# Database storage management with object-based storage devices

Sami Iren
Seagate Research

Steve Schlosser
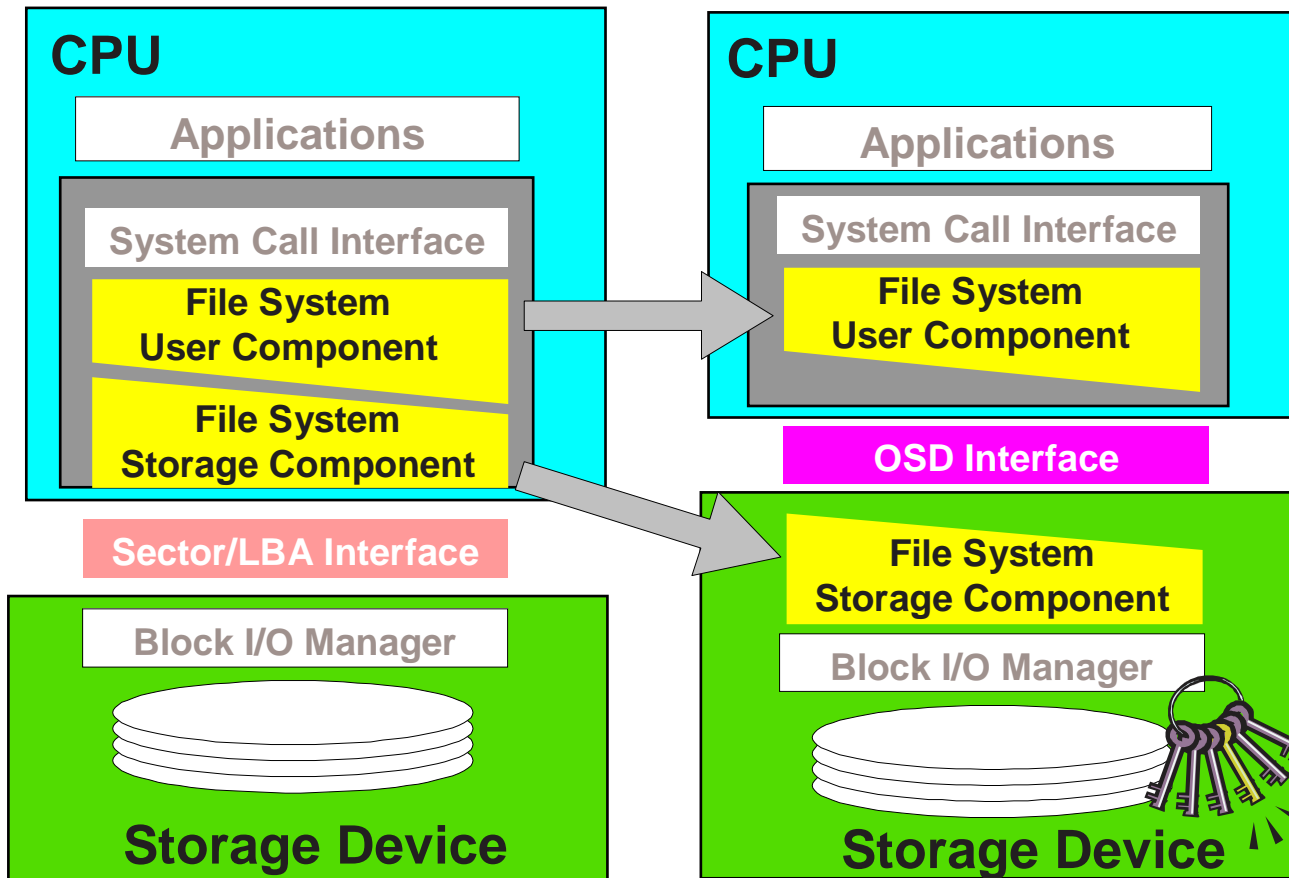Intel Research Pittsburgh

Seagate

intel.

# Outline

- New hardware: Object-based Storage Devices (OSD)
- Vision: Database-aware storage systems
- Example: Transparent device-specific data placement
- Moving forward: Issues with the current OSD specification

Seagate

intel.

# New hardware: Object-based Storage Devices (OSD)

- Same hardware, new interface ...
- Remember Network Attached Secure Disks?



- Disk array/server subsystem
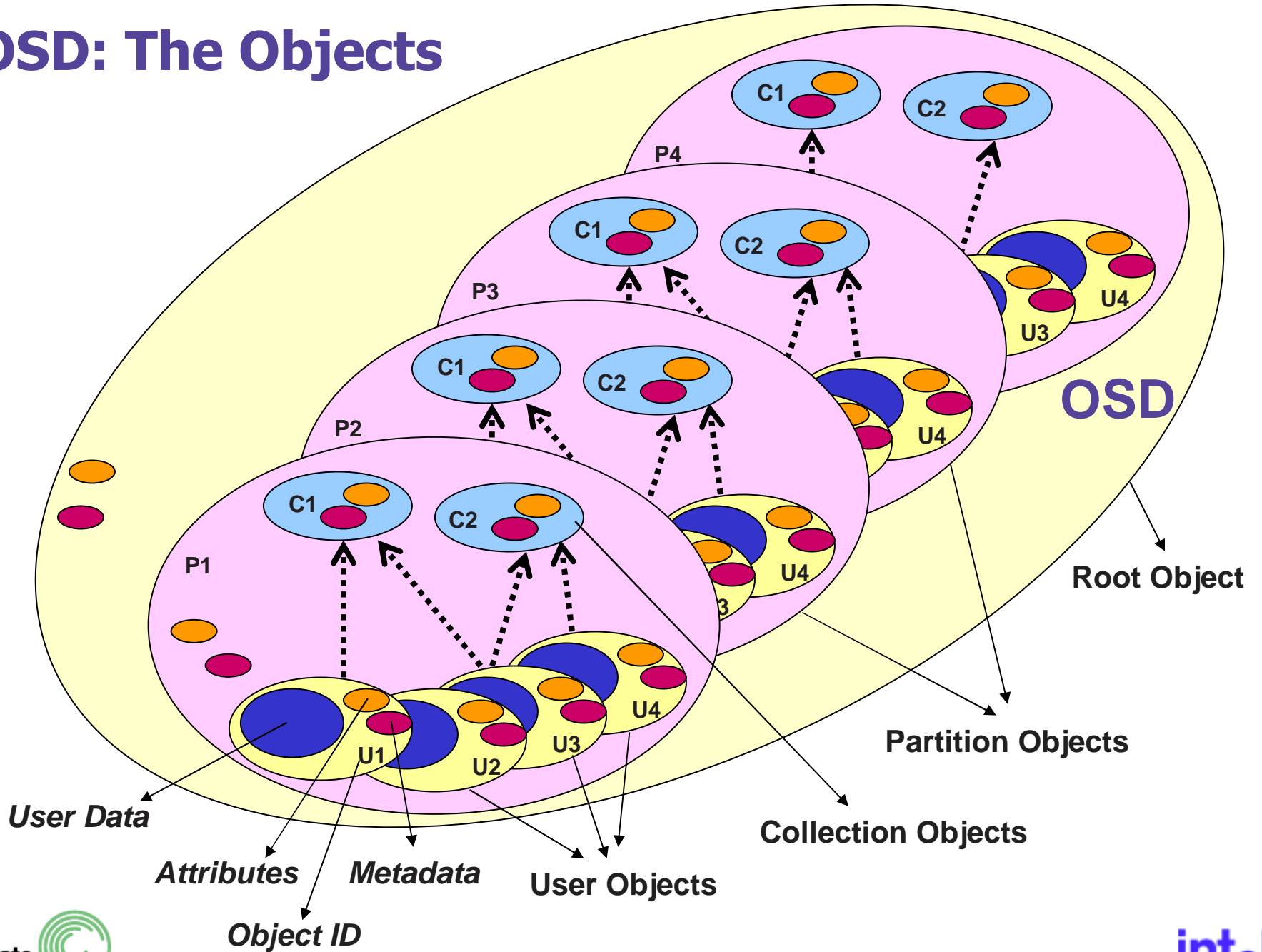- I.e. LLNL units with Lustre

