

Integrating SPDK in the NAS gateway

Erlang Li

Senior Storage Software Engineer

Integrated Storage Division

XSKY



Outline

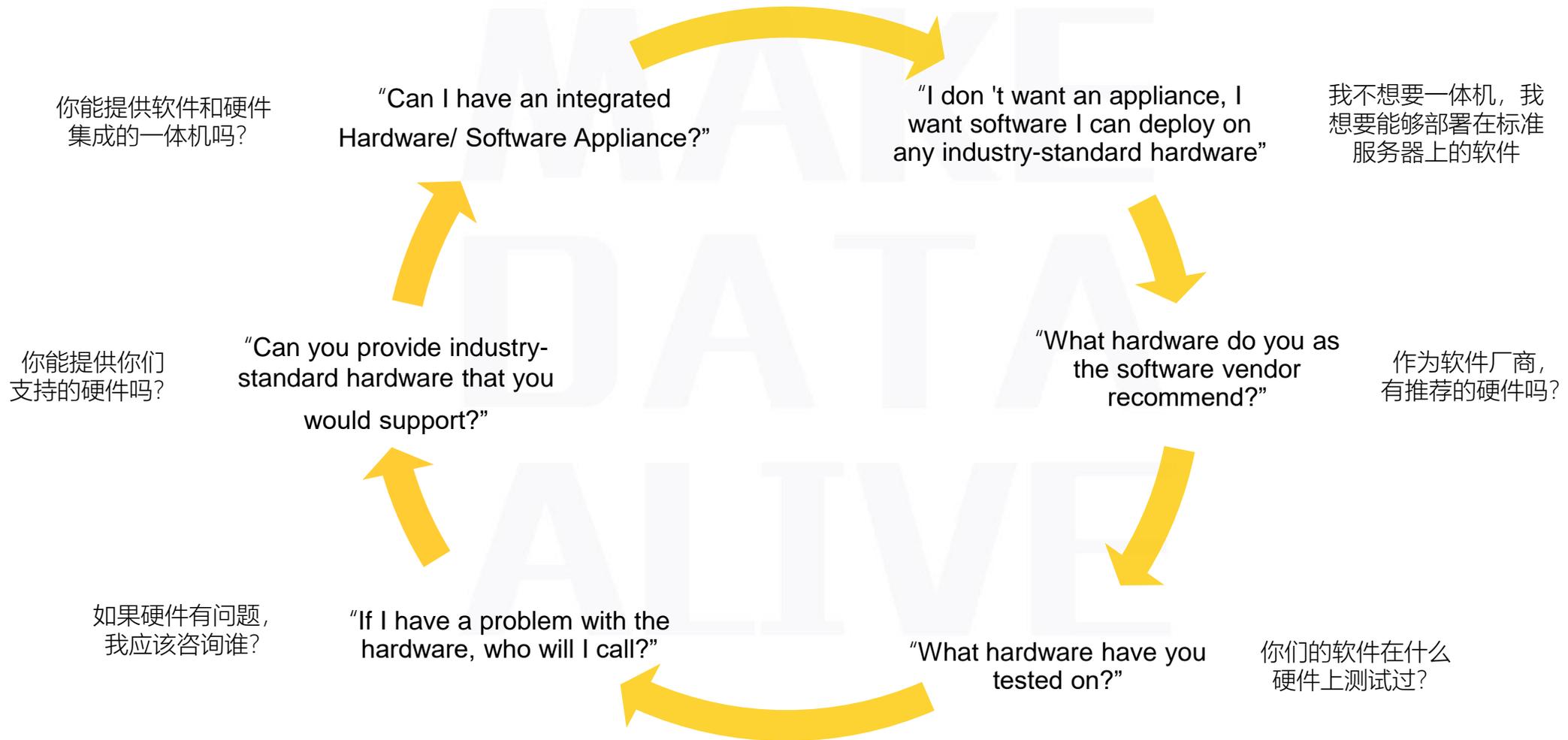
- Background
- Old Design of the NAS Gateway
- New Design of the NAS Gateway with SPDK
 - iSCSI based on SPDK
 - NVMe-of based on SPDK
- Performance comparison
 - Local Kernel VS NVMe-of VS iSCSI
 - Intel Optane VS NAND NVMe
 - Old System VS New System
- Add JSON-RPC Methods for NVMe SSD status

