# Overview of the DAOS Object Store

Johann Lombardi, Senior Principal Engineer, AXG, Intel
PMDK & DAOS Tutorial, Dallas, Nov 22, 2022

SC22

Dallas, TX | hpc accelerates.

# Notices and Disclaimers

Intel technologies' features and benefits depend on system configuration and may require enabled hardware, software or service activation. Performance varies depending on system configuration.

No product or component can be absolutely secure.

Tests document performance of components on a particular test, in specific systems. Differences in hardware, software, or configuration will affect actual performance. For more complete information about performance and benchmark results, visit http://www.intel.com/benchmarks .

Software and workloads used in performance tests may have been optimized for performance only on Intel microprocessors. Performance tests, such as SYSmark and MobileMark, are measured using specific computer systems, components, software, operations and functions. Any change to any of those factors may cause the results to vary. You should consult other information and performance tests to assist you in fully evaluating your contemplated purchases, including the performance of that product when combined with other products.   For more complete information visit http://www.intel.com/benchmarks .

Intel Advanced Vector Extensions (Intel AVX) provides higher throughput to certain processor operations. Due to varying processor power characteristics, utilizing AVX instructions may cause a) some parts to operate at less than the rated frequency and b) some parts with Intel® Turbo Boost Technology 2.0 to not achieve any or maximum turbo frequencies. Performance varies depending on hardware, software, and system configuration and you can learn more at http://www.intel.com/go/turbo.

Intel's compilers may or may not optimize to the same degree for non-Intel microprocessors for optimizations that are not unique to Intel microprocessors. These optimizations include SSE2, SSE3, and SSSE3 instruction sets and other optimizations. Intel does not guarantee the availability, functionality, or effectiveness of any optimization on microprocessors not manufactured by Intel. Microprocessor-dependent optimizations in this product are intended for use with Intel microprocessors. Certain optimizations not specific to Intel microarchitecture are reserved for Intel microprocessors. Please refer to the applicable product User and Reference Guides for more information regarding the specific instruction sets covered by this notice.