**CPU**

- Applications
- System Call Interface
  - File System User Component
  - File System Storage Component

**Sector/LBA Interface**

**Storage Device**
- Block I/O Manager

**CPU**

- Applications
- System Call Interface
  - File System User Component

**OSD Interface**

**Storage Device**
- File System Storage Component
- Block I/O Manager

- "Smart" disk for objects
- I.e. Panasas storage blade

- Highly integrated, single disk
- I.e. prototype Seagate OSD

Seagate

intel.

# OSD: The Objects

OSD

Root Object

Partition Objects

Collection Objects

User Objects

Metadata

Object ID

Attributes

User Data

P1 P2 P3 P4

C1 C2

U1 U2 U3 U4

Seagate  intel.

# OSD: The Command Set

- **Basic Protocol**
  - READ
  - WRITE
  - CREATE
  - REMOVE
  - GET ATTR
  - SET ATTR

  **Very Basic**

  **Space Mgmt**

  **Attributes**
  - opaque
  - internal
  - shared

- **Specialized**
  - APPEND – write w/o offset
  - CREATE & WRITE – save msg
  - FLUSH OBJ – force to media
  - LIST – recovery of objects

- **Security**
  - Authorization – on each request
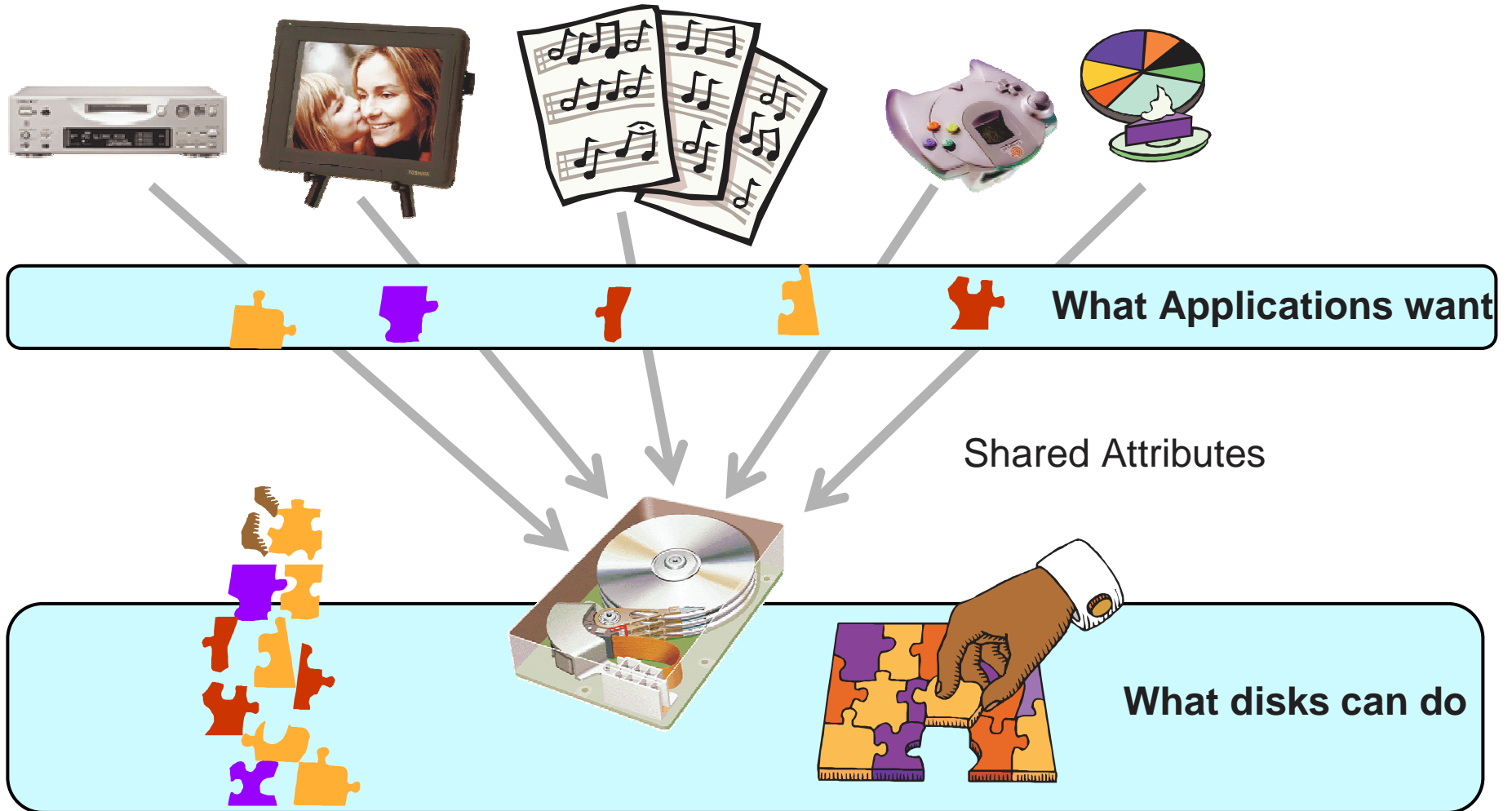  - Integrity – for args & data
  - SET KEY
  - SET MASTER KEY

- **Groups**
  - CREATE COLLECTION
  - REMOVE COLLECTION
  - LIST COLLECTION

- **Management**
  - FORMAT OSD
  - CREATE PARTITION
  - REMOVE PARTITION

Seagate

intel.

