

# Persistent Memory File System

## **PMFS(Persistent Memory File System)**

**LIU CHANGPENG && YI DONG**

INTEL CORPORATION | 3/3/2022

## **PMFS(Persistent Memory File System):**

*Single node userspace file system based on DAOS  
VOS and SPDK and relies on PMem.*

## **PMFS Application Scenario:**

**It's targeted scenarios including:**

- *Integrated with upper level application like Database as the optimal File Store against kernel based POSIX FS*
- *Support different kinds of backend block stores by SPDK like NVMe-oF target, Ceph backends and self developed backends*

# PMFS DETAILS

## Code files Introduction:

### Changes to be committed:

- # modified: src/SConscript
  - add pmfs to build.
- # new file: src/pmfs/README.md
  - file for introduction for pmfs mechanism.
- # new file: src/pmfs/SConscript
  - build file for compiling all files in pmfs.
- # new file: src/pmfs/include/pmfs.h
  - pmfs functions and some defines that can be used for external users.
- # new file: src/pmfs/pmfs.c
  - The realizations of pmfs APIs.
- # new file: src/pmfs/pmfs\_internal.h
  - Internal APIs that be called by pmfs APIs.
- # new file: src/pmfs/vos\_client.c
  - The realizations for client side  
fetch/update/punch/get\_num\_dkeys/list\_dkeys.
- # new file: src/pmfs/vos\_target\_engine.c
  - The realization of pools and containers and objects  
initialization and binding mechanism.
- # new file: src/pmfs/vos\_target\_engine.h
  - Target engine header file for redefine the linked list persistent  
memory pool.
- # new file: src/pmfs/vos\_target\_fs.c
  - Target file for processing tasks in spdk ring queue.
  - Drains the task queues and executes the tasks in callbacks.
  - It provides APIs that test APPs or external users can call:
    - ✓ vos\_task\_process\_init
    - ✓ vos\_task\_process
    - ✓ vos\_task\_process\_fini
    - ✓
- # new file: src/pmfs/vos\_target\_fs.h

➤ Header file for construct vos\_fs\_cmd\_args and APIs output.

# new file: src/pmfs/vos\_tasks.c

➤ Task using spdk\_ring\_create and vos\_target\_free\_tasks.

# new file: src/pmfs/vos\_tasks.h

➤ Task related structs defines :

- ✓ VOS\_OBJ\_UPDATE
- ✓ VOS\_OBJ\_FETCH,
- ✓ VOS\_OBJ\_PUNCH,
- ✓ VOS\_OBJ\_GET\_NUM\_DKEYS,
- ✓ VOS\_OBJ\_LIST\_DKEYS,

# modified: src/tests/SConscript

➤ Add tests APPs to build.

# new file: src/tests/pmfs\_demo.c

➤ Demo app for demonstrating mkfs/mount/createdir/...

# new file: src/tests/pmfs\_unittest.c

➤ Unittest that can be added to CI.

# new file: src/include/pmfs/pmfs.h

➤ PMFS functions and some defines that can be used for external users.

# new file: src/include/pmfs/vos\_target\_engine.h

➤ Copy header file for test apps.

# new file: src/include/pmfs/vos\_target\_fs.h

➤ Copy header file for test apps.