Outline

- Background

MAKE
DATA
ALIVE

The “Appliance” Vs. “Software-only” crazy illogical circle



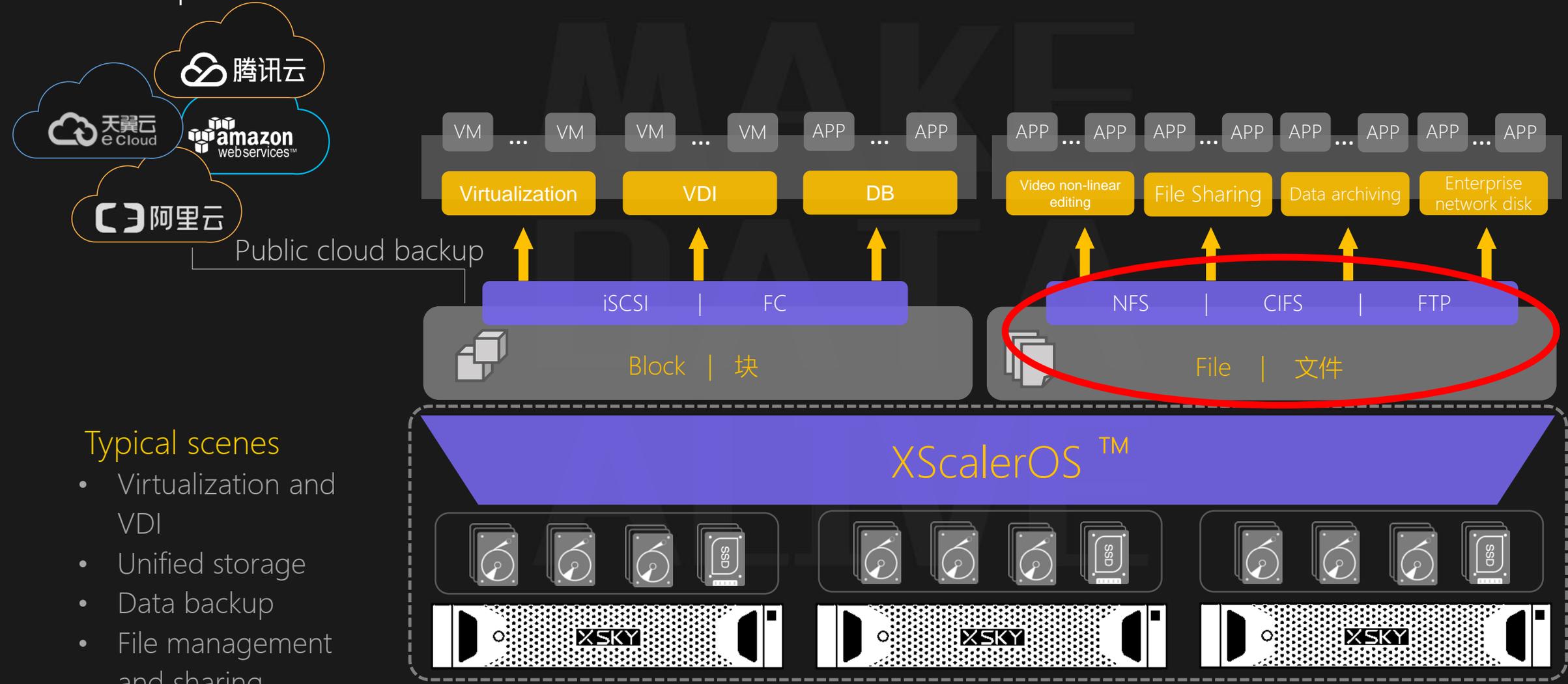
通向敏捷 IT 的理想选择

XSCALER EXPRESS 2000系列 软件定义存储一体机

- ➔ 软件定义存储引领IT变革
- ➔ 无缝对接主流企业级应用
- ➔ 经济易捷助力IT快速转型



XE2000 series focuses on the storage needs of multi-scenarios in growing companies



Typical scenes

- Virtualization and VDI
- Unified storage
- Data backup
- File management and sharing