# OSD: Shared attributes

- Mechanism to push application information into storage

**What Applications want**

Shared Attributes

**What disks can do**
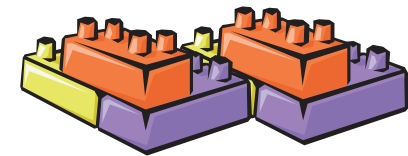
**Traditional**

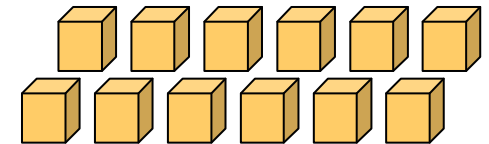**OSD**

# OSD: Benefits

- Improved performance
  - Hints, QoS, Differentiated Services
  - Data can be differentiated at the device
- Improved device and data sharing
  - Platform-dependent metadata moved to device
  - Systems need only agree on naming
- Improved scalability & security
  - Devices directly handle client requests
  - Object security w/ application-level granularity
  - Finer granularity than LUN-based security
- Improved storage management
  - Self-managed, policy-driven storage
  - Storage devices become more autonomous
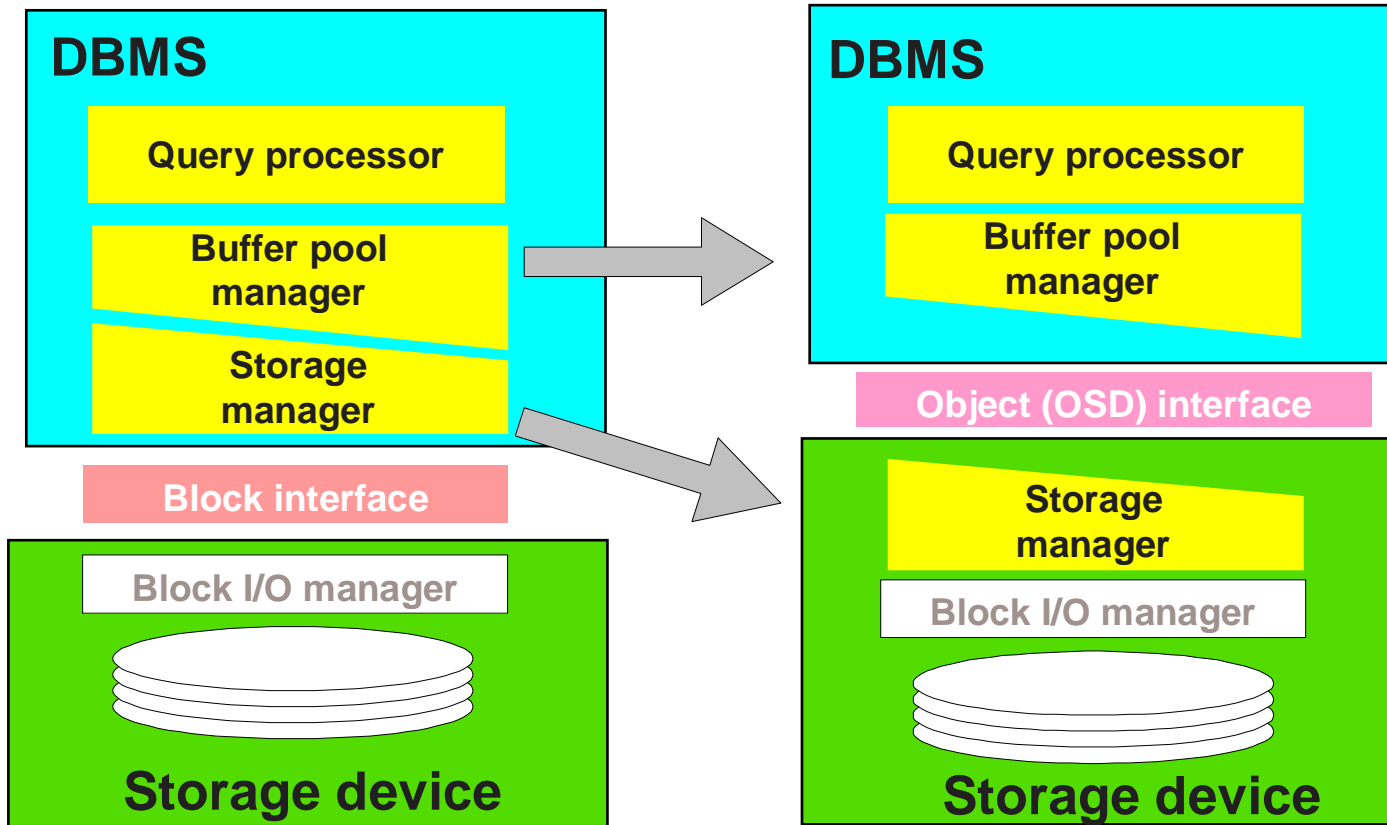
**Volumes**

