# GENESIS

Michal Szymaniak
mips@infradepot.com

# Genesis

- Zero-dependency storage platform

    · Replication
        · Goal: reliability and availability
        · How: Raft over RocksDBs

    · Sharding
        · Goal: capacity and throughput
        · How: BigTable-style range splitting

    · Change notification
        · Goal: consistency and integration
        · How: Zookeeper-style watches

# Genesis

- Google-inspired storage layer
  - Critical problems solved once
    - Less code – Fewer bugs – Higher reliability
    - Re-usable service – Easier maintenance

  - Higher layers simpler and focused
    - Less complexity – Faster development
    - Less specialized expertise – Easier staffing

  - Lower layers abstract and disaggregated
    - Same API for physical / virtual / cloud storage – Natural data mobility
    - Decoupled resource lifetimes – Smooth HW / DC / Cloud operations

# Genesis Use Cases

- Infrastructure metadata store
  - Distributed filesystem, package management, ...

- Massively shardable NoSQL
  - User metadata, event processing, ...

- Alternative implementation of popular APIs
  - DynamoDB, BigTable, ...

- Geo-distributed data storage
  - Multi-AZ, multi-region, multi-cloud, …

- Vehicle for physical data migrations
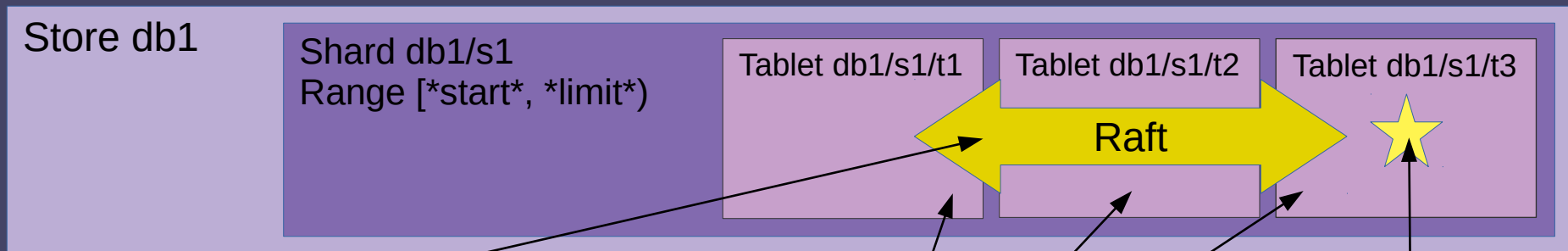  - Zero-downtime moves between clouds, regions, datacenters, ...

# Key-Value Store API

```
01| // Keys: N-tuples of binary strings, with key[0] determining the shard
02|
03| // Read: point reads or range scans, optionally with key filtering
04| ReadResult = { Key, Value, Stat }
05| Status Read(ReadOptions, Store, Key, ReadResult)
06| Iterator NewIterator(ReadOptions, Store)
07|
08| // Commit (aka MultiOp): atomic multi-key batch of operations
09| Op = CheckExists | CheckNotFound | CheckValue | CheckVersion |
10|     SetCounter | IncCounter | DecCounter |
11|     Write | WriteWithCounter | Delete | DeleteRange
12| Mutation = [ Op1, Op2, Op3, .. ]
13| Status Commit(CommitOptions, Store, Mutation)
14|
15| // Watch: individual keys or key "patterns" across entire store
16| Status WatchKey(WatchOptions, Store, Key, WatchCallback)
17| Status WatchStore(WatchOptions, Store, KeyFilter, WatchCallback)
```

# Hello World

```cpp
01| // Connect to Genesis
02| GrpcNetworkEnv env;
03| Client client(&env, FLAGS_genesis_bootstrap_servers);
04|
05| // Open the store
06| StoreHandle store;
07| CHECK_OK(client.OpenStore("db1", &store));
08|
09| // Write something
10| Mutation mutation;
11| const Key key("user:alice", "email");
12| mutation.Write(key, "alice@foo.com");
13| CHECK_OK(client.Commit(CommitOptions(), store, mutation));
14|
15| // Read it back
16| std::string email;
17| CHECK_OK(client.Read(ReadOptions(), store, key, &email));
```