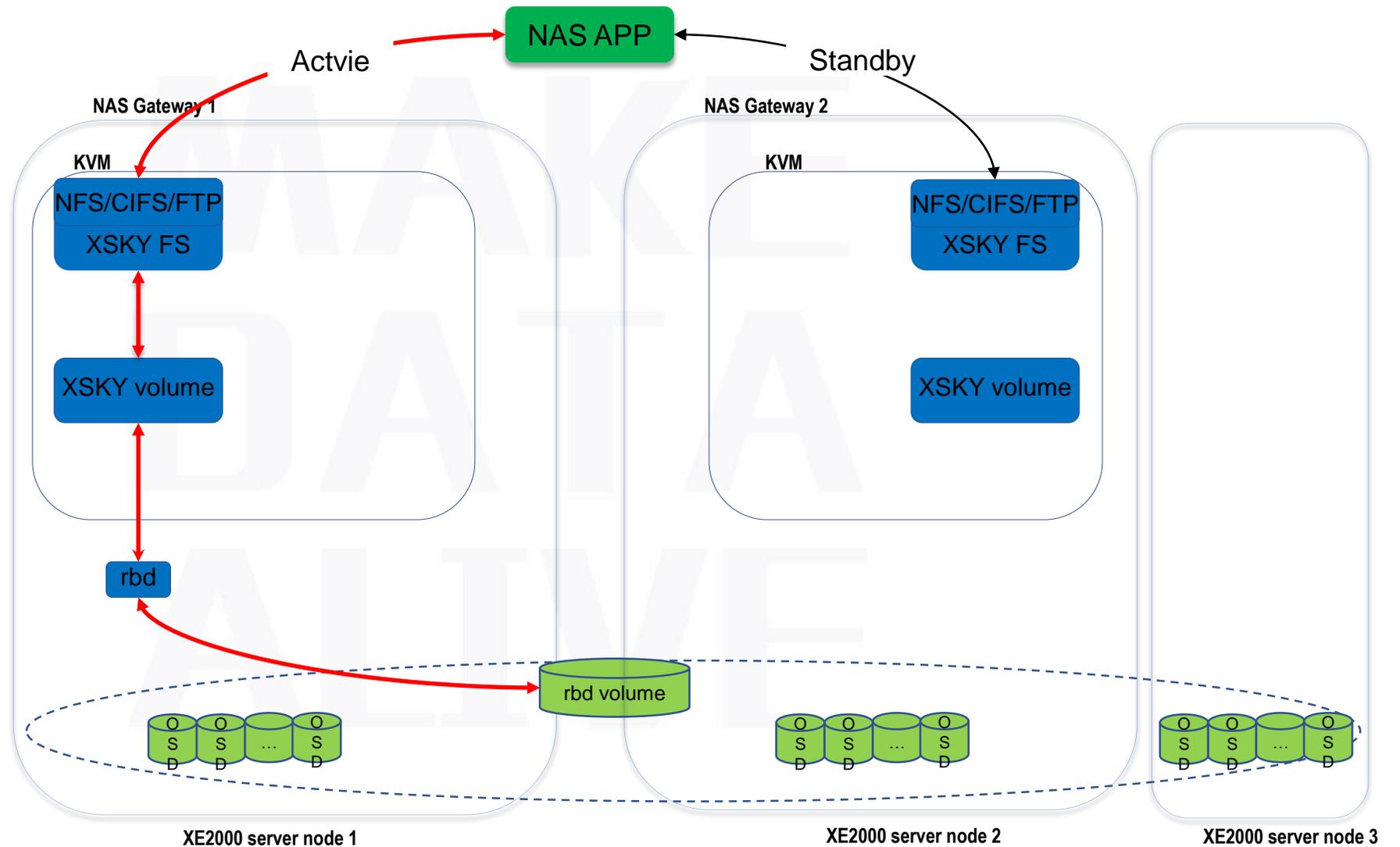
Outline

- Background
- Old Design of the NAS Gateway

MAKE
DATA
ALIVE

Old Design of the NAS Gateway

- AS-NAS (Active-Standby)
- KVM based
- Performance based on cluster



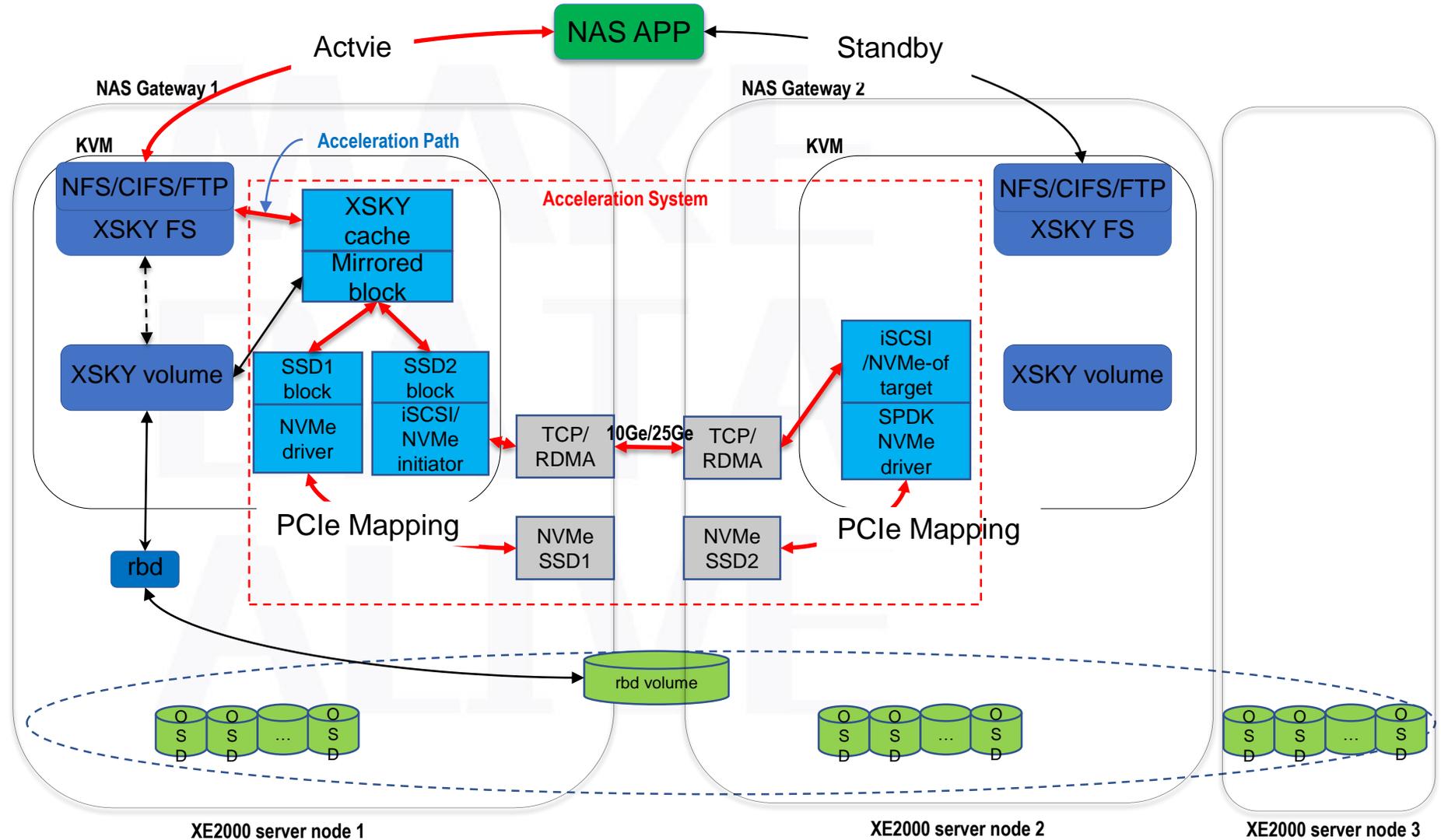
Outline

- Background
- Old Design of the NAS Gateway
