



Query best practices for Amazon Redshift

AWS Prescriptive Guidance



AWS Prescriptive Guidance: Query best practices for Amazon Redshift

Copyright © 2026 Amazon Web Services, Inc. and/or its affiliates. All rights reserved.

Amazon's trademarks and trade dress may not be used in connection with any product or service that is not Amazon's, in any manner that is likely to cause confusion among customers, or in any manner that disparages or discredits Amazon. All other trademarks not owned by Amazon are the property of their respective owners, who may or may not be affiliated with, connected to, or sponsored by Amazon.

Table of Contents

Introduction	1
Overview	1
Intended audience	1
Objectives	1
Architecture components	2
Query performance factors	7
Table properties	7
Sort keys	7
Data compression	8
Data distribution	8
Table maintenance	8
Cluster configuration	9
Node type	10
Node size, number of nodes, and slices	10
Workload management	10
Short query acceleration	10
SQL query	11
Query structure	11
Code compilation	11
Best practices for tables	12
Understand how sort keys work	12
Query tuning tips	12
Evaluate sort key effectiveness	13
Know your table	13
Choose the right table distribution style	14
Best practices for queries	15
Avoid using the SELECT * FROM statement	15
Identify query issues	15
Get summary information on your query	15
Avoid cross-joins	15
Avoid functions in query predicates	16
Avoid unnecessary cast conversions	16
Use CASE expressions for complex aggregations	17
Use subqueries	17

Use predicates	18
Add predicates to filter tables with joins	18
Use the least expensive operators for predicates	19
Use sort keys in GROUP BY clauses	19
Take advantage of materialized views	19
Be careful with columns in GROUP BY and ORDER BY clauses	19
Best practices for Redshift Spectrum	21
Predicate pushdown in Redshift Spectrum	22
Query tuning tips for Redshift Spectrum	23
Resources	24
Document history	25
Glossary	26
#	26
A	27
B	30
C	32
D	35
E	39
F	41
G	43
H	44
I	46
L	48
M	49
O	53
P	56
Q	59
R	59
S	62
T	66
U	67
V	68
W	68
Z	69

Query best practices for Amazon Redshift

Ethan Stark, Amazon Web Services (AWS)

June 2024 ([document history](#))

Overview

This guide provides recommendations and best practices for optimizing query and table performance in [Amazon Redshift](#). You can use Amazon Redshift to query petabytes of structured and semi-structured data across your data warehouse and your data lake by using standard SQL. This guide also provides an overview of the core architecture components of an Amazon Redshift data warehouse. This knowledge—along with an understanding of query performance factors such as table properties, cluster configuration, and query structure—can help you design efficient and effective tables and queries for your Amazon Redshift data warehouse.

Intended audience

This guide is intended for data engineers, data architects, and data analysts who design or use tables and queries in Amazon Redshift.

Objectives

This guide can help you and your organization achieve the following objectives:

- Design tables for optimal data storage and retrieval operations
- Design queries for optimal performance and cost savings
- Optimize the performance of [Amazon Redshift Spectrum](#) to query data directly from files on [Amazon Simple Storage Service \(Amazon S3\)](#)

Architecture components of an Amazon Redshift data warehouse

We recommend that you have a basic understanding of the core architecture components in an Amazon Redshift data warehouse. This knowledge can help you better understand how to design your queries and tables for optimal performance.

A data warehouse in Amazon Redshift consists of the following core architecture components:

- **Clusters** – A cluster, which is composed of one or more compute nodes, is the core infrastructure component of an Amazon Redshift data warehouse. Compute nodes are transparent to external applications, but your client application interacts directly with the leader node only. A typical cluster has two or more compute nodes. The compute nodes are coordinated through the leader node.
- **Leader node** – A leader node manages the communications for client programs and all compute nodes. A leader node also prepares the plans for running a query whenever a query is submitted to a cluster. When the plans are ready, the leader node compiles code, distributes the compiled code to the compute nodes, and then assigns slices of data to each compute node to process the query results.
- **Compute node** – A compute node runs a query. The leader node compiles code for individual elements of the plan to run the query and assigns the code to individual compute nodes. The compute nodes run the compiled code and send intermediate results back to the leader node for final aggregation. Each compute node has its own dedicated CPU, memory, and attached disk storage. As your workload grows, you can increase the compute capacity and storage capacity of a cluster by increasing the number of nodes, upgrading the node type, or both.
- **Node slice** – A compute node is partitioned into units called slices. Every slice in a compute node is allocated a portion of the node's memory and disk space where it processes a portion of the workload assigned to the node. The slices then work in parallel to complete the operation. Data is distributed among slices on the basis of the [distribution style](#) and distribution key of a particular table. An even distribution of data makes it possible for Amazon Redshift to evenly assign workloads to slices and maximizes the benefit of parallel processing. The number of slices per compute node is decided on the basis of the type of node. For more information, see [Clusters and nodes in Amazon Redshift](#) in the Amazon Redshift documentation.

- **Massively parallel processing (MPP)** – Amazon Redshift uses MPP architecture to quickly process data, even complex queries and vast amounts of data. Multiple compute nodes run the same query code on portions of data to maximize parallel processing.
- **Client application** – Amazon Redshift integrates with various data loading, extract, transform, and load (ETL), business intelligence (BI) reporting, data mining, and analytics tools. All client applications communicate with the cluster through the leader node only.

The following diagram shows how the architecture components of an Amazon Redshift data warehouse work together to accelerate queries.



There are seven stages of the query lifecycle:

1. Query reception and parsing:

- The leader node receives the query and parses the SQL.
- The parser produces an initial query tree, which represents the logical structure of the original query.
- Amazon Redshift feeds this query tree into the query optimizer.

2. Query optimization:

- The optimizer evaluates the query and, if necessary, rewrites it to maximize efficiency.
- This optimization process might involve creating multiple related queries to replace a single one.

3. Query plan generation:

- The optimizer generates a query plan (or multiple plans, if needed) for execution.
- The query plan specifies execution options, such as join types, join order, aggregation methods, and data distribution requirements.

4. Execution engine translation:

- The execution engine translates the query plan into discrete steps, segments, and streams:
 - **Step** – Represents an individual operation required during query execution. Steps can be combined to allow compute nodes to perform queries, joins, or other database operations.
 - **Segment** – Combines several steps that a single process can execute. It is the smallest compilation unit executable by a compute node slice. (A slice is the unit of parallel processing in Amazon Redshift.)
 - **Stream** – A collection of segments distributed across available compute node slices.
- The execution engine generates compiled code based on these steps, segments, and streams. Compiled code runs faster than interpreted code and consumes less compute capacity.
- The leader node broadcasts the compiled code to the compute nodes.

5. Parallel execution:

- This step occurs once for each stream.
- Compute node slices run query segments in parallel.
- During this process, Amazon Redshift optimizes network communication, memory usage, and disk management to pass intermediate results from one query plan step to the next.
- This optimization contributes to faster query execution.

6. Stream processing:

- This step occurs once for each stream.
- The engine creates executable segments for each stream, for efficient parallel processing.

7. Final sorting and aggregation:

- The leader node addresses any final sorting or aggregation that the query requires.
- Once completed, the leader node returns the results to the client.

For information on architecture components, see [Data warehouse system architecture](#) in the Amazon Redshift documentation.

Query performance factors for Amazon Redshift

A number of factors can affect query performance. The following aspects of your data, cluster, and database operations all play a part in how quickly your queries are processed:

- [Table properties](#)
 - [Sort keys](#) (Amazon Redshift Advisor)
 - [Data compression](#) (automated)
 - [Data distribution](#) (automated)
 - [Table maintenance](#) (automated)
- [Cluster configuration](#)
 - [Node type](#)
 - [Node size, number of nodes, and slices](#)
 - [Workload management](#) (automated)
 - [Short query acceleration](#) (automated)
- [SQL query](#)
 - [Query structure](#)
 - [Code compilation](#)

Table properties

Amazon Redshift tables are the fundamental units for storing data in Amazon Redshift, and each table has a set of properties that determine its behavior and accessibility. These properties include sorting, distribution style, compression encoding, and many others. Understanding these properties is crucial for optimizing the performance, security, and cost-effectiveness of Amazon Redshift tables.

Sort keys

Amazon Redshift stores data on disk in sorted order according to a table's sort keys. The query optimizer and the query processor use the information about where the data is located within a compute node to reduce the number of blocks that must be scanned. This improves query speed significantly by reducing the amount of data to process. We recommend that you use sort keys

to facilitate filters in the WHERE clause. For more information, see [Working with sort keys](#) in the Amazon Redshift documentation.

Data compression

Data compression reduces storage requirements, which reduces disk I/O and improves query performance. When you run a query, the compressed data is read into memory and then uncompressed when the query runs. By loading less data into memory, Amazon Redshift can allocate more memory to analyzing the data. Because columnar storage stores similar data sequentially, Amazon Redshift can apply adaptive compression encodings specifically tied to columnar data types. The best way to enable data compression on table columns is by using the AUTO option in Amazon Redshift to apply optimal compression encodings when you load the table with data. To learn more about using automatic data compression, see [Loading tables with automatic compression](#) in the Amazon Redshift documentation.

Data distribution

Amazon Redshift stores data on the compute nodes according to a table's distribution style. When you run a query, the query optimizer redistributes the data to the compute nodes as needed to perform any joins and aggregations. Choosing the right distribution style for a table helps minimize the impact of the redistribution step by locating the data where it needs to be before the joins are performed. We recommend that you use distribution keys to facilitate the most common joins. For more information, see [Working with data distribution styles](#) in the Amazon Redshift documentation.

Table maintenance

Although Amazon Redshift provides industry-leading performance out of the box for most workloads, keeping Amazon Redshift clusters running well requires maintenance. Updating and deleting data creates dead rows that must be vacuumed, and even append-only tables must be resorted if the append order isn't consistent with the sort key.

Vacuum

The vacuuming process in Amazon Redshift is essential for the health and maintenance of your Amazon Redshift cluster. It also affects the performance of queries. Because deletes and updates both flag the old data but don't actually remove it, you must use vacuuming to reclaim the disk space that was occupied by table rows that were marked for deletion by the previous UPDATE

and DELETE operations. Amazon Redshift can automatically sort and perform a VACUUM DELETE operation on tables in the background.

To clean up tables after a load or a series of incremental updates, you can also run the VACUUM command, either against the entire database or against individual tables. If tables have sort keys and table loads aren't optimized to sort as they insert, then you must use vacuums to resort the data (which can be crucial for performance). For more information, see [Vacuuming tables](#) in the Amazon Redshift documentation.