**Objects**

**Blocks**

# Outline

- New hardware: Object-based Storage Devices (OSD)
- Vision: Database-aware storage systems
- Example: Transparent device-specific data placement
- Moving forward: Issues with the current OSD specification
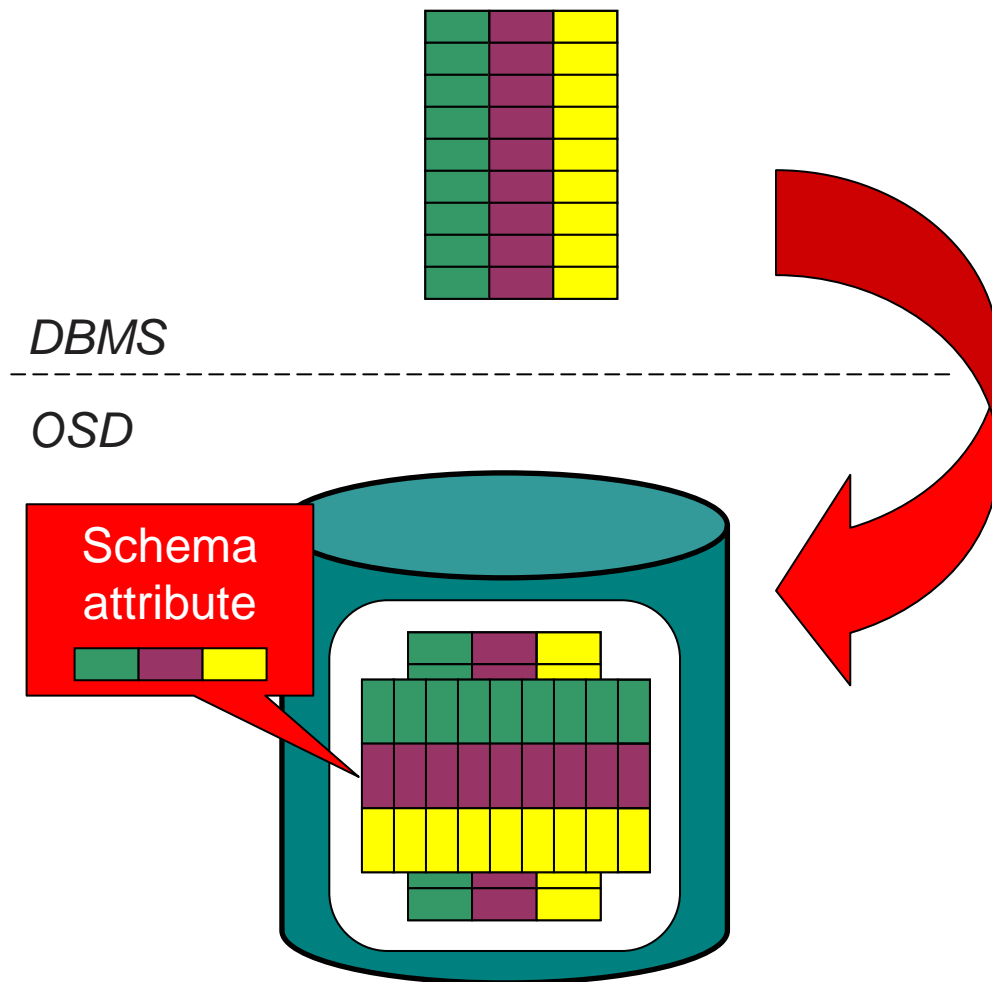
Seagate

intel.

# Vision: Storage that cooperates with databases

- DBMS storage managers do a lot to optimize storage access, with insufficient information

- Storage subsystem front-ends do a lot to optimize under the covers, without enough information about the application

- Wouldn't it be better if we could just get along?

- OSD can provide this mechanism

- **We have a unique opportunity to make DB software and storage more cooperative**
  - If you remember nothing else from this talk...

