



Xilinx Storage Services

Jamon Bowen
Product Marketing Director



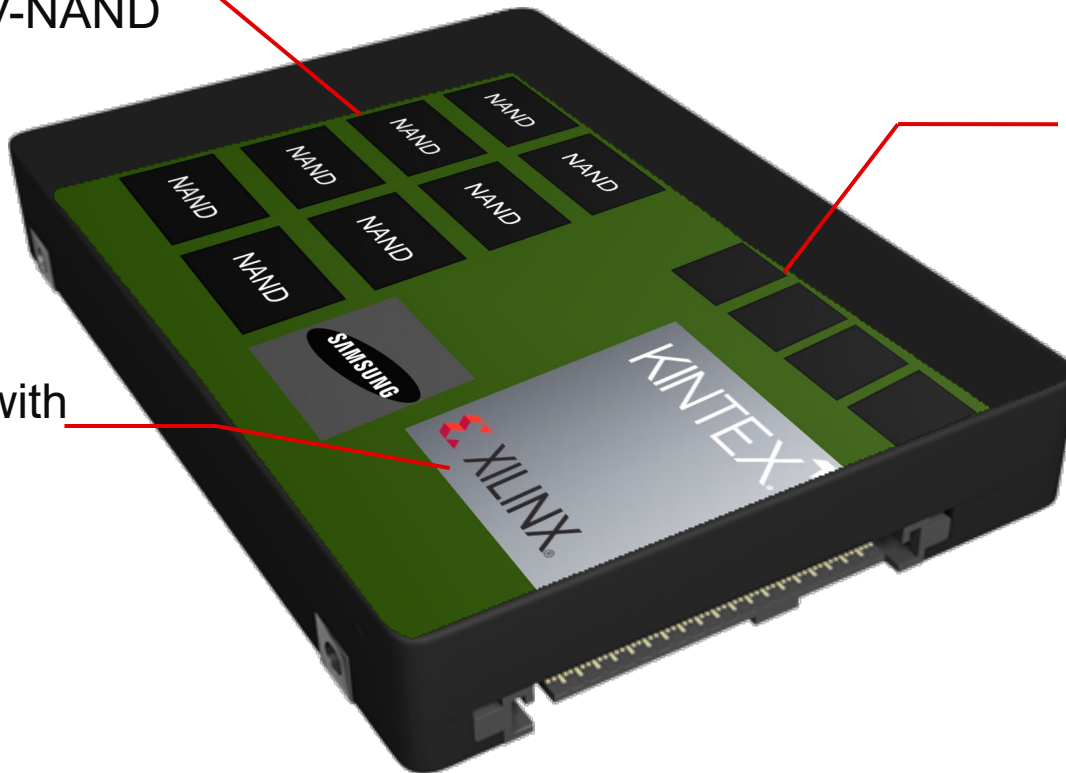
SmartSSD® CSD



From the outside it looks
just like a standard NVMe
SSD

SmartSSD® CSD

4 TB
5th Generation
Samsung V-NAND



4 GBs accelerator
memory

Xilinx FPGA with
customizable
accelerator

SmartSSD® CSD



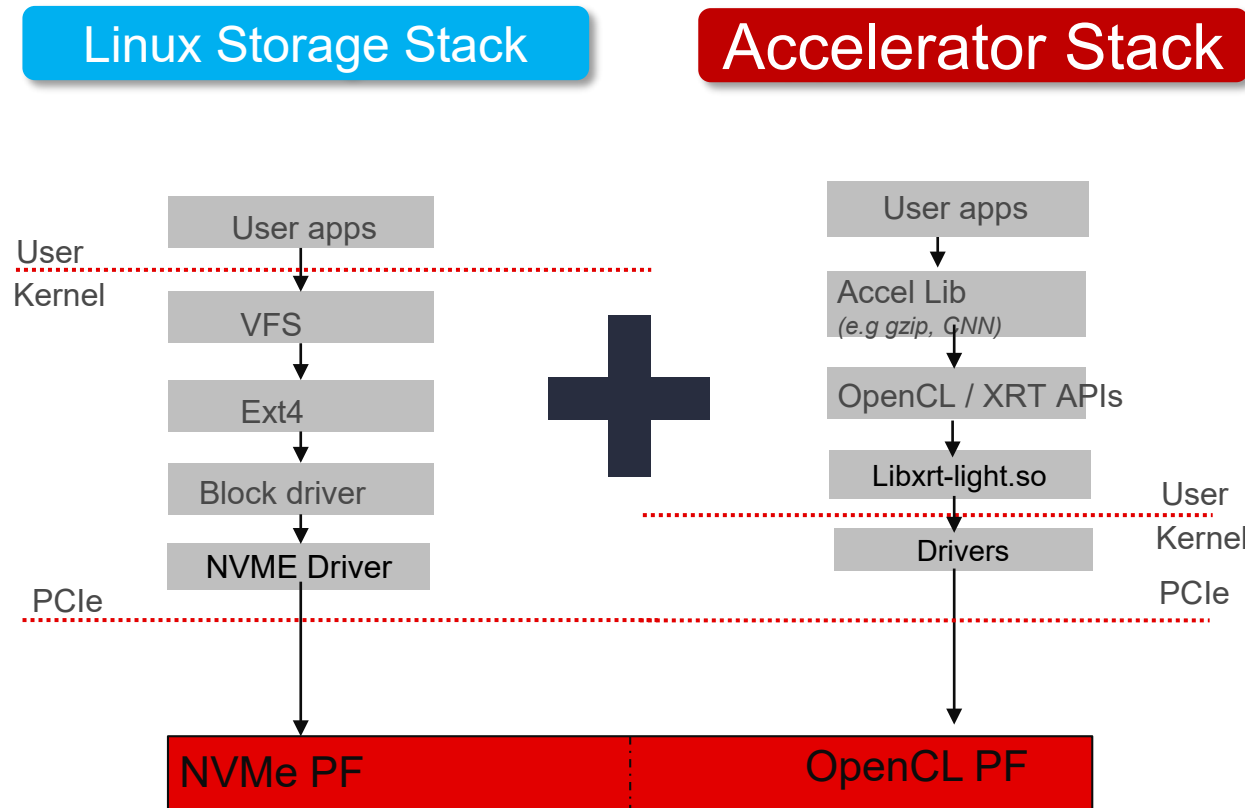
A data processing command is sent

The data is locally processed

Only the processed results are returned

Computational Storage API

Runtime Stack



- › Storage Accessed via NVMe Stack
- › Computational Storage / Accelerator Discovered, Managed, Orchestrated by XRT Stack
- › Shared Memory Space in the Compute Function Glues the Datapaths together

P2P example (Open CL)

<https://xilinx.github.io/XRT/master/html/p2p.html#p2p-data-transfer-between-fpga-card-and-nvme-device>

OpenCL coding style

Typical coding style

1. Create P2P buffer
2. Map P2P buffer to the host space
