Distributed Key-Value Store

System Design Document

Document Version: 1.0 Classification: Internal Technical Documentation

Executive Summary

leverages BadgerDB as the storage engine and implements the Raft consensus algorithm for distributed coordination. **Key Characteristics**

• **High Performance:** Sub-millisecond latencies with throughput exceeding 100K operations/second • **Strong Consistency:** ACID transactions with linearizable consistency guarantees

• **Cloud Native:** Kubernetes-optimized with comprehensive observability

Table of Contents

1. Architecture Overview 2. Core Components

3. Data Model 4. Consensus Algorithm 5. Storage Engine 6. Network Architecture

7. API Design 8. Performance Characteristics 9. Security Model 10. Operational Concerns 11. Future Roadmap 12. Conclusion 1. Architecture Overview **System Architecture Diagram** Client Applications

Raft Raft Raft Engine Engine Engine

Key Design Principles 1. **Separation of Concerns:** Clear separation between consensus, storage, and API layers 2. **Fault Isolation:** Component failures don't cascade across system boundaries 3. **Horizontal Scalability:** Linear performance scaling with cluster size 4. **Operational Simplicity:** Self-healing with minimal manual intervention 5. **Performance First:** Optimized for high-throughput, low-latency workloads 2. Core Components 1. API Gateway Layer **Responsibilities:** Client request routing and load balancing • Protocol translation (REST ↔ gRPC) • Authentication and authorization • Rate limiting and circuit breaking • Request/response transformation

• Intelligent leader routing for write operations Read load distribution across healthy nodes • Automatic retry with exponential backoff • Comprehensive request/response logging 2. Consensus Engine (Raft) **Responsibilities:** • Distributed consensus and leader election • Log replication across cluster members • Membership changes and configuration updates • Split-brain prevention and network partition handling **Implementation Details:**

Key Features:

• Multi-protocol support (gRPC, HTTP/REST)

3. Data Model **Key Structure** type Key struct { Value string // UTF-8 string, max 1KB Namespace string // Optional namespace for isolation Metadata Metadata // System and user metadata } type Metadata struct { CreatedAt time.Time UpdatedAt time.Time ExpiresAt *time.Time // Optional TTL Version uint64 // Optimistic concurrency control Tags map[string]string }

Compression encoding Checksum string // Data integrity verification }

The system supports ACID transactions with the following guarantees:

Atomicity: All operations in a transaction succeed or fail together

Isolation: Concurrent transactions don't interfere with each other

"user:123:last_login" value: "2024-01-15T10:30:00Z" ttl_seconds: 86400

• Write Batching: Multiple client operations batched into single log entries

• **Read Batching:** Follower read requests batched for improved throughput

• **Configurable Batch Size:** Adaptive batching based on load patterns

• **Priority-based Elections:** Preferred leaders based on node capabilities

← Direct memtable flushes

← Large values stored separately

• Write Amplification: Minimized through level-based compaction

• Value Threshold: Small values stored in LSM, large values in value log

• **Bloom Filters:** Reduce read amplification for non-existent keys

• **Block Cache:** LRU cache for frequently accessed blocks

← Compacted, sorted

← Higher levels

operation: "add_node" target: node_id: "node-4" address: "10.0.1.13:7000" role: "voter"

• **Pre-vote Phase:** Prevents unnecessary leader elections

• **Fast Recovery:** Optimized recovery from network partitions

• Consistency: Transactions maintain data integrity constraints

• Durability: Committed transactions survive system failures

Modified Raft algorithm with optimizations for key-value workloads

• Batched log entries for improved throughput

• Persistent data storage with ACID guarantees

• Automatic compaction and garbage collection

• Efficient key-value operations with prefix scanning

• LSM-tree based storage for write-optimized performance

• Memory-mapped value log for read optimization

• Configurable compression (Snappy, Zstandard)

3. Storage Engine (BadgerDB)

Backup and restore capabilities

• Encryption at rest with AES-256

Responsibilities:

Key Features:

Value Structure

Transaction Model

Example Transaction

Metadata map[string]string }

3. Leader Election Optimization

4. Membership Changes

5. Storage Engine

BadgerDB Architecture

1. LSM-Tree Structure

Disk Components:

Level 0

(Unsorted)

Level 1 (Sorted)

Level N

(Sorted)

Value Log

(Append-only)

Value Log:

2. Key Features

3. Performance Tuning

HTTP API for broader compatibility:

8. Performance Characteristics

Performance metrics

1. Latency Metrics

Operation Type

Single PUT

Single GET

Batch PUT (10)

Batch GET (10)

Cluster Size

3 nodes

5 nodes

7 nodes

Transaction (5 ops)

2. Throughput Capabilities

3. Scalability Characteristics

9. Security Model

TLS Configuration

RBAC System

1. Authentication & Authorization

"TLS_CHACHA20_POLY1305_SHA256"

• **Key Rotation:** Automatic daily rotation

Encryption in Transit

3. Audit & Compliance

\${AUDIT_TOKEN}"

10. Operational Concerns

1. Deployment Architecture

requests: storage: 100Gi

2. Monitoring & Observability

'kvstore-2:2112'] scrape_interval: 15s metrics_path: /metrics

• **Performance:** Latency percentiles, throughput, error rates

• **Storage:** Disk usage, compaction stats, cache hit rates

• Consensus: Leader elections, log replication lag, commit latency

1. **Point-in-time Recovery:** Restore from snapshot + replay logs

resources: cpu: 0.5 cores memory: 1GB storage: 10GB SSD network: 100 Mbps

resources: cpu: 4-8 cores memory: 16-32GB storage: 1TB NVMe SSD network: 10 Gbps

2. **Cross-region Recovery:** Replicate backups across regions

3. **Disaster Recovery:** Automated failover to backup region

Metrics Collection

Key Metrics

Recovery Procedures

4. Capacity Planning

Resource Requirements

Production Configuration

Minimum Configuration (Development)

Phase 2: Advanced Capabilities (Q2 2024)

Phase 3: Enterprise Features (Q3 2024)

12. Conclusion

contemporary applications.

high throughput

1. **Performance:** Sub-millisecond latencies with

2. **Reliability:** Strong consistency with automatic

Key Strengths

• **Cross-region Replication:** Active-active multi-region deployments

• Schema Evolution: Structured data with schema validation

• Time-series Support: Optimized storage for time-series data

• LDAP Integration: Enterprise directory service integration

• Automated Tuning: ML-based performance optimization

• Edge Computing: Lightweight edge node deployment

• Advanced Analytics: Built-in analytics and reporting

• **Graph Queries:** Basic graph traversal capabilities

Growth Projections

Kubernetes Deployment

• **Protocol:** TLS 1.3 minimum

• Cipher Suites: AEAD ciphers only

• **Scope:** All stored data including logs and snapshots

• **Certificate Management:** Automated with cert-manager

• mTLS: Enforced for inter-node communication

• **Read Scaling:** Linear with follower count

2. Batching Strategy

• Pre-vote mechanism to reduce leader election disruption

Configuration checkpointing for fast cluster recovery

4. Consensus Algorithm **Raft Implementation Details** Our Raft implementation includes several optimizations for key-value store workloads: 1. Log Structure Optimization type LogEntry struct { Index uint64 Term uint64 Timestamp time.Time Type EntryType Operations []Operation Checksum

string } type Operation struct { Type OpType // PUT, DELETE, BATCH Key string Value []byte Conditions []Condition

transaction: operations: - type: PUT key: "user:123:profile" value: "{\"name\": \"John Doe\", \"email\":

\"john@example.com\"}" conditions: - if_version_match: 5 - type: DELETE key: "user:123:temp_data" - type: PUT key:

type Value struct { Data []byte // Binary data, max 16MB ContentType string // MIME type hint Encoding string //

Memory Components: Memtable ← Active writes (Skip List) Immutable ← Sealed memtable Memtable

cluster_membership: current_config: - node_id: "node-1" address: "10.0.1.10:7000" role: "voter" - node_id: "node-2" address: "10.0.1.11:7000" role: "voter" - node_id: "node-3" address: "10.0.1.12:7000" role: "voter" pending_config:

7. API Design 1. gRPC API The primary API uses Protocol Buffers for type safety and performance: service KVStore { // Core operations rpc Put(PutRequest) returns (PutResponse); rpc Get(GetRequest) returns (GetResponse); rpc Delete(DeleteRequest) returns (DeleteResponse); // Batch operations rpc BatchPut(BatchPutRequest) returns (BatchPutResponse); rpc BatchGet(BatchGetRequest) returns (BatchGetResponse); // Advanced operations rpc List(ListRequest) returns (ListResponse); rpc Watch(WatchRequest) returns (stream WatchResponse); rpc Transaction(TransactionRequest) returns (TransactionResponse); // Administrative operations rpc Backup(BackupRequest) returns (BackupResponse); rpc Restore(RestoreRequest) returns (RestoreResponse); rpc Status(StatusRequest) returns (StatusResponse); } 2. REST API

endpoints: # Core operations PUT /api/v1/keys/ $\{key\}$ # Create/update key GET /api/v1/keys/ $\{key\}$ # Retrieve key DELETE /api/v1/keys/{key} # Delete key HEAD /api/v1/keys/{key} # Check existence # Batch operations POST /api/v1/batch/put # Batch create/update POST /api/v1/batch/get # Batch retrieve POST /api/v1/batch/delete # Batch delete # Advanced

/api/v1/admin/restore # Restore from backup GET /api/v1/admin/status # Cluster status GET /api/v1/admin/metrics #

2ms

1ms

4ms

12ms

Read QPS

200K

400K

600K

P99.9

10ms

8ms

30ms

20ms

50ms

Total QPS

250K

480K

700K

5ms

3ms

15ms

10ms

25ms

operations GET /api/v1/keys # List keys (with prefix) POST /api/v1/transaction # Execute transaction GET /api/v1/watch/{prefix} # Server-sent events # Administrative POST /api/v1/admin/backup # Create backup POST

0.5ms

0.3ms

2ms

1ms

Write QPS

50K

80K

100K

• **Write Scaling:** Limited by leader capacity, typically 50-100K writes/sec

• **Memory:** Configurable cache sizes, typically 10-50% of working set

• **Storage:** Supports PB-scale data with appropriate hardware

["developer"] auth_method: "jwt" namespace_prefix: "app:"

badger_config: # Memory settings memtable_size: 64MB block_cache_size: 256MB index_cache_size: 128MB # Compaction settings num_memtables: 5 num_level_zero_tables: 5 level_size_multiplier: 10 # Value log settings value_threshold: 1024 value_log_file_size: 1GB # Performance settings sync_writes: false num_compactors: 4 compression: snappy

2. Encryption **Encryption at Rest** • Algorithm: AES-256-GCM • **Key Management:** External KMS or Vault integration

tls: enabled: true cert_file: "/etc/certs/server.crt" key_file: "/etc/certs/server.key" ca_file: "/etc/certs/ca.crt"

client_auth: "require_and_verify_client_cert" min_version: "1.3" cipher_suites: - "TLS_AES_256_GCM_SHA384" -

roles: - name: "admin" permissions: - "kvstore:*" - "cluster:*" - "backup:*" - name: "developer" permissions: -

"kvstore:list" users: - username: "admin" roles: ["admin"] auth_method: "certificate" - username: "app-service" roles:

audit: enabled: true events: - "authentication_failure" - "authorization_failure" - "admin_operations" - "data_access" destinations: - type: "file" path: "/var/log/kvstore/audit.log" format: "json" - type: "syslog" facility: "LOG_AUTH" severity: "LOG_INFO" - type: "webhook" url: "https://siem.company.com/api/events" headers: Authorization: "Bearer

apiVersion: apps/v1 kind: StatefulSet metadata: name: kvstore spec: serviceName: kvstore replicas: 3 template: spec: containers: - name: kvstore image: kvstore:latest resources: requests: cpu: 2 memory: 8Gi storage: 100Gi limits: cpu:

volumeClaimTemplates: - metadata: name: data spec: accessModes: ["ReadWriteOnce"] storageClassName: "ssd" resources:

prometheus: scrape_configs: - job_name: 'kvstore' static_configs: - targets: ['kvstore-0:2112', 'kvstore-1:2112',

4 memory: 16Gi volumeMounts: - name: data mountPath: /data - name: config mountPath: /etc/kvstore

"kvstore:read" - "kvstore:write" - "kvstore:list:app:*" - name: "readonly" permissions: - "kvstore:read" -

