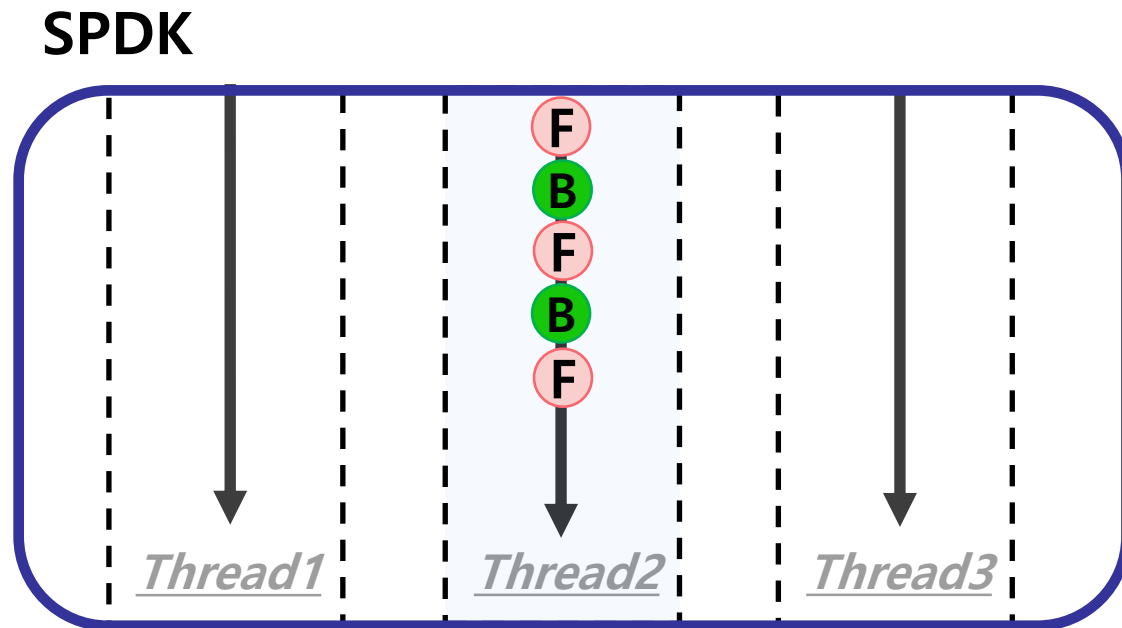
# 2) Core Selection Module

❑ Selects a core with low utilization and move background task to that core

❑ Builds an idle core group (G) based on the CPU utilization of each core

# 2) Core Selection Module

❑ Selects a core with low utilization and move background task to that core

❑ Builds an idle core group (G) based on the CPU utilization of each core

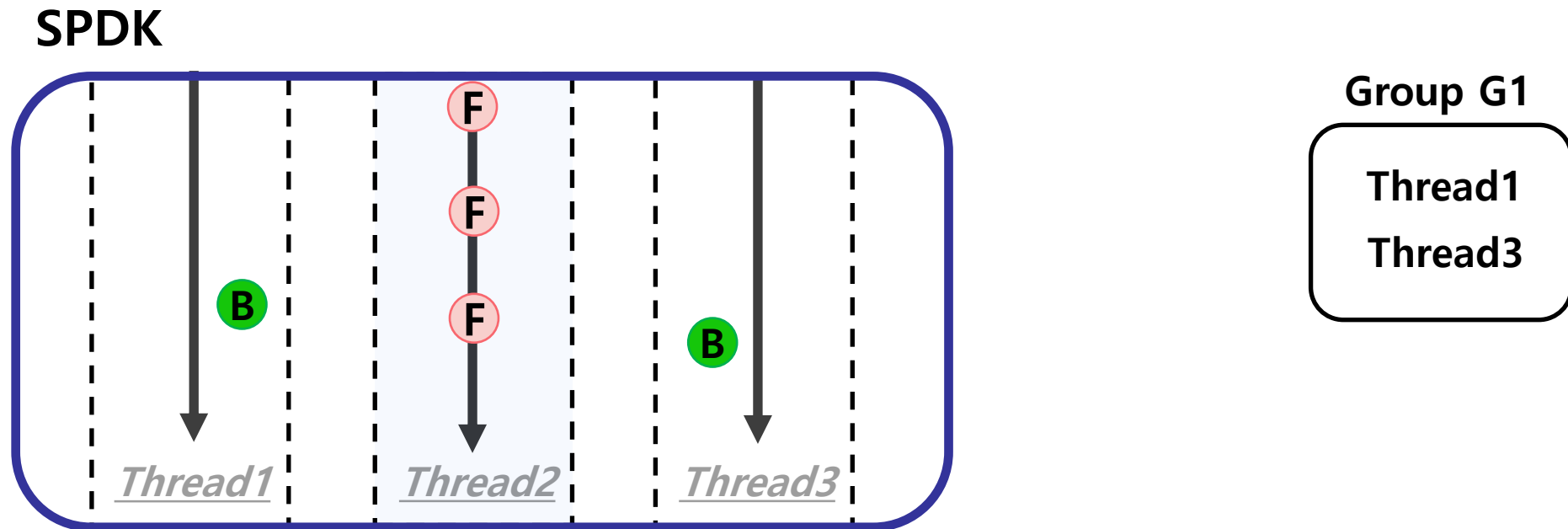❑ Selects idle cores to execute background tasks from the idle core group (G)
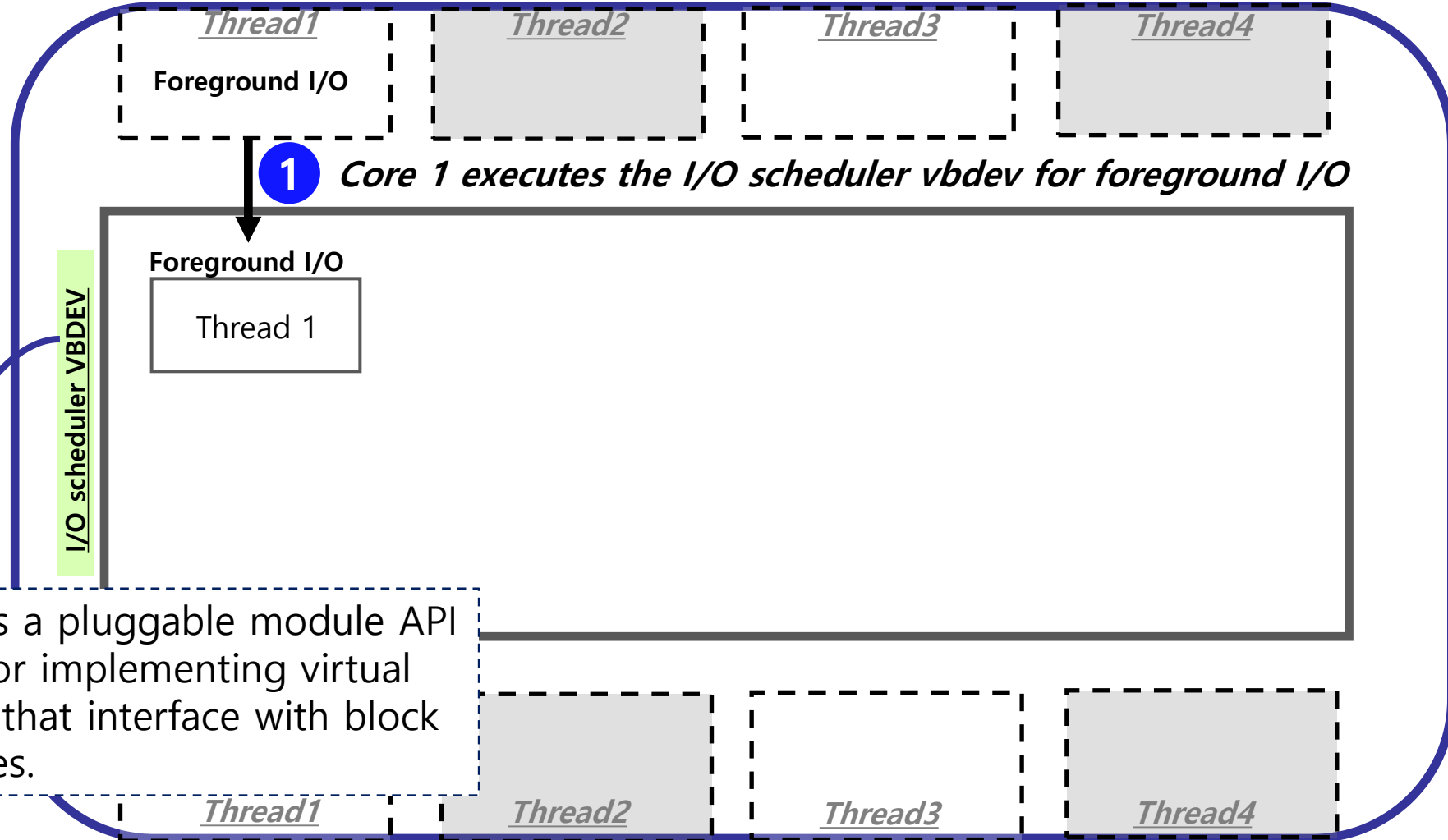
**SPDK**