# Logical View

**Store db1**

**Shard db1/s1**
Range [*start*, *limit*)

Tablet db1/s1/t1

Tablet db1/s1/t2

Tablet db1/s1/t3

Raft

**Raft**
– Consensus protocol
– Simplified Paxos
– Replicates mutation log
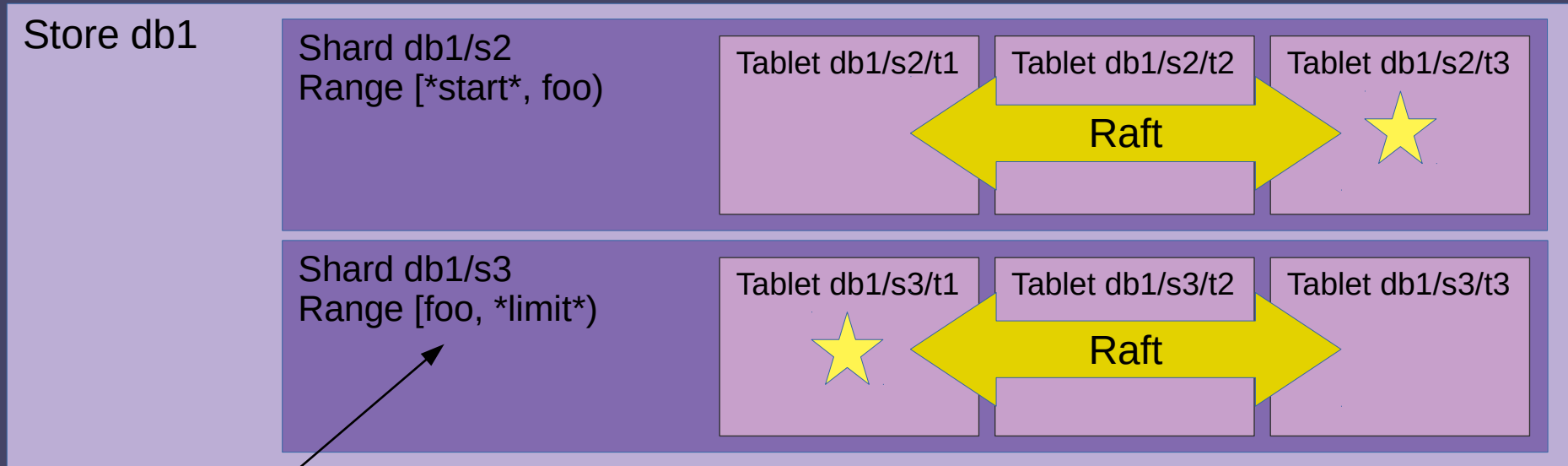– Needs only majority to work

**Tablet == Shard Replica**
– Voter or Observer
– Added / removed at will
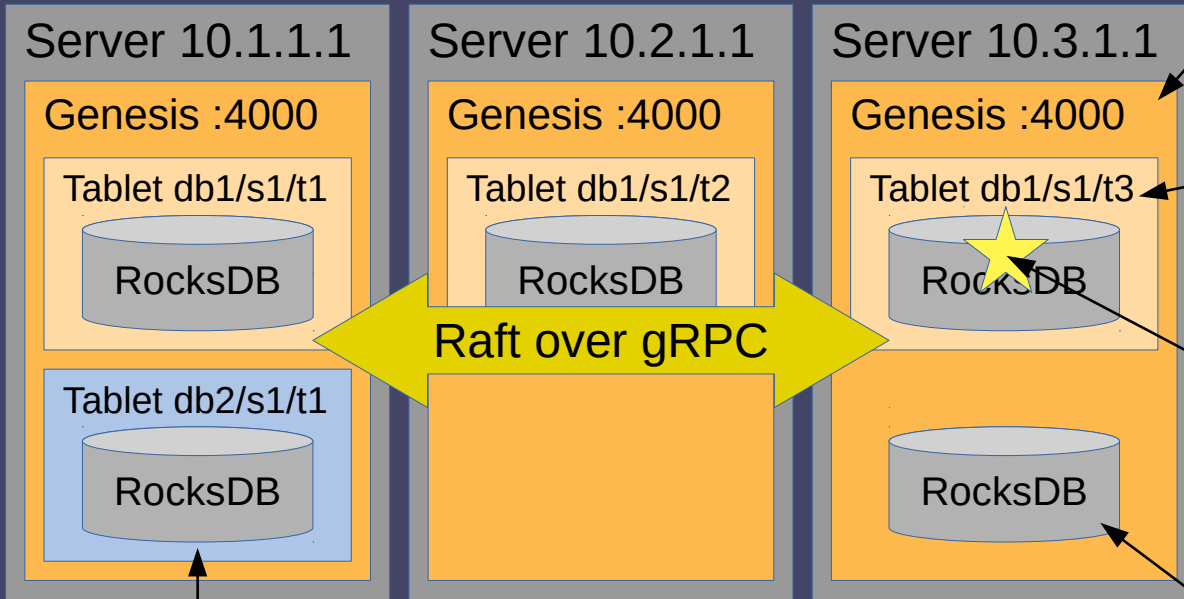  ("shard re-configuration")