3. Access the SSD location through Linux File System, the file needs to be opened with *O_DIRECT*.
4. Read/Write through Linux *pread/pwrite* function

```
// Creating P2P buffer
```

```
cl_mem_ext_ptr_t p2pBOExt = {0};  
p2pBOExt.flags = XCL_MEM_EXT_P2P_BUFFER;  
p2pBO = clCreateBuffer(context, CL_MEM_READ_ONLY | CL_MEM_EXT_PTR_XILINX, chunk_size, &p2pBOExt, NULL);  
clSetKernelArg(kernel, 0, sizeof(cl_mem), p2pBO);
```

```
// Map P2P Buffer into the host space
```

```
p2pPtr = (char *) clEnqueueMapBuffer(command_queue, p2pBO, CL_TRUE, CL_MAP_WRITE | CL_MAP_READ, 0, chunk_size, 0, NULL, NULL, NULL);  
filename = <full path to SSD>  
fd = open(filename, O_RDWR | O_DIRECT);
```

```
// Read chunk_size bytes starting at offset 0 from fd into p2pPtr
```

```
pread(fd, p2pPtr, chunk_size, 0);
```

```
// Write chunk_size bytes starting at offset 0 from p2pPtr into fd
```

```
pwrite(fd, p2pPtr, chunk_size, 0)
```



