Aerospike @ PayPal

Athreya Gopalakrishna NoSQL Engineering Lead

Aerospike NoSQL DB

Key/Values

Hardware Evolution

NoSQL Database Evolution

Case Study – Fraud Detection System

NoSQL DB Storage Architectures

Designing a 50TB NoSQL DB

System Efficiency, Performance etc.

Q & A

Key Values

Key Value, Computer Languages

KV is a storage paradigm to store and retrieve data

key	value
k1	1
k2	Aa, b1, 1c
k3	1,a,5/17/2018

KV Store in primitive forms

```
Map Container in C++ STL[map <int, int> kvstore;]
HashMap in Java [HashMap<Integer, Integer> kvstore = new HashMap(..)]
Dictionary in Python [kvstore = {}]
```

Simple interface

```
put(key, value), put(key, value, TTL), get(key), delete(key)
```

Hardware Evolution





NoSQL Database Evolution





Case Study

Fraud Detection System

- Analytical system
- Built on Relational and KV databases
- Requires Low Latency and High Throughput
- 1-200KB avg. object size
- 1ms@99 Percentile
- 4-5 Millions of transactions/second
- Trillions of keys
- 100+ of Terabytes of Storage

KV STORAGE GROWTH TREND

-Data Growth (TB)



KEY SPACE GROWTH TREND



Storage Architecture

High Performance NoSQL

Cache	In-Memory	Memory-First	Hybrid-Memory
 Index, Data Stored and Served from Memory No Data on Disk 	 Index, Data Stored and Served from Memory Data is persisted on Disk 	 Index, Data Stored in Memory and Disk OR Index in Memory and Data on Disk Data is persisted on Disk 	 Index in Memory and Data on Disk Data is persisted on Disk



Write Path – Latency and Consistency



Why Aerospike ?

Aerospike Architecture

Key Differentiators in NoSQL space

Aerospike



Samsuns

Ground up, Designed for SSDs. (Achieves – Even wear and tear on Device)

Proprietary file system (Achieves – Consistent Device Latency, Follows Device throughput)



Hybrid Storage – Predictable capacity. (Achieves – Enables huge storage on SSD)

Aerospike as NoSQL Database



- Written in C
- Simple KV database
- Distributed shared nothing architecture
- AP and CP (Strong consistency) modes
- Operates In-Memory or Hybrid-Memory Modes
- Low write amplification
- SSD optimized for consistent performance
- High storage density
- Low CPU utilization
- **UDF** for server side computations

Designing a 50TB A/A Database 50T 50T В В C2 C1 DC1, RF=2 X-DC replication X-DC replication A/A 50T 50T 50T 50T В В В В C3 C4 C6 C5 **X-DC** replication DC3, RF=2 DC2, RF=2

In-Memory Database for 50TB

(Predictable performance)



Servers ~ 1024

Price ~\$12M



Memory-First Database for 50TB

(Predictable performance)



Servers ~ 1024

Price ~\$15M



Memory-First Database for 50TB

(Unpredictable performance)



Price ~\$1.8M



Hybrid Memory Database for 50TB

(Predictable performance)



System configuration





Load 250M 1KB Value Size Raw Data = 250GB

Max Write = 100K TPS CPU = 5% Mem = 15GB Disk = 268GB

Aerospike



Write Throughput (in-mem vs hma)



HW Profiles

Package	M(Samsung)		XL(Samsung)		L(Intel)		XL(HP-OEM)		
Server	· # # # # #.								
	Dell Power Ha	Edge R730XD	Dell Power	Edge R730XD	Dell Power Ha	Edge R730XD	HP DL3	60P Gen 9	
Cost/Lipit			Haswell						
Cost/GR	528,000		\$45,382		35,000		\$23,518		
	256GB		512GB		512GB		384GB		
	3.4TB		12.8TB		8TB		12TB		
SSD Model	Samsung SM1715		Samsung SM1715		Intel P3700		Intel \$3700		
SSD Form Factor	PCle	e NVMe	PCIe NVMe		PCIe NVMe		SAS		
SSD Storage Density	3.2TB		3.2TB		2ТВ		1.92TB		
SSD #	2x		4x		4x		6x		
Random Read IOPS (4KB)	750,000		750,000		450,000		75,000		
Random Write IOPS (4KB)	130,000		130,000		175,000		36,000		
Sequential Read (MB/s)	3	8000	3000		2800		500		
Sequential Write (MB/s)	2200		2200		2000		460		
Aerospike Max Read throughput (1KB)	38	0,000	760,000		780,000		200,000		
Aerospike Max Write throughput (1KB)	42	0,000	840,000		800,000		542,000		
Aerospike Read Latencies(Min/Avg/95th/99th/Max)	Percentile	Latency(µs)	Percentile	Latency(µs)	Percentile	Latency(µs)	Percentile	Latency(µs)	
80K-R/200K-W @1KB	Min	113	Min	113	Min	112	Min	271	
	Avg	115	Avg	115	Avg	113	Avg	402	
	95th	108	95th	108	95th	114	95th	407	
	99th	108	99th	108	99th	114	99th	645	
	Max	108	Max	108	Max	114	Max	856	
Aerospike Write Latencies(Min/Avg/95th/99th/Max) 80K-R/200K-W @1KB) 215/264/288/288/288(μs)		215/264/288/288/288(µs)		233/236/240/240(µs)		189/196/202/221/277(μs)		
Cold restart (Minutes)	23m		1hr 45m 52s		2hr 10m 16s		1h 11m 53s		
Warm restart(Minutes)	1r	n 57s	3m 14s		7m 5s		3r	3m 10s	



Aerospike Object size vs. Latency

2x = 80K read + 200K writes



Client response time





propriotory

Thank you

