# The Story So Far

**2014:** Comparing Existing DBMSs

**2015:** Evaluating Storage Architectures

# 2014: Existing DBMSs

Comparison of disk vs. main-memory DBMSs running on Intel NVM SDV.

Found that logging is (still) the main bottleneck in both systems.

Paper: ADMS @ VLDB'14

# 2015: Storage Architectures

Evaluated storage and recovery methods for OLTP DBMSs.

Developed NVM-optimized methods that achieve 5.5x better throughput with 2x fewer writes.
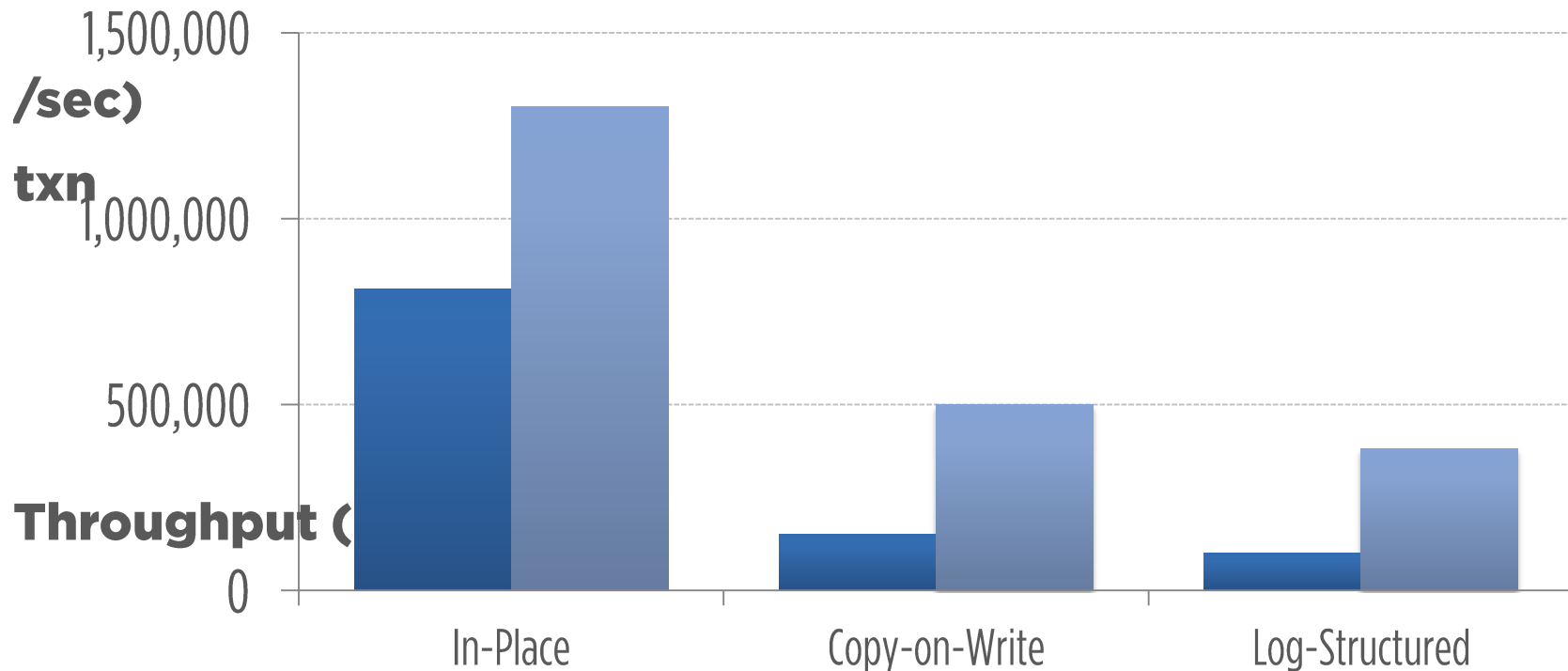
Paper: SIGMOD'15

# 2015: Storage Architectures

|  | Table Storage | Logging | Example |
|---|---|---|---|
| **In-Place** | Yes | Yes | VOLTDB |
| **Copy-on-Write** | Yes | No | LMDB |
| **Log-based** | No | Yes | RocksDB |

# YCSB :: 10/90 RW :: 2x Latency

# PCOMMIT Evaluation

Weakly-ordered sync primitive that retains data in the flushed cached lines.

