

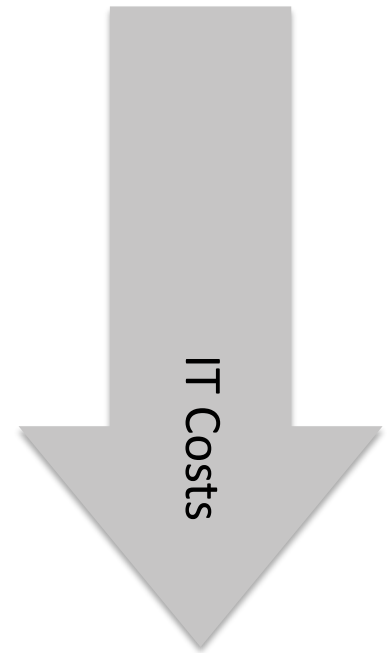
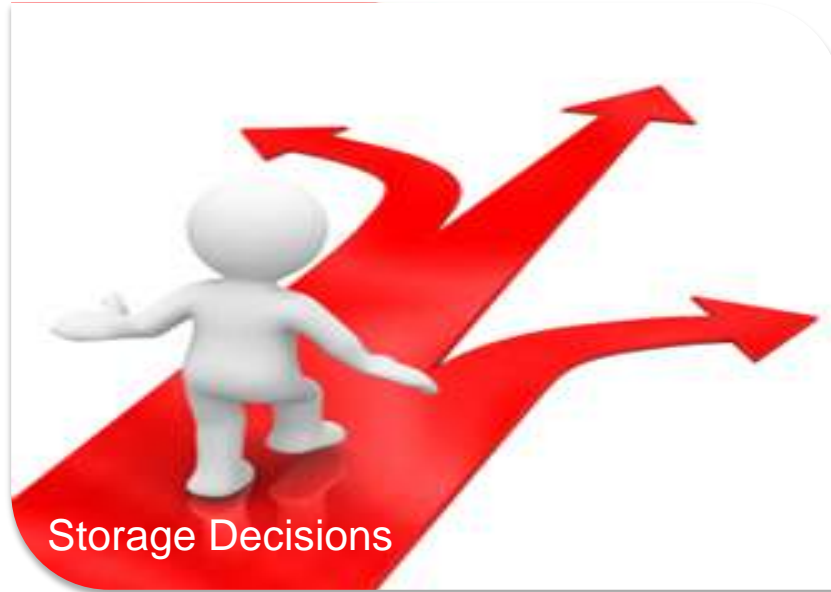
Achieving Extraordinary  
Workload Performance  
via Hyperscale  
Solid State Storage

