

Chapter 7: Databases on AWS

Rahul Daware

June 6, 2018

Contents

1	RDS - Backups, Multi-AZ and Read Replicas	3
1.1	Automated Backups	3
1.2	Database Snapshots	3
1.3	Restoring Backups	3
1.4	Encryption	3
1.5	What is Multi-AZ RDS	3
1.6	Read Replicas	4
2	DynamoDB	4
2.1	Eventual Consistent Reads	4
2.2	Strongly Consistent Reads	4
3	Amazon Redshift	4
3.1	Columnar Data Storage	4
3.2	Advanced Compression	5
3.3	Massively Parallel Processing (MPP)	5
3.4	Redshift Security	5
4	Elasticache	5
4.1	ElastiCache Exam Tips	5
5	Aurora	5
5.1	Aurora Scaling	6
6	DynamoDB vs RDS	6

1 RDS - Backups, Multi-AZ and Read Replicas

There are two types of backups for AWS : Automated Backups and Database Snapshots

1.1 Automated Backups

Automated Backups allow you to recover your database to any point in time within a "retention period". The retention period can be between one and 35 days. Automated backups will take a full daily snapshot and will also store transaction logs throughout the day. When you do a recovery, AWS will first choose the most recent daily back up , and then apply transaction logs relevant to that day. This allows you to do a point in time recovery down to a second, within the retention period.

Automated backups are enabled by default. The backup data is stored in S3 and you get free storage equal to the size of your database. So if you have an RDS instance of 10GB, you will get 10GB worth of storage. Backups are taken within a defined window. During the backup window, storage I/O may be suspended while your data is being backed up and you may experience elevated latency.

1.2 Database Snapshots

DB Snapshots are done manually(i.e. they are user initiated). They are stored even after you delete the original RDS instance, unlike automated backups.

1.3 Restoring Backups

Whenever you restore either an automatic backup or a manual snapshot, the restored version of the database will be a new RDS instance with a new DNS endpoint.

1.4 Encryption

Encryption at rest is supported for MySQL, Oracle, SQL Server, PostgreSQL, MariaDB and Aurora. Encryption is done using the AWS key management service (KMS) service. Once your RDS instance is encrypted, the data stored at rest in the underlying storage is encrypted, as are its automated backups, read replicas and snapshots. At the present time, encrypting an existing DB instance is not supported. To use Amazon RDS encryption for an existing database, you must first create a snapshot make a copy of that snapshot and encrypt the copy.

1.5 What is Multi-AZ RDS

Multi-AZ allows you to have an exact copy of your production database in another availability zone. AWS handles the replication for you, so when your production database is written to, this write will automatically be synchronized to the stand by database. In the event of planned database maintenance, DB instance failure, or an availability zone failure, Amazon RDS will automatically failover to the standby so that database operations can resume quickly without administrative intervention. Multi-AZ is for disaster recovery only. It is not primarily used for improving performance, For performance improvement, you need read replicas. Multi AZ is available for:

- SQL server
- Oracle
- MySQL Server
- PostgreSQL
- MariaDB

1.6 Read Replicas

Read replicas allow you to have a read-only copy of your production database. This is achieved by using asynchronous replication from the primary RDS instance to the read replica. You use read replicas primarily for very read-heavy database workloads. Read replicas are available for :

- MySQL Server
- PostgreSQL
- MariaDB
- Aurora

Read replicas are used for scaling and not for DR. Must have automatic backups turned on in order to deploy a read replica. You can have up to 5 read replica copies of any database. You can have read replicas (but watch out for latency). Each read replica will have its own DNS endpoint. You can have read replicas that have Multi-AZ. You can create read replicas of Multi-AZ source databases. Read replicas can be promoted to be their own databases. This breaks the replication. You can have a read replica in a second region.

2 DynamoDB

Amazon DynamoDB is a fast and flexible NoSQL database service for all applications that need consistent, single digit millisecond latency at any scale. It is a fully managed database and supports both document and key-value data models. Its flexible data model and reliable performance make it a great fit for mobile, web, gaming, ad-tech, IoT, and many other applications.