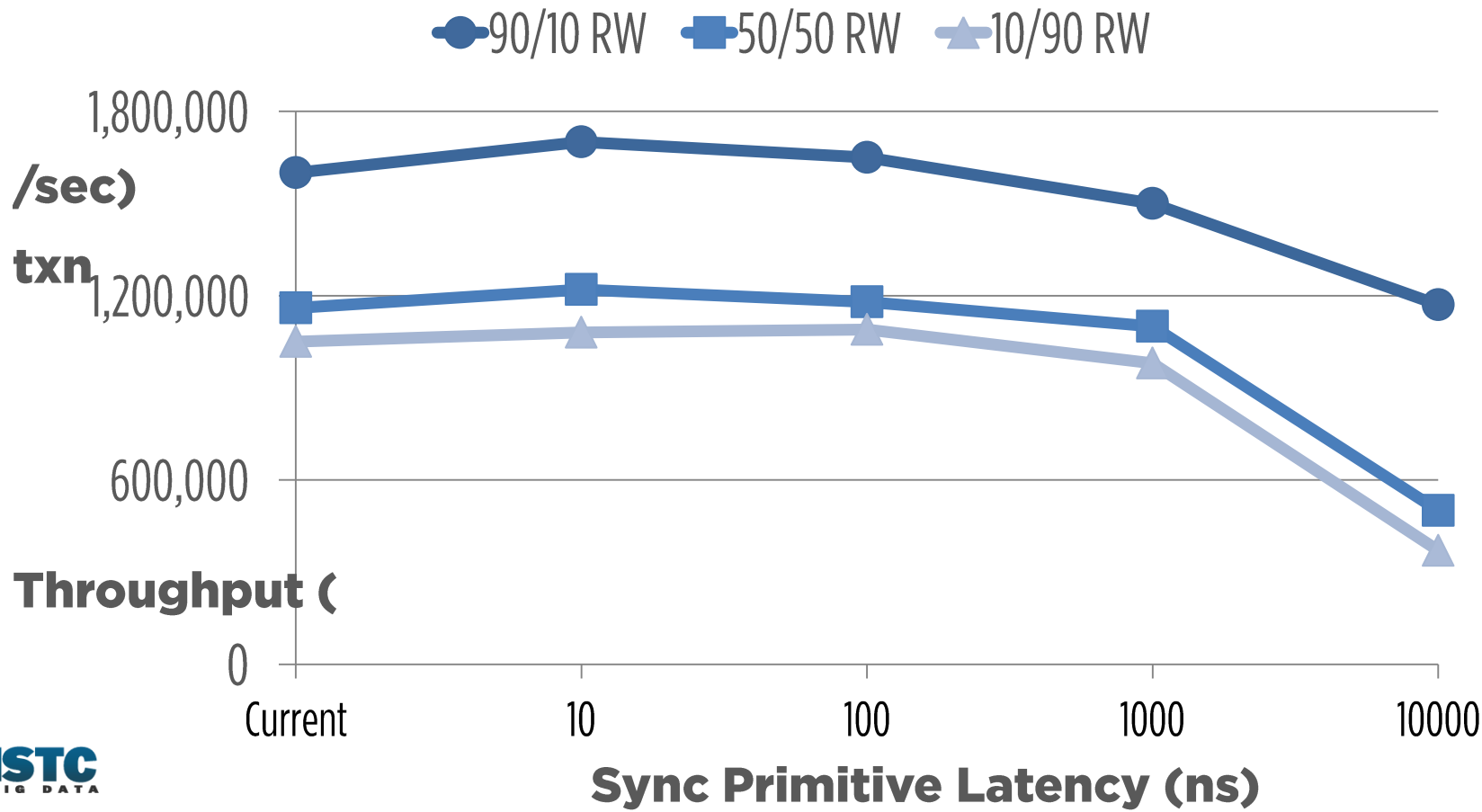Emulated with **RDTSC** and **PAUSE** instructions on NVM SDV.