Xilinx Storage Services Reference Design

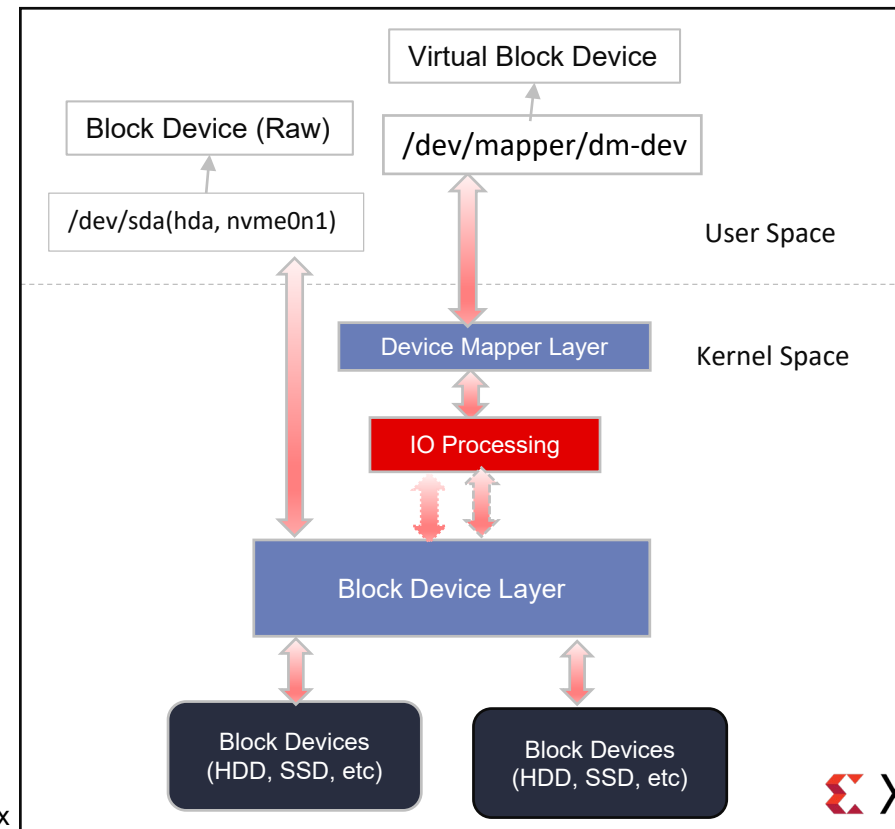
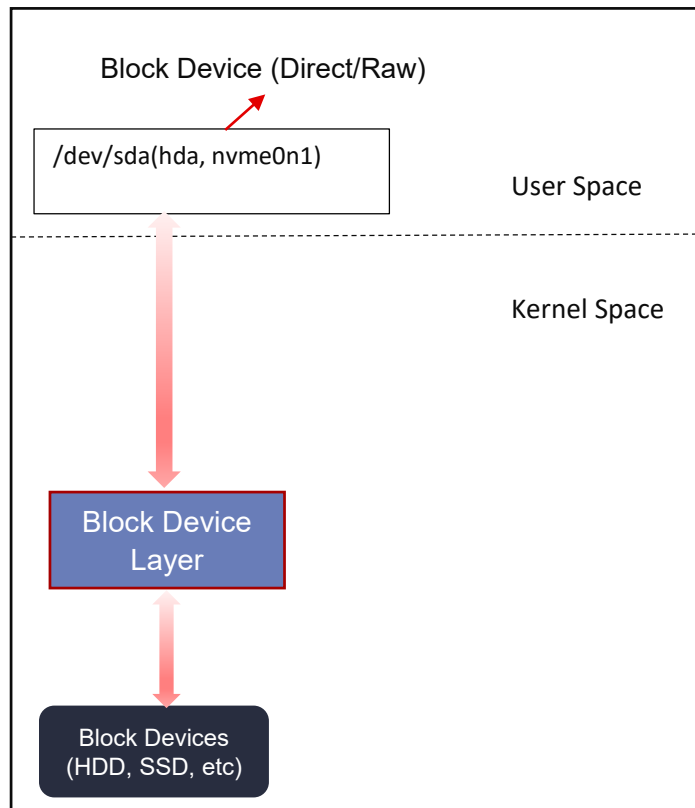
Xilinx Storage Services – Introduction

- The standard API is unaware of accelerator and SSD colocation
 - Requires application to map NVMe and Accelerators to one another.
- Many Linux Software tools run in kernel space and need to call kernel libraries
- Memory allocation overhead and FPGA program times create challenges for storage block level applications

Xilinx Storage Services provide an easy-to-use API solution to accelerate storage kernel applications

XSS Example: Accelerating device mapper

- The device mapper is a framework provided by the Linux kernel for mapping physical block devices onto higher-level virtual block devices.
- Device mapper works by passing data from a virtual block device, which is provided by the device mapper itself, to another block device. Data can be also modified/processed in transition, which is performed, for example, in the case of device mapper providing disk encryption.



