# Handling Real-Time Use Cases at Scale Using a Hybrid-Memory Architecture

Srini V. Srinivasan

Founder, Chief Development Officer
Aerospike

*In-Memory Computing Summit*
*Silicon Valley*
*October 25, 2017*

AEROSPIKE

# Use Cases

# Fraud Prevention for Interactive Payments
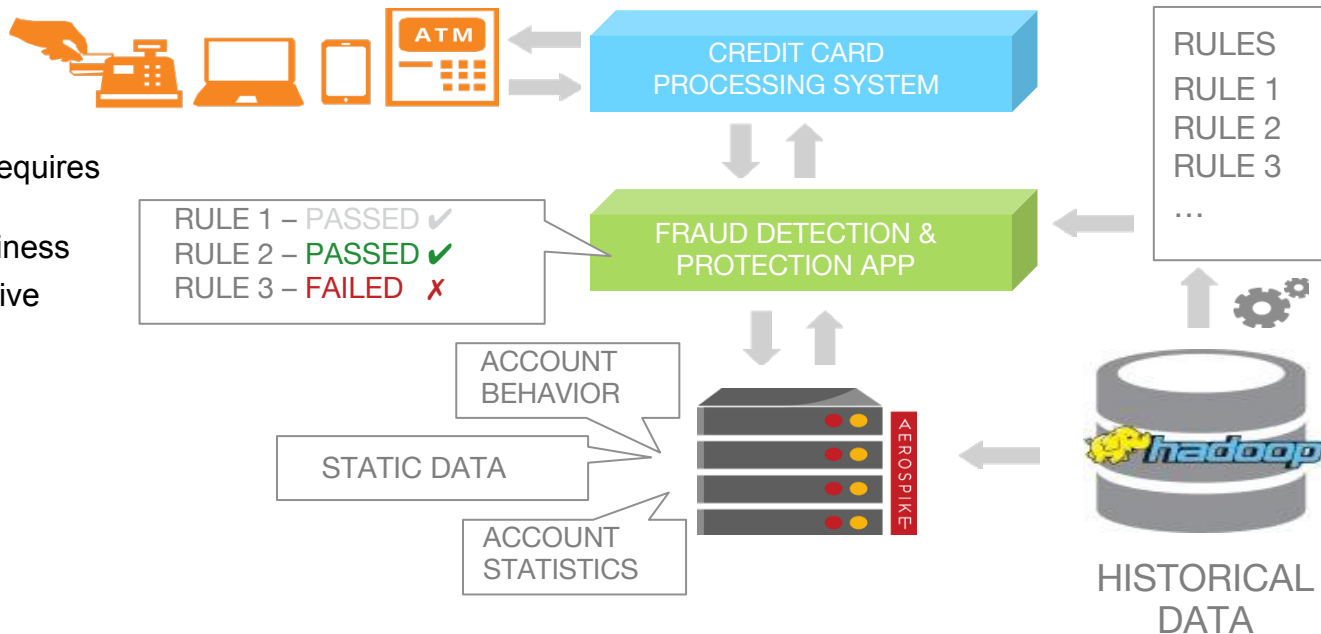
## Business Challenge

- Every payment transaction requires hundreds of DB reads/writes
- Missed latency SLA lost business
- Caching solution too expensive

## Need to scale up

- 10 → 100 TB
- 10B → 100 B objects
- 200k → I Million+ TPS

## Selected Hybrid Memory

- Built for Flash – eliminated inconsistencies
- Predictable Low latency at High Throughput
- 20 Server Cluster reduced from 150 in-memory cache servers

CREDIT CARD PROCESSING SYSTEM

RULES
RULE 1
RULE 2
RULE 3
…

RULE 1 – PASSED ✔
RULE 2 – PASSED ✔
RULE 3 – FAILED ✗

FRAUD DETECTION & PROTECTION APP

ACCOUNT BEHAVIOR

STATIC DATA

ACCOUNT STATISTICS

HISTORICAL DATA

# Retail Banking Positions – Risk Management

## Business Challenge

- Must update stock prices, show balances on 300 positions
- process 250M transactions, 2 M updates/day
- Calculate risk metrics on portfolios on a continuous basis
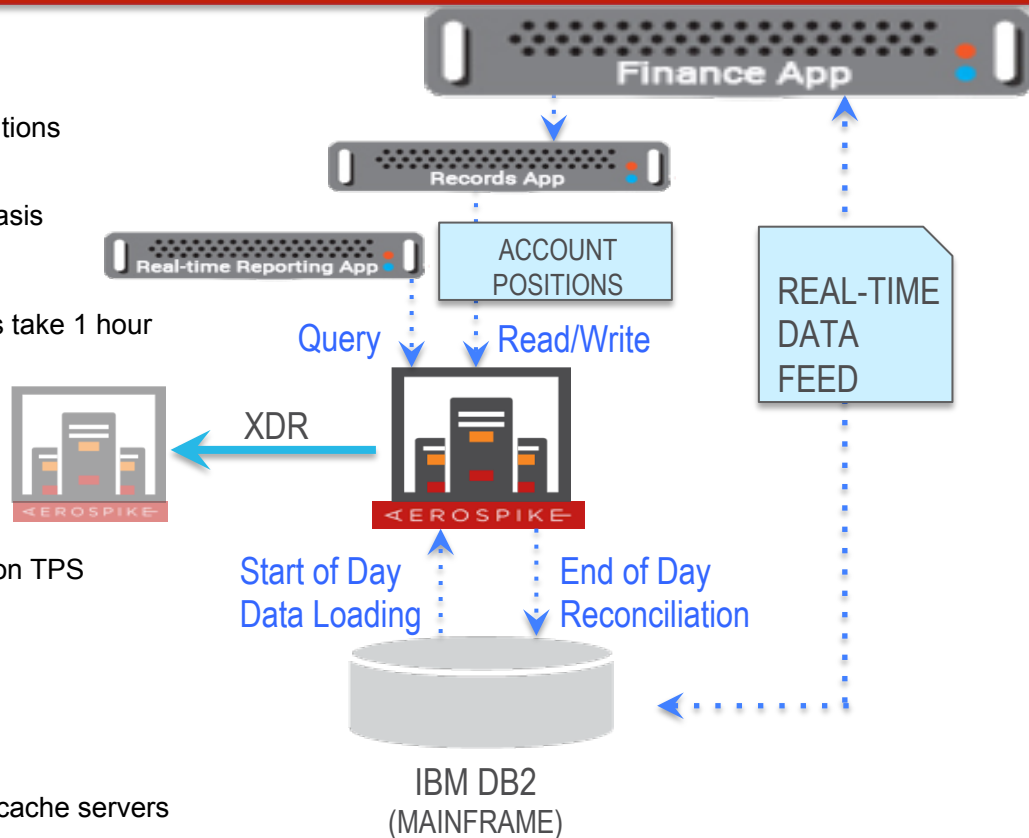
## Caching solution failed

- Running out of memory, data inconsistencies, restarts take 1 hour
- 150 Servers -> Growing to 1000

## Need to scale business

- 3 → 13 TB, 100 → 400 Million objects, 200k → I Million TPS

## Hybrid Memory Advantage

- Built for Flash – eliminated inconsistencies
- Predictable Low latency at High Throughput
- 10-12 Server Cluster – reduced from 150 in-memory cache servers

Finance App

Records App

Real-time Reporting App

ACCOUNT POSITIONS

REAL-TIME DATA FEED

Query

Read/Write

XDR

Start of Day Data Loading

End of Day Reconciliation

IBM DB2 (MAINFRAME)