Cost reduction scenarios described are intended as examples of how a given Intel-based product, in the specified circumstances and configurations, may affect future costs and provide cost savings.  Circumstances will vary.  Intel does not guarantee any costs or cost reduction.
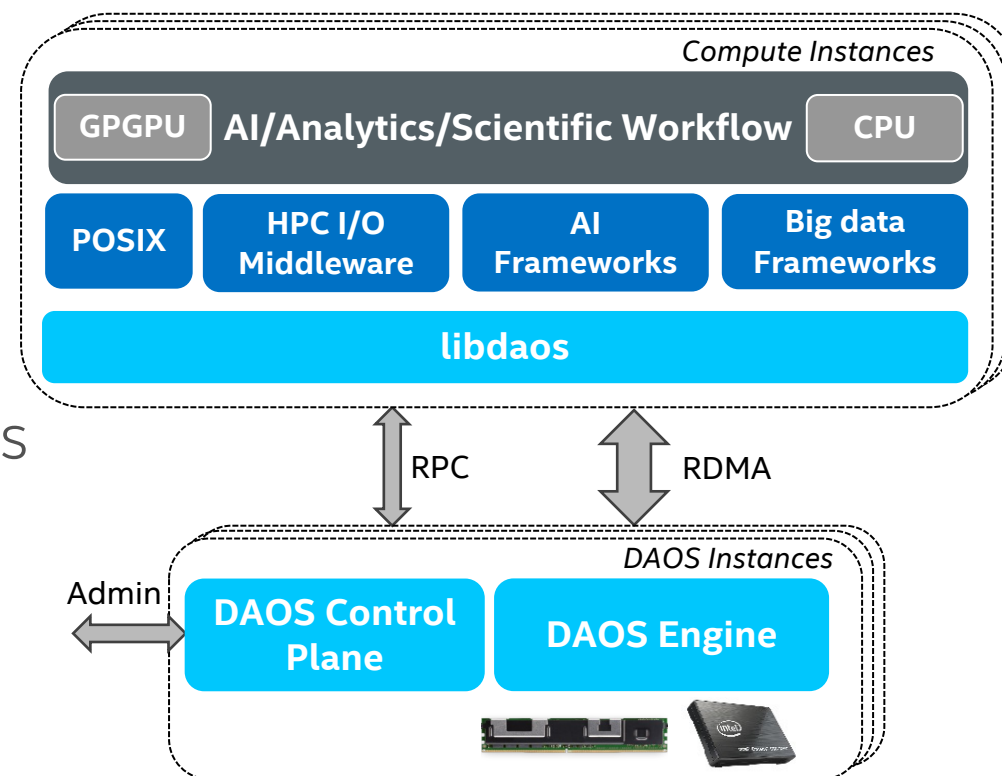
Intel does not control or audit third-party benchmark data or the web sites referenced in this document. You should visit the referenced web site and confirm whether referenced data are accurate.
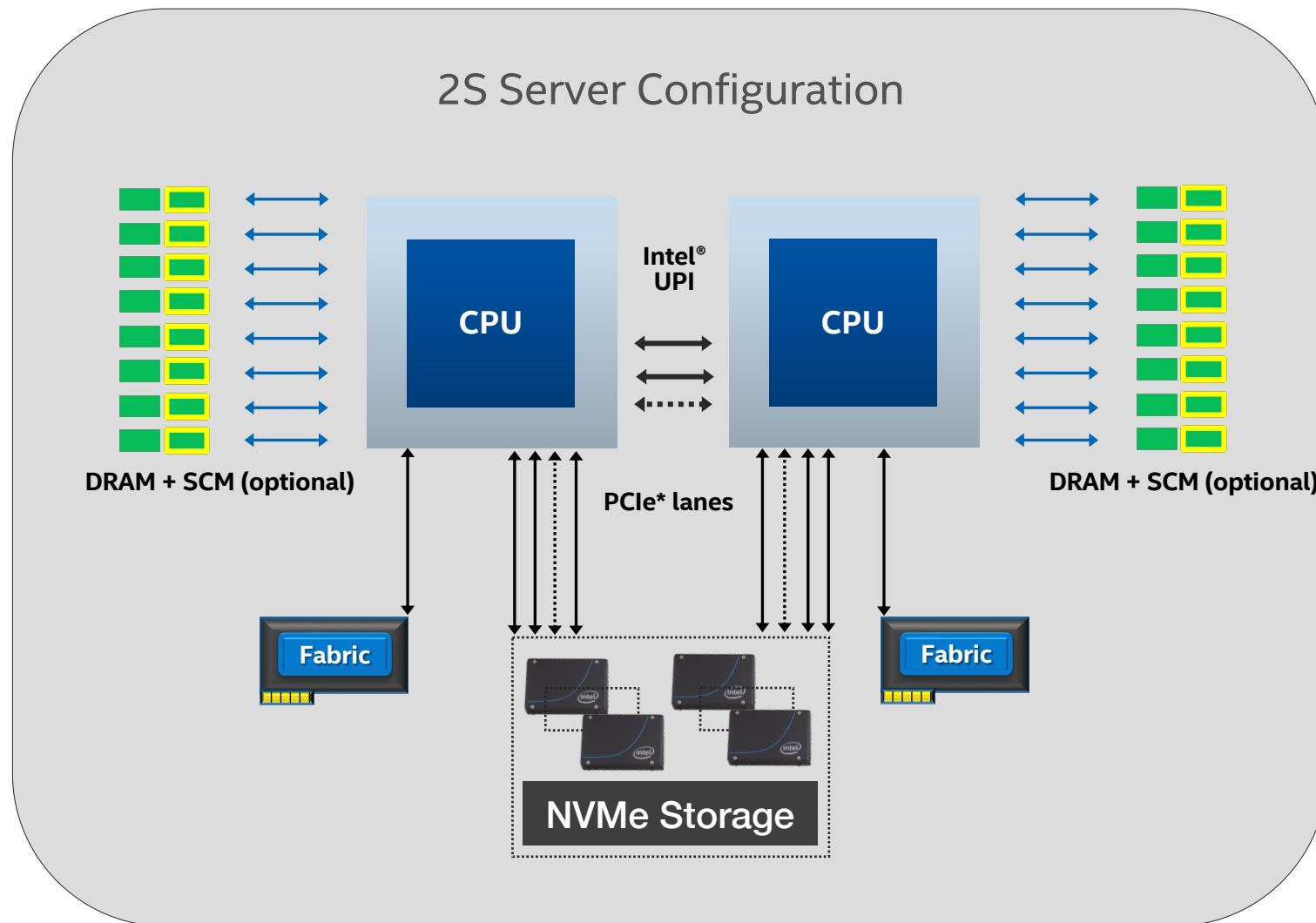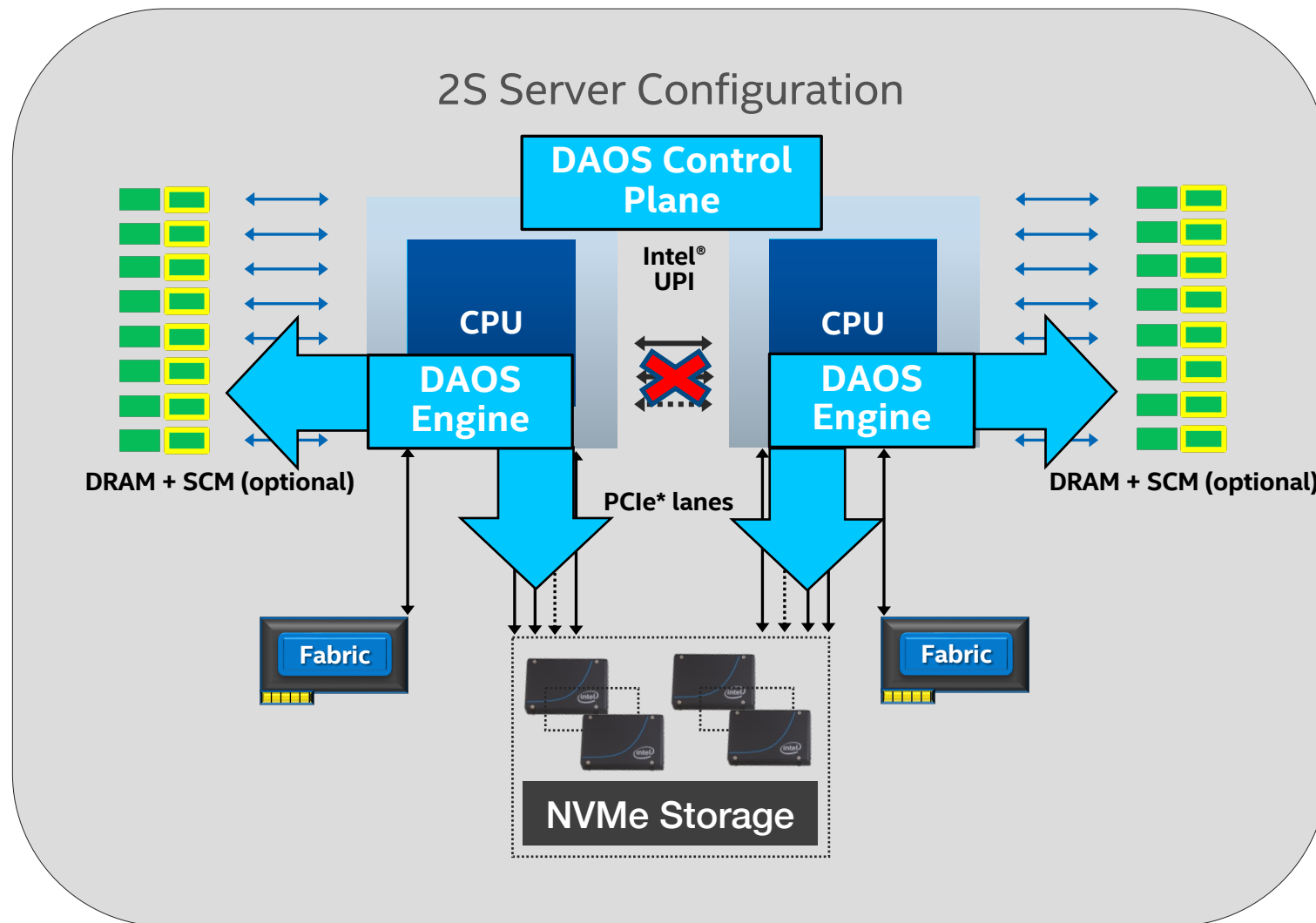
# DAOS: Nextgen Open Storage Platform

