



WHAT **NON-VOLATILE  
MEMORY**  
MEANS FOR THE FUTURE OF  
**DATABASE  
SYSTEMS**

# STORAGE LATENCY



**RAMAC 350**  
(600 ms)

**1956**

$10^5 \times \downarrow$



**NAND SSD**  
(60 us)

**2016**

# COMPUTE LATENCY



**RAMAC 305**  
(100 Hz)

**1956**

**$10^8$ x ↓**

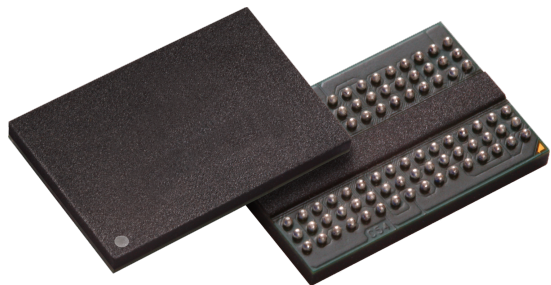
**1000x**



**CORE I7**  
(1 GHz)

**2016**

# NON-VOLATILE MEMORY



**3D XPOINT**  
(60 ns)

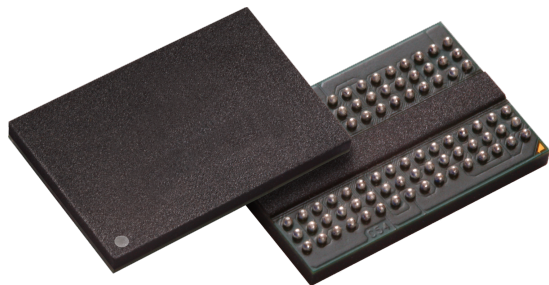
**2017**

**1000x faster than NAND**

**10x denser than DRAM**

**1000x endurance of NAND**

# NON-VOLATILE MEMORY



**3D XPOINT**  
(60 ns)

**2017**

*Support in Linux 4.3+*

*New assembly instructions*

*SNIA programming model*

# WHAT NVM MEANS FOR DATABASES?

- Option 1: Treat NVM like a faster SSD

DBMS



*Use NVM as a logging and backing store*

DRAM



*Improves performance*

NVM



# WHAT NVM MEANS FOR DATABASES?

- Option 2: Treat NVM as extended memory

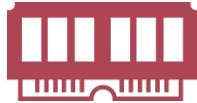
DBMS



*Redesign the key algorithms  
in a database system*



DRAM

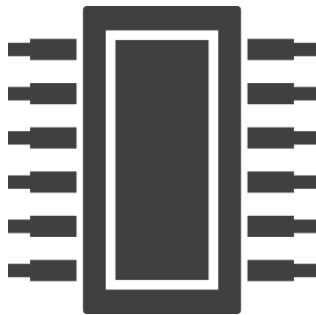


NVM

*Improves not only performance,  
but also availability  
and device utilization*



**PAST:  
EXISTING  
SYSTEMS**



**PRESENT:  
NVM-AWARE  
DBMS**



**FUTURE:  
ANALYTICS  
ON NVM**



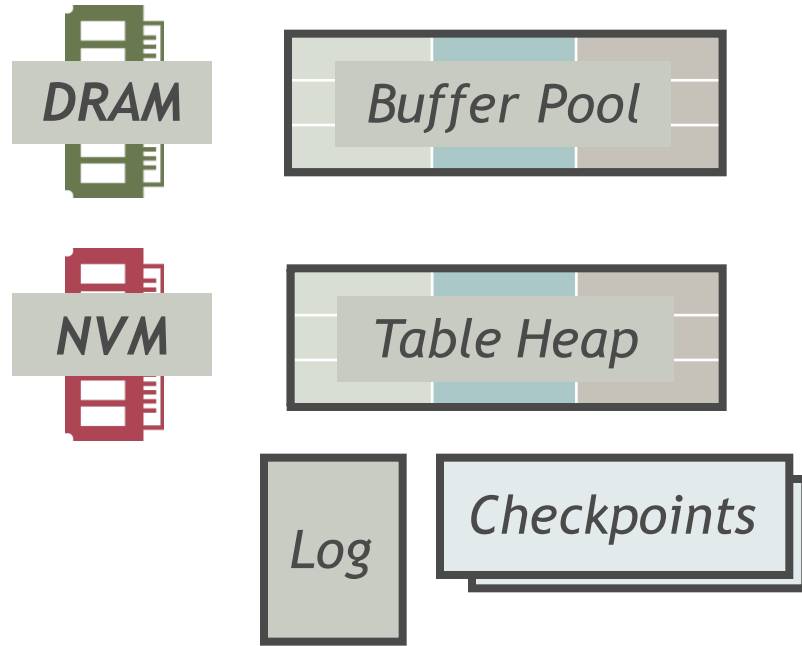
# PAST – EXISTING SYSTEMS

- Investigate how existing systems perform on NVM
  - *Option 1: Treat NVM like a faster SSD*
- Evaluate two types of DBMSs
  - *Disk-oriented DBMS*
  - *In-memory DBMS*