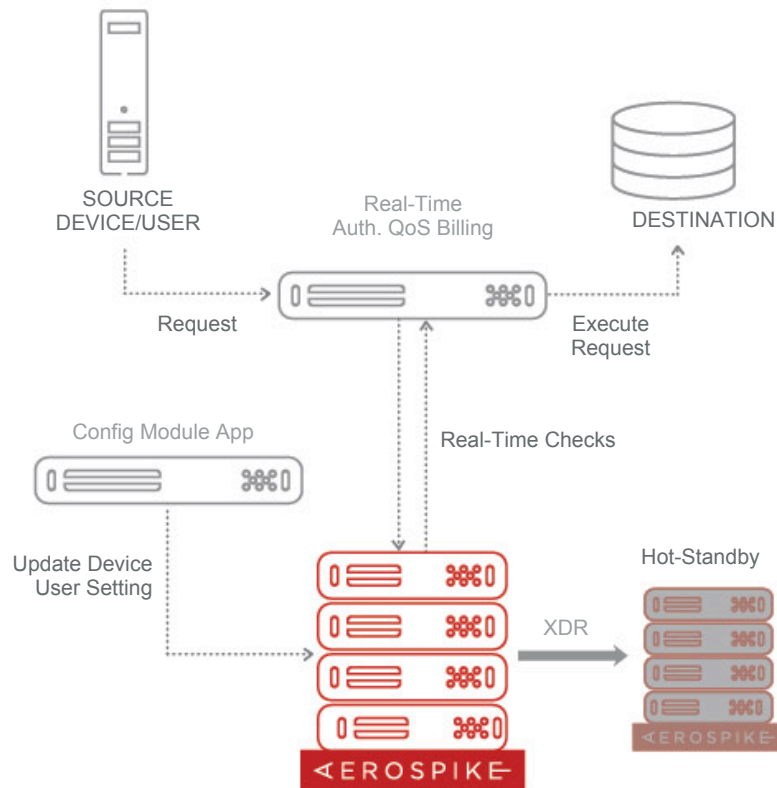# Telco – Billing and Charging

## Challenge

- Edge access to regulate traffic
- Accessible using provisioning applications (self-serve and through support personnel)
- Ensure accuracy in billing and charging
- Quick turn around for provisioning changes

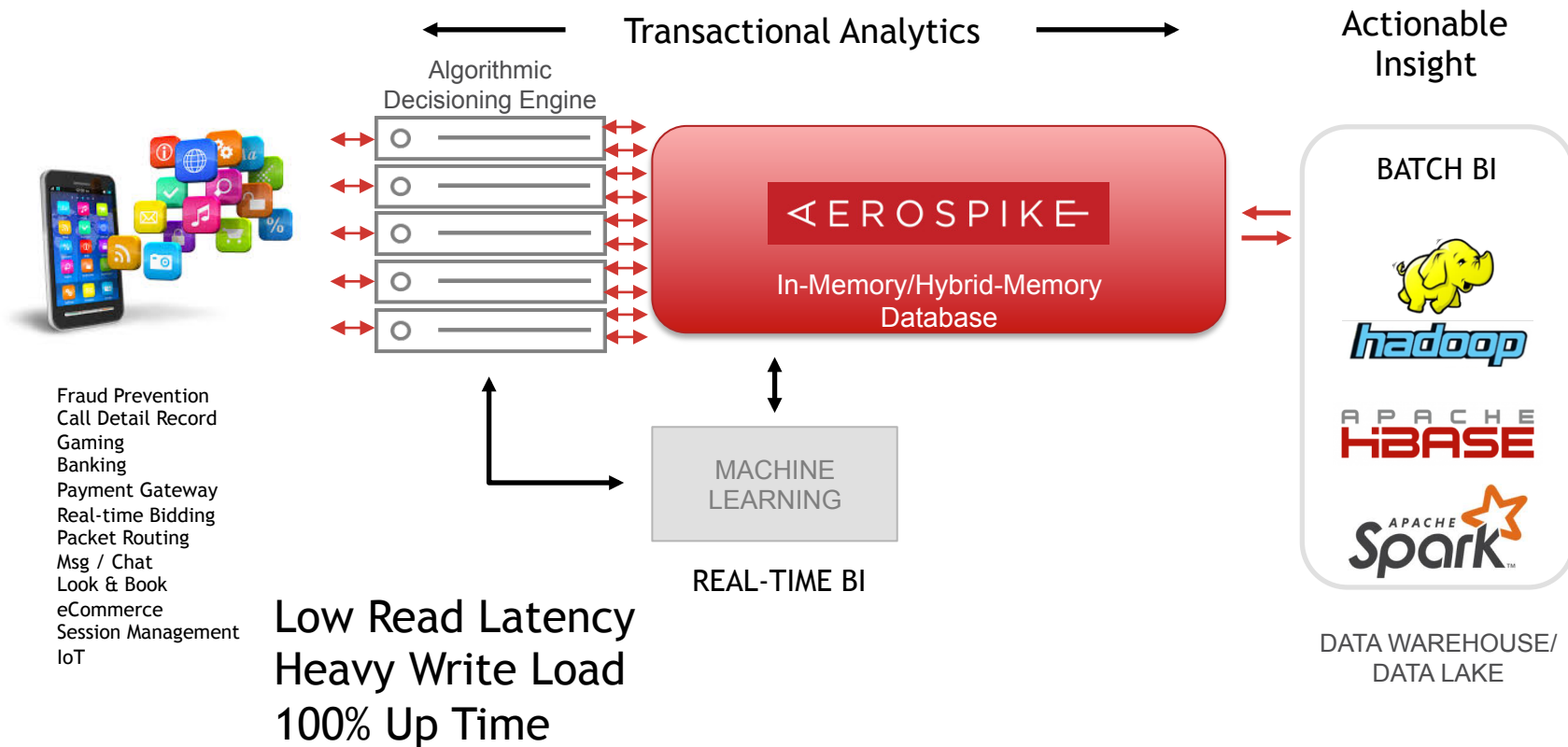## Need Extremely High Availability, Reliability, Low latency

