

Microsoft Azure Data Fundamentals SECOND EDITION Exam Ref DP-900

Nicola Farquharson

FREE SAMPLE CHAPTER







Exam Ref DP-900 Microsoft Azure Data Fundamentals

Nicola Farquharson

Exam Ref DP-900 Microsoft Azure Data Fundamentals

Published with the authorization of Microsoft Corporation by: Pearson Education, Inc.

Copyright © 2024 by Pearson Education, Inc.

Hoboken, New Jersey

All rights reserved. This publication is protected by copyright, and permission must be obtained from the publisher prior to any prohibited reproduction, storage in a retrieval system, or transmission in any form or by any means, electronic, mechanical, photocopying, recording, or likewise. For information regarding permissions, request forms, and the appropriate contacts within the Pearson Education Global Rights & Permissions Department, please visit www.pearson.com/permissions.

No patent liability is assumed with respect to the use of the information contained herein. Although every precaution has been taken in the preparation of this book, the publisher and author assume no responsibility for errors or omissions. Nor is any liability assumed for damages resulting from the use of the information contained herein.

ISBN-13: 978-0-13-826190-0 ISBN-10: 0-13-826190-3

Library of Congress Control Number: 2024932734

\$PrintCode

TRADEMARKS

Microsoft and the trademarks listed at http://www.microsoft.com on the "Trademarks" webpage are trademarks of the Microsoft group of companies. All other marks are property of their respective owners.

WARNING AND DISCLAIMER

Every effort has been made to make this book as complete and as accurate as possible, but no warranty or fitness is implied. The information provided is on an "as is" basis. The author, the publisher, and Microsoft Corporation shall have neither liability nor responsibility to any person or entity with respect to any loss or damages arising from the information contained in this book or from the use of the programs accompanying it.

SPECIAL SALES

For information about buying this title in bulk quantities, or for special sales opportunities (which may include electronic versions; custom cover designs; and content particular to your business, training goals, marketing focus, or branding interests), please contact our corporate sales department at corpsales@pearsoned.com or (800) 382-3419.

For government sales inquiries, please contact governmentsales@pearsoned.com.

For questions about sales outside the U.S., please contact intlcs@pearson.com.

CREDITS

EDITOR-IN-CHIEF Brett Bartow

EXECUTIVE EDITOR Loretta Yates

ASSOCIATE EDITOR Shourav Bose

DEVELOPMENT EDITOR Songlin Qiu

MANAGING EDITOR Sandra Schroeder

SENIOR PROJECT EDITOR Tracey Croom

COPY EDITOR Kim Wimpsett

INDEXER Timothy Wright

PROOFREADER Barbara Mack

TECHNICAL EDITOR Owen Auger

EDITORIAL ASSISTANT Cindy Teeters

COVER DESIGNER Twist Creative, Seattle

COMPOSITOR codeMantra

GRAPHICS codeMantra

FIGURE CREDITS Figure 4-13: Umamfals/ Shutterstock, Figure 4-14: Ma8/123RF, Figure 4-14: Yeliena Brovko/ Shutterstock

Contents at a glance

	Introduction	viii
	About the author	Х
CHAPTER 1	Describe core data concept	1
CHAPTER 2	Identify considerations for relational data on Azure	29
CHAPTER 3	Describe considerations for working with non-relational	
	data on Azure	73
CHAPTER 4	Describe an analytics workload on Azure	103
CHAPTER 5	DP-900 Microsoft Azure Data Fundamentals	
	Exam Updates	149
	Index	155

Contents

	Introduction	viii
	Organization of this book	viii
	Microsoft certifications	ix
	Errata, updates, and book support	ix
	Stay in touch	ix
Chapter 1	Describe core data concept	1
	Skill 1.1: Describe ways to represent data	1
	Describe features of structured data	2
	Describe features of semi-structured data	3
	Describe features of unstructured data	4
	Skill 1.2: Identify options for data storage	5
	Describe common formats for data files	6
	Describe types of databases	13
	Skill 1.3: Describe common data workloads.	17
	Describe features of transactional workloads	18
	Describe features of analytical workloads	18
	Skill 1.4: Identify roles and responsibilities for data workloads	19
	Describe responsibilities for database administrators	20
	Describe responsibilities for data engineers	21
	Describe responsibilities for data analysts	22
	Chapter summary	23
	Thought experiment	24
	Thought experiment answers	27
Chapter 2	Identify considerations for relational data on Azure	29
	Skill 2.1: Describe relational concepts	29
	Identify features of relational data	30
	Describe normalization and how it is used	32