# 2) Core Selection Module

❑ Selects a core with low utilization and move background task to that core

❑ Builds an idle core group (G) based on the CPU utilization of each core

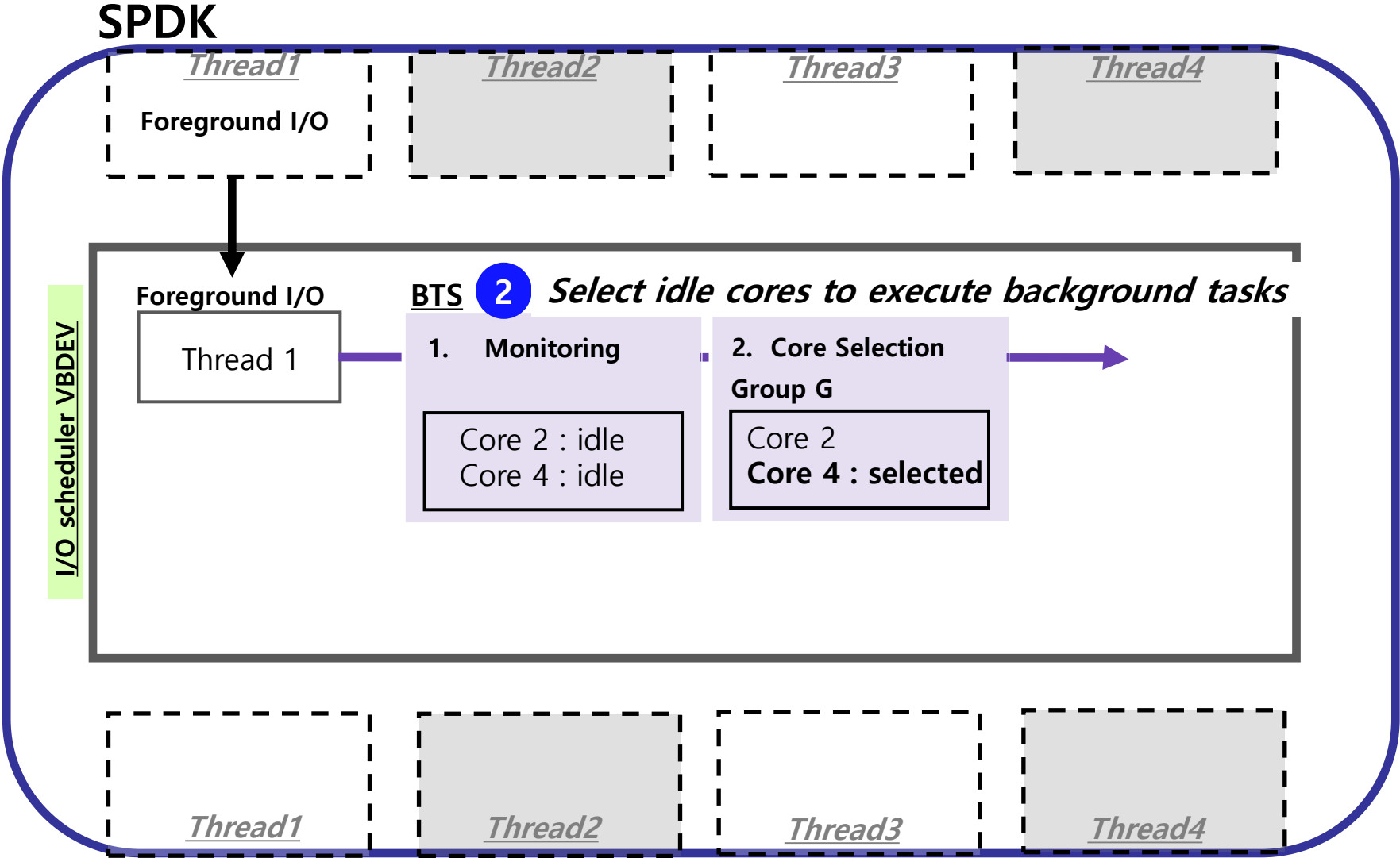❑ Selects idle cores to execute background tasks from the idle core group (G)

