# Data-centric OSes

## NVM and the Death of the Process

HPTS '19

**Daniel Bittman**        Peter Alvaro        Ethan Miller

UC Santa Cruz

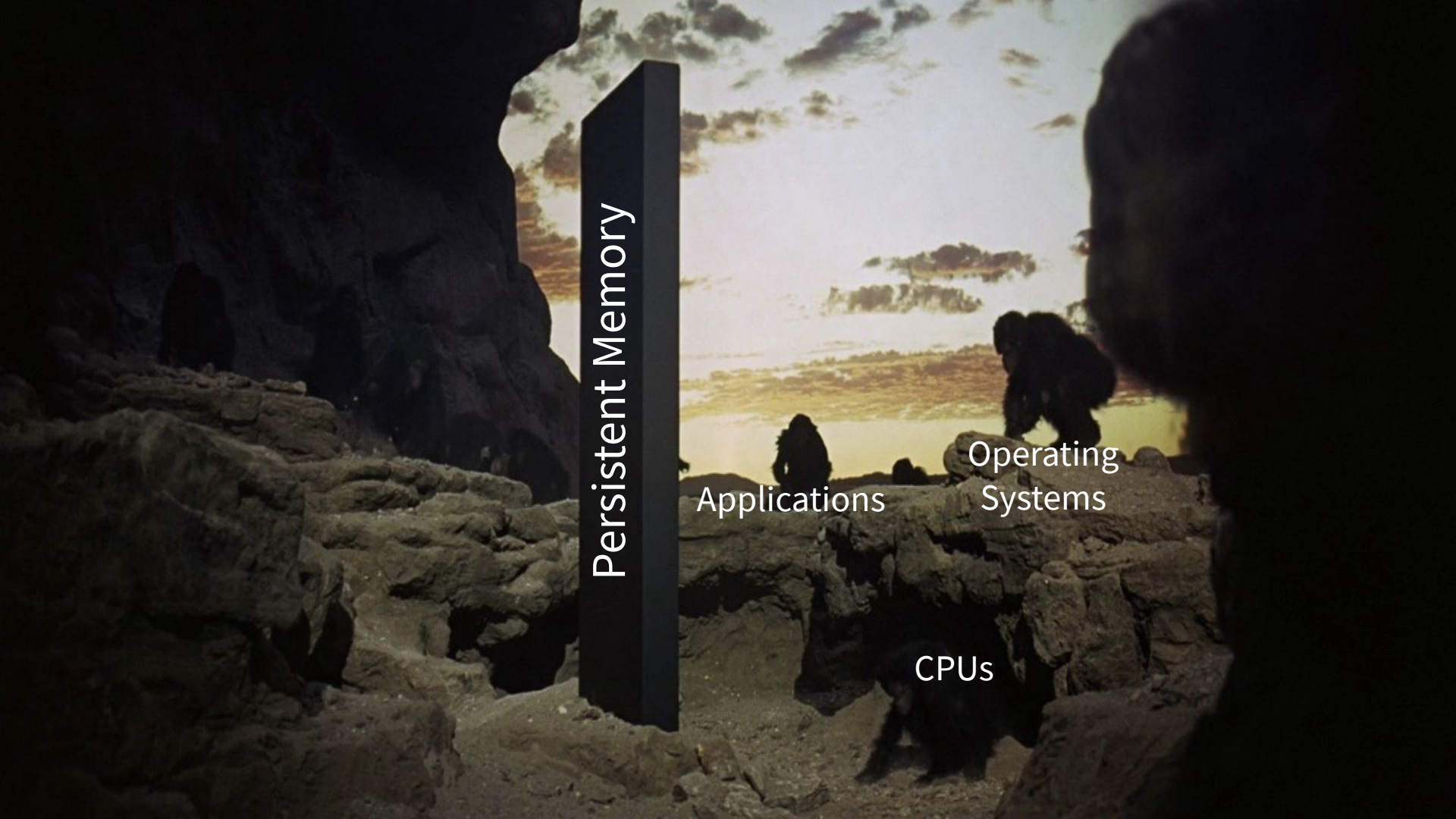# Data-centric OSes
## NVM and the Death of the Process

HPTS '19

**Daniel Bittman**        Peter Alvaro        Ethan Miller

UC Santa Cruz

# Hardware Trends



(artistic rendering;
actual implementation may vary)