# STORAGE DECISIONS


# BUSINESS THE WORKLOAD CHALLENGE



Strategic Objectives




# ACCELERATING WORKLOADS MEANS ACCELERATING BUSINESS




**Analytics & Intelligence**

High IOPS



**Batch Processing**

High IOPS  
High Bandwidth



**Email**




**Online Transaction Processing**

High IOPS  
Low Latency




**Video Transcoding**

High IOPS & Bandwidth




**Virtual Desktops**

High Write IOPS



**Databases Loads**

High IOPS

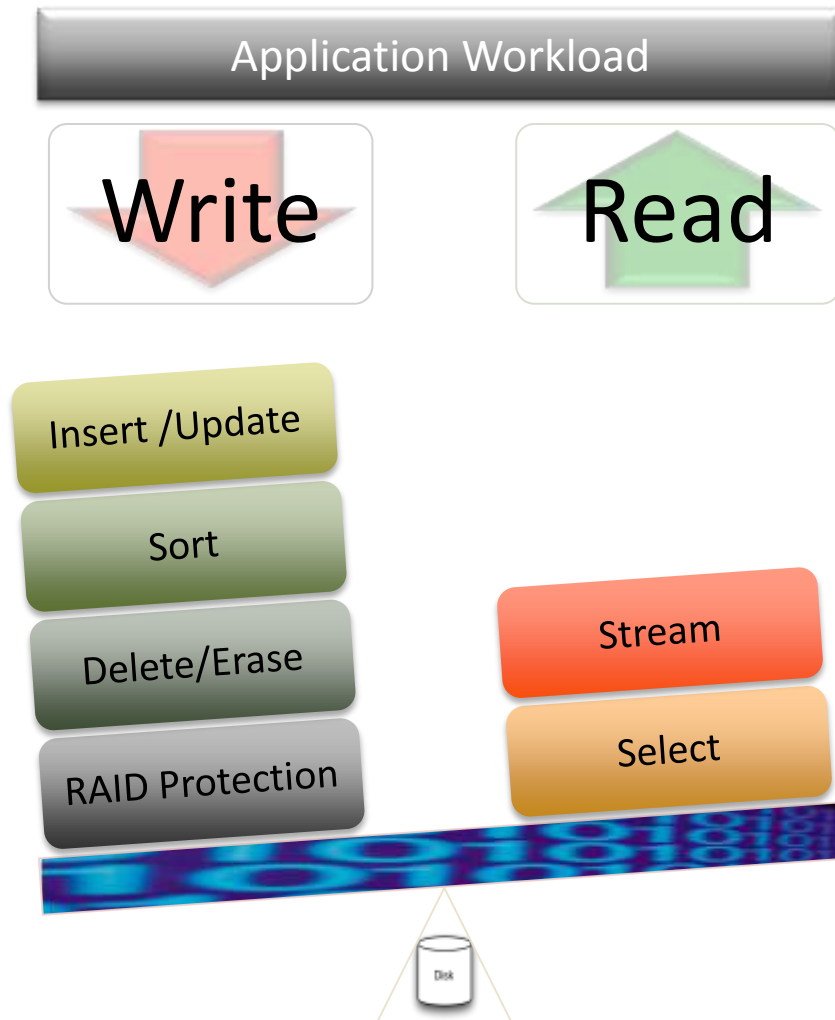


**High Performance Computing**

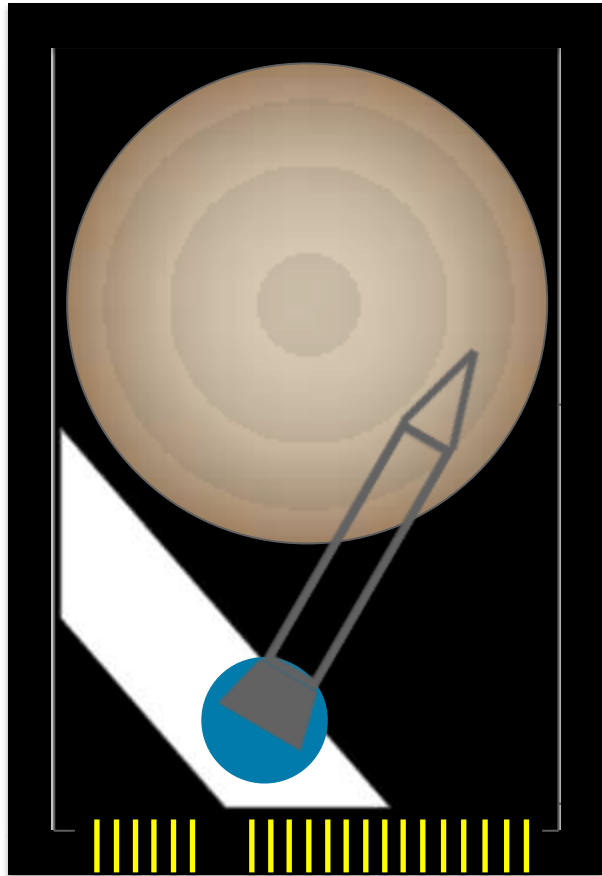
Low Latency

WRITE TO THE POINT

# THE I/O BALANCING ACT.



# HISTORICALLY HDDS DRIVE APPLICATION PERFORMANCE



## Speed

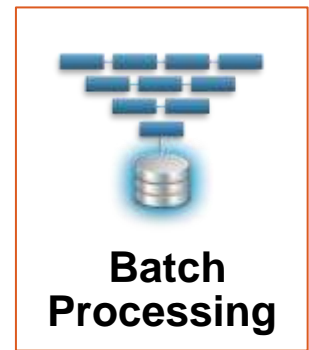
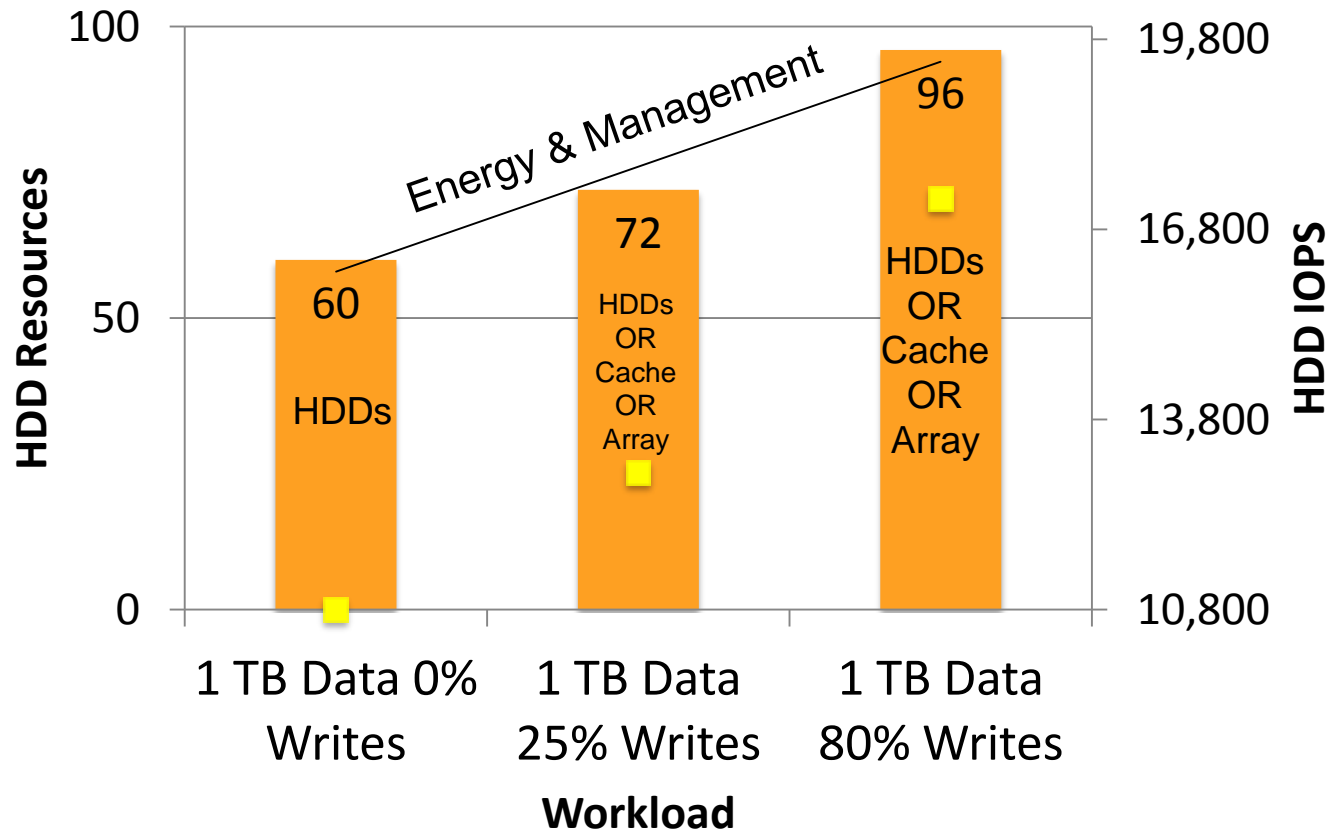
- 10s of MB/s Data Transfer Rates
- 100s of Write / Read operation per second
- .001s Latency (ms)

## Design

- Motors
- Spindles
- High Energy Consumption

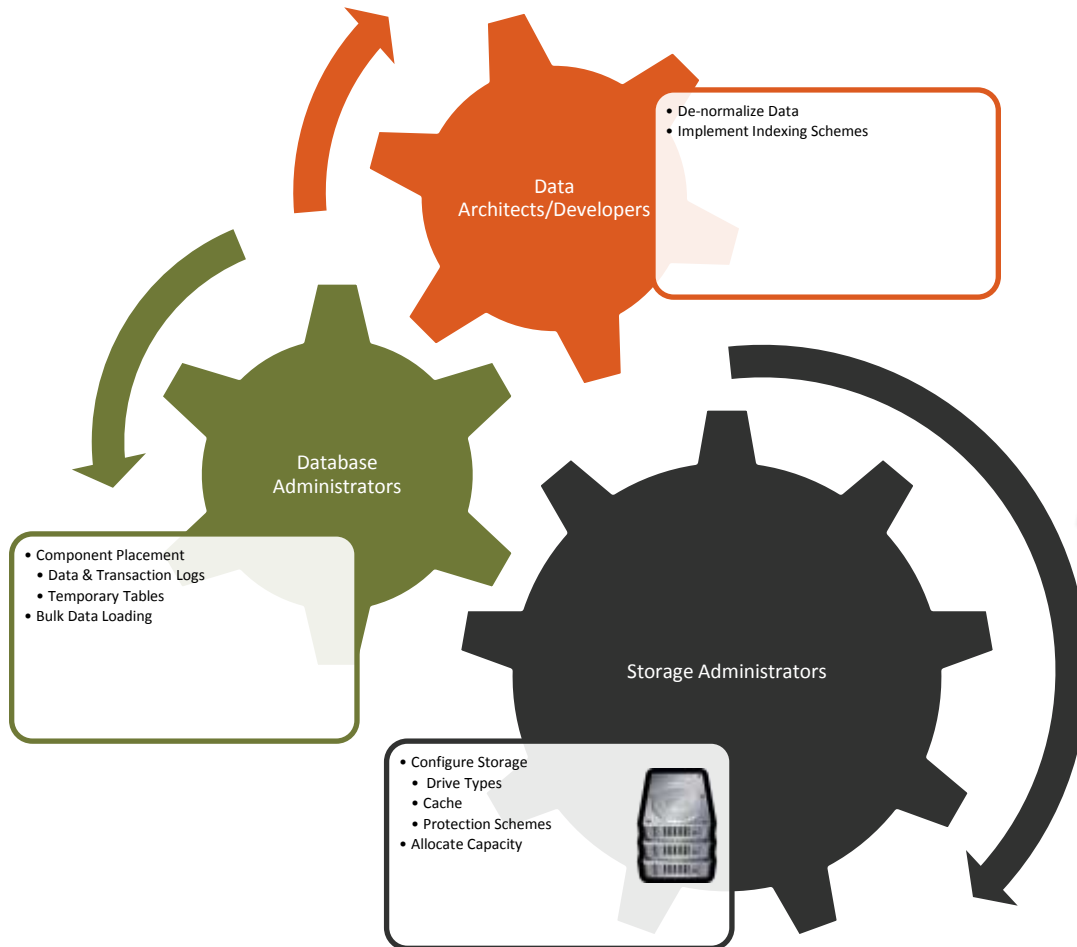
# TODAY, WORKLOAD ACCELERATION IS CONSTRAINED BY WRITES