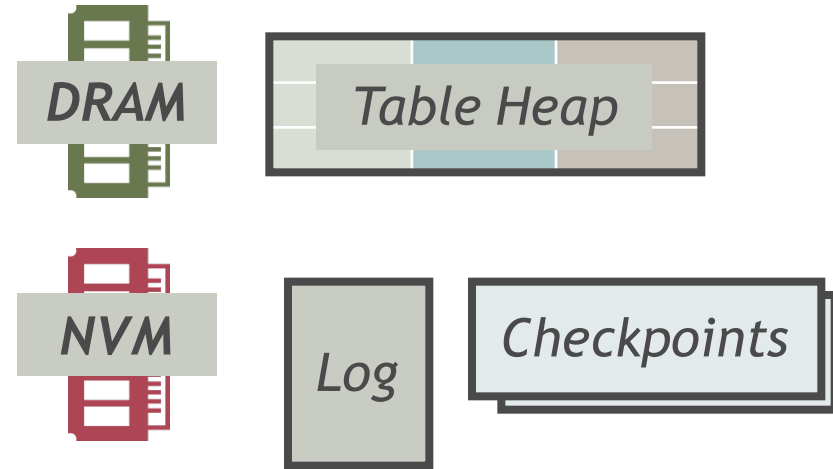


# DBMS ARCHITECTURES

## DISK-ORIENTED DBMS

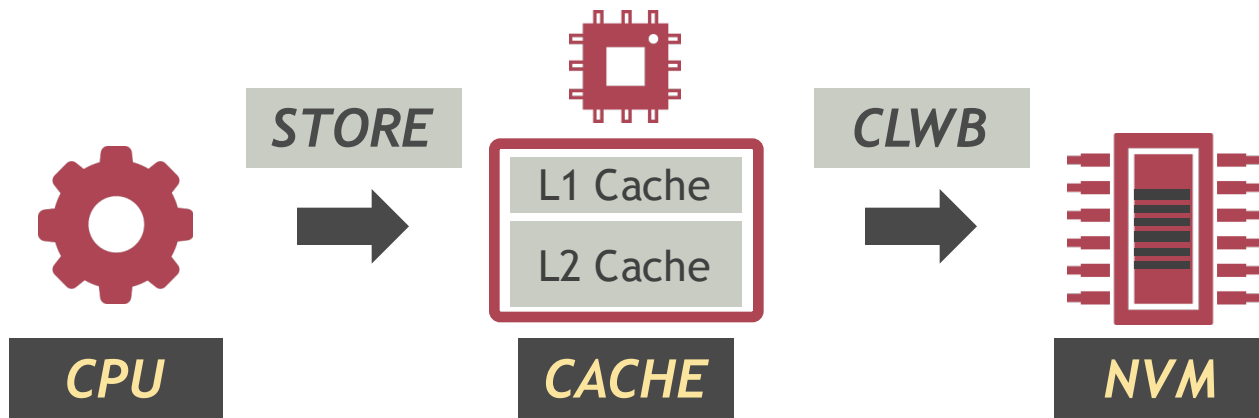


## IN-MEMORY DBMS



# NVM HARDWARE EMULATOR

- Tunable DRAM latency for emulating NVM
- Special assembly instructions for NVM
  - *Cache line write-back*



# EVALUATION

---

- Compare existing DBMSs on NVM emulator
  - *MySQL (Disk-oriented DBMS)*
  - *H-Store (In-memory DBMS)*
- TPC-C benchmark
  - *1/8<sup>th</sup> of database fits in DRAM, rest on NVM*

# PERFORMANCE

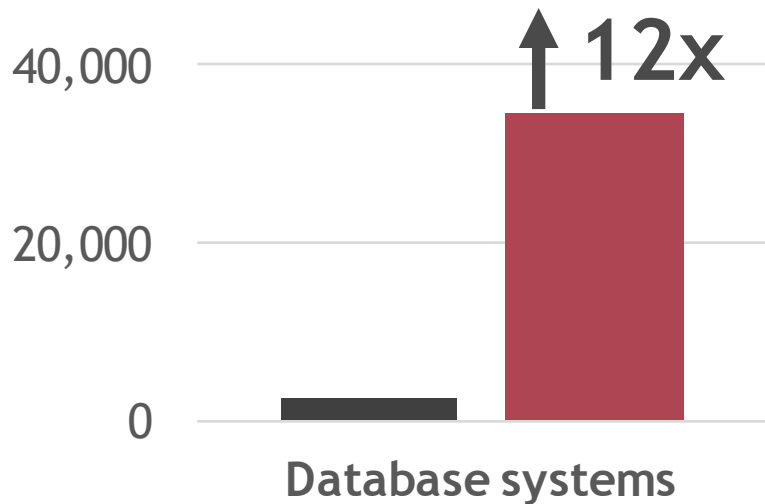


Disk-Oriented DBMS



In-memory DBMS

*8x DRAM Latency*



Throughput  
(txn/sec)

# PERFORMANCE

■ Disk-Oriented DBMS      ■ In-memory DBMS

*2x DRAM Latency*

↓ 4x

*Legacy database systems are not prepared for NVM*

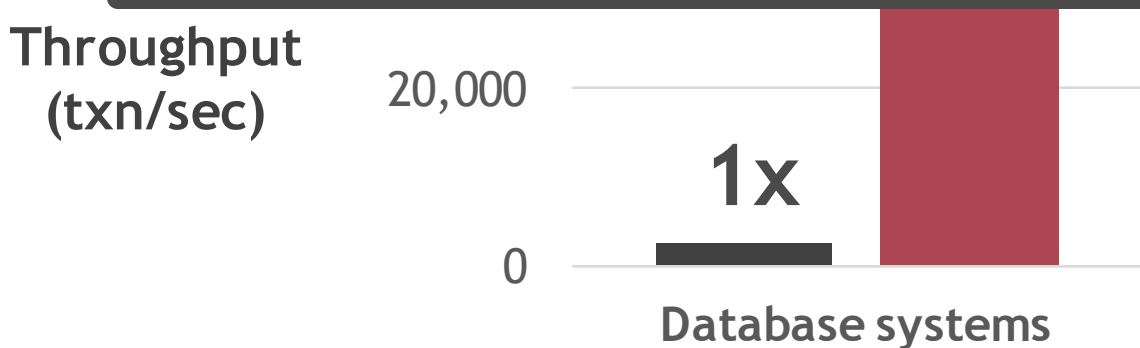
Throughput  
(txn/sec)

