



*BY Developers FOR Developers*

**Storage Developer Conference**  
**September 22-23, 2020**

# Unlocking the New Performance and QoS Capabilities of the **Software-Enabled Flash™ API**

**Rory Bolt**  
**Principal Architect and Sr. Fellow**  
**KIOXIA America, Inc.**



**01**

# **What is Software-Enabled Flash?**

## Software-Enabled Flash™ Technology

**Fundamentally  
redefines the  
relationship  
between host  
and solid-state  
storage**

- **Brings control of media to the host**
- **Host applications have complete control over storage functionality and behavior**
- **Solves legacy overhead problems and enables new features**
- **Maximizes flash flexibility, performance and parallelism...**

**in other words, its value.**

## A Software Enabled API

- 1 A software enabling technology built around an open source flash-native API
- 2 It delivers a rich interface of functions and tools to simplify storage innovation
- 3 The API abstracts flash details enabling future generations of flash to work without code changes
- 4 Allows any flash vendor to build and optimize their flash to the API
- 5 KIOXIA will provide sample source code and libraries in the future

**Software-Enabled  
Flash Technology**

**is not  
software  
but it is**

**software  
enabled**

# Hardware and Software Working Together

Purpose-built, media-centric flash controllers focused on **hyperscaler** requirements.

- Flash vendor/generational abstraction
- Advanced die time scheduling
- Access to entire media
- Host CPU offload
- Flexible DRAM configurations

Open sourced API and libraries providing functionality **hyperscalers** demand.

- Data placement
- Workload / tenant isolation
- Latency control
- Buffer management
- Adaptable to new workloads

# Provided By the API

## Hardware Structure

- Channels/Dies
- Parallelism
- Flash structure
  - Plane/Block/Page
  - Defect management
  - Lifetime management