## A More Assets Problem



- ↓ PERFORMANCE
- ↓ PRODUCTIVITY
- ↑ TOTAL COSTS

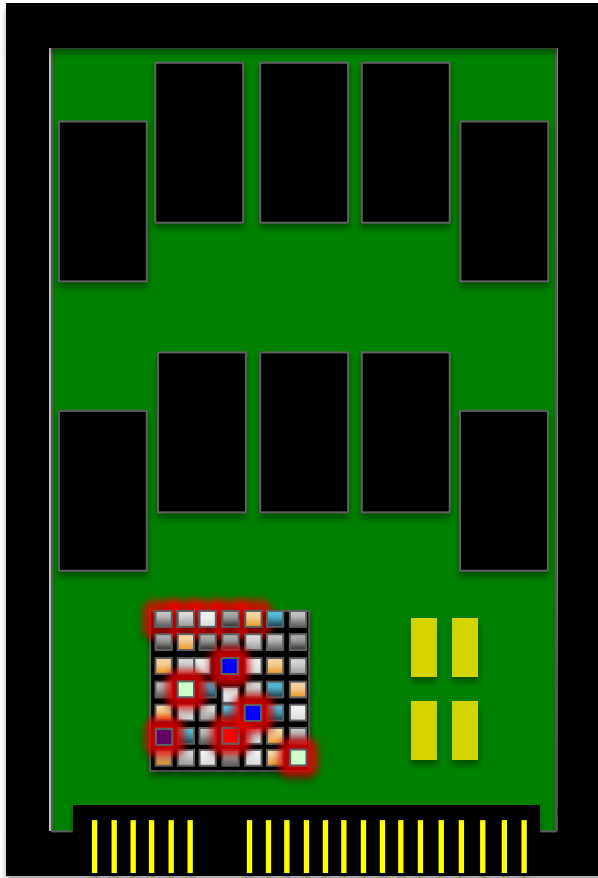
# ACCELERATING WORKLOADS DRIVES IT EFFORT



IT is focused on managing Application Performance

DO THE WRITE THING

# FLASH IMPROVES WORKLOAD PERFORMANCE



## Speed

- 100s of MB/s data transfer rates
- 1000s of Write or Read operations per second
- .000001 Latency ( $\mu$ s)

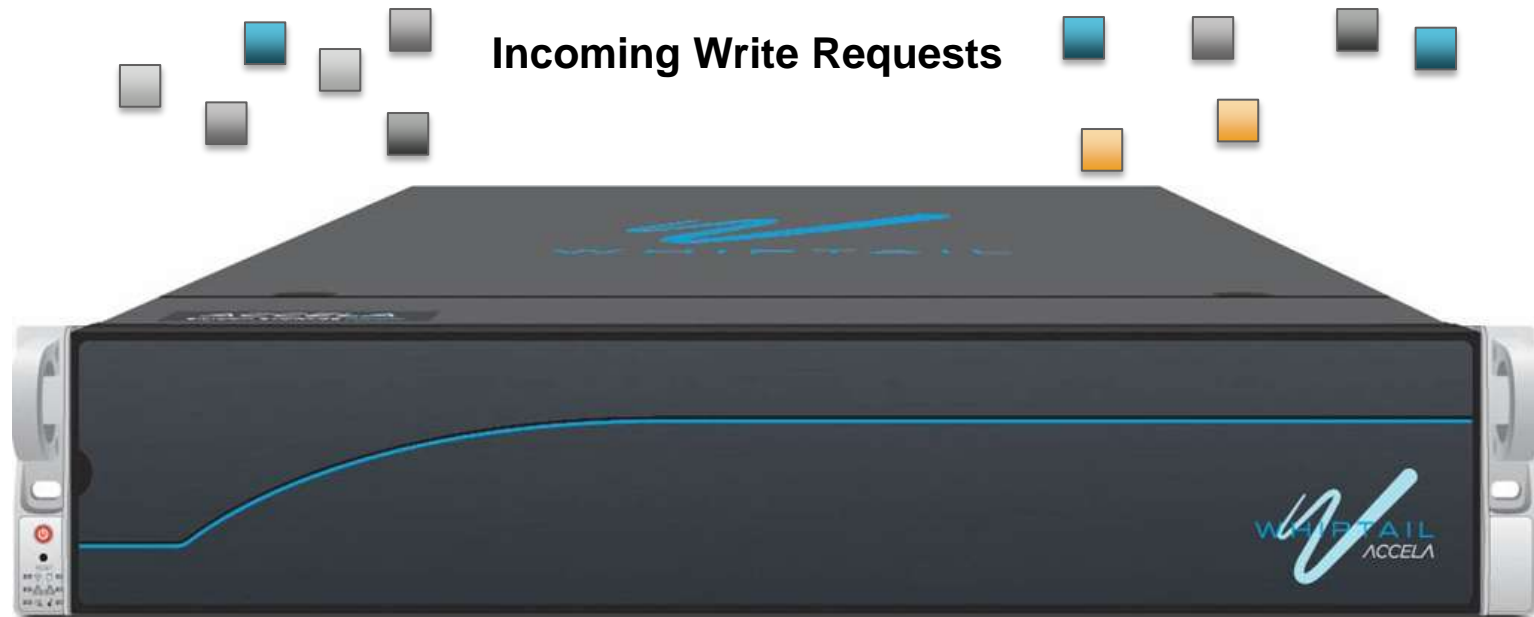
## Design

- Silicon
- NAND Flash
- Low energy consumption

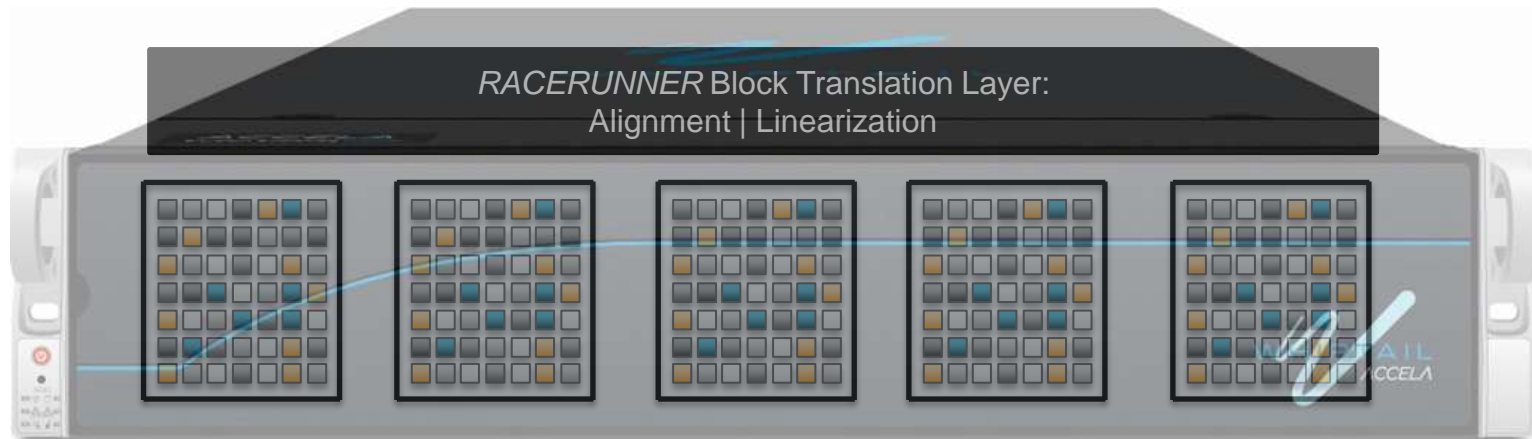
# THE BLOCK TRANSLATION LAYER MANAGING FLASH WRITES

- Mask Asymmetric Write Behavior (Page Erase)
- Logical block Addressing to Physical Mapping
- Garbage Collection
- Bad Block Management
- Align & Distribute Writes

