

AWS Academy Data Engineering Practice Test (Sample)

Study Guide



Everything you need from our exam experts!

Copyright © 2026 by Examzify - A Kaluba Technologies Inc. product.

ALL RIGHTS RESERVED.

No part of this book may be reproduced or transferred in any form or by any means, graphic, electronic, or mechanical, including photocopying, recording, web distribution, taping, or by any information storage retrieval system, without the written permission of the author.

Notice: Examzify makes every reasonable effort to obtain accurate, complete, and timely information about this product from reliable sources.

SAMPLE

Table of Contents

- Copyright** 1
- Table of Contents** 2
- Introduction** 3
- How to Use This Guide** 4
- Questions** 5
- Answers** 8
- Explanations** 10
- Next Steps** 16

SAMPLE

Introduction

Preparing for a certification exam can feel overwhelming, but with the right tools, it becomes an opportunity to build confidence, sharpen your skills, and move one step closer to your goals. At Examzify, we believe that effective exam preparation isn't just about memorization, it's about understanding the material, identifying knowledge gaps, and building the test-taking strategies that lead to success.

This guide was designed to help you do exactly that.

Whether you're preparing for a licensing exam, professional certification, or entry-level qualification, this book offers structured practice to reinforce key concepts. You'll find a wide range of multiple-choice questions, each followed by clear explanations to help you understand not just the right answer, but why it's correct.

The content in this guide is based on real-world exam objectives and aligned with the types of questions and topics commonly found on official tests. It's ideal for learners who want to:

- Practice answering questions under realistic conditions,
- Improve accuracy and speed,
- Review explanations to strengthen weak areas, and
- Approach the exam with greater confidence.

We recommend using this book not as a stand-alone study tool, but alongside other resources like flashcards, textbooks, or hands-on training. For best results, we recommend working through each question, reflecting on the explanation provided, and revisiting the topics that challenge you most.

Remember: successful test preparation isn't about getting every question right the first time, it's about learning from your mistakes and improving over time. Stay focused, trust the process, and know that every page you turn brings you closer to success.

Let's begin.

How to Use This Guide

This guide is designed to help you study more effectively and approach your exam with confidence. Whether you're reviewing for the first time or doing a final refresh, here's how to get the most out of your Examzify study guide:

1. Start with a Diagnostic Review

Skim through the questions to get a sense of what you know and what you need to focus on. Your goal is to identify knowledge gaps early.

2. Study in Short, Focused Sessions

Break your study time into manageable blocks (e.g. 30 - 45 minutes). Review a handful of questions, reflect on the explanations.

3. Learn from the Explanations

After answering a question, always read the explanation, even if you got it right. It reinforces key points, corrects misunderstandings, and teaches subtle distinctions between similar answers.

4. Track Your Progress

Use bookmarks or notes (if reading digitally) to mark difficult questions. Revisit these regularly and track improvements over time.

5. Simulate the Real Exam

Once you're comfortable, try taking a full set of questions without pausing. Set a timer and simulate test-day conditions to build confidence and time management skills.

6. Repeat and Review

Don't just study once, repetition builds retention. Re-attempt questions after a few days and revisit explanations to reinforce learning. Pair this guide with other Examzify tools like flashcards, and digital practice tests to strengthen your preparation across formats.

There's no single right way to study, but consistent, thoughtful effort always wins. Use this guide flexibly, adapt the tips above to fit your pace and learning style. You've got this!

Questions

SAMPLE

- 1. Which service is used for real-time analytics on streaming data?**
 - A. Amazon Kinesis Data Analytics**
 - B. Amazon EMR**
 - C. AWS Glue**
 - D. Amazon RDS**

- 2. Which statement is NOT correct regarding Apache Hadoop?**
 - A. Hadoop is designed for distributed storage and processing**
 - B. Hadoop is best suited to batch processing**
 - C. Hadoop is best suited to real-time analytics applications**
 - D. Hadoop promotes high-throughput access to application data**

- 3. What is one of the main advantages of using partitioning in database management?**
 - A. Increased data duplication**
 - B. Improved ease of maintenance**
 - C. Optimized query speeds**
 - D. Elimination of data redundancy**