## Buffer Management

- Write buffer
- Number of open blocks

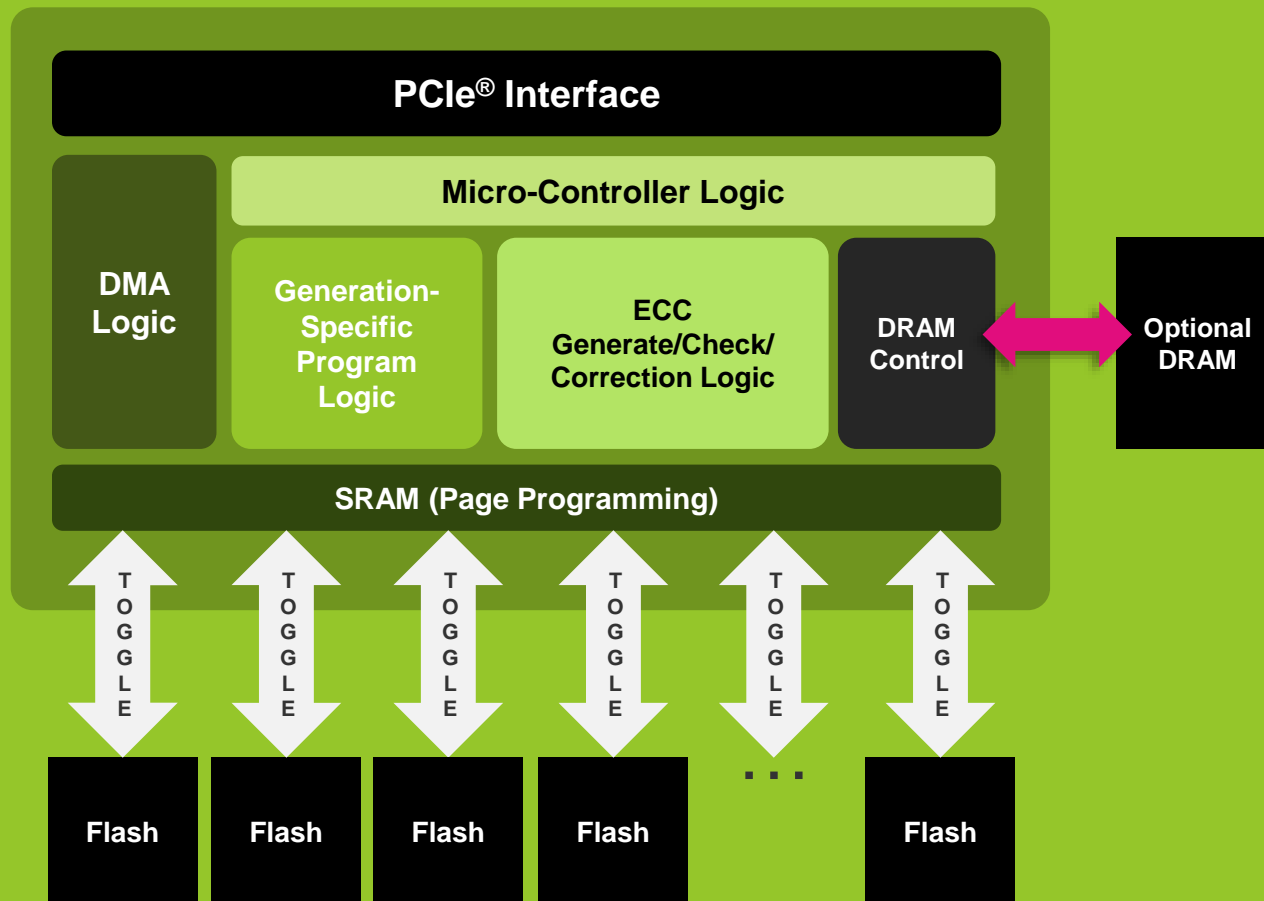
## Programming Algorithm

- Physical placement
- Number of passes
- Voltage threshold
- Thermal management

## Error Management

- Error reduction
- ECC deadline

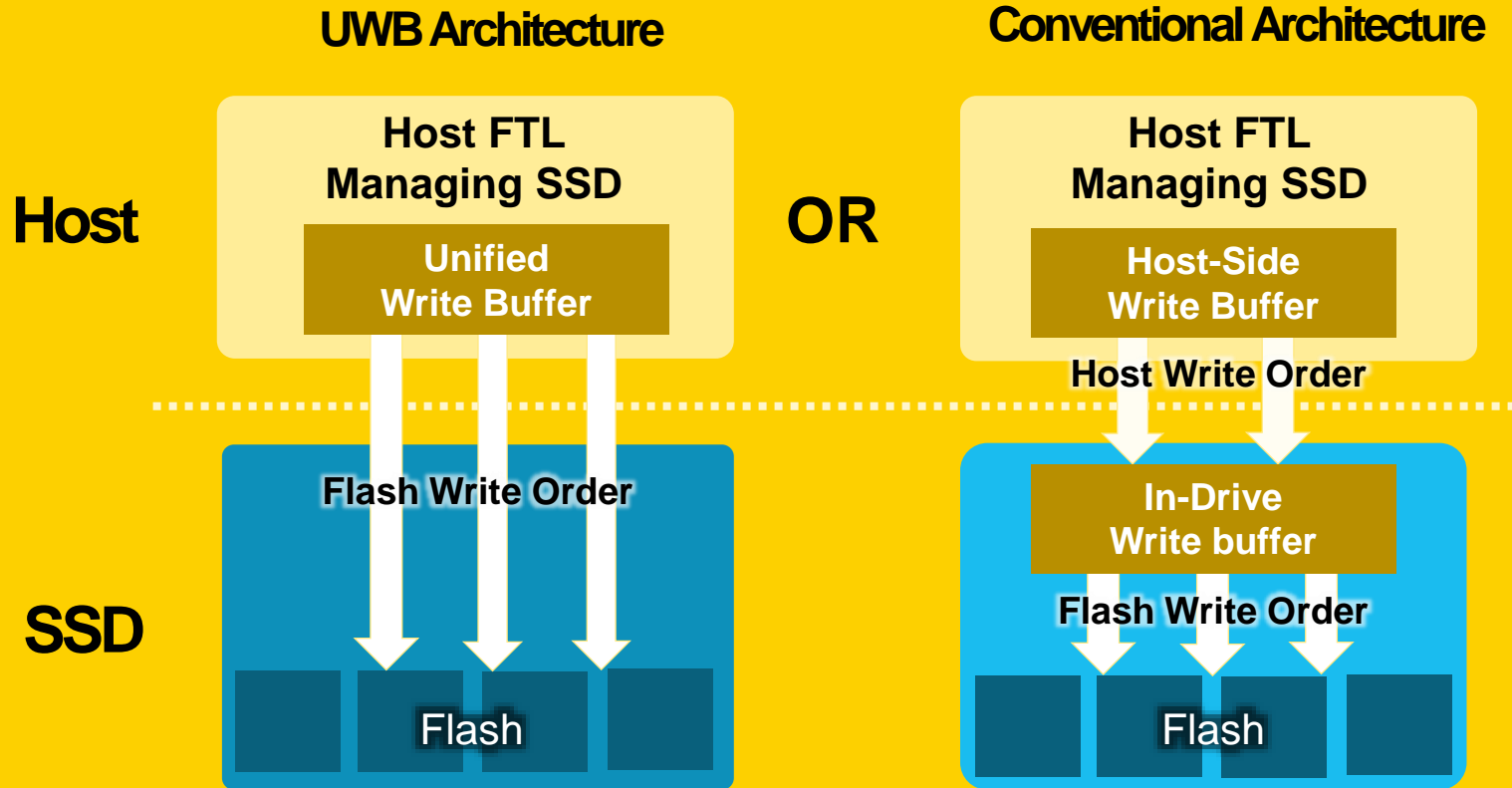
# Purpose-Built Hardware to Manage Media



PCIe is a registered trademark of PCI-SIG.

- **Controller handles all media details**
  - Programming algorithms
  - Lifetime extension
  - Defect management
- **Multiple prioritized and weighted queues per die**
- **Copy offload functionality**
- **ECC hardware**
- **Zero DRAM configurations (SRAM only)**

# Unified Write Buffers (UWB) offers Flexible use of DRAM



**Supports both architectures as well as hybrid implementations!**

## Unified Write Buffer Benefits:

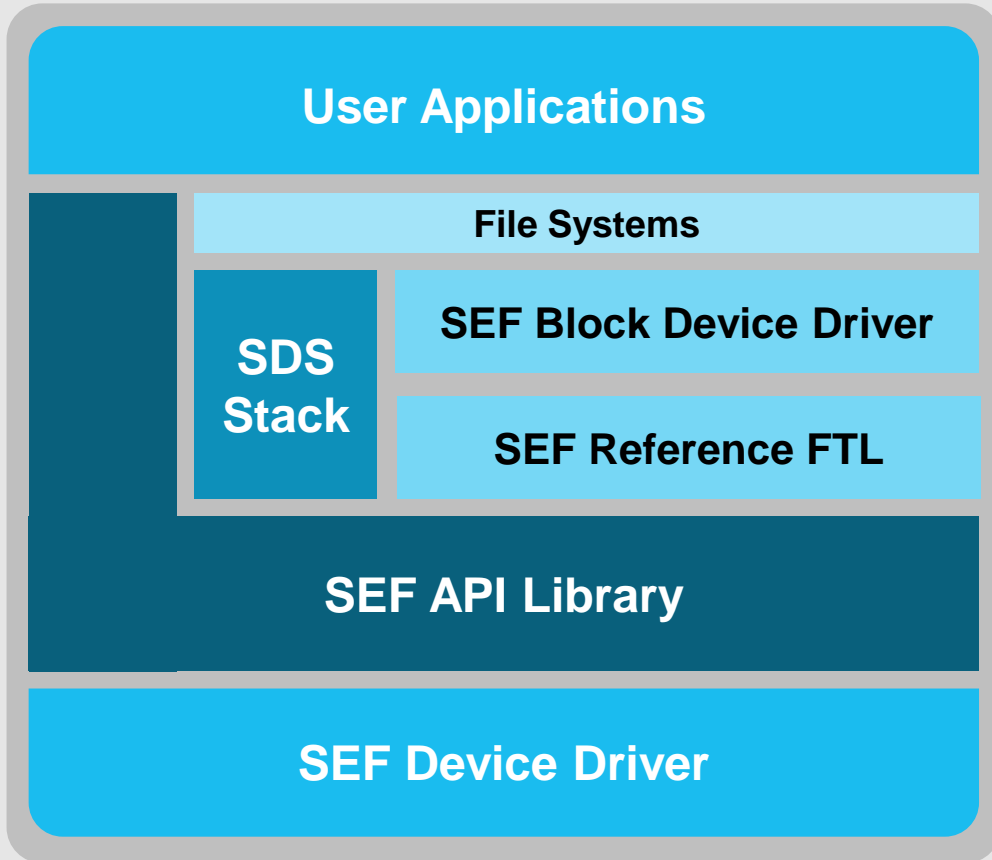
- ✓ Can dynamically adapt to changing workloads
- ✓ Can optimize DRAM at the system level