- **New Design of the NAS Gateway with SPDK**
 - iSCSI based on SPDK
 - NVMe-of based on SPDK

MAKE
DATA
ALIVE

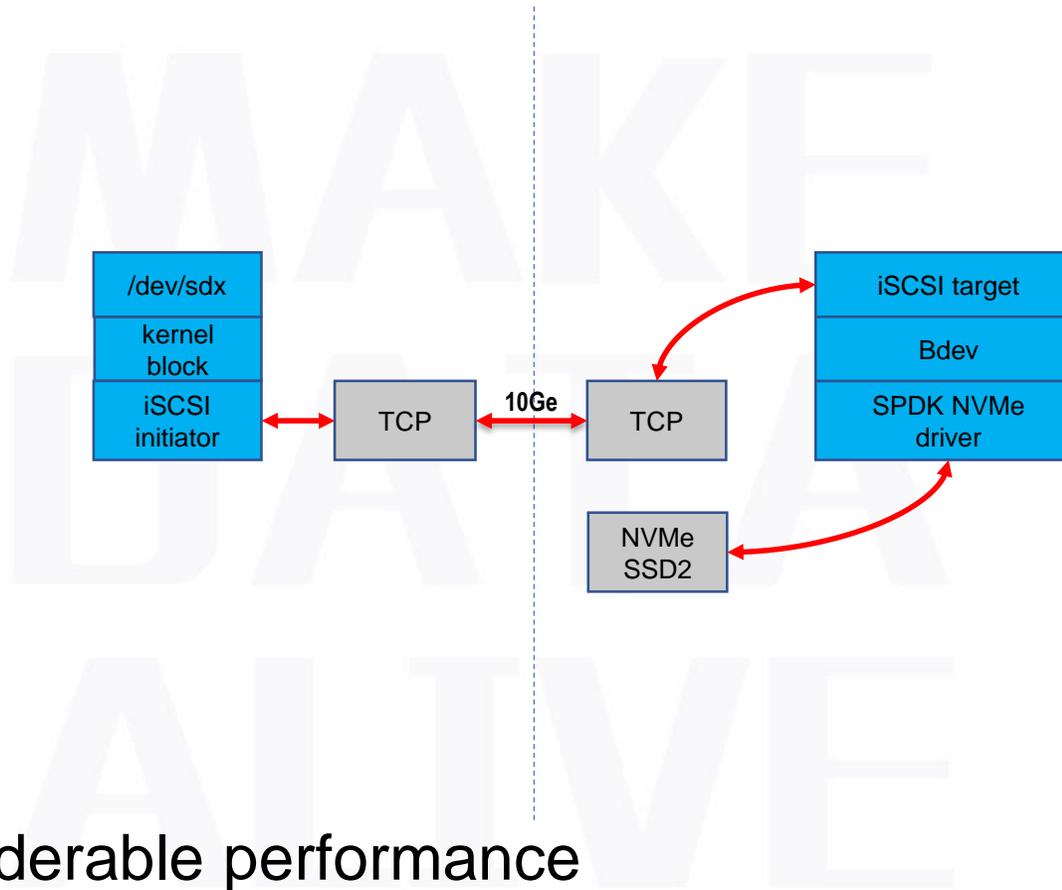
New Design of the NAS Gateway with SPDK

- Performance based on NVMe&&Network
- HA
- Optional protocol stacks(iSCSI or NVMe-of) with different hardware



