

### インテルIPC最新技術アップデート

Technical Solution Specialist Fumiyasu Ishibashi

### パラダイムチェンジを牽引 - HPC & AI融合



HPC のワークフローに 組み込まれる AI HPC のシミュレーション を高速化する AI HPC のシミュレーション にとって代わる AI



### 第2世代 インテル<sup>®</sup> Xeon<sup>®</sup> スケーラブル・プロセッサー



http://www.intel.co.jp/xeonscalable/

### インテル<sup>®</sup> Xeon<sup>®</sup> プロセッサーさらなる発展

#### 唯一、コンバージェンス向けに最適化されたデータセンター CPU インテル® アドバンスト・ベクトル・エクステンション 512 インテル® ディープラーニング・ブースト (インテル® DL ブースト) インテル® Optane™ DC パーシステント・メモリー

2019年

#### 2020年

#### 2021年

<u>次世代テクノロジー</u>

### **COOPER LAKE**

CASCADE LAKE

14NM AI の新たな加速化 (VNNI) メモリーストレージの新しい階層 14NM 次世代インテル® DL ブースト (BFLOAT16)

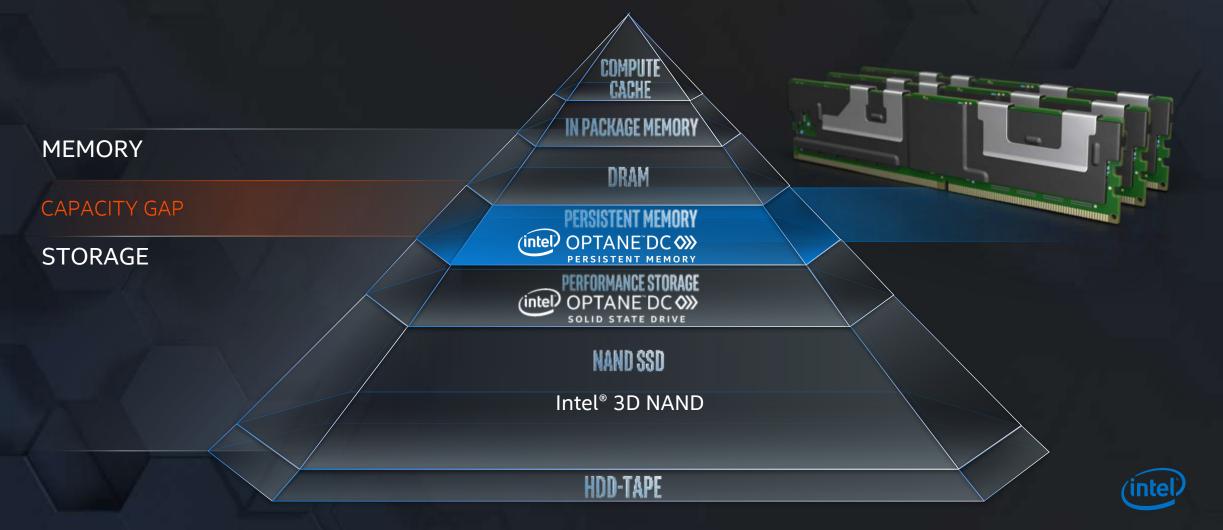
**ICE LAKE** 

10NM 現在サンプルを出荷開始

#### 業界最先端のパフォーマンス



### MEMORY AND STORAGE HIERARCHY GAPS THE CAPACITY GAP



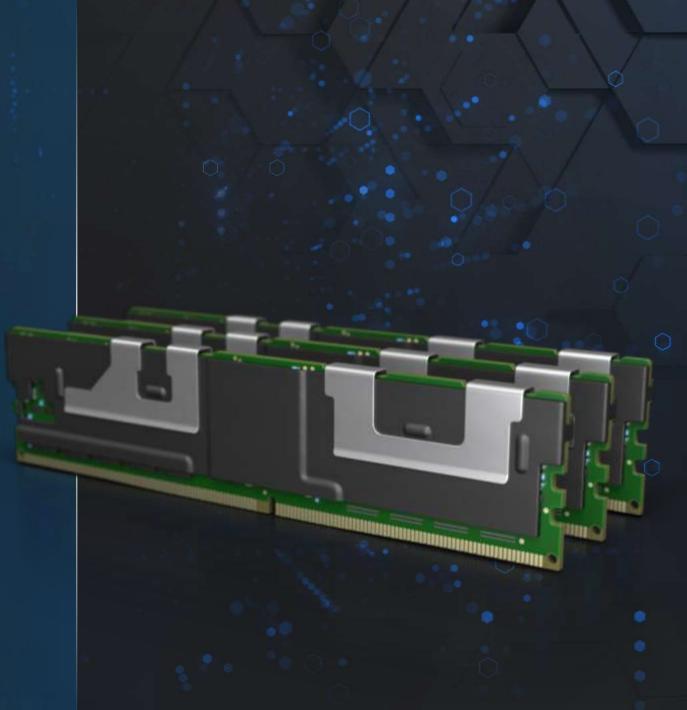


**BIG AND AFFORDABLE MEMORY** 128, 256, 512GB MODULES DDR4 PIN COMPATIBLE

**BYTE ADDRESSABLE** DIRECT LOAD/STORE ACCESS

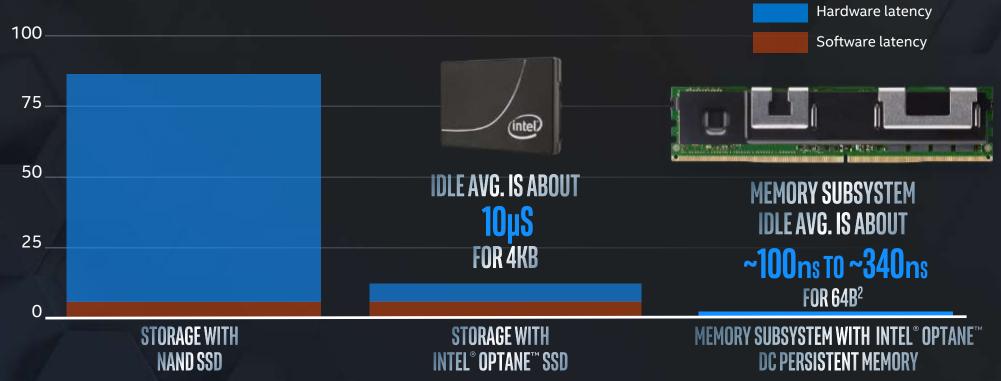
HIGH PERFORMANCE STORAGE NATIVE PERSISTENCE

