

# SPDK Engineering Practices of Software and Hardware Integrated Virtualization for JD.com Cloud

Qiming zhang

# Agenda

01

Background

02

Integrated Virtualization platform

03

Engineering optimization practices

04

Summary

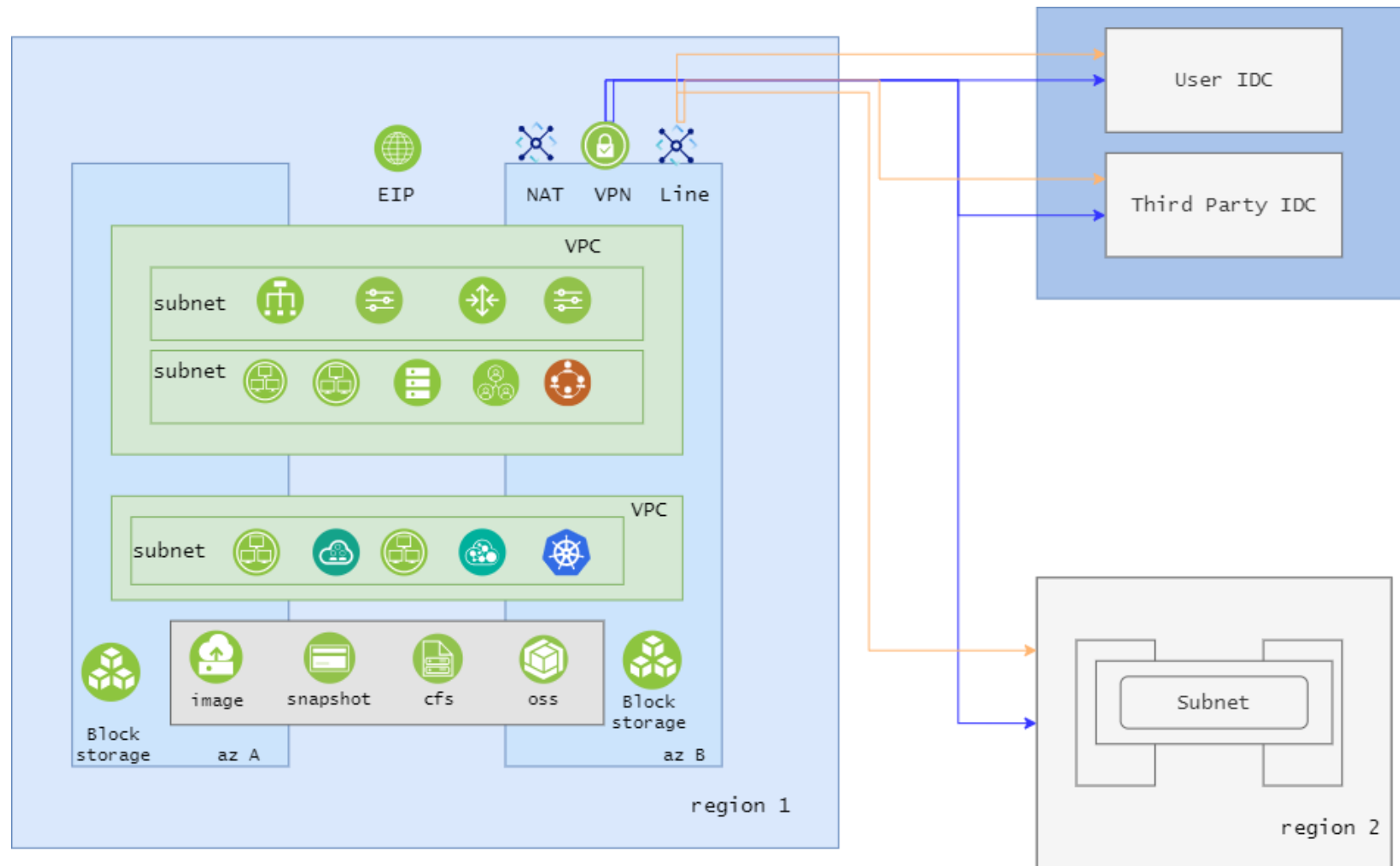
# 01

## Background

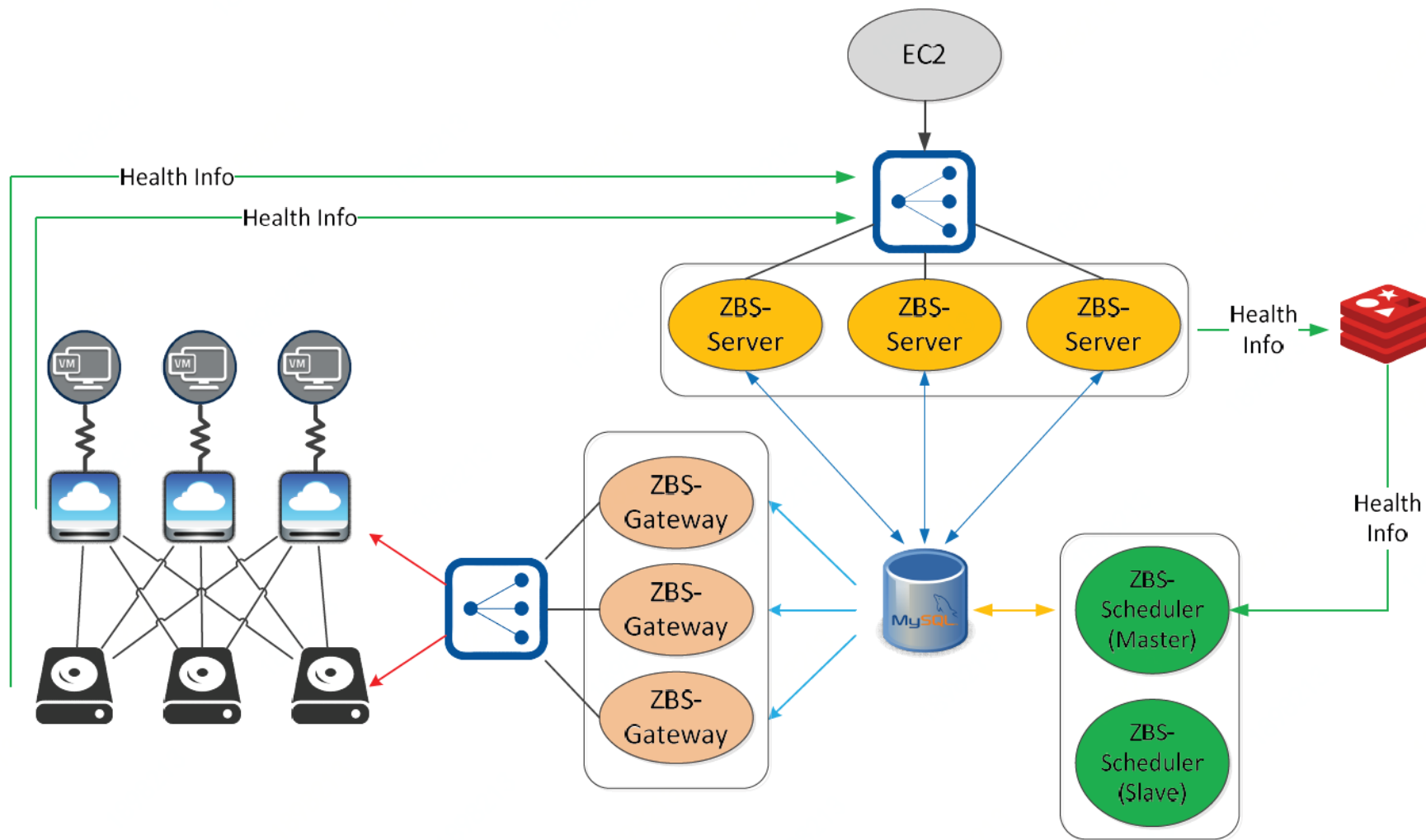
# The challenge of public cloud storage

Faced with the growing storage needs of internal and external users on the cloud, storage is in urgent need of