20,000

0

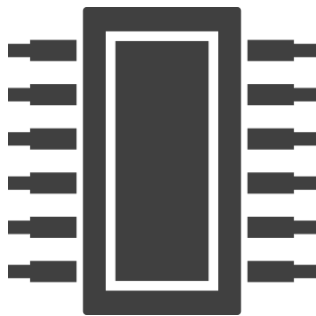
1x

Database systems





**PAST:  
EXISTING  
SYSTEMS**



**PRESENT:  
NVM-AWARE  
DBMS**



**FUTURE:  
ANALYTICS  
ON NVM**

# PRESENT – NVM-AWARE DBMS

- Understand changes required in DBMSs to leverage NVM
  - *Option 2: Treat NVM as extended memory*
- Rethink key algorithms in database systems
  - *Logging and recovery protocol*
- Designed for real-time analytics
  - *Multi-versioned database*





# MULTI-VERSIONED DBMS

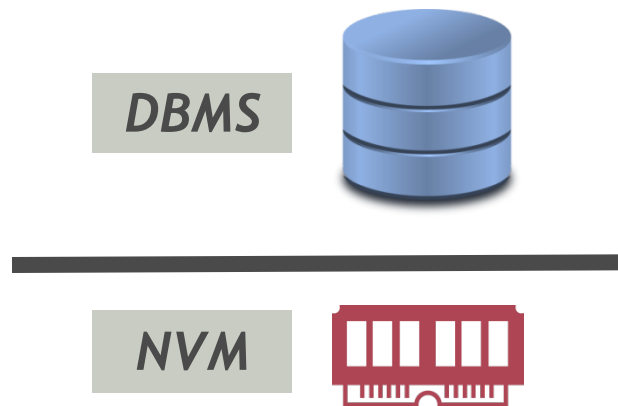
---

TUPLE ID	BEGIN TIMESTAMP	END TIMESTAMP	PREVIOUS VERSION	TUPLE DATA
1	10	$\infty$	—	X
2	10	20	—	Y
3	20	$\infty$	2	Y'

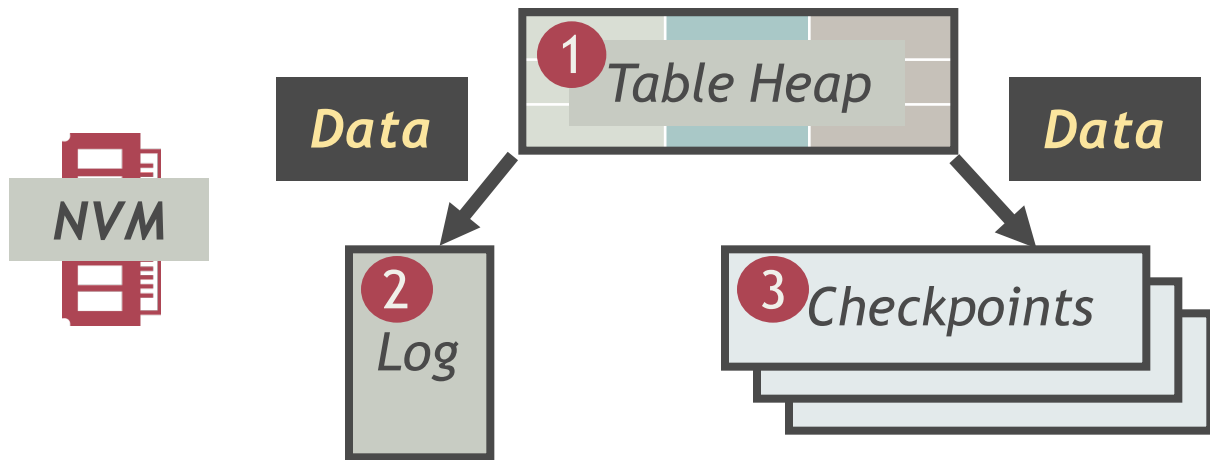
# THOUGHT EXPERIMENT

---

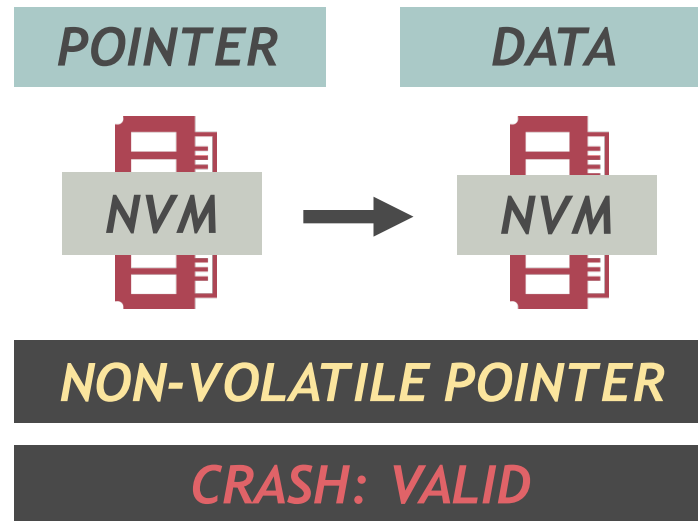
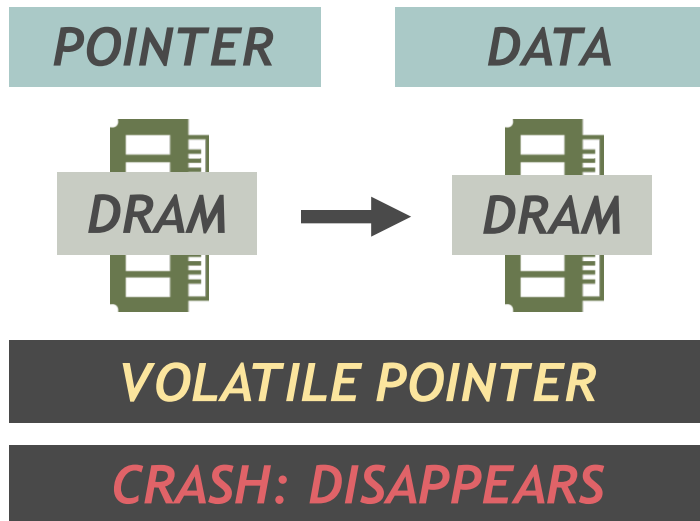
- NVM-only storage hierarchy
  - *No volatile DRAM*