**SPDK**



**Group G1**

Thread1
Thread3

# Implementation

# Implementation

**SPDK**

*Thread1*

Foreground I/O

*Thread2*

*Thread3*

*Thread4*

**1** *Core 1 executes the I/O scheduler vbdev for foreground I/O*
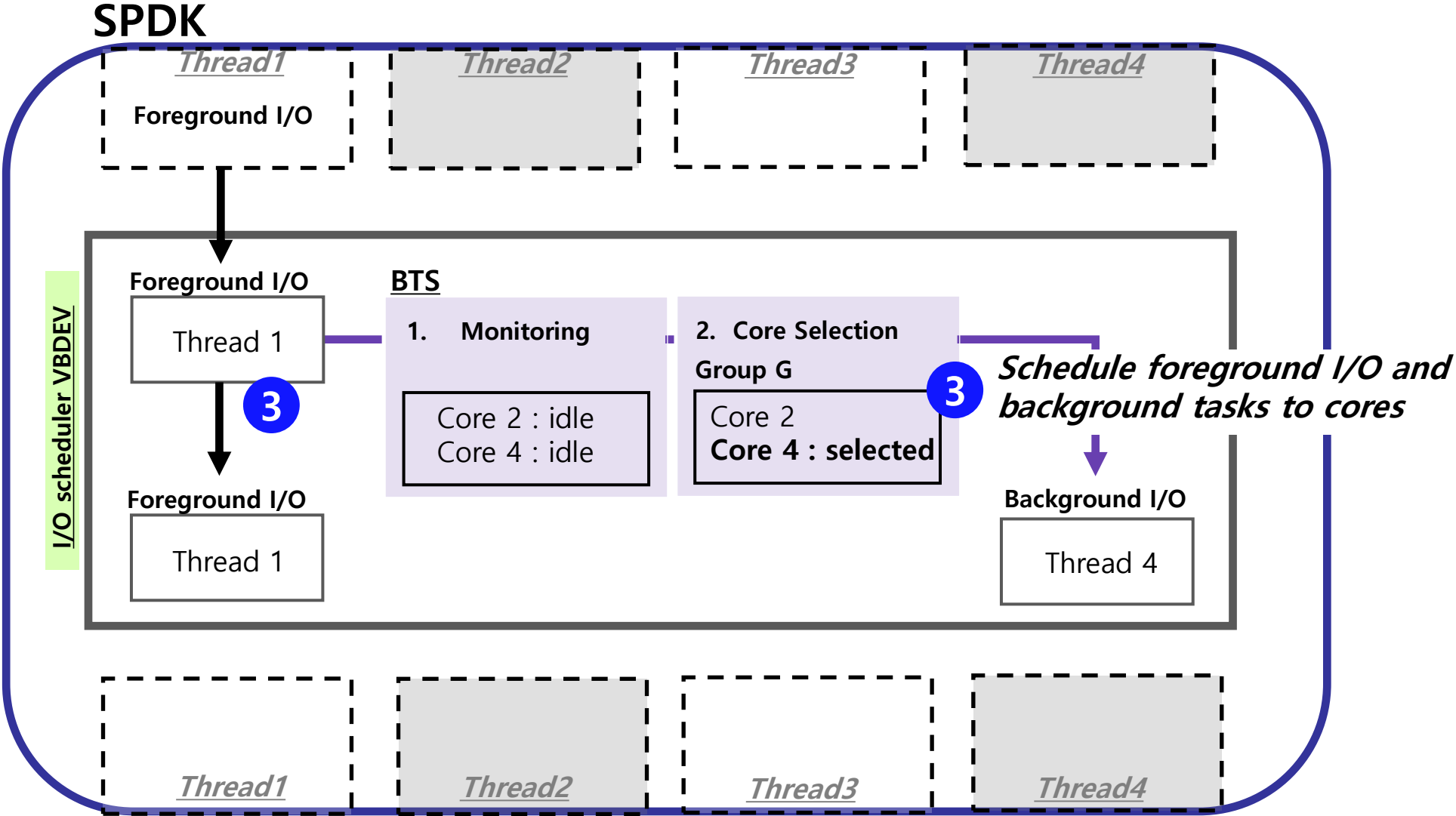
**Foreground I/O**