- > TBs of data
- 10-100M objects
- 10-200K TPS

## Selected Aerospike

- Clustered system
- Predictable low latency at high throughput
- Highly-available and reliable on failure
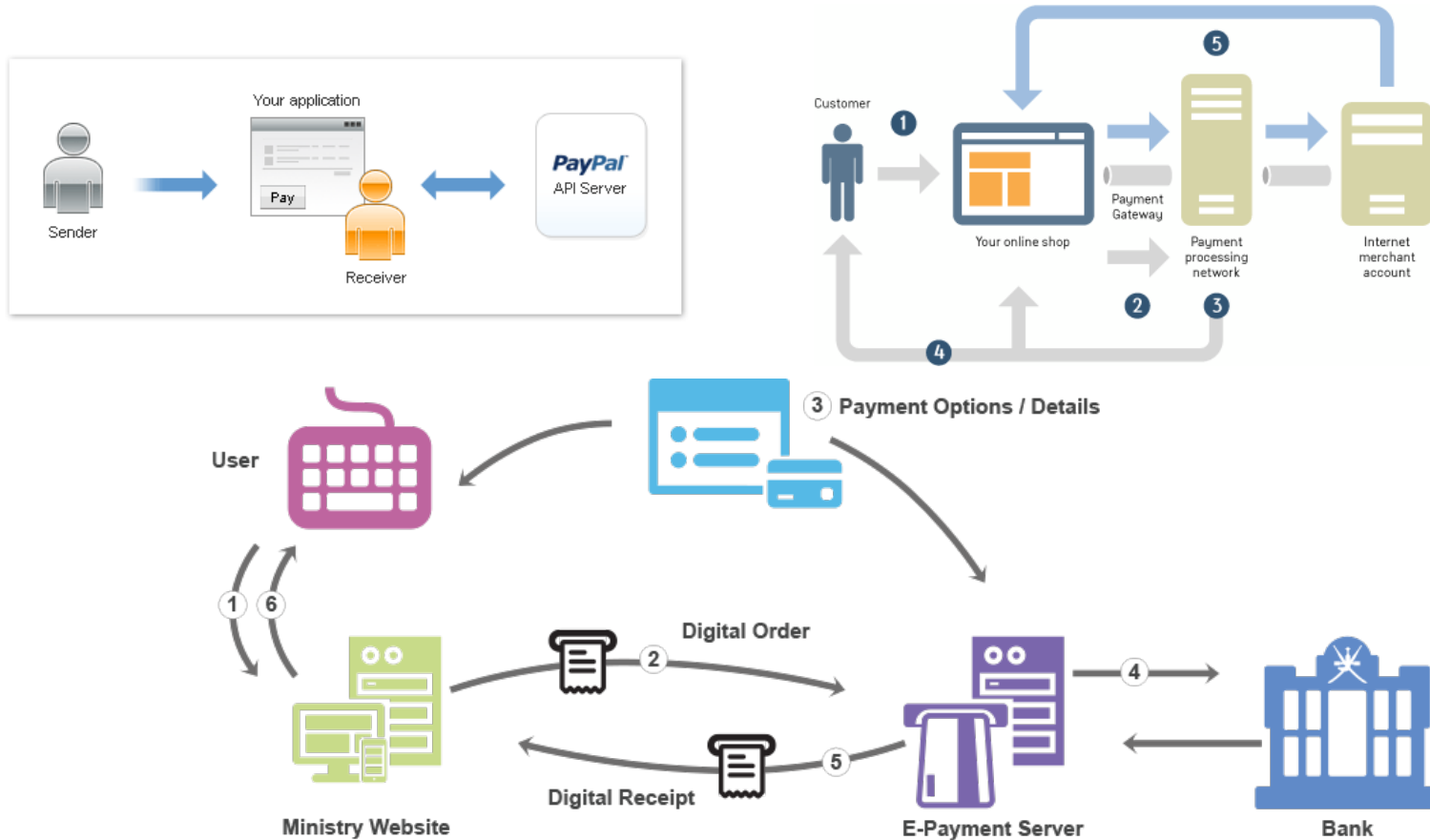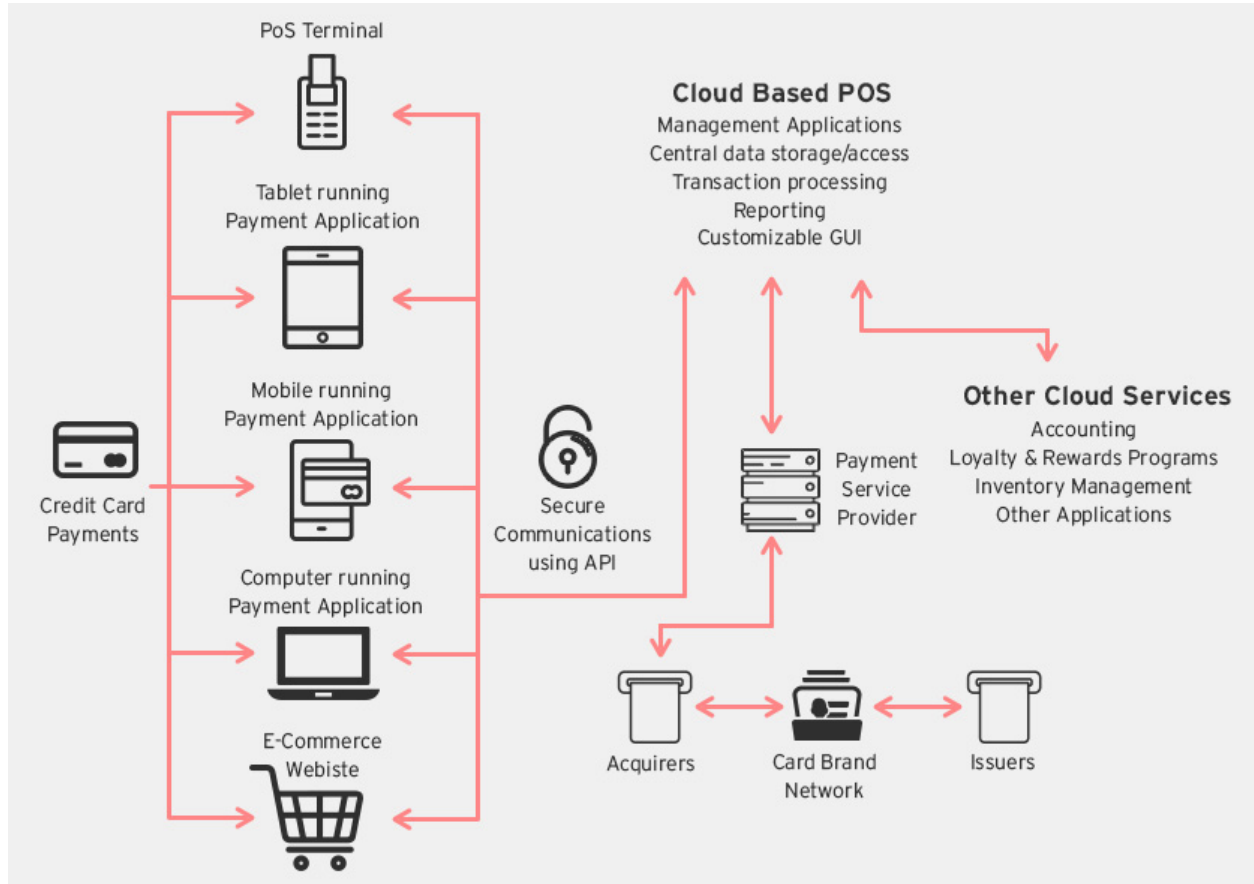- Cross data center (XDR) support



SOURCE DEVICE/USER

Real-Time Auth. QoS Billing

DESTINATION

Request

Execute Request

Config Module App

Real-Time Checks

Update Device User Setting

Hot-Standby

XDR

# Systems of Engagement

Transactional Analytics →

Actionable Insight

Algorithmic Decisioning Engine

**AEROSPIKE**

In-Memory/Hybrid-Memory Database

BATCH BI

hadoop

APACHE HBASE

APACHE Spark™

DATA WAREHOUSE/ DATA LAKE

Fraud Prevention
Call Detail Record
Gaming
Banking
Payment Gateway
Real-time Bidding
Packet Routing
Msg / Chat
Look & Book
eCommerce
Session Management
IoT

MACHINE LEARNING

REAL-TIME BI

Low Read Latency
Heavy Write Load
100% Up Time

AEROSPIKE

# Example

# The Scale Problem in Payment Fraud Detection

# Payment systems are evolving fast

# Operational Scale Explosion

**BUSINESS TRANSACTIONS**

Web views

( Payments )
( Mobile Queries )
( Recommendation )
( And More )

Decisioning Engine

High Performance NoSQL

**LEGACY RDBMS HDFS BASED**

XDR

"REAL-TIME BIG DATA"
"DECISIONING"

LEGACY DATABASE
(Mainframe)

DATA WAREHOUSE/
DATA LAKE

| 500 | X | 5000 | = | 2.5 M |
|---|---|---|---|---|
| Business Trans per sec | | Calculations per sec | | Database Transactions per sec |