- Higher performance
- Lower costs
- More stable service



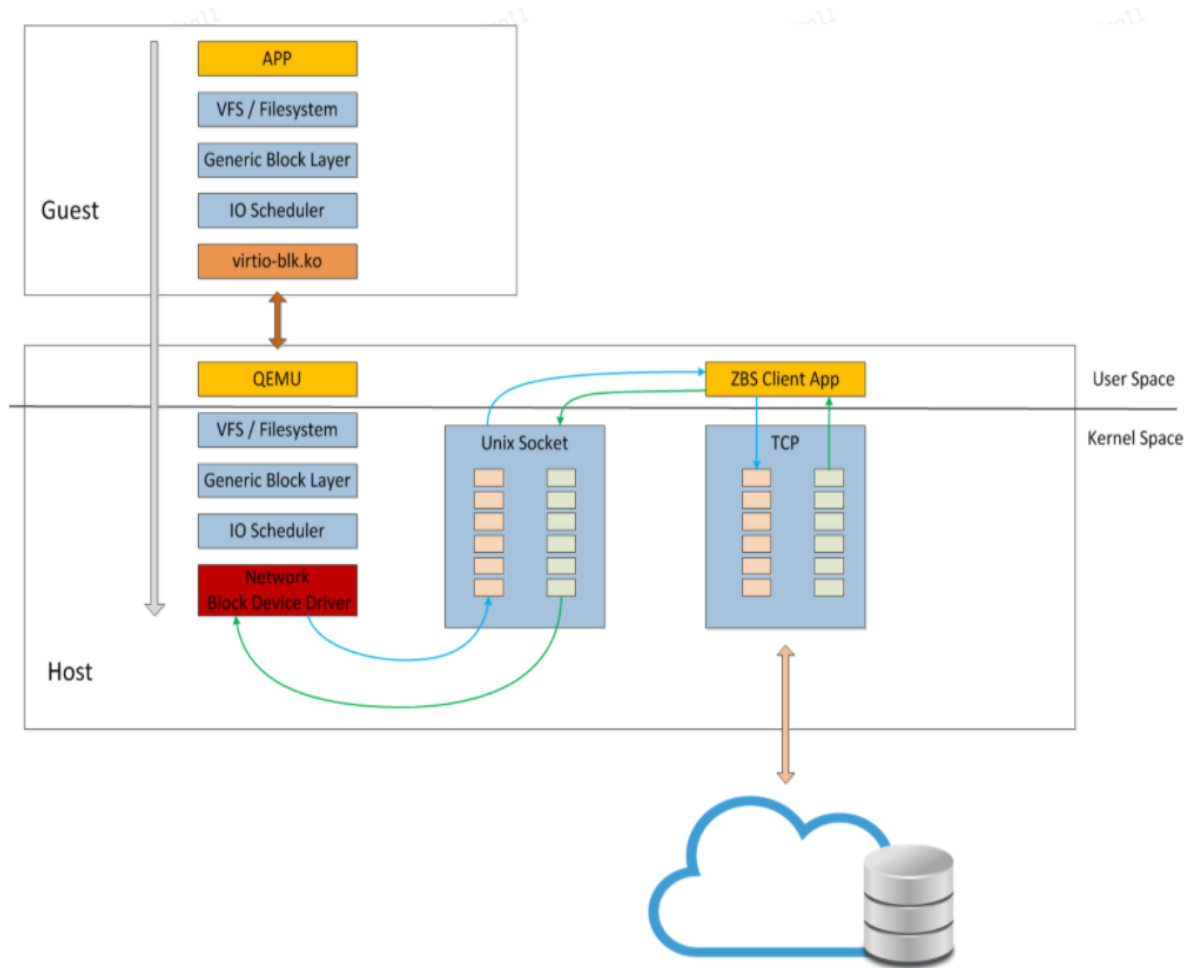
# Storage Architecture



# Architectural defects

qemu && nbd

- Longer io path
- Expensive virtualization overhead
- Service stability impact

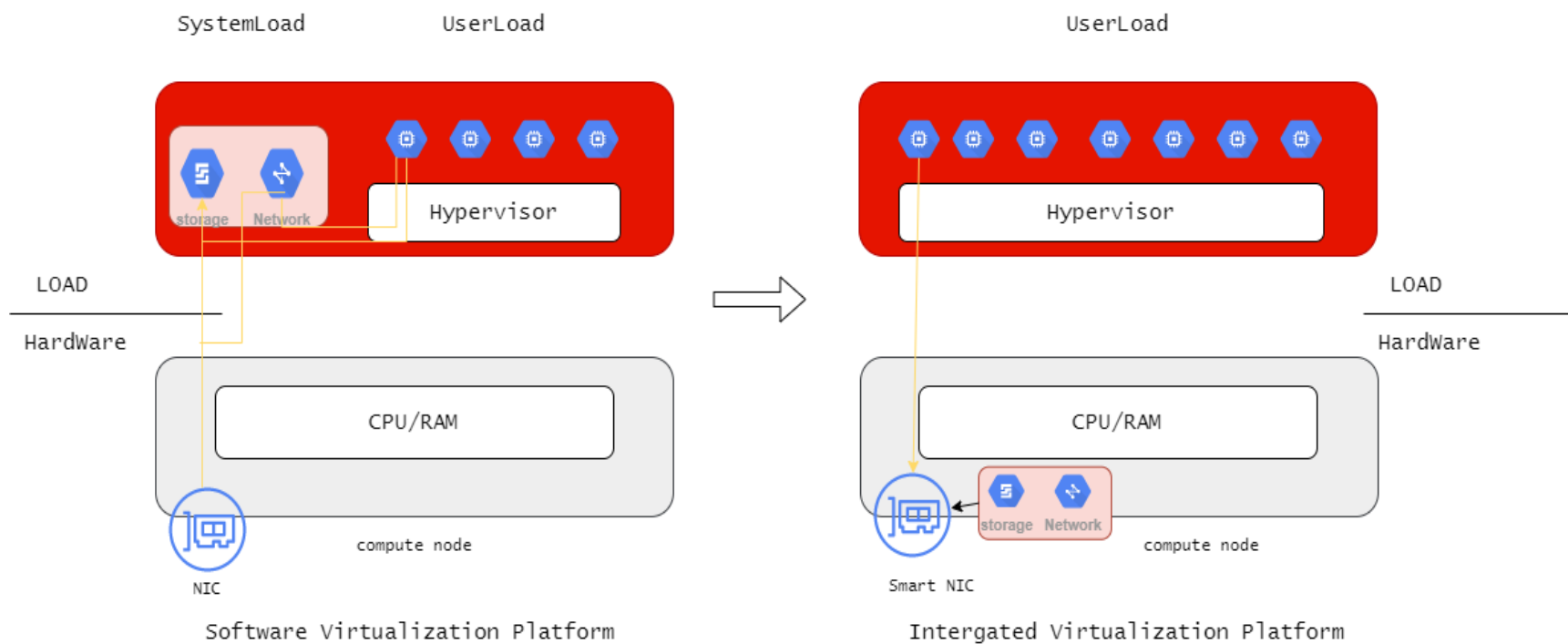