- 4. Which option is NOT a flow state in AWS Step Functions?**
 - A. Pass**
 - B. Choice**
 - C. Loop**
 - D. Task**

- 5. Which AWS service is best suited for ingesting data from a software as a service (SaaS) application?**
 - A. Amazon S3**
 - B. Amazon AppFlow**
 - C. Amazon Kinesis**
 - D. Amazon RDS**

- 6. What is the main use case for AWS IoT Core?**
- A. Web hosting**
 - B. Mobile application development**
 - C. Secure data ingestion from edge devices**
 - D. Cloud storage management**
- 7. What does the AWS Direct Connect service provide?**
- A. A dedicated network connection to on-premises**
 - B. A cloud storage solution**
 - C. A virtual private cloud**
 - D. A managed database service**
- 8. What is the main advantage of using Apache Spark for machine learning?**
- A. Cost-effectiveness**
 - B. Built-in workflow management**
 - C. Better support for batch processing**
 - D. High performance for iterative algorithms**
- 9. Which service can automatically scale resources in response to demand?**
- A. Amazon EC2 Auto Scaling**
 - B. Amazon RDS**
 - C. AWS Lambda**
 - D. Amazon EFS**
- 10. Why is unstructured data considered more flexible?**
- A. It follows a strict schema.**
 - B. It can be stored in any format.**
 - C. It is easier to query than structured data.**
 - D. It has no fixed format and adapts to various use cases.**

Answers

SAMPLE

1. A
2. C
3. C
4. C
5. B
6. C
7. A
8. D
9. A
10. D

SAMPLE

Explanations

SAMPLE

1. Which service is used for real-time analytics on streaming data?

- A. Amazon Kinesis Data Analytics**
- B. Amazon EMR**
- C. AWS Glue**
- D. Amazon RDS**

Amazon Kinesis Data Analytics is specifically designed for real-time analytics on streaming data. It enables users to process and analyze real-time data from sources like Amazon Kinesis Data Streams and can also integrate with other AWS services. This service allows you to run SQL queries on streaming data, which means you can monitor, process, and gain insights from data as it is generated, making it essential for applications that require immediate analytics, such as live monitoring or real-time dashboards. Other services mentioned serve different purposes: Amazon EMR is primarily used for big data processing and batch analytics using frameworks like Apache Hadoop and Apache Spark, which are not focused on real-time data analysis. AWS Glue is primarily an ETL (Extract, Transform, Load) service designed for preparing and transforming data for analysis, rather than for real-time analysis. Amazon RDS is a managed relational database service that handles structured data and is optimized for transactional workloads, not specifically for real-time streaming data analytics.

2. Which statement is NOT correct regarding Apache Hadoop?

- A. Hadoop is designed for distributed storage and processing**
- B. Hadoop is best suited to batch processing**
- C. Hadoop is best suited to real-time analytics applications**
- D. Hadoop promotes high-throughput access to application data**

The assertion that Hadoop is best suited to real-time analytics applications is not correct. Apache Hadoop is primarily designed for batch processing rather than real-time analytics. Its architecture focuses on storing and processing large datasets in a distributed manner, making it highly effective for tasks that can be performed in multiple stages over a significant duration, such as ETL processes, large-scale data transformations, and extensive report generation. While there are components within the Hadoop ecosystem, such as Apache Storm or Apache Kafka, that can handle real-time processing, the core Hadoop framework (particularly the MapReduce component) is optimized for batch workloads, leading to higher throughput rather than low-latency, real-time response. Thus, emphasizing Hadoop's strengths aligns more with batch processing and high-throughput scenarios rather than immediate data query responses typically associated with real-time analytics applications.

3. What is one of the main advantages of using partitioning in database management?

- A. Increased data duplication
- B. Improved ease of maintenance
- C. Optimized query speeds**
- D. Elimination of data redundancy

Using partitioning in database management significantly boosts query speeds, which is one of its primary advantages. Partitioning divides a large database into smaller, more manageable segments called partitions. When queries are executed, the database system can focus on only the relevant partitions rather than scanning the entire dataset. This selective reading minimizes the amount of data processed, reduces scan times, and enhances overall query performance. In addition to optimizing query speeds, partitioning can also improve the efficiency of data management tasks, as maintenance operations can be performed on individual partitions rather than the entire dataset. However, the most notable benefit for query performance remains the ability to quickly access smaller, targeted subsets of data. This leads to faster response times for end-users and improves the overall performance of applications relying on the database. The other choices address aspects that can be beneficial but are not the primary advantage of partitioning. For instance, while partitioning may help with maintenance and organization, such benefits are not as directly impactful on performance as the speed improvements resulting from optimized data access.