- Stored on SSD storage
- Spread across 3 geographically distinct data centers
- Eventual consistent reads (default)
- Strongly consistent reads

2.1 Eventual Consistent Reads

Consistency across all copies of data is usually reached within a second. Repeating a read after a short time should return the updated data. (Best Read Performance)

2.2 Strongly Consistent Reads

A strongly consistent read returns a result that reflects all writes that received a successful response prior to the read

3 Amazon Redshift

Amazon Redshift is a fast and powerful, fully managed, petabyte-scale data warehouse service in the cloud. Customers can start small for just \$0.25 per hour with no commitments or upfront costs and scale to a petabyte or more for \$1,000 per terabyte per year, less than a tenth of most other data warehousing solutions

3.1 Columnar Data Storage

Instead of storing data as a series of rows, Amazon Redshift organizes the data by column. Unlike row-based systems, which are ideal for transaction processing, column-based systems are ideal for data warehousing and analytics, where queries often involve aggregates performed over large data sets. Since only the columns involved in the queries are processed and columnar data is stored sequentially on the storage media, column-based systems require far fewer I/Os, greatly improving query performance.

3.2 Advanced Compression

Columnar data stores can be compressed much more than row-based data stores because similar data is stored sequentially on disk. Amazon Redshift employs multiple compression techniques and can often achieve significant compression relative to traditional relational data stores. In addition, Amazon Redshift doesn't require indexes of materialized views and so uses less space than traditional relational database systems. When loading data into an empty table, Amazon Redshift automatically samples your data and selects the most appropriate compression scheme.

3.3 Massively Parallel Processing (MPP)

Amazon Redshift automatically distributes data and query load across all nodes. Amazon Redshift makes it easy to add nodes to your data warehouse and enables you to maintain fast query performance as your data warehouse grows.

3.4 Redshift Security

- Encrypted in transit using SSL
- Encrypted at rest using AES-256 encryption
- By default RedShift takes care of key management
- You can also manage your own keys through HSM
- Also through AWS Key Management
- Currently only available in 1 AZ
- Can restore snapshots to new AZs in the event of an outage

4 Elasticache

ElastiCache is a web service that makes it easy to deploy, operate, and scale an in-memory cache in the cloud. The service improves the performance of web applications by allowing you to retrieve information from fast, managed, in-memory caches, instead of relying entirely on slower disk-based databases

Types of ElastiCache

- Memcached - A widely adopted memory object caching system. ElastiCache is protocol compliant with Memcached, so popular tools that you use today with existing Memcached environments will work seamlessly with the service.
- Redis - A popular open-source in-memory key value store that supports data structures such as sorted sets and lists. ElastiCache supports Master/Slave replication and Multi-AZ which can be used to achieve cross AZ redundancy.

4.1 ElastiCache Exam Tips

Typically you will be given a scenario where a particular database is under a lot of stress/load. You may be asked which service you should use to alleviate this. Elasticache is a good choice if your database is particularly read heavy and not prone to frequent changing. Redshift is a good answer if the reason your database is feeling stress is because management keep running OLAP transactions on it etc.

5 Aurora

Amazon Aurora is a MySQL-compatible, relational database engine that combines the speed and availability of high-end commercial databases with the simplicity and cost effectiveness of open source databases. Amazon Aurora provides up to five times better performance than MySQL at a price point one tenth that of a commercial database while delivering similar performance and availability.

5.1 Aurora Scaling

- Start with 10 GB, scales in 10GB increments to 64TB (Storage autoscaling)
- Compute resources can scale up to 32 vCPUs and 244GB of memory
- 2 copies of your data is contained in each availability zone, with minimum of 3 availability zones. 6 copies of your data
- Aurora is designed to transparently handle the loss of up to two copies of data without affecting database write availability and up to three copies without affecting read availability
- Aurora storage is also self-healing, Data blocks and disks are continuously scanned for errors and repaired automatically

6 DynamoDB vs RDS

DynamoDB offers "push button" scaling, meaning that you can scale your database on the fly, without any down time. RDS is not so easy.- you usually have to use a bigger instance size or to add a read replica.