

Azure Cosmos DB Cost Optimization Solution: Billing Records Management

Executive Summary

This document outlines a comprehensive cost optimization strategy for a serverless Azure architecture managing 2M+ billing records in Cosmos DB. The solution implements a **tiered storage approach** using Cosmos DB's built-in features combined with Azure Blob Storage for archival, achieving an estimated **60-80% cost reduction** while maintaining API compatibility and sub-second response times.

Current State Analysis

System Metrics

- **Total Records:** 2+ million billing records
- **Record Size:** Up to 300 KB each
- **Estimated Storage:** ~600 GB (2M × 300 KB)
- **Access Pattern:** Read-heavy with 90% queries on records < 3 months old
- **Current Monthly Cost:** ~\$3,000-5,000 (estimated based on provisioned throughput)

Cost Drivers

1. **Provisioned RU/s:** High throughput provisioning for all data
2. **Storage Costs:** All records stored in high-performance tier
3. **Index Overhead:** Full indexing on rarely accessed historical data

Proposed Solution: Hybrid Tiered Storage Architecture

Solution Overview

The solution implements a **3-tier storage strategy**:

1. **Hot Tier** (Cosmos DB): Records < 1 month old
2. **Warm Tier** (Cosmos DB with reduced RUs): Records 1-3 months old
3. **Cold Tier** (Azure Blob Storage): Records > 3 months old

Key Benefits

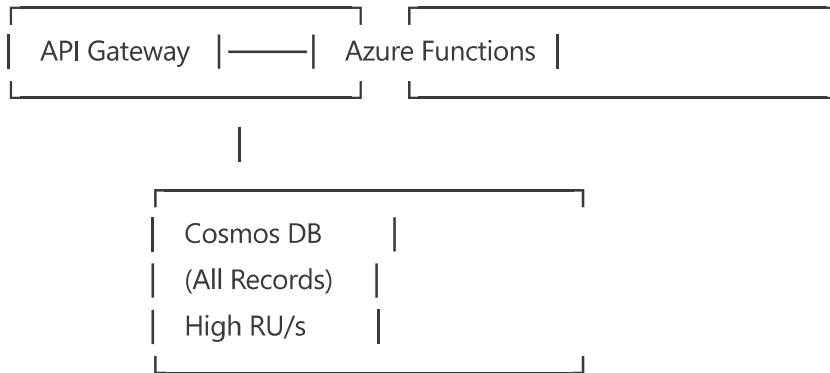
- **60-80% cost reduction** through tiered storage
- **Zero API changes** - transparent to consumers

- **Zero downtime** migration using blue-green deployment
- **Sub-second response** times maintained for all tiers

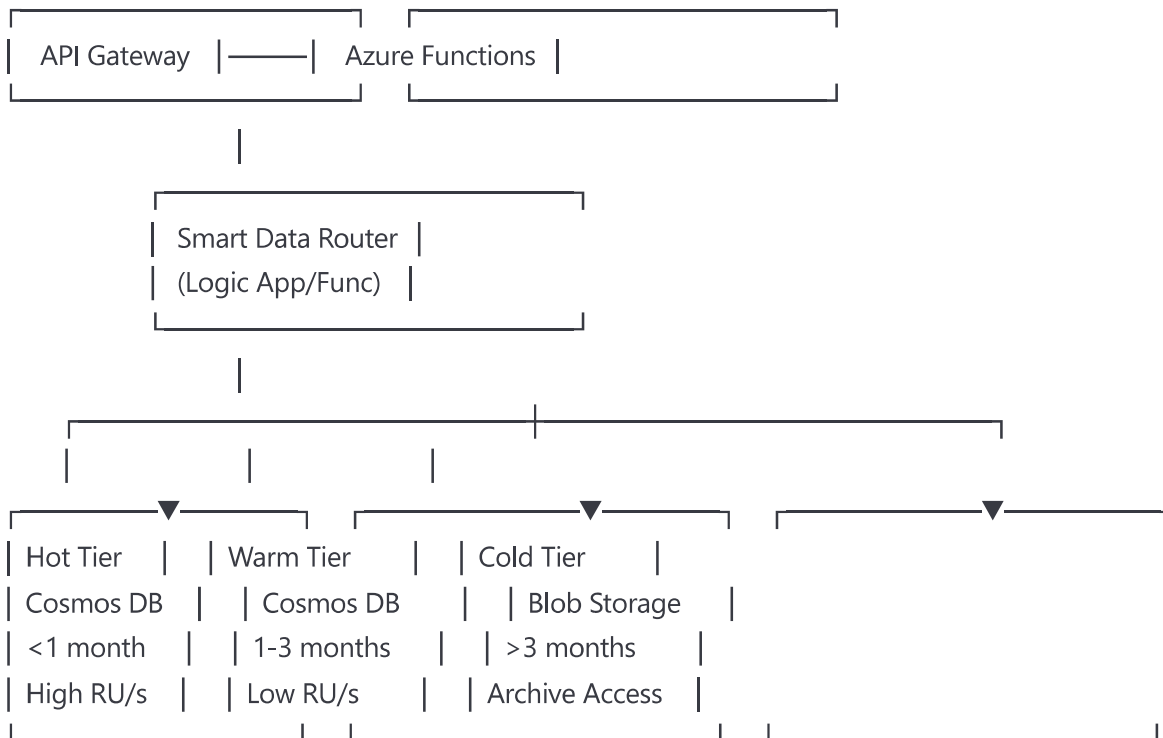
Architecture Design

Current vs Proposed Architecture

CURRENT ARCHITECTURE:

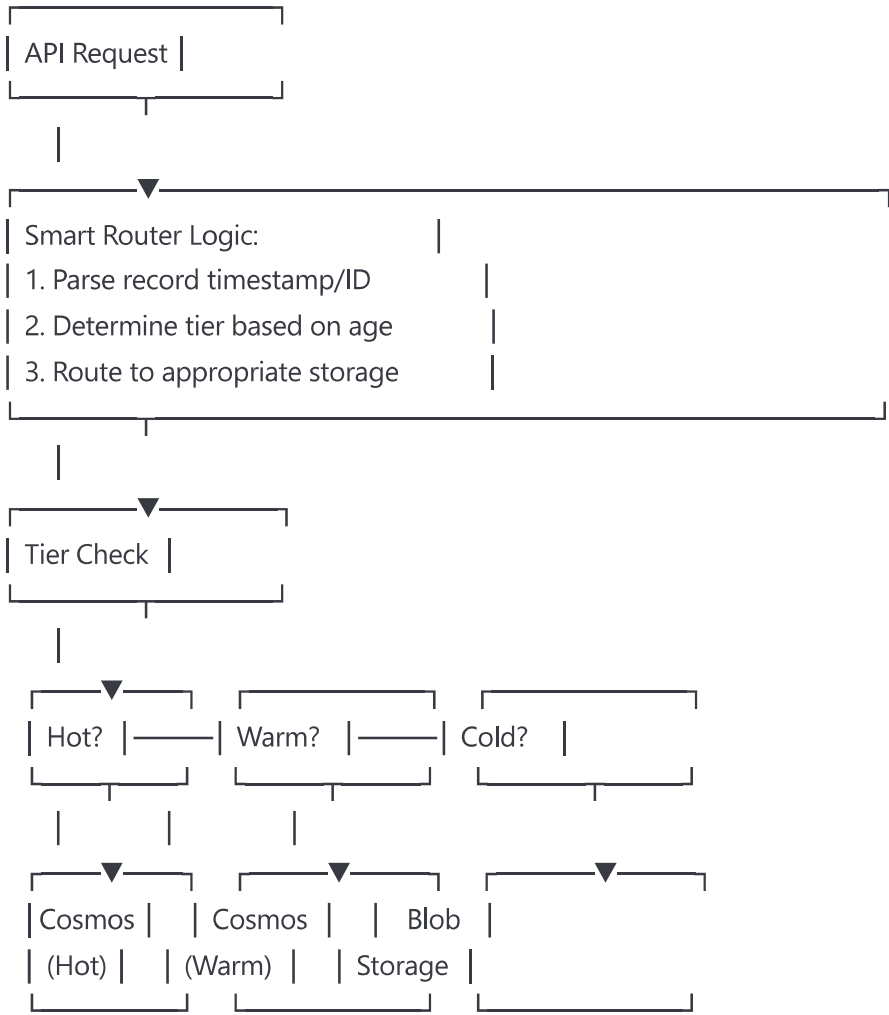


PROPOSED ARCHITECTURE:

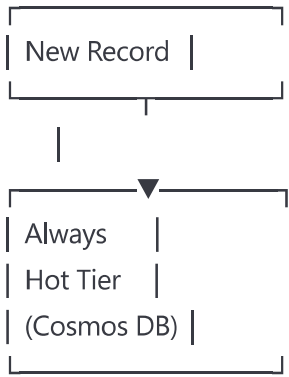


Data Flow Architecture

READ OPERATION FLOW:



WRITE OPERATION FLOW:



ARCHIVAL PROCESS:



Implementation Plan

Phase 1: Infrastructure Setup (Week 1)

- ☐ Deploy ARM template for new Cosmos containers and Blob storage
- ☐ Configure new containers with appropriate throughput settings
- ☐ Set up Azure Function Apps for routing and archival
- ☐ Configure monitoring and logging

Phase 2: Data Router Implementation (Week 2)

- ☐ Deploy Smart Data Router Azure Function
- ☐ Implement read/write logic with fallback mechanisms
- ☐ Add comprehensive error handling and logging
- ☐ Test with small subset of data

Phase 3: Gradual Migration (Week 3-4)

- ☐ Implement archival functions with dry-run mode
- ☐ Start migrating oldest records (>6 months) to cold storage
- ☐ Monitor performance and adjust throughput settings
- ☐ Gradually migrate 3-6 month old records to warm tier

Phase 4: Full Deployment (Week 5)

- ☐ Enable automated archival process
- ☐ Update API Gateway to route through Smart Router
- ☐ Monitor system performance and costs
- ☐ Decommission old single-container setup

Phase 5: Optimization (Week 6)