4. Which option is NOT a flow state in AWS Step Functions?

- A. Pass
- B. Choice
- C. Loop**
- D. Task

In AWS Step Functions, a flow state refers to different types of states that control the flow of execution within a state machine. Each state type serves distinct purposes in orchestrating complex workflows. The Pass state allows you to pass input to output without performing work, effectively enabling the transition of data between states without executing any action. The Choice state provides branching logic, allowing the workflow to take different paths based on conditions specified in the state. The Task state is utilized to execute a specific task, such as invoking a Lambda function or making a call to an API. The term "Loop" does not correspond to a defined flow state in AWS Step Functions. While looping behavior can be achieved through iterative constructs like chaining states or using a combination of Choice and Task states, there is no dedicated state type named "Loop." Instead, developers implement loops through other means within the state machine design, emphasizing that AWS Step Functions does not explicitly provide a flow state called Loop.

5. Which AWS service is best suited for ingesting data from a software as a service (SaaS) application?

A. Amazon S3

B. Amazon AppFlow

C. Amazon Kinesis

D. Amazon RDS

Amazon AppFlow is the most appropriate service for ingesting data from a software as a service (SaaS) application due to its specialized design for seamlessly connecting SaaS applications with AWS services. This managed integration service allows users to securely transfer data between various SaaS applications and AWS services with minimal coding and setup. It supports the connection to a wide range of popular SaaS applications, making data ingestion straightforward and efficient. By providing built-in connectors, AppFlow can handle data transformations and facilitate features like filtering and mapping before it moves the data to its destination, such as Amazon S3 or other AWS services. This integration capability simplifies the process of automating the flow of data, which can be particularly useful for businesses looking to analyze and utilize data from multiple sources without heavy lifting or complicated ETL (Extract, Transform, Load) processes. In contrast, while Amazon S3 is primarily used for object storage and can receive data, it does not have the specialized capabilities needed for direct integration with SaaS applications like AppFlow. Similarly, Amazon Kinesis is more suited for real-time data streaming and processing and is not specifically designed for SaaS data ingestion. Amazon RDS is a relational database service that can store data but does not provide the direct,

6. What is the main use case for AWS IoT Core?

A. Web hosting

B. Mobile application development

C. Secure data ingestion from edge devices

D. Cloud storage management

The primary use case for AWS IoT Core is to enable secure data ingestion from edge devices. AWS IoT Core provides a platform for connecting Internet of Things (IoT) devices to the cloud, allowing these devices to collect and send data securely for further processing, analysis, or storage. This capability is essential for applications that rely on real-time data from distributed devices, such as sensors, cameras, and other smart devices, as it facilitates efficient data flow and integration with other AWS services for analytics and machine learning. This use case leverages features such as device authentication, secure communication protocols, and message routing, making it suitable for various IoT applications across numerous industries, including smart homes, agriculture, industrial automation, and healthcare. In contrast, web hosting, mobile application development, and cloud storage management focus on different aspects of cloud computing and do not specifically leverage the capabilities of IoT device management and secure data ingestion that AWS IoT Core is designed for.

7. What does the AWS Direct Connect service provide?

- A. A dedicated network connection to on-premises**
- B. A cloud storage solution**
- C. A virtual private cloud**
- D. A managed database service**

AWS Direct Connect is a service that enables users to establish a dedicated network connection from their on-premises data center or office to AWS. This service is particularly beneficial for organizations that require a stable, high-speed connection that can reduce latency and improve performance compared to standard internet connections. By providing a direct physical connection, AWS Direct Connect facilitates data transfer to and from the cloud more efficiently and can help eliminate the variability of internet-based connections. The dedicated connection allows businesses to connect their internal networks to a range of AWS services, making it a vital service for enterprises that rely on consistent and secure data transfer between their on-premises environments and AWS environments. As a significant advantage, it can also lead to cost savings on data transfer charges, especially for large-scale data transfers. This has important implications for data-heavy applications and operations, such as big data processing, backup, and disaster recovery strategies. Other available options, while also beneficial services within the AWS ecosystem, do not pertain to establishing a network connection directly to on-premises infrastructure.