# 02

## Integrated Virtualization platform

# Integrated Virtualization Platform

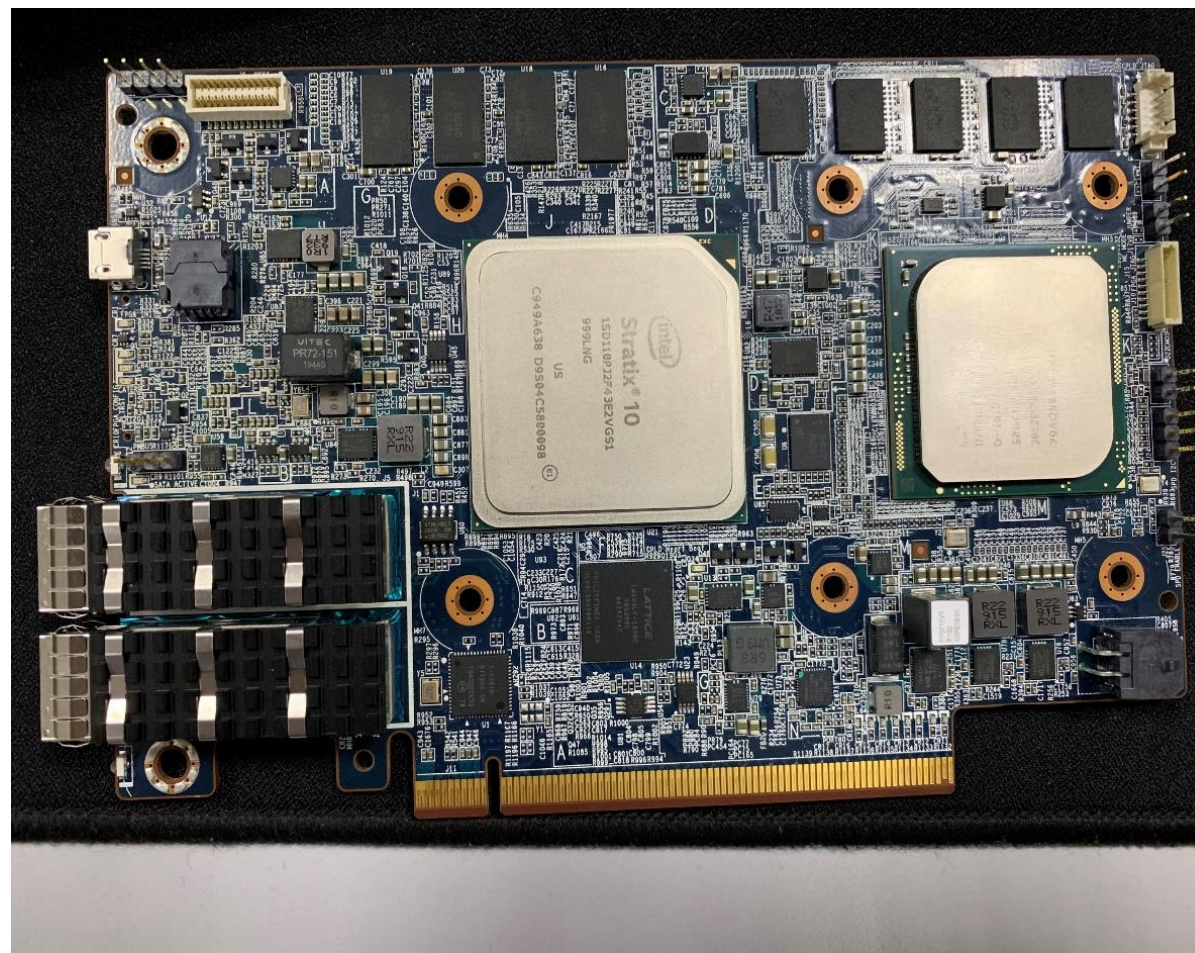
- Reduce virtualization overhead
- Reduce interference and improve stability
- Unified virtualization platform that supports heterogeneous
- Multi-vendor support