# The billions of objects problem – Streaming vs DBMS

Streaming system can only store a limited number of objects in memory

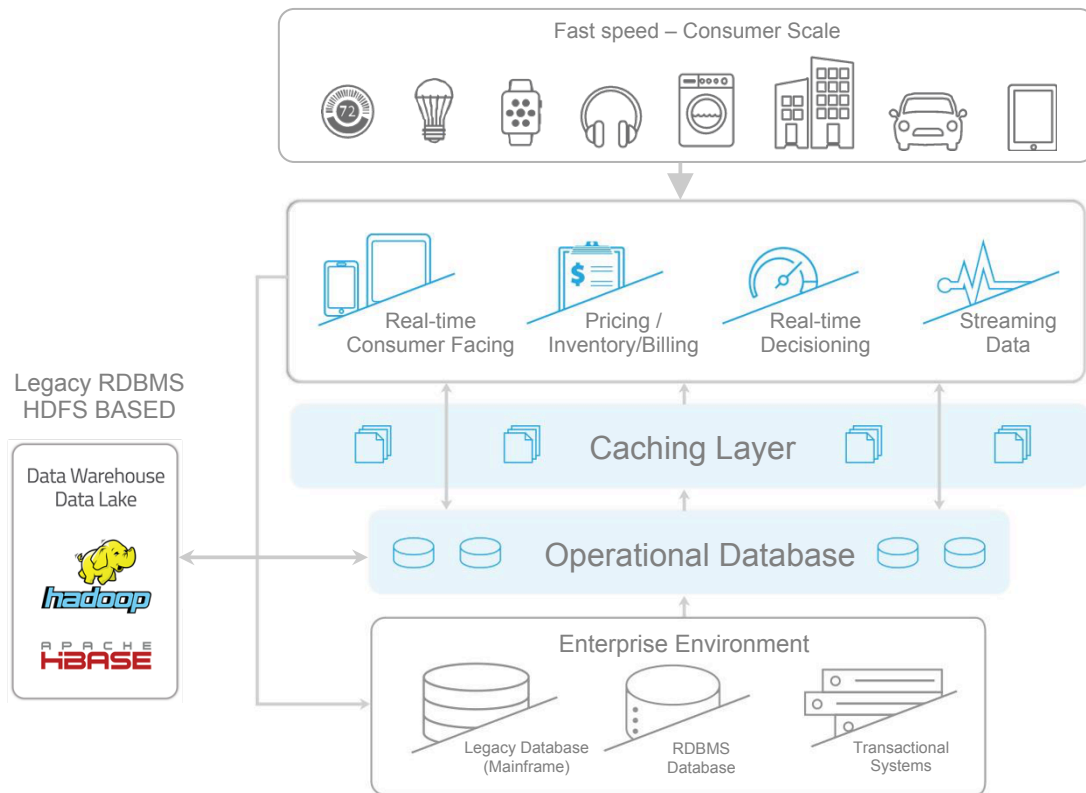Joining the active objects (millions) to database objects (billions) is best done using a distributed KVS

Active Data for Real-time Fraud

Real-Time

Read/Update

Database with historical data

- Millions of objects in DRAM
- 1-100 GB
- Fraud and Risk processing
  - bespoke/Spark/TickDB etc.

- Billions of objects on Storage
- 10-100TB
- Cassandra, Aerospike, etc.

# Hybrid Memory Database

# Systems of Engagement … what is required ?

- **>> TPS (speed) … greater than 1 million tps**

- **>> Scale … greater than 5 to 10 TB**

- **<< Low Latency …  ~ 1 msec per transaction**

- **Reliability … ~ five 9s**

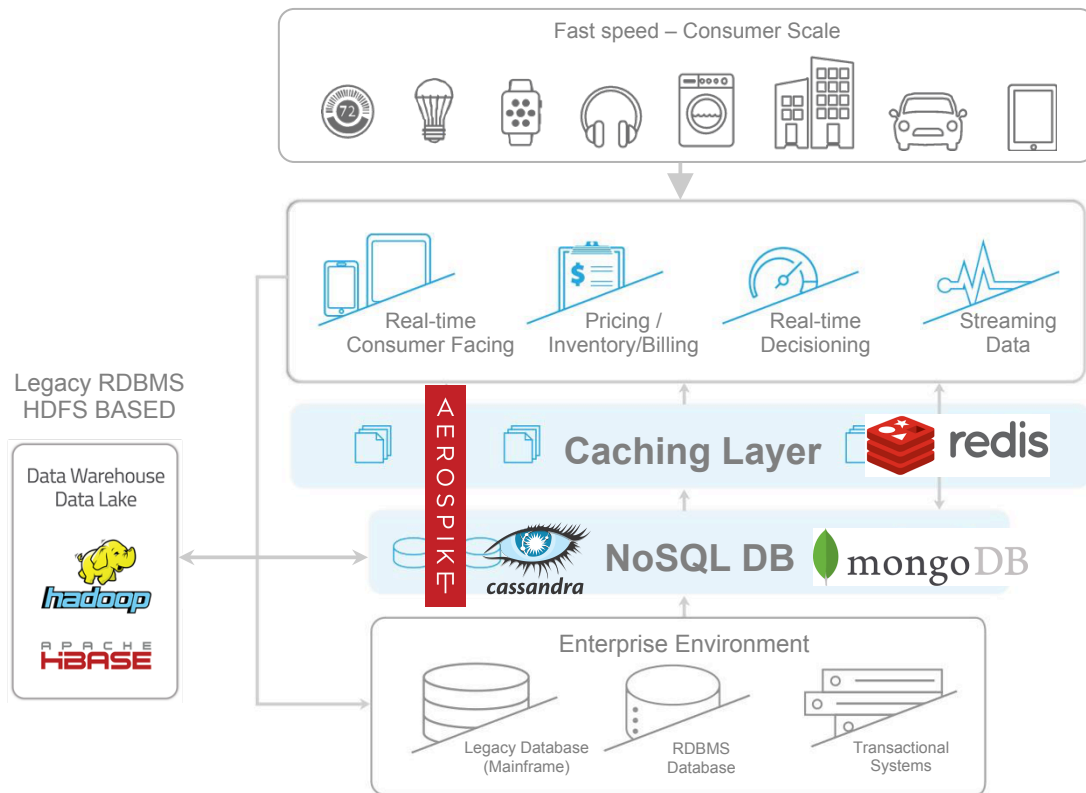- **TCO … the lowest reasonable cost**

# Traditional architecture has significant limitations

Fast speed – Consumer Scale

Real-time Consumer Facing

Pricing / Inventory/Billing

Real-time Decisioning

Streaming Data

Legacy RDBMS HDFS BASED

Data Warehouse Data Lake

Caching Layer

Operational Database

Enterprise Environment

Legacy Database (Mainframe)