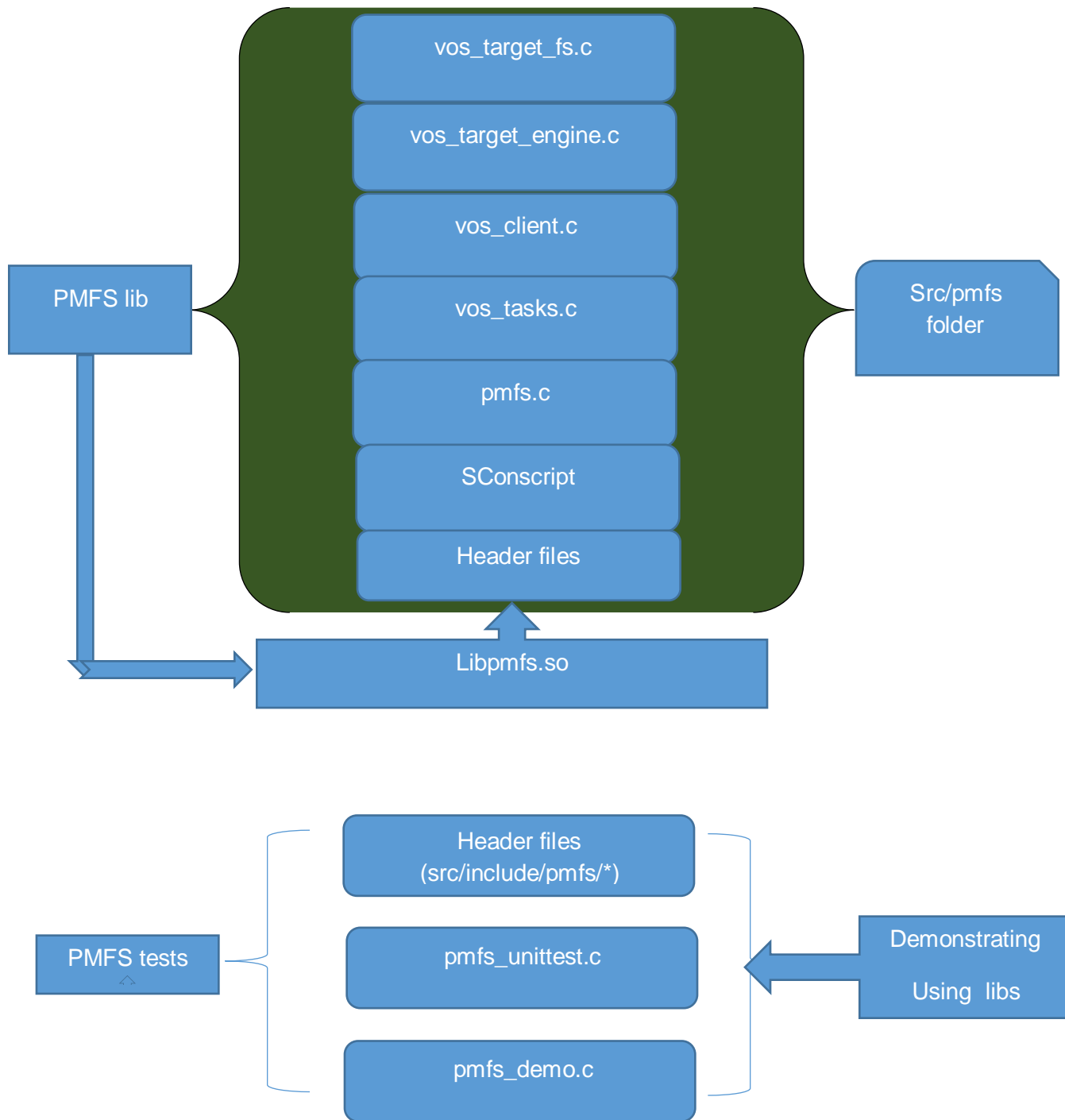
# new file: src/include/pmfs/vos\_tasks.h

➤ Copy header file for test apps.