`sys_read`

~300 ns

~1 us
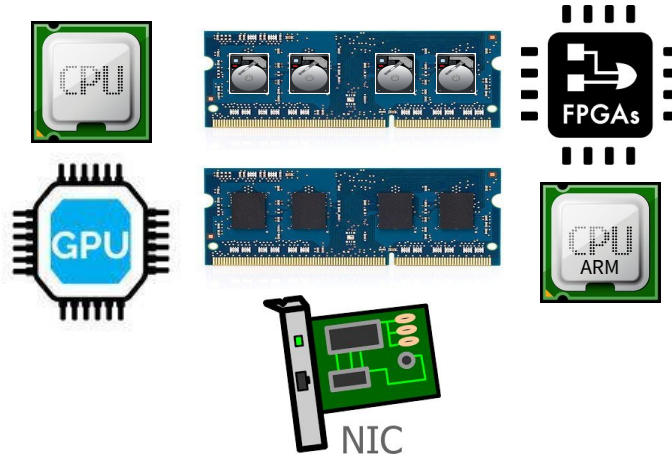
~1-10 ms

Growing, becoming persistent

Outdated interface

Cannot compute on directly

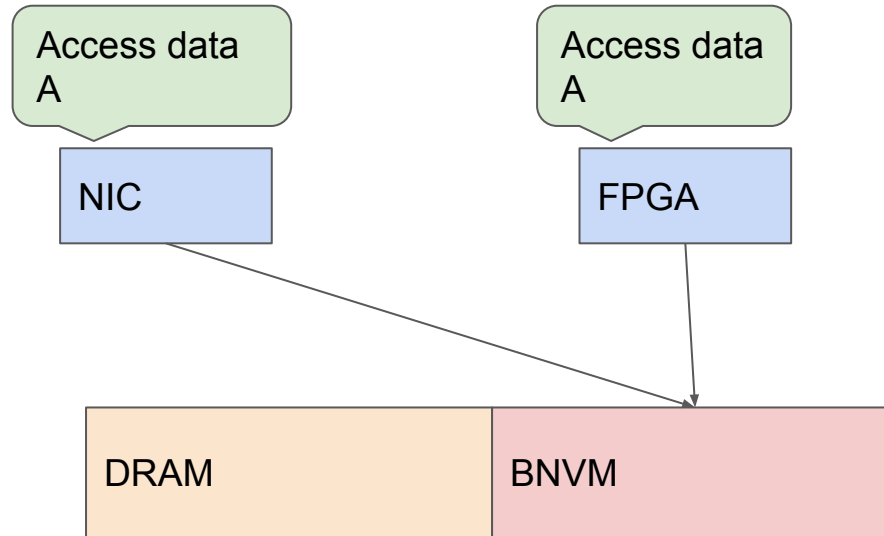**Persistent data should be operated on *directly* and *like memory***

# Hardware Trends

Multiplicity of Computing Devices and
Heterogeneous Memory

# Hardware's Needs vs. Software's Needs

| Consideration | Hardware | Software |
|---|:---:|:---:|
| Latency | ✓ | ✓ |
| In-memory Data Structures | ✗ | ✓ |
| Data Lifetime and Persistent Data References | ✗ | ✓ |
| Memory Heterogeneity and Data Movement | ✓ | ✗ |

# Heterogeneity and Autonomy

# Data Movement

Access data A

**Move** data A

NIC

FPGA

DRAM

BNVM

# In short…

**Software** cares about
**long-lived data relationships**,
even across program runs.

**Hardware** must provide
**consistent data access**, even
if it moves in memory.

Virtual memory is the **wrong** abstraction.

Virtual memory is fine.

Software is easier to change than hardware

# Twizzler: A new OS

The kernel is "out of the way"