**Raft Leader**
– Elected among voters
– Replication driver
– Consistent state
– Automatic failover

# Logical View: Multiple Shards

**Store db1**

**Shard db1/s2**
**Range [*start*, foo)**

| Tablet db1/s2/t1 | Tablet db1/s2/t2 | Tablet db1/s2/t3 ⭐ |
|---|---|---|

Raft

**Shard db1/s3**
**Range [foo, *limit*)**

| Tablet db1/s3/t1 ⭐ | Tablet db1/s3/t2 | Tablet db1/s3/t3 |
|---|---|---|

Raft

- Old shard s1 split at key "foo"
- New shards own adjacent key ranges
- Separate tablets, Raft state, leaders
- Sharding hidden from clients

# Physical View



Server 10.1.1.1

Genesis :4000

Tablet db1/s1/t1

RocksDB

Tablet db2/s1/t1

RocksDB

Server 10.2.1.1

Genesis :4000

Tablet db1/s1/t2

RocksDB

Server 10.3.1.1

Genesis :4000

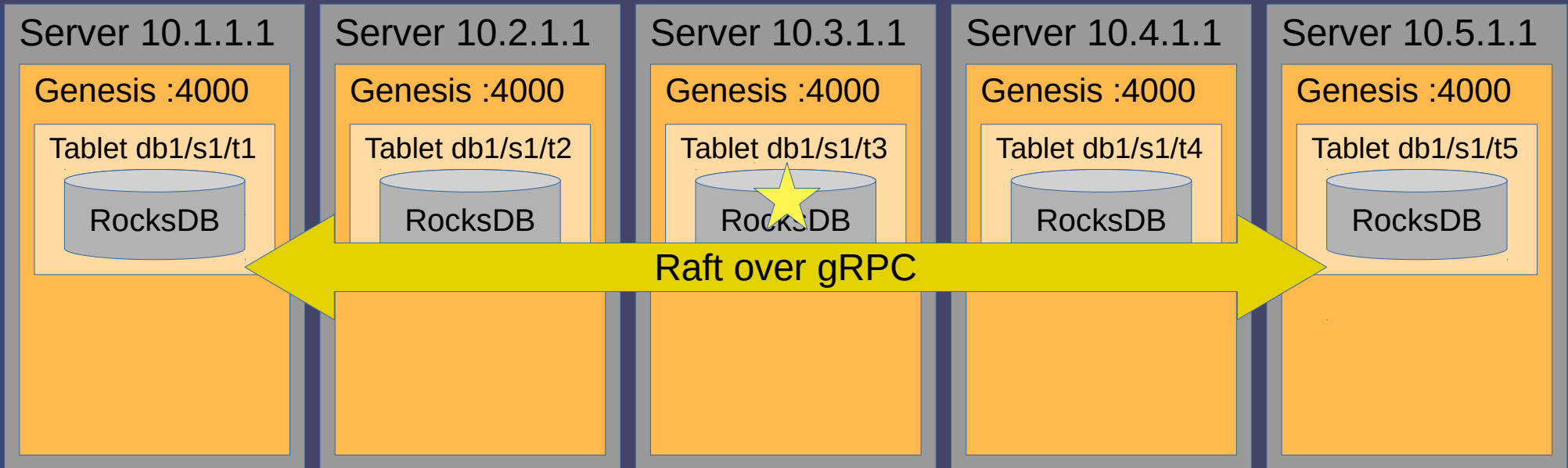Tablet db1/s1/t3

RocksDB

RocksDB

Raft over gRPC

**Tablet Server**
– Owner of disk space
– Platform for tablets
– RPC dispatcher

**Tablet**
– Owner of RocksDB
– Raft code and buffers
– RPC destination

**Raft Leader**
– Hot tablet handling all writes
– Commit calls + replication
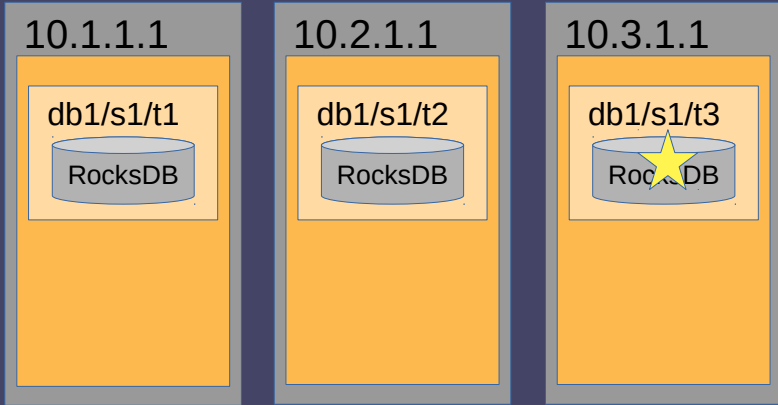– Leader flip == traffic switch

**Tablet Storage**
– "Unloaded" == inaccessible
– Tablet either being created, moving, or pending deletion

**Colocated Tablet**
– Many tablets per server