Analyze

The ANALYZE operation updates statistical metadata on the tables in an Amazon Redshift database. Keeping statistics current improves query performance by enabling the query planner to choose optimal plans. Amazon Redshift continuously monitors your database and automatically performs analyze operations in the background. To minimize impact to your system performance, the ANALYZE operation automatically runs during periods when workloads are light. If you choose to explicitly run ANALYZE, do the following:

- Run the ANALYZE command before running queries.
- Run the ANALYZE command on the database routinely at the end of every regular load or update cycle.
- Run the ANALYZE command on new tables that you create and existing tables or columns that undergo significant change.
- Consider running ANALYZE operations on different schedules for different types of tables and columns, depending on their use in queries and their propensity to change.
- To save time and cluster resources, use the PREDICATE COLUMNS clause when you run the ANALYZE command.

Cluster configuration

A cluster is a collection of nodes which perform the actual storing and processing of data. Setting up your Amazon Redshift cluster the right way is critical if you want to achieve the following:

- High scalability and concurrency
- Efficient use of Amazon Redshift
- Better performance

- Lower cost

Node type

An Amazon Redshift cluster can use one of several node types (RA3, DC2, and DS2). Each node type offers different sizes and limits to help you scale your cluster appropriately. The node size determines the storage capacity, memory, CPU, and price of each node in the cluster. Cost and performance optimization starts with choosing the right node type and size. For more information about node types, see [Overview of Amazon Redshift clusters](#) in the Amazon Redshift documentation.

Node size, number of nodes, and slices

A compute node is partitioned into slices. More nodes means more processors and slices, which enables your queries to process faster by running portions of the query concurrently across the slices. However, more nodes also means greater expense. This means that you must find the balance of cost and performance that is appropriate for your system. For more information on Amazon Redshift cluster architecture, see [Data warehouse system architecture](#) in the Amazon Redshift documentation.

Workload management

Amazon Redshift workload management (WLM) enables users to flexibly manage workload queues with priorities so that short, fast-running queries won't get stuck in queues behind long-running queries. Automatic WLM uses machine learning (ML) algorithms to profile queries and place them in the appropriate queue with the appropriate resources, while managing query concurrency and memory allocation. For more information about WLM, see [Implementing workload management](#) in the Amazon Redshift documentation.

Short query acceleration

Short query acceleration (SQA) prioritizes short-running queries ahead of long-running queries. SQA runs queries in a dedicated space so that SQA queries aren't forced to wait in queues behind longer queries. SQA only prioritizes queries that are short-running and are in a user-defined queue. If you use SQA, short-running queries begin running more quickly and you can see results sooner. If you enable SQA, you can reduce or eliminate WLM queues that are dedicated to short-running queries. In addition, long-running queries don't need to contend for slots in a WLM queue. This means that you can configure your WLM queues to use fewer query slots. If you use lower

concurrency, query throughput is increased and overall system performance is improved for most workloads. For more information about SQA, see [Working with short query acceleration](#) in the Amazon Redshift documentation.

SQL query

A database query is a request for data from a database. The request should come in an Amazon Redshift cluster using SQL. Amazon Redshift supports SQL client tools that connect through Java Database Connectivity (JDBC) and Open Database Connectivity (ODBC). You can use most SQL client tools that support JDBC or ODBC drivers.

Query structure

How your query is written greatly affects its performance. We recommend that you write queries to process and return as little data as necessary to meet your needs. For more information on how to structure your queries, see the [Best practices for designing Amazon Redshift queries](#) section of this guide.

Code compilation

Amazon Redshift generates and compiles optimized code for each query execution plan. The compiled code runs faster because it eliminates the overhead of using an interpreter. To minimize latency for new queries while preserving the performance benefits of compiled code, Amazon Redshift uses a technique called *composition*. Composition generates a lightweight arrangement of pre-existing logic to process new queries immediately, while simultaneously compiling highly optimized, query-specific code in the background. This removes compilation from the critical path of query execution. This means that new queries start faster and deliver performance consistent with subsequent runs.

Amazon Redshift also uses a serverless compilation service to scale query compilations beyond the compute resources of an Amazon Redshift cluster. Compiled code segments are cached both locally on the cluster and in a virtually unlimited remote cache that persists after cluster reboots. Subsequent executions of the same query run faster because they can skip the compilation phase. By using a scalable compilation service, Amazon Redshift compiles code in parallel to provide consistently fast performance.

Best practices for designing Amazon Redshift tables

This section provides an overview of best practices for designing database tables. We recommend that you follow these best practices to achieve optimal query performance and efficiency.

Understand how sort keys work

Amazon Redshift stores your data on disk in sorted order according to the sort key. The Amazon Redshift query optimizer uses sort order when it determines optimal query plans. To use sort keys effectively, we recommend that you do the following:

- Keep the table sorted as much as possible.
- Use VACUUM sort to restore optimal performance.
- Avoid compressing the sort key column.
- If the sort key is compressed and if the `sortkey1_skew` ratio is significantly high, then recreate the table without enabling compression on the sort key.
- Avoid applying a function to the sort key columns. For example, in the following query, the `trans_dt` : TIMESTAMPTZ sort key column isn't used if you cast it to DATE:

```
select order_id, order_amt
from sales
where trans_dt::date = '2021-01-08'::date
```

- Perform INSERT operations in sort key order.
- Use sort keys in the GROUP BY clause when possible.

Query tuning tips

We recommend that you do the following to tune your queries:

- Always order compound sort keys from lowest-cardinality to highest-cardinality for optimal effectiveness.
- If the leading key in a compound sort key is relatively unique (that is, it has high cardinality), then avoid adding additional columns to your sort key. Adding additional columns has little impact on query performance but does add maintenance costs.

Evaluate sort key effectiveness

To optimize your queries, you must be able to evaluate the effectiveness of your queries. We recommend that you use the [SVL_QUERY_SUMMARY](#) view to find general information about the execution of a query. In this view, you can use the attribute `IS_RRSCAN` to determine if an EXPLAIN plan step uses a range-restricted scan. You can also use the attribute `rows_pre_filter` to determine the selectivity the of a sort key.

You can also use an admin view from GitHub called [v_my_last_query_summary](#). The view displays information for the last query that ran.

The following statement shows how to find general information about the execution of a query.

```
select lpad(' ',stm+seg+step) || label as label,
       rows,
       bytes,
       is_diskbased,
       is_rrscan,
       rows_pre_filter
from svl_query_summary
where query = pg_last_query_id()
order by stm, seg, step;
```

The preceding query returns the following sample output.

label	<input type="checkbox"/> rows	bytes	is_diskbased	is_rrscan	rows_pre_filter
scan tbl=163860 name=orders	<input type="checkbox"/> 1500000	24000000	f	f	1500000
project	<input type="checkbox"/> 1500000	0	f	f	0
project	<input type="checkbox"/> 1500000	0	f	f	0
hash tbl=968	<input type="checkbox"/> 1500000	24000000	f	f	0
scan tbl=163852 name=lineitem	<input type="checkbox"/> 6001215	144029160	f	t	6001215
project	<input type="checkbox"/> 6001215	0	f	f	0
project	<input type="checkbox"/> 6001215	0	f	f	0
hjoin tbl=968	<input type="checkbox"/> 6001215	0	f	f	0
project	<input type="checkbox"/> 6001215	0	f	f	0
project	<input type="checkbox"/> 6001215	0	f	f	0

Know your table

It's important to understand the critical properties of your table. To learn more about your table, do the following:

- Use [PG_TABLE_DEF](#) to view information about table columns.
- Use [SVV_TABLE_INFO](#) to view more comprehensive information about a table, including data distribution skew, key distribution skew, table size, and statistics.

Choose the right table distribution style

When you run a query, the query optimizer redistributes the rows to the compute nodes as needed to perform any joins and aggregations. The goal in selecting a table distribution style is to minimize the impact of the redistribution step by locating the data where it needs to be before you run the query.

We recommend the following approach to choosing the right table distribution style:

- Avoid broadcasting and redistribution in a query execution plan by collocating the rows within the same node. For example, by selecting a DISTKEY, you can distribute the fact table and one-dimension table on their common columns. Choose the largest dimension based on the size of the filtered dataset. Only the rows that are used in the join must be distributed, so consider the size of the dataset after filtering, not the size of the table.
- Make sure there is no skewness on the column where the distribution key is created. Otherwise, one compute node could perform more heavy lifting than others. If you notice skewness, then consider changing the distribution key column. A column can be considered as a candidate for a distribution key if its values are uniformly distributed or high cardinal values.
- If the table used in the join condition is small (less than 1 GB), then consider the distribution style ALL.
- You can compress the distribution key, but you must avoid compressing the sort key column (especially the first column of the sort key).

Note

If you use automatic table optimization, you don't need to choose the distribution style of your table. For more information, see [Working with automatic table optimization](#) in the Amazon Redshift documentation. To have Amazon Redshift choose the appropriate distribution style, specify AUTO for the distribution style.

Best practices for designing Amazon Redshift queries

This section provides an overview of best practices for designing queries. We recommend that you follow the best practices in this section to achieve optimal query performance and efficiency.

Avoid using the `SELECT * FROM` statement

We recommend that you avoid using the `SELECT * FROM` statement. Instead, always list out columns for analysis. This reduces query execution time and scans costs for Amazon Redshift Spectrum queries.

Example of what to avoid

```
select *  
from sales;
```

Best practice example

```
select sales_date, sales_amt  
from sales;
```

Identify query issues

We recommend that you check the [STL_ALERT_EVENT_LOG](#) view to identify and correct possible issues with your query.

Get summary information on your query

We recommend that you use the [SVL_QUERY_SUMMARY](#) and [SVL_QUERY_REPORT](#) views to get summary information on your queries. You can use this information to optimize your queries.

Avoid cross-joins

We recommend that you avoid using cross-joins unless absolutely necessary. Without a join condition, cross-joins result in the Cartesian product of two tables. Cross-joins are typically run as nested-loop joins (the slowest of the possible join types).

Example of what to avoid

```
select c.c_name,  
       n.n_name  
from tpch.customer c,  
     tpch.nation n;
```

Best practice example

```
select c.c_name,  
       n.n_name  
from tpch.customer c,  
     join tpch.nation n  
     on n.n_nationkey = c.c_nationkey;
```

Avoid functions in query predicates

We recommend that you avoid using functions in query predicates. Using functions in query predicates can negatively impact performance because functions typically add extra processing overhead to each row and slow down the overall execution of the query.

Example of what to avoid

```
select sum(o_totalprice)  
from tpch.orders  
where datepart(year, o_orderdate) = 1992;
```

Best practice example

```
select sum(o_totalprice)  
from tpch.orders  
where o_orderdate between '1992-01-01' and '1992-12-31';
```

Avoid unnecessary cast conversions

We recommend that you avoid using unnecessary cast conversion on the queries because casting data types takes time and resources and slows down query execution.