Presents a unified interface for data sharing, security, and persistent pointers
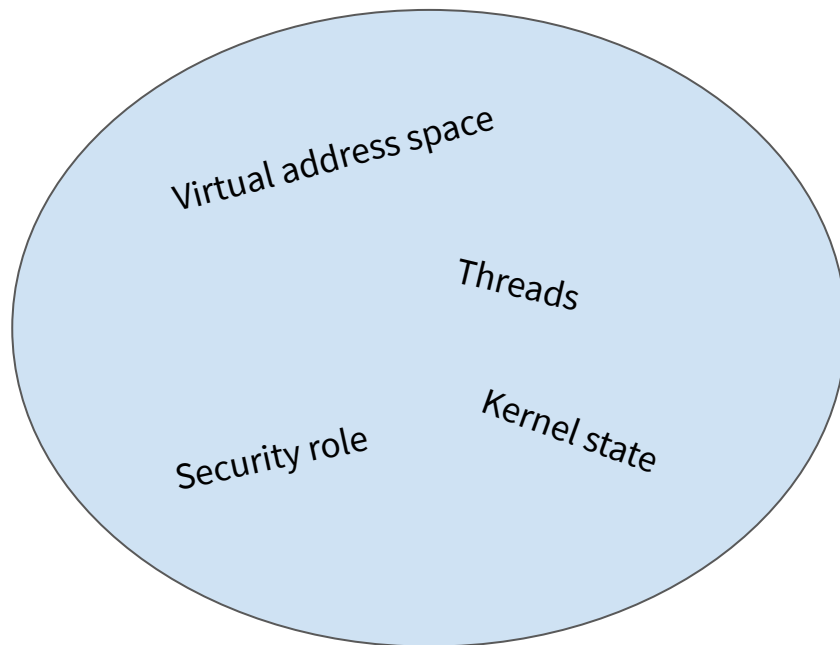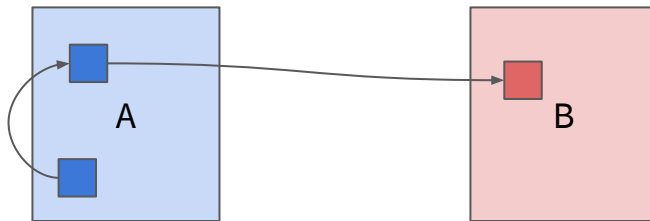
**OS Community**

**DB Community**

# The Death of the Process

Virtual address space

Threads

Security role

Kernel state

# A global object space

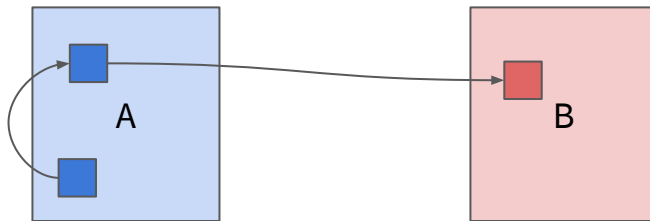**Persistent data should be operated on *directly* and *like memory***



An object is a unit of semantically similar information

E.g. a b-tree, or part of one.

# A global object space

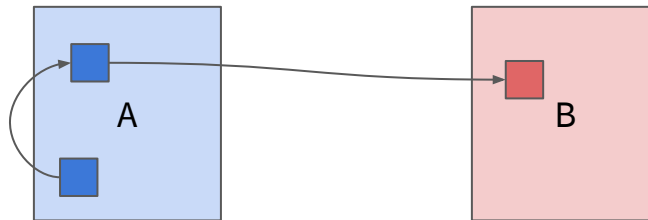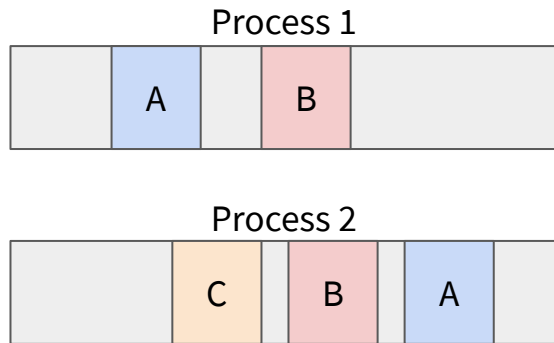**Persistent data should be operated on *directly* and *like memory***



An object is a unit of semantically similar information

E.g. a b-tree, or part of one.

**Pointers may be *cross-object*: referring to data within a different object**

# Persistent pointers in Twizzler

**Virtual addresses are the *wrong* abstraction**

Process 1

Process 2

object-id offset

# Twizzler's pointers

| FOT entry | offset |
|:---:|:---:|

64-bits

## Foreign Object Table

| | object ID or Name | Name Resolver | flags |
|:---:|:---:|:---:|:---:|
| 1 | object ID or Name | Name Resolver | flags |
| 2 | object ID or Name | Name Resolver | flags |

...