Thread 1

I/O scheduler VBDEV

SPDK provides a pluggable module API called BDEV for implementing virtual block devices that interface with block storage devices.
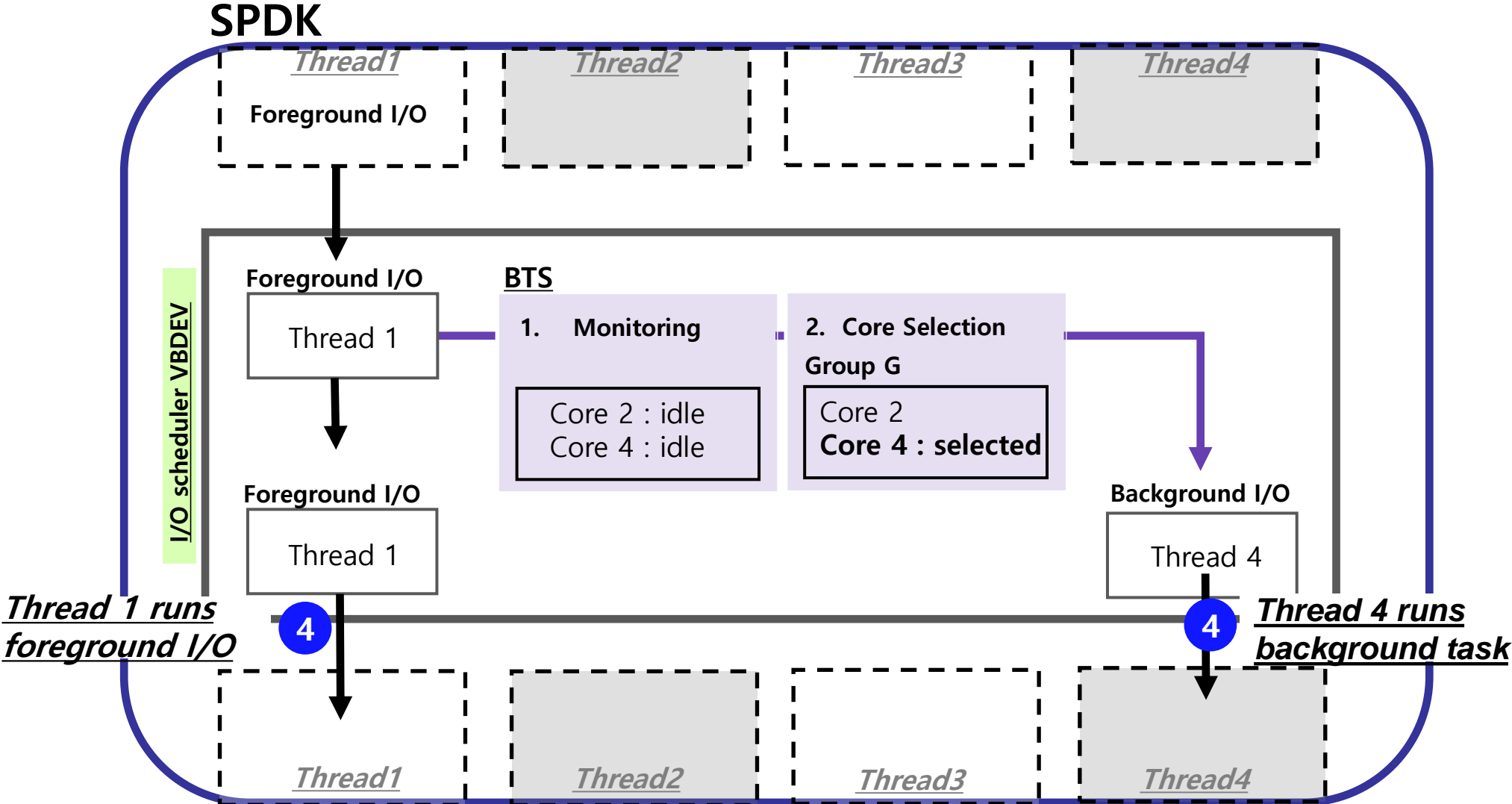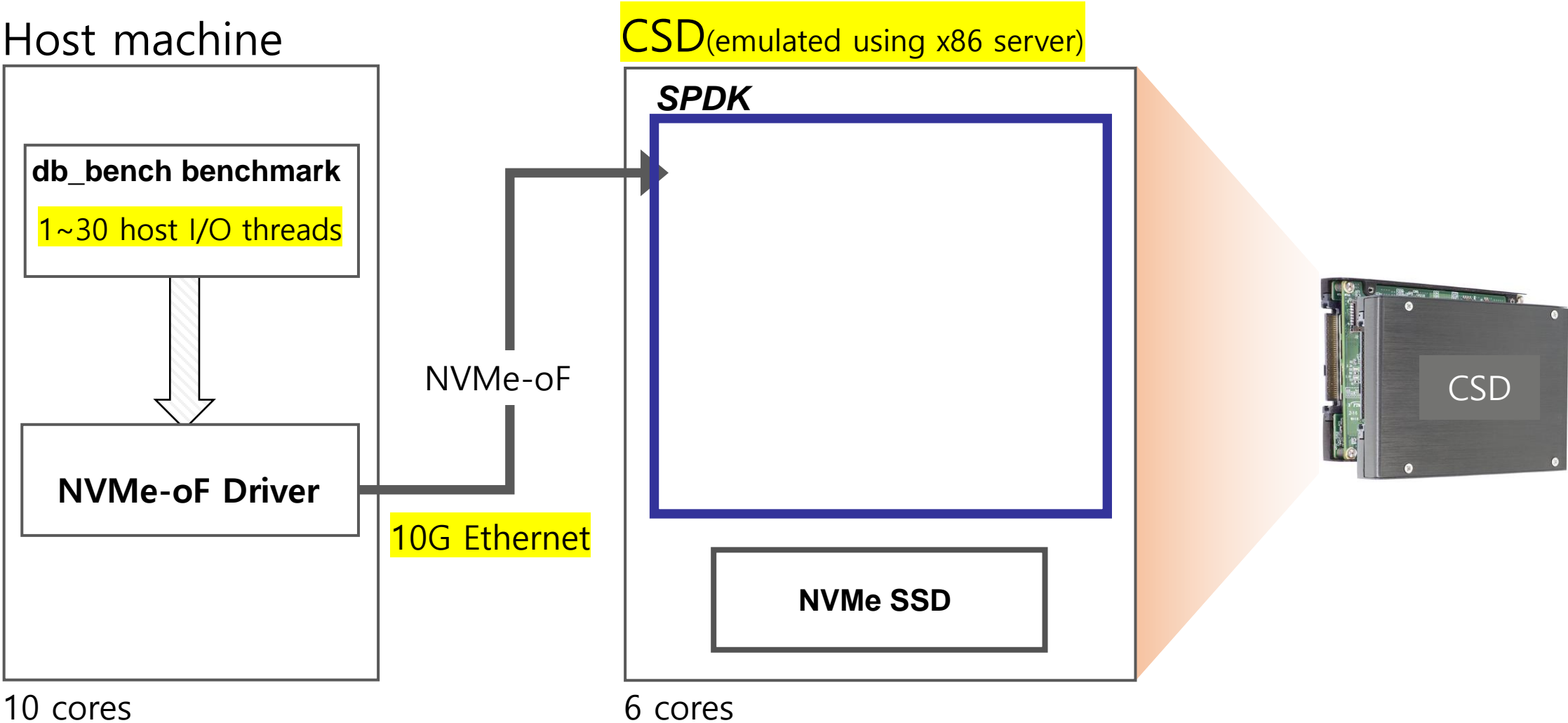
*Thread1*