# RACERUNNER OS: ALIGNS WRITES TO NAND & RAID

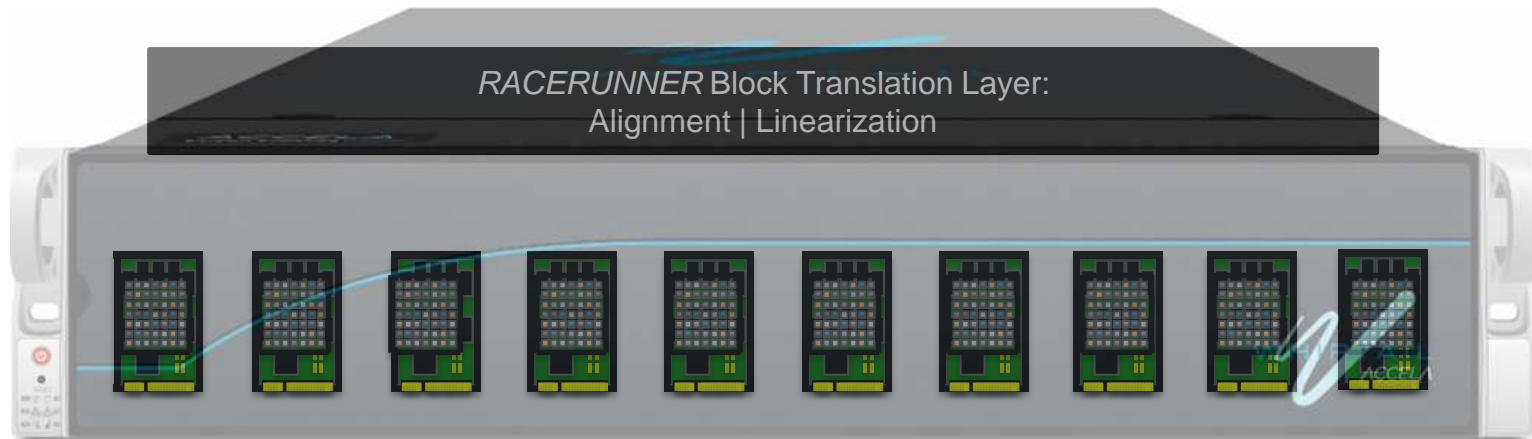


# RACERUNNER OS: ALIGNS WRITES TO NAND & RAID



Write requests are aligned in native NAND page size before passed to RAID.

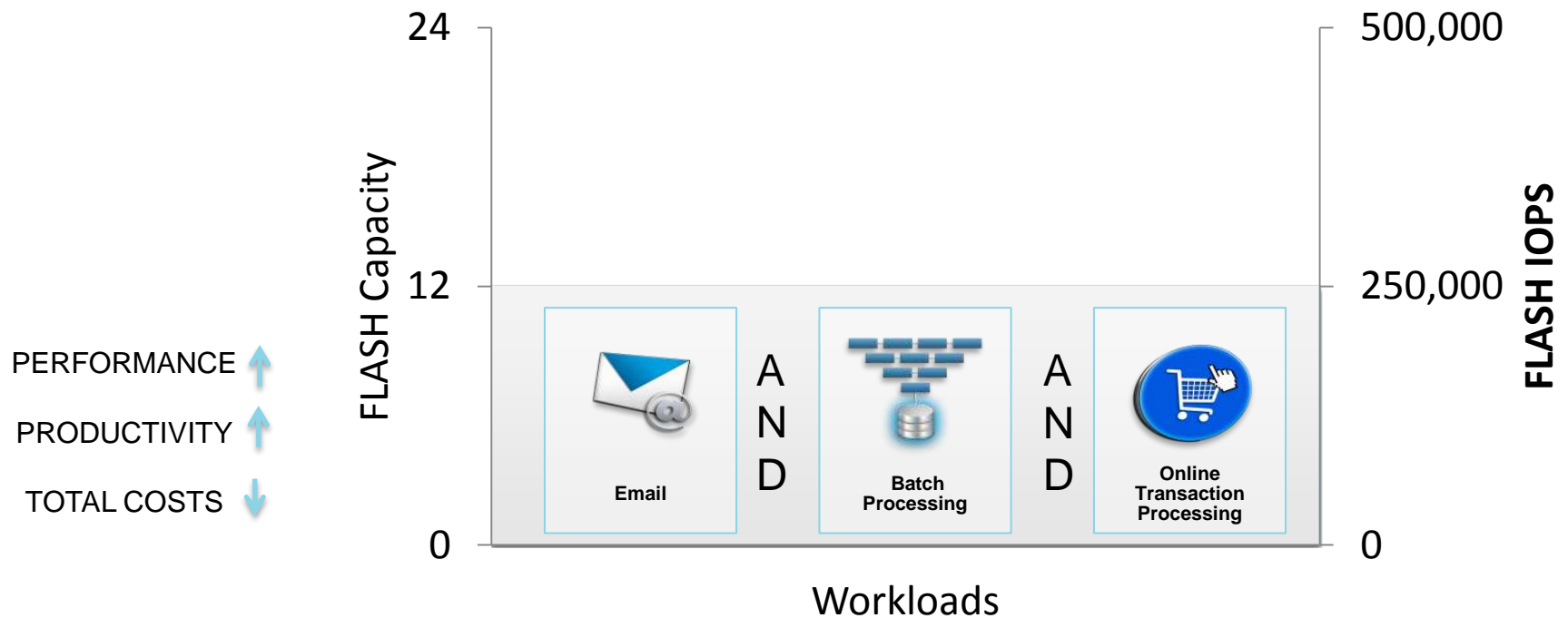
# RACERUNNER OS: ALIGNS WRITES TO NAND & RAID



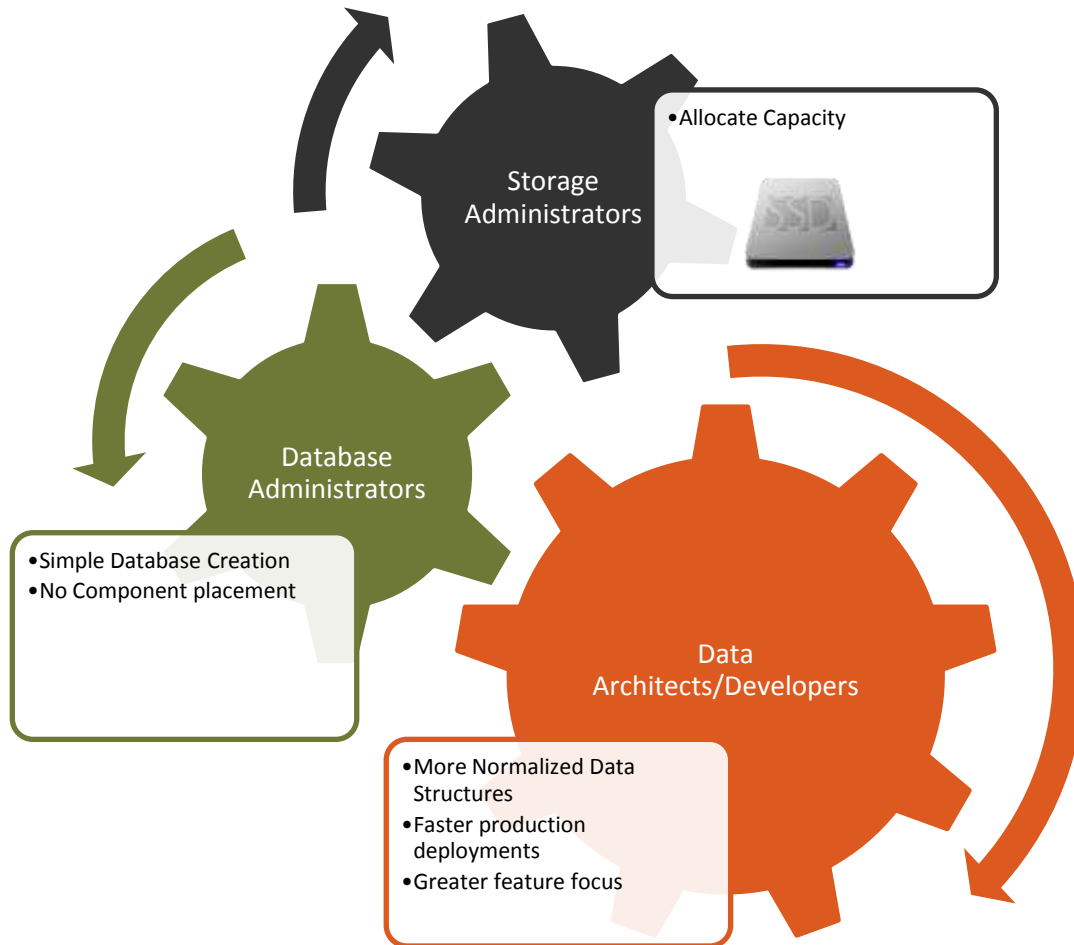
**Write requests are flushed to media as full RAID stripes.**