# Smart NIC Choice

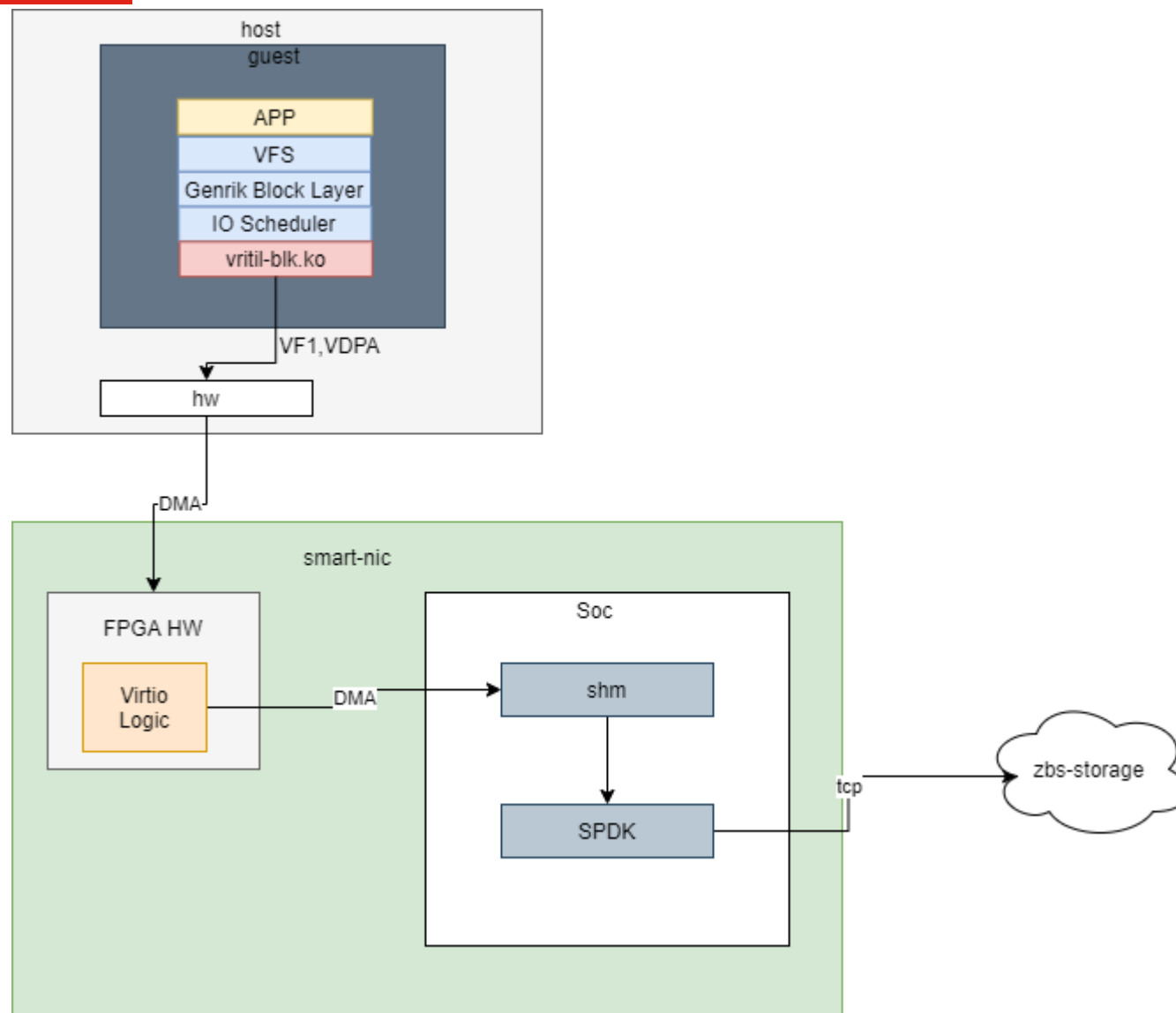
Our choice is working on Intel BSC



# New Storage Architecture

qemu + spdk

- Shorter io path
- Lower virtualization overhead
- Better storage performance
- Improved service stability