	Identify common structured query language (SQL) statements	34
	Identify common database objects	42
	Skill 2.2: Describe relational Azure data services	44
	Describe the Azure SQL family of products including Azure SQL Database, Azure SQL, Azure Managed Instance, and SQL Server on Azure Virtual Machines Identify Azure database services for open-source	45
	database systems	55
	Chapter summary	63
	Thought experiment	63
	Thought experiment answers	67
	Scenario Exercises	69
	Exercise 1: Create a Database and Table	70
	Exercise 2: Explore the Scalability of Azure SQL Database	72
	Exercise 3: Explore the Security Features of Azure SQL Database	72
	Exercise 4: Explore Data Recovery in Azure SQL Database	72
Chapter 3	Describe considerations for working	72
Chapter 3	with non-relational data on Azure	73
Chapter 3	with non-relational data on Azure Skill 3.1: Describe capabilities of Azure Storage	74
Chapter 3	with non-relational data on Azure Skill 3.1: Describe capabilities of Azure Storage Describe Azure Blob storage	74 74
Chapter 3	with non-relational data on Azure Skill 3.1: Describe capabilities of Azure Storage Describe Azure Blob storage Describe Azure Data Lake Storage Gen2	74 74 81
Chapter 3	with non-relational data on Azure Skill 3.1: Describe capabilities of Azure Storage Describe Azure Blob storage	74 74
Chapter 3	with non-relational data on Azure Skill 3.1: Describe capabilities of Azure Storage Describe Azure Blob storage Describe Azure Data Lake Storage Gen2 Describe Azure File storage Describe Azure Table storage	74 74 81 83 86
Chapter 3	with non-relational data on Azure Skill 3.1: Describe capabilities of Azure Storage Describe Azure Blob storage Describe Azure Data Lake Storage Gen2 Describe Azure File storage	74 74 81 83 86
Chapter 3	with non-relational data on Azure Skill 3.1: Describe capabilities of Azure Storage Describe Azure Blob storage Describe Azure Data Lake Storage Gen2 Describe Azure File storage Describe Azure Table storage Skill 3.2: Describe capabilities and features of Azure Cosmos DB	74 74 81 83 86 87
Chapter 3	with non-relational data on Azure Skill 3.1: Describe capabilities of Azure Storage Describe Azure Blob storage Describe Azure Data Lake Storage Gen2 Describe Azure File storage Describe Azure Table storage Skill 3.2: Describe capabilities and features of Azure Cosmos DB Identify use cases for Azure Cosmos DB	74 74 81 83 86 87 88 90
Chapter 3	with non-relational data on Azure Skill 3.1: Describe capabilities of Azure Storage Describe Azure Blob storage Describe Azure Data Lake Storage Gen2 Describe Azure File storage Describe Azure Table storage Skill 3.2: Describe capabilities and features of Azure Cosmos DB Identify use cases for Azure Cosmos DB Describe Azure Cosmos DB APIs	74 74 81 83 86 87 88 90 93
Chapter 3	with non-relational data on Azure Skill 3.1: Describe capabilities of Azure Storage Describe Azure Blob storage Describe Azure Data Lake Storage Gen2 Describe Azure File storage Describe Azure Table storage Skill 3.2: Describe capabilities and features of Azure Cosmos DB Identify use cases for Azure Cosmos DB Describe Azure Cosmos DB APIs Chapter summary	74 74 81 83 86 87 88 90 93 94
Chapter 3	with non-relational data on AzureSkill 3.1: Describe capabilities of Azure StorageDescribe Azure Blob storageDescribe Azure Data Lake Storage Gen2Describe Azure File storageDescribe Azure Table storageSkill 3.2: Describe capabilities and features of Azure Cosmos DB.Identify use cases for Azure Cosmos DBDescribe Azure Cosmos DB APIsChapter summaryThought experiment	74 74 81 83 86 87 88 90 93 94 97

Chapter 4	Describe an analytics workload on Azure	103
	Skill 4.1 Describe common elements of large-scale analytics	104
	Describe large-scale data warehousing architecture	104
	Describe considerations for data ingestion and processing	107
	Describe options for analytical data stores	110
	Describe Azure services for data warehousing	113
	Skill 4.2 Describe consideration for real-time data analytics	
	Describe the difference between batch and streaming data	118
	Identify Microsoft cloud services for real-time analytics	120
	Skill 4.3 Describe data visualization in Microsoft Power Bl	
	Identify capabilities of Power BI	128
	Describe features of data models in Power Bl	129
	Chapter summary	139
	Thought experiment	140
	Thought experiment answers	143
	Scenario Exercises	145
	Task 1: Large-scale data ingestion and processing	146
	Task 2: An analytical data store	146
	Task 3: Real-time analytics	147
	Task 4: Power BI visualization	147
	Task 5: Azure data warehousing services	147
Chapter 5	DP-900 Microsoft Azure Data Fundamentals	
•	Exam Updates	149
	The purpose of this chapter	149
	About possible exam updates	150
	Impact on you and your study plan	150
	News and commentary about the exam objective updates	150
	Updated technical content	152
	Objective mapping	152

Introduction

Many study resources for technical exams focus on detailed, low-level tasks, teaching how to use specific functionalities and accomplish granular objectives. However, the DP-900 exam, which covers Microsoft Azure data fundamentals, adopts a more high-level approach. This exam is designed to build upon your basic understanding of data concepts and extend it to strategic application within Microsoft Azure's data services. It emphasizes a broad understanding of data management, storage, and processing in the Azure cloud environment. The DP-900 exam, and the materials aligned with it, are geared more towards conceptual understanding rather than intricate coding. Most of the code samples and discussions in this context aim to illustrate broader data principles and cloud data service applications, highlighting the strategic aspect of data handling in Azure rather than focusing on detailed coding techniques.