– Same or different store

# Physical View: More Tablets

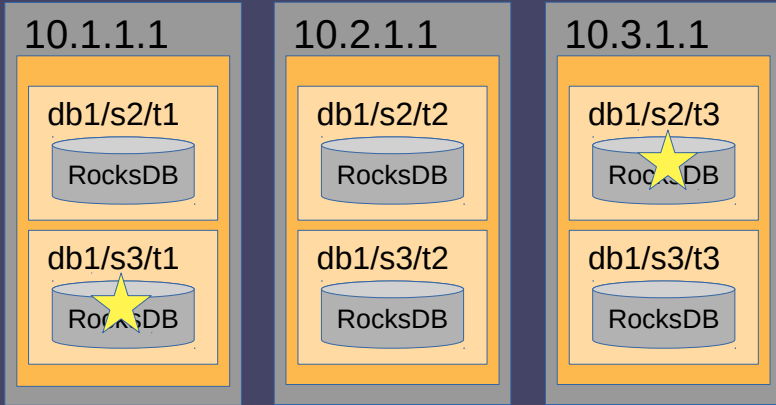| Server 10.1.1.1 | Server 10.2.1.1 | Server 10.3.1.1 | Server 10.4.1.1 | Server 10.5.1.1 |
|---|---|---|---|---|
| Genesis :4000 | Genesis :4000 | Genesis :4000 | Genesis :4000 | Genesis :4000 |
| Tablet db1/s1/t1 | Tablet db1/s1/t2 | Tablet db1/s1/t3 | Tablet db1/s1/t4 | Tablet db1/s1/t5 |
| RocksDB | RocksDB | RocksDB | RocksDB | RocksDB |

Raft over gRPC

– More tablets == higher availability + higher read throughput
– Any odd number of voters, plus any number of observers
– Configured separately per shard

# Physical View: More Shards

**10.1.1.1**

db1/s1/t1

RocksDB

**10.2.1.1**

db1/s1/t2

RocksDB

**10.3.1.1**

db1/s1/t3

RocksDB

– More shards == more capacity + more throughput

# Physical View: More Shards

| 10.1.1.1 | 10.2.1.1 | 10.3.1.1 |
|----------|----------|----------|
| db1/s2/t1 RocksDB | db1/s2/t2 RocksDB | db1/s2/t3 RocksDB ★ |
| db1/s3/t1 RocksDB ★ | db1/s3/t2 RocksDB | db1/s3/t3 RocksDB |

– More shards == more capacity + more throughput

– Shard split:

   1) split each tablet into N >= 2 new ones

# Physical View: More Shards

| 10.1.1.1 | 10.2.1.1 | 10.3.1.1 |
|---|---|---|
| db1/s2/t1 | db1/s2/t2 | db1/s2/t3 |
| RocksDB | RocksDB | RocksDB ⭐ |