Example of what to avoid

```
select sum(o_totalprice)
from tpch.orders
where o_ordertime::date = '1992-01-01';
```

Best practice example

```
select sum(o_totalprice)
from tpch.orders
where o_ordertime between '1992-01-01 00:00:00' and '1992-12-31 23:59:59';
```

Use CASE expressions for complex aggregations

We recommend that you use a [CASE expression](#) to perform complex aggregations instead of selecting from the same table multiple times.

Example of what to avoid

```
select sum(sales_amt) as us_sales
from sales
where country = 'US';

select sum(sales_amt) as ca_sales
from sales
where country = 'CA';
```

Best practice example

```
select sum(case when country = 'US' then sales_amt end) as us_sales,
       sum(case when country = 'CA' then sales_amt end) as ca_sales
from sales;
```

Use subqueries

We recommend that you use subqueries in cases where one table in the query is used only for predicate conditions and the subquery returns a small number of rows (less than about 200).

Example of what to avoid

If a subquery returns less than 200 rows:

```
select sum(order_amt) as total_sales
from sales
where region_key IN
    (select region_key
     from regions
     where state = 'CA');
```

Best practice example

If a subquery returns greater than or equal to 200 rows:

```
select sum(o.order_amt) as total_sales
from sales o
join regions r
  on r.region_key = o.region_key
  and r.state = 'CA';
```

Use predicates

We recommend that you use predicates to restrict the dataset as much as possible. Predicates are used in SQL to filter and restrict the data that's returned in a query. By specifying conditions in a predicate, you can specify which rows must be included in the query results based on specified conditions. This allows you to retrieve only the data that you're interested in and improves the efficiency and accuracy of your queries. For more information, see [Conditions](#) in the Amazon Redshift documentation.

Add predicates to filter tables with joins

We recommend that you add predicates to filter tables that participate in joins, even if the predicates apply the same filters. Using predicates to filter tables with joins in SQL can improve query performance by reducing the amount of data that must be processed and reducing the size of the intermediate result set. By specifying the conditions for the join operation in the WHERE clause, the query execution engine can eliminate rows that don't match the conditions before they are joined. This results in a smaller result set and faster query execution.

Example of what to avoid

```
select p.product_name, sum(o.order_amt)
```

```
from sales o
join product p
  on r.product_key = o.product_key
where o.order_date > '2022-01-01';
```

Best practice example

```
select p.product_name, sum(o.order_amt)
from sales o
join product p
  on p.product_key = o.product_key
  and p.added_date > '2022-01-01'
where o.order_date > '2022-01-01';
```

Use the least expensive operators for predicates

In the predicate, use the least expensive operators that you can. [Comparison condition](#) operators are preferable to [LIKE](#) operators. LIKE operators are still preferable to [SIMILAR TO](#) or [POSIX](#) operators.

Use sort keys in GROUP BY clauses

Use sort keys in the GROUP BY clause so that the query planner can use more efficient aggregation. A query might qualify for one-phase aggregation when its GROUP BY list contains only sort key columns, one of which is also the distribution key. The sort key columns in the GROUP BY list must include the first sort key, followed by other sort keys that you want to use in sort key order.

Take advantage of materialized views

If possible, rewrite the query by replacing the complex code with a materialized view, which will significantly improve the performance of the query. For more information, see [Creating materialized views in Amazon Redshift](#) in the Amazon Redshift documentation.

Be careful with columns in GROUP BY and ORDER BY clauses

If you use both GROUP BY and ORDER BY clauses, make sure that you put the columns in the same order in both GROUP BY and ORDER BY clauses. GROUP BY implicitly requires the data to be sorted. If your ORDER BY clause is different, then the data must be sorted twice.

Example of what to avoid

```
select a, b, c, sum(d)
from a_table
group by b, c, a
order by a, b, c
```

Best practice example

```
select a, b, c, sum(d)
from a_table
group by a, b, c
order by a, b, c
```

Best practices for using Amazon Redshift Spectrum

This section provides an overview of best practices for using [Amazon Redshift Spectrum](#). We recommend that you follow these best practices to achieve optimal performance when you use Redshift Spectrum:

- Consider that file types have a significant influence on Redshift Spectrum query performance. To improve performance, use columnar encoded files such as ORC or Parquet, and use CSV format only for very small dimension tables.
- Use prefix-based partitioning to take advantage of partition pruning. This means using filters that are keyed to the partitions in your data lake.
- Redshift Spectrum scales automatically to process large requests, so do as much as possible in Redshift Spectrum (for example, predicate pushdown).
- Pay attention to partition files on frequently filtered columns. If data is partitioned by one or more filtered columns, Redshift Spectrum can take advantage of partition pruning and skip scanning unneeded partitions and files. A common practice is to partition the data based on time.
- You can check the effectiveness of your partitions and efficiency of your Redshift Spectrum query by using the following query.

```
Select query,  
       segment,  
       max(assigned_partitions) as total_partitions,  
       max(qualified_partitions) as qualified_partitions  
From svl_s3partition  
Where query=pg_last_query_id()  
Group by 1,2;
```

The preceding query shows the following:

- **total_partitions** – The number of partitions recognized by the AWS Glue Data Catalog
- **qualified_partitions** – The number of prefixes on Amazon Simple Storage Service (Amazon S3) that are accessed for the Redshift Spectrum query
- You can also check the SVL_S3QUERY_SUMMARY system table to learn about the effectiveness of your partitions and efficiency of your Redshift Spectrum query. To do so, use the following statement.

```
Select *
From svl_s3query_summary
Where query=pg_last_query_id();
```

The preceding query returns even more information, including `is_partitioned`, `s3_scanned_rows/bytes`, and `s3_returned_rows/bytes` values in addition to files showing the efficiency of partition pruning.

Predicate pushdown in Redshift Spectrum

Using predicate pushdown avoids consuming resources in the Amazon Redshift cluster. You can push many SQL operations down to the Redshift Spectrum layer. We recommend taking advantage of this wherever possible.

Keep in mind the following:

- You can evaluate some types of SQL operations completely within the Redshift Spectrum layer, including the following:
 - GROUP BY clauses
 - Comparison and pattern-matching conditions (for example, LIKE)
 - Aggregate functions (for example, COUNT, SUM, AVG, MIN, and MAX)
 - `regex_replace`, `to_upper`, `date_trunc`, and other functions
- You can't push some operations to the Redshift Spectrum layer, including DISTINCT and ORDER BY. Perform ORDER BY only at the top level of the query if possible, since sorting is done in the leader node.
- Examine your query EXPLAIN plan to verify if predicate pushdown is effective. To find Redshift Spectrum portions in an EXPLAIN command, look for these steps:
 - S3 Seq Scan
 - S3 HashAggregate
 - S3 Query Scan
 - Seq Scan PartitionInfo
 - Partition Loop
- Use the fewest number of columns in your query. Redshift Spectrum can eliminate columns for scanning if data is in Parquet or ORC format.

- Make extensive use of partitions for parallel processing and partition elimination, and keep file sizes to at least 64 MB if possible.
- Set the `TABLE PROPERTIES 'numRows' = 'nnn'` if you use `CREATE EXTERNAL TABLE` or `ALTER TABLE`. Amazon Redshift doesn't analyze external tables to generate table statistics that the query optimizer uses to generate a query plan. If the statistics aren't set, then Amazon Redshift assumes that external tables are the larger tables and local tables are the smaller tables.

Query tuning tips for Redshift Spectrum

We recommend that you keep the following in mind when you tune your queries:

- The number of Redshift Spectrum nodes that your Amazon Redshift cluster can engage for a query is tied to the number of slices in your cluster.
- Resizing your cluster up can benefit your cluster's local compute profiles, storage profiles, and the query capabilities of the Amazon S3 data lake query.
- The Amazon Redshift query planner pushes predicates and aggregations to the Redshift Spectrum query layer whenever possible.
- When large amounts of data are returned from Amazon S3, the processing is limited by your cluster's resources.
- Because Redshift Spectrum scales automatically to process large requests, your overall performance improves whenever you can push processing to the Redshift Spectrum layer.

Resources

- [Amazon Redshift best practices](#) (Amazon Redshift documentation)
- [Amazon Redshift best practices for designing queries](#) (Amazon Redshift documentation)
- [Tuning query performance](#) (Amazon Redshift documentation)
- [Query plan](#) (Amazon Redshift documentation)
- [Improving Amazon Redshift Spectrum query performance](#) (Amazon Redshift documentation)
- [Understanding the query lifecycle in Amazon Redshift](#) (AWS Prescriptive Guidance)

Document history

The following table describes significant changes to this guide. If you want to be notified about future updates, you can subscribe to an [RSS feed](#).

Change	Description	Date
Removed AQUA	We removed information about Advanced Query Accelerator (AQUA).	June 14, 2024
Initial publication	—	February 3, 2023

AWS Prescriptive Guidance glossary

The following are commonly used terms in strategies, guides, and patterns provided by AWS Prescriptive Guidance. To suggest entries, please use the **Provide feedback** link at the end of the glossary.

Numbers

7 Rs

Seven common migration strategies for moving applications to the cloud. These strategies build upon the 5 Rs that Gartner identified in 2011 and consist of the following:

- Refactor/re-architect – Move an application and modify its architecture by taking full advantage of cloud-native features to improve agility, performance, and scalability. This typically involves porting the operating system and database. Example: Migrate your on-premises Oracle database to the Amazon Aurora PostgreSQL-Compatible Edition.
- Replatform (lift and reshape) – Move an application to the cloud, and introduce some level of optimization to take advantage of cloud capabilities. Example: Migrate your on-premises Oracle database to Amazon Relational Database Service (Amazon RDS) for Oracle in the AWS Cloud.
- Repurchase (drop and shop) – Switch to a different product, typically by moving from a traditional license to a SaaS model. Example: Migrate your customer relationship management (CRM) system to Salesforce.com.
- Rehost (lift and shift) – Move an application to the cloud without making any changes to take advantage of cloud capabilities. Example: Migrate your on-premises Oracle database to Oracle on an EC2 instance in the AWS Cloud.
- Relocate (hypervisor-level lift and shift) – Move infrastructure to the cloud without purchasing new hardware, rewriting applications, or modifying your existing operations. You migrate servers from an on-premises platform to a cloud service for the same platform. Example: Migrate a Microsoft Hyper-V application to AWS.
- Retain (revisit) – Keep applications in your source environment. These might include applications that require major refactoring, and you want to postpone that work until a later time, and legacy applications that you want to retain, because there's no business justification for migrating them.

- Retire – Decommission or remove applications that are no longer needed in your source environment.

A

A2A (Agent-to-Agent)

A stateful protocol for agent-to-agent collaboration supporting task delegation and state transfer.

ABAC

See [attribute-based access control](#).

abstracted services

See [managed services](#).