• System: CPU, memory, network, disk I/O **Alerting Rules** groups: - name: kvstore rules: - alert: KVStoreHighLatency expr: histogram_quantile(0.95, kvstore_request_duration_seconds) > 0.1 for: 5m labels: severity: warning annotations: summary: "KVStore high latency $\label{lem:detected} \mbox{ detected" - alert: KVStoreLeaderElection expr: increase(kvstore_leader_elections_total[5m]) > 0 for: 1m labels: \\ \mbox{ detected" - alert: KVStoreLeaderElection expr: increase(kvstore_leader_elections_total[5m]) > 0 for: 1m labels: \\ \mbox{ detected" - alert: KVStoreLeaderElection expr: increase(kvstore_leader_elections_total[5m]) > 0 for: 1m labels: \\ \mbox{ detected" - alert: KVStoreLeaderElection expr: increase(kvstore_leader_elections_total[5m]) > 0 for: 1m labels: \\ \mbox{ detected" - alert: KVStoreLeaderElection expr: increase(kvstore_leader_elections_total[5m]) > 0 for: 1m labels: \\ \mbox{ detected" - alert: KVStoreLeaderElection expr: increase(kvstore_leader_elections_total[5m]) > 0 for: 1m labels: \\ \mbox{ detected" - alert: KVStoreLeaderElection expr: increase(kvstore_leader_elections_total[5m]) > 0 for: 1m labels: \\ \mbox{ detected" - alert: KVStoreLeaderElection expr: increase(kvstore_leader_elections_total[5m]) > 0 for: 1m labels: \\ \mbox{ detected" - alert: KVStoreLeaderElection expr: increase(kvstore_leader_elections_total[5m]) > 0 for: \\ \mbox{ detected" - alert: KVStoreLeaderElection expr: increase(kvstore_leader_elections_total[5m]) > 0 for: \\ \mbox{ detected" - alert: KVStoreLeaderElection expr: increase(kvstore_eleader_elections_total[5m]) > 0 for: \\ \mbox{ detected" - alert: KVStoreLeaderElection expr: increase(kvstore_eleader_elead$ severity: critical annotations: summary: "KVStore leader election occurred" 3. Backup & Recovery **Backup Strategy**

backup: schedule: "0 2 * * * " # Daily at 2 AM retention: daily: 7 weekly: 4 monthly: 12 storage: type: "s3" bucket: "kvstore-backups" encryption: "AES256" compression: "zstd" verification: enabled: true sample_rate: 0.1 # Verify 10%

capacity_planning: metrics: - name: "storage_growth" current: "500GB" growth_rate: "10GB/month" projection_months: 12 threshold_warning: "80%" threshold_critical: "90%" - name: "request_rate_growth" current: "10K QPS" growth_rate: "15%/quarter" projection_quarters: 4 threshold_warning: "70% of capacity" threshold_critical: "85% of capacity"

11. Future Roadmap Phase 1: Enhanced Features (Q1 2024) • Secondary Indexes: Support for custom indexing strategies • **Streaming API:** Real-time change streams • **Multi-tenancy:** Namespace isolation with resource quotas • **Compression:** Advanced compression algorithms (LZ4, Zstandard)

Phase 4: Ecosystem Integration (Q4 2024) • Kafka Integration: Change data capture to Kafka • **Spark Connector:** Direct integration with Apache Spark • Kubernetes Operator: Advanced cluster lifecycle management • Service Mesh: Native Istio/Envoy integration

The Distributed Key-Value Store represents a modern approach to distributed data storage, combining proven algorithms

Design Trade-offs

1. Consistency vs. Availability: Strong consistency

2. Write Scalability: Single leader limits write

may impact availability during network partitions

(Raft consensus) with high-performance storage engines (BadgerDB) to deliver a robust, scalable solution for

failover throughput 3. **Scalability:** Linear read scaling with efficient 3. **Memory Usage:** LSM trees require significant write handling memory for optimal performance 4. Operability: Comprehensive monitoring and 4. Complexity: Distributed consensus adds self-healing capabilities operational complexity 5. **Security:** Defense-in-depth with encryption and **RBAC** This system design provides a solid foundation for building reliable, high-performance distributed applications while maintaining operational simplicity and strong consistency guarantees. @ 2025 System Architecture Team | Distributed Key-Value Store Design Document v1.0 This document contains confidential and proprietary information. Distribution is restricted to authorized personnel only.

This document provides a comprehensive technical overview of the Distributed Key-Value Store (KVStore) system - a high-performance, fault-tolerant distributed database designed for modern cloud-native applications. The system • **Fault Tolerance:** Automatic failover with no data loss in minority failure scenarios • Horizontal Scalability: Dynamic cluster membership with intelligent load balancing

Load Balancer (gRPC/HTTP Multiplexing) KVStore Cluster Node 2 Node 3 Node 1 (Leader) (Follower) |← (Follower)

API API Server Server Server |BadgerDB | BadgerDB BadgerDB Storage Storage Storage

Date: August 2025

Status: Final Draft