# TRADITIONAL STORAGE ENGINE



# NON-VOLATILE POINTER



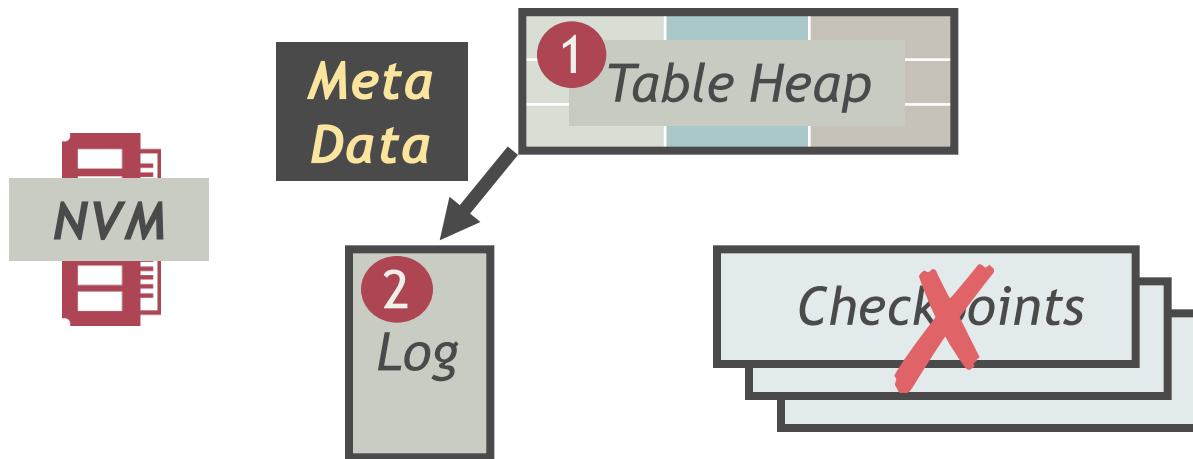
# NVM-AWARE STORAGE ENGINE

- Instead of duplicating tuple data in log
  - *Store non-volatile tuple pointers in log records*

TRADITIONAL LOG
INSERT TUPLE 1 (DATA)
UPDATE TUPLE 1 (DATA)

NVM-AWARE LOG
INSERT TUPLE 1 (POINTER)
UPDATE TUPLE 1 (POINTER)

# NVM-AWARE STORAGE ENGINE



# EVALUATION

---

- Compare storage engines on NVM emulator
  - *Traditional engine*
  - *NVM-aware engine*
- Yahoo! Cloud Serving Benchmark
  - *Database fits in NVM*

# PERFORMANCE



Traditional Engine



NVM-Aware Engine

**8x DRAM Latency**



**3x**

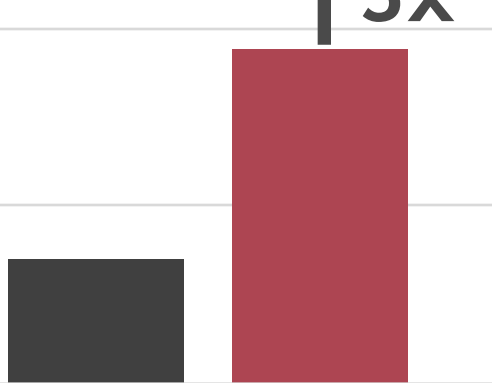
600,000

300,000

0

Throughput  
(txn/sec)