HIGH RELIABILITY AND SECURITY TWO OPERATIONAL MODES



## MORE TO BE GAINED BY BEING ON MEMORY BUS

#### **IDLE AVERAGE RANDOM READ LATENCY**<sup>1</sup>



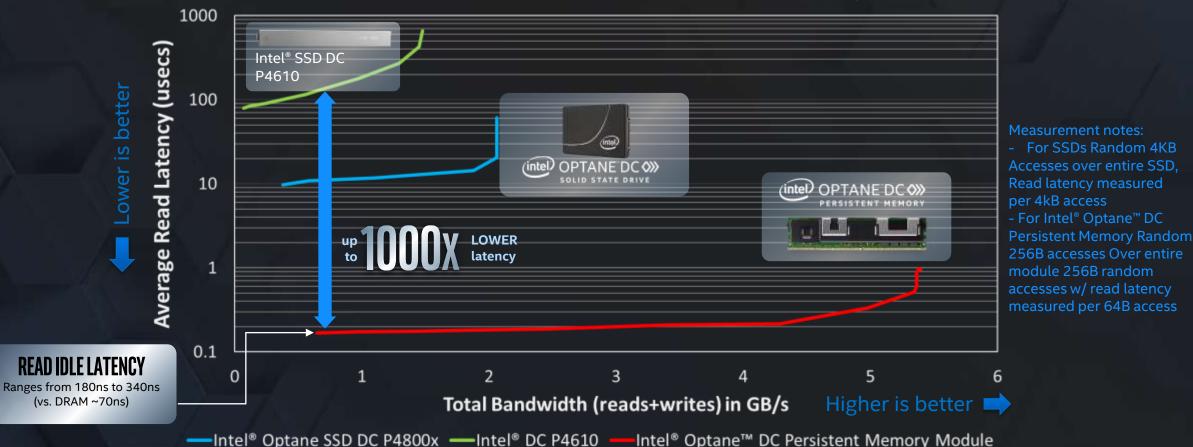


#### See Aopendix J

Performance results are based on testing as of July 24, 2018 set forth in the configurations and may not reflect all publicly available security updates. See configuration disclosure for details. No product can be absolutely secure. For more complete information about performance and benchmark results, visit www.intel.com/benchmarks.

#### **PERFORMANCE** Latency vs. Load

(70Read/30Write Random, 4kB for SSD and 256B for Memory)



#### See Appendix K

Performance results are based on testing as of February 22, 2019 set forth in the Configurations and may not reflect all publicly available security updates. See configuration disclosure for details. No product can be absolutely secure. For more complete information about performance and benchmark results, visit <u>www.intel.com/benchmarks</u>.