## Conventional Architecture

- ✓ Appropriate for systems that rely on drive power loss protection



# Software Components and Layering



## A Software Development Kit (SDK)

- Provides an Open Source Linux<sup>®</sup> reference block driver
- Provides an Open Source Linux<sup>®</sup> reference FTL

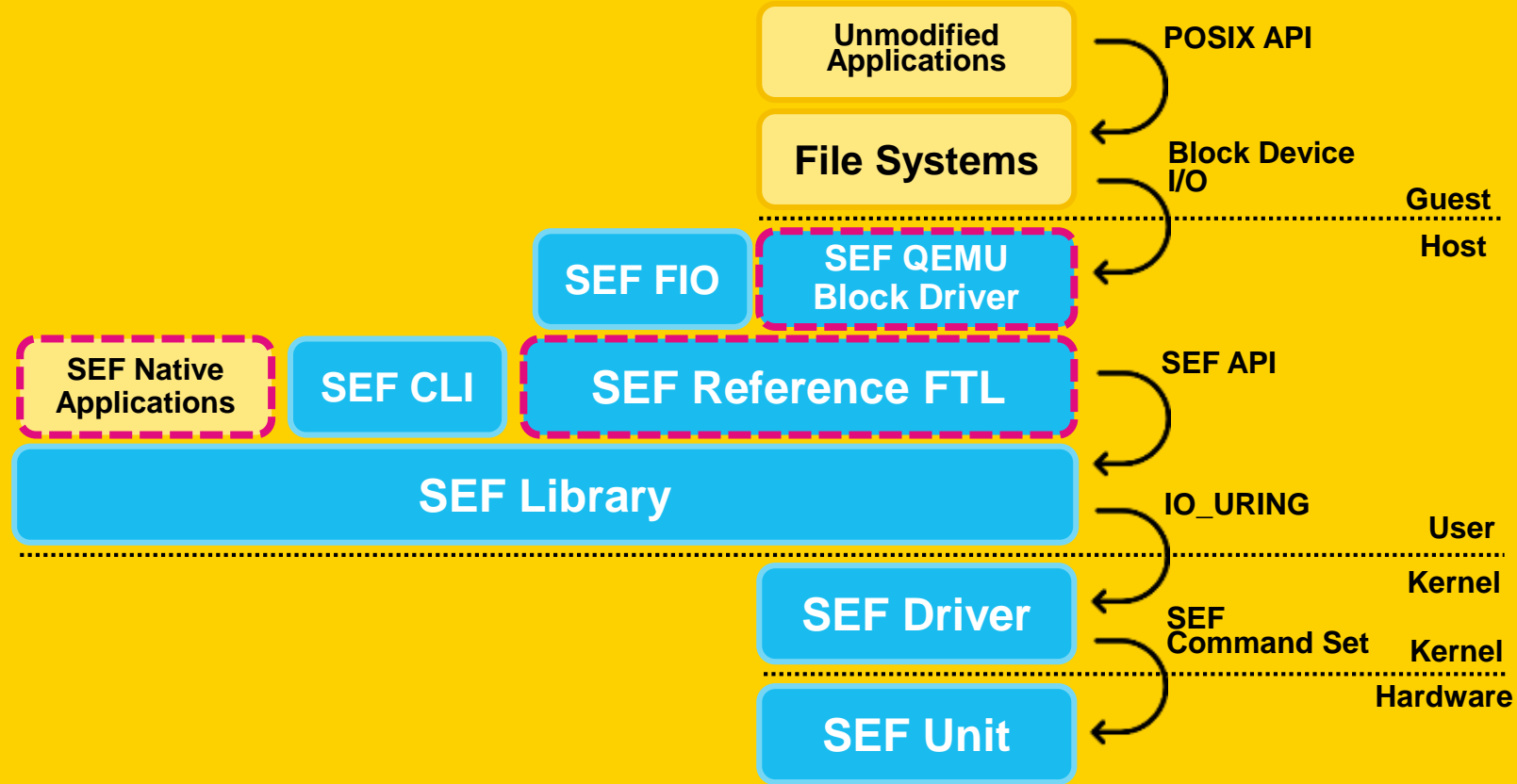
Linux is the registered trademark of Linus Torvalds in the U.S. and other countries.

## An API Library

- Open Source Linux interface library that implements the API (Application Programming Interface)
- C language with C++ wrappers, data structures, enumerated types, utility functions
- User space and kernel implementations
- Hides flash generational differences such as programming algorithms
- Controls advanced scheduling