RDBMS Database

Transactional Systems

## Challenges

- Complex
- Maintainability
- Durability
- Consistency
- Scalability
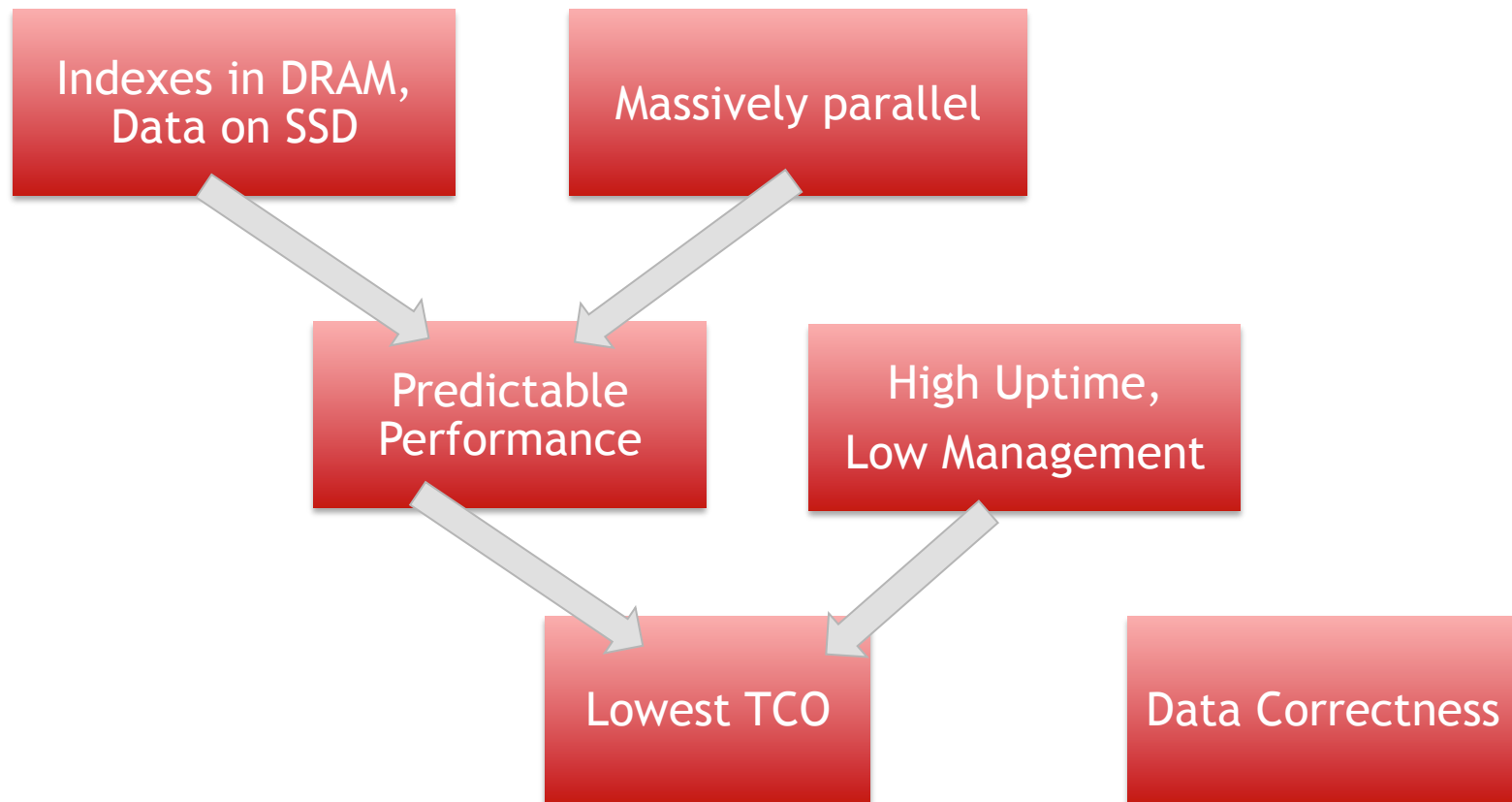- Cost ($)
- Data Lag

# Hybrid Memory Architecture



Fast speed – Consumer Scale

Real-time Consumer Facing

Pricing / Inventory/Billing

Real-time Decisioning

Streaming Data

Legacy RDBMS HDFS BASED

Data Warehouse Data Lake

hadoop

APACHE HBASE

AEROSPIKE

Caching Layer

redis

NoSQL DB

cassandra

mongoDB

Enterprise Environment

Legacy Database (Mainframe)

RDBMS Database

Transactional Systems

## Benefits

- Fast App Development

- Richer data schema

- No need for SQL

- In-memory performance

- High scale

- Lower latency

- Distributed

- Tradeoff: Consistency versus Availability

AEROSPIKE

# Aerospike System Overview
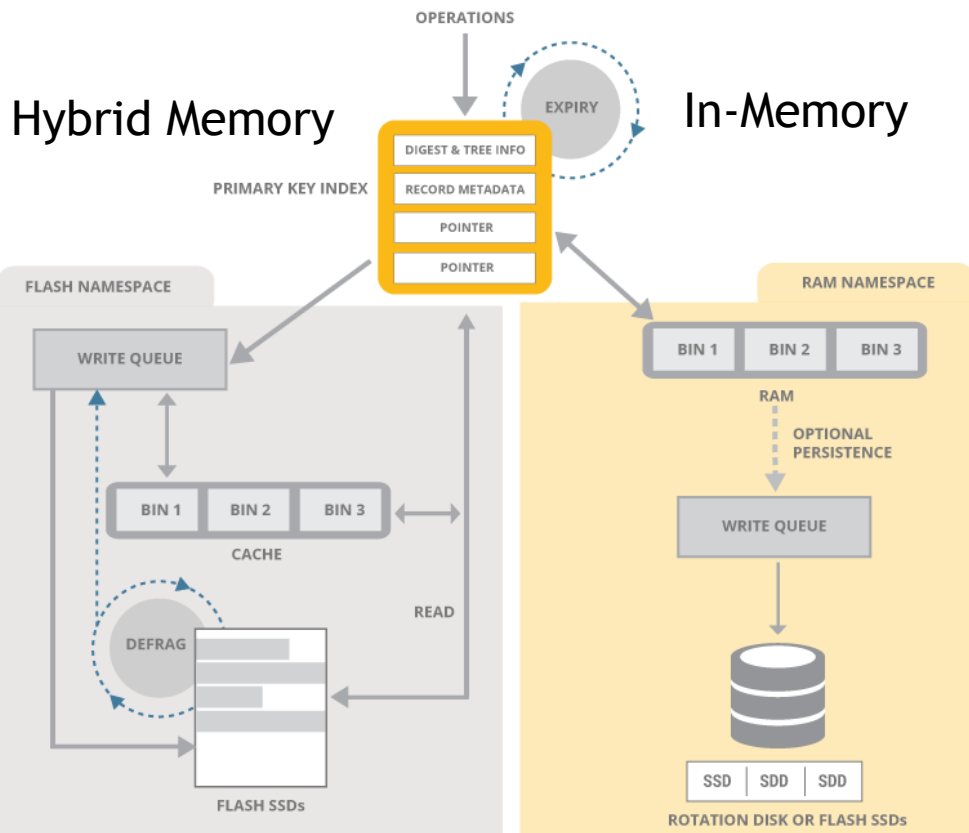


1) **No Hotspots** – Distributed Hash Table simplifies data partitioning

2) **Smart Client** – **1 hop** to data, no load balancers

3) **Shared Nothing Architecture**, every node is identical

4) **Smart Cluster, Zero Touch** – auto-failover, rebalancing, rack aware, rolling upgrades

5) **Transactions and long-running tasks prioritized in real-time**

6) **XDR** – async replication across data centers ensures **Zero Downtime**

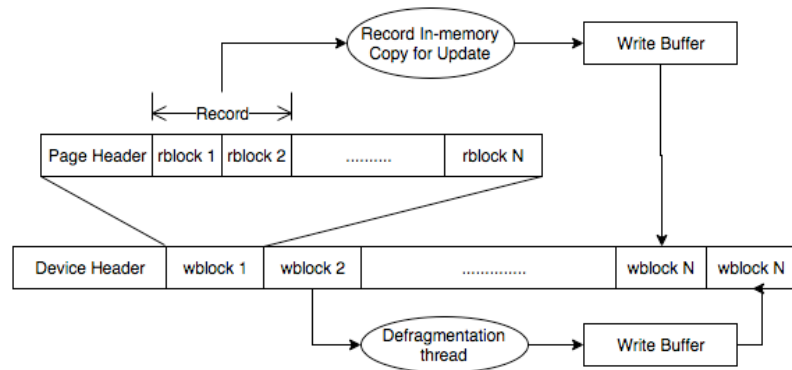# Attributes of a Hybrid Memory Architecture