# REMOVE WRITE CONSTRAINTS AND ACCELERATE WORKLOADS

## A Demand Solution



# ACCELERATED WORKLOADS DRIVES BUSINESS EFFECTIVENESS



IT is Focused on Application Features

# MEASURING PERFORMANCE

# HOW DO WE ACCELERATE APPLICATION WORKLOADS?



## Increase

- **BANDWIDTH**: Sustained WRITE or READ bit-rate (GB/s or TB/h)
- **IOPS**: Input/Output Operations Per Second or Transactions
- **SYMMETRY** between Writes and Reads
- **WORKLOAD CONSOLIDATION**:



## Decrease

- **LATENCY**: Milliseconds to Microseconds
- **ENERGY** : Power and Cooling.
- **CAPACITY**: eliminate overprovisioning
- **WORKLOAD TUNING**:

## Velocity

- Near time
- Real time
- Streams
- Batch

## WHIPTAIL Labs

Reengineer Workloads for Solid State Storage

Workload Reference Architectures

Workload Benchmarks

Workload Best Practices

## Product Engineering

Solid State Storage platforms

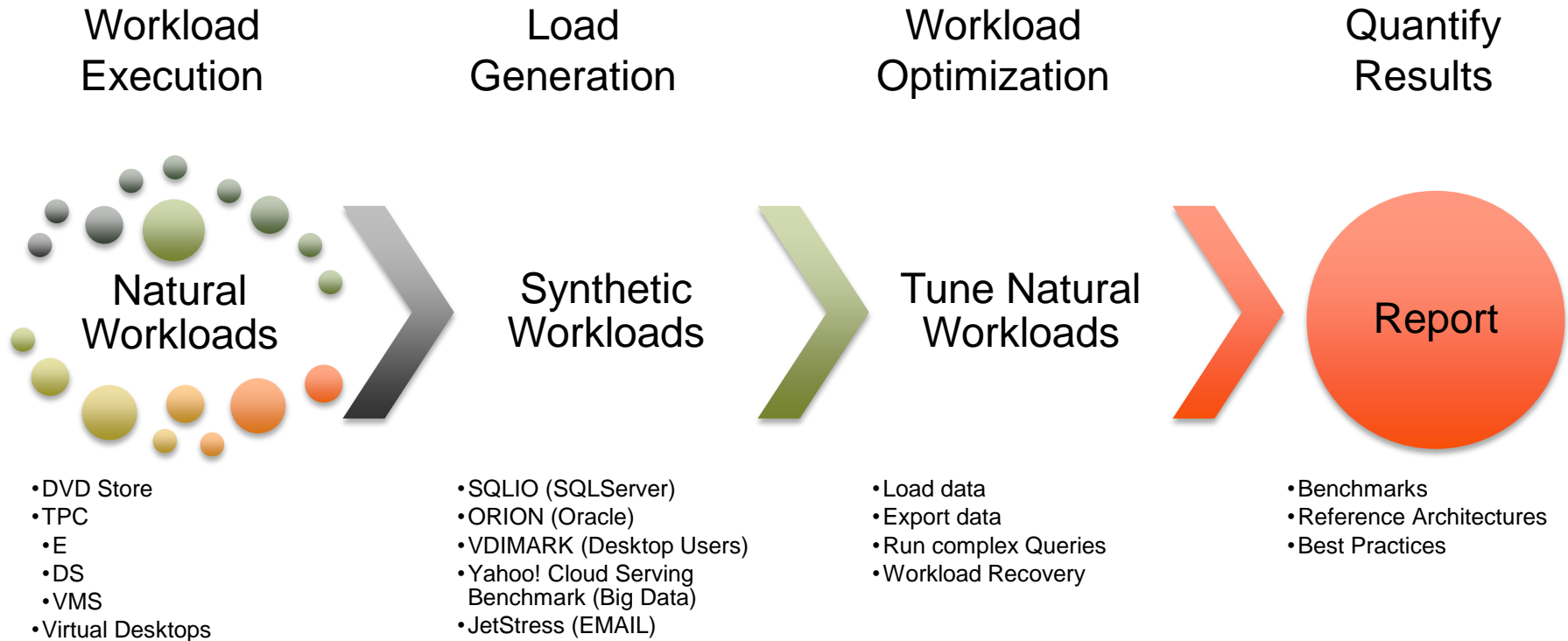
Data management features

## Research & Development

New advancements in Solid State Storage

New data management techniques

# WORKLOAD TESTING PROCESS



## Develop Prototypes



- Reference Architectures
- Benchmarks

## Scale Testing



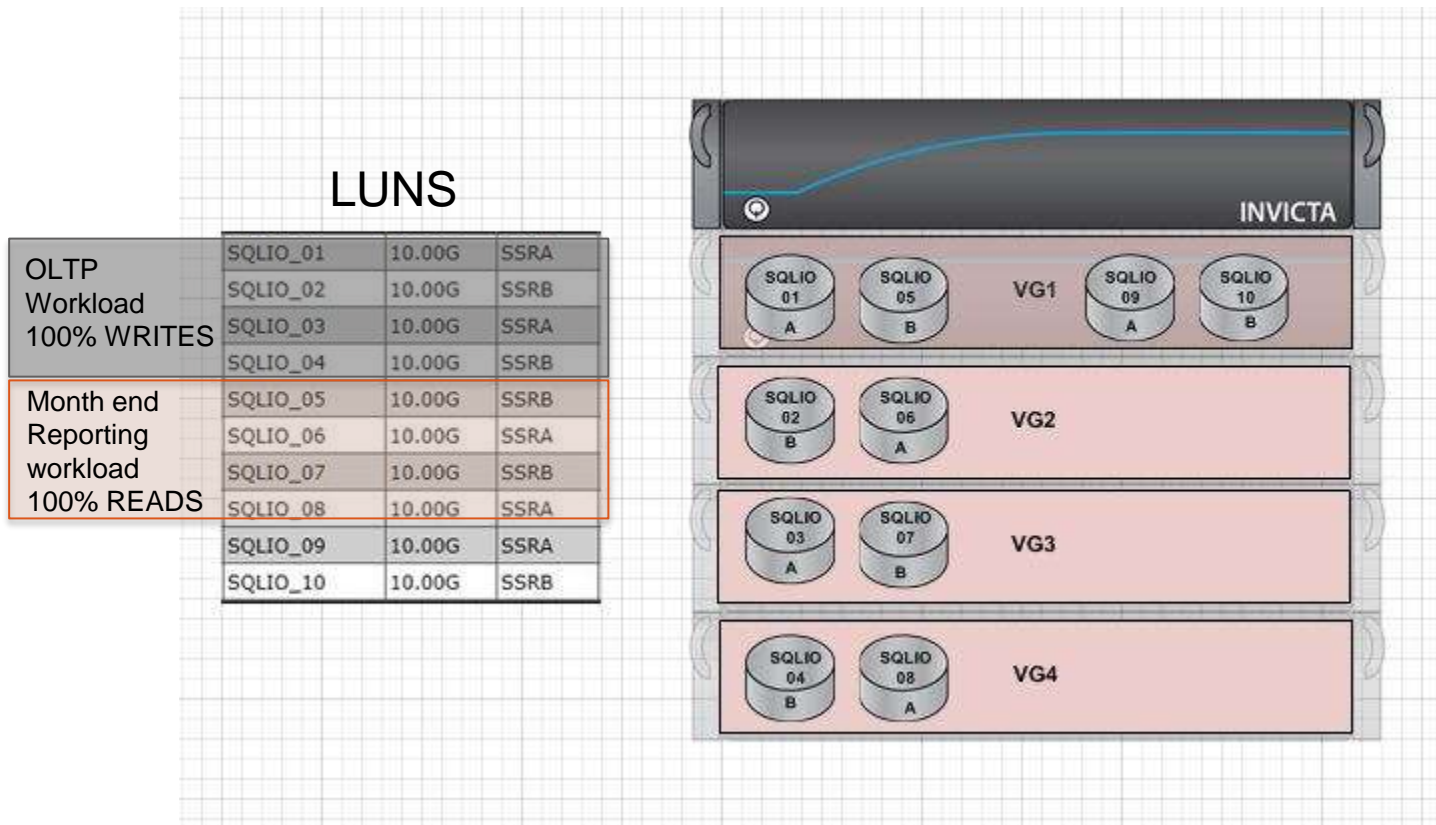
- Architectures
- Benchmarks
- Compare / Contrast