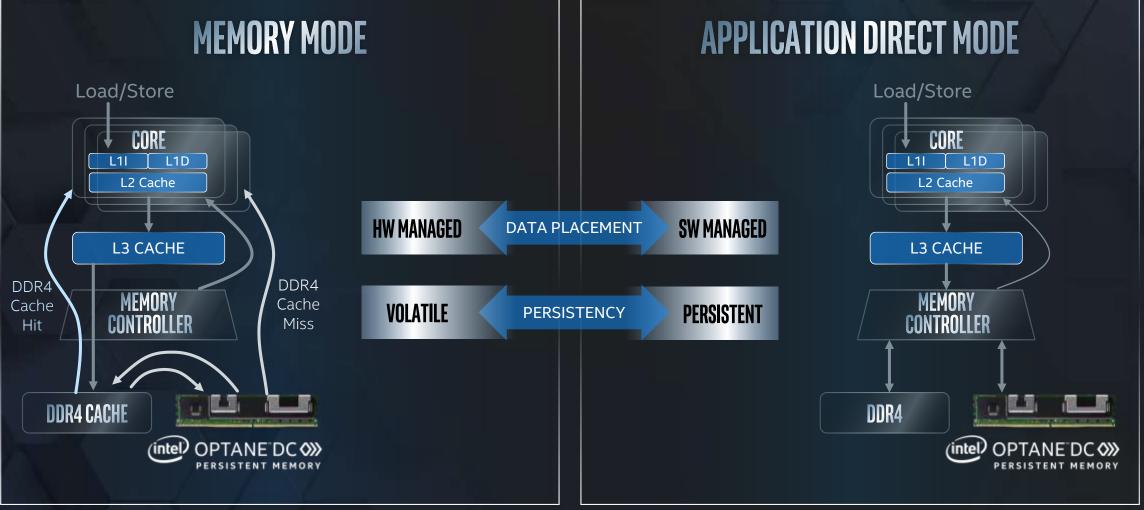
### **PERFORMANCE DETAILS**

- Intel<sup>®</sup> Optane<sup>™</sup> DC persistent memory is programmable for different power limits for power/performance optimization
  - 12W 18W, in 0.25 watt granularity for example: 12.25W, 14.75W, 18W
  - Higher power settings give best performance
- Performance varies based on traffic pattern
  - Contiguous 4 cacheline (256B) granularity vs. single random cacheline (64B) granularity
  - Read vs. writes

Granularity	Traffic	Module	Bandwidth
256B (4x64B)	Read		8.3 GB/s
256B (4x64B)	Write		3.0 GB/s
256B (4x64B)	2 Read/1 Write		5.4 GB/s
64B	Read	256GB, 18W	2.13 GB/s
64B	Write		0.73 GB/s
64B	2 Read/1 Write		1.35 GB/s



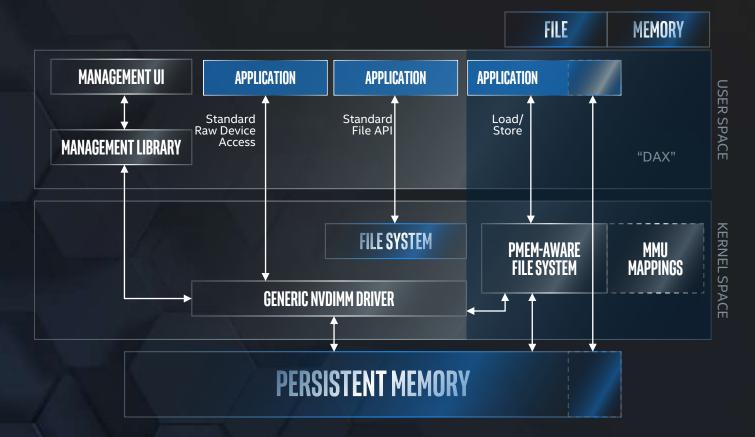
### **TWO OPERATIONAL MODES**





# **PERSISTENCY: EARLY ENGAGEMENT WITH INDUSTRY**

#### PERSISTENT MEMORY PROGRAMING MODEL developed through SNIA



#### PERSISTENT MEMORY DESIGN KIT (PMDK) available on <u>http://pmem.io</u>

a presenta his	ia a					- 0
-> C (	() Not secure (	name 20			9)6	9 6
		pme	em.io			
P	ersiste	nt Memo	ory Pr	ogram	ming	g
Home	Glossary	Documents	PMDK	ndcti i	Blog /	About
If you're j Area for li Guide, and Here are s • Persis • Work- book • Dersig • Intel I • PIRL C	ust getting star nks to backgrou d lots of additio some of the top tent Memory Di In-Progress Pro tent Memory Si Developer Zone conference (Per	for persistent men sistent Programmi	formation: ent Memory	ri rted Rec Hulti Intro vmer have desci close	ent Blog level vm level vm duction mcache wi recently ribed perfi- tibed perfi- to optimur all keys	encache hich we orms um when
The term, technologi memory, o non-volati that are li but it doe	persistent mem les which allow directly byte-ad le, preserved a ke memory, an sn't typically re	nt Memory? iory is used to design programs to access idressable, while the cross power cycless d aspects that are place either memo provise a their there.	ss data as he contents i i. It has aspe like storage, жу or storage	appro likely cts libver base	oximately y to minicachie d URU cac mig UT. CEU duction	equally