Seagate

intel.
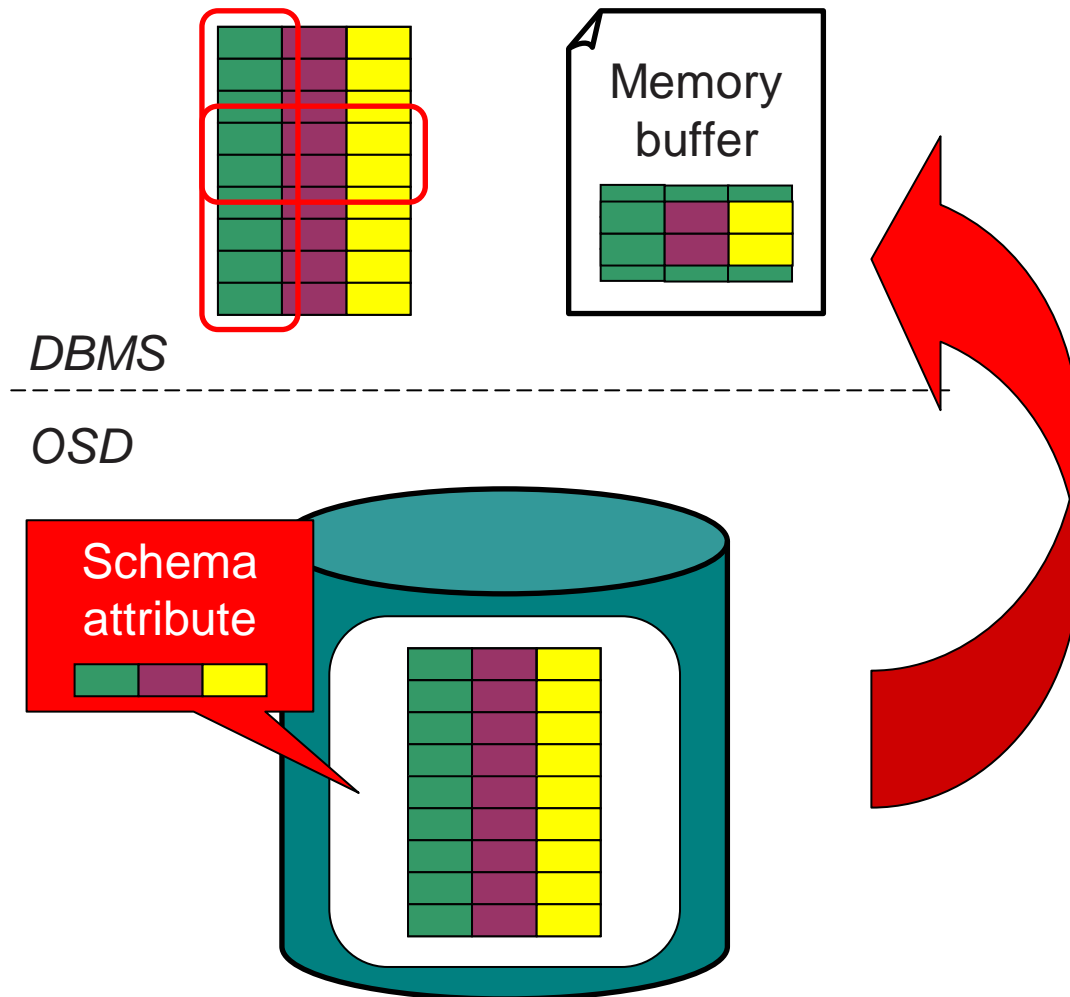
# Solution: OSD for databases

# Database-aware storage systems

1. Create object for relation
2. Attach schema attribute
3. Populate relation

- Goal 1: OSD uses its parameters to optimize layout under the covers

- Goal 2: Access method should be independent of storage optimizations

*DBMS*

*OSD*

Schema attribute

# Database-aware storage systems



Memory buffer

DBMS

OSD

Schema attribute

- Goal 2: Access method should be independent of storage optimizations

- DBMS should just specify the data it wants, and the destination address

- DBMS shouldn't worry about
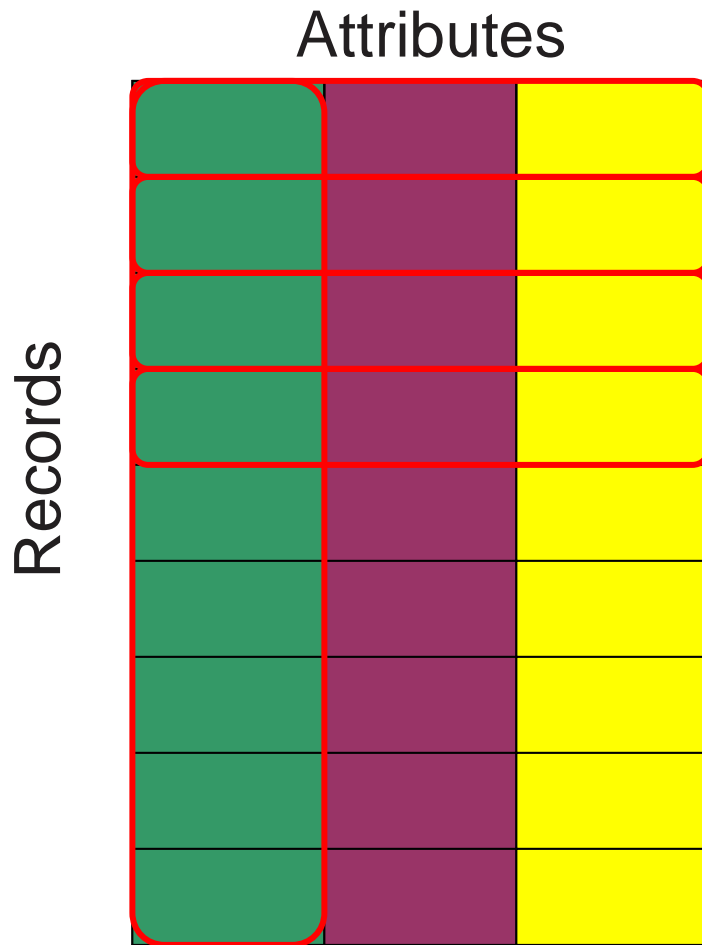  - Data placement
  - Addressing blocks/bytes

Seagate

intel.

# Outline

- New hardware: Object-based Storage Devices (OSD)
- Vision: Database-aware storage systems
- **Example: Transparent device-specific data placement**
- Moving forward: Issues with the current OSD specification
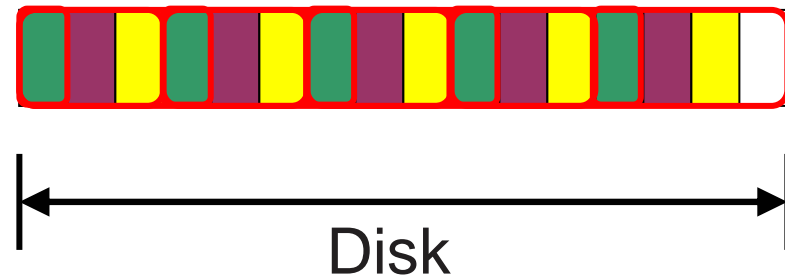
Seagate

intel.

# The Fates project at CMU

- Collaboration between Database Group and Parallel Data Lab
  - Ongoing involvement of researchers at Intel Research and EMC
  - Lachesis: Improved communication between storage subsystem and DB storage manager (VLDB 2003)
  - Clotho: Retool database software to fetch query-specific data (VLDB 2004)
  - Atropos: Leverage detailed disk information for improved 2D data placement (FAST 2004)

- Key idea: Storage informs database of its characteristics
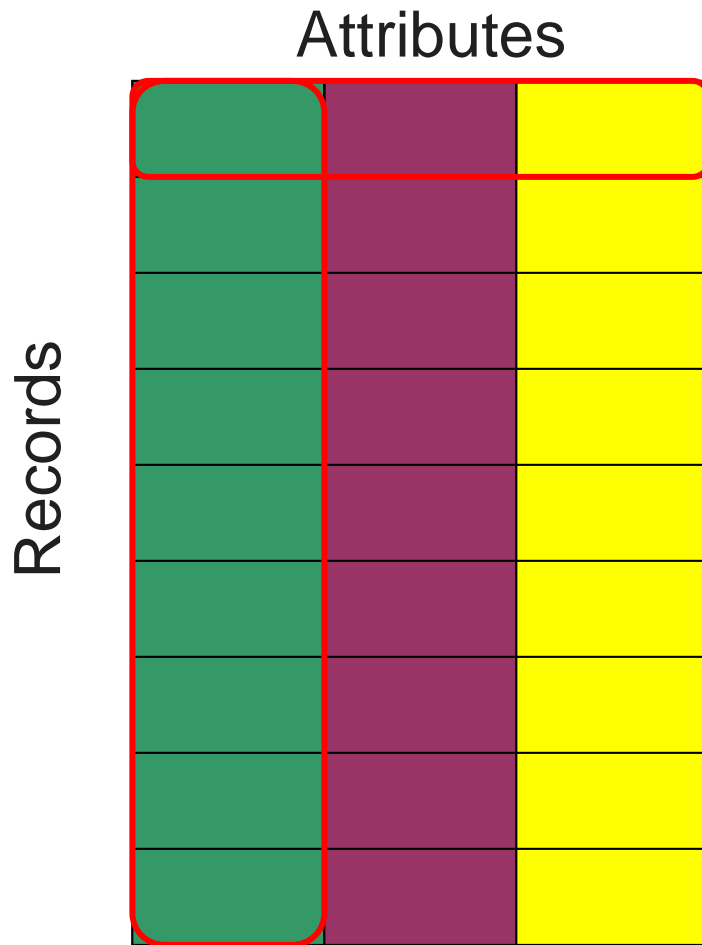  - Combination of schema and disk parameters yields 2D placement

**Seagate**

**intel.**

# 2D data structure access
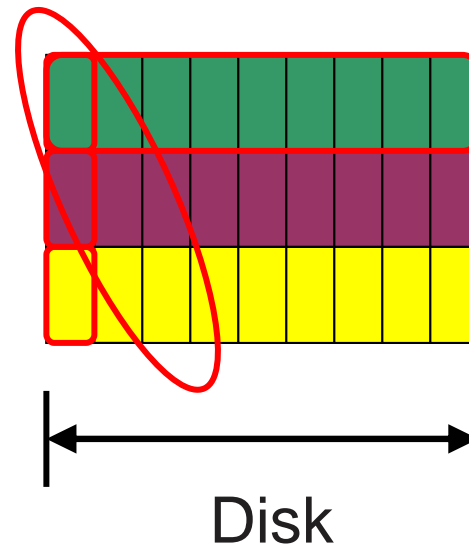
Attributes

Records



- On-disk storage requires serialization
- Access along one dimension is efficient
  - i.e., sequential
- Access along the other is inefficient
  - i.e., random I/O
  - Or, read entire relation and discard unwanted attributes

Disk

Seagate

intel.