iSCSI based on SPDK

- iSCSI protocol based on TCP
- Hardware:
 - bonding 10Ge Network
 - NAND NVMe SSD
- Soft:
 - SPDK iscsi tgt
 - Kernel iSCSI initiator



```
[root@nas-edbzztxdbj ~]# cat /etc/spdk/conf/iscsi.conf
[Global]
[Bdev]
[iSCSI]
  NodeRase "iqn.2019-05.XS.spdk"
  AuthFile /usr/local/etc/spdk/auth.conf
  MinConnectionsPerCore 4
  Timeout 30
  DiscoveryAuthMethod Auto
  DefaultTimeWait 2
  DefaultTime2Retain 60
  FirstBurstLength 8192
  ImmediateData Yes
  ErrorRecoveryLevel 0
[PortalGroup1]
  Portal DA1 10.10.10.1:3260
[InitiatorGroup1]
  InitiatorName ANY
  Netmask 10.10.10.2/32
[Nvme]
  TransportID "trtype:PCIe traddr:0000:01:01:0" Nvme0
  RetryCount 4
  TimeoutUseC 0
  ActionOnTimeout None
  AdminPollRate 100000
  HotplugEnable No
  HotplugPollRate 0
[TargetNode2]
  TargetName disk2
  TargetAlias "Data Disk2"
  Mapping PortalGroup1 InitiatorGroup1
  AuthMethod Auto
  AuthGroup AuthGroup1
  UseDigest Auto
  LUN0 Nvme0n1
```

Limited cost reach considerable performance

NVMe-of based on SPDK

➤ NVMe-of protocol based on RDMA(RoCE v2)

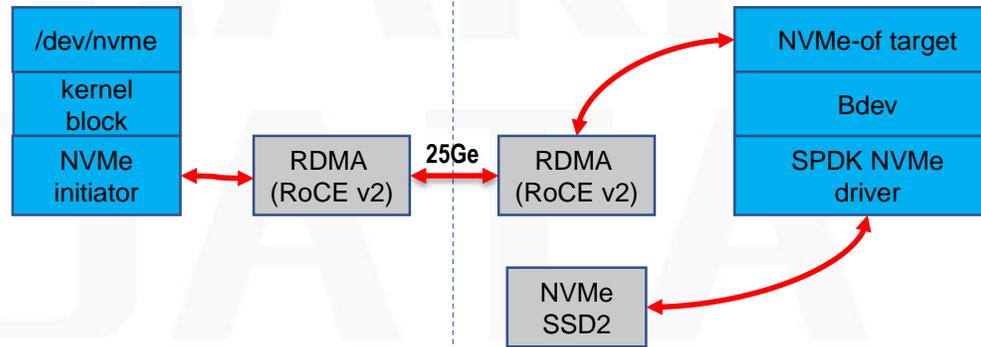
➤ Hardware:

- bonding 25Ge/40Ge Network
- NAND NVMe SSD or Optane SSD

➤ Soft:

- SPDK nvme-of tgt
- Kernel nvme-of initiator

Pursue better performance



```
[root@ns-n3r3d4yuk ~]# cat nvm.conf
(Global)
# SetterPollRate 10000
(Nvme)
# TransportId "trtype:PCIe traddr:0000:01:01:0" Nvme0
RetryCount 4
Timeout 0
ActionOnTimeout None
AdminPollRate 10000
NvmeInEnable No
(Subsystem)
NVMe nqn.2016-06.io.spdk.cnodel
Listen RDMA 10.10.10.2:4420
AttachHostType
Host nqn.2016-06.io.spdk:init
SN SPDK0000000000000001
Namespace Nvme0n1
(Transport)
# Set RDMA transport type.
Type RDMA
# Set the maximum number of outstanding I/O per queue.
#MaxQueueDepth 128
# Set the maximum number of submission and completion queues per session.
# Setting this to '8', for example, allows for 8 submission and 8 completion queues
# per session.
#MaxQueuesPerSession 4
# Set the maximum in-capsule data size. Must be a multiple of 16.
# 0 is a valid choice.
#(Capsule)CapsuleSize 4096
# Set the maximum I/O size. Must be a multiple of 4096.
#MaxIOSize 131072
# Set the I/O unit size, and this value should not be larger than MaxIOSize
#IOUnitSize 131072
# Set the maximum number of IO for admin queue
#MaxAQDepth 32
# Set the number of pooled data buffers available to the transport
# It is used to provide the read/write data buffers for the pairs on this transport.
#NumSharedBuffers 512
# Set the number of shared buffers to be cached per poll group
#BufCacheSize 32
# Set the maximum number outstanding I/O per shared receive queue. Relevant only for RDMA transport
#MaxSRQDepth 4096
```