Instead, persistent memory is a third tier, used in

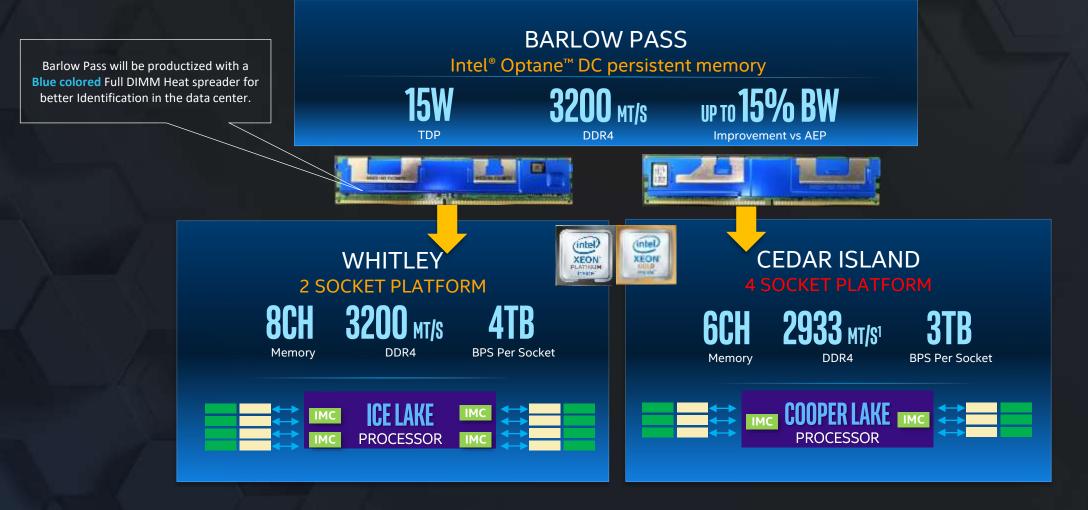


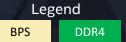
# **A STRONG MEMORY & STORAGE FUTURE**

TODAY	FUTURE		
2 <sup>ND</sup> GEN INTEL XEON SCALABLE (CASCADE LAKE)	COOPER LAKE / ICE LAKE	SAPPHIRE RAPIDS	FUTURE INTEL XEON PROCESSOR
OPTANE DC APACHE PASS	ARLOW PASS	3 <sup>RD</sup> GEN DC PERSISTENT MEMORY	4 <sup>TH</sup> GEN DC PERSISTENT MEMORY
Intel® SSD DC P4800X (COLDSTREAM)AL	LDER STREAM	NEXT GENERATION	NEXT GENERATION
	ALE-R/ARBORDALE (96-L, 144-L)	NEXT GENERATION	NEXT GENERATION