- Fully Distributed multi-tenant global namespace

- Platform for innovation
  - Modular API and layering
  - Can leverage latest HW & SW technology

- Built for high performance
  - 10's µs latency, billions of IOPS, TB/s to PB/s

- Full userspace model
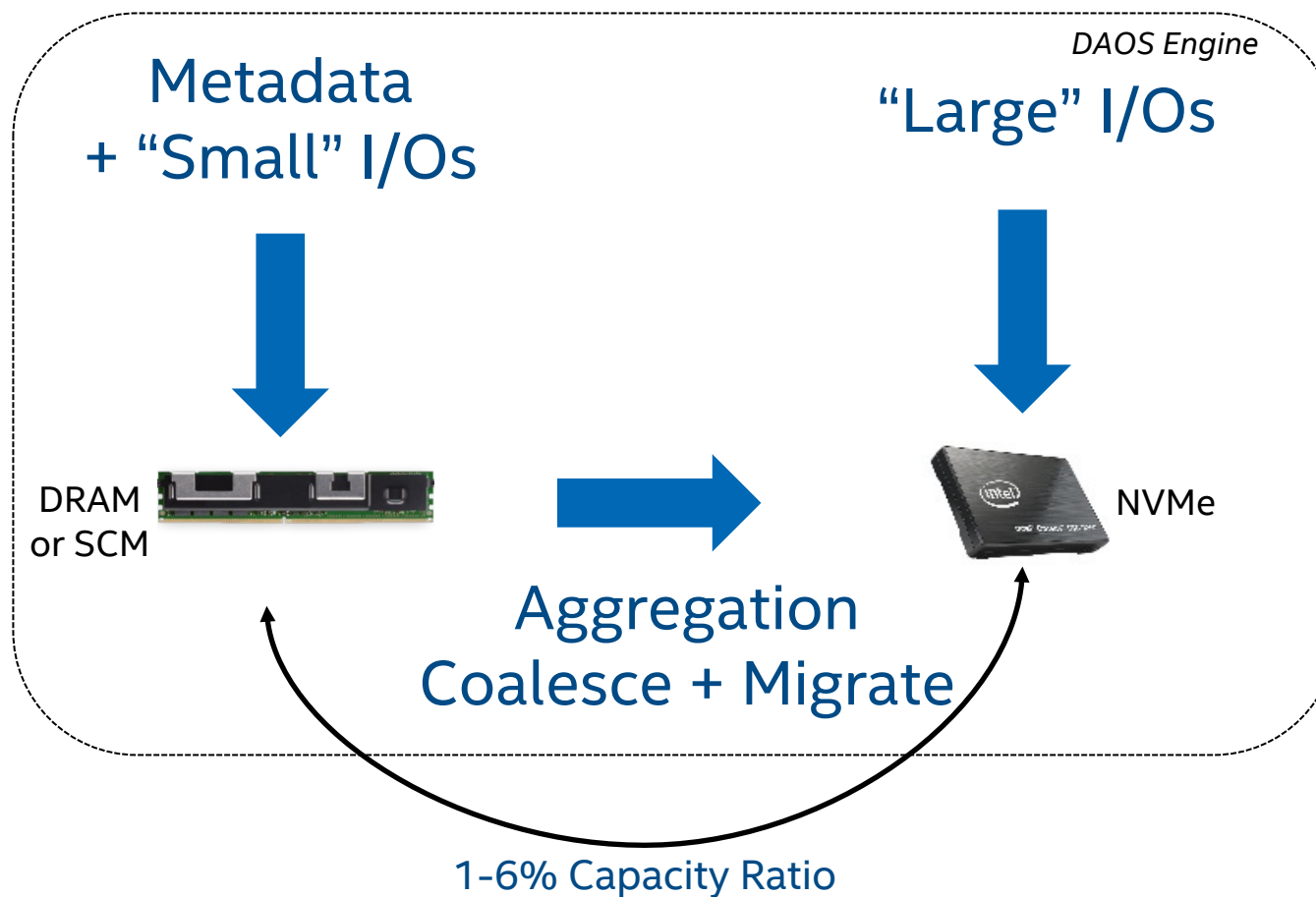  - Run on-prem or in the cloud
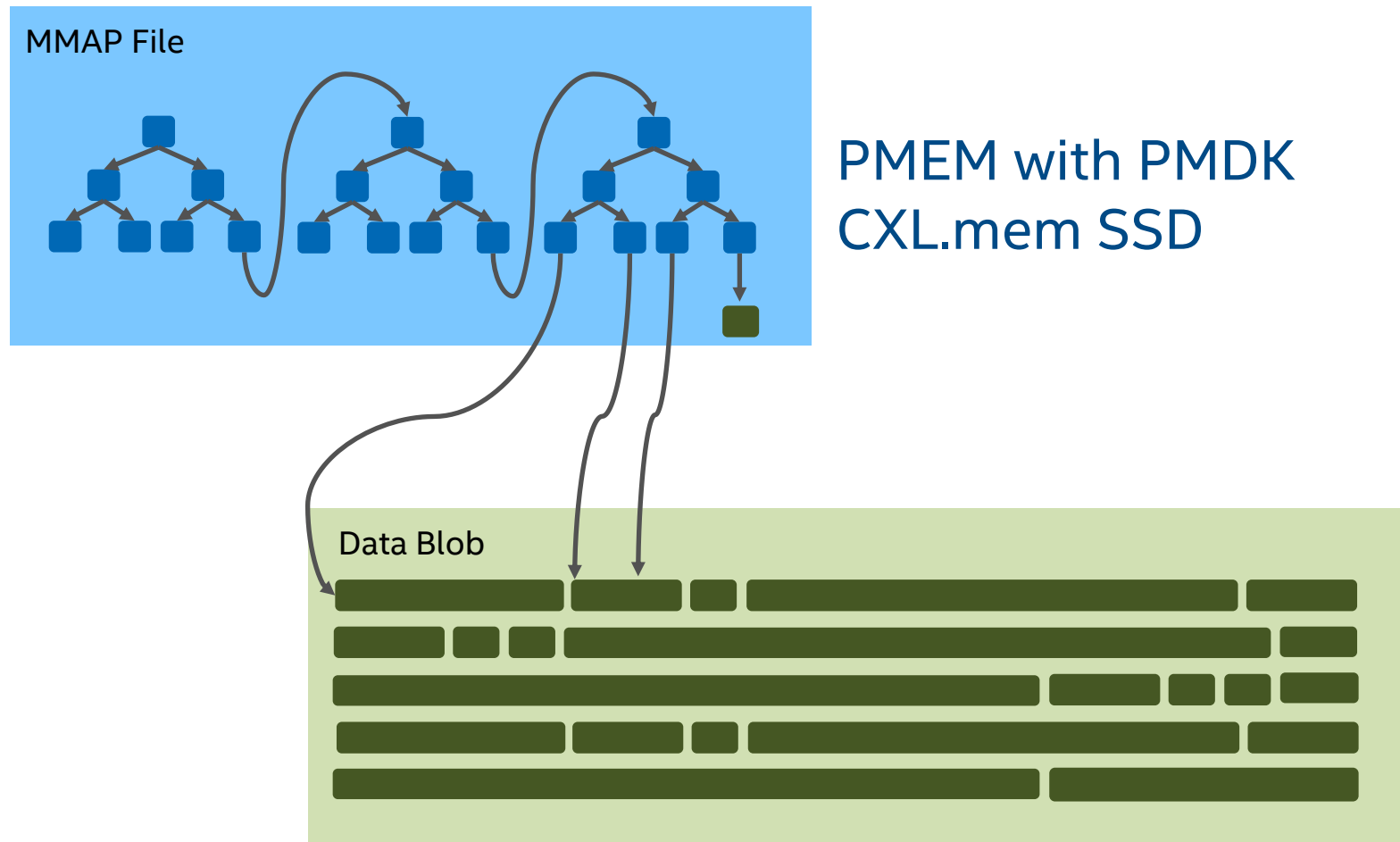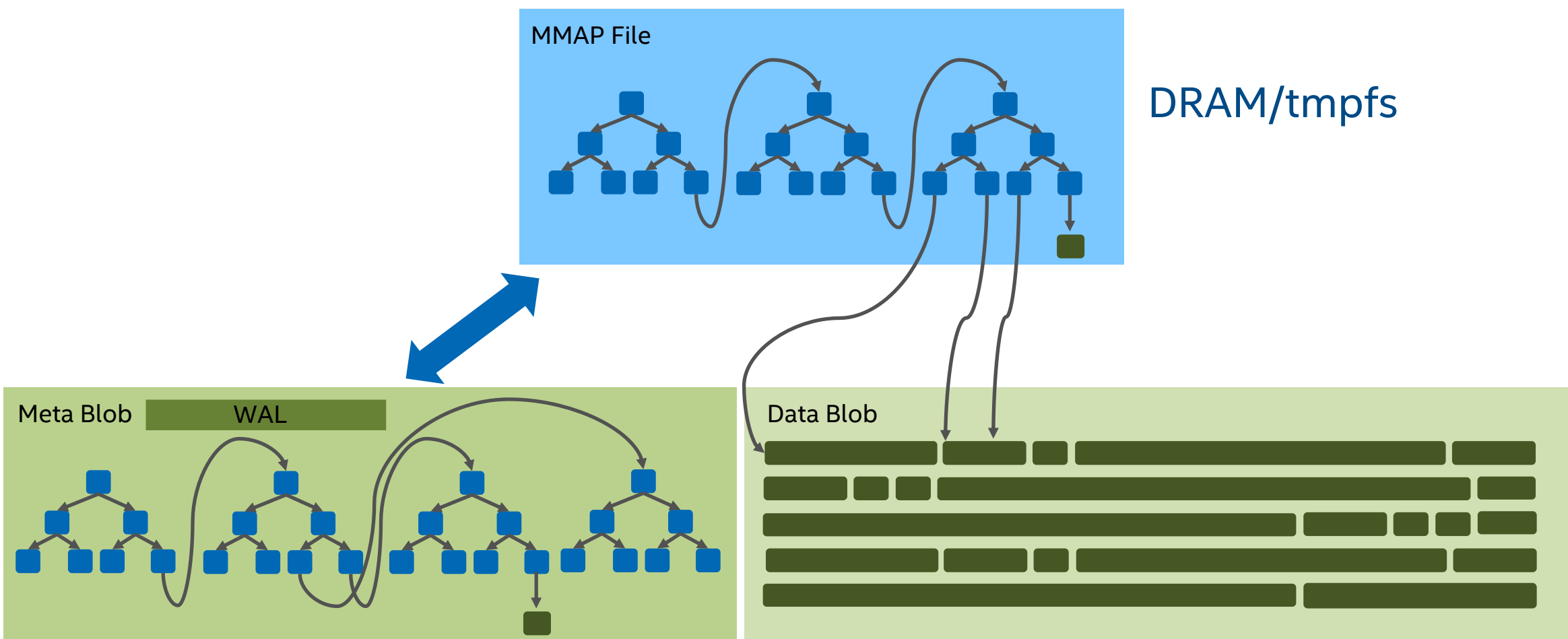
- Growing open-source community