Outline

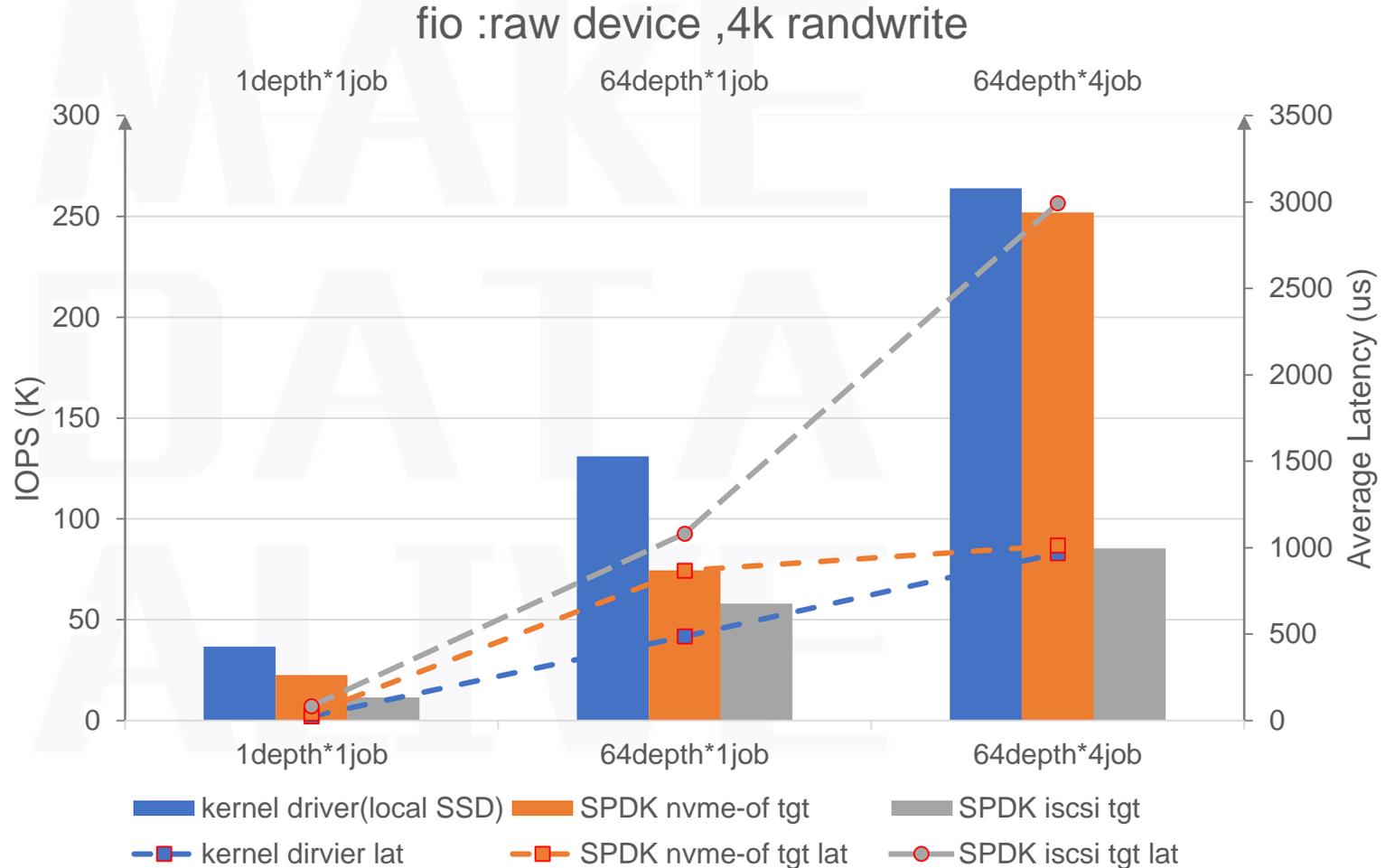
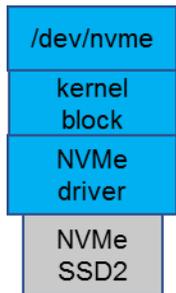
- Background
- Old Design of the NAS Gateway
- New Design of the NAS Gateway with SPDK
 - iSCSI based on SPDK
 - NVMe-of based on SPDK
- Performance comparison
 - Local Kernel VS NVMe-of VS iSCSI
 - Intel Optane VS NAND NVMe
 - Old System VS New System

Local Kernel VS NVMe-of VS iSCSI

➤ Hardware:

- 25Ge Ethernet controller
- NAND NVMe SSD

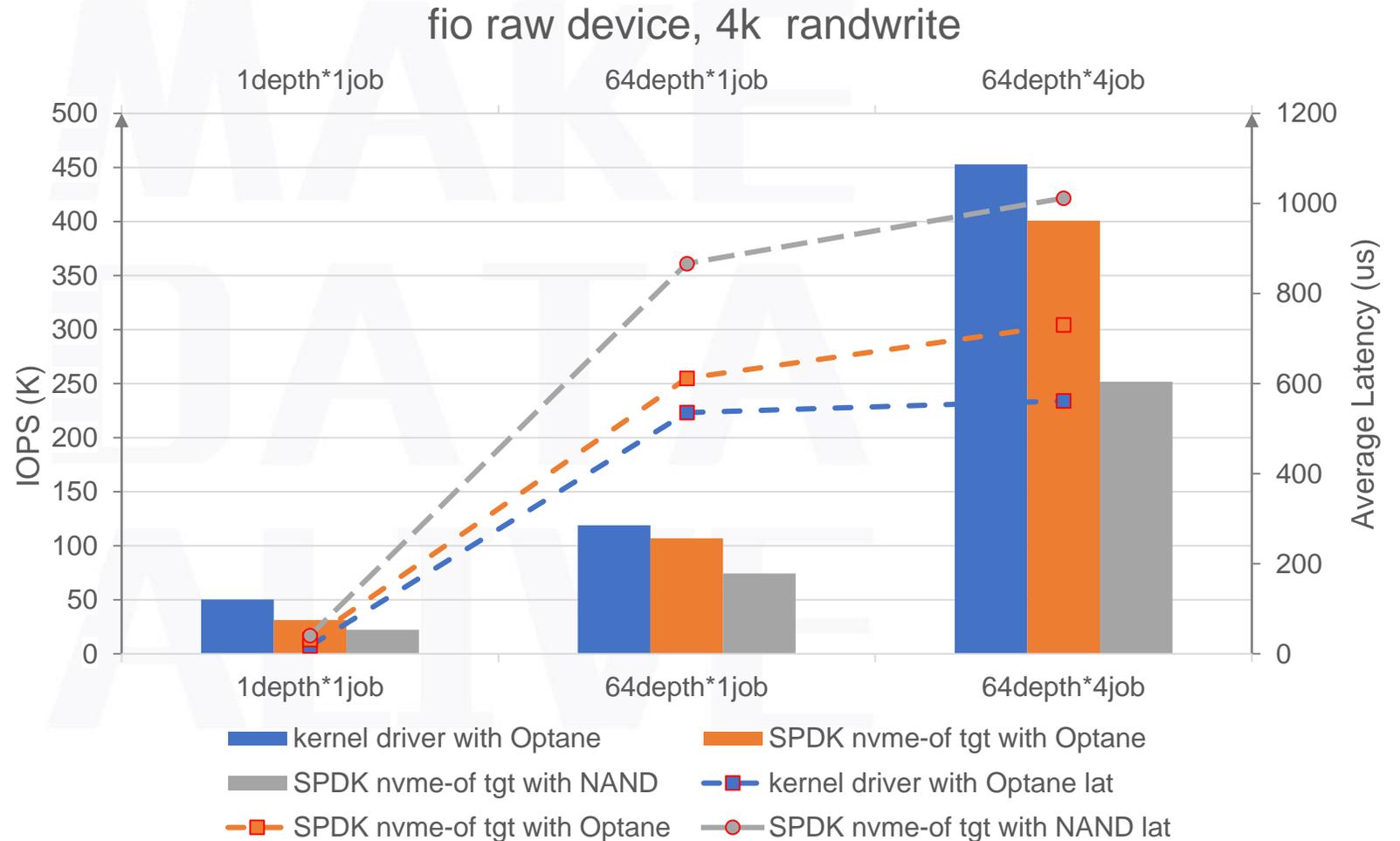
Local kernel is just baseline:



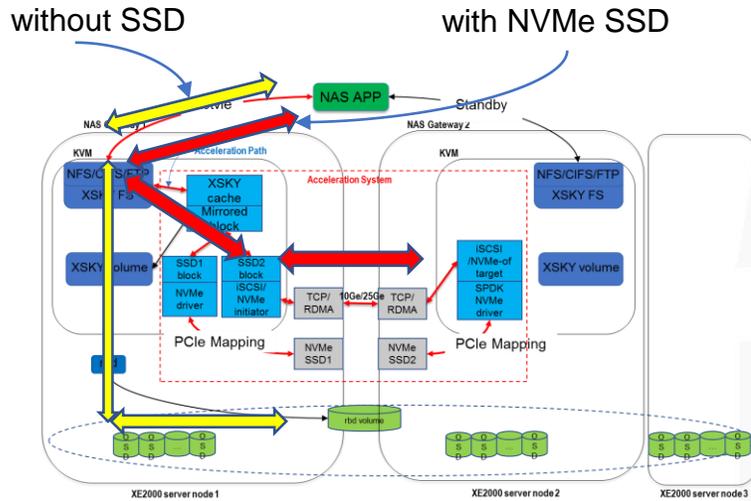
Intel Optane VS NAND NVMe

➤ Hardware:

- 25Ge Ethernet controller
- NAND NVMe SSD
- Intel Optane SSD



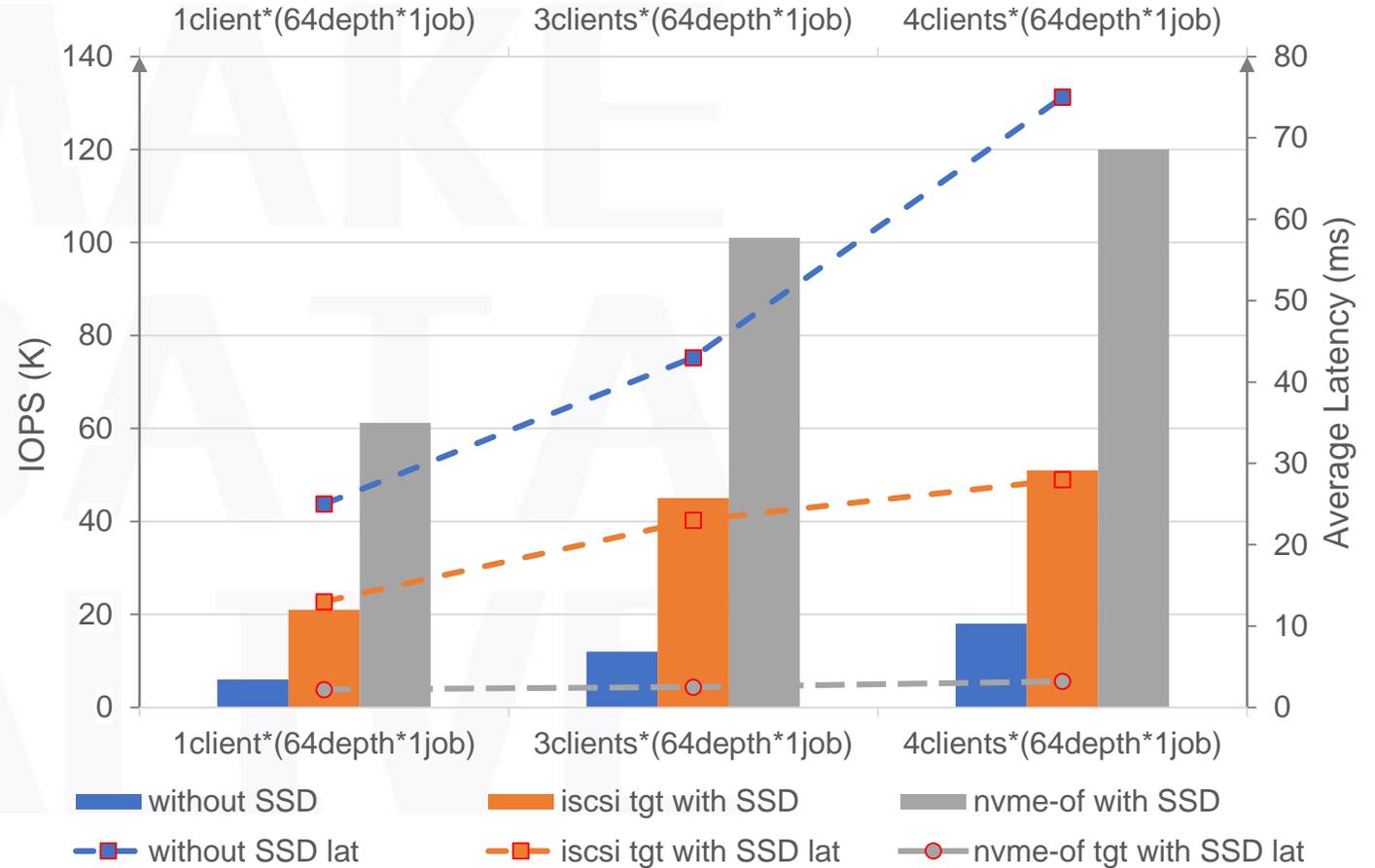
Old System VS New System



End-to-end performance (NFS):

- ① without NVMe SSD(Old system)
- ② iscsi tgt with NVMe SSD
- ③ nvme-of tgt with NVMe SSD

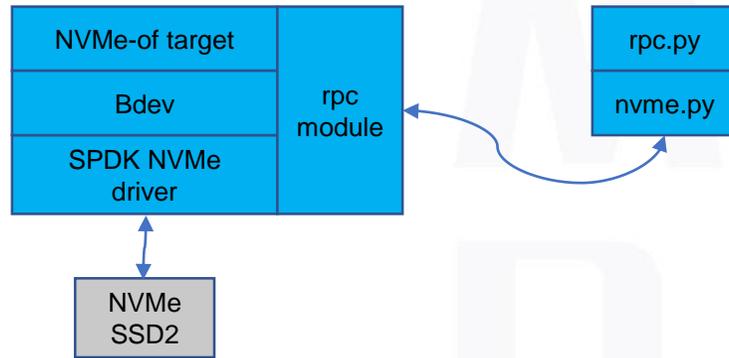
fio: raw device,4k randwrite



Outline

- Background
- Old Design of the NAS Gateway
- New Design of the NAS Gateway with SPDK
 - iSCSI based on SPDK
 - NVMe-of based on SPDK
- Performance comparison
 - Local Kernel VS NVMe-of VS iSCSI
 - intel Optane VS NAND NVMe
 - Old System VS New System
- Add JSON-RPC Methods for NVMe SSD status

Add JSON-RPC Methods for NVMe SSD status



```
[root@nas-n3r3dt4yak ~]# spdk-rpc get_nvme_info -r 'trtype:PCIe traddr:0000:01:01.0' -i health
{"Available Spare Space": "OK",
 "Temperature": "OK",
 "Current Temperature": "303 Kelvin (30 Celsius)",
 "Volatile Memory Backup": "OK",
 "Available Spare": "100%",
 "Available Spare Threshold": "5%",
 "Temperature Threshold": "343 Kelvin (70 Celsius)",
 "Read Only": "No",
 "Device Reliability": "OK",
 "Life Percentage Used": "1%"}

[root@nas-n3r3dt4yak ~]# ps -ef |grep tgt
root      10577 10981 99 06:48 ttyS0    00 03:39 ./app/nvmf_tgt/nvmf_tgt -c /root/nvme.conf
root      10683 10641 0 06:52 pts/0    00.00.00 grep  color=auto tgt
[root@nas-n3r3dt4yak ~]#
```

- Add rpc API to get information of the nvme SSD attached in iscsi/nvme-of tgt.
- Just like “identify” example in SPDK, but for running process.

Beijing | Shanghai | Shenzhen | Chengdu | Nanjing | Wuhan

Thank You



MAKE
D
A
ALIVE