### BARLOW PASS OVERVIEW – 2<sup>ND</sup> GENERATION



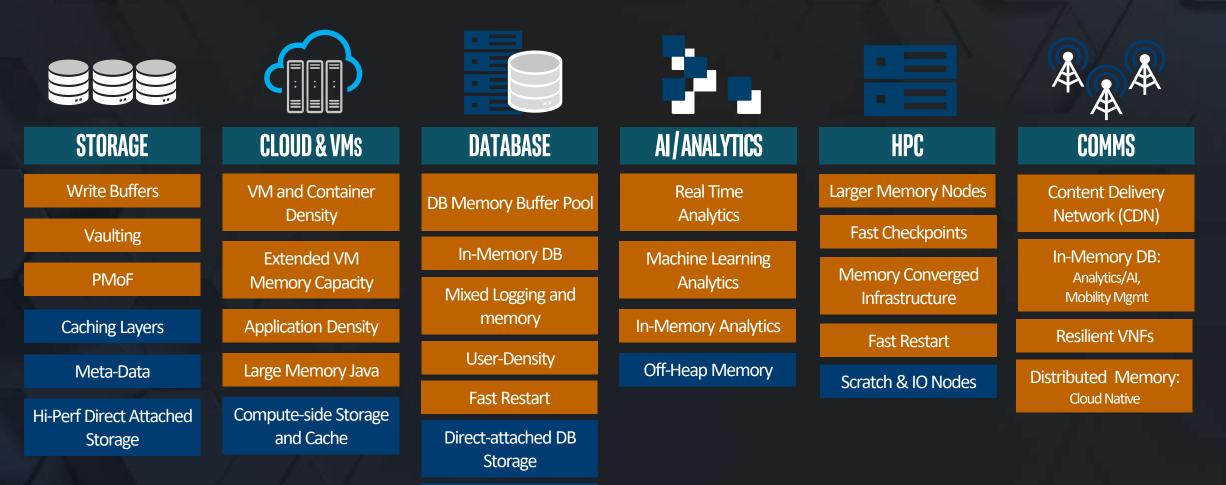


1. For Cedar Island, 2933 2DPC & 3200 1DPC is POR. Stretch target is 3200 2DPC



#### INTEL OPTANE DC DESIGN IN ADVOCACY

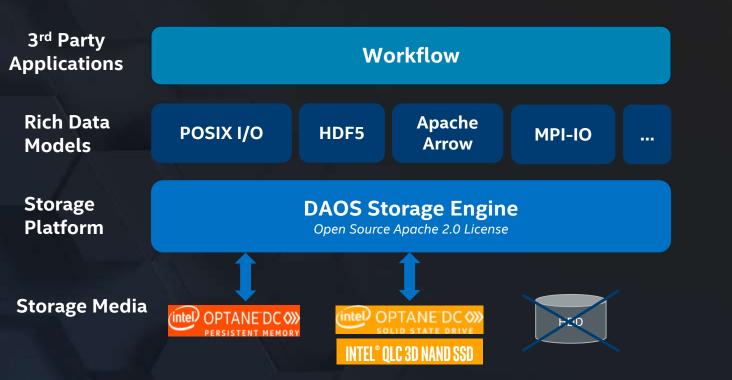




**Dedicated Logging** 

VNF – Virtual Network Function PMoF - Persistent Memory over Fabric

# <u>DISTRIBUTED ASYNCHRONOUS OBJECT STORAGE</u>



#### Benefits

- Built natively over new userspace PMEM/NVMe software stack
- Rich storage semantics
- High throughput/IOPS @arbitrary alignment/size
- Fine-grained, low-latency & True zero-copy I/Os
- Scalable communications
- Software-managed redundancy
- Rely on COTS hardware



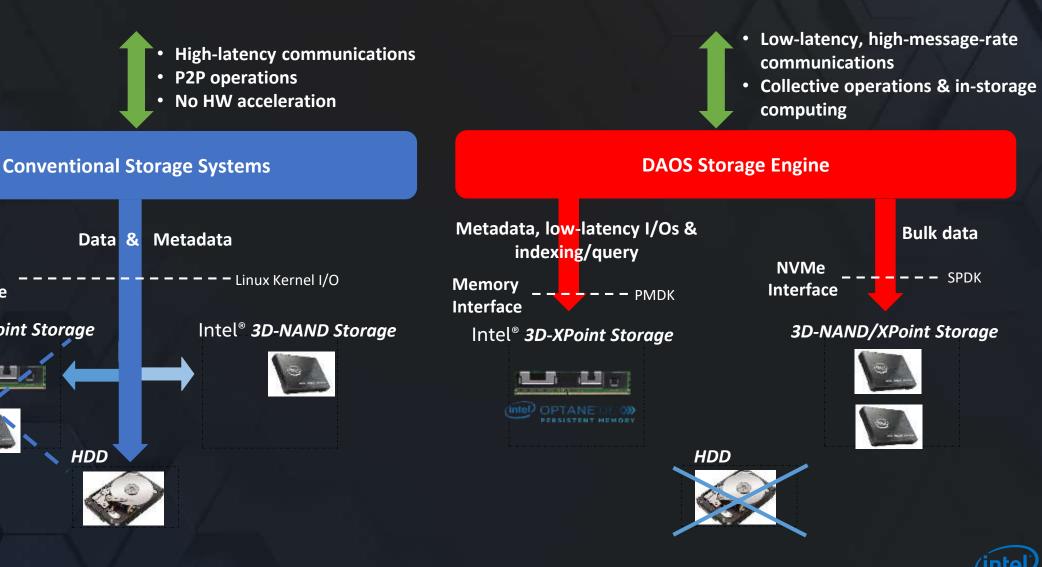
### **DAOS ARCHITECTURE**

Block

Interface

Intel<sup>®</sup> **3D-XPoint Storage** 

HDD



### **DEPLOYMENT OPTIONS**

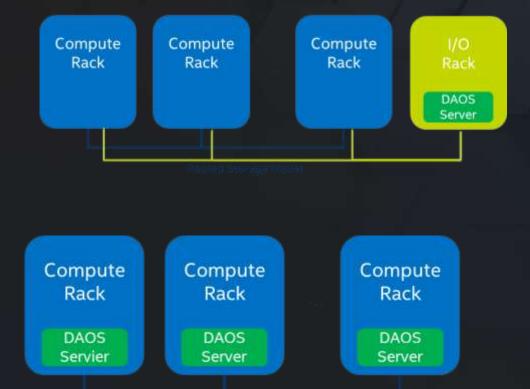
DAOS system can be deployed in two different ways:

#### **Pooled Storage Model :**

The DAOS servers can run on dedicated storage nodes in separate racks. This is a traditional pool model where storage is uniformly accessed by all compute nodes. In order to minimize the number of I/O racks and to optimize floor space, this approach usually requires high density storage servers.

#### **Disaggregated Storage Model :**

In the disaggregated model, the storage nodes are integrated into compute racks and can be either dedicated or shared (e.g. in a hyper-converged infrastructure) nodes. The DAOS servers are thus massively distributed and storage access is non-uniform and must take locality into account.