*Compute Instances*

| GPGPU | AI/Analytics/Scientific Workflow | CPU |

| POSIX | HPC I/O Middleware | AI Frameworks | Big data Frameworks |

**libdaos**

RPC          RDMA

*DAOS Instances*

Admin — **DAOS Control Plane**   **DAOS Engine**

intel

# DAOS Node Design



2S Server Configuration

CPU — Intel® UPI — CPU

DRAM + SCM (optional)

PCIe* lanes

Fabric

NVMe Storage

Fabric

# DAOS Node Design



2S Server Configuration

DAOS Control Plane

Intel® UPI

CPU

DAOS Engine

CPU

DAOS Engine

DRAM + SCM (optional)

DRAM + SCM (optional)

PCIe* lanes

Fabric

Fabric

NVMe Storage

# Engine: Media Management



Metadata + "Small" I/Os

DAOS Engine

"Large" I/Os

DRAM or SCM

NVMe

Aggregation
Coalesce + Migrate

1-6% Capacity Ratio

# Metadata: Persistent Device

MMAP File

PMEM with PMDK
CXL.mem SSD

Data Blob

# Metadata: Volatile Device

# Engine Software Stack

DAOS Engine

RPC

VOS

Mercury

UMEM

BIO

OFI

PMDK

AD-MEM

SPDK

Fabric  Fabric

intel OPTANE ›››
PERSISTENT MEMORY

intel

# I/O Operation Flow



SCM
NVMe

# I/O Operation Flow

SCM

NVMe

## DAOS XStream

- pmemobj_reserve() new buffer
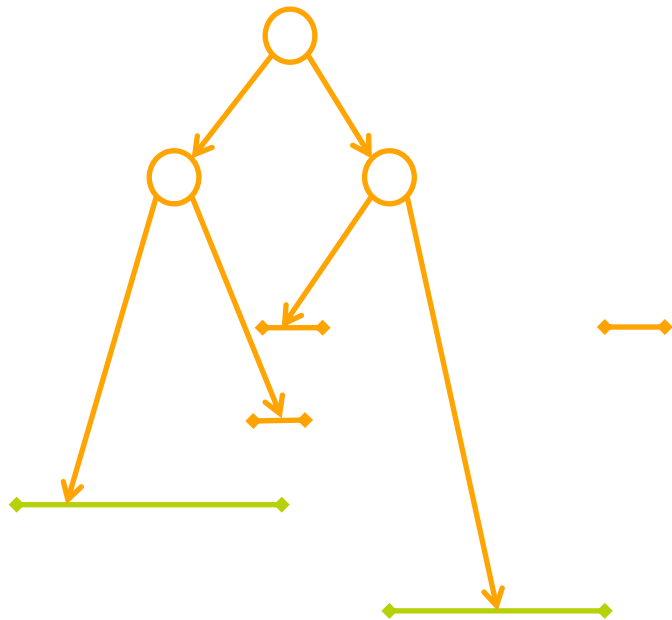
intel.