XSS API

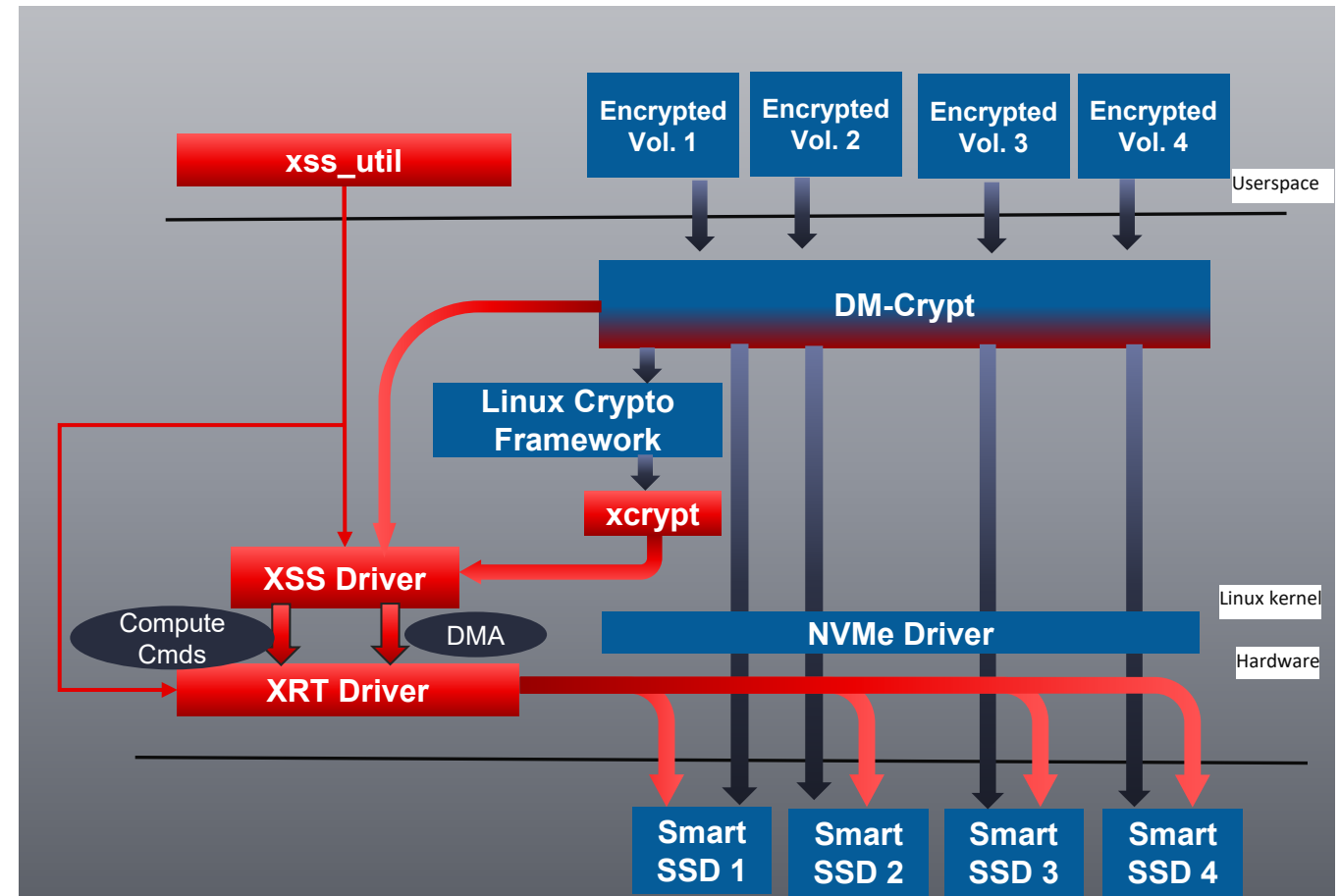
xss.h

- ▶ Provides an API to kernel applications to leverage XSS
- ▶ Can be Extended to create new kernel storage applications

```
86 //The functions below enable allocation of p2p memory
87 struct page* xss_alloc_p2p_page(int ctx_id, gfp_t gfp_mask);
88 void xss_free_p2p_page(int ctx_id, struct page *page);
89
90 //stand alone api to just use HW aes as library to enc/dec and return results to host/ p2p buffer.
91 // for p2p - leverage the p2p allocated handles
92 int xss_crypt(int ctx_id, struct xss_cipher_req *req);
93
94 //Lists the devices that support a particular service (assumes that the xclbin is
95 // programmed on the device with the service)
96 int xss_get_devs_with_service(int service_id, int dev_id[]);
97 bool xss_is_service_avail_on_dev(int dev_id, int service_id);
98 int xss_get_pci_addr(int dev_id, unsigned *pci_addr);
99 int xss_get_dev_id(unsigned pci_addr);
100 //For the smartSSD give the accelerator device id associated with the NVMe device.
101 int xss_get_peer_dev(const char * blkdev_path);
102 //Creates a context on the accelerator device for an application using XSS.
103 // Client is a user name for the context that can be referenced in utilities
104 // dev_id - accelerator device ID
105 // Services - list from the supported services
106 // num_services - length of the services array.
107 // ret - handle to the context.
108 int xss_create_ctx(const char * client, int dev_id, int services[], int num_services);
109 void xss_delete_ctx(int ctx_id);
110 bool xss_is_user_ptr_p2p(struct page *page, int ctx_id);
```

Acceleration Example: XSS/dm-crypt

- DM-Crypt is an existing Linux solution that
 - Provides inline full disk encryption (FDE)
 - leverages Linux device-mapper layer and Linux Crypto Framework to accomplish FDE.
- XSS Driver
 - Provides easy to use APIs for kernel applications matching kernel storage library interfaces.
 - Creates accelerator contexts for applications and manages P2P buffers.
 - Manages complexity from typical in mult-accelerator SmartSSD deployments.
 - Contains services Xilinx developed.
- `xss_util`
 - Programs the `xclbin` and passes configuration information to the XSS Driver.
- `xcrypt`
 - Linux crypto framework module that leverages a HW accelerator on the SmartSSD using XSS.



- DM-Crypt
 - Updated to use P2P Buffers for IO using XSS.

Xilinx Storage Services : XSS/dm-crypt Example

➤ XSS Configuration:

1. Update /etc/xss.conf

```
#<bdf_address> <xclbin_path>  
0000:05:00.1 /home/xss/fa_aes_xts2_rtl_enc_dec.xclbin
```

2. Load xss configuration on device

```
# xss_util load-config  
~~~~~  
Reading configuration from /etc/xss.conf:  
~~~~~  
Configuring device 0000:05:00.1... Loaded xclbin:/home/xss/fa_aes_xts2_rtl_enc_dec.xclbin, Added device info to XSS.
```

➤ *Dm-crypt VBD Creation*

1. Using dmsetup

```
# dmsetup create xss-dev --table "0 $(blockdev --getsz /dev/nvme0n1) crypt capi:xss_aes_xts_async-plain64 <128-byte-key> 0 /dev/nvme0n1 0 1  
sector_size:4096"
```

2. Using Cryptsetup

```
# cryptsetup --type luks2 --cipher capi:xss_aes_xts_async-plain64 --key-size 512 --sector-size 4096 luksFormat /dev/nvme0n1  
  
# cryptsetup luksOpen /dev/nvme0n1 xss-dev
```

➤ *Use VBD with a file-system*

```
# mkfs.ext4 /dev/mapper/xss-dev  
  
# mount /dev/mapper/xss-dev /mnt
```