# Intelligent data placement in Fates
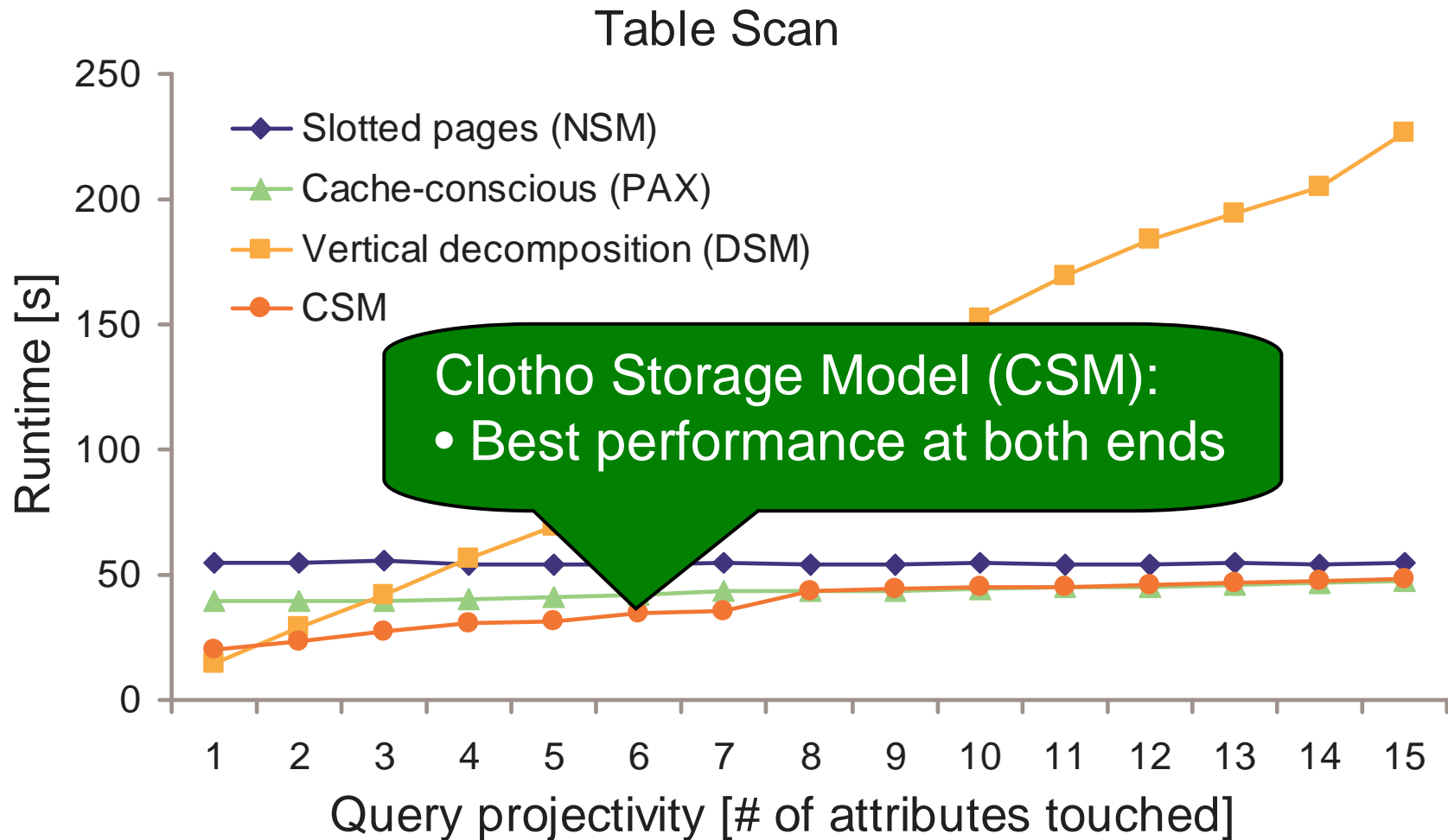
## Attributes



Records

- One column per disk track
- Column-major access is efficient
- Row-major access is inefficient
  - One block per rotation
- Semi-sequential layout enables efficient row-major access

Disk

Seagate

intel.

# High-level Fates result

Table: CREATE TABLE R (FLOAT a1, …, FLOAT a15) (1GB)
Query: SELECT a1, a2, …, FROM R WHERE a1 < Hi



Table Scan

Legend:
- ◆ Slotted pages (NSM)
- ▲ Cache-conscious (PAX)
- ■ Vertical decomposition (DSM)
- ● CSM

Clotho Storage Model (CSM):
- Best performance at both ends

Y-axis: Runtime [s]
X-axis: Query projectivity [# of attributes touched]
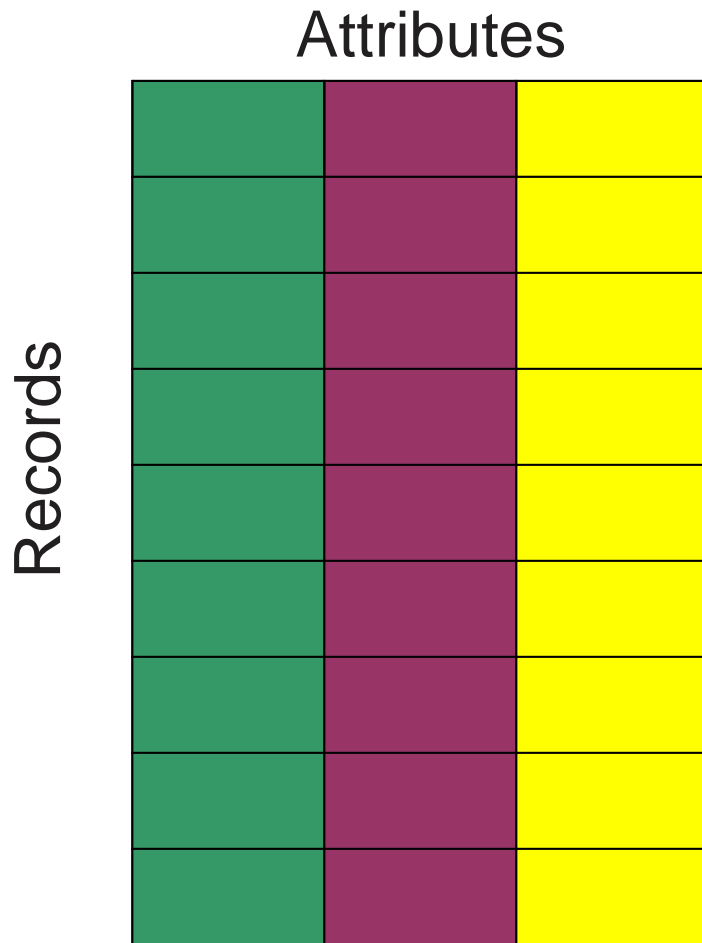
Seagate      intel.