8. What is the main advantage of using Apache Spark for machine learning?

- A. Cost-effectiveness**
- B. Built-in workflow management**
- C. Better support for batch processing**
- D. High performance for iterative algorithms**

The primary advantage of using Apache Spark for machine learning lies in its high performance for iterative algorithms. Apache Spark is designed to handle large-scale data processing effectively, and its architecture allows for execution of tasks in memory. This capability significantly enhances the efficiency of iterative algorithms, which are common in many machine learning applications. For instance, when training models, especially those that require multiple passes over the dataset, such as gradient descent methods, Spark's ability to keep data in memory reduces the overhead of reading and writing data to disk repeatedly, resulting in faster computations and smoother experience when dealing with iterative processes. Additionally, Spark's distributed computing framework allows it to scale out across multiple nodes, further improving the performance and speed of machine learning tasks. This makes Spark particularly useful for handling large datasets and complex computations inherent in machine learning workflows, leading to faster model training and evaluation cycles.

9. Which service can automatically scale resources in response to demand?

- A. Amazon EC2 Auto Scaling**
- B. Amazon RDS**
- C. AWS Lambda**
- D. Amazon EFS**

Amazon EC2 Auto Scaling is designed specifically to automatically adjust the number of Amazon EC2 instances in response to changes in demand. This service monitors the specified metrics and conditions, such as CPU utilization or memory usage, and dynamically scales the number of instances up or down based on these metrics. This capability ensures that applications have the necessary resources during peak demand periods while minimizing costs during low usage times. Choosing Amazon EC2 Auto Scaling enables businesses to maintain performance and availability without manual intervention, providing an efficient way to manage resource allocation. Its functionality to set policies for scaling actions also integrates well with other AWS services, allowing for more complex architectures where resource optimization is critical. In contrast, while AWS Lambda also scales automatically, it does so in a different context, operating on a serverless model where individual functions execute in response to events without the need for direct server management. Amazon RDS offers managed relational database services, and while it can handle scaling to an extent through read replicas and instance resizing, it doesn't automatically scale in the same adaptive manner as EC2 Auto Scaling when it comes to varying server workloads. Lastly, Amazon EFS provides scalable file storage and can increase storage capacity automatically, but it doesn't scale compute resources directly in response to demand like EC2 Auto

10. Why is unstructured data considered more flexible?

- A. It follows a strict schema.**
- B. It can be stored in any format.**
- C. It is easier to query than structured data.**
- D. It has no fixed format and adapts to various use cases.**

Unstructured data is considered more flexible primarily because it has no fixed format and can be easily adapted to a variety of use cases. This characteristic allows organizations to collect and utilize a wide range of data types, such as text, images, audio, and video, without the constraints that often come with structured data. Structured data typically conforms to a predefined schema, which can limit how the data can be used and analyzed. In contrast, unstructured data's lack of rigid structure means that it can be stored in any format, making it easier to accommodate new types of data as they arise and to integrate disparate data sources. Another aspect of unstructured data's flexibility is its potential for insights. Since it can contain valuable information that isn't confined to specific fields or formats, organizations can derive insights through processing techniques such as natural language processing, image analysis, and more. This adaptability can be instrumental in today's rapidly evolving data landscape.

Next Steps

Congratulations on reaching the final section of this guide. You've taken a meaningful step toward passing your certification exam and advancing your career.

As you continue preparing, remember that consistent practice, review, and self-reflection are key to success. Make time to revisit difficult topics, simulate exam conditions, and track your progress along the way.

If you need help, have suggestions, or want to share feedback, we'd love to hear from you. Reach out to our team at hello@examzify.com.

Or visit your dedicated course page for more study tools and resources:

<https://awsacademydataengr.examzify.com>

We wish you the very best on your exam journey. You've got this!

SAMPLE