ACID

See [atomicity, consistency, isolation, durability](#).

active-active migration

A database migration method in which the source and target databases are kept in sync (by using a bidirectional replication tool or dual write operations), and both databases handle transactions from connecting applications during migration. This method supports migration in small, controlled batches instead of requiring a one-time cutover. It's more flexible but requires more work than [active-passive migration](#).

active-passive migration

A database migration method in which the source and target databases are kept in sync, but only the source database handles transactions from connecting applications while data is replicated to the target database. The target database doesn't accept any transactions during migration.

Agent

An AI system that can autonomously reason, plan, and take actions using tools to achieve goals.

Agent Ops

Operational practices for building, testing, deploying, and running AI agents in production at scale.

aggregate function

A SQL function that operates on a group of rows and calculates a single return value for the group. Examples of aggregate functions include SUM and MAX.

AI

See [artificial intelligence](#).

AIOps

See [artificial intelligence operations](#).

anonymization

The process of permanently deleting personal information in a dataset. Anonymization can help protect personal privacy. Anonymized data is no longer considered to be personal data.

anti-pattern

A frequently used solution for a recurring issue where the solution is counter-productive, ineffective, or less effective than an alternative.

application control

A security approach that allows the use of only approved applications in order to help protect a system from malware.

application portfolio

A collection of detailed information about each application used by an organization, including the cost to build and maintain the application, and its business value. This information is key to [the portfolio discovery and analysis process](#) and helps identify and prioritize the applications to be migrated, modernized, and optimized.

artificial intelligence (AI)

The field of computer science that is dedicated to using computing technologies to perform cognitive functions that are typically associated with humans, such as learning, solving problems, and recognizing patterns. For more information, see [What is Artificial Intelligence?](#)

artificial intelligence operations (AIOps)

The process of using machine learning techniques to solve operational problems, reduce operational incidents and human intervention, and increase service quality. For more information about how AIOps is used in the AWS migration strategy, see the [operations integration guide](#).

asymmetric encryption

An encryption algorithm that uses a pair of keys, a public key for encryption and a private key for decryption. You can share the public key because it isn't used for decryption, but access to the private key should be highly restricted.

atomicity, consistency, isolation, durability (ACID)

A set of software properties that guarantee the data validity and operational reliability of a database, even in the case of errors, power failures, or other problems.

attribute-based access control (ABAC)

The practice of creating fine-grained permissions based on user attributes, such as department, job role, and team name. For more information, see [ABAC for AWS](#) in the AWS Identity and Access Management (IAM) documentation.

authoritative data source

A location where you store the primary version of data, which is considered to be the most reliable source of information. You can copy data from the authoritative data source to other locations for the purposes of processing or modifying the data, such as anonymizing, redacting, or pseudonymizing it.

Availability Zone

A distinct location within an AWS Region that is insulated from failures in other Availability Zones and provides inexpensive, low-latency network connectivity to other Availability Zones in the same Region.

AWS Cloud Adoption Framework (AWS CAF)

A framework of guidelines and best practices from AWS to help organizations develop an efficient and effective plan to move successfully to the cloud. AWS CAF organizes guidance into six focus areas called perspectives: business, people, governance, platform, security, and operations. The business, people, and governance perspectives focus on business skills and processes; the platform, security, and operations perspectives focus on technical skills and processes. For example, the people perspective targets stakeholders who handle human resources (HR), staffing functions, and people management. For this perspective, AWS CAF provides guidance for people development, training, and communications to help ready the organization for successful cloud adoption. For more information, see the [AWS CAF website](#) and the [AWS CAF whitepaper](#).

AWS Workload Qualification Framework (AWS WQF)

A tool that evaluates database migration workloads, recommends migration strategies, and provides work estimates. AWS WQF is included with AWS Schema Conversion Tool (AWS SCT). It analyzes database schemas and code objects, application code, dependencies, and performance characteristics, and provides assessment reports.

B

bad bot

A [bot](#) that is intended to disrupt or cause harm to individuals or organizations.

BCP

See [business continuity planning](#).

behavior graph

A unified, interactive view of resource behavior and interactions over time. You can use a behavior graph with Amazon Detective to examine failed logon attempts, suspicious API calls, and similar actions. For more information, see [Data in a behavior graph](#) in the Detective documentation.

big-endian system

A system that stores the most significant byte first. See also [endianness](#).

binary classification

A process that predicts a binary outcome (one of two possible classes). For example, your ML model might need to predict problems such as "Is this email spam or not spam?" or "Is this product a book or a car?"

bloom filter

A probabilistic, memory-efficient data structure that is used to test whether an element is a member of a set.

blue/green deployment

A deployment strategy where you create two separate but identical environments. You run the current application version in one environment (blue) and the new application version in the other environment (green). This strategy helps you quickly roll back with minimal impact.

bot

A software application that runs automated tasks over the internet and simulates human activity or interaction. Some bots are useful or beneficial, such as web crawlers that index information on the internet. Some other bots, known as *bad bots*, are intended to disrupt or cause harm to individuals or organizations.

botnet

Networks of [bots](#) that are infected by [malware](#) and are under the control of a single party, known as a *bot herder* or *bot operator*. Botnets are the best-known mechanism to scale bots and their impact.

branch

A contained area of a code repository. The first branch created in a repository is the *main branch*. You can create a new branch from an existing branch, and you can then develop features or fix bugs in the new branch. A branch you create to build a feature is commonly referred to as a *feature branch*. When the feature is ready for release, you merge the feature branch back into the main branch. For more information, see [About branches](#) (GitHub documentation).

break-glass access

In exceptional circumstances and through an approved process, a quick means for a user to gain access to an AWS account that they don't typically have permissions to access. For more information, see the [Implement break-glass procedures](#) indicator in the AWS Well-Architected guidance.

brownfield strategy

The existing infrastructure in your environment. When adopting a brownfield strategy for a system architecture, you design the architecture around the constraints of the current systems and infrastructure. If you are expanding the existing infrastructure, you might blend brownfield and [greenfield](#) strategies.

buffer cache

The memory area where the most frequently accessed data is stored.

business capability

What a business does to generate value (for example, sales, customer service, or marketing). Microservices architectures and development decisions can be driven by business capabilities.

For more information, see the [Organized around business capabilities](#) section of the [Running containerized microservices on AWS](#) whitepaper.

business continuity planning (BCP)

A plan that addresses the potential impact of a disruptive event, such as a large-scale migration, on operations and enables a business to resume operations quickly.

C

CAF

See [AWS Cloud Adoption Framework](#).

canary deployment

The slow and incremental release of a version to end users. When you are confident, you deploy the new version and replace the current version in its entirety.

CCoE

See [Cloud Center of Excellence](#).

CDC

See [change data capture](#).

change data capture (CDC)

The process of tracking changes to a data source, such as a database table, and recording metadata about the change. You can use CDC for various purposes, such as auditing or replicating changes in a target system to maintain synchronization.

chaos engineering

Intentionally introducing failures or disruptive events to test a system's resilience. You can use [AWS Fault Injection Service \(AWS FIS\)](#) to perform experiments that stress your AWS workloads and evaluate their response.

CI/CD

See [continuous integration and continuous delivery](#).

classification

A categorization process that helps generate predictions. ML models for classification problems predict a discrete value. Discrete values are always distinct from one another. For example, a model might need to evaluate whether or not there is a car in an image.

Citizen Developer

A business user who creates AI applications using no-code/low-code platforms without specialized technical skills.

client-side encryption

Encryption of data locally, before the target AWS service receives it.

Cloud Center of Excellence (CCoE)

A multi-disciplinary team that drives cloud adoption efforts across an organization, including developing cloud best practices, mobilizing resources, establishing migration timelines, and leading the organization through large-scale transformations. For more information, see the [CCoE posts](#) on the AWS Cloud Enterprise Strategy Blog.

cloud computing

The cloud technology that is typically used for remote data storage and IoT device management. Cloud computing is commonly connected to [edge computing](#) technology.

cloud operating model

In an IT organization, the operating model that is used to build, mature, and optimize one or more cloud environments. For more information, see [Building your Cloud Operating Model](#).

cloud stages of adoption

The four phases that organizations typically go through when they migrate to the AWS Cloud:

- Project – Running a few cloud-related projects for proof of concept and learning purposes
- Foundation – Making foundational investments to scale your cloud adoption (e.g., creating a landing zone, defining a CCoE, establishing an operations model)
- Migration – Migrating individual applications
- Re-invention – Optimizing products and services, and innovating in the cloud

These stages were defined by Stephen Orban in the blog post [The Journey Toward Cloud-First & the Stages of Adoption](#) on the AWS Cloud Enterprise Strategy blog. For information about how they relate to the AWS migration strategy, see the [migration readiness guide](#).

CMDB

See [configuration management database](#).

code repository

A location where source code and other assets, such as documentation, samples, and scripts, are stored and updated through version control processes. Common cloud repositories include GitHub or Bitbucket Cloud. Each version of the code is called a *branch*. In a microservice structure, each repository is devoted to a single piece of functionality. A single CI/CD pipeline can use multiple repositories.

cold cache

A buffer cache that is empty, not well populated, or contains stale or irrelevant data. This affects performance because the database instance must read from the main memory or disk, which is slower than reading from the buffer cache.

cold data

Data that is rarely accessed and is typically historical. When querying this kind of data, slow queries are typically acceptable. Moving this data to lower-performing and less expensive storage tiers or classes can reduce costs.

computer vision (CV)

A field of [AI](#) that uses machine learning to analyze and extract information from visual formats such as digital images and videos. For example, Amazon SageMaker AI provides image processing algorithms for CV.

configuration drift

For a workload, a configuration change from the expected state. It might cause the workload to become noncompliant, and it's typically gradual and unintentional.

configuration management database (CMDB)

A repository that stores and manages information about a database and its IT environment, including both hardware and software components and their configurations. You typically use data from a CMDB in the portfolio discovery and analysis stage of migration.

conformance pack

A collection of AWS Config rules and remediation actions that you can assemble to customize your compliance and security checks. You can deploy a conformance pack as a single entity in

an AWS account and Region, or across an organization, by using a YAML template. For more information, see [Conformance packs](#) in the AWS Config documentation.

continuous integration and continuous delivery (CI/CD)

The process of automating the source, build, test, staging, and production stages of the software release process. CI/CD is commonly described as a pipeline. CI/CD can help you automate processes, improve productivity, improve code quality, and deliver faster. For more information, see [Benefits of continuous delivery](#). CD can also stand for *continuous deployment*. For more information, see [Continuous Delivery vs. Continuous Deployment](#).

CV

See [computer vision](#).