Kernel application modification

- ▶ Dm-crypt – Updated to allocate buffers from P2P region of SmartSSD CSD.

```
2125 static void *crypt_page_alloc(gfp_t gfp_mask, void *pool_data)
2126 {
2127     struct crypt_config *cc = pool_data;
2128     struct page *page;
2129     struct crypto_tfm *base;
2130     struct xss_aes_ctx *xctx;
2131
2132     if (unlikely(percpu_counter_compare(&cc->n_allocated_pages, dm_crypt_pages_per_client) >= 0) &&
2133         likely(gfp_mask & __GFP_NORETRY))
2134         return NULL;
2135
2136     if (!test_bit(DM_CRYPT_USE_P2P_PAGES, &cc->flags))
2137         page = alloc_page(gfp_mask);
2138     else {
2139         base = crypto_skcipher_tfm(any_tfm(cc));
2140         xctx = crypto_tfm_ctx(base);
2141         page = xss_alloc_p2p_page(xctx->ctx_id, gfp_mask);
2142     }
2143
2144     if (likely(page != NULL))
2145         percpu_counter_add(&cc->n_allocated_pages, 1);
2146
2147     return page;
2148 }
2149
2150 static void crypt_page_free(void *page, void *pool_data)
2151 {
2152     struct crypt_config *cc = pool_data;
2153     struct crypto_tfm *base = crypto_skcipher_tfm(any_tfm(cc));
2154     struct xss_aes_ctx *xctx = crypto_tfm_ctx(base);
2155
2156     if (!test_bit(DM_CRYPT_USE_P2P_PAGES, &cc->flags))
2157         __free_page(page);
2158     else {
2159         xss_free_p2p_page(xctx->ctx_id, page);
2160     }
2161     percpu_counter_sub(&cc->n_allocated_pages, 1);
2162 }
```

Modified to have the bio buffers allocated in the P2P address space. Once that is done the reads/ writes to NVMe device happen via P2P – automatically!

Life of a Block: Write

▶ Application Writes to storage:

```
filename = <full path dm-crypt block device>
fd = open(filename, O_RDWR);
// Write chunk_size bytes starting at offset 0 from p2pPtr into fd
pwrite(fd, data_pointer, size, offset)
```

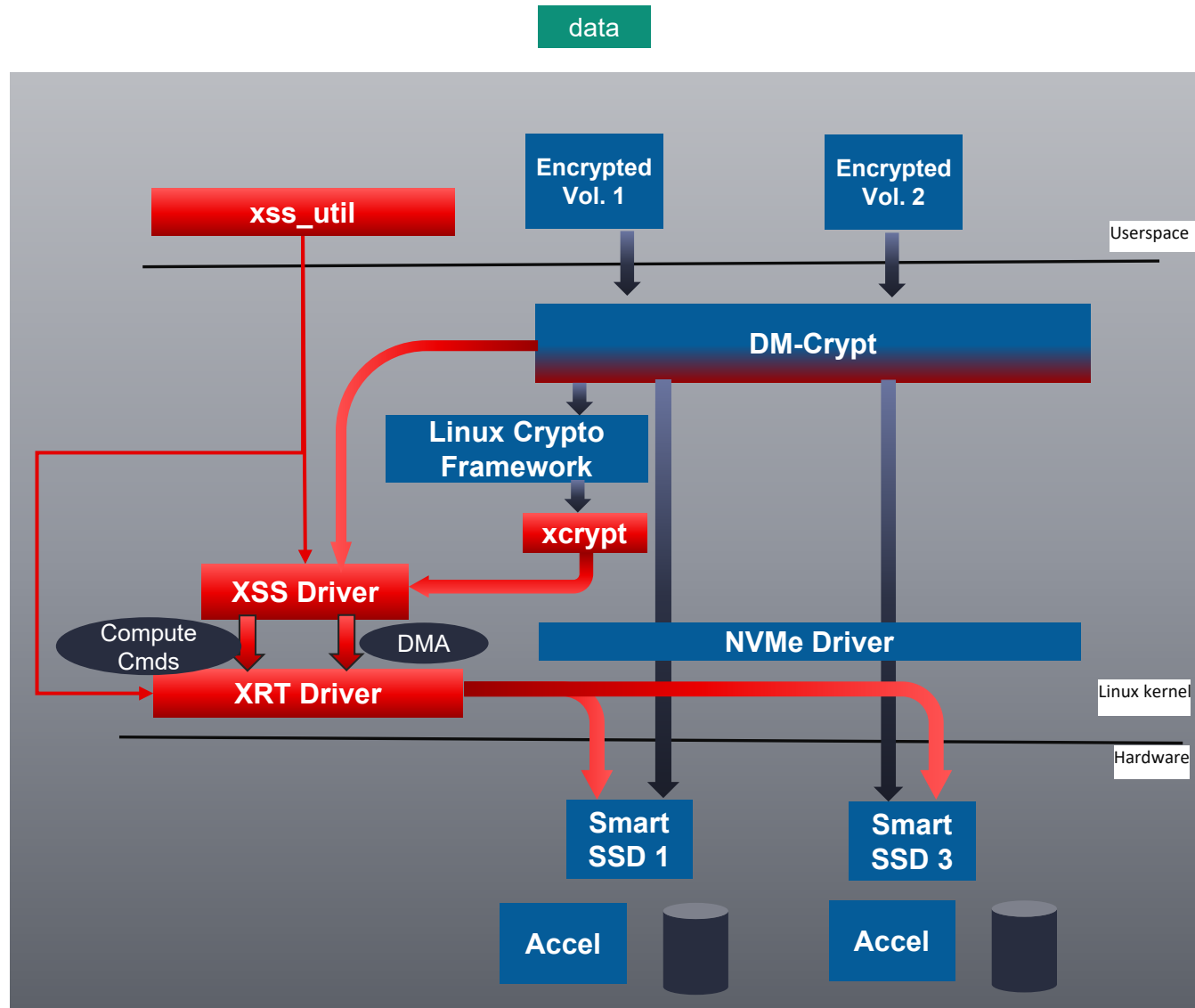
- ▶ BIO Layer in Linux forwards the IO to DM-crypt.
- ▶ DM-Crypt creates a crypto work item to encrypt the write data
- ▶ In dm-crypt-5.4.c:

```
1066 static int crypt_convert_block_skcipher(struct crypt_config *cc,
1067                                         struct convert_context *ctx,
1068                                         struct skcipher_request *req,
1069                                         unsigned int tag_offset)
1070 {
1071     struct bio_vec bv_in = bio_iter_iovec(ctx->bio_in, ctx->iter_in);
1072     struct bio_vec bv_out = bio_iter_iovec(ctx->bio_out, ctx->iter_out);
1073     struct scatterlist *sg_in, *sg_out;
```

```
1126     if (bio_data_dir(ctx->bio_in) == WRITE)
1127         r = crypto_skcipher_encrypt(req);
1128     else
1129         r = crypto_skcipher_decrypt(req);
```