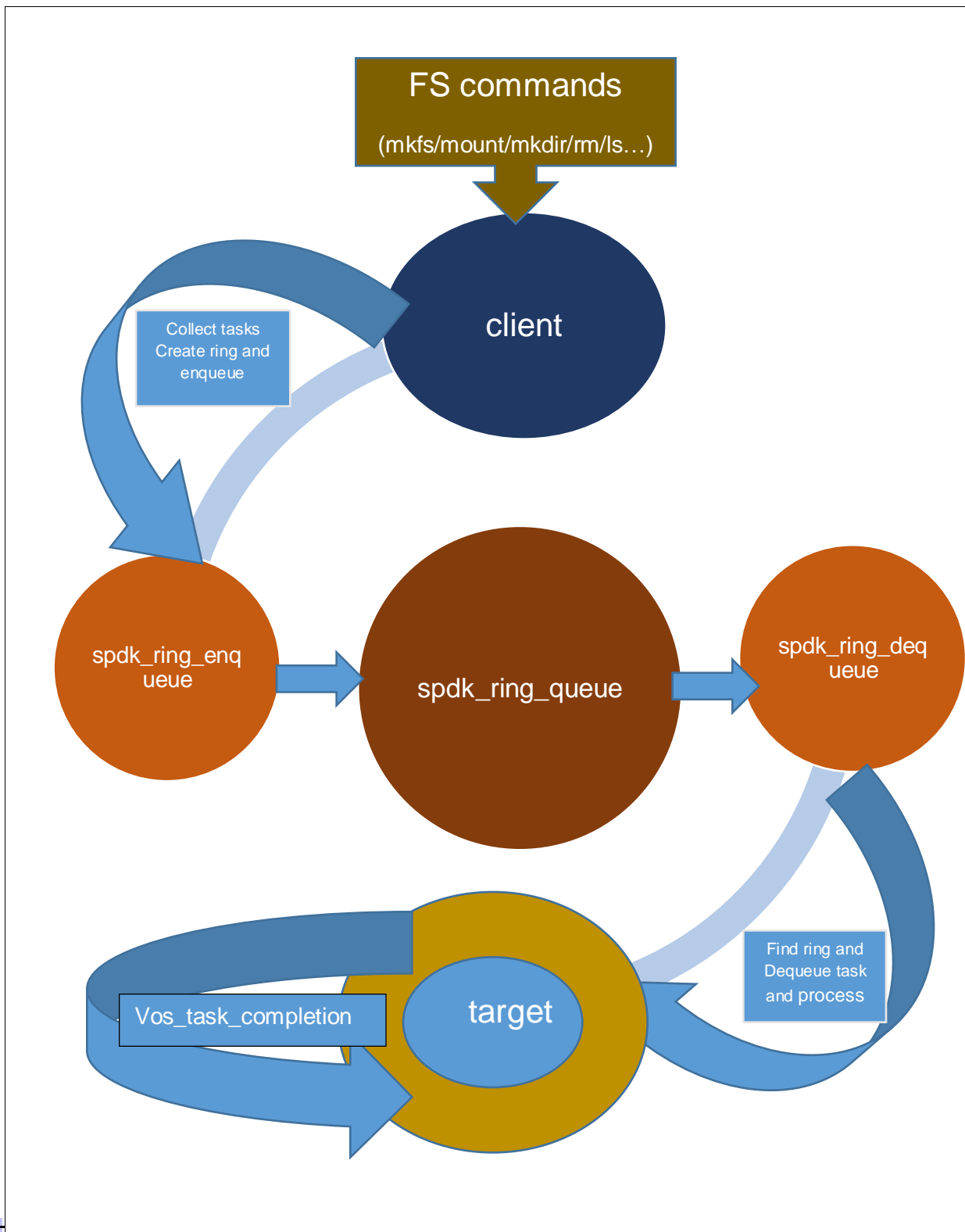
# modified: utils/run\_test.sh

➤ Add pmfs\_unittest to CI script.