D

data at rest

Data that is stationary in your network, such as data that is in storage.

data classification

A process for identifying and categorizing the data in your network based on its criticality and sensitivity. It is a critical component of any cybersecurity risk management strategy because it helps you determine the appropriate protection and retention controls for the data. Data classification is a component of the security pillar in the AWS Well-Architected Framework. For more information, see [Data classification](#).

data drift

A meaningful variation between the production data and the data that was used to train an ML model, or a meaningful change in the input data over time. Data drift can reduce the overall quality, accuracy, and fairness in ML model predictions.

data in transit

Data that is actively moving through your network, such as between network resources.

data mesh

An architectural framework that provides distributed, decentralized data ownership with centralized management and governance.

data minimization

The principle of collecting and processing only the data that is strictly necessary. Practicing data minimization in the AWS Cloud can reduce privacy risks, costs, and your analytics carbon footprint.

data perimeter

A set of preventive guardrails in your AWS environment that help make sure that only trusted identities are accessing trusted resources from expected networks. For more information, see [Building a data perimeter on AWS](#).

data preprocessing

To transform raw data into a format that is easily parsed by your ML model. Preprocessing data can mean removing certain columns or rows and addressing missing, inconsistent, or duplicate values.

data provenance

The process of tracking the origin and history of data throughout its lifecycle, such as how the data was generated, transmitted, and stored.

data subject

An individual whose data is being collected and processed.

data warehouse

A data management system that supports business intelligence, such as analytics. Data warehouses commonly contain large amounts of historical data, and they are typically used for queries and analysis.

database definition language (DDL)

Statements or commands for creating or modifying the structure of tables and objects in a database.

database manipulation language (DML)

Statements or commands for modifying (inserting, updating, and deleting) information in a database.

DDL

See [database definition language](#).

deep ensemble

To combine multiple deep learning models for prediction. You can use deep ensembles to obtain a more accurate prediction or for estimating uncertainty in predictions.

deep learning

An ML subfield that uses multiple layers of artificial neural networks to identify mapping between input data and target variables of interest.

defense-in-depth

An information security approach in which a series of security mechanisms and controls are thoughtfully layered throughout a computer network to protect the confidentiality, integrity, and availability of the network and the data within. When you adopt this strategy on AWS, you add multiple controls at different layers of the AWS Organizations structure to help secure resources. For example, a defense-in-depth approach might combine multi-factor authentication, network segmentation, and encryption.

delegated administrator

In AWS Organizations, a compatible service can register an AWS member account to administer the organization's accounts and manage permissions for that service. This account is called the *delegated administrator* for that service. For more information and a list of compatible services, see [Services that work with AWS Organizations](#) in the AWS Organizations documentation.

deployment

The process of making an application, new features, or code fixes available in the target environment. Deployment involves implementing changes in a code base and then building and running that code base in the application's environments.

development environment

See [environment](#).

detective control

A security control that is designed to detect, log, and alert after an event has occurred. These controls are a second line of defense, alerting you to security events that bypassed the preventative controls in place. For more information, see [Detective controls](#) in *Implementing security controls on AWS*.

development value stream mapping (DVSM)

A process used to identify and prioritize constraints that adversely affect speed and quality in a software development lifecycle. DVSM extends the value stream mapping process originally designed for lean manufacturing practices. It focuses on the steps and teams required to create and move value through the software development process.

digital twin

A virtual representation of a real-world system, such as a building, factory, industrial equipment, or production line. Digital twins support predictive maintenance, remote monitoring, and production optimization.

dimension table

In a [star schema](#), a smaller table that contains data attributes about quantitative data in a fact table. Dimension table attributes are typically text fields or discrete numbers that behave like text. These attributes are commonly used for query constraining, filtering, and result set labeling.

disaster

An event that prevents a workload or system from fulfilling its business objectives in its primary deployed location. These events can be natural disasters, technical failures, or the result of human actions, such as unintentional misconfiguration or a malware attack.

disaster recovery (DR)

The strategy and process you use to minimize downtime and data loss caused by a [disaster](#). For more information, see [Disaster Recovery of Workloads on AWS: Recovery in the Cloud](#) in the AWS Well-Architected Framework.

DML

See [database manipulation language](#).

domain-driven design

An approach to developing a complex software system by connecting its components to evolving domains, or core business goals, that each component serves. This concept was introduced by Eric Evans in his book, *Domain-Driven Design: Tackling Complexity in the Heart of Software* (Boston: Addison-Wesley Professional, 2003). For information about how you can use domain-driven design with the strangler fig pattern, see [Modernizing legacy Microsoft ASP.NET \(ASMX\) web services incrementally by using containers and Amazon API Gateway](#).

DR

See [disaster recovery](#).

drift detection

Tracking deviations from a baselined configuration. For example, you can use AWS CloudFormation to [detect drift in system resources](#), or you can use AWS Control Tower to [detect changes in your landing zone](#) that might affect compliance with governance requirements.

DVSM

See [development value stream mapping](#).

E

EDA

See [exploratory data analysis](#).

EDI

See [electronic data interchange](#).

edge computing

The technology that increases the computing power for smart devices at the edges of an IoT network. When compared with [cloud computing](#), edge computing can reduce communication latency and improve response time.

electronic data interchange (EDI)

The automated exchange of business documents between organizations. For more information, see [What is Electronic Data Interchange](#).

encryption

A computing process that transforms plaintext data, which is human-readable, into ciphertext.

encryption key

A cryptographic string of randomized bits that is generated by an encryption algorithm. Keys can vary in length, and each key is designed to be unpredictable and unique.

endianness

The order in which bytes are stored in computer memory. Big-endian systems store the most significant byte first. Little-endian systems store the least significant byte first.

endpoint

See [service endpoint](#).

endpoint service

A service that you can host in a virtual private cloud (VPC) to share with other users. You can create an endpoint service with AWS PrivateLink and grant permissions to other AWS accounts or to AWS Identity and Access Management (IAM) principals. These accounts or principals can connect to your endpoint service privately by creating interface VPC endpoints. For more information, see [Create an endpoint service](#) in the Amazon Virtual Private Cloud (Amazon VPC) documentation.

enterprise resource planning (ERP)

A system that automates and manages key business processes (such as accounting, [MES](#), and project management) for an enterprise.

envelope encryption

The process of encrypting an encryption key with another encryption key. For more information, see [Envelope encryption](#) in the AWS Key Management Service (AWS KMS) documentation.

environment

An instance of a running application. The following are common types of environments in cloud computing:

- development environment – An instance of a running application that is available only to the core team responsible for maintaining the application. Development environments are used to test changes before promoting them to upper environments. This type of environment is sometimes referred to as a *test environment*.
- lower environments – All development environments for an application, such as those used for initial builds and tests.
- production environment – An instance of a running application that end users can access. In a CI/CD pipeline, the production environment is the last deployment environment.

- upper environments – All environments that can be accessed by users other than the core development team. This can include a production environment, preproduction environments, and environments for user acceptance testing.

epic

In agile methodologies, functional categories that help organize and prioritize your work. Epics provide a high-level description of requirements and implementation tasks. For example, AWS CAF security epics include identity and access management, detective controls, infrastructure security, data protection, and incident response. For more information about epics in the AWS migration strategy, see the [program implementation guide](#).

ERP

See [enterprise resource planning](#).

exploratory data analysis (EDA)

The process of analyzing a dataset to understand its main characteristics. You collect or aggregate data and then perform initial investigations to find patterns, detect anomalies, and check assumptions. EDA is performed by calculating summary statistics and creating data visualizations.

F

fact table

The central table in a [star schema](#). It stores quantitative data about business operations. Typically, a fact table contains two types of columns: those that contain measures and those that contain a foreign key to a dimension table.

fail fast

A philosophy that uses frequent and incremental testing to reduce the development lifecycle. It is a critical part of an agile approach.

fault isolation boundary

In the AWS Cloud, a boundary such as an Availability Zone, AWS Region, control plane, or data plane that limits the effect of a failure and helps improve the resilience of workloads. For more information, see [AWS Fault Isolation Boundaries](#).

feature branch

See [branch](#).

features

The input data that you use to make a prediction. For example, in a manufacturing context, features could be images that are periodically captured from the manufacturing line.

feature importance

How significant a feature is for a model's predictions. This is usually expressed as a numerical score that can be calculated through various techniques, such as Shapley Additive Explanations (SHAP) and integrated gradients. For more information, see [Machine learning model interpretability with AWS](#).

feature transformation

To optimize data for the ML process, including enriching data with additional sources, scaling values, or extracting multiple sets of information from a single data field. This enables the ML model to benefit from the data. For example, if you break down the "2021-05-27 00:15:37" date into "2021", "May", "Thu", and "15", you can help the learning algorithm learn nuanced patterns associated with different data components.

few-shot prompting

Providing an [LLM](#) with a small number of examples that demonstrate the task and desired output before asking it to perform a similar task. This technique is an application of in-context learning, where models learn from examples (*shots*) that are embedded in prompts. Few-shot prompting can be effective for tasks that require specific formatting, reasoning, or domain knowledge. See also [zero-shot prompting](#).

FGAC

See [fine-grained access control](#).

fine-grained access control (FGAC)

The use of multiple conditions to allow or deny an access request.

flash-cut migration

A database migration method that uses continuous data replication through [change data capture](#) to migrate data in the shortest time possible, instead of using a phased approach. The objective is to keep downtime to a minimum.

FM

See [foundation model](#).

foundation model (FM)

A large deep-learning neural network that has been training on massive datasets of generalized and unlabeled data. FMs are capable of performing a wide variety of general tasks, such as understanding language, generating text and images, and conversing in natural language. For more information, see [What are Foundation Models](#).

FM gateway

A centralized intermediary that controls and normalizes access to [foundation models](#). Also known as an *LLM gateway*.

G

generative AI

A subset of [AI](#) models that have been trained on large amounts of data and that can use a simple text prompt to create new content and artifacts, such as images, videos, text, and audio. For more information, see [What is Generative AI](#).

geo blocking

See [geographic restrictions](#).

geographic restrictions (geo blocking)

In Amazon CloudFront, an option to prevent users in specific countries from accessing content distributions. You can use an allow list or block list to specify approved and banned countries. For more information, see [Restricting the geographic distribution of your content](#) in the CloudFront documentation.