*Thread2*

*Thread3*

*Thread4*

# Implementation

**SPDK**



*Thread1*

Foreground I/O

*Thread2*

*Thread3*

*Thread4*

**I/O scheduler VBDEV**

Foreground I/O

Thread 1

**BTS** ② *Select idle cores to execute background tasks*

1. **Monitoring**

Core 2 : idle
Core 4 : idle

2. **Core Selection**

Group G

Core 2
**Core 4 : selected**

*Thread1*

*Thread2*

*Thread3*

*Thread4*

# Implementation

**SPDK**

# Implementation

# Evaluation

# Experimental Setup

Host machine

CSD(emulated using x86 server)

SPDK

**db_bench benchmark**

1~30 host I/O threads

NVMe-oF Driver

NVMe-oF

10G Ethernet

NVMe SSD

CSD

10 cores

6 cores

# Experimental Setup



**Host machine**

**db_bench benchmark**

1~30 host I/O threads

**NVMe-oF Driver**

10 cores

NVMe-oF

10G Ethernet

**CSD**(emulated using x86 server)

*SPDK*

NVMe-oF Target

I/O scheduler ( BTS)

KV

Background (dedup)
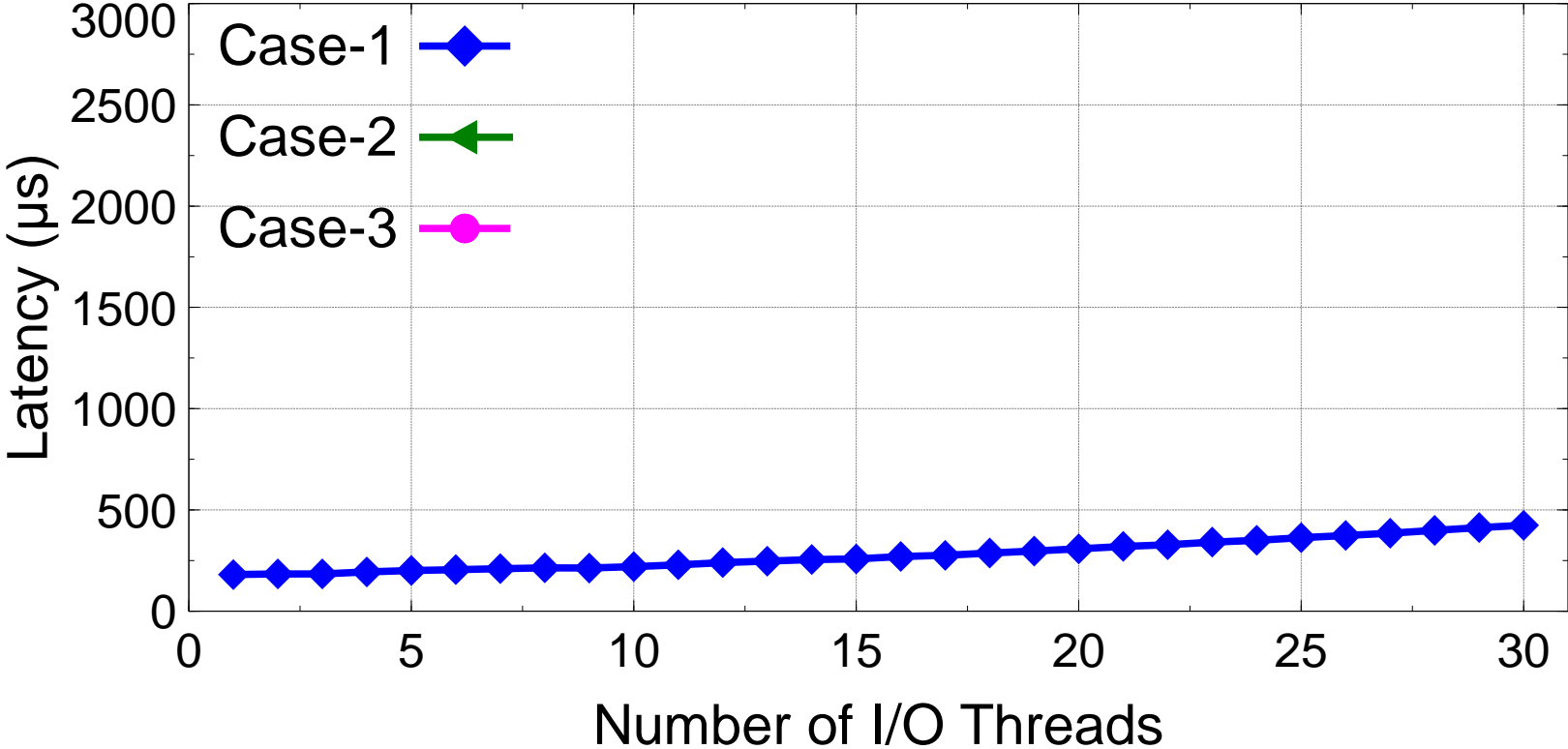
NVMe Driver

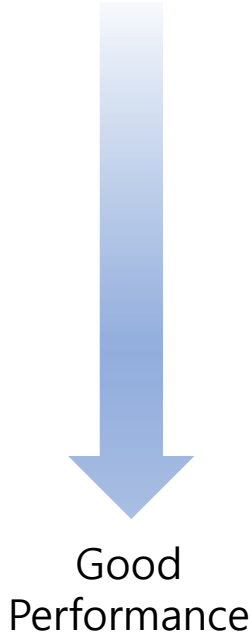**NVMe SSD**

6 cores

CSD

# Experimental Setup

❑ Comparisons

➢ Case1: Only foreground I/O

➢ Case2: Foreground I/O + Deduplication without BTS
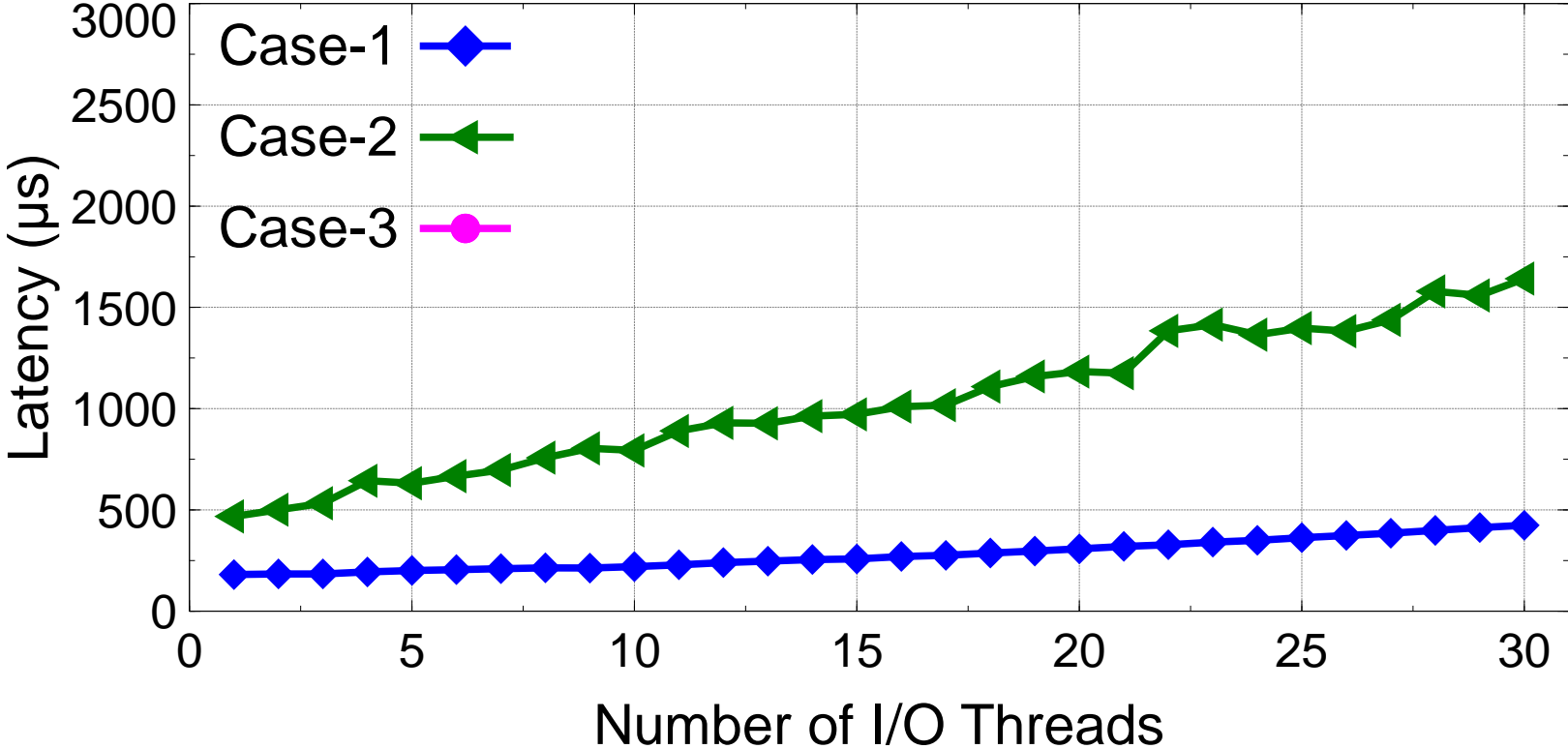
➢ Case3: Foreground I/O + Deduplication with BTS

# Put() Performance



Case 1 : Only foreground I/O
Case 2 : Foreground I/O & Dedup without BTS
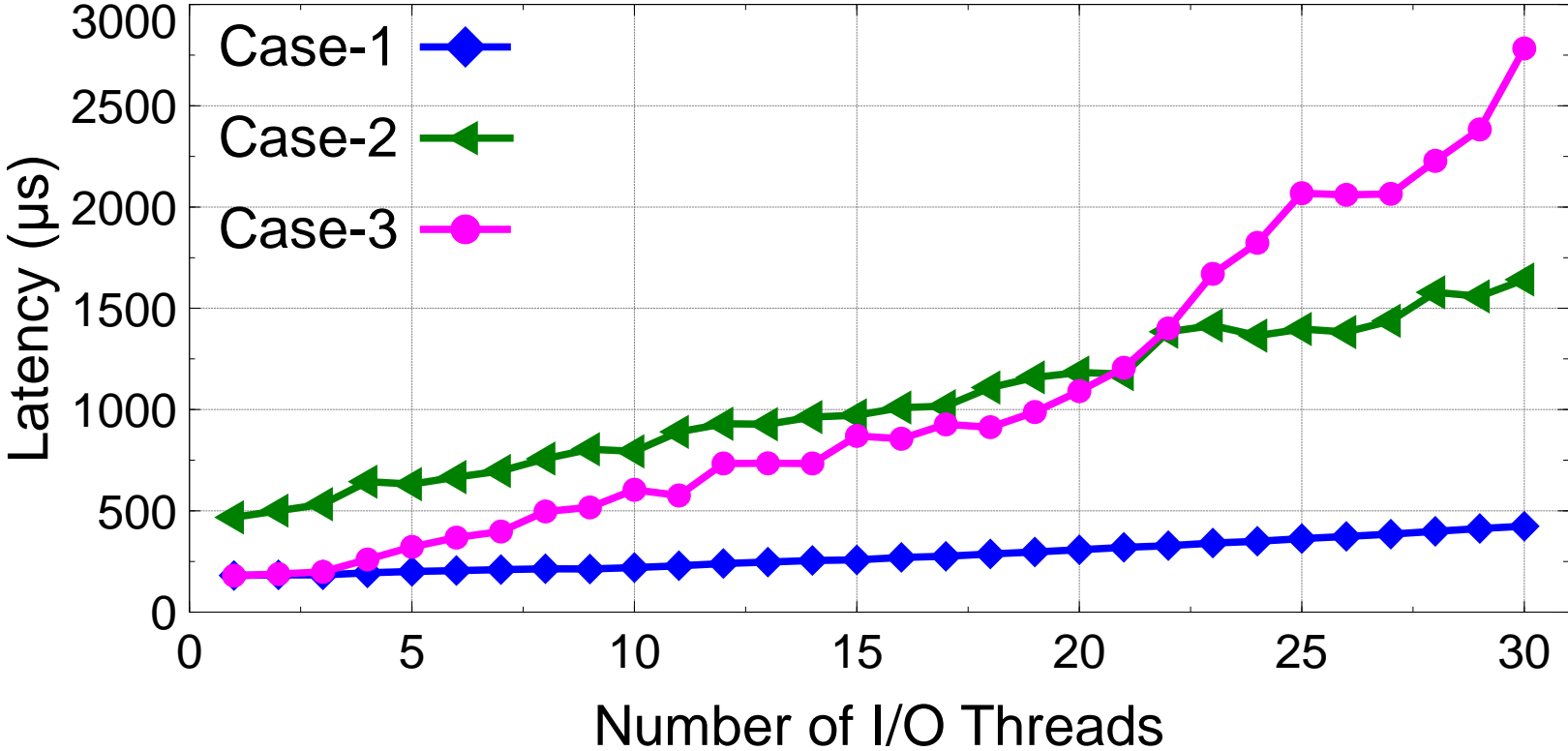Case 3 : Foreground I/O & Dedup with BTS
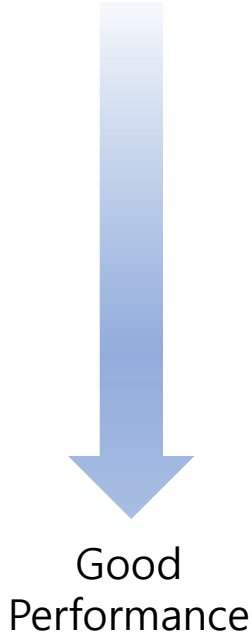
Good Performance