## Validate



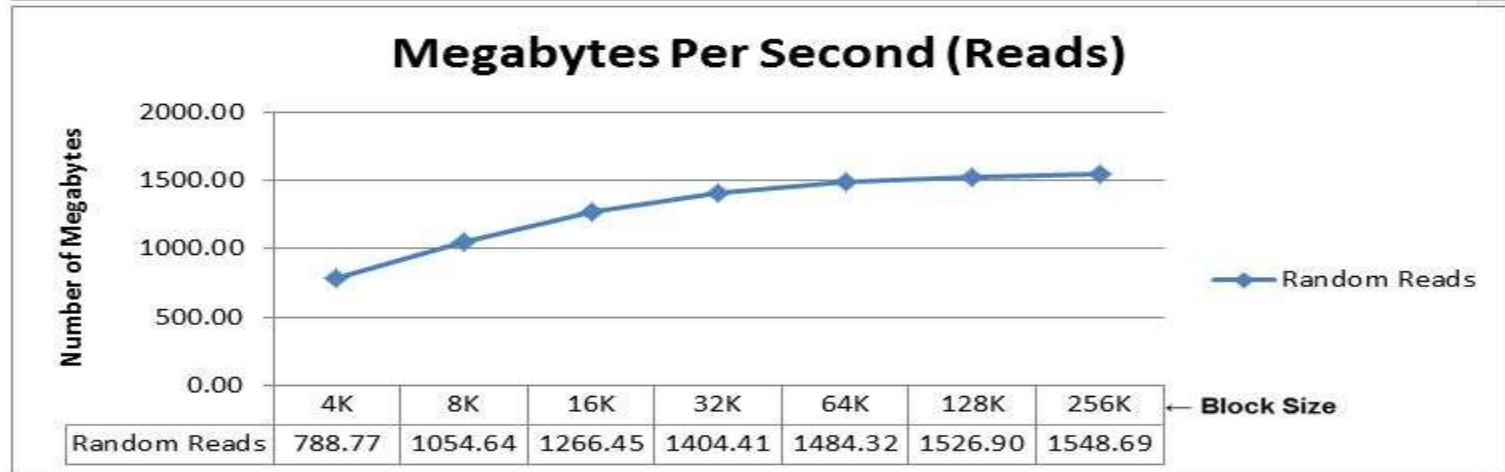
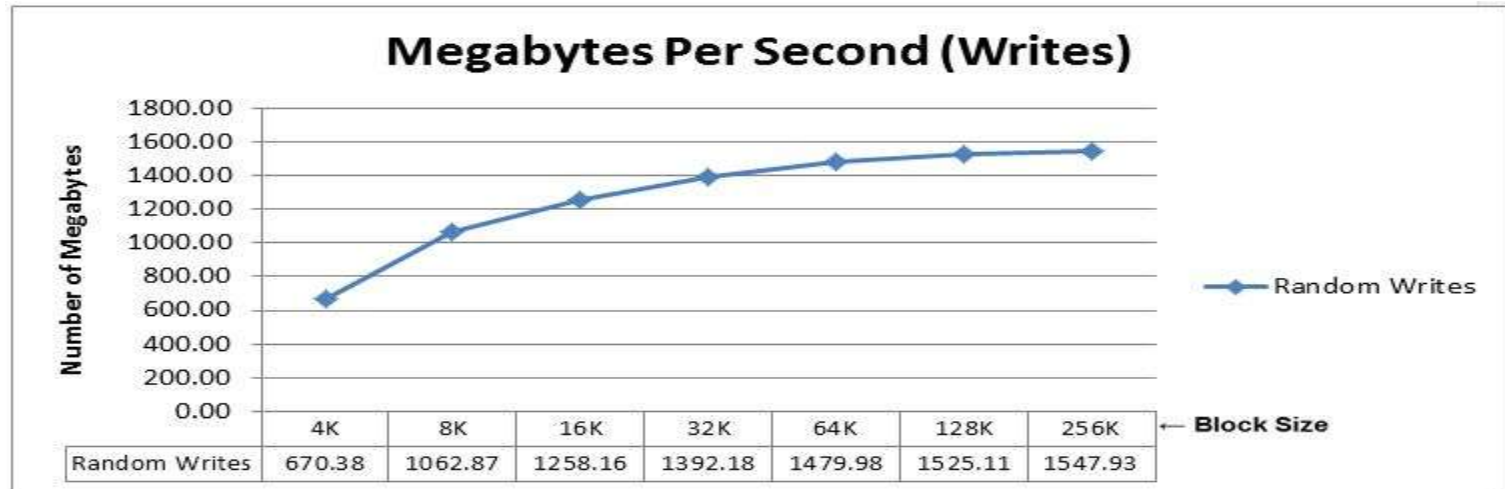
**Audit  
Results**