Gitflow workflow

An approach in which lower and upper environments use different branches in a source code repository. The Gitflow workflow is considered legacy, and the [trunk-based workflow](#) is the modern, preferred approach.

golden image

A snapshot of a system or software that is used as a template to deploy new instances of that system or software. For example, in manufacturing, a golden image can be used to provision

software on multiple devices and helps improve speed, scalability, and productivity in device manufacturing operations.

greenfield strategy

The absence of existing infrastructure in a new environment. When adopting a greenfield strategy for a system architecture, you can select all new technologies without the restriction of compatibility with existing infrastructure, also known as [brownfield](#). If you are expanding the existing infrastructure, you might blend brownfield and greenfield strategies.

guardrail

A high-level rule that helps govern resources, policies, and compliance across organizational units (OUs). *Preventive guardrails* enforce policies to ensure alignment to compliance standards. They are implemented by using service control policies and IAM permissions boundaries. *Detective guardrails* detect policy violations and compliance issues, and generate alerts for remediation. They are implemented by using AWS Config, AWS Security Hub CSPM, Amazon GuardDuty, AWS Trusted Advisor, Amazon Inspector, and custom AWS Lambda checks.

guardrails (AI)

Safety mechanisms that filter, validate, and constrain [agent](#) inputs and outputs to help ensure responsible and safe AI behavior.

H

HA

See [high availability](#).

heterogeneous database migration

Migrating your source database to a target database that uses a different database engine (for example, Oracle to Amazon Aurora). Heterogeneous migration is typically part of a re-architecting effort, and converting the schema can be a complex task. [AWS provides AWS SCT](#) that helps with schema conversions.

high availability (HA)

The ability of a workload to operate continuously, without intervention, in the event of challenges or disasters. HA systems are designed to automatically fail over, consistently deliver

high-quality performance, and handle different loads and failures with minimal performance impact.

historian modernization

An approach used to modernize and upgrade operational technology (OT) systems to better serve the needs of the manufacturing industry. A *historian* is a type of database that is used to collect and store data from various sources in a factory.

holdout data

A portion of historical, labeled data that is withheld from a dataset that is used to train a [machine learning](#) model. You can use holdout data to evaluate the model performance by comparing the model predictions against the holdout data.

human-in-the-loop (HitL)

A workflow pattern where [agent](#) execution pauses for human review and approval at critical decision points.

homogeneous database migration

Migrating your source database to a target database that shares the same database engine (for example, Microsoft SQL Server to Amazon RDS for SQL Server). Homogeneous migration is typically part of a rehosting or replatforming effort. You can use native database utilities to migrate the schema.

hot data

Data that is frequently accessed, such as real-time data or recent translational data. This data typically requires a high-performance storage tier or class to provide fast query responses.

hotfix

An urgent fix for a critical issue in a production environment. Due to its urgency, a hotfix is usually made outside of the typical DevOps release workflow.

hypercare period

Immediately following cutover, the period of time when a migration team manages and monitors the migrated applications in the cloud in order to address any issues. Typically, this period is 1–4 days in length. At the end of the hypercare period, the migration team typically transfers responsibility for the applications to the cloud operations team.

I

laC

See [infrastructure as code](#).

identity-based policy

A policy attached to one or more IAM principals that defines their permissions within the AWS Cloud environment.

idle application

An application that has an average CPU and memory usage between 5 and 20 percent over a period of 90 days. In a migration project, it is common to retire these applications or retain them on premises.

IIoT

See [industrial Internet of Things](#).

immutable infrastructure

A model that deploys new infrastructure for production workloads instead of updating, patching, or modifying the existing infrastructure. Immutable infrastructures are inherently more consistent, reliable, and predictable than [mutable infrastructure](#). For more information, see the [Deploy using immutable infrastructure](#) best practice in the AWS Well-Architected Framework.

inbound (ingress) VPC

In an AWS multi-account architecture, a VPC that accepts, inspects, and routes network connections from outside an application. The [AWS Security Reference Architecture](#) recommends setting up your Network account with inbound, outbound, and inspection VPCs to protect the two-way interface between your application and the broader internet.

incremental migration

A cutover strategy in which you migrate your application in small parts instead of performing a single, full cutover. For example, you might move only a few microservices or users to the new system initially. After you verify that everything is working properly, you can incrementally move additional microservices or users until you can decommission your legacy system. This strategy reduces the risks associated with large migrations.

Industry 4.0

A term that was introduced by [Klaus Schwab](#) in 2016 to refer to the modernization of manufacturing processes through advances in connectivity, real-time data, automation, analytics, and AI/ML.

infrastructure

All of the resources and assets contained within an application's environment.

infrastructure as code (IaC)

The process of provisioning and managing an application's infrastructure through a set of configuration files. IaC is designed to help you centralize infrastructure management, standardize resources, and scale quickly so that new environments are repeatable, reliable, and consistent.

industrial Internet of Things (IIoT)

The use of internet-connected sensors and devices in the industrial sectors, such as manufacturing, energy, automotive, healthcare, life sciences, and agriculture. For more information, see [Building an industrial Internet of Things \(IIoT\) digital transformation strategy](#).

inspection VPC

In an AWS multi-account architecture, a centralized VPC that manages inspections of network traffic between VPCs (in the same or different AWS Regions), the internet, and on-premises networks. The [AWS Security Reference Architecture](#) recommends setting up your Network account with inbound, outbound, and inspection VPCs to protect the two-way interface between your application and the broader internet.

Internet of Things (IoT)

The network of connected physical objects with embedded sensors or processors that communicate with other devices and systems through the internet or over a local communication network. For more information, see [What is IoT?](#)

interpretability

A characteristic of a machine learning model that describes the degree to which a human can understand how the model's predictions depend on its inputs. For more information, see [Machine learning model interpretability with AWS](#).

IoT

See [Internet of Things](#).

IT information library (ITIL)

A set of best practices for delivering IT services and aligning these services with business requirements. ITIL provides the foundation for ITSM.

IT service management (ITSM)

Activities associated with designing, implementing, managing, and supporting IT services for an organization. For information about integrating cloud operations with ITSM tools, see the [operations integration guide](#).

ITIL

See [IT information library](#).

ITSM

See [IT service management](#).

L

label-based access control (LBAC)

An implementation of mandatory access control (MAC) where the users and the data itself are each explicitly assigned a security label value. The intersection between the user security label and data security label determines which rows and columns can be seen by the user.

landing zone

A landing zone is a well-architected, multi-account AWS environment that is scalable and secure. This is a starting point from which your organizations can quickly launch and deploy workloads and applications with confidence in their security and infrastructure environment. For more information about landing zones, see [Setting up a secure and scalable multi-account AWS environment](#).

large language model (LLM)

A deep learning [AI](#) model that is pretrained on a vast amount of data. An LLM can perform multiple tasks, such as answering questions, summarizing documents, translating text into other languages, and completing sentences. For more information, see [What are LLMs](#).

large migration

A migration of 300 or more servers.

LBAC

See [label-based access control](#).

least privilege

The security best practice of granting the minimum permissions required to perform a task. For more information, see [Apply least-privilege permissions](#) in the IAM documentation.

lift and shift

See [7 Rs](#).

little-endian system

A system that stores the least significant byte first. See also [endianness](#).

LLM

See [large language model](#).

lower environments

See [environment](#).

M

machine learning (ML)

A type of artificial intelligence that uses algorithms and techniques for pattern recognition and learning. ML analyzes and learns from recorded data, such as Internet of Things (IoT) data, to generate a statistical model based on patterns. For more information, see [Machine Learning](#).

main branch

See [branch](#).

malware

Software that is designed to compromise computer security or privacy. Malware might disrupt computer systems, leak sensitive information, or gain unauthorized access. Examples of malware include viruses, worms, ransomware, Trojan horses, spyware, and keyloggers.

managed services

AWS services for which AWS operates the infrastructure layer, the operating system, and platforms, and you access the endpoints to store and retrieve data. Amazon Simple Storage

Service (Amazon S3) and Amazon DynamoDB are examples of managed services. These are also known as *abstracted services*.

manufacturing execution system (MES)

A software system for tracking, monitoring, documenting, and controlling production processes that convert raw materials to finished products on the shop floor.

MAP

See [Migration Acceleration Program](#).

MCP

See [Model Context Protocol](#).

Model Context Protocol (MCP)

A stateless protocol for [agent](#)-to-[tool](#) communication.

MCP server

A service that exposes one or more [tools](#) through the [Model Context Protocol](#).

mechanism

A complete process in which you create a tool, drive adoption of the tool, and then inspect the results in order to make adjustments. A mechanism is a cycle that reinforces and improves itself as it operates. For more information, see [Building mechanisms](#) in the AWS Well-Architected Framework.

member account

All AWS accounts other than the management account that are part of an organization in AWS Organizations. An account can be a member of only one organization at a time.

MES

See [manufacturing execution system](#).

Message Queuing Telemetry Transport (MQTT)

A lightweight, machine-to-machine (M2M) communication protocol, based on the [publish/subscribe](#) pattern, for resource-constrained [IoT](#) devices.

microservice

A small, independent service that communicates over well-defined APIs and is typically owned by small, self-contained teams. For example, an insurance system might include

microservices that map to business capabilities, such as sales or marketing, or subdomains, such as purchasing, claims, or analytics. The benefits of microservices include agility, flexible scaling, easy deployment, reusable code, and resilience. For more information, see [Integrating microservices by using AWS serverless services](#).

microservices architecture

An approach to building an application with independent components that run each application process as a microservice. These microservices communicate through a well-defined interface by using lightweight APIs. Each microservice in this architecture can be updated, deployed, and scaled to meet demand for specific functions of an application. For more information, see [Implementing microservices on AWS](#).

Migration Acceleration Program (MAP)

An AWS program that provides consulting support, training, and services to help organizations build a strong operational foundation for moving to the cloud, and to help offset the initial cost of migrations. MAP includes a migration methodology for executing legacy migrations in a methodical way and a set of tools to automate and accelerate common migration scenarios.

migration at scale

The process of moving the majority of the application portfolio to the cloud in waves, with more applications moved at a faster rate in each wave. This phase uses the best practices and lessons learned from the earlier phases to implement a *migration factory* of teams, tools, and processes to streamline the migration of workloads through automation and agile delivery. This is the third phase of the [AWS migration strategy](#).

migration factory

Cross-functional teams that streamline the migration of workloads through automated, agile approaches. Migration factory teams typically include operations, business analysts and owners, migration engineers, developers, and DevOps professionals working in sprints. Between 20 and 50 percent of an enterprise application portfolio consists of repeated patterns that can be optimized by a factory approach. For more information, see the [discussion of migration factories](#) and the [Cloud Migration Factory guide](#) in this content set.

migration metadata

The information about the application and server that is needed to complete the migration. Each migration pattern requires a different set of migration metadata. Examples of migration metadata include the target subnet, security group, and AWS account.

migration pattern

A repeatable migration task that details the migration strategy, the migration destination, and the migration application or service used. Example: Rehost migration to Amazon EC2 with AWS Application Migration Service.

Migration Portfolio Assessment (MPA)

An online tool that provides information for validating the business case for migrating to the AWS Cloud. MPA provides detailed portfolio assessment (server right-sizing, pricing, TCO comparisons, migration cost analysis) as well as migration planning (application data analysis and data collection, application grouping, migration prioritization, and wave planning). The [MPA tool](#) (requires login) is available free of charge to all AWS consultants and APN Partner consultants.