} No modifications – calls the HW accelerated crypto based on the dm-crypt VBD specifying that module under the linux crypto framework

Write to Linux Crypt



Life of a Block: Write

- ▶ Xcrypt.c – accelerated linux crypto framework module. Provides all of the linux crypto framework interfaces. Actual hw accelerated work done by xss_crypt() in XSS.

```
▼ S xcrypt_requeue_request
  ○ list : struct list_head
  ○ req : struct skcipher_request*
  ○ encrypt : bool
  ○ xss_work : struct work_struct
  ⚡ S xcrypt_handle_requeue_reqs(struct work_struct*) : void
  ⚡ S xcrypt_enqueue_request(int, struct skcipher_request*) : void
  ● S xss_alg_skcipher_init(struct crypto_skcipher*) : int
  ● S xss_alg_skcipher_exit(struct crypto_skcipher*) : void
  ● S xss_alg_sync_skcipher_encrypt(struct skcipher_request*) : int
  ● S xss_alg_sync_skcipher_decrypt(struct skcipher_request*) : int
  ● S xss_alg_sync_skcipher_xts_setkey(struct crypto_skcipher*, const u8*, unsigned int) : int
  ● S callback(void*, int) : void
  ● xcrypt_enqueue_request(int, struct skcipher_request*) : void
  ● S xss_alg_async_skcipher_encrypt(struct skcipher_request*) : int
  ● S xss_alg_async_skcipher_decrypt(struct skcipher_request*) : int
  ● S xcrypt_handle_requeue_reqs(struct work_struct*) : void
  ● S xss_alg_async_skcipher_xts_setkey(struct crypto_skcipher*, const u8*, unsigned int) : int
  ● S xss_crypto_algs : struct skcipher_alg[]
  ● xss_crypt_init(void) : int
  ● xss_crypt_exit(void) : void
  ⚡ module_init()
  ⚡ module_exit()
```

```
144 static int xss_alg_async_skcipher_encrypt(struct skcipher_request *req)
145 {
146     int ret;
147     struct crypto_skcipher *cipher = crypto_skcipher_reqtfm(req);
148     struct crypto_tfm *tfm = crypto_skcipher_tfm(cipher);
149     struct xss_aes_ctx *xctx = crypto_tfm_ctx(tfm);
150
151     struct xss_cipher_req xss_req;
152     xss_req.src = req->src;
153     xss_req.dst = req->dst;
154     xss_req.encrypt = true;
155     xss_req.iv = (uint8_t *)req->iv;
156     xss_req.datalen = req->cryptlen;
157     xss_req.datatype = SGL_TYPE;
158
159     xss_req.xss_cb.func = (void *)callback;
160     xss_req.xss_cb.data = (void *)&req->base;
161
162     xss_req.key_len = xctx->key_len;
163     xss_req.key = xctx->key;
164     ret = xss_crypt(xctx->ctx_id, &xss_req);
165     if(ret == -EBUSY) {
166         xcrypt_enqueue_request(true, req);
167     }
168     return ret;
169 }
```

Life of a Block: Write

► xss_crypt:

```
347 int xss_crypt(int ctx_id, struct xss_cipher_req *cipher_req)
348 {
349     struct xss_context *ctx = xss_get_ctx_struct(ctx_id);
350     int ret=0;
351     int cu_index;
352     int svc_offset;
353     struct xss_request *xss_req;
354     struct xss_bo **bo_list;
355 }
```