Indexes in DRAM, Data on SSD

Massively parallel

Predictable Performance

High Uptime, Low Management

Lowest TCO

Data Correctness

# Aerospike Storage Architecture (HMA+)



## Highlights

1. Direct device access
2. Large Block Writes
3. Indexes in DRAM
4. Highly Parallelized
5. Log-structured FS "copy-on-write"
6. Fast restart with shared memory
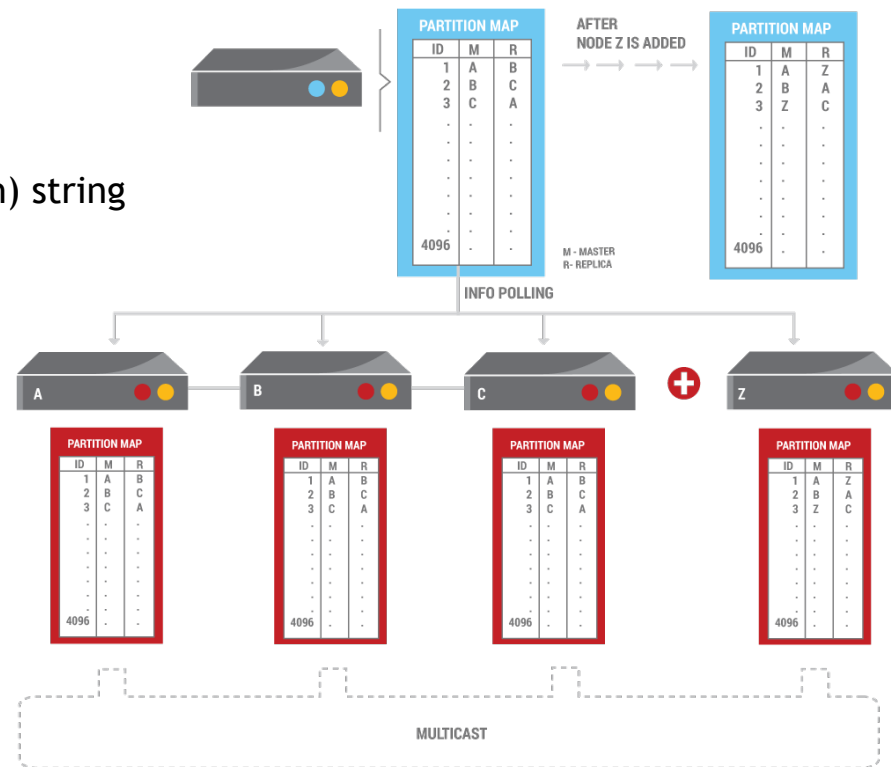
## Storage Layout

# Hybrid Memory Characteristics

**Indexes in DRAM, Data on SSD**

- **Small amount of DRAM**
  - Avoid cost and server sprawl

- **No cache, so no cache misses**
  - Predictable, low latency performance on NVMe/SSD

- **Optimized for SSDs**
  - Reads done in parallel
  - Writes done optimally for SSD to reduce wear-and-tear

# Distributed Hash Based Partitioning

- **Distributed Hashing with No Hotspots**
  - Every key **hashed** with **RIPEMD160** into an ultra efficient 20 byte (fixed length) string
  - Hash + additional (fixed 64 bytes) data forms **index entry** in RAM
  - **Some bits** from hash value are used to calculate the **Partition ID** (4096 partitions)
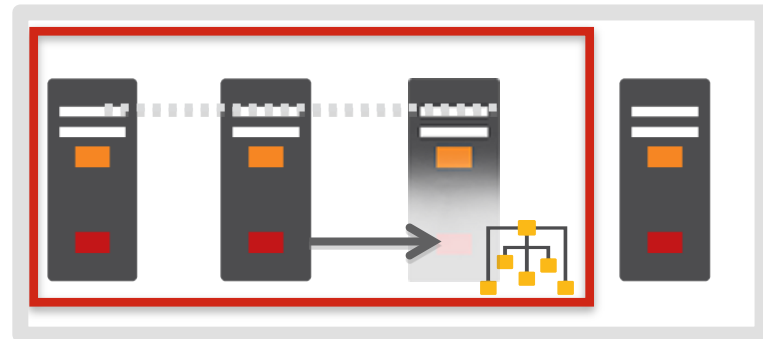  - Partition ID maps to Node ID in the cluster

# Automatic rebalancing

**Adding, or removing a node, the cluster automatically rebalances**

1. **Cluster discovers new node via gossip protocol**
2. **Paxos vote determines new data organization**
3. **Partition migrations occur**

**After migration is complete, the cluster is evenly balanced.**
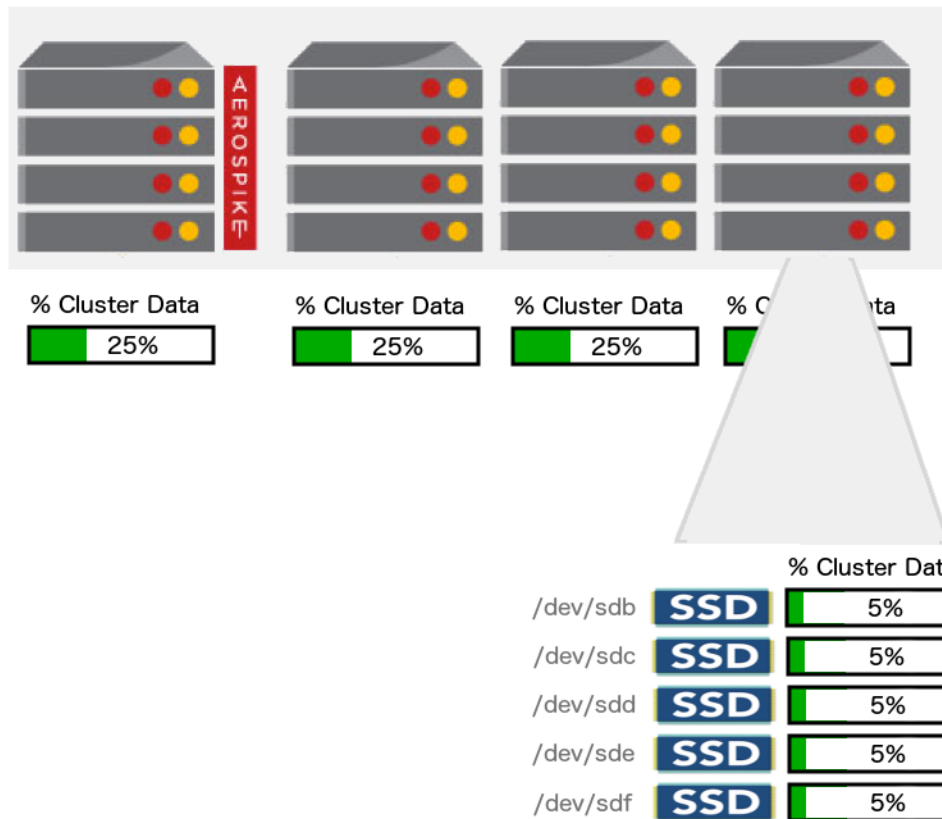
**Clients keep working during rebalancing.**

Massively parallel

- **Take full advantage of all the hardware**
  - Scaling up

- **Scale linearly with number of nodes**
  - Scaling out

## Automatic Distribution of Data using Smart Partitions™ algorithm

- Even amount data on every node and on every flash device
- All hardware used equally
- Load on all servers is balanced
- No "hot spots"
- No configuration changes as workload or use case changes