Migration Readiness Assessment (MRA)

The process of gaining insights about an organization's cloud readiness status, identifying strengths and weaknesses, and building an action plan to close identified gaps, using the AWS CAF. For more information, see the [migration readiness guide](#). MRA is the first phase of the [AWS migration strategy](#).

migration strategy

The approach used to migrate a workload to the AWS Cloud. For more information, see the [7 Rs](#) entry in this glossary and see [Mobilize your organization to accelerate large-scale migrations](#).

ML

See [machine learning](#).

modernization

Transforming an outdated (legacy or monolithic) application and its infrastructure into an agile, elastic, and highly available system in the cloud to reduce costs, gain efficiencies, and take advantage of innovations. For more information, see [Strategy for modernizing applications in the AWS Cloud](#).

modernization readiness assessment

An evaluation that helps determine the modernization readiness of an organization's applications; identifies benefits, risks, and dependencies; and determines how well the organization can support the future state of those applications. The outcome of the assessment is a blueprint of the target architecture, a roadmap that details development phases and

milestones for the modernization process, and an action plan for addressing identified gaps. For more information, see [Evaluating modernization readiness for applications in the AWS Cloud](#).

monolithic applications (monoliths)

Applications that run as a single service with tightly coupled processes. Monolithic applications have several drawbacks. If one application feature experiences a spike in demand, the entire architecture must be scaled. Adding or improving a monolithic application's features also becomes more complex when the code base grows. To address these issues, you can use a microservices architecture. For more information, see [Decomposing monoliths into microservices](#).

MPA

See [Migration Portfolio Assessment](#).

MQTT

See [Message Queuing Telemetry Transport](#).

multiclass classification

A process that helps generate predictions for multiple classes (predicting one of more than two outcomes). For example, an ML model might ask "Is this product a book, car, or phone?" or "Which product category is most interesting to this customer?"

mutable infrastructure

A model that updates and modifies the existing infrastructure for production workloads. For improved consistency, reliability, and predictability, the AWS Well-Architected Framework recommends the use of [immutable infrastructure](#) as a best practice.

O

OAC

See [origin access control](#).

OAI

See [origin access identity](#).

OCM

See [organizational change management](#).

offline migration

A migration method in which the source workload is taken down during the migration process. This method involves extended downtime and is typically used for small, non-critical workloads.

OI

See [operations integration](#).

OLA

See [operational-level agreement](#).

online migration

A migration method in which the source workload is copied to the target system without being taken offline. Applications that are connected to the workload can continue to function during the migration. This method involves zero to minimal downtime and is typically used for critical production workloads.

OPC-UA

See [Open Process Communications - Unified Architecture](#).

Open Process Communications - Unified Architecture (OPC-UA)

A machine-to-machine (M2M) communication protocol for industrial automation. OPC-UA provides an interoperability standard with data encryption, authentication, and authorization schemes.

operational-level agreement (OLA)

An agreement that clarifies what functional IT groups promise to deliver to each other, to support a service-level agreement (SLA).

operational readiness review (ORR)

A checklist of questions and associated best practices that help you understand, evaluate, prevent, or reduce the scope of incidents and possible failures. For more information, see [Operational Readiness Reviews \(ORR\)](#) in the AWS Well-Architected Framework.

operational technology (OT)

Hardware and software systems that work with the physical environment to control industrial operations, equipment, and infrastructure. In manufacturing, the integration of OT and information technology (IT) systems is a key focus for [Industry 4.0](#) transformations.

operations integration (OI)

The process of modernizing operations in the cloud, which involves readiness planning, automation, and integration. For more information, see the [operations integration guide](#).

organization trail

A trail that's created by AWS CloudTrail that logs all events for all AWS accounts in an organization in AWS Organizations. This trail is created in each AWS account that's part of the organization and tracks the activity in each account. For more information, see [Creating a trail for an organization](#) in the CloudTrail documentation.

organizational change management (OCM)

A framework for managing major, disruptive business transformations from a people, culture, and leadership perspective. OCM helps organizations prepare for, and transition to, new systems and strategies by accelerating change adoption, addressing transitional issues, and driving cultural and organizational changes. In the AWS migration strategy, this framework is called *people acceleration*, because of the speed of change required in cloud adoption projects. For more information, see the [OCM guide](#).

origin access control (OAC)

In CloudFront, an enhanced option for restricting access to secure your Amazon Simple Storage Service (Amazon S3) content. OAC supports all S3 buckets in all AWS Regions, server-side encryption with AWS KMS (SSE-KMS), and dynamic PUT and DELETE requests to the S3 bucket.

origin access identity (OAI)

In CloudFront, an option for restricting access to secure your Amazon S3 content. When you use OAI, CloudFront creates a principal that Amazon S3 can authenticate with. Authenticated principals can access content in an S3 bucket only through a specific CloudFront distribution. See also [OAC](#), which provides more granular and enhanced access control.

ORR

See [operational readiness review](#).

OT

See [operational technology](#).

outbound (egress) VPC

In an AWS multi-account architecture, a VPC that handles network connections that are initiated from within an application. The [AWS Security Reference Architecture](#) recommends

setting up your Network account with inbound, outbound, and inspection VPCs to protect the two-way interface between your application and the broader internet.

P

permissions boundary

An IAM management policy that is attached to IAM principals to set the maximum permissions that the user or role can have. For more information, see [Permissions boundaries](#) in the IAM documentation.

personally identifiable information (PII)

Information that, when viewed directly or paired with other related data, can be used to reasonably infer the identity of an individual. Examples of PII include names, addresses, and contact information.

PII

See [personally identifiable information](#).

playbook

A set of predefined steps that capture the work associated with migrations, such as delivering core operations functions in the cloud. A playbook can take the form of scripts, automated runbooks, or a summary of processes or steps required to operate your modernized environment.

PLC

See [programmable logic controller](#).

PLM

See [product lifecycle management](#).

policy

An object that can define permissions (see [identity-based policy](#)), specify access conditions (see [resource-based policy](#)), or define the maximum permissions for all accounts in an organization in AWS Organizations (see [service control policy](#)).

polyglot persistence

Independently choosing a microservice's data storage technology based on data access patterns and other requirements. If your microservices have the same data storage technology, they can encounter implementation challenges or experience poor performance. Microservices are more easily implemented and achieve better performance and scalability if they use the data store best adapted to their requirements.

portfolio assessment

A process of discovering, analyzing, and prioritizing the application portfolio in order to plan the migration. For more information, see [Evaluating migration readiness](#).

predicate

A query condition that returns true or false, commonly located in a WHERE clause.

predicate pushdown

A database query optimization technique that filters the data in the query before transfer. This reduces the amount of data that must be retrieved and processed from the relational database, and it improves query performance.

preventative control

A security control that is designed to prevent an event from occurring. These controls are a first line of defense to help prevent unauthorized access or unwanted changes to your network. For more information, see [Preventative controls](#) in *Implementing security controls on AWS*.

principal

An entity in AWS that can perform actions and access resources. This entity is typically a root user for an AWS account, an IAM role, or a user. For more information, see *Principal* in [Roles terms and concepts](#) in the IAM documentation.

privacy by design

A system engineering approach that takes privacy into account through the whole development process.

private hosted zones

A container that holds information about how you want Amazon Route 53 to respond to DNS queries for a domain and its subdomains within one or more VPCs. For more information, see [Working with private hosted zones](#) in the Route 53 documentation.

proactive control

A [security control](#) designed to prevent the deployment of noncompliant resources. These controls scan resources before they are provisioned. If the resource is not compliant with the control, then it isn't provisioned. For more information, see the [Controls reference guide](#) in the AWS Control Tower documentation and see [Proactive controls](#) in *Implementing security controls on AWS*.

product lifecycle management (PLM)

The management of data and processes for a product throughout its entire lifecycle, from design, development, and launch, through growth and maturity, to decline and removal.

production environment

See [environment](#).

programmable logic controller (PLC)

In manufacturing, a highly reliable, adaptable computer that monitors machines and automates manufacturing processes.

prompt chaining

Using the output of one [LLM](#) prompt as the input for the next prompt to generate better responses. This technique is used to break down a complex task into subtasks, or to iteratively refine or expand a preliminary response. It helps improve the accuracy and relevance of a model's responses and allows for more granular, personalized results.

pseudonymization

The process of replacing personal identifiers in a dataset with placeholder values. Pseudonymization can help protect personal privacy. Pseudonymized data is still considered to be personal data.

publish/subscribe (pub/sub)

A pattern that enables asynchronous communications among microservices to improve scalability and responsiveness. For example, in a microservices-based [MES](#), a microservice can publish event messages to a channel that other microservices can subscribe to. The system can add new microservices without changing the publishing service.

Q

query plan

A series of steps, like instructions, that are used to access the data in a SQL relational database system.

query plan regression

When a database service optimizer chooses a less optimal plan than it did before a given change to the database environment. This can be caused by changes to statistics, constraints, environment settings, query parameter bindings, and updates to the database engine.

R

RACI matrix

See [responsible, accountable, consulted, informed \(RACI\)](#).

RAG

See [Retrieval Augmented Generation](#).

ransomware

A malicious software that is designed to block access to a computer system or data until a payment is made.

RASCI matrix

See [responsible, accountable, consulted, informed \(RACI\)](#).

RCAC

See [row and column access control](#).

read replica

A copy of a database that's used for read-only purposes. You can route queries to the read replica to reduce the load on your primary database.

re-architect

See [7 Rs](#).

recovery point objective (RPO)

The maximum acceptable amount of time since the last data recovery point. This determines what is considered an acceptable loss of data between the last recovery point and the interruption of service.

recovery time objective (RTO)

The maximum acceptable delay between the interruption of service and restoration of service.

refactor

See [7 Rs](#).

Region

A collection of AWS resources in a geographic area. Each AWS Region is isolated and independent of the others to provide fault tolerance, stability, and resilience. For more information, see [Specify which AWS Regions your account can use](#).

regression

An ML technique that predicts a numeric value. For example, to solve the problem of "What price will this house sell for?" an ML model could use a linear regression model to predict a house's sale price based on known facts about the house (for example, the square footage).

rehost

See [7 Rs](#).

release

In a deployment process, the act of promoting changes to a production environment.

relocate

See [7 Rs](#).

replatform

See [7 Rs](#).

repurchase

See [7 Rs](#).

resiliency

An application's ability to resist or recover from disruptions. [High availability](#) and [disaster recovery](#) are common considerations when planning for resiliency in the AWS Cloud. For more information, see [AWS Cloud Resilience](#).

resource-based policy

A policy attached to a resource, such as an Amazon S3 bucket, an endpoint, or an encryption key. This type of policy specifies which principals are allowed access, supported actions, and any other conditions that must be met.

responsible, accountable, consulted, informed (RACI) matrix

A matrix that defines the roles and responsibilities for all parties involved in migration activities and cloud operations. The matrix name is derived from the responsibility types defined in the matrix: responsible (R), accountable (A), consulted (C), and informed (I). The support (S) type is optional. If you include support, the matrix is called a *RASCI matrix*, and if you exclude it, it's called a *RACI matrix*.

responsive control

A security control that is designed to drive remediation of adverse events or deviations from your security baseline. For more information, see [Responsive controls](#) in *Implementing security controls on AWS*.

retain

See [7 Rs](#).

retire

See [7 Rs](#).