# I/O Operation Flow

SCM

NVMe

## DAOS XStream

- pmemobj_reserve() new buffer

- Start RDMA transfer to newly allocated buffer
  - Switch to other ULT until completion is reported
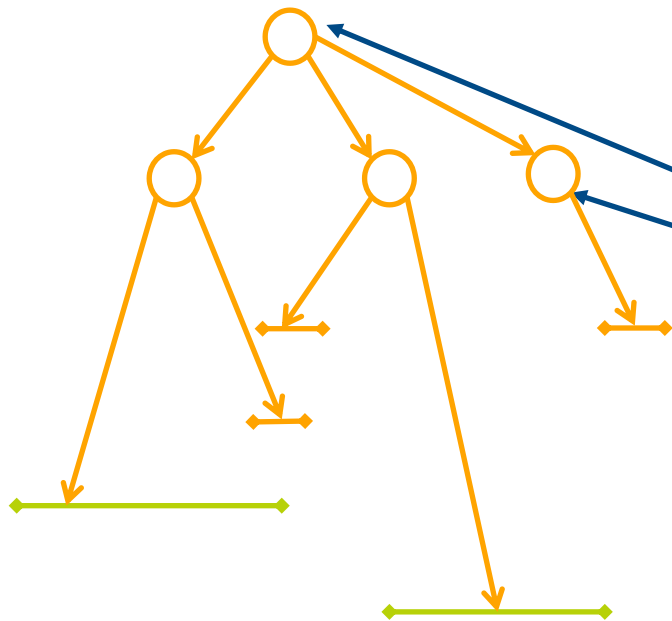
# I/O Operation Flow

SCM

NVMe

## DAOS XStream

- pmemobj_reserve() new buffer

- Start RDMA transfer to newly allocated buffer
  - Switch to other ULT until completion is reported

- If RDMA transfer failed, free buffer with pmemobj_cancel()

# I/O Operation Flow

SCM

NVMe

## DAOS XStream

- pmemobj_reserve() new buffer

- Start RDMA transfer to newly allocated buffer
  - Switch to other ULT until completion is reported

- If RDMA transfer failed, free buffer with pmemobj_cancel()

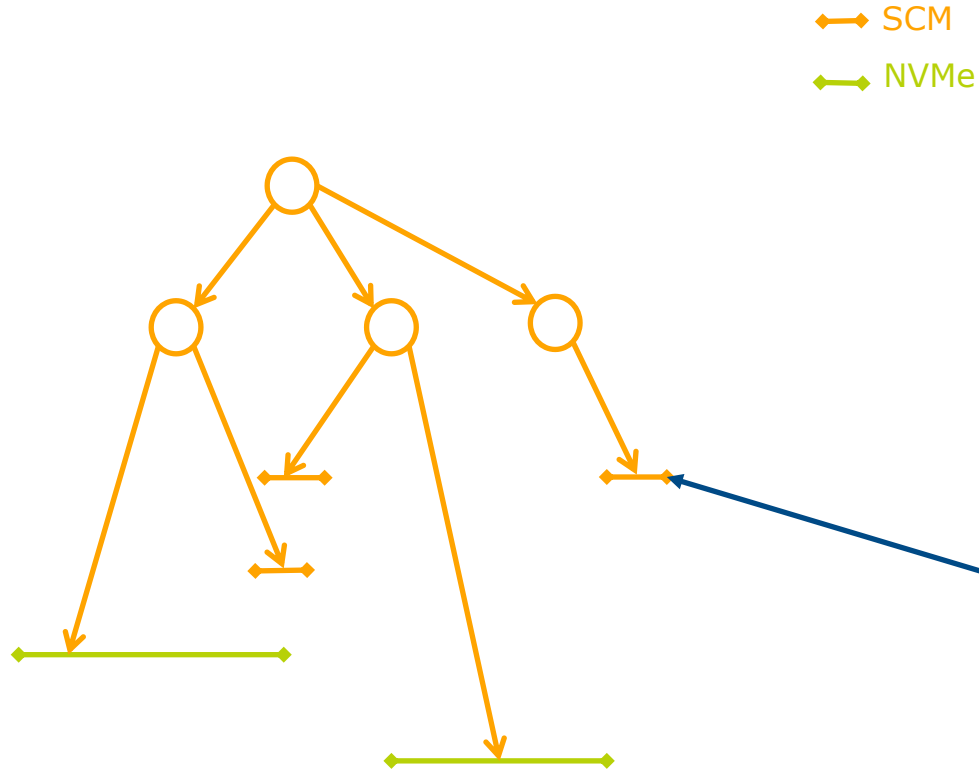- Otherwise, start pmemobj transaction

# I/O Operation Flow

SCM

NVMe

## DAOS XStream

- pmemobj_reserve() new buffer
- Start RDMA transfer to newly allocated buffer
  - Switch to other ULT until completion is reported
- If RDMA transfer failed, free buffer with pmemobj_cancel()
- Otherwise, start pmemobj transaction
- Modify index to insert new extent

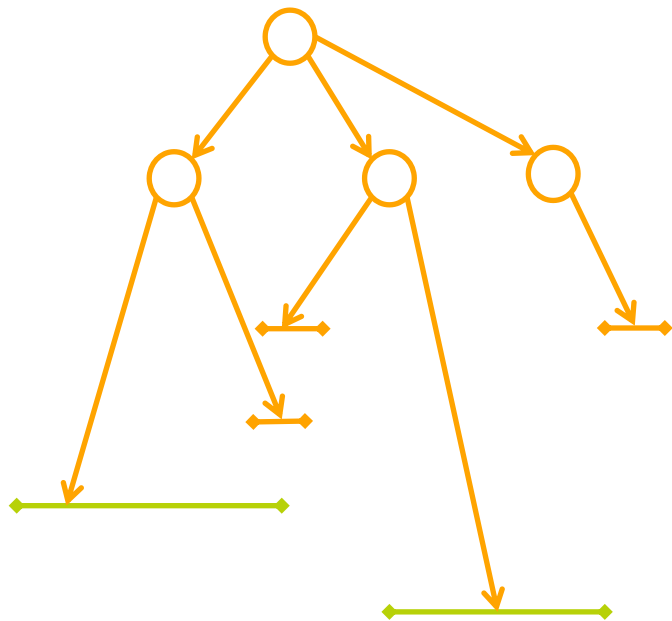intel.