## Files relationships

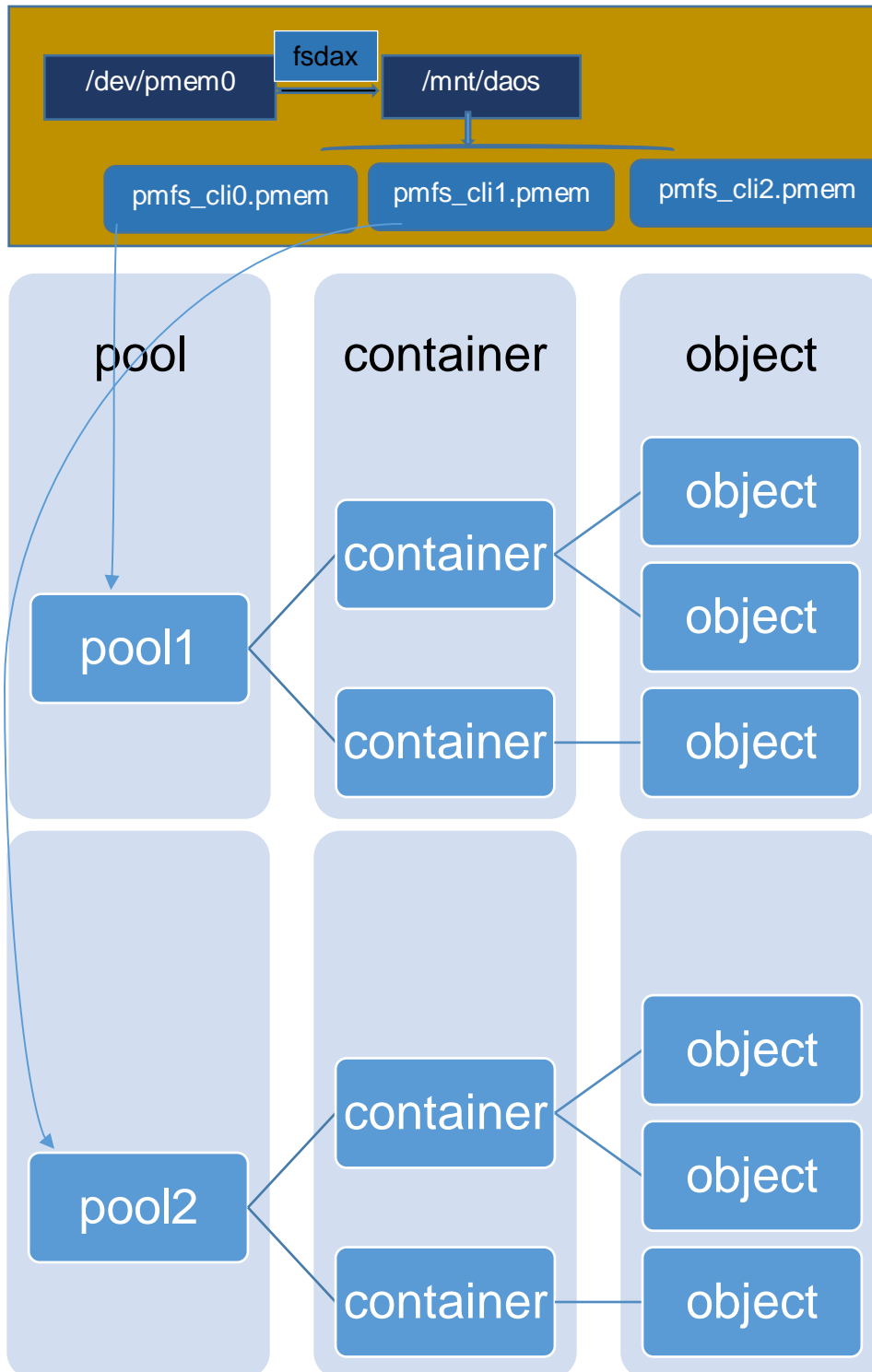


## PMFS threading model



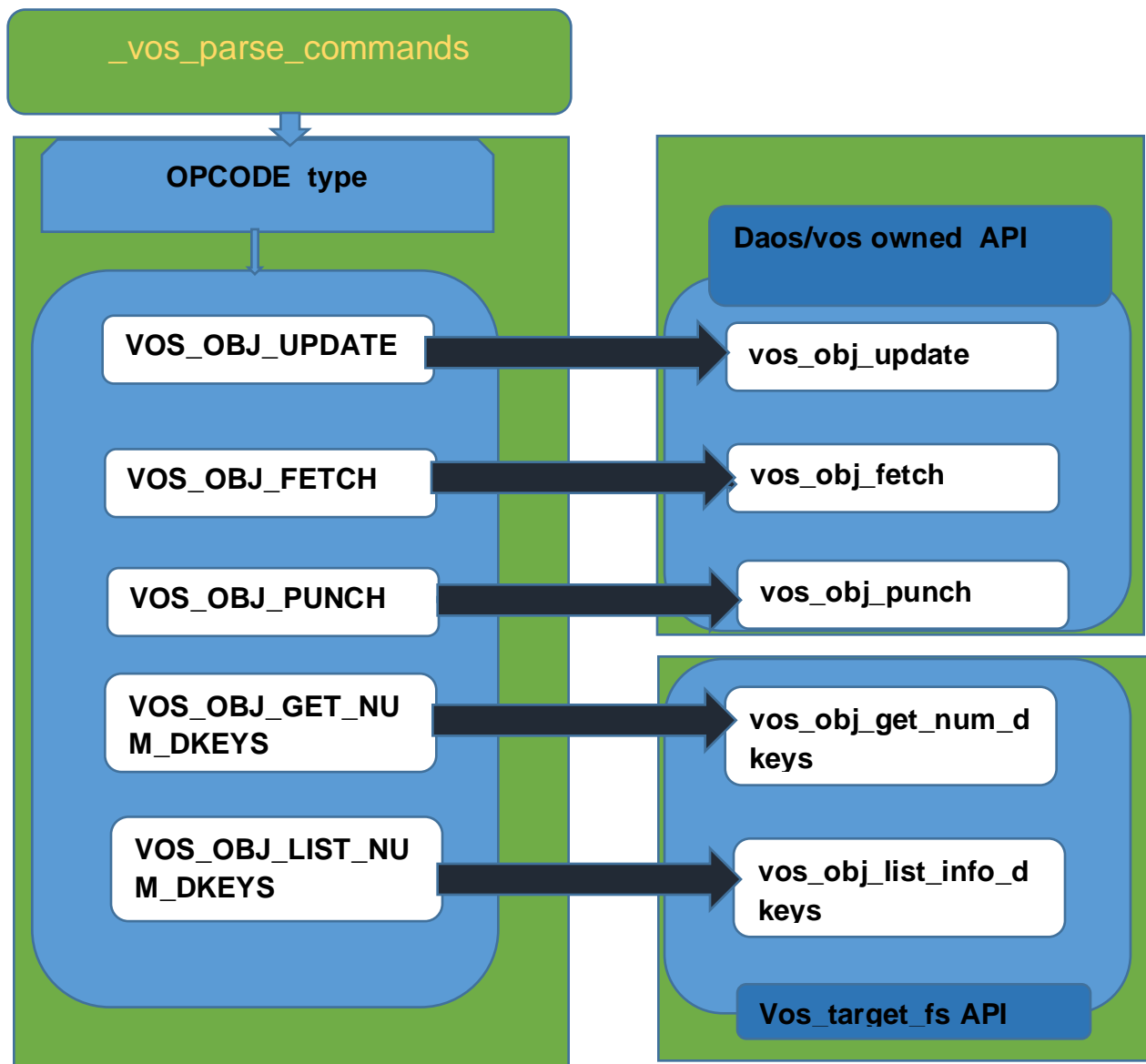
# Storage hierarchy

- Persistent memory pool
- Container
- Objects



## + VOS target processes FS command calls

- vos\_task\_process\_init
- vos\_task\_process
  - ✓ while(1){ vos\_task\_completion}
    - vos\_task\_dequeue
    - vos\_task\_ult
    - ❖ \_vos\_parse\_commands
- vos\_task\_process\_fini





## PMFS vos APIs

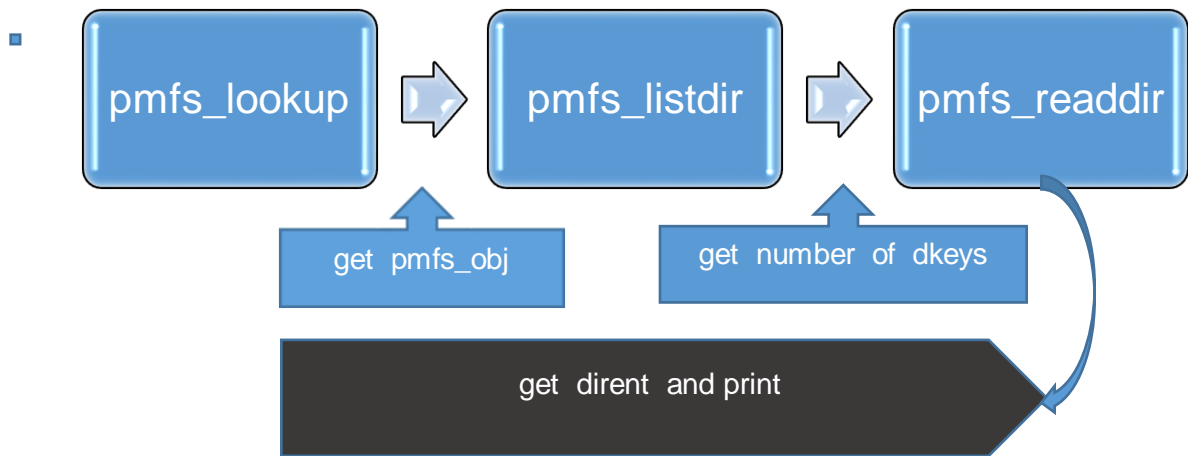
- *int pmfs\_mkfs(daos\_handle\_t poh, uuid\_t uuid);*
- *int pmfs\_mount(daos\_handle\_t poh, daos\_handle\_t coh, int flags, pmfs\_t \*\*pmfs);*
- *int pmfs\_umount(pmfs\_t \*pmfs);*
- *int pmfs\_mkdir(pmfs\_t \*pmfs, pmfs\_obj\_t \*parent, const char \*name, mode\_t mode);*
- *int pmfs\_listdir(pmfs\_t \*pmfs, pmfs\_obj\_t \*obj, uint32\_t \*nr);*
- *int pmfs\_remove(pmfs\_t \*pmfs, pmfs\_obj\_t \*parent, const char \*name, bool force, daos\_obj\_id\_t \*oid);*