# MULTI-WORKLOAD TESTING



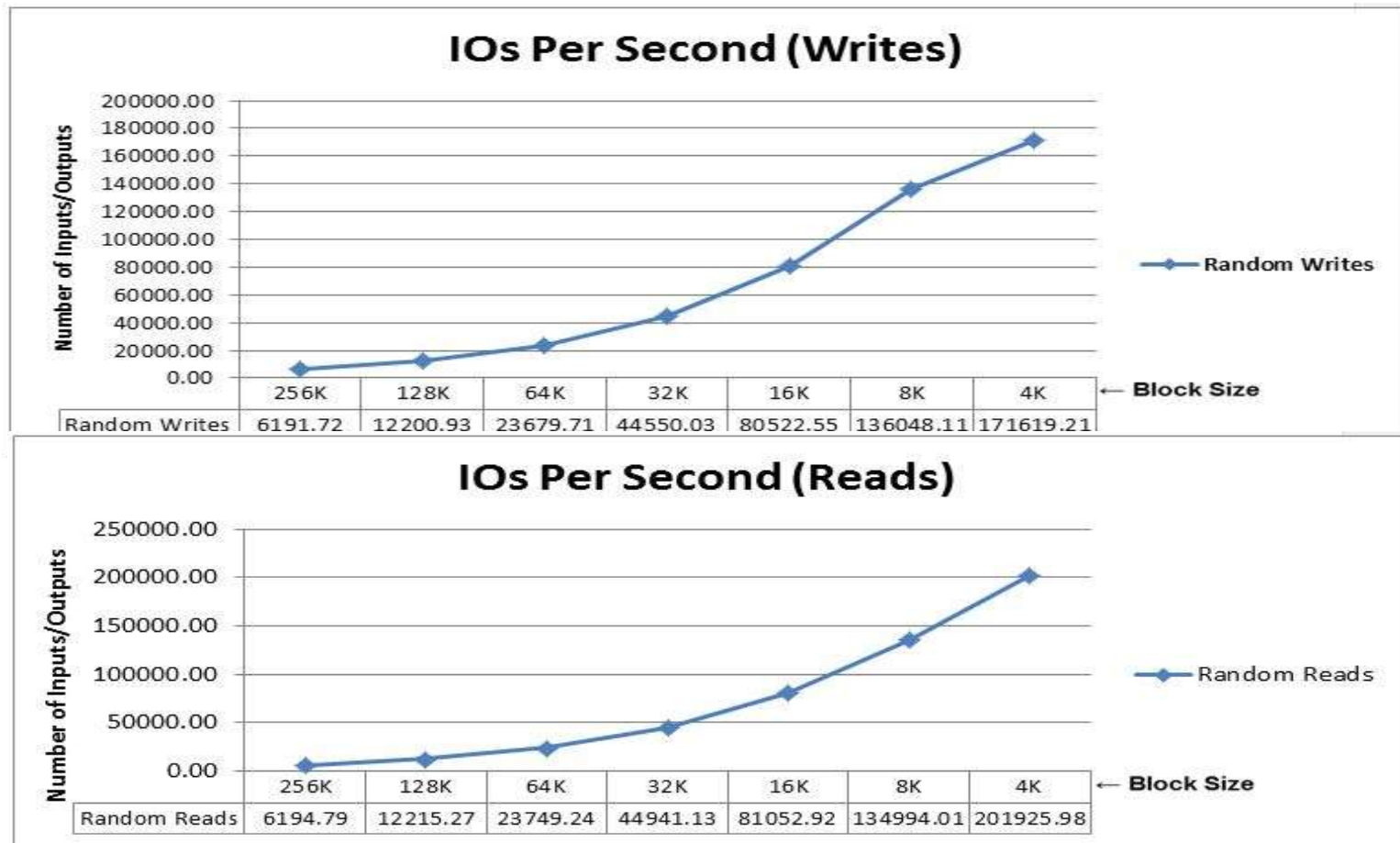
# BANDWIDTH TEST RESULTS

## OLTP & REPORTING WORKLOADS



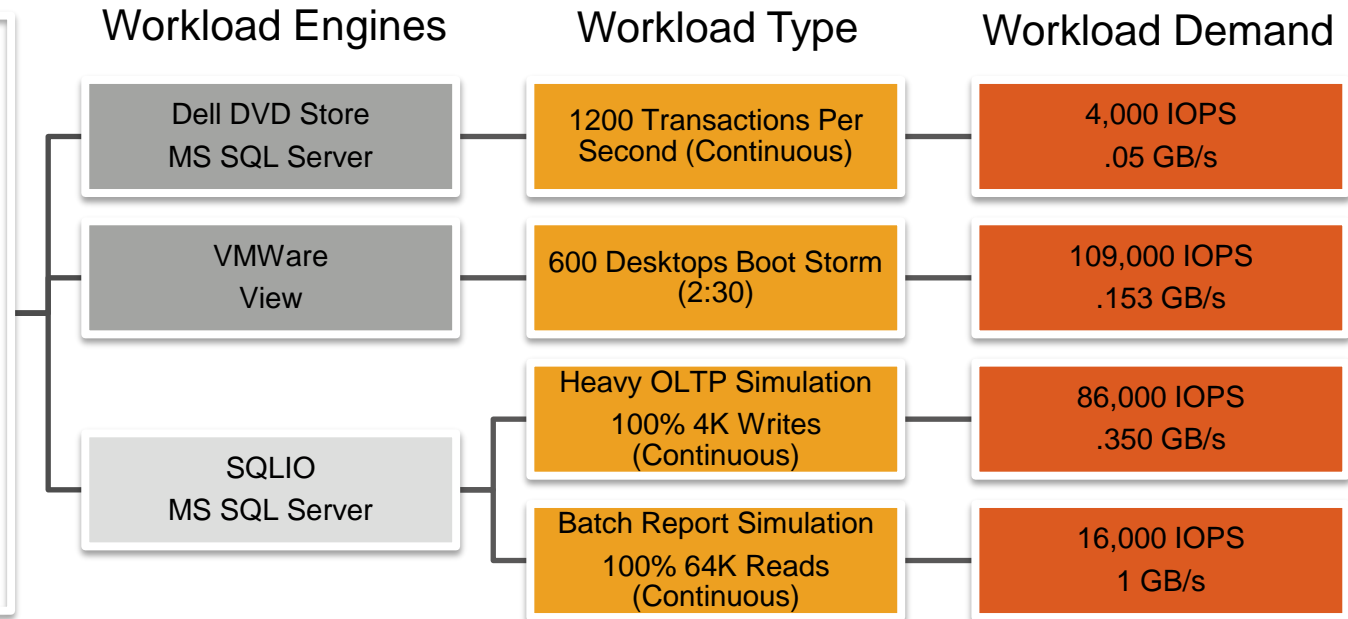
# IO TEST RESULTS

## OLTP & REPORTING WORKLOADS



# MULTI-WORKLOAD REFERENCE ARCHITECTURE

## Mercury



- **INVICTA**
  - 350,000 IOPS
  - 3.5 GB/s
  - 18 TB
- 8 Servers

In 2012 Mercury traveled to Barcelona, New York, San Francisco, Santa Clara, and Seattle demonstrating the ability to accelerate multiple workloads on to Solid State Storage.

**215,000 IOPS**  
**1.553 GB/s**

Raid 5 HDD Equivalent = 3,800  
RAID 10 HDD Equivalent = 2,000



# PERFORMANCE BENCHMARK

## THE CARL SAGAN RUN



### Create 4.225 Billion Records

- Randomly generated census like data
- 10 fields
- 85 bytes per record



### Load records into Oracle Database

- A single 1.5TB LUN for table, temp and logs
- Tablespace .365 TB
- No Indexes
- 76,000 Writes per second
- 1.2 GB/s



### Generate additional loads via ORION

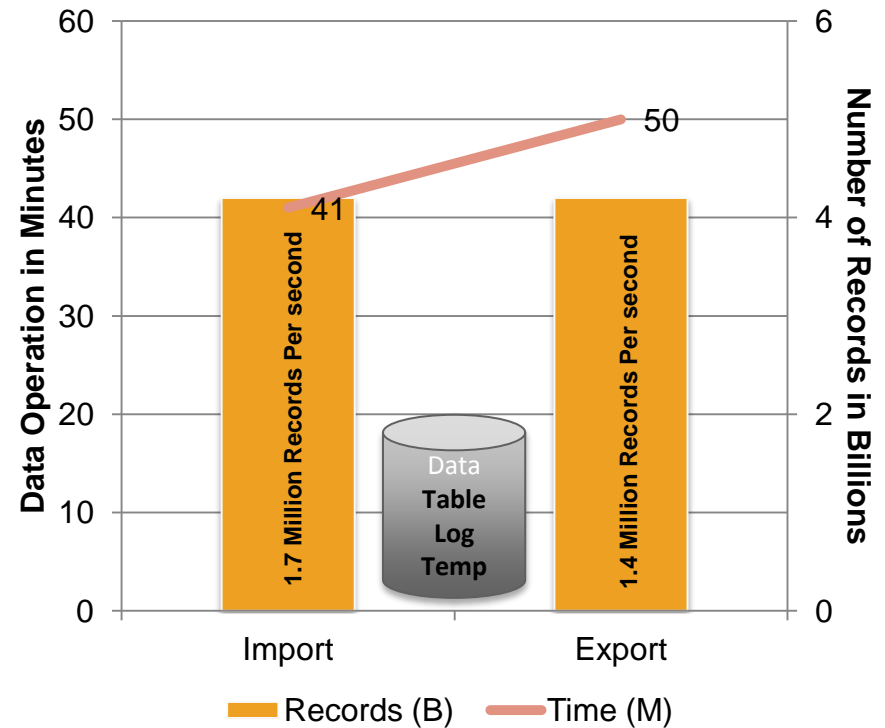
- Run Heavy OLTP Simulation
- 70,000 Writes/Reads per second (40/60)
- 1 GB/s