# I/O Operation Flow

SCM

NVMe

## DAOS Xstream

- pmemobj_reserve() new buffer
- Start RDMA transfer to newly allocated buffer
  - Switch to other ULT until completion is reported
- If RDMA transfer failed, free buffer with pmemobj_cancel()
- Otherwise, start pmemobj transaction
- Modify index to insert new extent
- pmemobj_tx_publish()

intel.
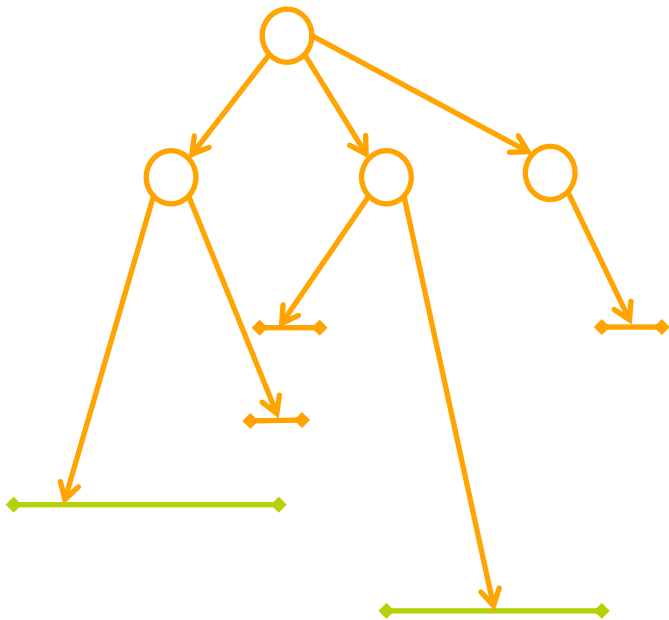
# I/O Operation Flow



Legend:
- SCM (orange arrow)
- NVMe (green arrow)

## DAOS Xstream

- pmemobj_reserve() new buffer
- Start RDMA transfer to newly allocated buffer
  - Switch to other ULT until completion is reported
- If RDMA transfer failed, free buffer with pmemobj_cancel()
- Otherwise, start pmemobj transaction
- Modify index to insert new extent
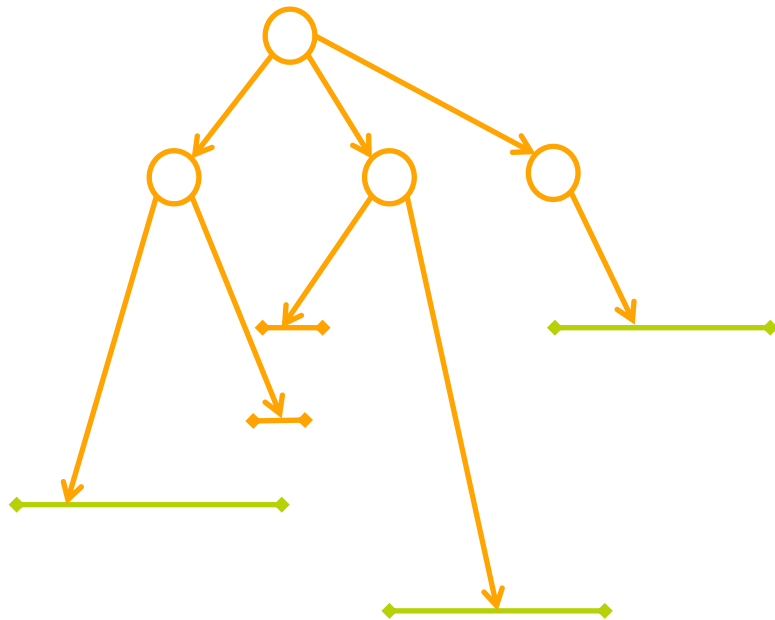- pmemobj_tx_publish()
- Commit pmemobj transaction

# I/O Operation Flow

SCM

NVMe

## DAOS Xstream

- pmemobj_reserve() new buffer

- Start RDMA transfer to newly allocated buffer
  - Switch to other ULT until completion is reported

- If RDMA transfer failed, free buffer with pmemobj_cancel()