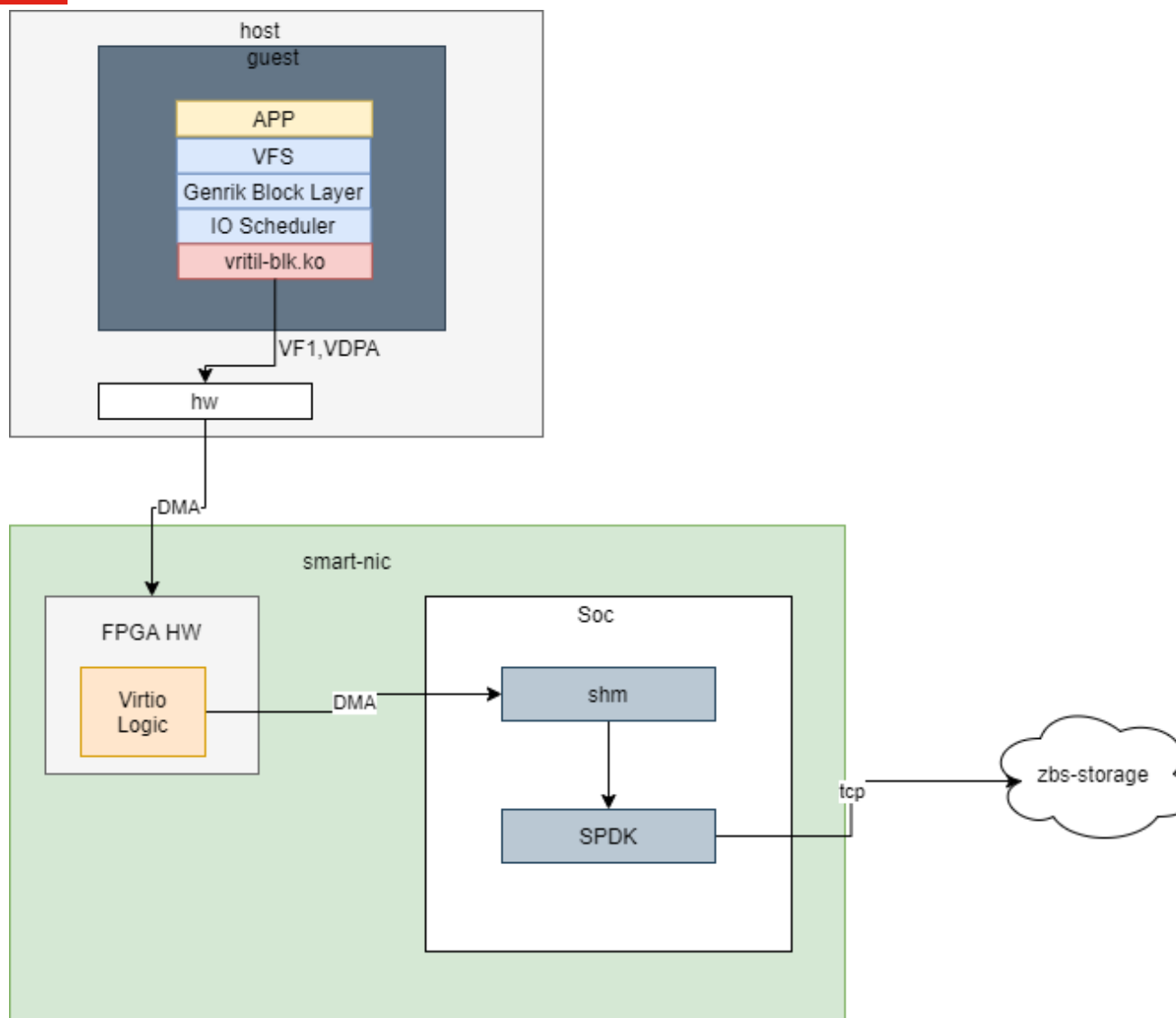


# 03

## Engineering optimization practices

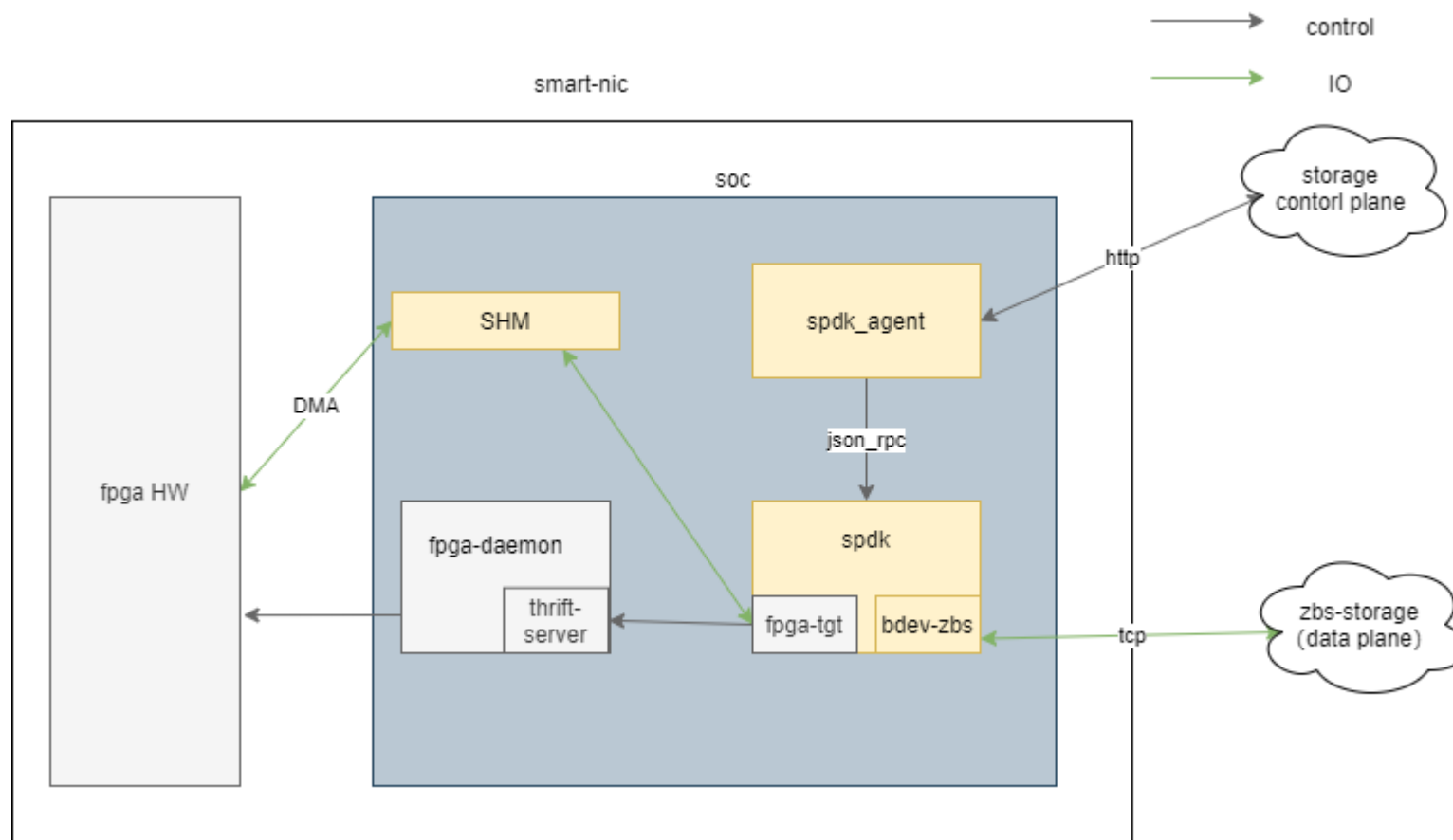
# Problems on the new platform

- R&D effectiveness trade-offs
  - Ability reuse
  - Technical selection
- Hardware capability limitations
  - Longer IO path
  - Host/SOC Status synchronization
  - Resource restrictions