# Solving the "parameter problem"

- Relying on storage vendors to expose parameters is problematic
  - Measuring parameters is difficult and fragile (but not impossible)

- Key idea: OSD can solve this problem
  - Expose schema to storage subsystem via shared attributes
  - Storage subsystem can do data placement under the covers
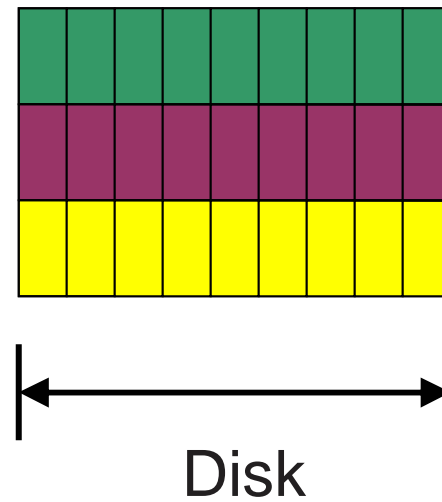
Seagate

intel.

# Outline

- New hardware: Object-based Storage Devices (OSD)
- Vision: Database-aware storage systems
- Example: Transparent device-specific data placement
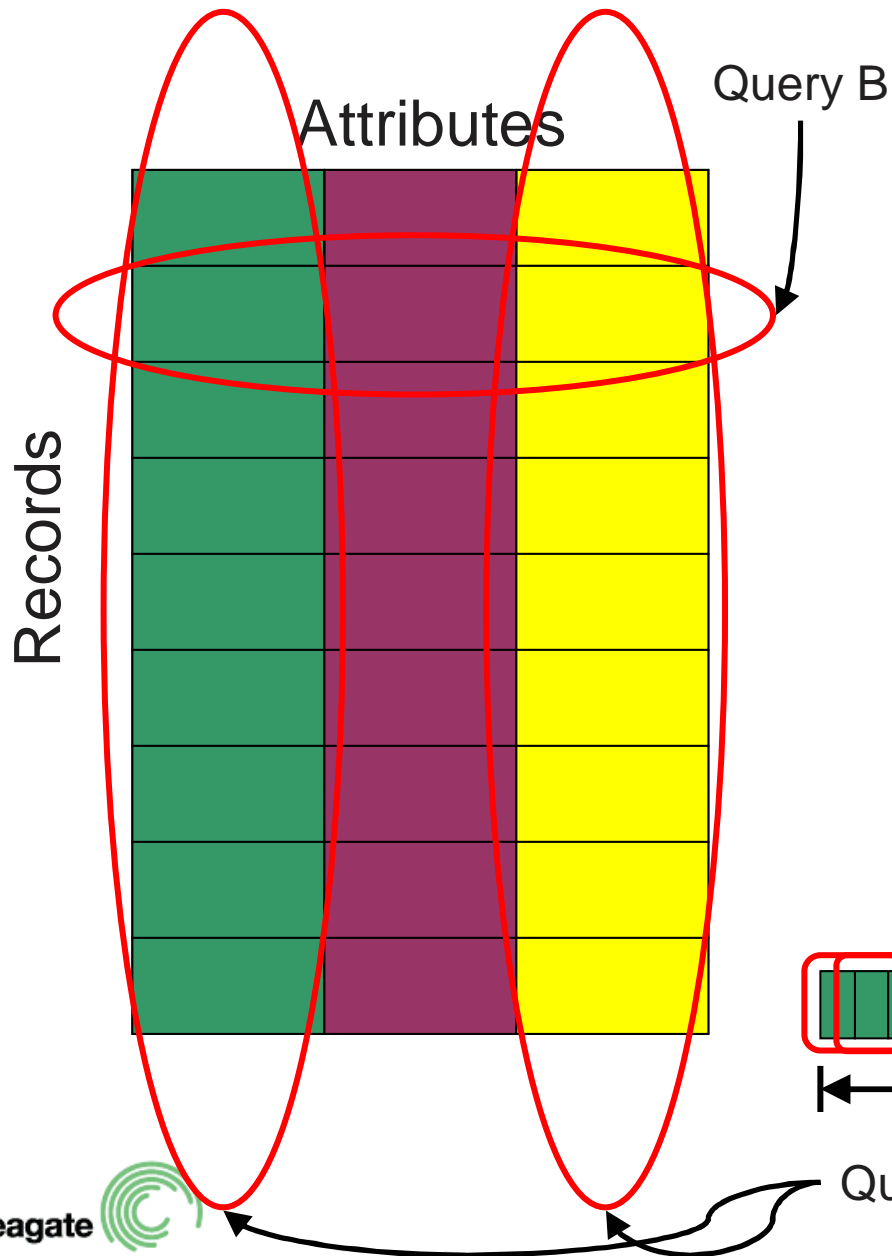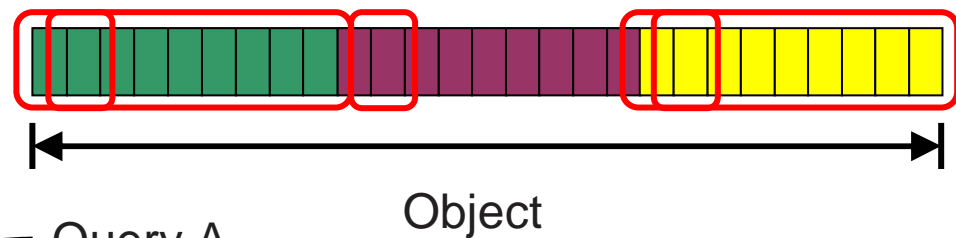- **Moving forward: Issues with the current OSD specification**

Seagate

intel.