- Otherwise, start pmemobj transaction

- Modify index to insert new extent

- pmemobj_tx_publish()
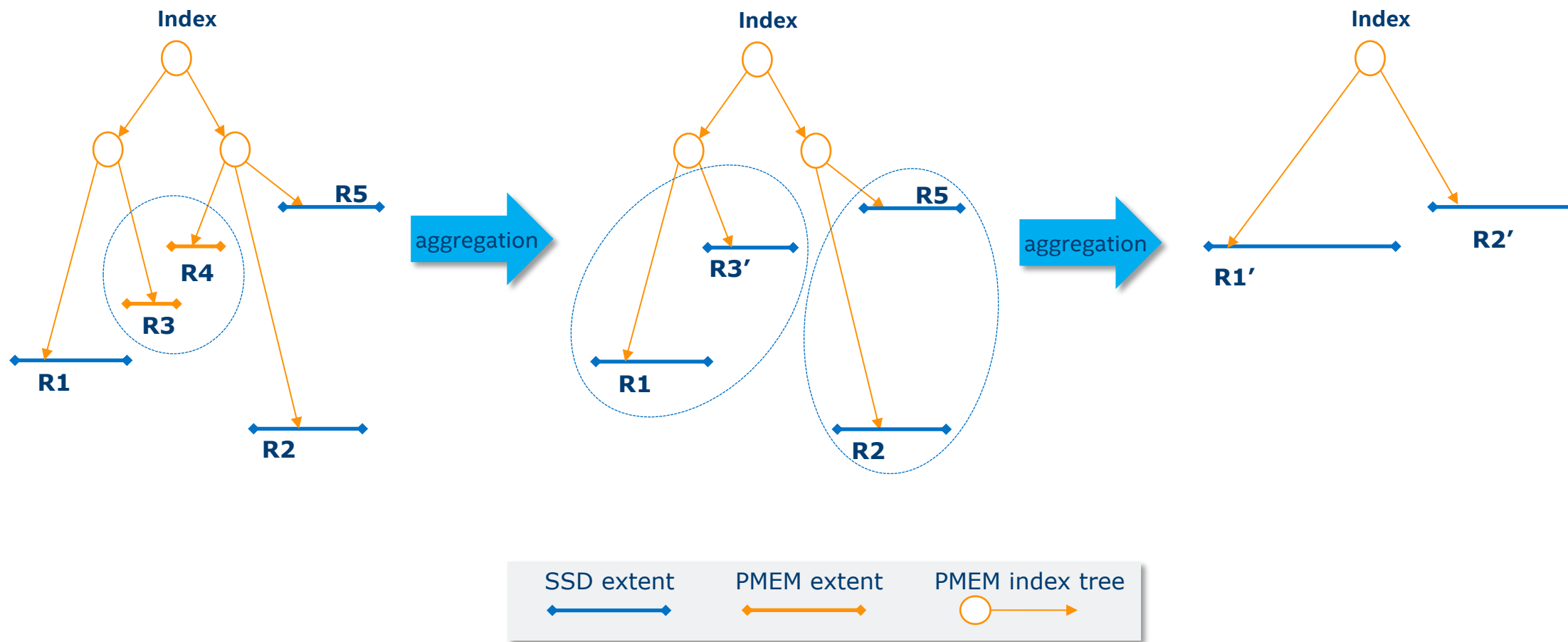
- Commit pmemobj transaction

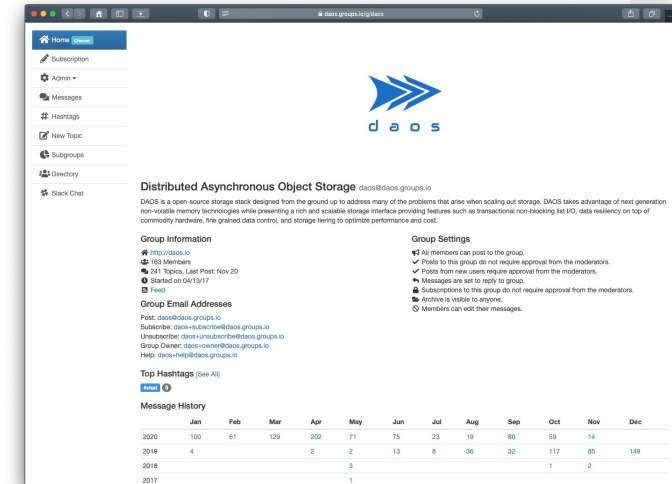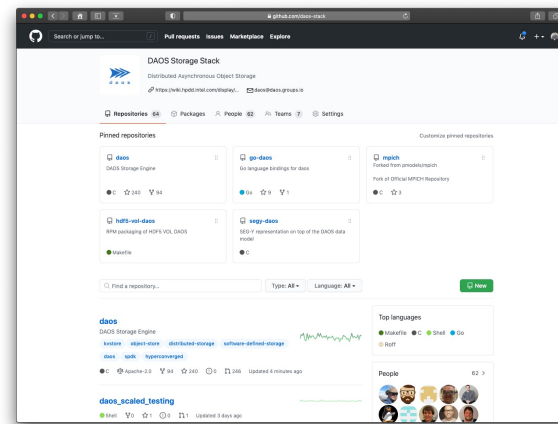- Reply to the client

intel.

# I/O Operation Flow

SCM
NVMe

## Extent in NVMe

- Track NVMe allocated/free space in data structures maintained in pmemobj pool
- Same processing flow

intel.

# Data Aggregation



SSD extent   PMEM extent   PMEM index tree

# Resources



- Open-source Community
  - Github: https://github.com/daos-stack/daos
  - Online doc: http://daos.io
  - Mailing list & slack: https://daos.groups.io
  - YouTube channel: http://video.daos.io
- 6th DAOS User Group (DUG'22)
  - Recordings will be available at http://dug.daos.io
- DAOS BoF Community at SC'22
- Intel landing page
  - https://www.intel.com/content/www/us/en/high-performance-computing/daos.html





intel.