- *int pmfs\_open(pmfs\_t \*pmfs, pmfs\_obj\_t \*parent, const char \*name, mode\_t mode, int flags, daos\_size\_t chunk\_size, const char \*value, pmfs\_obj\_t \*\*\_obj);*
- *int pmfs\_readdir(pmfs\_t \*pmfs, pmfs\_obj\_t \*obj, uint32\_t \*nr, struct dirent \*dirs);*
- *int pmfs\_release(pmfs\_obj\_t \*obj);*
- *int pmfs\_lookup(pmfs\_t \*pmfs, const char \*path, int flags, pmfs\_obj\_t \*\*\_obj, mode\_t \*mode, struct stat \*stbuf);*
- *int pmfs\_punch(pmfs\_t \*pmfs, pmfs\_obj\_t \*obj, daos\_off\_t offset, daos\_size\_t len);*
- *int pmfs\_write\_sync(pmfs\_t \*pmfs, pmfs\_obj\_t \*obj, d\_sg\_list\_t \*usr\_sgl, daos\_off\_t off);*
- *int pmfs\_read\_sync(pmfs\_t \*pmfs, pmfs\_obj\_t \*obj, d\_sg\_list\_t \*usr\_sgl, daos\_off\_t off, daos\_size\_t \*read\_size);*
- *int pmfs\_stat(pmfs\_t \*pmfs, pmfs\_obj\_t \*parent, const char \*name, struct stat \*stbuf);*