Storage Engines





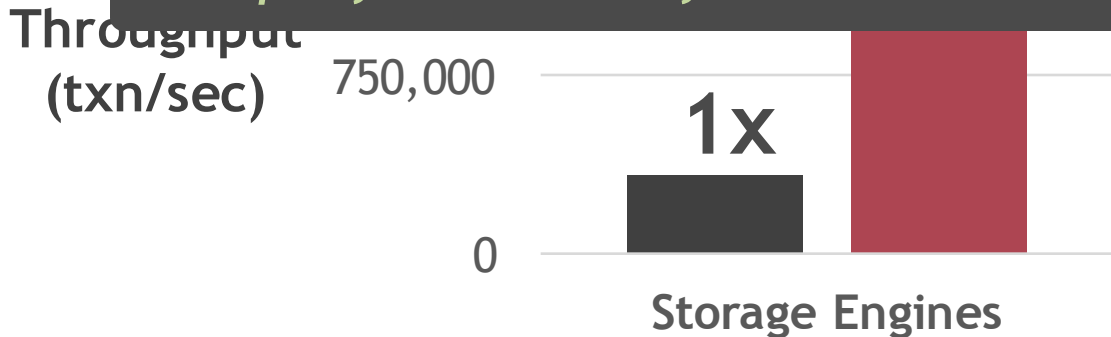
# PERFORMANCE

■ Traditional Engine    ■ NVM-Aware Engine

**2x DRAM Latency**

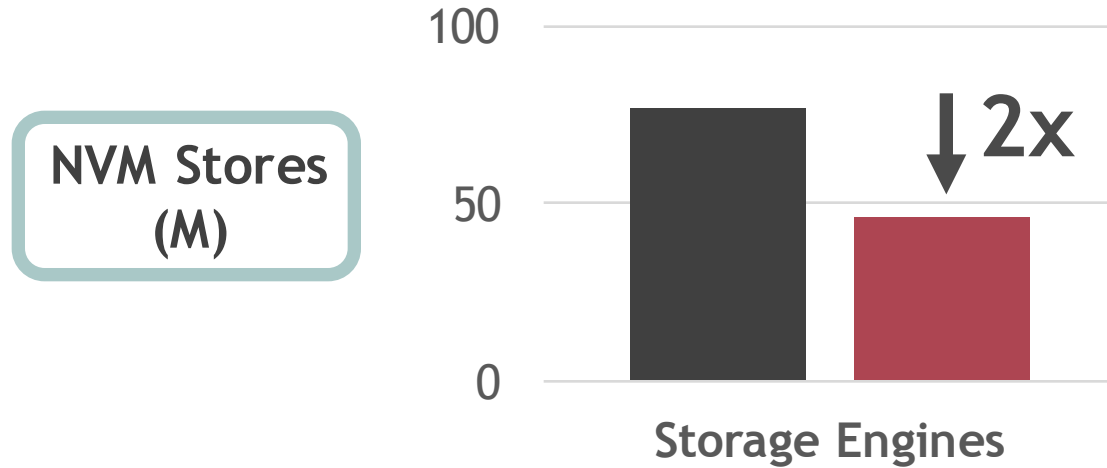
↓ 4x

*NVM latency has a significant impact on the performance of NVM-aware storage engine*



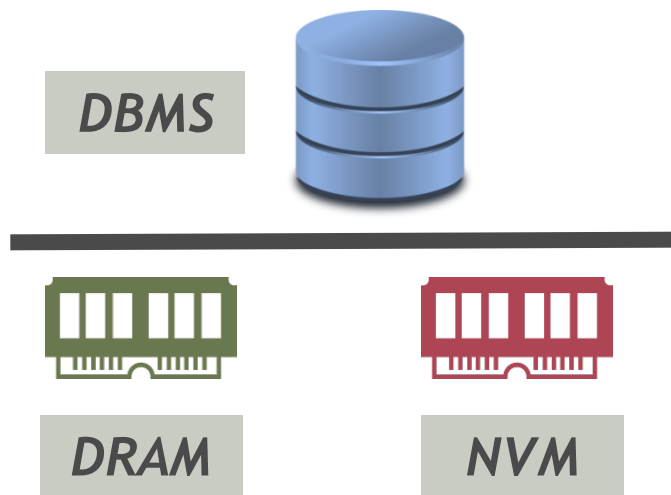
# DEVICE LIFETIME

■ Traditional Engine    ■ NVM-Aware Engine



# THE PELOTON DBMS

- A new database system for NVM research
  - *Designed for a two-tier storage hierarchy*