# Put() Performance



Case 1 : Only foreground I/O
Case 2 : Foreground I/O & Dedup without BTS
Case 3 : Foreground I/O & Dedup with BTS

# Put() Performance

Case 1 : Only foreground I/O
Case 2 : Foreground I/O & Dedup without BTS
Case 3 : Foreground I/O & Dedup with BTS
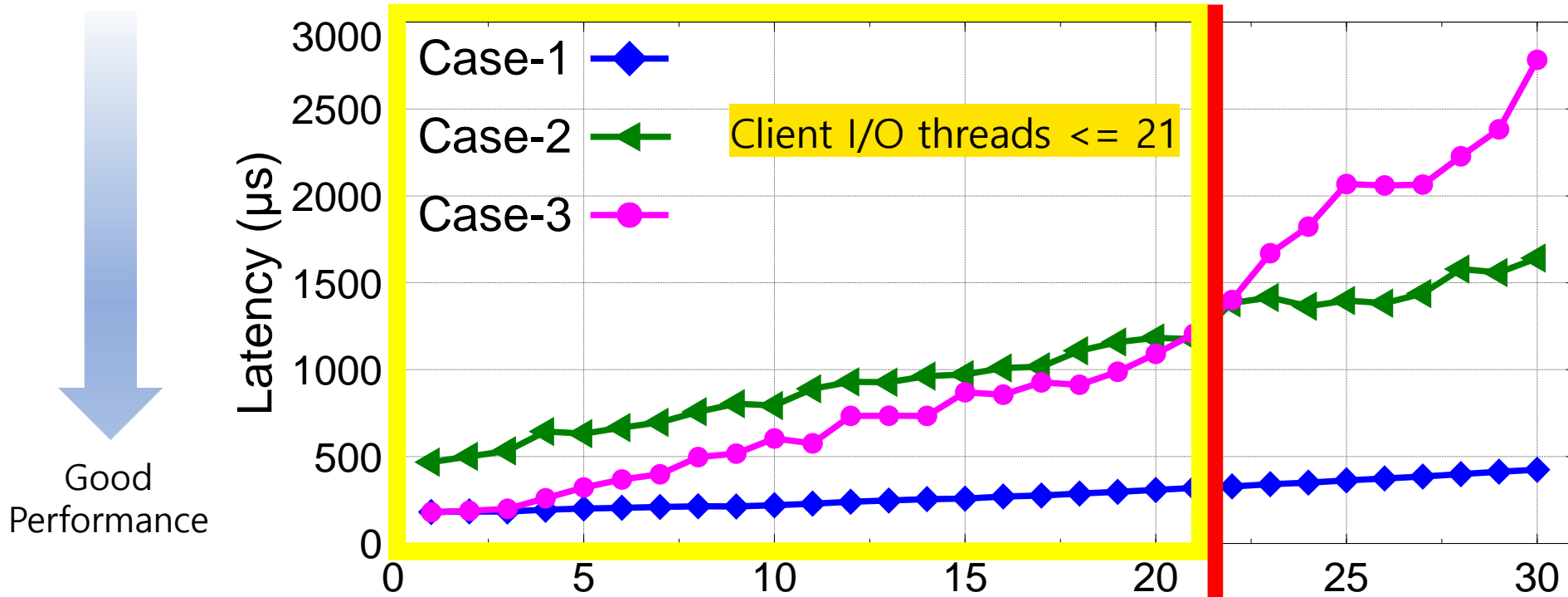