# R&D Efficiency Tradeoff

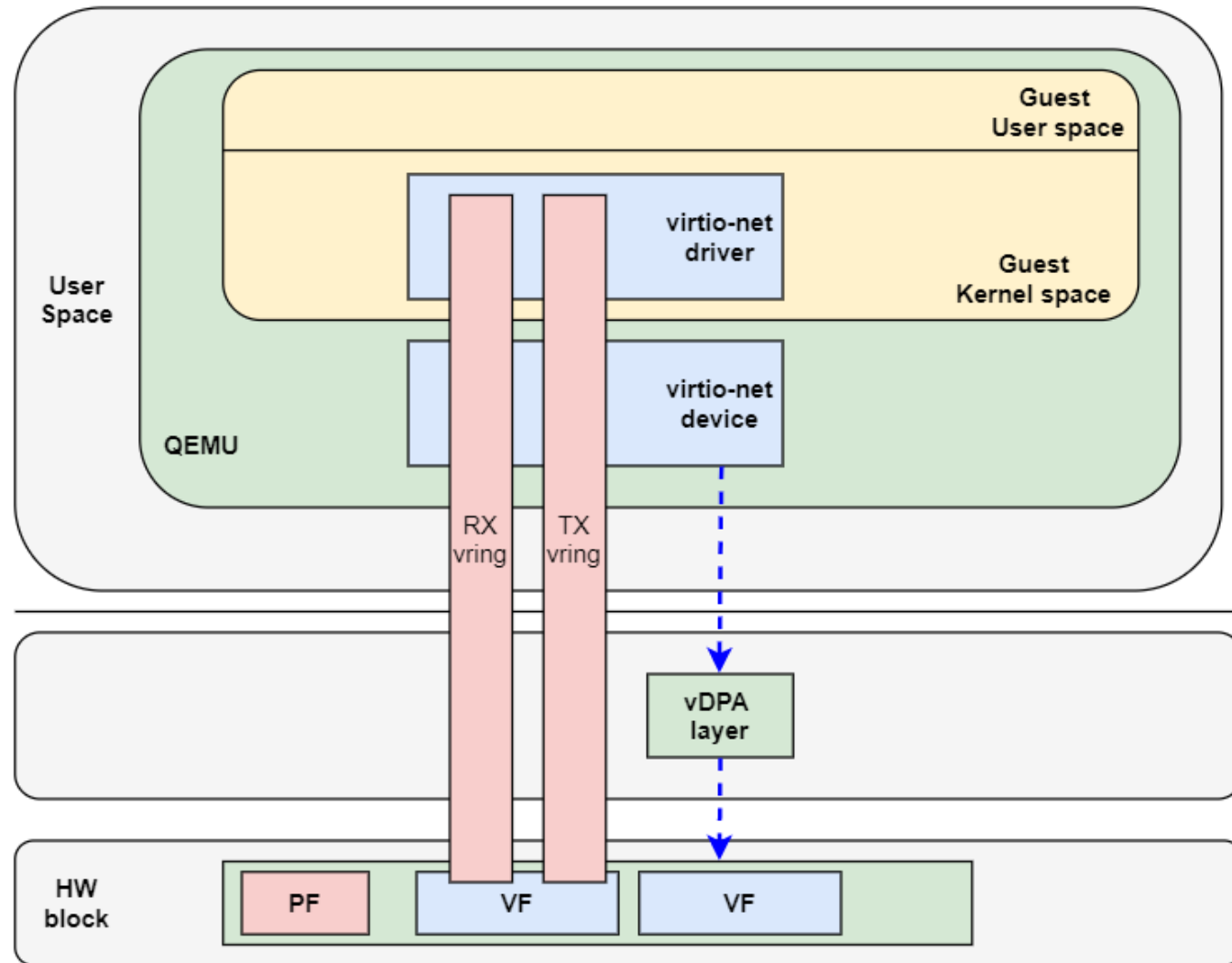
- Re-use control plane capability
- Re-implement data plane: VDPA & SPDK



# VDPA

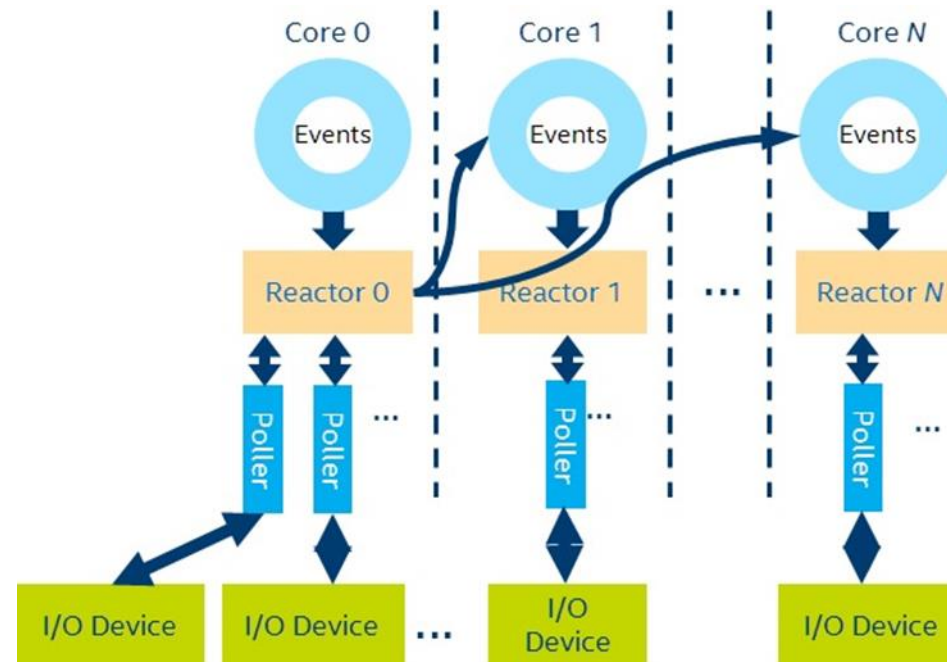
Virtio-net -> vdpas + sriov

- zero copy
- zero interrupt overhead
- Shorten io path



## High Performance

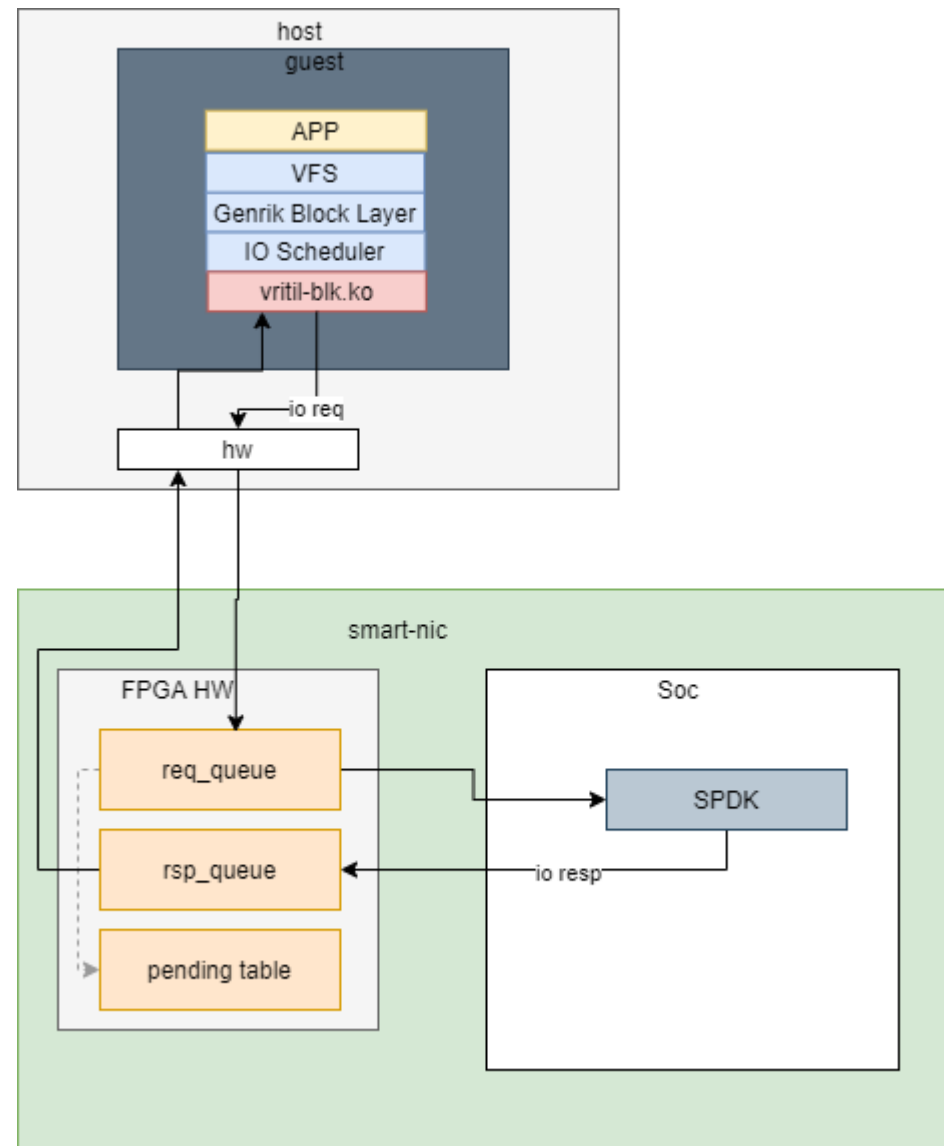
- PMD
- User Space Driver
- Huge page mem
- CPU affinity
- lockless