## The Software-Enabled Flash Device Driver

# A Possible Software-Enabled Flash deployment

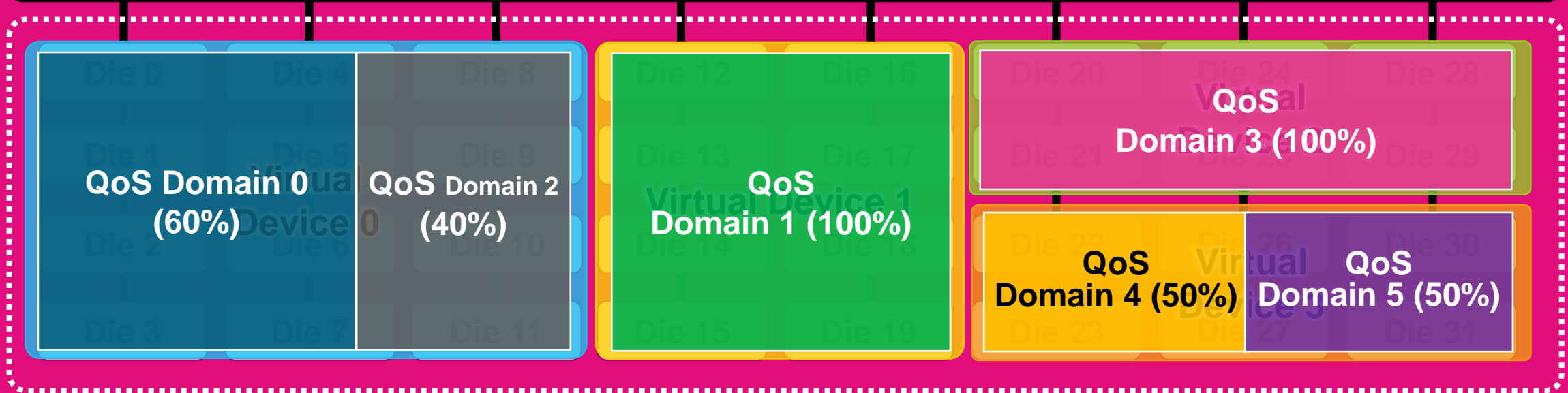


**02**

# **Software-Enabled Flash Concepts and Technologies**

# Placement Control and Isolation

## SEF Controller



### CONCEPTS

**Virtual Devices** – a set of one or more flash die providing hardware isolation

**QoS Domains** – a mechanism to impose capacity quotas and scheduling policy to provide software based isolation

**Placement IDs** – a mechanism to group data at the superblock level

# Data Placement Mechanism for Data Isolation

Isolation is done via the data placement mechanisms through definition of virtual devices and QoS domains.

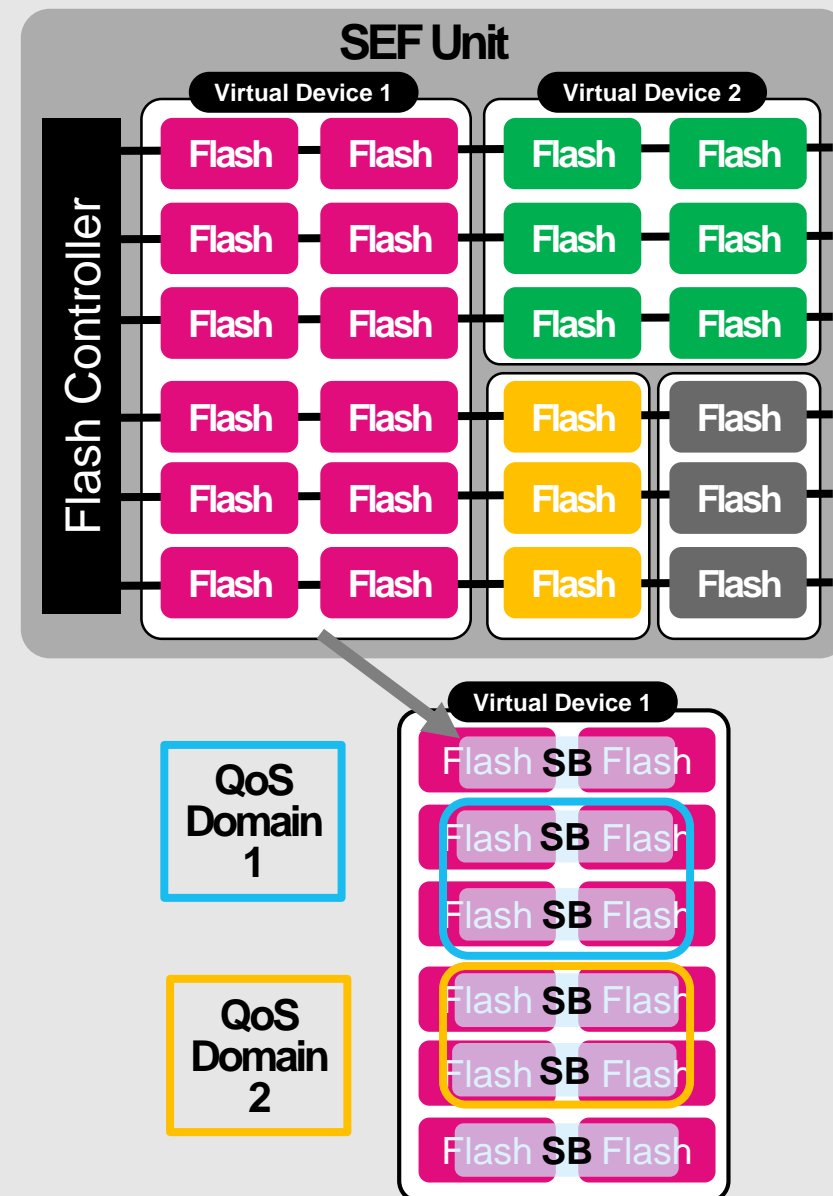
## Virtual Device: physical isolation at die level

- Capacity and I/O throughput can be defined
- Configures “Superblock (\*)” which is the minimum data set size for internal data management.

(\*) Superblocks(SB) are striped across each die in the virtual device

## QoS Domain: logical isolation at flash block level

- QoS Domain is a data container within the virtual device
- Superblocks are assigned to a QoS domain
- API manages life-cycle of superblocks



# Isolation Capability

## Die Level Isolation (Virtual Devices)

**Most effective,  
least scalable.**

- Full isolation is possible for small numbers (tens) of tenants
- Practical number of die is limited, so flash device cannot service large numbers of tenants with die level isolation

## Block Level Isolation (QoS Domains)

**Most scalable,  
not total isolation.**

- To enable predictable performance and extend lifetime, block level isolation is very effective for large number (thousands) of tenants

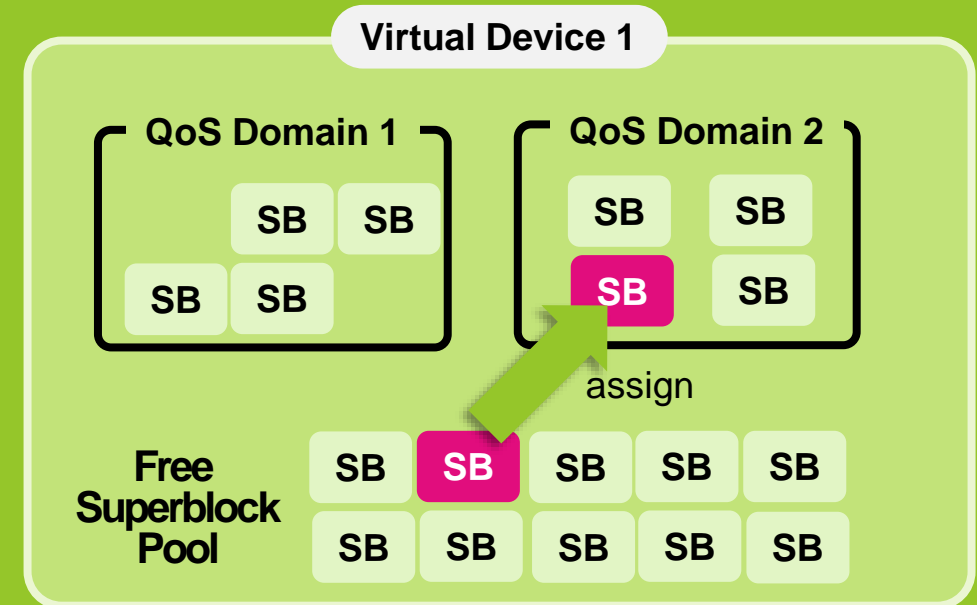
# Nameless Write Mechanism (1) - “Automatically assigns a block to write”