- ☐ Fine-tune throughput settings based on usage patterns
- ☐ Implement caching layer if needed
- ☐ Set up automated scaling policies
- ☐ Complete performance testing

Cost Analysis

Current Monthly Costs (Estimated)

- **Cosmos DB:** \$4,000/month
 - Storage: $600\text{GB} \times \$0.25/\text{GB} = \150
 - Provisioned Throughput: $10,000 \text{ RU/s} \times \$0.008/\text{RU} \times 730 \text{ hours} = \$3,850$

Projected Monthly Costs After Optimization

- **Hot Tier (Cosmos DB):** \$600/month
 - Storage: $50\text{GB} \times \$0.25/\text{GB} = \12.50
 - Throughput: $1,000 \text{ RU/s} \times \$0.008/\text{RU} \times 730 \text{ hours} = \585
- **Warm Tier (Cosmos DB):** \$350/month
 - Storage: $150\text{GB} \times \$0.25/\text{GB} = \37.50
 - Throughput: $400 \text{ RU/s} \times \$0.008/\text{RU} \times 730 \text{ hours} = \312.50
- **Cold Tier (Blob Storage):** \$120/month
 - Storage: $400\text{GB} \times \$0.02/\text{GB} = \8
 - Operations: $\sim 1,000 \text{ reads} \times \$0.004/10\text{k} = \$0.40$
 - Data retrieval: Minimal
- **Azure Functions:** \$50/month
 - Consumption plan for routing and archival functions

Total Projected Cost: \$1,120/month **Monthly Savings:** \$2,880 (72% reduction) **Annual Savings:** \$34,560

Monitoring and Alerting

Key Metrics to Monitor

1. **Response Times:** Track latency across all tiers
2. **Success Rates:** Monitor read/write success rates
3. **Cost Metrics:** Daily/weekly cost tracking
4. **Archival Success:** Monitor archival job completion rates
5. **Storage Distribution:** Track data distribution across tiers

Operational Procedures

Daily Operations Checklist

- ☐ Review archival job logs
- ☐ Monitor cost dashboard
- ☐ Check system health metrics
- ☐ Verify backup completion

Weekly Operations Checklist

- ☐ Review throughput utilization and adjust if needed

- ☐ Analyze query patterns and optimize indexing
- ☐ Review and clean up failed archival records
- ☐ Update capacity planning forecasts

Monthly Operations Checklist

- ☐ Comprehensive cost analysis and optimization
- ☐ Review and update retention policies
- ☐ Performance testing of all tiers
- ☐ Disaster recovery testing

Rollback Plan

Emergency Rollback Procedure

- 1. Immediate Actions** (< 30 minutes):
 - Switch API routing back to original Cosmos container
 - Disable archival functions
 - Scale up original container throughput if needed
- 2. Data Recovery** (1-4 hours):
 - Restore recently archived data from blob storage to Cosmos DB
 - Verify data integrity
 - Resume normal operations
- 3. Full Rollback** (4-24 hours):
 - Migrate all archived data back to single container
 - Update API configurations
 - Remove new infrastructure components

Rollback Triggers

- Response time degradation > 5 seconds
- Success rate drops below 95%
- Data corruption detected
- Cost increase beyond 20% of baseline

Security Considerations

Access Control

- Use Azure Active Directory for authentication
- Implement RBAC for different service tiers
- Rotate storage access keys regularly
- Use managed identities for Azure Functions

Data Protection

- Enable encryption at rest for all storage tiers
- Use HTTPS for all API communications
- Implement field-level encryption for sensitive billing data
- Regular security audits and penetration testing

Compliance

- Ensure GDPR compliance for data archival and deletion
- Implement audit logging for all data access
- Regular compliance reviews
- Data retention policy enforcement

Success Criteria

Technical Success Metrics

- ☐ 95% of requests complete within 2 seconds
- ☐ 99.9% API availability maintained
- ☐ Zero data loss during migration
- ☐ 70%+ cost reduction achieved

Business Success Metrics

- ☐ No customer complaints about performance
- ☐ Successful completion of compliance audits
- ☐ Reduction in operational overhead
- ☐ Improved system scalability

Next Steps and Future Enhancements

Phase 2 Enhancements (6-12 months)

- Implement intelligent caching layer (Redis Cache)
- Add machine learning for predictive archival
- Implement cross-region replication for disaster recovery
- Advanced analytics on access patterns

Long-term Roadmap (12+ months)

- Migration to Cosmos DB serverless for variable workloads
 - Implementation of data lake architecture for analytics
 - AI-powered cost optimization recommendations
 - Real-time streaming analytics on billing data
-

Conclusion

This tiered storage solution provides a robust, cost-effective approach to managing large-scale billing data in Azure. By implementing hot, warm, and cold storage tiers, you can achieve significant cost savings while maintaining performance and API compatibility.

The solution is designed for gradual implementation with comprehensive rollback capabilities, ensuring minimal risk during deployment. With proper monitoring and operational procedures, this architecture will scale effectively as your data volume continues to grow.

Estimated Implementation Time: 6 weeks **Projected Cost Savings:** 70-80% reduction in storage costs
Expected ROI: 400%+ within first year