## Smart Clients

- Single "hop" from client to server
- Cluster-spanning operations (scan, query, batch) sent to all processing nodes for parallel processing.



| % Cluster Data | % Cluster Data | % Cluster Data | % Cluster Data |
|---|---|---|---|
| 25% | 25% | 25% | |

| | | % Cluster Data |
|---|---|---|
| /dev/sdb | SSD | 5% |
| /dev/sdc | SSD | 5% |
| /dev/sdd | SSD | 5% |
| /dev/sde | SSD | 5% |
| /dev/sdf | SSD | 5% |

# Aerospike's Predictable Performance

## Performance Built In

- Written in C with memory-optimized libraries => No garbage collection
- Continual defragmentation of storage => No compactions
- Known master for any piece of data => No quorum reads
- Designed as a distributed database => Networking primary consideration

## Storage Optimizations

- Writes done to memory buffer => Avoid storage slowdown
- Storage used in "block" mode => No file system overhead
- Reads and writes striped across devices => Concurrent use of hardware
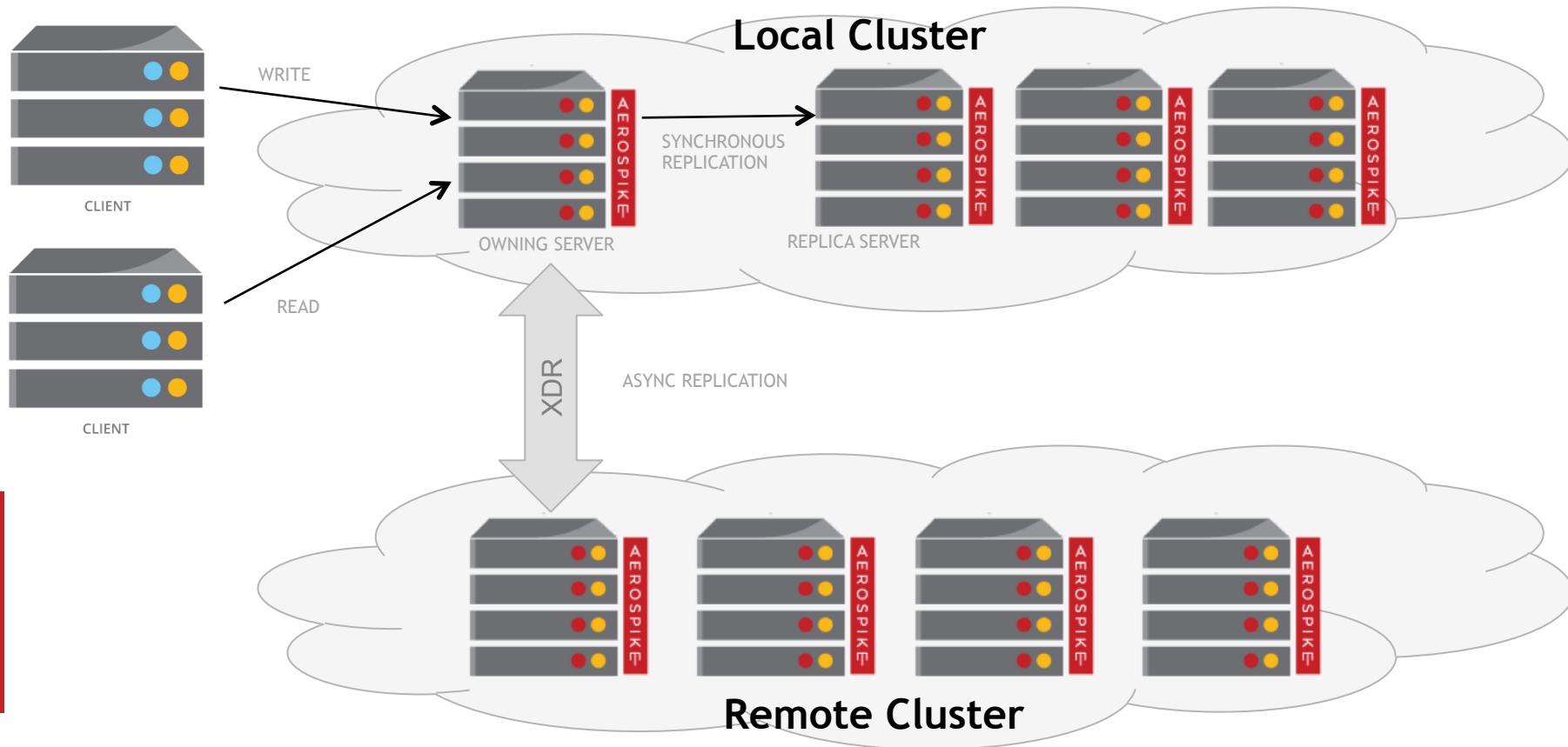
## Smart Clients

- Single "hop" from client to server
- Partition map stored on client
- Automatic load balancing – no external load balancers!

AEROSPIKE

**Data Correctness**

- **Reads should return the latest copy of the data**
  - With no latency penalty

- **Caches should not be necessary**
  - Eliminates stale data reads

- **Mixed workloads should not cause issues**
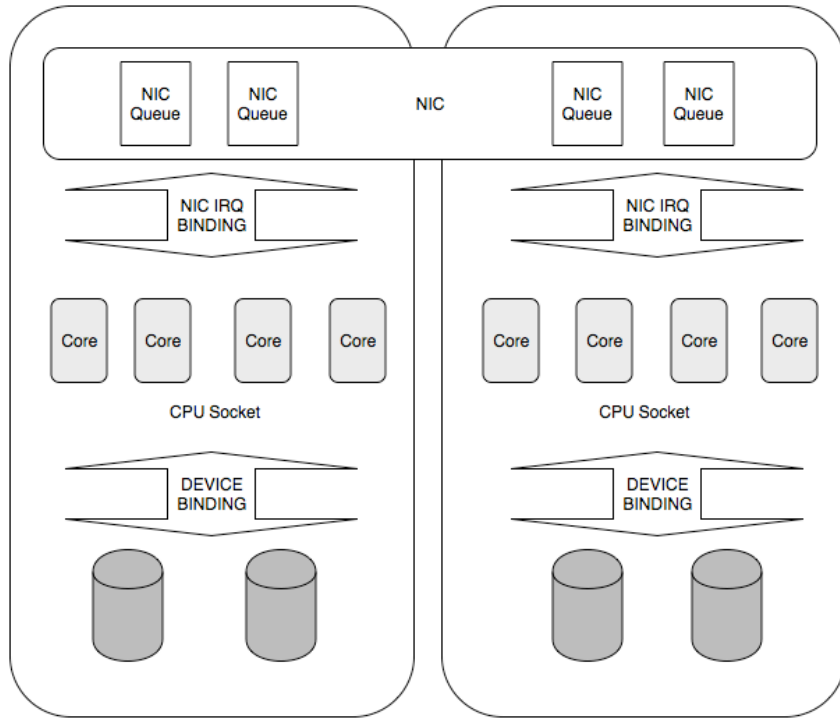  - True concurrent reads/writes

AEROSPIKE

**Local Cluster**

CLIENT

WRITE

READ

CLIENT

OWNING SERVER

SYNCHRONOUS REPLICATION

REPLICA SERVER

XDR

ASYNC REPLICATION

**Remote Cluster**

# High Uptime, Low Management

**High Uptime,
Low Management**

- **High Uptime**
  - "Shared Nothing" Architecture
  - No single points of failure
  - No cascading failures
  - Seamless loss of nodes with self-heal capability