Host wants control over data placement to minimize write amplification (WAF).



Optimal block selection is best left to the device to ensure media life and health.

- SEF Unit will automatically select most appropriate superblock and assign it to the QoS domain (applies to both manual and automatic allocation)
- Ownership of the superblock may change after it is released to the free pool

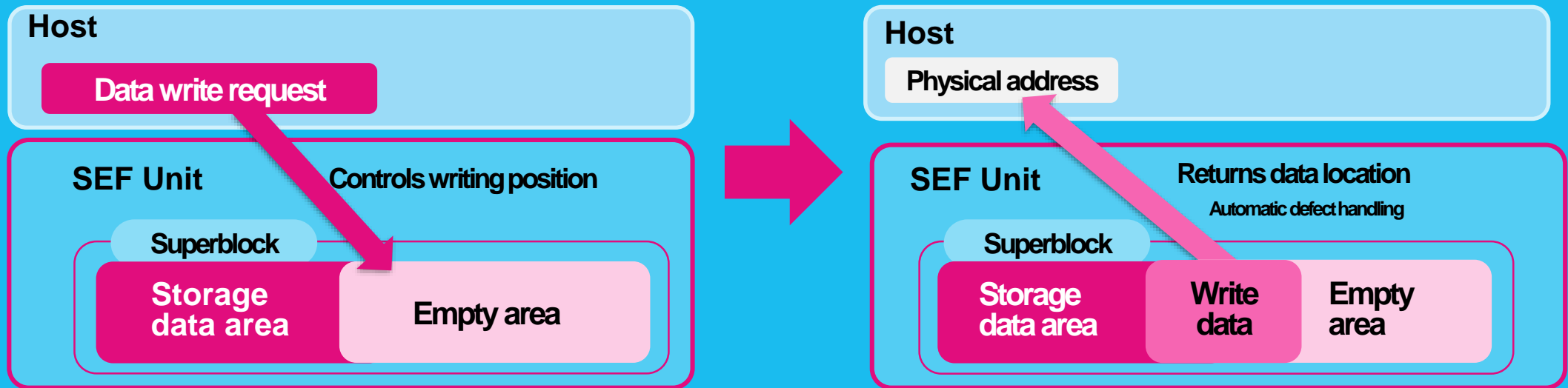


# Nameless Write Mechanism (2) “Bounded Automatic Data Placement”

## Write Operation:

**QoSDomainID** and **PlacementID** or **FlashAddress** bounds selection of superblock

- Host does not specify physical address for write
- Physical address returned upon write completion



## Read Operation:

**Direct physical address read**

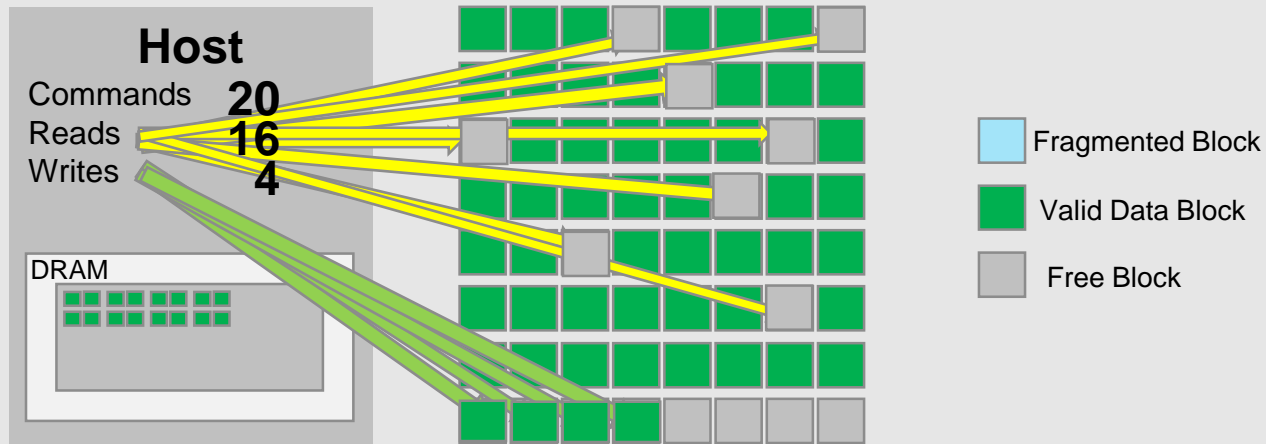
- Logical-physical address conversion process is not required by device
- Minimal read overhead and latency



# Nameless Copy Mechanism – “Offloading Garbage Collection(GC)”

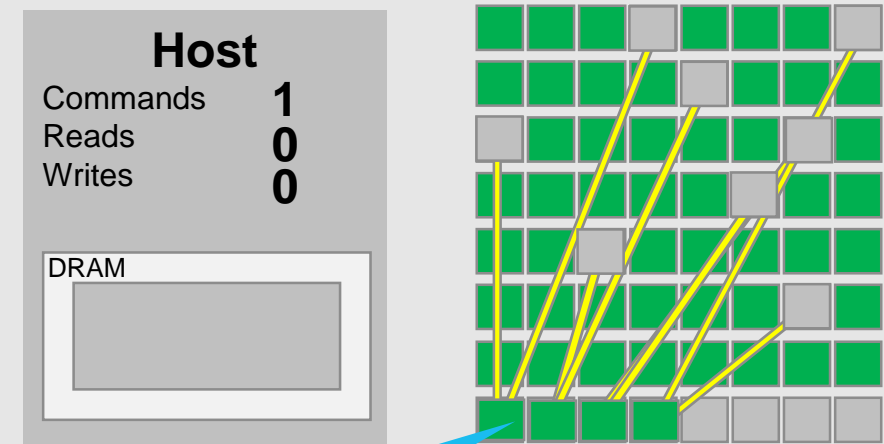
“Nameless Copy” is a copy offload function that implements internal data copy from die to die without the host moving data. Garbage collection becomes fast and efficient.