## Object Layout

| FOT | Data |
|:---:|:---:|

# Example pointer resolution
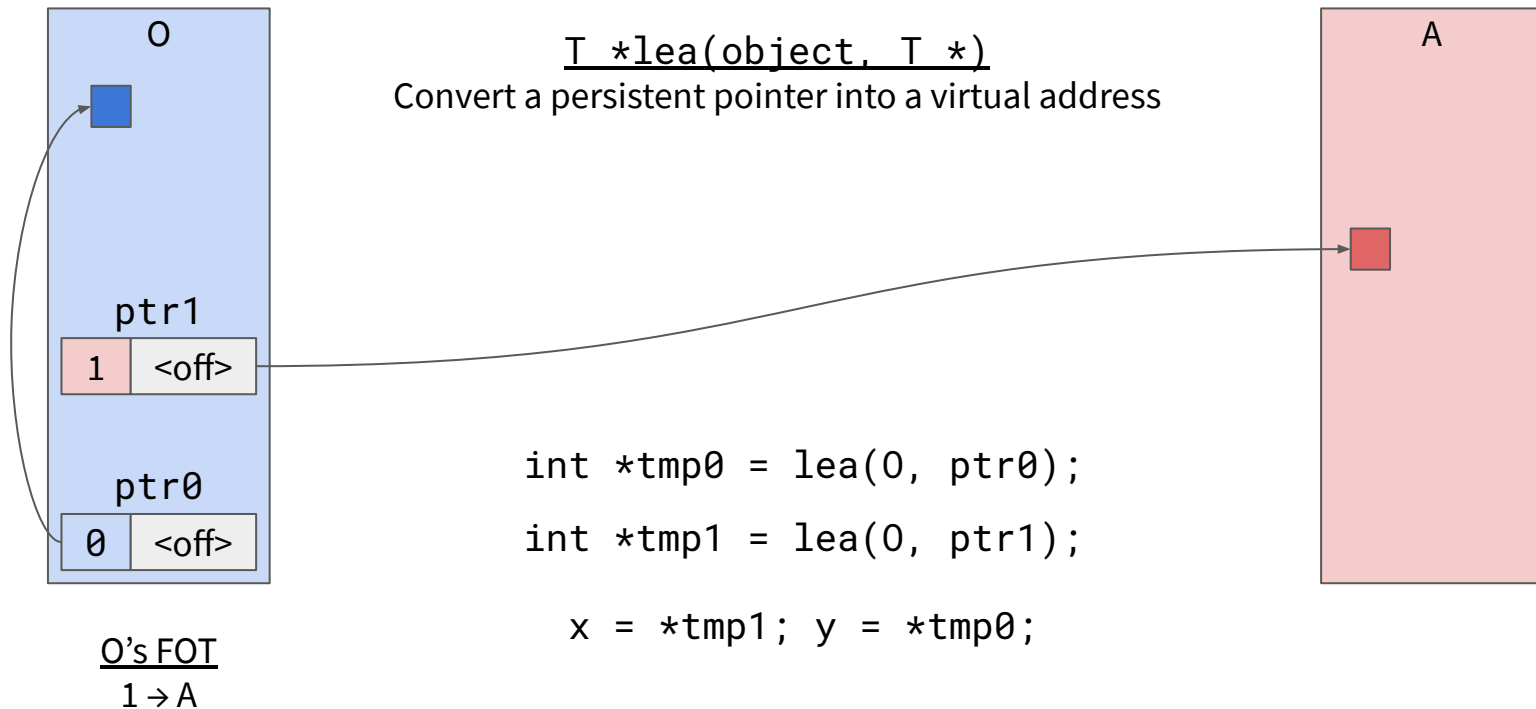


FOT entry of >0 means "cross-object"—points to a different object.

# Pointer implementation



O

T *lea(object, T *)
Convert a persistent pointer into a virtual address

A

ptr1

| 1 | <off> |

ptr0

| 0 | <off> |

O's FOT
1 → A

```
int *tmp0 = lea(O, ptr0);

int *tmp1 = lea(O, ptr1);

x = *tmp1; y = *tmp0;
```

# Two-level Mapping



Virtual Space

| Object A r-x | Object B rw- | |
|---|---|---|

Object Space

| | Object B r-- | | Object A rwx | | Object C r-- |
|---|---|---|---|---|---|

Physical Memory

DRAM          NVRAM

n+m page tables!
(instead of n*m)