# Hardware optimization

Our FPGA team provides the ability to solve soc, host status synchronization issues

- Hold io
- Host status check

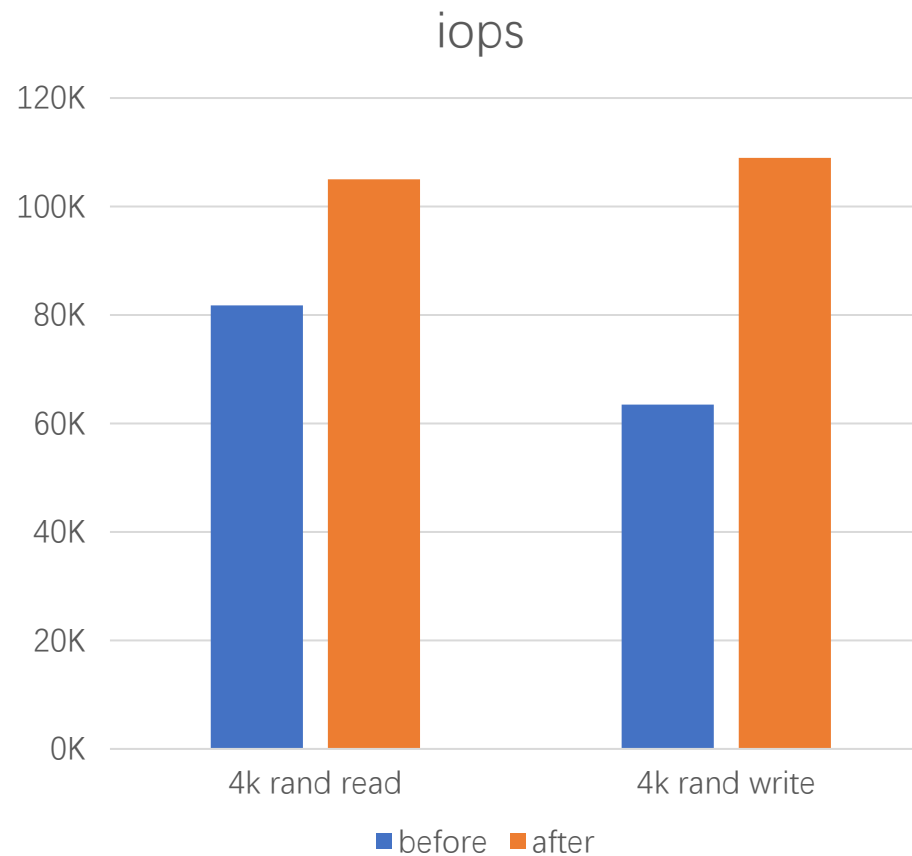




# System optimization

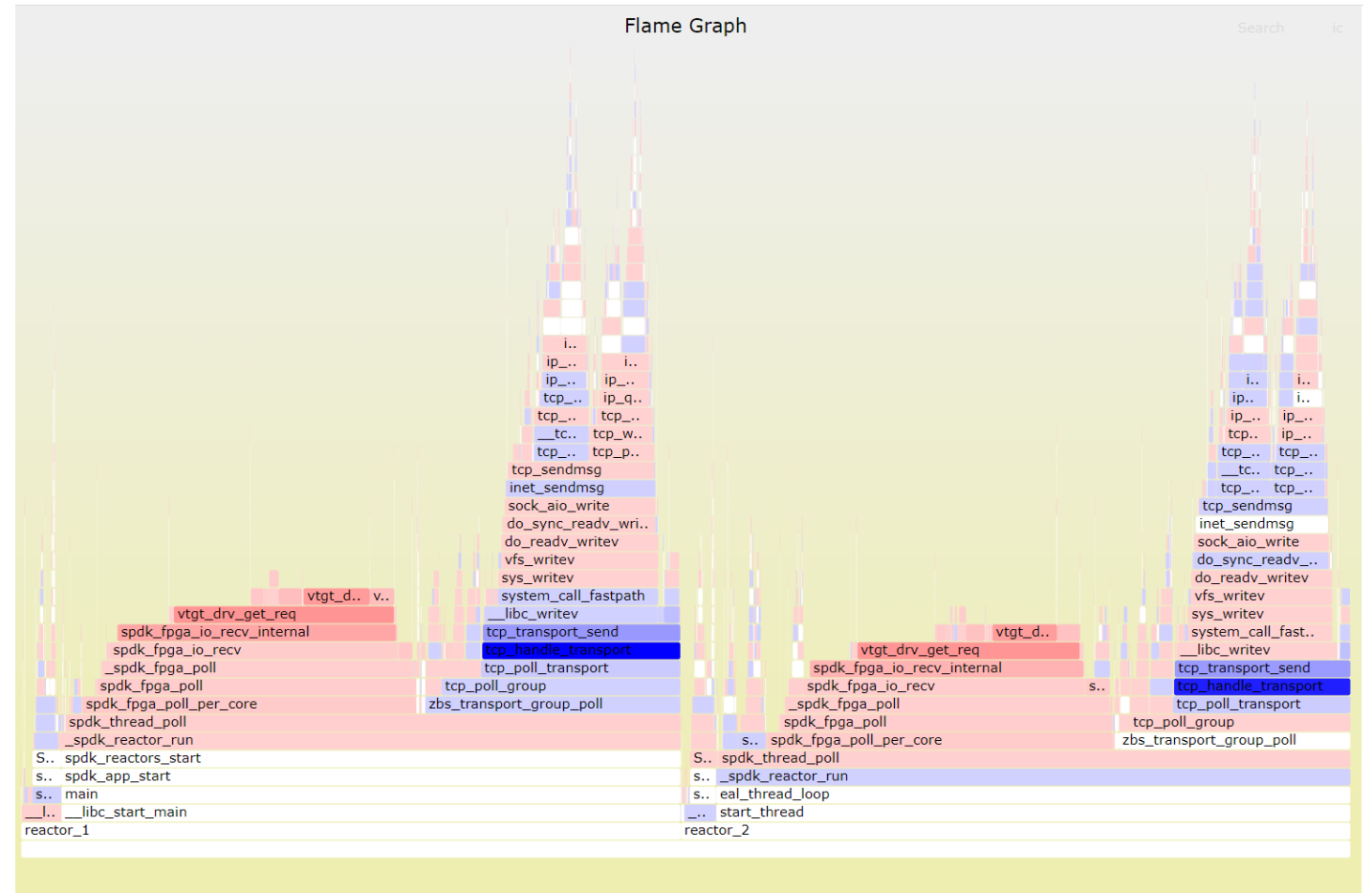
Reduce CPU losses and increase throughput