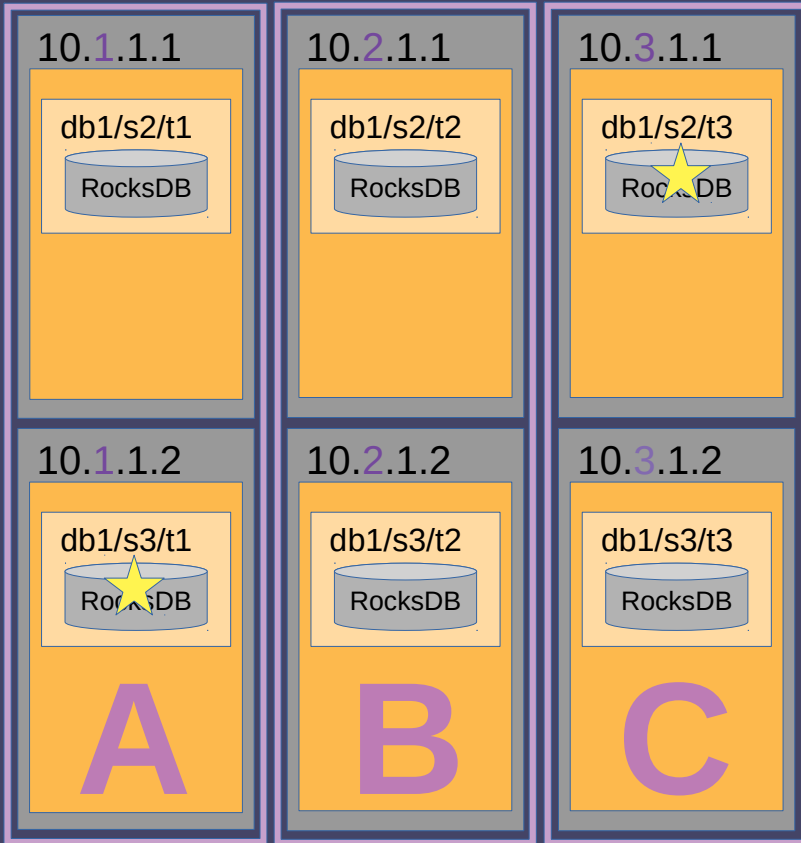| 10.1.1.2 | 10.2.1.2 | 10.3.1.2 |
|---|---|---|
| db1/s3/t1 | db1/s3/t2 | db1/s3/t3 |
| RocksDB ⭐ | RocksDB | RocksDB |

– More shards == more capacity + more throughput

– Shard split:

   1) split each tablet into N >= 2 new ones

   2) keep 1 new tablet in place, move other(s)

# Physical View: More Shards

| | | |
|---|---|---|
| 10.1.1.1 | 10.2.1.1 | 10.3.1.1 |
| db1/s2/t1 RocksDB | db1/s2/t2 RocksDB | db1/s2/t3 RocksDB ⭐ |
| 10.1.1.2 | 10.2.1.2 | 10.3.1.2 |
| db1/s3/t1 RocksDB ⭐ | db1/s3/t2 RocksDB | db1/s3/t3 RocksDB |
| **A** | **B** | **C** |

– More shards == more capacity + more throughput

– Shard split:

    1) split each tablet into N >= 2 new ones

    2) keep 1 new tablet in place, move other(s)

– New shards replicated across same server pools (typically, different availability zones)