*Security Contexts!*

# Hey look it's a Venn Diagram

Persistent Pointers

Data sharing

Twizzler

PMDK

Security model

# Benchmark: SQLite, throughput

# Benchmark: SQLite, latency

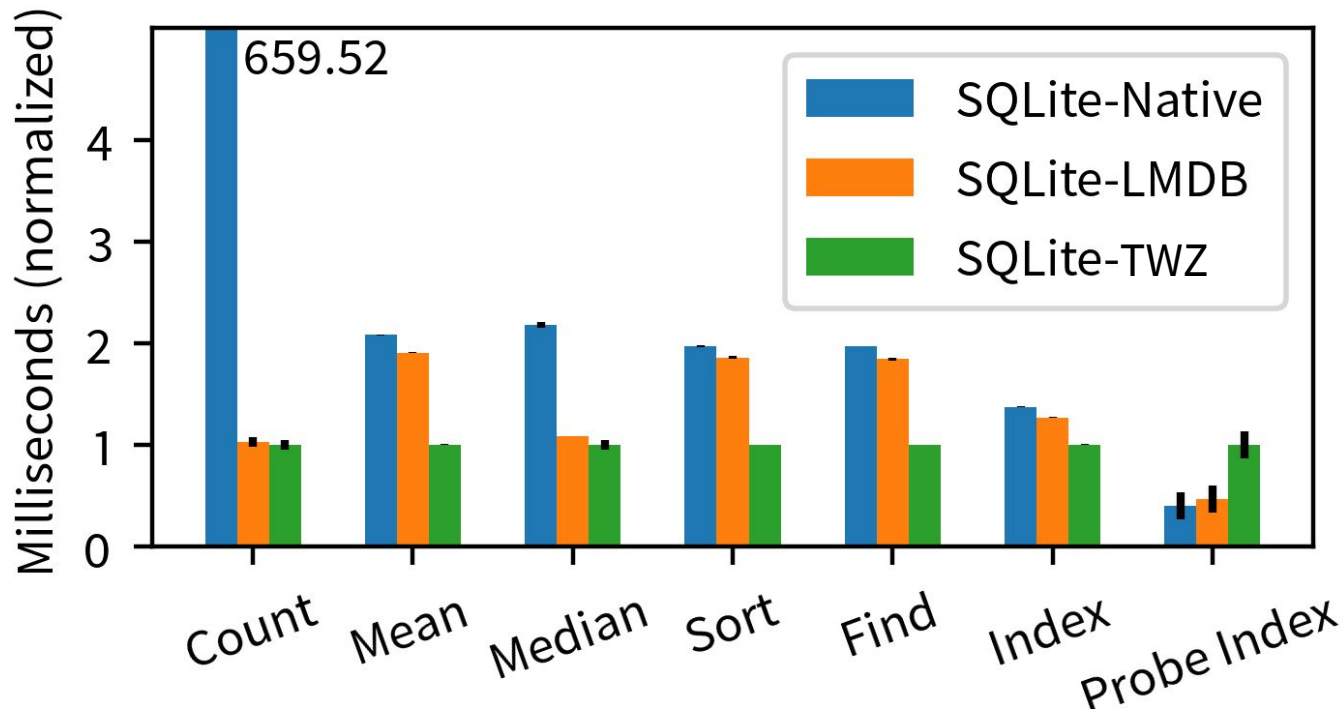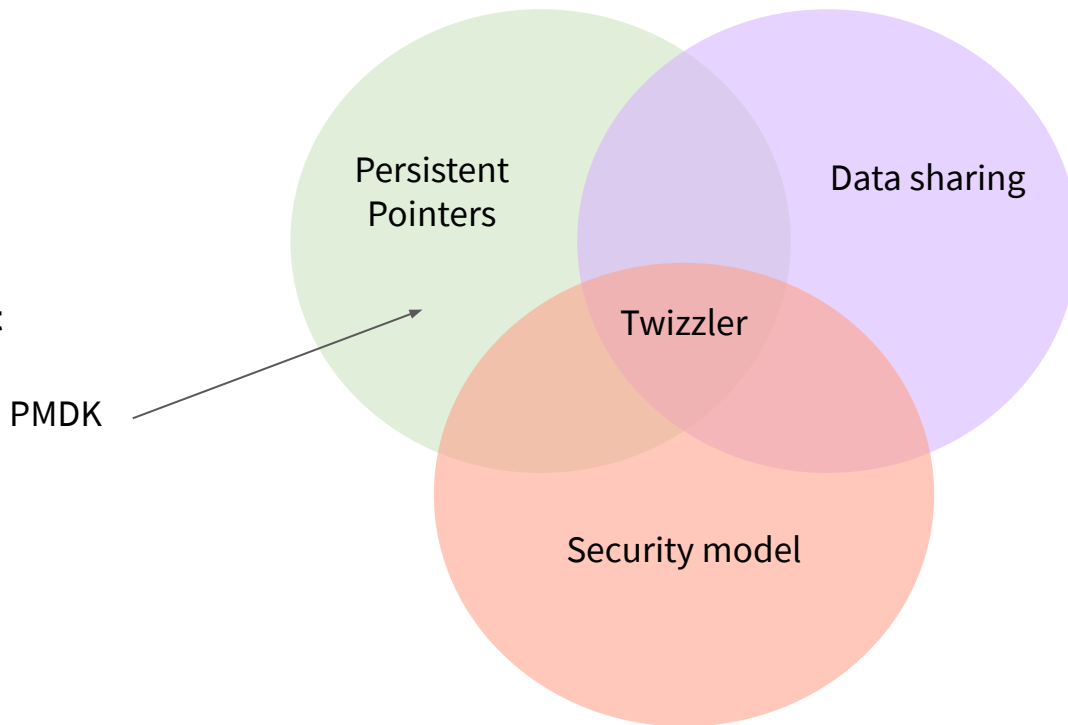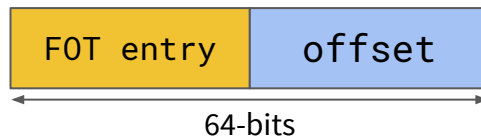# Takeaways - 1