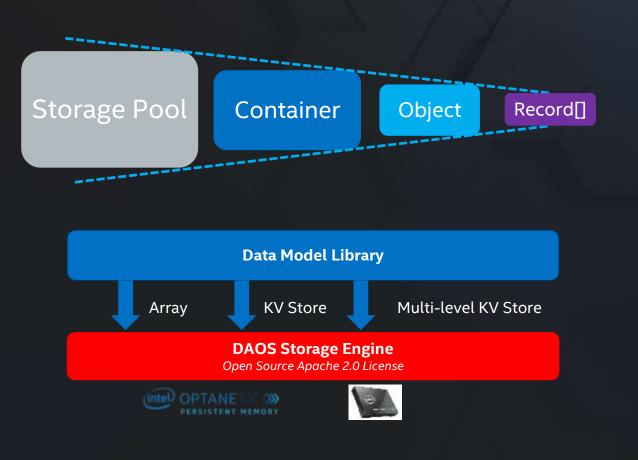


Disaggregated Storage Model



# **DAOS DATA MODEL**

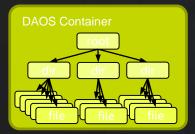
- Non-POSIX rich storage API as the new foundation
  - Scalable storage model suitable for both structured & unstructured data
    - key-value stores, multi-dimensional arrays, columnar databases, ...
    - Accelerate data analytic/AI frameworks
  - Non-blocking data & metadata operations
  - Extendable through microservice architecture





## **DATASET MANAGEMENT**

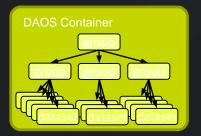
- Aggregate related datasets into manageable and coherent entities
  - Distributed consistency & automated recovery
  - Full Versioning
  - Simplified data management
    - Snapshot
    - Cross-tier Migration
    - Indexing





Encapsulated POSIX Namespace

File-per-process

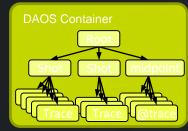


HDF5 « File »

Columnar Database



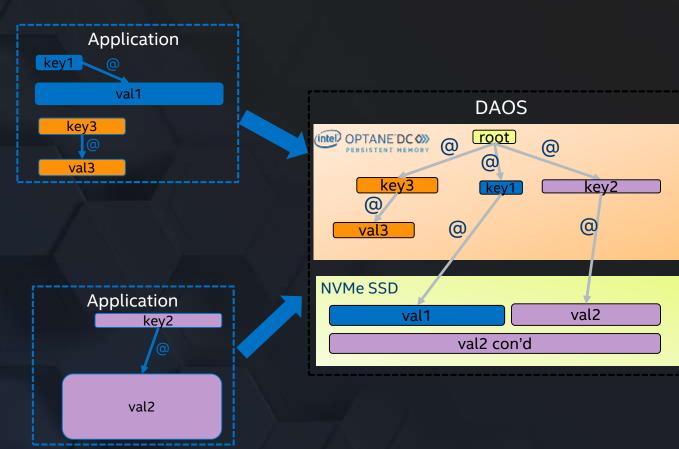
Key-value store



SEGY



## **ADVANCED STORAGE API**



#### Fast data retrieval

- Avoid file serialization and offset management
- Keys can be of any size/type
- Keys can be ordered with range query support Scalable Insert & Fetch
- Allow concurrent access/update
- Unconstrained by POSIX serialization
- Non-blocking
- Distributed transactions keep KV store always consistent