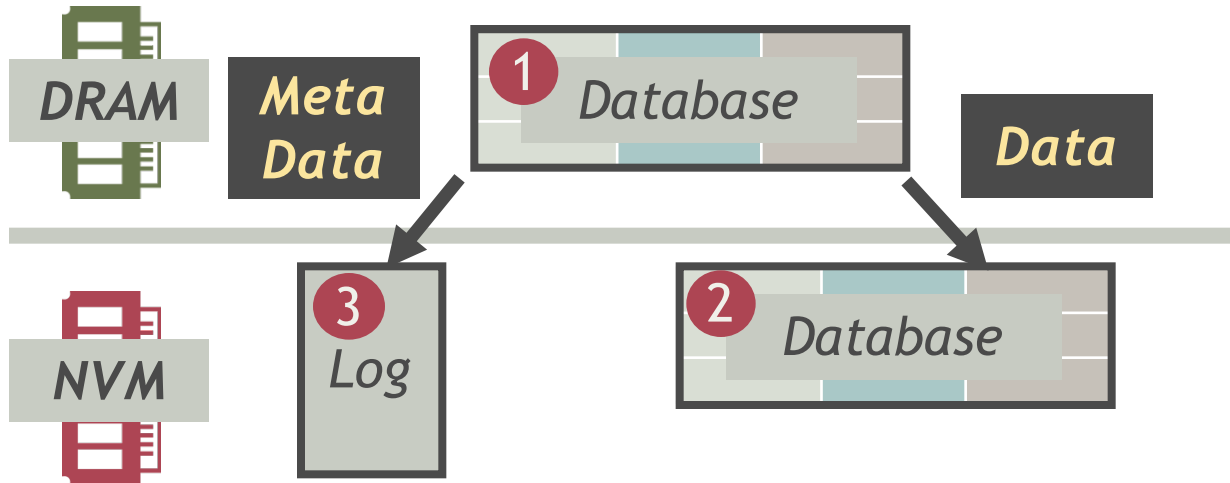
WRITE-BEHIND LOGGING  
UNDER REVIEW

# WRITE-BEHIND LOGGING

---

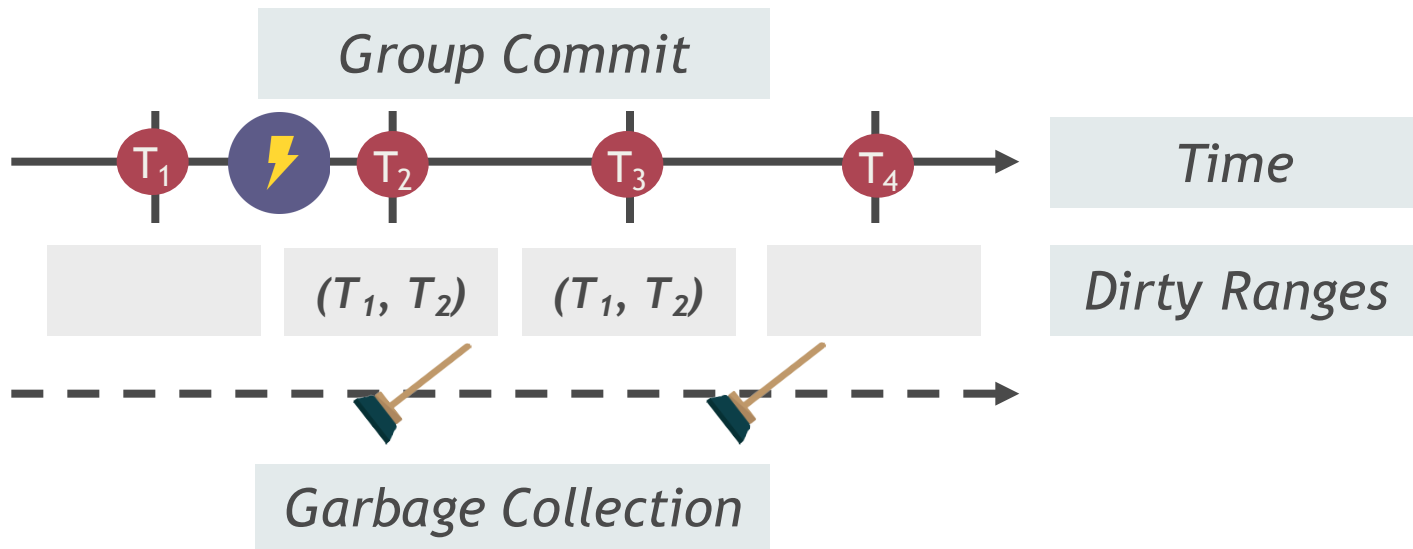
- Write-ahead log serves two purposes
  - *Transform random database writes into sequential log writes*
  - *Support transaction rollback*
- NVM supports fast random writes
  - *Directly write data to the database*
  - *Later, record metadata in write-behind log*

# WRITE-BEHIND LOGGING



# METADATA FOR INSTANT RECOVERY

- Record failed timestamp ranges in log
  - Use it to ignore effects of uncommitted transactions



# EVALUATION

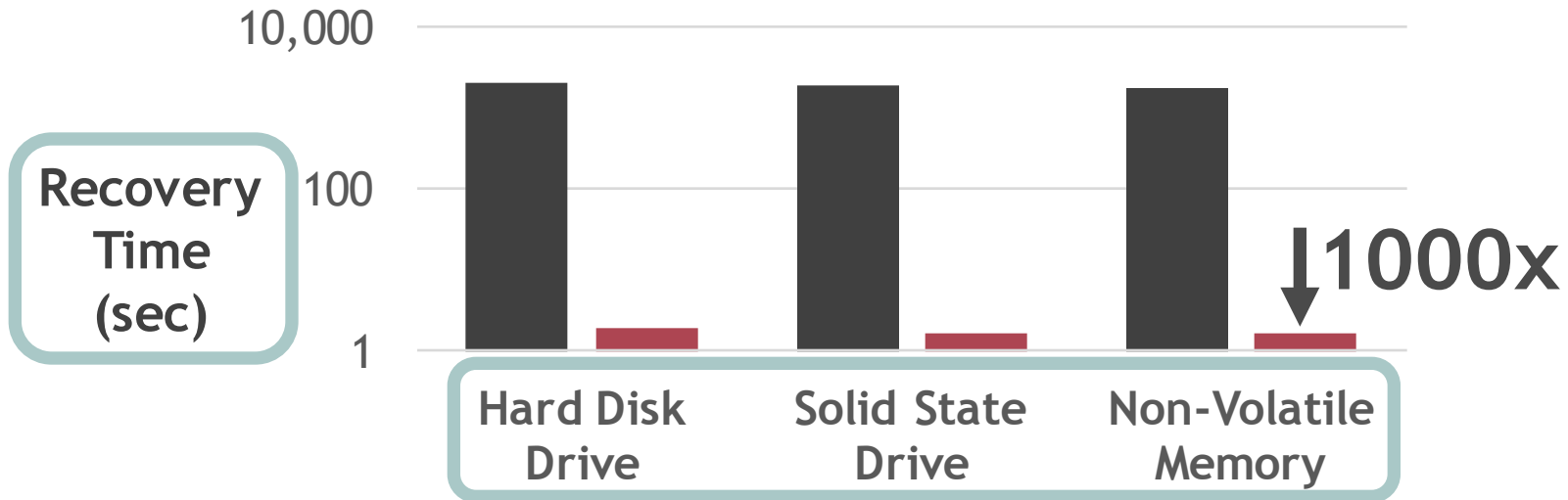
---

- Compare logging protocols in Peloton
  - *Write-Ahead logging*
  - *Write-Behind logging*
- Storage devices in Intel's NVM hardware emulator
  - *Hard-disk drive*
  - *Solid-state drive*
  - *Non-volatile memory emulator*
- TPC-C benchmark