Good Performance

# Put() Performance



Case 1 : Only foreground I/O
Case 2 : Foreground I/O & Dedup without BTS
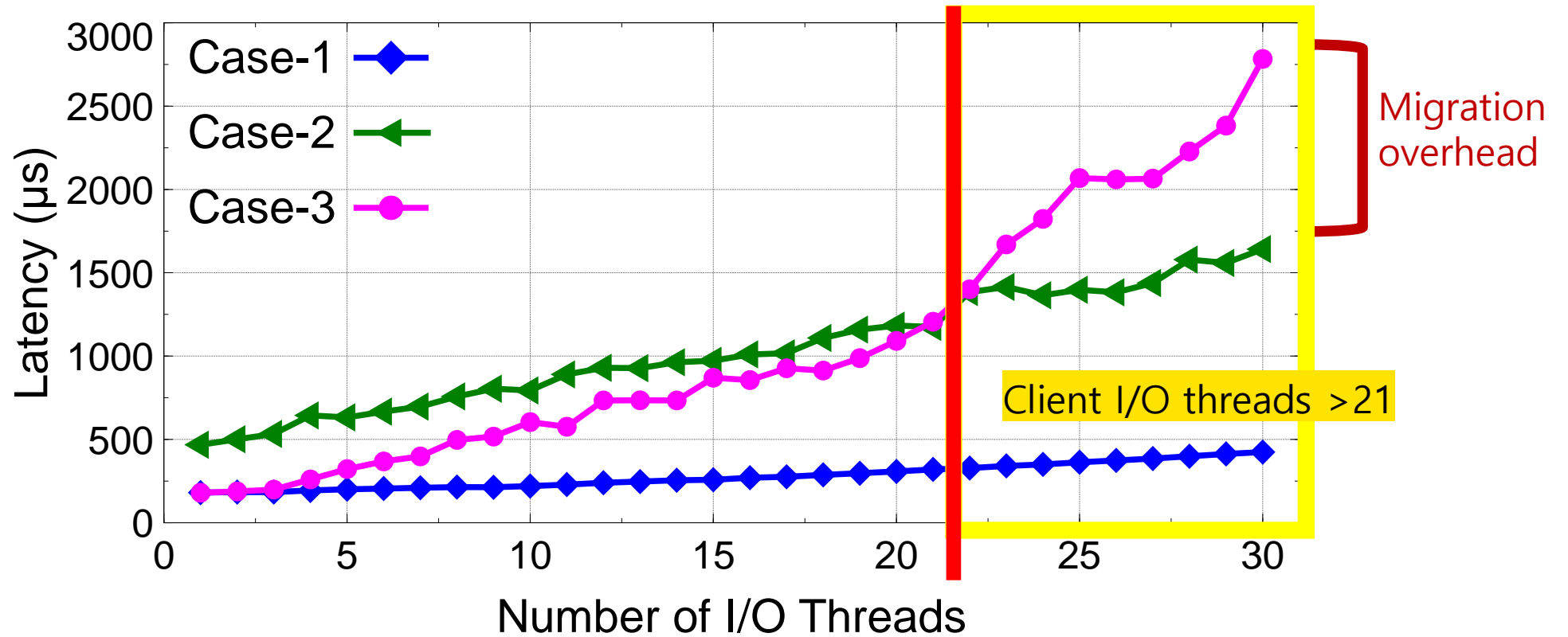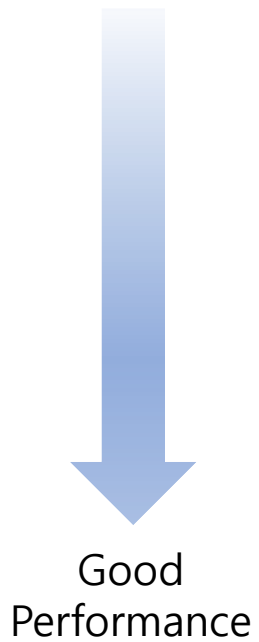Case 3 : Foreground I/O & Dedup with BTS

BTS reduced latency by an average of 47.8% when the number of host I/O threads was less than 10.

# Put() Performance



Case 1 : Only foreground I/O
Case 2 : Foreground I/O & Dedup without BTS
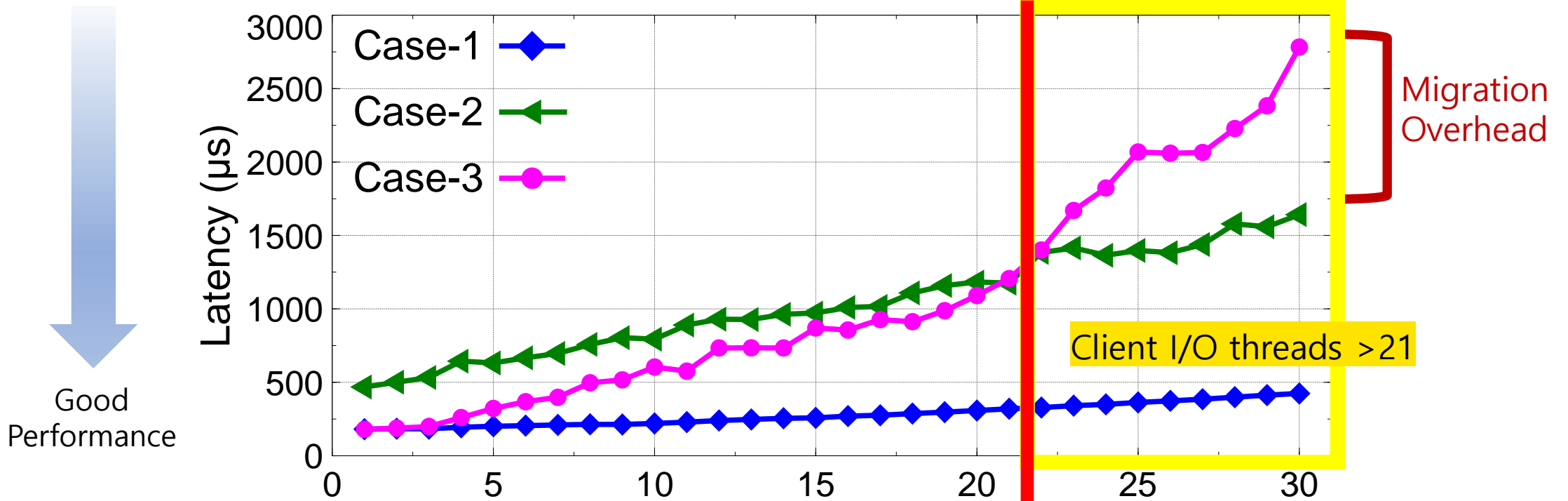Case 3 : Foreground I/O & Dedup with BTS

Good Performance

Migration overhead

Client I/O threads >21

Case-1
Case-2
Case-3

Latency (μs)

Number of I/O Threads

# Put() Performance



Case 1 : Only foreground I/O
Case 2 : Foreground I/O & Dedup without BTS
Case 3 : Foreground I/O & Dedup with BTS

Migration Overhead

Client I/O threads >21

Good Performance

Case-1
Case-2
Case-3

Latency (μs)

Under heavy load, the migration overhead for background tasks outweighs the performance gains.

# Conclusion

❑ We have identified a problem with SPDK where background tasks increase the response time of foreground I/O in CSD using SPDK

❑ We proposed a **Background Task-Aware Scheduler (BTS)** in SPDK for CSD

❑ Comprehensive evaluation showed that the BTS scheduler is effective when core utilization is rather low

# Conclusion

❑ We have identified a problem with SPDK where background tasks increase the response time of foreground I/O in CSD using SPDK

❑ We proposed a **Background Task-Aware Scheduler (BTS)** in SPDK for CSD

❑ Comprehensive evaluation showed that the BTS scheduler is effective when core utilization is rather low

❑ **BTS can be applied to any storage system using SPDK**

# Thank you!

Yeohyeon Park
yeohyeon@sogang.ac.kr