# Linearization problem remains

Attributes

Records

- Placement within an object can be optimized

Disk

Seagate

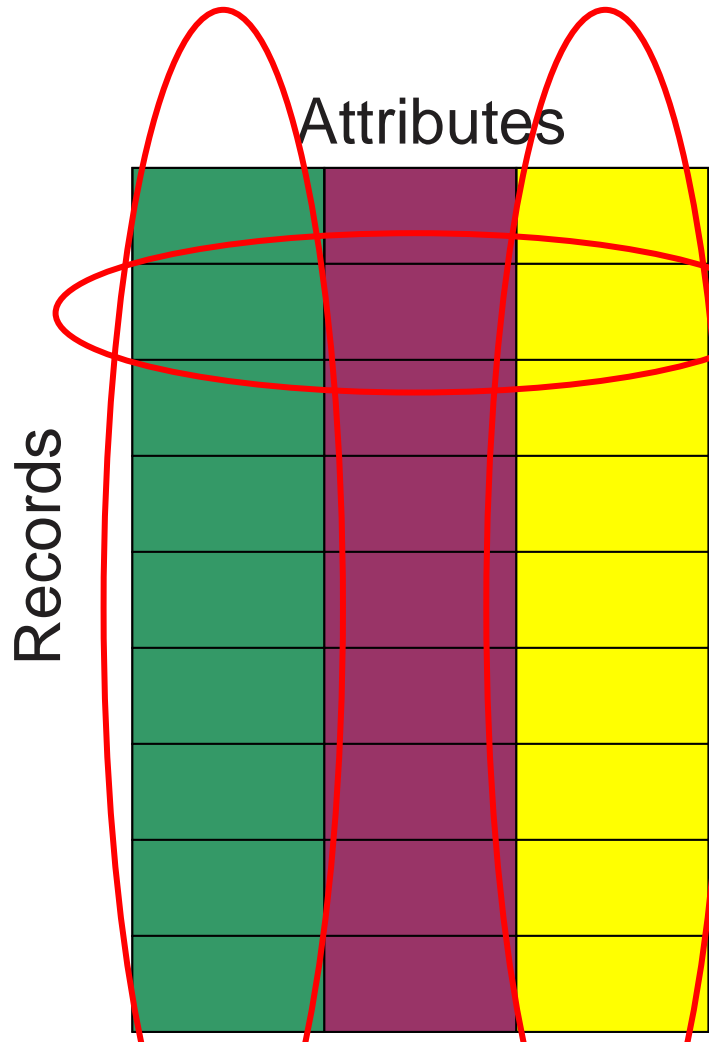intel.

# Linearization problem remains

Attributes

Query B

Records

Object

Query A

- Placement within an object can be optimized
- However, OSD objects are still addressed as a linear array of bytes
  - Query A: Two requests
  - Query B: Many requests, DBMS must calculate byte offsets

- Violates our goal!
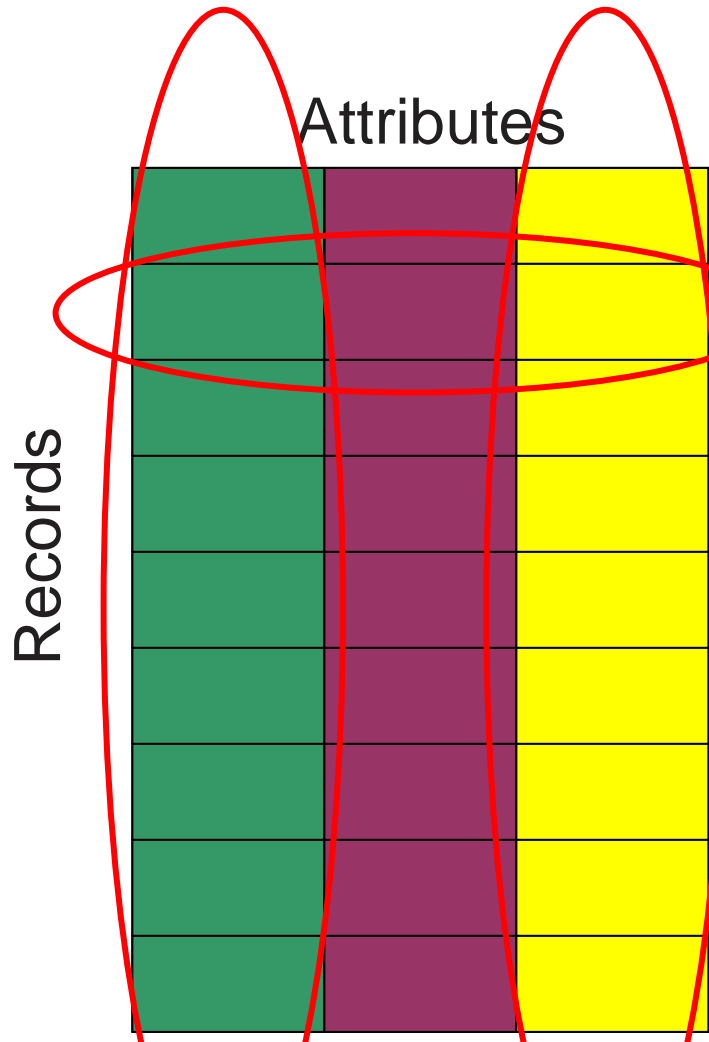  - DBMS should be independent of storage optimizations

Seagate

intel.

# Solution #1: Two-dimensional object interface

Attributes

Query B

Records

- Provide two-dimensional READ and WRITE commands
  - READ (offsetX, lenX, offsetY, lenY)
  - WRITE (offsetX, lenX, offsetY, lenY)

- Query A: Two requests
- Query B: One request

- Violates our goal!
  - DBMS must still translate byte offsets

Seagate

Query A

intel.

# Solution #2: Relational object interface

Attributes

Records

Query B

Query A

- Provide relational READ and WRITE commands
  - READ (record offset, len, field bitmap)
  - WRITE (record offset, len, field bitmap)

- Bitmap specifies fields to be accessed
  - Query A: READ(0, 9, 101)
  - Query B: READ(1, 1, 111)

- Meets our goal
  - DBMS addresses object in terms of the relation's schema

Seagate

intel.

# Conclusion

- Object-based storage devices allow applications to provide storage subsystems with high-level knowledge about data

- Demonstrated how database systems can take advantage of modern placement techniques using OSD interfaces

- Placement is just the first step

- In the future, how can database systems cooperate better with storage systems?