– Tablets move within each server pool

# Physical View: More Stores



Orange store "db1":
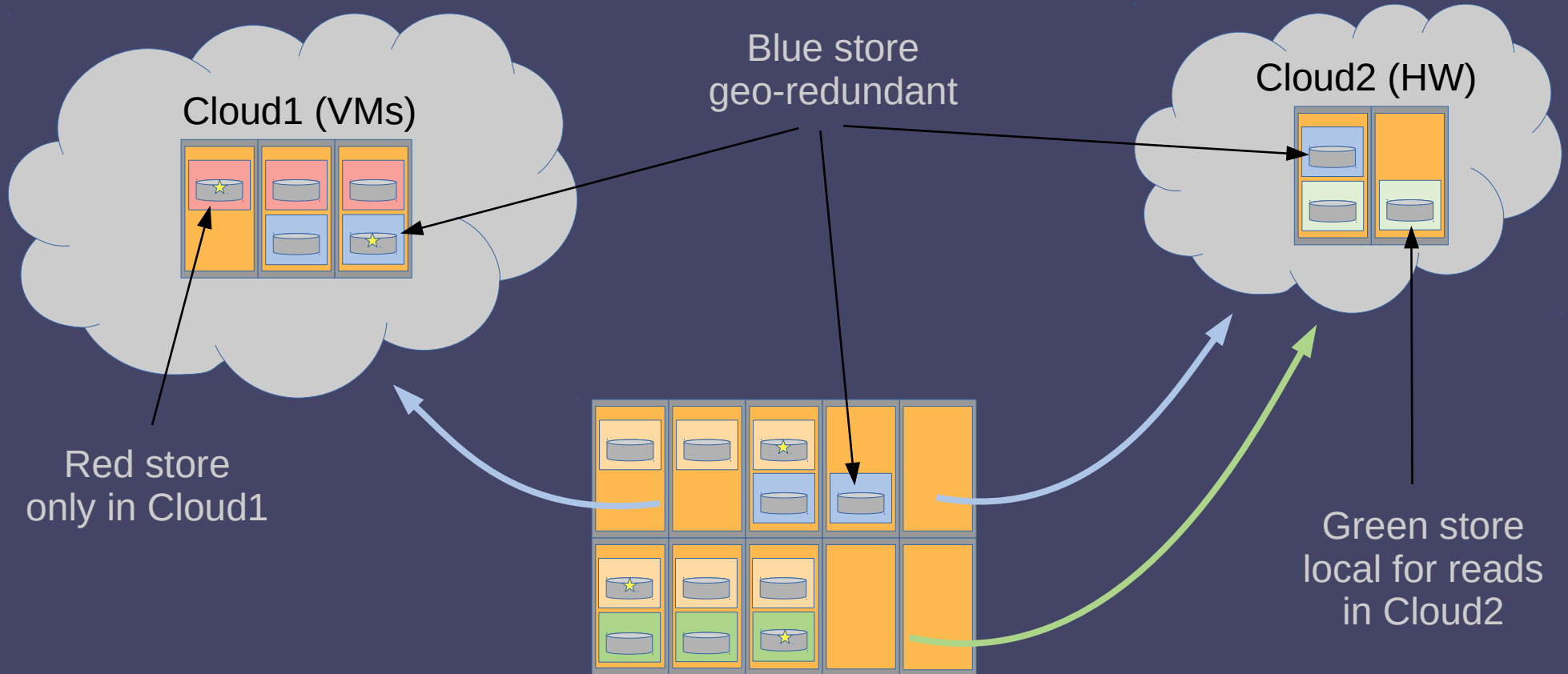– 2 shards
– sharded for capacity and throughput

Blue store "db2":
– 1 shard, 5 tablets
– higher availability than with 3 tablets

Green store "db3":
– 1 shard, 5 tablets, including 2 observers
– higher read throughput than with 3 tablets

# Physical View: More Locations



Cloud1 (VMs)

Blue store
geo-redundant

Cloud2 (HW)

Red store
only in Cloud1

Green store
local for reads
in Cloud2

# Maintenance

- Store
    - Create / Delete
    - Backup / Restore
    - Update – change store configuration (ACLs, rate limits, ...)

- Shard
    - Split / Merge / AddObserver / RemoveObserver
    - Update – promote observers to voters or vice versa
    - SetLeader – force leader election for traffic rebalancing

- Tablet
    - Move – move tablets between servers for disk space rebalancing
    - Checkpoint – clone tablet storage for backup or offline processing
    - Recover – re-create lost tablet from another, in emergency

# Performance

- Read Latency
    - 128b / 1kb / 16kb / 128kb / 1MB
    - 170 / 174 / 195 / 217 / 594 usec
    - Mostly leader RPC time

- Write Latency
    - 128b / 1kb / 16kb / 128kb / 1MB
    - 882 / 894 / 1148 / 2477 / 15429 usec
    - Mostly time of replication to quorum

- Throughput
    - Scales linearly with shard count
    - Reads per shard: up to 600,000 rps / 50 Gbps
    - Writes per shard: up to 30,000 rps / storage saturation