- Tso
- NIC interrupted
- NIC ring buffer
- IrqBalance
- IPtable
- CPU patch



# Software optimization

- SIMD
- Message batch size tune
- Hotspot code optimization

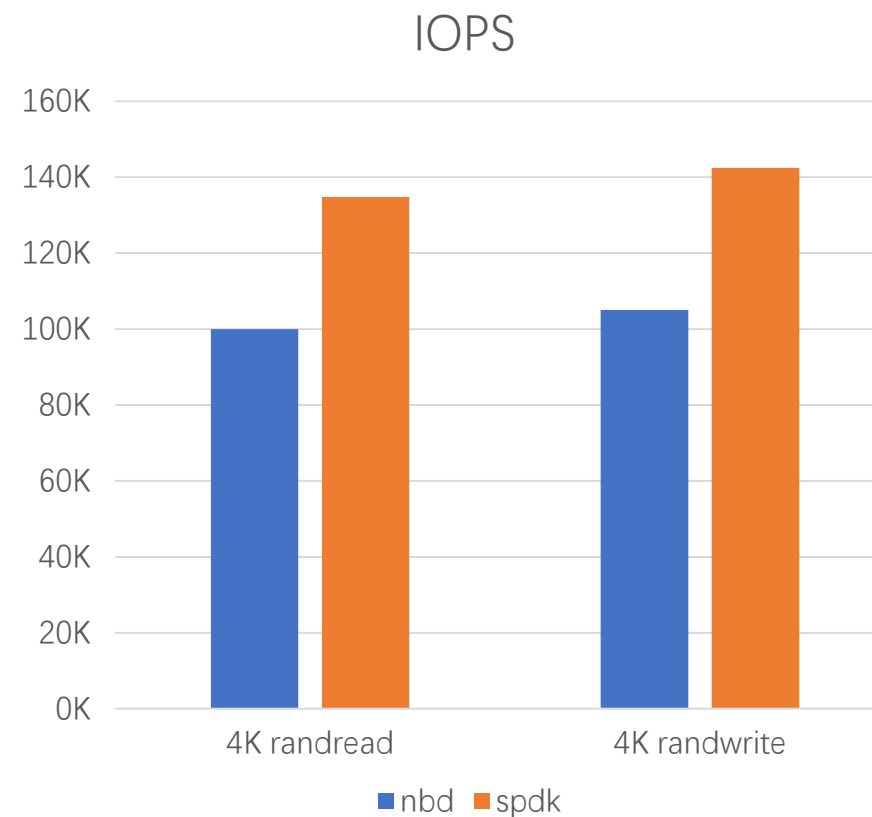


# 04

## Summary

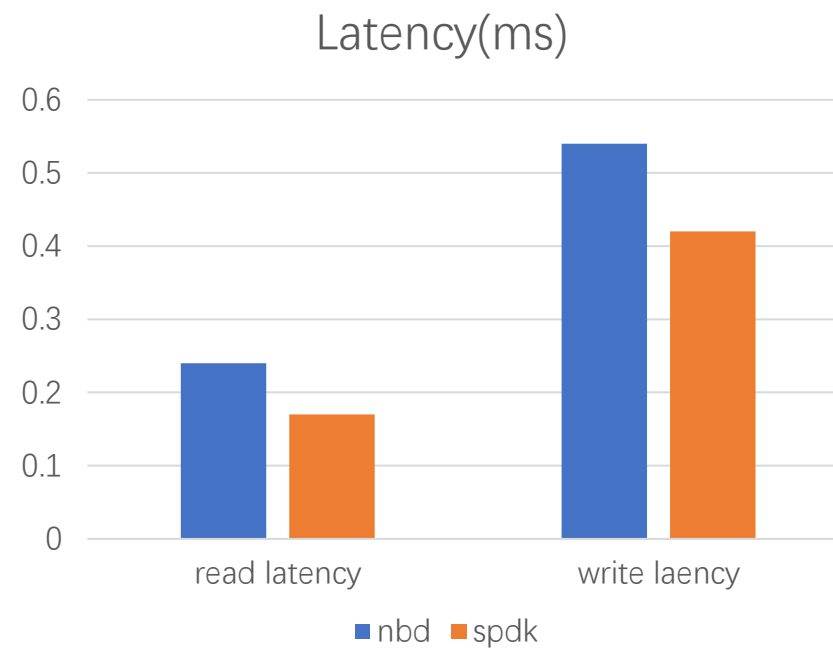
## Summary

- 30% iops performance improvement on single volume
- Upper limit of IOPs of the whole SOC is lower than nbd-client, due to SOC performance limit



# Summary

- 15% reduction in latency



## Future work

- User-level TCP
- RDMA
- Next-generation storage engine

# Q & A