## Manual Copy Operation



Garbage collection “read and write” data flows through the host

## Nameless Copy Operation



Host controls when to garbage collect, data is copied internally by device

## Manual Copy




## CPU Utilization



## Nameless Copy



## User Workload

	User Write Data	0
	Elapsed Time	0
	CPU Utilization (Sec)	0

## Garbage Collection

	Copied SuperBlocks	0
	Released SuperBlocks	0
	Copy Commands	0
	GC Read Commands	0
	GC Write Commands	0

 User Workload  
 Garbage Collection

## User Workload

	User Write Data	0
	Elapsed Time	0
	CPU Utilization (Sec)	0

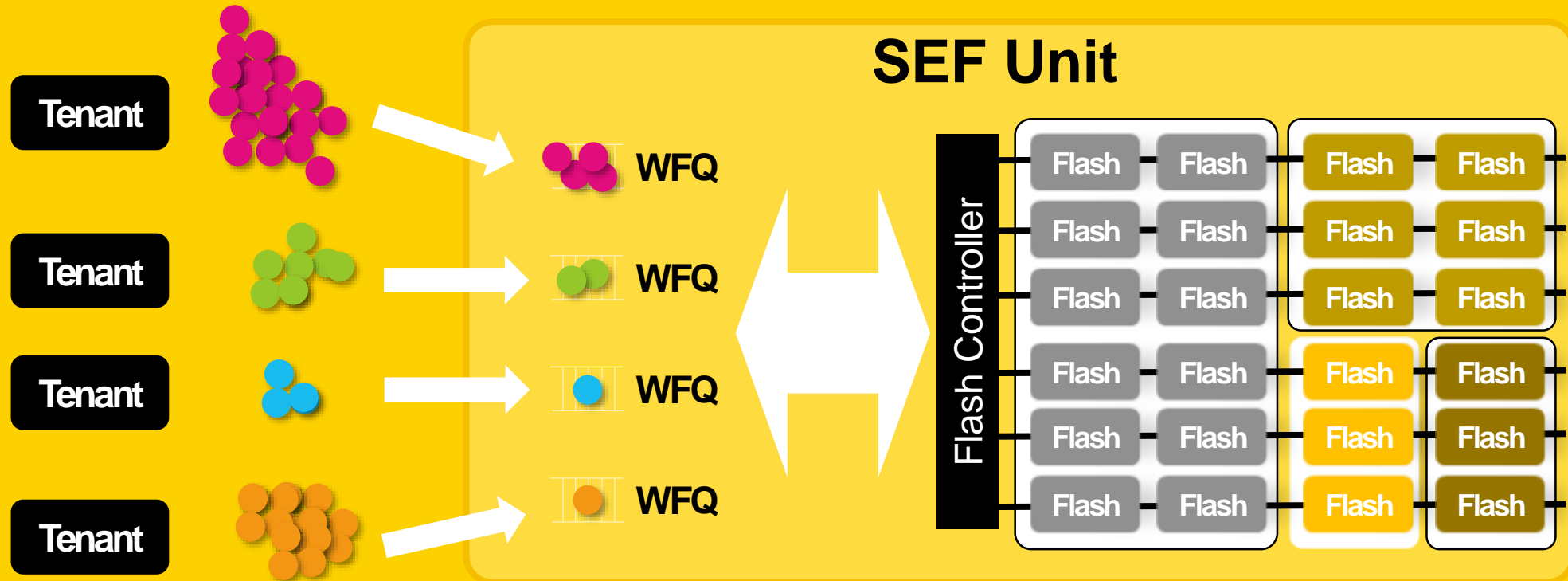
## Garbage Collection

	Copied SuperBlocks	0
	Released SuperBlocks	0
	Copy Commands	0
	GC Read Commands	0
	GC Write Commands	0

# Weighted Fair Queuing (WFQ)

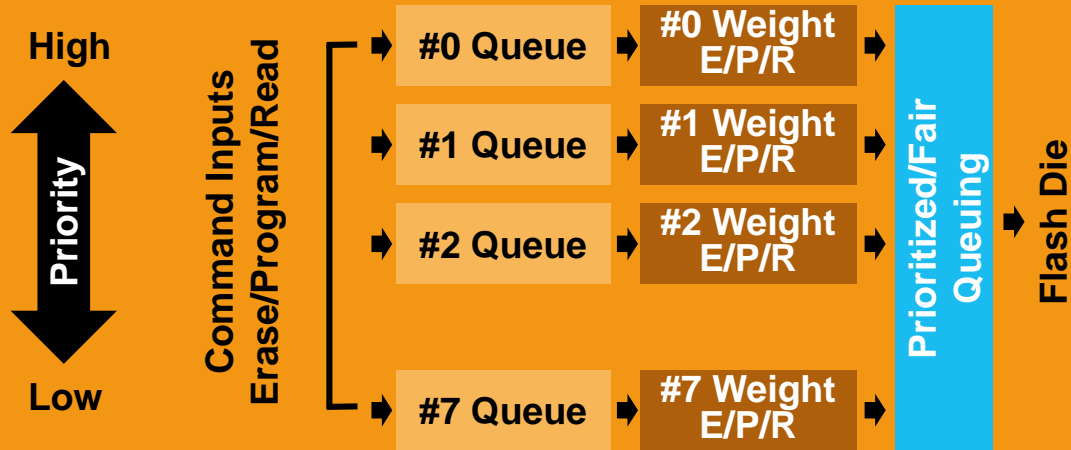
## Control of processing fairness provides adjustable service levels for each tenant

- Host can control the priority ratio of write and read requests per QoS domain
- SEF unit implements a die-based IO fairness control