### INVICTA

- 4 Nodes
- 24 TB Capacity
- 425,000 IOPS

### The Carl Sagan Run Results Minutes



# WORKLOADS IN THE QUEUE



Document Oriented



Key Value Pairs



Reference Architecture



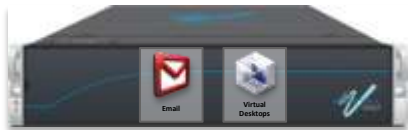
Data Warehouse  
Reference Architecture





# PRODUCT LINEUP TO INFINITY AND BEYOND

1-9 Application Workloads



**ACCELA**  
1.5TB – 12TB  
250,000 IOPS  
1.9 GB/s Bandwidth



→  
**Scalability Path**

1-60 Application Workloads



**INVICTA**  
2-6 Nodes  
6TB-72TB  
650,000 IOPS  
7GB/s Bandwidth

1-300 Application Workloads

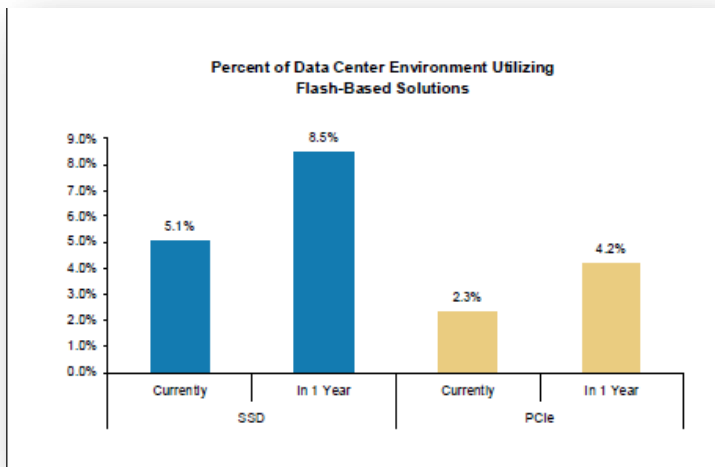


**INVICTA – INFINITY (Q1/13)**  
7-30 Nodes  
86TB-360TB  
800,000 – 4 Million IOPS  
40GB/s Bandwidth



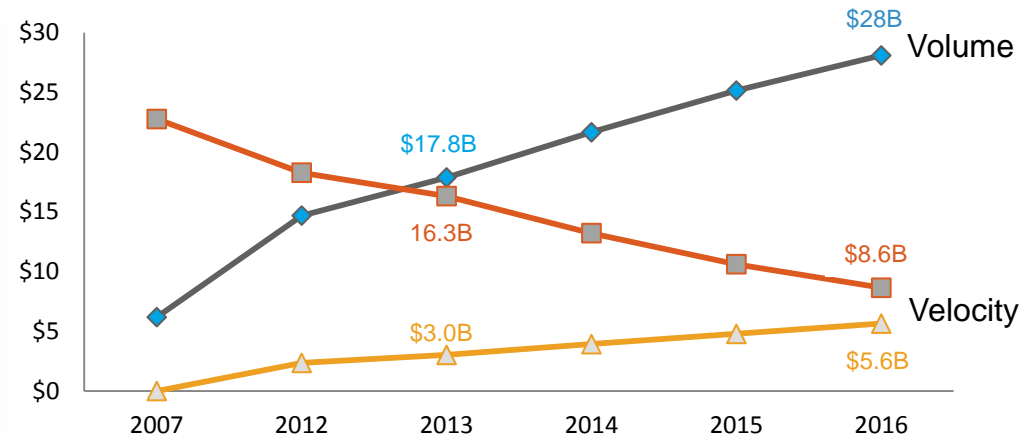
# SOLID STATE IS QUICKLY GAINING FAVOR WORKLOAD PERFORMANCE

## Data Center Adoption



Source Morgan Stanley  
April 2012

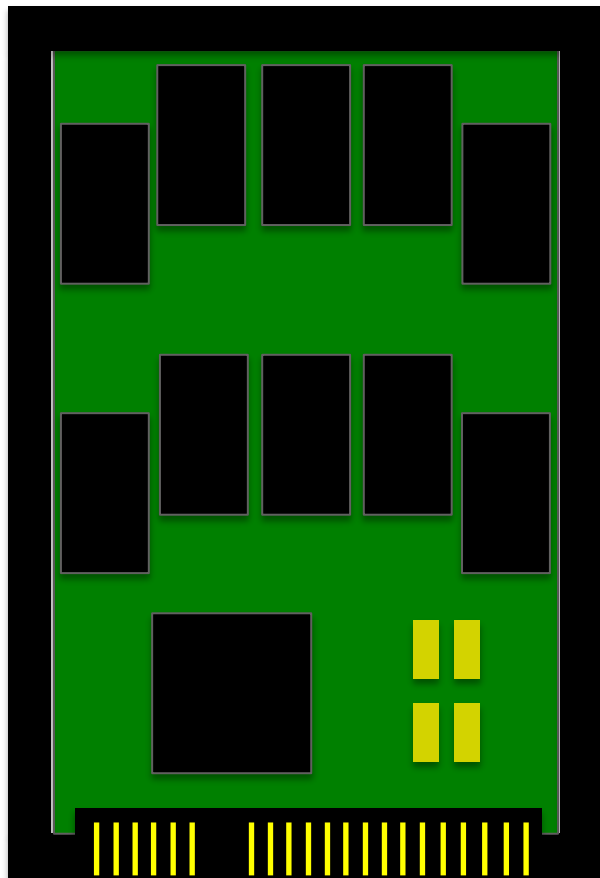
## Storage Systems Revenue (\$B)



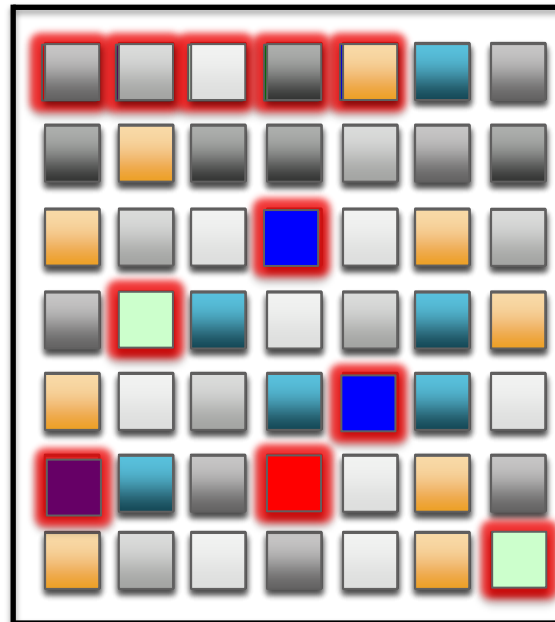
Source IDC  
March 2012

—◆— Slow HDD Revenue (\$B)  
—■— Fast HDD Revenue (\$B)  
—▲— SSD Revenue (\$B)

# NAND FLASH FUNDAMENTALS: *FLASH WRITE PROCESS*



2MB NAND Page

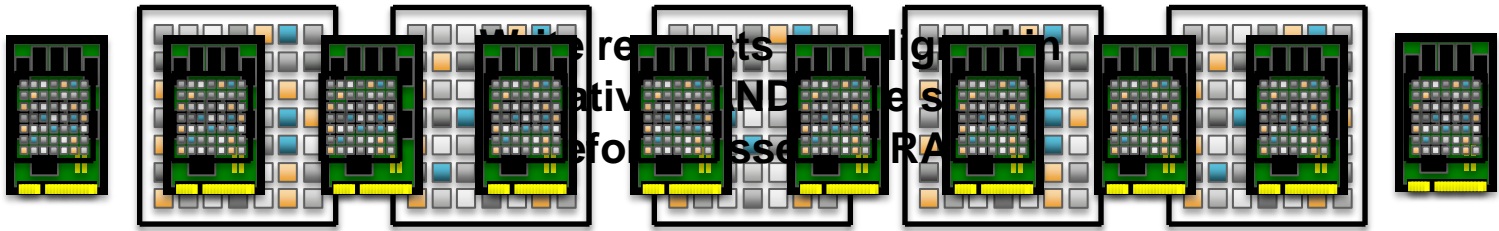


1. NAND Page contents are read to a buffer.
2. NAND Page is erased (aka, "flashed").
3. Buffer is written back with previous data and any changed or new blocks – including zeroes.

# RACERUNNER OS: *ALIGNING WRITES TO NAND*



RaceRunner Block Translation Layer:  
Alignment | Linearization



Write requests are flushed to media as full RAID stripes.