Retrieval Augmented Generation (RAG)

A [generative AI](#) technology in which an [LLM](#) references an authoritative data source that is outside of its training data sources before generating a response. For example, a RAG model might perform a semantic search of an organization's knowledge base or custom data. For more information, see [What is RAG](#).

rotation

The process of periodically updating a [secret](#) to make it more difficult for an attacker to access the credentials.

row and column access control (RCAC)

The use of basic, flexible SQL expressions that have defined access rules. RCAC consists of row permissions and column masks.

RPO

See [recovery point objective](#).

RTO

See [recovery time objective](#).

runbook

A set of manual or automated procedures required to perform a specific task. These are typically built to streamline repetitive operations or procedures with high error rates.

S

SAML 2.0

An open standard that many identity providers (IdPs) use. This feature enables federated single sign-on (SSO), so users can log into the AWS Management Console or call the AWS API operations without you having to create user in IAM for everyone in your organization. For more information about SAML 2.0-based federation, see [About SAML 2.0-based federation](#) in the IAM documentation.

SCADA

See [supervisory control and data acquisition](#).

SCP

See [service control policy](#).

secret

In AWS Secrets Manager, confidential or restricted information, such as a password or user credentials, that you store in encrypted form. It consists of the secret value and its metadata. The secret value can be binary, a single string, or multiple strings. For more information, see [What's in a Secrets Manager secret?](#) in the Secrets Manager documentation.

security by design

A system engineering approach that takes security into account through the whole development process.

security control

A technical or administrative guardrail that prevents, detects, or reduces the ability of a threat actor to exploit a security vulnerability. There are four primary types of security controls: [preventative](#), [detective](#), [responsive](#), and [proactive](#).

security hardening

The process of reducing the attack surface to make it more resistant to attacks. This can include actions such as removing resources that are no longer needed, implementing the security best practice of granting least privilege, or deactivating unnecessary features in configuration files.

security information and event management (SIEM) system

Tools and services that combine security information management (SIM) and security event management (SEM) systems. A SIEM system collects, monitors, and analyzes data from servers, networks, devices, and other sources to detect threats and security breaches, and to generate alerts.

security response automation

A predefined and programmed action that is designed to automatically respond to or remediate a security event. These automations serve as [detective](#) or [responsive](#) security controls that help you implement AWS security best practices. Examples of automated response actions include modifying a VPC security group, patching an Amazon EC2 instance, or rotating credentials.

server-side encryption

Encryption of data at its destination, by the AWS service that receives it.

service control policy (SCP)

A policy that provides centralized control over permissions for all accounts in an organization in AWS Organizations. SCPs define guardrails or set limits on actions that an administrator can delegate to users or roles. You can use SCPs as allow lists or deny lists, to specify which services or actions are permitted or prohibited. For more information, see [Service control policies](#) in the AWS Organizations documentation.

service endpoint

The URL of the entry point for an AWS service. You can use the endpoint to connect programmatically to the target service. For more information, see [AWS service endpoints](#) in *AWS General Reference*.

service-level agreement (SLA)

An agreement that clarifies what an IT team promises to deliver to their customers, such as service uptime and performance.

service-level indicator (SLI)

A measurement of a performance aspect of a service, such as its error rate, availability, or throughput.

service-level objective (SLO)

A target metric that represents the health of a service, as measured by a [service-level indicator](#).

shared responsibility model

A model describing the responsibility you share with AWS for cloud security and compliance. AWS is responsible for security *of* the cloud, whereas you are responsible for security *in* the cloud. For more information, see [Shared responsibility model](#).

Shadow AI

Unauthorized [AI](#) applications built or used outside of governed channels within an organization.

SIEM

See [security information and event management system](#).

single point of failure (SPOF)

A failure in a single, critical component of an application that can disrupt the system.

SLA

See [service-level agreement](#).

SLI

See [service-level indicator](#).

SLO

See [service-level objective](#).

split-and-seed model

A pattern for scaling and accelerating modernization projects. As new features and product releases are defined, the core team splits up to create new product teams. This helps scale your organization's capabilities and services, improves developer productivity, and supports rapid innovation. For more information, see [Phased approach to modernizing applications in the AWS Cloud](#).

SPOF

See [single point of failure](#).

star schema

A database organizational structure that uses one large fact table to store transactional or measured data and uses one or more smaller dimensional tables to store data attributes. This structure is designed for use in a [data warehouse](#) or for business intelligence purposes.

strangler fig pattern

An approach to modernizing monolithic systems by incrementally rewriting and replacing system functionality until the legacy system can be decommissioned. This pattern uses the analogy of a fig vine that grows into an established tree and eventually overcomes and replaces its host. The pattern was [introduced by Martin Fowler](#) as a way to manage risk when rewriting monolithic systems. For an example of how to apply this pattern, see [Modernizing legacy Microsoft ASP.NET \(ASMX\) web services incrementally by using containers and Amazon API Gateway](#).

subnet

A range of IP addresses in your VPC. A subnet must reside in a single Availability Zone.

supervisory control and data acquisition (SCADA)

In manufacturing, a system that uses hardware and software to monitor physical assets and production operations.

symmetric encryption

An encryption algorithm that uses the same key to encrypt and decrypt the data.

synthetic testing

Testing a system in a way that simulates user interactions to detect potential issues or to monitor performance. You can use [Amazon CloudWatch Synthetics](#) to create these tests.

system prompt

A technique for providing context, instructions, or guidelines to an [LLM](#) to direct its behavior. System prompts help set context and establish rules for interactions with users.

T

tags

Key-value pairs that act as metadata for organizing your AWS resources. Tags can help you manage, identify, organize, search for, and filter resources. For more information, see [Tagging your AWS resources](#).

target variable

The value that you are trying to predict in supervised ML. This is also referred to as an *outcome variable*. For example, in a manufacturing setting the target variable could be a product defect.

task list

A tool that is used to track progress through a runbook. A task list contains an overview of the runbook and a list of general tasks to be completed. For each general task, it includes the estimated amount of time required, the owner, and the progress.

test environment

See [environment](#).

training

To provide data for your ML model to learn from. The training data must contain the correct answer. The learning algorithm finds patterns in the training data that map the input data attributes to the target (the answer that you want to predict). It outputs an ML model that captures these patterns. You can then use the ML model to make predictions on new data for which you don't know the target.

tool

A function or API that an [agent](#) can invoke to perform operations in external systems.

transit gateway

A network transit hub that you can use to interconnect your VPCs and on-premises networks. For more information, see [What is a transit gateway](#) in the AWS Transit Gateway documentation.

trunk-based workflow

An approach in which developers build and test features locally in a feature branch and then merge those changes into the main branch. The main branch is then built to the development, preproduction, and production environments, sequentially.

trusted access

Granting permissions to a service that you specify to perform tasks in your organization in AWS Organizations and in its accounts on your behalf. The trusted service creates a service-linked role in each account, when that role is needed, to perform management tasks for you. For more information, see [Using AWS Organizations with other AWS services](#) in the AWS Organizations documentation.

tuning

To change aspects of your training process to improve the ML model's accuracy. For example, you can train the ML model by generating a labeling set, adding labels, and then repeating these steps several times under different settings to optimize the model.

two-pizza team

A small DevOps team that you can feed with two pizzas. A two-pizza team size ensures the best possible opportunity for collaboration in software development.

U

uncertainty

A concept that refers to imprecise, incomplete, or unknown information that can undermine the reliability of predictive ML models. There are two types of uncertainty: *Epistemic uncertainty* is caused by limited, incomplete data, whereas *aleatoric uncertainty* is caused by the noise and randomness inherent in the data.

undifferentiated tasks

Also known as *heavy lifting*, work that is necessary to create and operate an application but that doesn't provide direct value to the end user or provide competitive advantage. Examples of undifferentiated tasks include procurement, maintenance, and capacity planning.

upper environments

See [environment](#).

V

vacuuming

A database maintenance operation that involves cleaning up after incremental updates to reclaim storage and improve performance.

version control

Processes and tools that track changes, such as changes to source code in a repository.

VPC peering

A connection between two VPCs that allows you to route traffic by using private IP addresses. For more information, see [What is VPC peering](#) in the Amazon VPC documentation.

vulnerability

A software or hardware flaw that compromises the security of the system.

W

warm cache

A buffer cache that contains current, relevant data that is frequently accessed. The database instance can read from the buffer cache, which is faster than reading from the main memory or disk.

warm data

Data that is infrequently accessed. When querying this kind of data, moderately slow queries are typically acceptable.

window function

A SQL function that performs a calculation on a group of rows that relate in some way to the current record. Window functions are useful for processing tasks, such as calculating a moving average or accessing the value of rows based on the relative position of the current row.

workload

A collection of resources and code that delivers business value, such as a customer-facing application or backend process.

workstream

Functional groups in a migration project that are responsible for a specific set of tasks. Each workstream is independent but supports the other workstreams in the project. For example, the portfolio workstream is responsible for prioritizing applications, wave planning, and collecting migration metadata. The portfolio workstream delivers these assets to the migration workstream, which then migrates the servers and applications.

WORM

See [write once, read many](#).

WQF

See [AWS Workload Qualification Framework](#).

write once, read many (WORM)

A storage model that writes data a single time and prevents the data from being deleted or modified. Authorized users can read the data as many times as needed, but they cannot change it. This data storage infrastructure is considered [immutable](#).

Z

zero-day exploit

An attack, typically malware, that takes advantage of a [zero-day vulnerability](#).

zero-day vulnerability

An unmitigated flaw or vulnerability in a production system. Threat actors can use this type of vulnerability to attack the system. Developers frequently become aware of the vulnerability as a result of the attack.

zero-shot prompting

Providing an [LLM](#) with instructions for performing a task but no examples (*shots*) that can help guide it. The LLM must use its pre-trained knowledge to handle the task. The effectiveness of zero-shot prompting depends on the complexity of the task and the quality of the prompt. See also [few-shot prompting](#).

zombie application

An application that has an average CPU and memory usage below 5 percent. In a migration project, it is common to retire these applications.