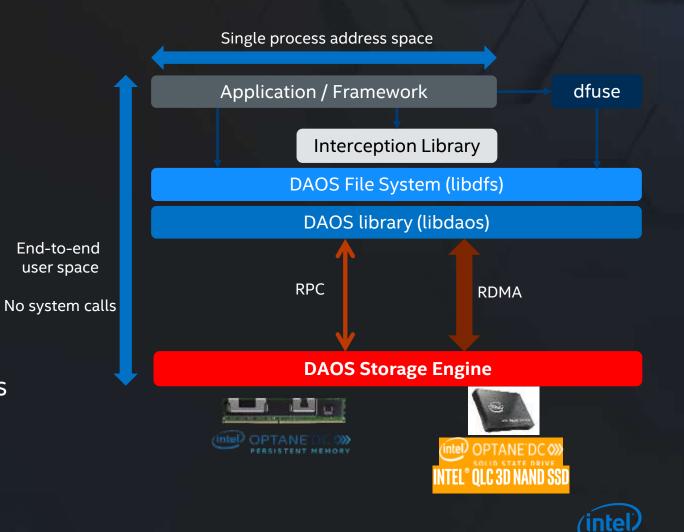
#### Data indexing

- Enable in-storage computing
- Query & custom index
- Data provenance



# **POSIX I/O SUPPORT**

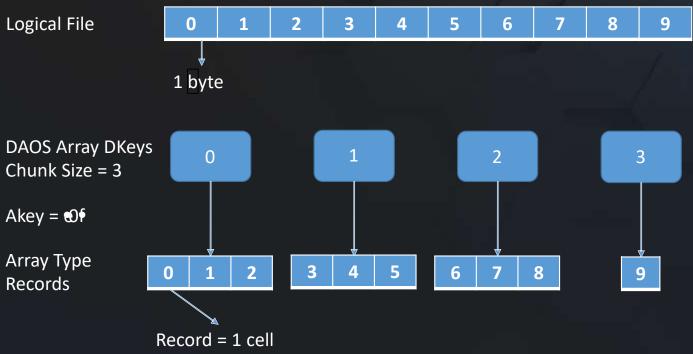
- DAOS File System (libdfs)
  - Encapsulated POSIX namespace
  - Application/framework can link directly with libdfs
    - ior/mdtest backend provided
    - MPI-IO driver leveraging collective open
    - TensorFlow, ...
- FUSE Daemon (dfuse)
  - Transparent access to DAOS
  - Involves system calls
- I/O interception library
  - OS bypass for read/write operations



## **MPI-IO DRIVER FOR DAOS**

The DAOS MPI-IO driver is implemented within the I/O library in MPICH (ROMIO).

- Added as an ADIO driver
- Portable to Open-MPI, Intel MPI, etc.
- <u>https://github.com/daos-stack/mpich</u>
- daos\_adio branch
- PR to mpich master in review
- 1 MPI File = 1 DAOS Array Object



Application works seamlessly by just specifying the use of the driver by appending "daos:" to the path.



## DAOS IN 10-500 2019-11

#### #1 in 10 Node Challenge

										D-Re	igister 🔒 Log
7	4	Virtual	Institut	e for I/O				10	Search		C
4		4								Recent C	hanges Sitema
are	here: Virt	tual Institute for I/O » I	O500 » Lists » :	2019-11 <b>» 10 Node</b> (	Challenge						
Ra	nked List	10 Node Challen	ge Full L	List Histor	ic data St	udent Clu	ister Co.	-			
This	is the offic	Challenge Ial list from Supercom				s the best	result for		C	)	00
This	is the offic		in/filesystem qual			s the best	result for		C	10500	00
This a giv	is the officience of the offic	al list from Supercom	in/filesystem qual	lifying for the 10 Node	Challenge.	client	client	data	Score		md
This a giv	is the officience of the offic	al list from Supercom ation of system/institutio	in/filesystem qual	information	Challenge					10500	
This a giv	is the officience of the offic	al list from Supercom ation of system/institutio	in/filesystem qual	information	Challenge.	client	client total			io500 bw	md
This a giv	is the officience of the offic	al list from Supercom ation of system/institution	system	illying for the 10 Node information storage vendor	Challenge filesystem type	client nodes	client total procs	data	score	lo500 bw GiB/s	md kIOP/s

#### #2 in Ranked List

0
-
cent Changes Silemag

This is the official list from Wapercomputing 2019. The list shows the best result for a given combination of system/institution/filesystem.

#### Please see also the 10 node challenge ranked list.

#	information								10500		
	list	institution	system	storage vendor	filesystem	client	client	data	score	DW	md
	id				type	nodes	total procs			GiB/s	kIOP/s
T.	9019	WekatO	WekałO on AWS	WekatO	WekaiO Matrix	345	8625	zip	838.95	174.74	5045.33
2	sc19	inter	Wolf	Intel	DAOS	- 26	728	zip	933.64	183.36	4753.79
	3413	Supercomputing Center in Changsha	Teams, 20	of Defense Technology	Cugoo		0200	sala	400.00	200.00	



## **DAOS: PRIMARY STORAGE ON AURORA**



#### **Aurora DAOS configuration**

- Capacity: 230PB
- Bandwidth: >25TB/s

"Combined in Aurora, the Intel compute system, Cray Slingshot network, and the Intel DAOS storage open new possibilities for accelerating the scientific research needed to solve critical human challenges such as cancer and disease. DAOS enables the creation of new storage data models tailored specifically to applications like the Cancer Distributed Learning Environment (CANDLE) which provide a powerful platform to advance a wide array of scientific challenges using deep learning."