## VOS client APIs

- *vos\_client\_obj\_update*
- *vos\_client\_obj\_update\_sync*
- *vos\_client\_obj\_fetch*
- *vos\_client\_obj\_fetch\_sync*
- *vos\_client\_obj\_punch*
- *vos\_client\_obj\_punch\_sync*
- *vos\_client\_obj\_get\_num\_dkys*
- *vos\_client\_obj\_get\_num\_dkys\_sync*
- *vos\_client\_obj\_list\_dkys*
- *vos\_client\_obj\_list\_dkys\_sync*

## + Example

### ***“ls” command:***



### ❖ Code flow path:

#### ➤ Client :

`pmfs_lookup->fetch_entry->vos_client_obj_fetch_sync->client_task_enqueue->spdk_ring_enqueue`

#### ➤ Target :

`vos_task_completion->vos_task_dequeue->vos_task_get_ring->spdk_ring_dequeue->vos_task_ult->vos_fs_execute_command->vos_parse_commands->vos_obj_fetch`

## + Demo APP

### ➤ Setup env script:

```
#!/usr/bin/env bash
umount -l /dev/pmem0
ndctl destroy-namespace namespace0.0 --force
ndctl create-namespace --size 120G --mode fsdax --map mem -
e namespace0.0 -f
ndctl list
mkfs.xfs -f -m reflink=0 /dev/pmem0
mkdir /mnt/daos
mount -o dax /dev/pmem0 /mnt/daos
```

### ➤ config file /etc/daos\_nvme.conf:

```
{
  "daos_data": {
    "config": []
  },
  "subsystems": [
    {
      "subsystem": "bdev",
      "config": [
        {
          "method": "bdev_nvme_set_options",
          "params": {
            "action_on_timeout": "none",
            "timeout_us": 0,
            "retry_count": 4,
            "arbitration_burst": 0,
            "low_priority_weight": 0,
            "medium_priority_weight": 0,
            "high_priority_weight": 0,
            "nvme_adminq_poll_period_us": 10000,
            "keep_alive_timeout_ms": 10000,
```

```
        "nvme_ioq_poll_period_us": 0,  
        "io_queue_requests": 0,  
        "delay_cmd_submit": true  
    }  
},  
{  
    "params": {  
        "name": "Nvme0",  
        "trtype": "PCle",  
        "traddr": "0000:68:00.0"  
    },  
    "method": "bdev_nvme_attach_controller"  
}  
]  
}  
]
```

`./install/bin/pmfs_demo`