# Controller Architecture: Die Scheduling and Queuing

## Functionality



- When all weights are 0, it works as an 8-level priority scheduler
- When all weights are same value, it works as a round-robin scheduler
- When weights are unique, it works as a die time weighted scheduler

## FEATURES

- Each virtual device has 8 FIFO command queues (up to 1:1 ratio with die!)
- Device scheduler handles suspend/resume for program and erase commands
- Host can specify a specific queue for each flash access command (read, program, erase) per QoS domain
- Every queue can specify die time weighting for read, program, and erase commands
- Host can override defaults queues and weights for individual commands

Start Test

### Read Weight

QoS Domain 1  
(198)

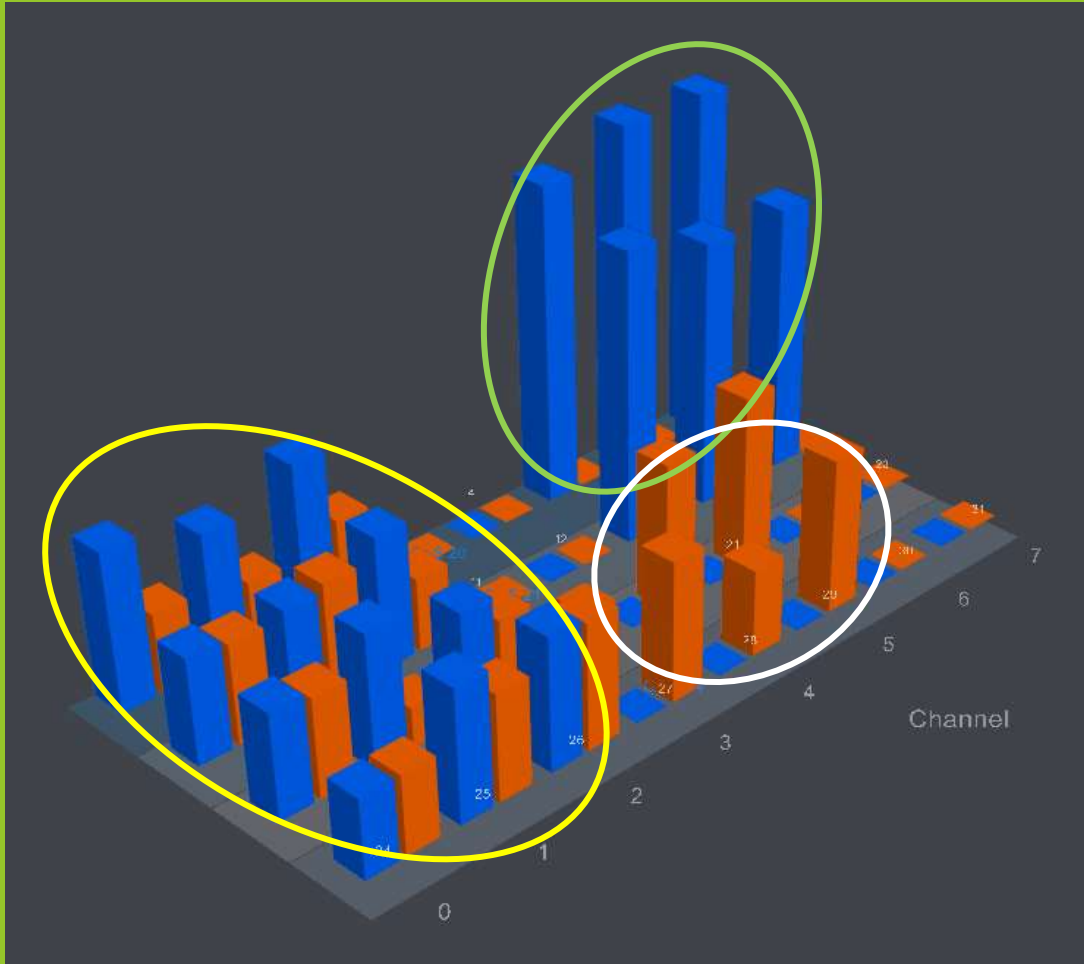
QoS Domain 2  
(202)



LATENCY

QoS Domain 1  QoS Domain 2

# Software-Enabled Flash Adaptability... Demonstrated



- A single device/SKU that is customizable and reconfigurable
- ZNS, standard SSD, Software-Enabled Flash native RocksDB, and a custom hyperscale FTL all running simultaneously on same device
- As technologies, applications, and workloads change over time, Software-Enabled Flash can be adapted to changing needs

**03**

# **API Examples**

# Sef-cli

```
01  ===== General Usage =====
02  ./sef-cli <target-specific-action> <target> -[hv] [target-specific-parameters]
03  sef-cli supports many different engines and each engine support multiple actions.
04  In order to print which arguments are supported by each target use the verbose flag.
05
06  ===== General Arguments =====
07  The following arguments are general and are applied to all SEF Cli Targets
08  -h, --help          Print help and exit
09  -, --version       Print version number and git commit hash and exit
10  -V, --verbose      Explain what is being done or show additional information regarding the action being taken
11
12  ===== Target =====
13  The following targets are supported by your sef-cli software:
14  virtual-device     Perform actions towards Virtual Devices
15  superblock        Perform actions towards superblocks
16  shell             A Python Shell that is extended with SEF functionalities in order to interact with the device
17  sef-unit          Perform actions towards mounted SEF units
18  sdk               Perform actions towards SEF SDK (block layer)
19  root-pointer      Perform actions towards QoS Domain's root pointers
20  qos-domain        Perform actions towards QoS Domains
21  adu              Perform direct IO towards the SEF
```



# Practical Example: Creating a Virtual Device

```
01  #!/bin/bash
02  # Practical Example: Mapping SEF unit 0 to a Single Virtual Device
03  # The following script assumes an initially un-configured SEF unit
04  #
05  # Note that this script should be run with sudo permissions
06
07  /usr/local/bin/sef-cli create virtual \
08      -s0 \
09      --start-channel=0 \
10      --num-channel=4 \
11      --start-bank=0 \
12      --num-bank=4 \
13      --virtual-device-id=0 \
14      --max-qos-domains=1
15
16  # To see all the virtual devices on unit 0:
17  /usr/local/bin/sef-cli list virtual -s
```