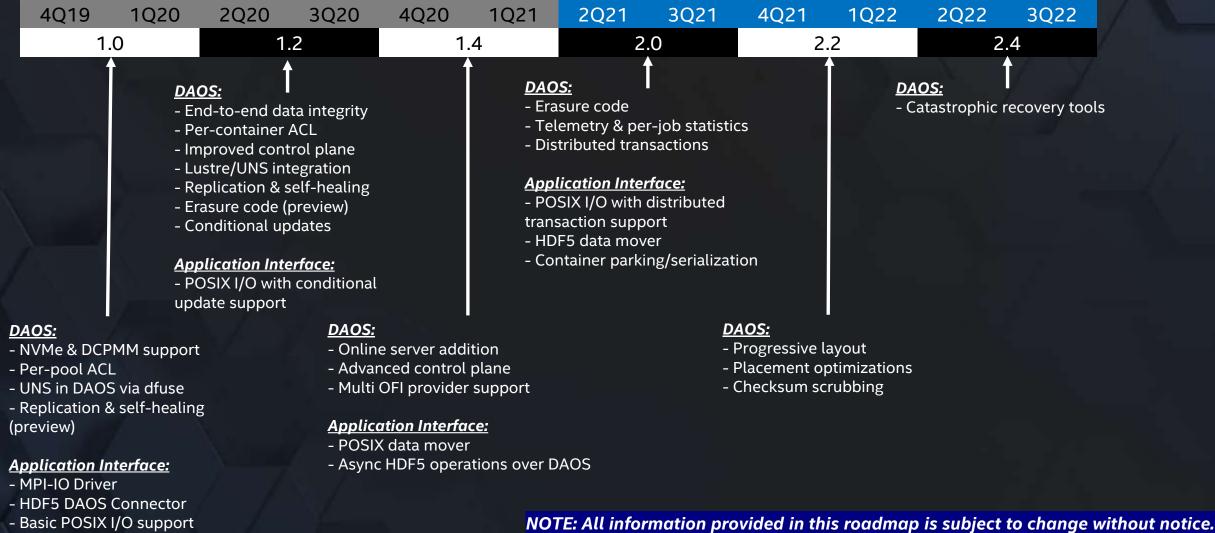
– Rick Stevens, Associate Laboratory Director for Computing, Environment and Life Sciences

"The Argonne Leadership Computing Facility is excited to be the first major production deployment of the DAOS storage system as part of Aurora, an US exascale system coming in 2021. As designed, it will provide us unprecedented levels of metadata operation rates and extremely high bandwidth for I/O intensive workloads."

– Susan Coghlan, ALCF-X Project Director/Exascale Computing Systems Deputy Director



### **DAOS COMMUNITY ROADMAP**



- Spark

# RESOURCES

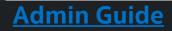
Why Githuh? - Dr	erptis Toplare – Marke	quines Michig -		- 111 s	yn in Sigri wi			
U date-stack/ date			Ø thesh	#3a 1	Viet 10			
43-Grafe (2) Income (K. 17)	Pott response AN I II Property	and the second states	SHC 2					
ADS Storage Engine Mite	Weenle							
	and service in supervision and	ter all services	Allow phat stop	Server and	200			
@-2.441 commits	p-02 branches	O B referation	AL 44 contributions		lyster 20			
Real and Post of Street	-			nets				
VL		Print's Sal		-	tillion i loss age			
	(DBC-28) mild (Doud	product services (#11.(#11)	eths.		J miriti apr			
and the state of the second se					to function approximately and the function of the second s			
in survivor a testint	20 Percei age							
Birth .	Data Still an abit, he	when the last proph (#1107)		( his ope				
	deally been an about a	the local sector of the local						

**Source Code on GitHub** 



#### Introduction The Disbellianed Assocharies Object Marige (DACK) a structure access regists were designed from the proved up for muscline/valid/factor Non-Valid in Mercory (NVM). DADS bries, showing early next generation /4V/Mastrolingy, Westcorage Class Memory (SDM) and NVM express (NVM)4. while presenting a key value computed or one cap of commuting further or that provides Halterscaultus, transcherative-blocking RD, adversidate protectiet with self-brailing, traitioend data megnity, fine-grained data control, and south interrupy, to sythmise performance and scart. This administration gasks seraises is emorized with DAOS v0.7. Additional Documentation Helice to the following documentation for excitine clare and description: Document Location DAY Durrait https:/gtts/sometikee.epsch/dam/historyame/eps/HEATME.ed IMOS Scorage Myeler Hitse Within Little and Characteria and Stateman randoms. Print, Solid patt intercommunity, DC Baselines Compare heatings

Lift testhere-skie Mater



Dress & Administration Gabbie introduction



#### **DAOS Solution Brief**

Community mailing list on Groups.io

daos@daos.groups.io

#### Support

https://jira.hpdd.intel.com



### **LEGAL DISCLAIMER**

All information provided here is subject to change without notice.

The products described in this document may contain design defects or errors known as errata which may cause the product to deviate from published specifications. Current characterized errata are available on request.

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 computer system can be absolutely secure. Check with your system manufacturer or retailer or learn more at intel.com.

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 <u>www.intel.com/benchmarks</u>.

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 data. You should review this content, consult other sources, and confirm whether referenced data are accurate.

Results have been estimated or simulated using internal Intel analysis or architecture simulation or modeling, and provided to you for informational purposes. Any differences in your system hardware, software or configuration may affect your actual performance.

© Intel Corporation. Intel, the Intel logo, and other Intel marks are trademarks of Intel Corporation or its subsidiaries. Other names and brands may be claimed as the property of others.