**We need to consider persistent memory programming in the context of sharing and security**

Persistent Pointers

Data sharing

Twizzler

PMDK

Security model

# Takeaways - 2

| FOT entry | offset |
|-----------|--------|

64-bits

**A *flexible* persistent pointer design enables sharing, upgrades, and late-binding**

Foreign Object Table

| | object ID or Name | Name Resolver | flags |
|---|-------------------|---------------|-------|
| 1 | object ID or Name | Name Resolver | flags |
| 2 | object ID or Name | Name Resolver | flags |

...

Object Layout

| FOT | Data |
|-----|------|

# Takeaways - 3

**We are building Twizzler to explore new programming models for NVM**

**We must evolve our storage models for new technology**
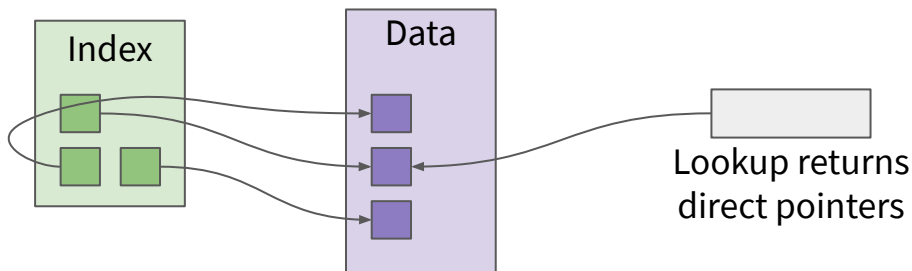
# Thank You!
## questions / discussion

**Daniel Bittman**
dbittman@ucsc.edu
@danielbittman

Peter Alvaro
palvaro@ucsc.edu

Ethan L. Miller
elm@ucsc.edu

# Case Study: KVS



Index

Data

Lookup returns
direct pointers

### insert(key, value)
```
bucket = get_bucket(key)
bucket.ptr = store(Index, value.ptr)
bucket.len = value.len
```
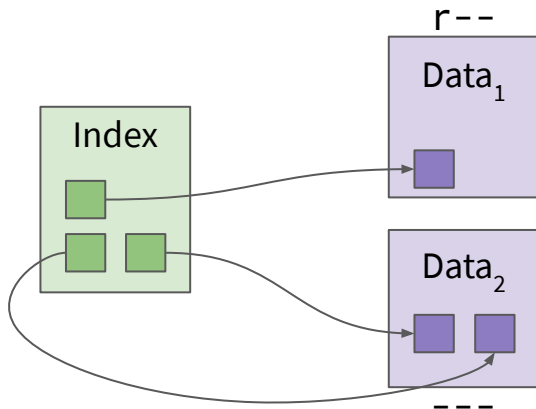
### lookup(key)
```
bucket = get_bucket(key)
item.ptr = lea(Index, bucket.ptr)
item.len = bucket.len
```

**250 lines of simple C code is *all you need***

(store is the reverse of lea: convert a virtual address into a persistent pointer)

# Cast Study: KVS

**Add access control to the existing design**



```
bucket = get_bucket(key)
item.ptr = lea(Index, bucket.ptr)
item.len = bucket.len
```

Index points to **different data objects** with **different access control**.

Can hand out pointers to these objects, which can **only be dereferenced with proper permissions**.

# Late-binding of access control

Super user

User
manager

User
Database

rw-

Normal user

User
manager

User
Database

r--