# Practical Example: Creating a QoS Domain

```
01  #!/bin/bash
02  # Practical Example: Create a QoS Domain on Virtual Device 0 on SEF unit 0
03  #
04  # Note that this script should be run with sudo permissions
05
06  /usr/local/bin/sef-cli create qos \
07    -s0 \
08    --virtual-device-id=0 \
09    --qos-domain-id=2 \
10    --flash-capacity=3000000 \
11    --adu-size=4096 \
12    --num-root-pointers=4 \
13    --num-placement-id=4
14
15  # To see all QoS Domains on SEF Unit 0:
16  /usr/local/bin/sef-cli list qos -s0
```

# API Example: Nameless Write

```
01 // Get the handle for the SEF Unit
02 sefHandle = SEFGetHandle(sefUnitIndex);
03
04 // Open QoS domain
05 status = SEFOpenQoSDomain(sefHandle, qosDomainId, notifyFunc, context, key, &qosHandle);
06
07 // Write data to SEF
08 status = SEFWriteWithoutPhysicalAddress1(qosHandle, SEFAutoAllocate, placementId, userAddress,
09 numADU, &iov, 1, &permanentAddress, &distanceToEndOfSuperblock, &overrides);
```

# API Example: Asynchronous Event Handler

```
01 // Open QoS domain
02 status = SEFOpenQoSDomain(sefHandle, qosDomainId, HandleSEFNotification, context, key, &qosHandle);
03
04 void HandleSEFNotification(void *context, struct SEFQoSNotification event) {
05     FTLContext *ctxt = (FTLContext *) context;
06     switch (event.type) {
07         case kAddressUpdate:
08             SEFFTLUpdate(ctxt, SEFGetUserAddressLba(event.changedUserAddress),
09                         event.oldFlashAddress, event.newFlashAddress);
10             break;
11         case kSuperblockStateChanged:
12             SSBStateEvent(ctxt, event.changedFlashAddress, kSSBEventClosed);
13             break;
14         ...
15     }
16 }
```

# Lockless Lookup Tables – Address Update

```
01 // This is the non-authoritative way to update the lut[]. First writer wins
02 // implies the data has been moved and not updated
03
04 int SEFFTLUpdate(FTLContext *ctxt, int64_t userLBA, struct SEFFlashAddress oldFlashAddress,
05                 struct SEFFlashAddress newFlashAddress) {
06     bool lutUpdated;
07     uint64_t expectedBits = oldFlashAddress.bits;
08     // Mark new flash address as used in superblock map
09     // Update lut[] with the new flash address as long as the old was there
10     // Mark old flash address as free in superblock map, if lut[] updated
11     // Mark new flash address as free in superblock map, if lut[] failed update
12     SSBSetAduValid(ctxt,newFlashAddress); // In case exchange works, must be set
13     lutUpdated = atomic_compare_exchange_strong(&ctxt->lut[userLBA], &expectedBits, newFlashAddress.bits);
14     if (lutUpdated)
15     {
16         assert(oldFlashAddress.bits);
17         SSBClearAduValid(ctxt,oldFlashAddress);
18         ctxt->lutDirty = true;
19     } else {
20         SSBClearAduValid(ctxt,newFlashAddress);
21         // Undo it if it didn't
22     #if CC_DEBUG_COUNTERS
23         atomic_fetch_add(&ctxt->ccCounter.lateUpdate, 1);
24     #endif
25     }
26     return 0;
27 }
```

# API Example: Direct Access Read

```
01 // Get the handle for the SEF Unit
02 sefHandle = SEFGetHandle(sefUnitIndex);
03
04 // Open QoS domain
05 status = SEFOpenQoSDomain(sefHandle, qosDomainId, notifyFunc, context, key, &qosHandle);
06
07 // Read data from the flash
08 status = SEFReadWithPhysicalAddress1(qosHandle, flashAddress, numADU, &iov, iovCount, iovOffset, userAddress, &overrides);
09
10 // NOTE: The error recovery mode for a QoS Domain is set at time of
11 // creation. This determines if the SEF unit will perform automatic
12 // error recovery on the QoS Domain. In manual mode SEFSetReadDeadline
13 // determines how quickly the SEF unit must respond when attempting
14 // recovery
```

# Data Path Commands Have Asynchronous versions

```
01 struct SEFReadWithPhysicalAddressIOCB {
02     struct SEFStatus status;           /**< Library sets error field to a non-zero value to indicate ... */
03     int16_t opcode;                   /**< Should never be accessed - for internal use by library */
04     int16_t done;                      /**< Flag for polled I/O - library sets this field */
05     ALIGN_FOR_POINTER(4);
06     void *param1;                     /**< Ignored by the library; the caller can store context ... */
07     void (*complete_func)(struct SEFReadWithPhysicalAddressIOCB *);
08     const struct SEFReadOverrides *overrides; /**< Override parameters for scheduling purposes, may be NULL */
09     struct SEFlashAddress flashAddress; /**< Physical address for the read command */
10     struct SEFUserAddress;            /**< Contains LBA information */
11     const struct iovec *iov;          /**< A pointer to the scatter gather list */
12     uint32_t iovOffset;               /**< Starting byte offset into iov array */
13     uint32_t numADU;                  /**< Number of ADUs to be read, maximum is superblockCapacity */
14     uint16_t iovcnt;                  /**< Number of elements in the scatter gather list */
15 };
16
17 void SEFReadWithPhysicalAddress1Async(void *SEFHandle, struct SEFReadWithPhysicalAddressIOCB *iocb);
```

**04**

# **Software-Enabled Flash Summary & Where to Learn More**

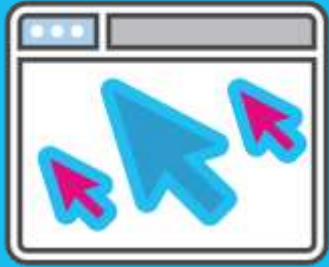


# Software-Enabled Flash™ technology fundamentally redefines the relationship between host and solid-state storage.

- ✓ Purpose-built hardware for managing flash media under host control
- ✓ Open-source, flash-native API
- ✓ Industry standards and protocols
- ✓ May be used as a building block for multiple drive protocols (*block, ZNS, custom*)
- ✓ Combines full host control with ease of use

# For More Information on Software-Enabled Flash™ Technology

1



**microsite**

[www.softwareenabledflash.com](http://www.softwareenabledflash.com)

2



**white paper**

Available on the microsite

3



**API**

<http://github.com/kioxiaamerica>

**KIOXIA**