- **Low Management**
  - Automatic sharding of data
  - No re-tuning of cluster for use-case changes
  - No requirement for caches
  - Smaller number of nodes for easier management
  - "Set and forget" DevOps management
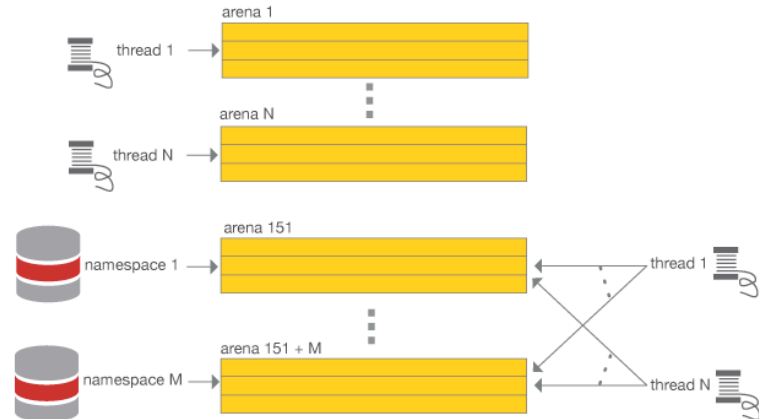
# Designed for Wire-Line Speed

## Multi-core architecture



## Optimized C based DB kernel

1. Multi-threaded data structures
2. Nested locking model for synchronization
3. Lockless data structures
4. Partitioned single threaded data structures
5. Index entries are aligned to cache line (64 bytes)
6. Custom memory management (arenas)
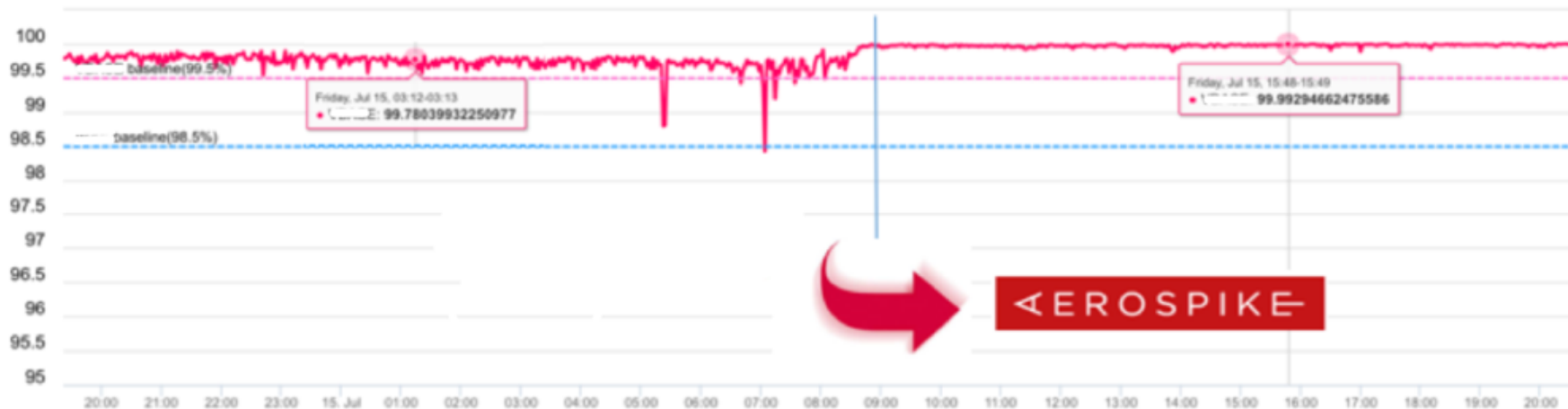
## Memory Arena Assignment

# Hybrid Memory Benefits

## 30X Improvement in SLA

In-Memory SLA 98.5%     Hybrid Memory SLA 99.95%



## Missed SLA is lost Revenue!!!

# Lowest TCO

Indexes in DRAM, Data on SSD

Massively parallel

Predictable Performance

High Uptime, Low Management

## Lowest TCO

- **Hybrid Memory Architectures offer**
  - Cacheless, consistent performance using NVMe/ Optane.
  - Server count reduced (3x or more)
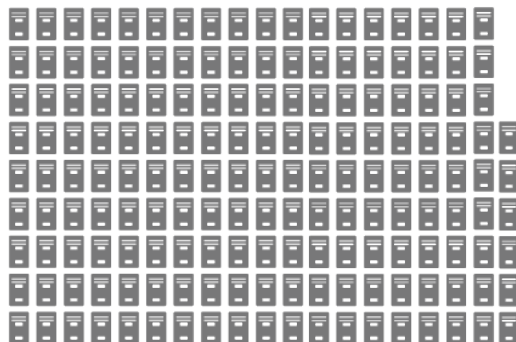  - Significant reduction in TCO (10x documented)

**⊿EROSPIKE**

UP TO

## 10x FASTER
## 10x FEWER

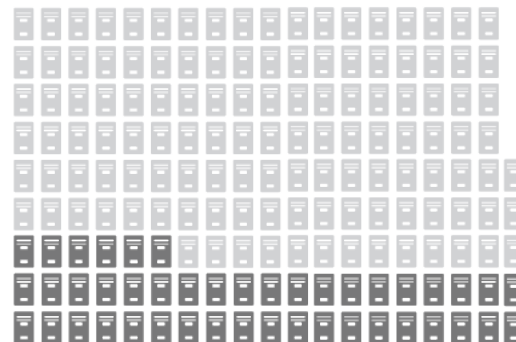ACTUAL CUSTOMER ANALYSIS

**168 SERVERS**

**OTHER DATABASE**
DRAM & HDD

**ONLY 44 SERVERS**

**⊿EROSPIKE**

SSD & DRAM

| | Year 1 (Millions USD) | Year 2 (Millions USD) | Year 3 (Millions USD) | Total (Millions USD) |
|---|---|---|---|---|
| Cassandra | $1.82 | $2.53 | $3.94 | $8.28 |
| Aerospike | $1.88 | $1.24 | $1.85 | $4.97 |
| **Total OpEx Savings** | **-$0.06** | **$1.28** | **$2.09** | **$3.32** |

**YoY SPEND ON OPERATIONS**  ● Cassandra  ● AEROSPIKE

$4M

$3M

$2M

$1M

$0M   Year 1          Year 2          Year 3

AEROSPIKE

# Case Studies: HMA - Lower TCO & better SLA

| Customer | Situation | Problem | Hybrid Memory System |
|---|---|---|---|
| **Trading Account Risk Management** | DB2+Gemfire cache | 150 Servers growing to 1000 | Single cluster – 12 servers |
| **Payments Fraud Detection** | 2 ORCL RAC clusters + Terracotta cache | System Stability & missing SLA's | 3 Clusters – 20 Servers each |
| **User Integrity Checking for Internet Transactions** | DataStax/Cassandra | 168 DataStax Servers growing to 450+ | 30 Servers – 2 clusters |
| **Telco Device and User Access** | ORCL Coherence / DataStax Cassandra | Existing SOE solutions unstable & Costly | 5 successful POC's |
| **Telco Revenue Assurance** | DataStax/Cassandra PostgreSQL + cache | Hundreds of cache & Cassandra Servers Scalability challenges | Significant reduction of server footprint – global deployment |

AEROSPIKE

# Next Generation Systems of Engagement –
## An Emerging Market with Multiple Technologies

**Systems of Engagement – Many Choices**

Unique Functional Capabilities and Hybrid Memory Solution

Significant functional overlap - Commodity DB problem set

Speed TPS

Scale TB

**Systems of Engagement - TCO**

Alternative TCO

Hybrid Memory TCO

TCO ($)

Scale TB

*Hybrid Memory Architecture Delivers Predictable Performance, Highest Availability, and Lowest TCO*

# Thank You
*Questions?*