**Summer 2015**: ~10,000 **PCOMMIT** invocations per second per CPU core.

# YCSB // In-Place Engine

# New Stuff

NVM vs. SSD

Multi-level Anti-Caching

DRAM+NVM storage manager

# NVM vs. SSD

Two-level Storage Hierarchy

Disk-oriented vs. Memory-Oriented
- *Caching (MySQL)*
- *Anti-caching (H-Store)*
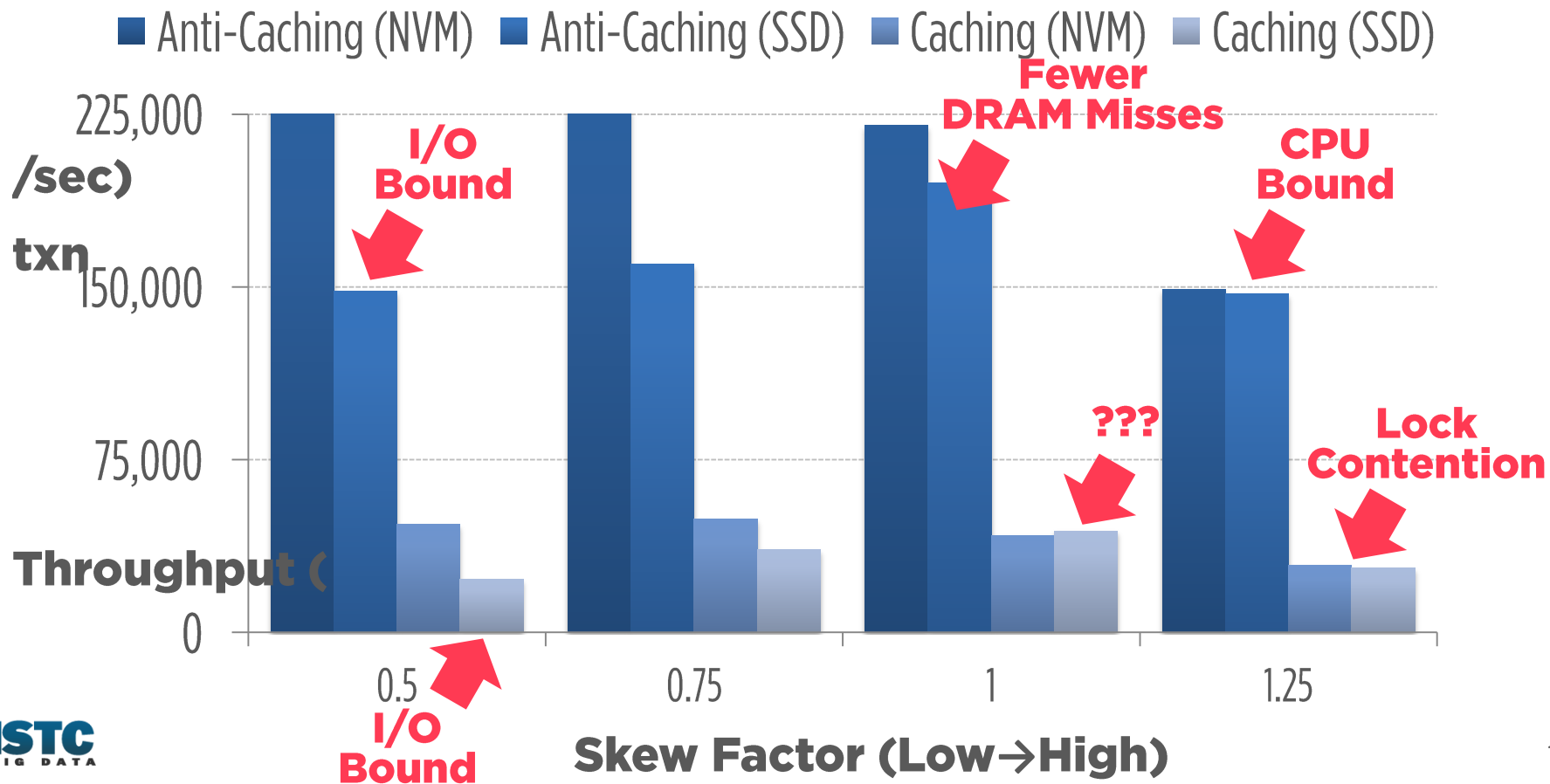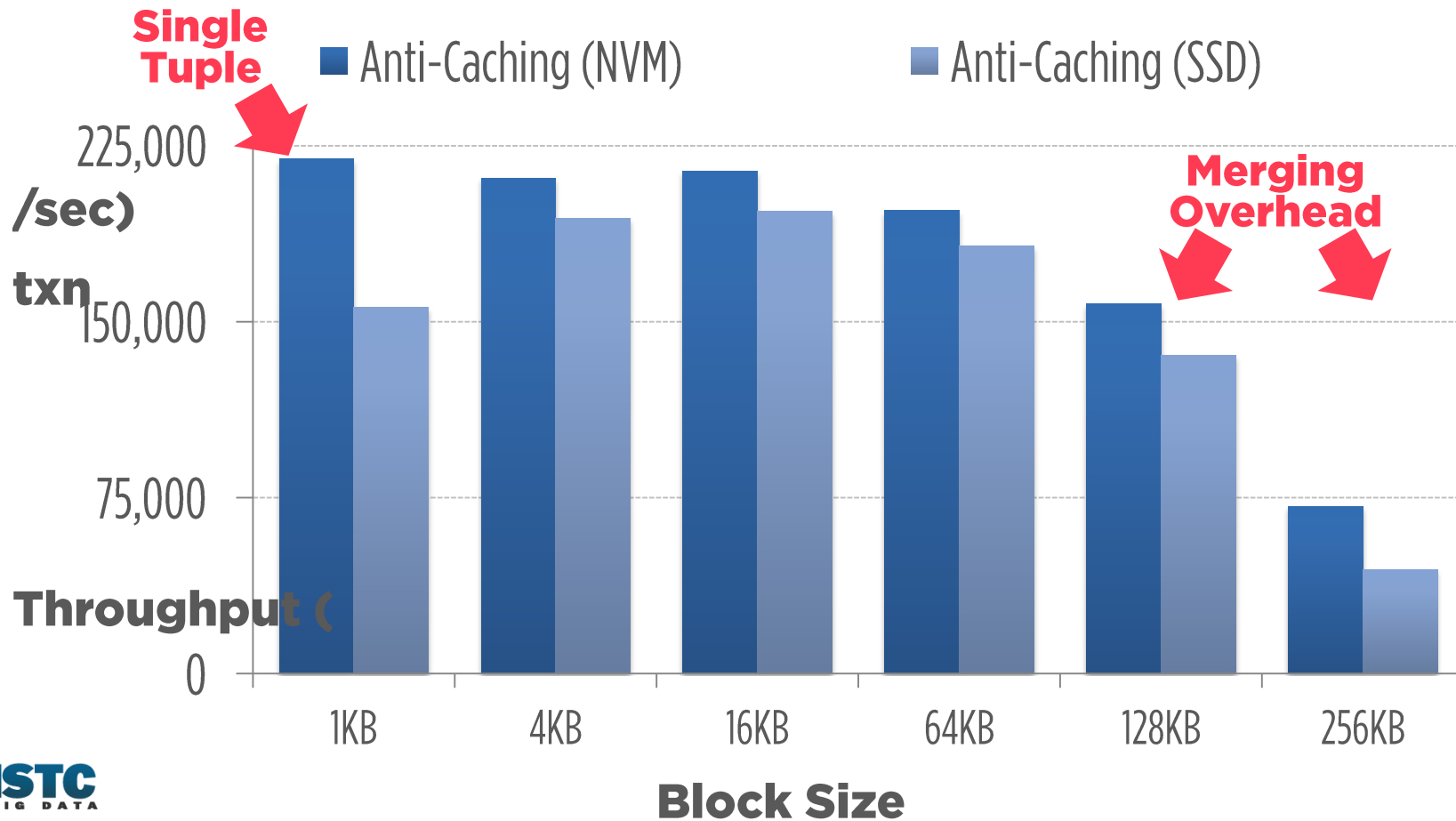
# Disk vs. Memory Oriented DBMSs