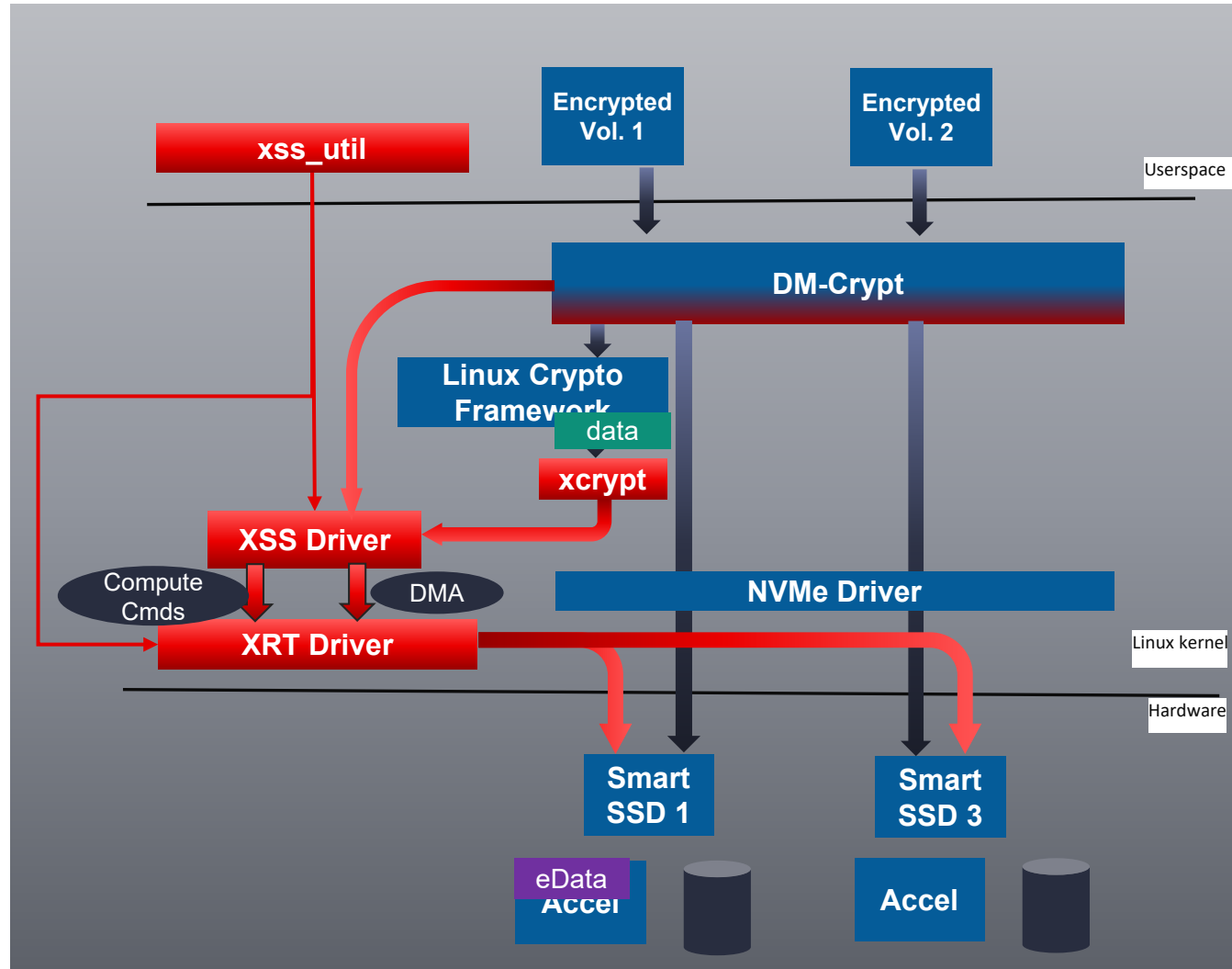
```
363     svc_offset = xss_service_offset_in_ctx(ctx, crypt_services[cipher_req->encrypt]);
364     if (svc_offset<0)
365         return svc_offset;
366
367     cu_index = xss_get_service_cu(ctx, svc_offset, crypt_kernels[cipher_req->encrypt]);
```

```
387     ret = crypt_setup_ert_packet(ctx, cipher_req, bo_list, cu_index);
388     if (ret)
389         goto out;
390
391     xss_req->ctx = ctx;
392     xss_req->svc_id = crypt_services[cipher_req->encrypt];
393     xss_req->cu_index = cu_index;
394     xss_req->bo_list = bo_list;
395     xss_req->num_bos = CRYPT_NUM_BOS;
396     xss_req->op_dir = cipher_req->encrypt ? XSS_OP_DIR_WRITE:XSS_OP_DIR_READ;
397     xss_req->complete = xss_crypt_req_cmplt;
398     xss_req->user_cb = cipher_req->xss_cb.func;
399     xss_req->user_cb_data = cipher_req->xss_cb.data;
400
401     if ((ret = xss_submit_request(xss_req)) < 0)
402         goto out;
```

Get the Compute unit from XSS

Submit the encryption request to XSS with the P2P buffers.