# APPLICATION AVAILABILITY

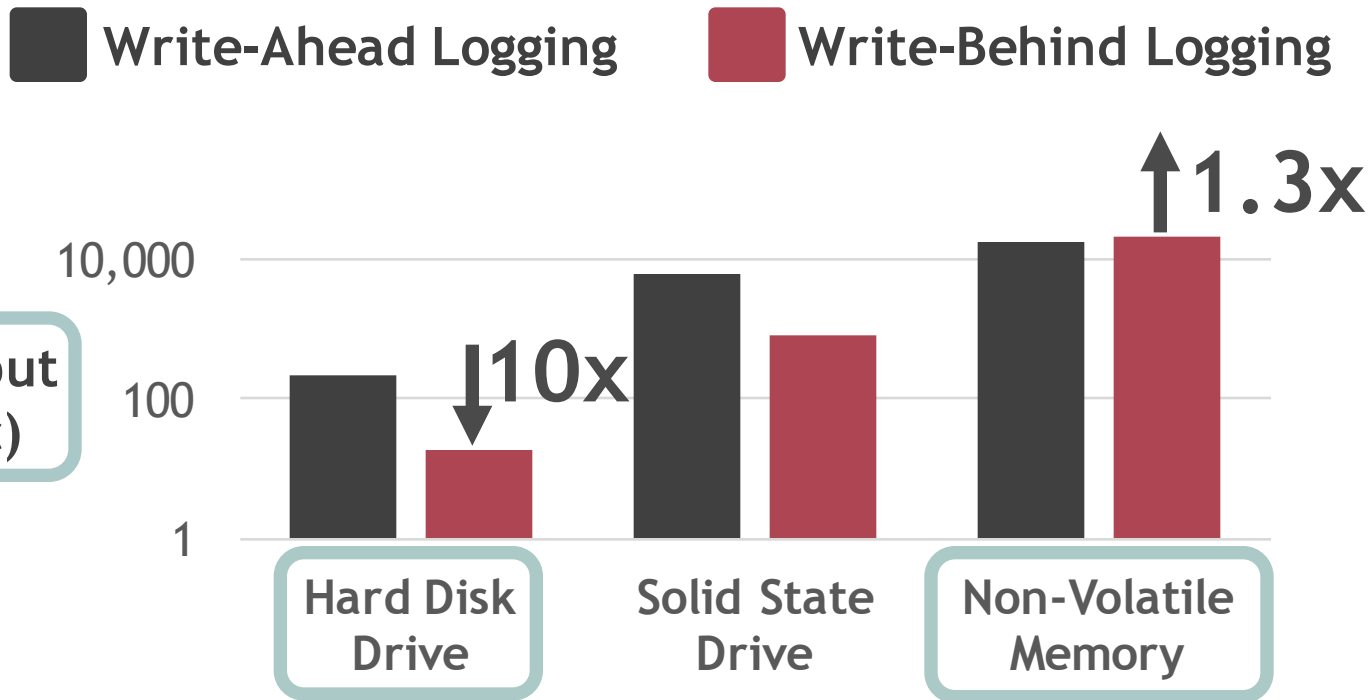
Write-Ahead Logging

Write-Behind Logging



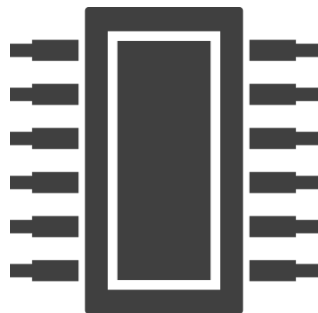


# PERFORMANCE





**PAST:  
EXISTING  
SYSTEMS**



**PRESENT:  
NVM-AWARE  
DBMS**



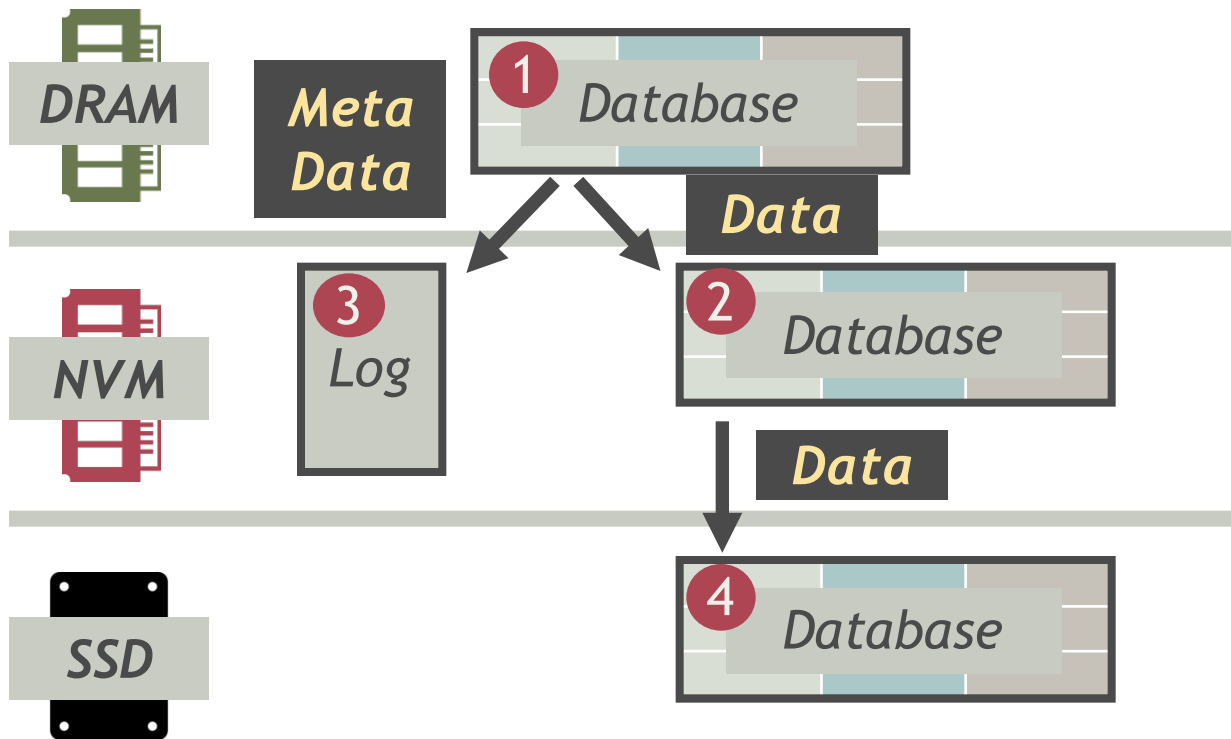
**FUTURE:  
ANALYTICS  
ON NVM**

# FUTURE — REAL-TIME ANALYTICS

- Support analytics on larger-than-NVM databases
  - *Cost of first-generation NVM devices*
- Three-tier storage hierarchy
  - *DRAM + NVM + SSD*
- When to migrate data between different layers?



# THREE-TIER STORAGE HIERARCHY

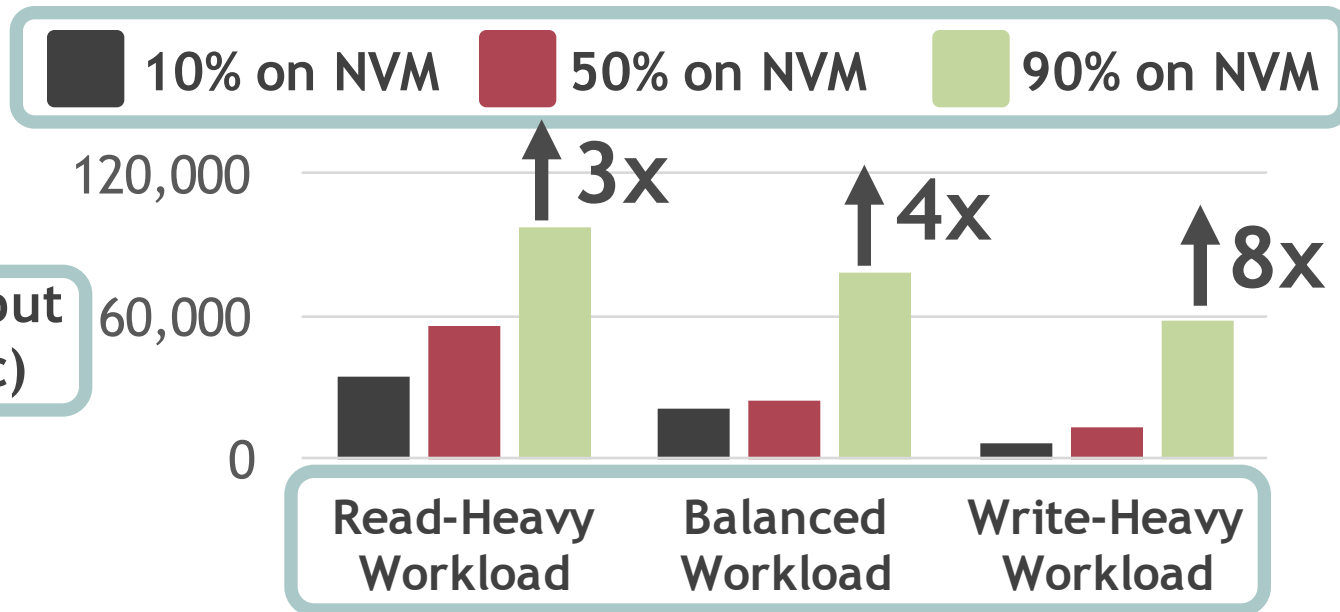


# DATA MIGRATION

---

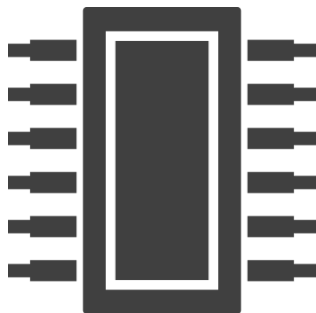
- Can directly read warm data from NVM
  - *No need to copy data over to DRAM for reading*
- Keep hottest data in DRAM
- Dynamically migrate cold data to SSD

# THREE-TIER STORAGE HIERARCHY





**PAST:  
EXISTING  
SYSTEMS**



**PRESENT:  
NVM-AWARE  
DBMS**



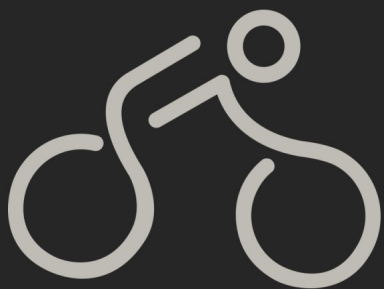
**FUTURE:  
ANALYTICS  
ON NVM**

# LESSONS LEARNED

---

- NVM outperforms SSD when correctly used
- Rethink key algorithms in database systems
  - *Write-behind logging enables instant recovery*
  - *Improves device utilization and extends its lifetime*
- Ongoing work
  - *Data placement policy*
  - *Replication*





**PELTON**

<http://pelotondb.org>



NVM Ready



Autonomous



Apache Licensed

END

@joy\_arulraj