# YCSB :: 90/10 RW :: 4x Latency

Legend: Anti-Caching (NVM) · Anti-Caching (SSD) · Caching (NVM) · Caching (SSD)

Y-axis: Throughput (txn /sec) — 0, 75,000, 150,000, 225,000

X-axis: Skew Factor (Low→High) — 0.5, 0.75, 1, 1.25

Annotations: I/O Bound · Fewer DRAM Misses · CPU Bound · ??? · Lock Contention · I/O Bound

# YCSB :: Byte-Addressable Access

# Voter :: 4x Latency



Anti-Caching (NVM)

225000

150000

/sec)

txn   75000

0

**DRAM Evictions**

Anti-Caching (SSD)

225000

150000

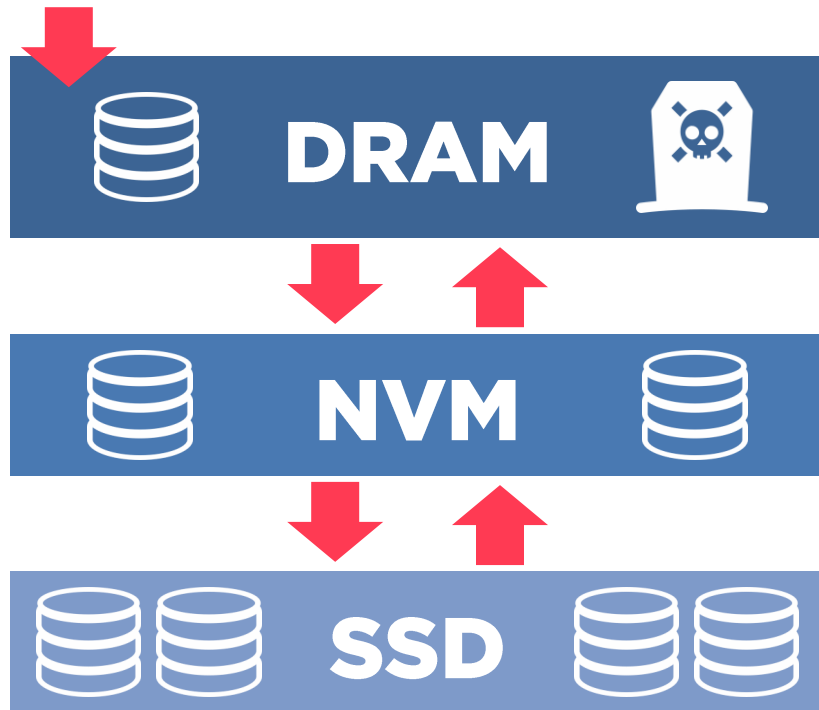**Throughput (**

75000

0

**Elapsed Time**

16

# Multi-Level Anti-Caching



OLTP Workload

DRAM

NVM

SSD

**Current Investigation:**
Eviction Policies
Retrieval Policies
Access Interfaces
Data Organization

# Multi-Level Anti-Caching

OLTP Workload

## Multi-Level Anti-Caching

Lin Ma[5], Michael Giardino[2], Sam Zhao[1], Joy Arulraj[5], Dana Van Aken[5], Prashanth Menon[5]
Ugur Cetintemel[1], Justin DeBrabant[1], Kshitij Doshi[2], Subramanya R. Dulloor[2], Aaron Elmore[6]
Samuel Madden[3], David Maier[4], Michael Kaminsky[2], Tim Kraska[1], Jeff Parkhurst[2]
Andrew Pavlo[5], Michael Stonebraker[3], Nesime Tatbul[2,3], Donald Trump
Kristin Tufte[4], Stanley Zdonik[1]

[1]Brown  [2]Intel Labs  [3]MIT  [4]Portland State University  [5]CMU  [6]Univ. of Chicago
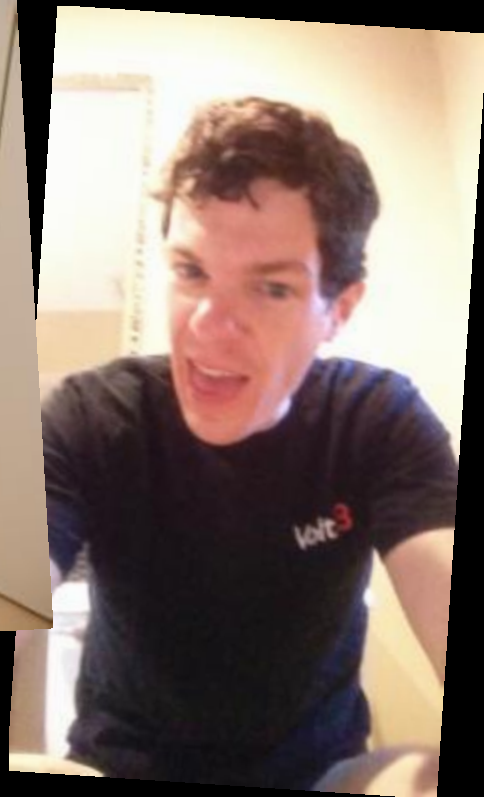
SSD

# Voter :: Multi-Level :: 2x Latency

# DRAM+NVM DBMS

Building a new storage manager for our new DBMS that will seamlessly incorporate NVM as an extension to its address space.

Upper-levels of the system are oblivious to "true" location of data.

# END

@ANDY_PAVLO