Write to Linux Crypt



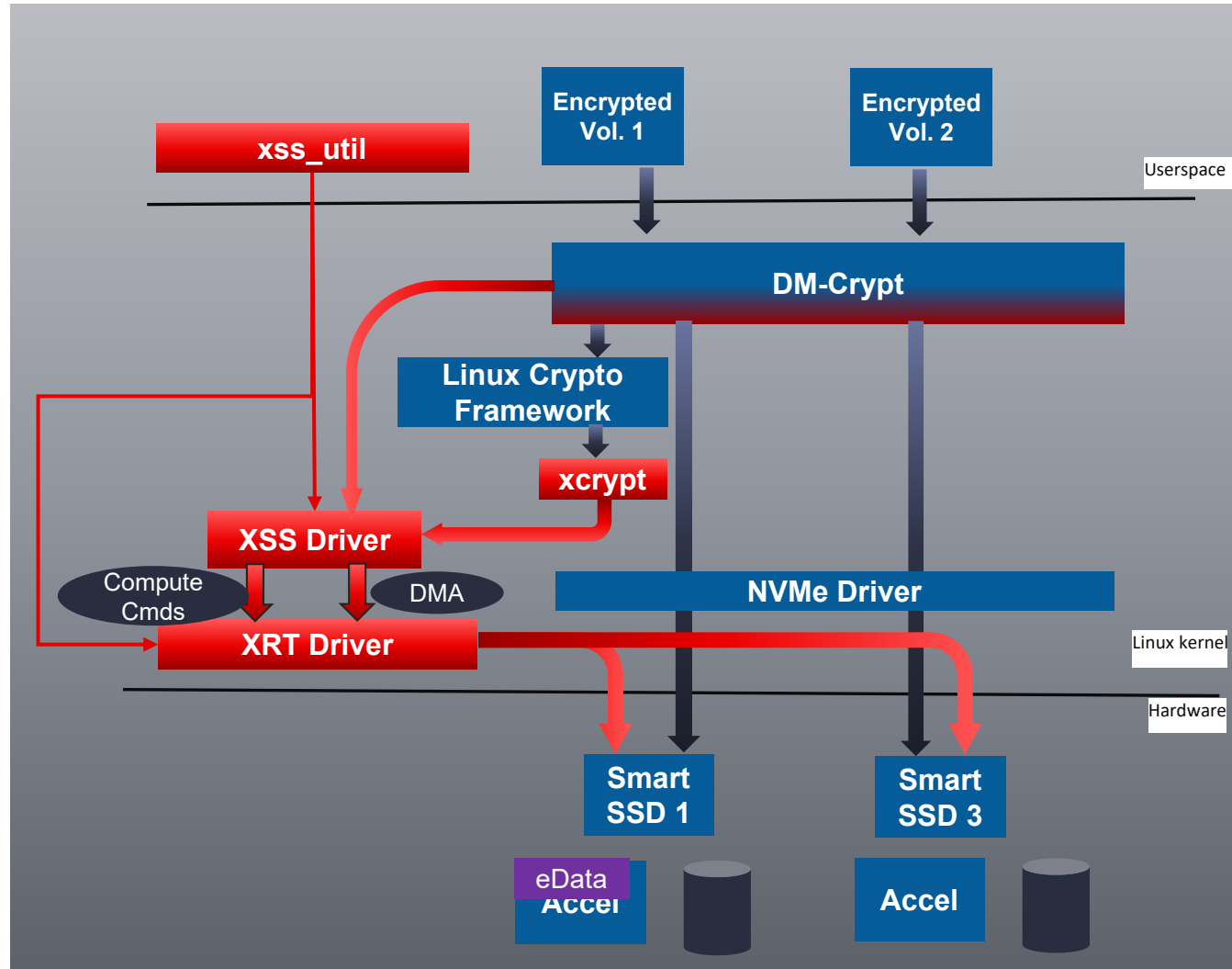
Dm-crypt

- ▶ Encrypted data written by Dm-crypt to backing storage via the BIO layer
- ▶ BIO Layer forwards the IO to NVMe device.

```
1536 static int dmccrypt_write(void *data)
1537 {
1538     struct crypt_config *cc = data;
1539     struct dm_crypt_io *io;
1540
1541     while (1) {
1542         struct rb_root write_tree;
1543         struct blk_plug plug;
1544
1545         spin_lock_irq(&cc->write_thread_lock);
1546     continue_locked:
1547
1548         if (!RB_EMPTY_ROOT(&cc->write_tree))
1549             goto pop_from_list;
1550
1551         set_current_state(TASK_INTERRUPTIBLE);
1552
1553         spin_unlock_irq(&cc->write_thread_lock);
1554
1555         if (unlikely(kthread_should_stop())) {
1556             set_current_state(TASK_RUNNING);
1557             break;
1558     }
```

No modifications – The encrypted data input is already in a P2P buffer, so that when the write goes to the NVMe DMA it is pulled from the buffer and written to disk automatically.

Write to Linux Crypt



Next steps

- ▶ Build accelerated storage applications!
- ▶ Register for access to the Xilinx Storage Services reference design:

<https://www.xilinx.com/products/intellectual-property/xss.html>



Thank You