Ideal candidates include those at the outset of their data-related career path, IT professionals seeking to broaden their expertise into cloud data services, and academicians involved in teaching or learning data and cloud technologies. Even professionals from non-technical backgrounds will find the DP-900 exam valuable for understanding data management in the cloud context.

As for the required knowledge and experience, the DP-900 exam is tailored for individuals with a basic understanding of data concepts, including the creation, storage, and processing of data. It expects familiarity with both relational and non-relational data types, as well as an introductory knowledge of data analytics and reporting concepts. A general grasp of cloud concepts, especially around Microsoft Azure, is beneficial, though extensive prior experience is not necessary. The exam content covers Azure data services, data processing, and data storage solutions, so even a fundamental level of exposure to these areas can be advantageous for aspirants.

This book covers every major topic area on the exam, but it does not cover every exam question. Only the Microsoft exam team has access to the exam questions, and Microsoft regularly adds new questions to the exam, making it impossible to cover specific questions. You should consider this book a supplement to your relevant real-world experience and other study materials. If you encounter a topic in this book you are not completely comfortable with, use the "Need more review?" links in the text to find more information and take the time to research and study the topic. Great information is available on MSDN, on TechNet, and in blogs and forums.

Organization of this book

This book is organized by the "Skills measured" list published for the exam. The "Skills measured" list is available for each exam on the Microsoft Learn website: *Microsoft.com/learn*.

Each chapter in this book corresponds to a major topic area in the list, and the technical tasks in each topic area determine a chapter's organization. If an exam covers six major topic areas, for example, the book will contain six chapters.

Microsoft certifications

Microsoft certifications distinguish you by proving your command of a broad set of skills and experience with current Microsoft products and technologies. The exams and corresponding certifications are developed to validate your mastery of critical competencies as you design and develop, or implement and support, solutions with Microsoft products and technologies both on-premises and in the cloud. Certification brings a variety of benefits to the individual and to employers and organizations.

NEED MORE REVIEW? ALL MICROSOFT CERTIFICATIONS

For information about Microsoft certifications, including a full list of available certifications, go to *microsoft.com/learn*.

Check back often to see what is new!

Errata, updates, and book support

We've made every effort to ensure the accuracy of this book and its companion content. You can access updates to this book—in the form of a list of submitted errata and their related corrections—at *MicrosoftPressStore.com/ERDP9002e/errata*.

If you discover an error that is not already listed, please submit it to us at the same page.

For additional book support and information, visit *MicrosoftPressStore.com/Support*.

Please note that product support for Microsoft software and hardware is not offered through the previous addresses. For help with Microsoft software or hardware, go to *support.microsoft.com*.

Stay in touch

Let's keep the conversation going! We're on Twitter: twitter.com/MicrosoftPress.

About the author

NICOLA FARQUHARSON has more than two decades of experience in networking infrastructure and Microsoft technologies, such as artificial intelligence, Microsoft SQL, Microsoft Power BI, data science, Dynamics 365, machine learning, Microsoft Azure, and Azure DevOps. Her extensive background in these technologies ensures a deep and comprehensive coverage of the key concepts necessary for the Microsoft DP-900 certification exam.

Her career journey includes roles as a Microsoft technical trainer and professor, focusing on data, machine learning, and artificial intelligence. These roles have equipped Nicola with valuable insights into the skills and knowledge required in these rapidly advancing fields. Her experience as a cybersecurity/infrastructure senior analyst and IT consultant has also given her a profound understanding of data security and risk management, vital elements in contemporary data management.

Nicola's dedication to professional growth is reflected in her acquisition of multiple Microsoft technology certifications, including Microsoft AI Engineer, Microsoft Data Engineer, Microsoft DevOps Engineer Expert, Microsoft Teams Administrator Associate, Azure Security Engineer Associate, Microsoft 365 Enterprise Administrator, and Microsoft Cloud Solution Expert. These qualifications underscore her commitment to staying updated with the latest technological advancements. This second edition of *Exam Ref DP-900 Microsoft Azure Data Fundamentals* aims to provide readers with a thorough understanding of Azure data fundamentals, preparing them for the DP-900 exam and enhancing their career prospects in the fields of cloud and data technologies.

Describe core data concept

Data has played an integral role in human endeavors throughout history, but its significance has surged exponentially in recent times, driven by technological advancements and the digital revolution. From ancient records inscribed on stone tablets to the vast digital repositories we have today, the evolution of data has profoundly shaped the trajectory of human progress. The field of data management and analytics has witnessed remarkable break-throughs, fueled by the exponential growth of computing power and the escalating complexity of modern business challenges.

The structured organization of data originated during the emergence of relational databases in the 1970s. Visionaries like Edgar F. Codd revolutionized the field with their pioneering work on the relational model, laying the groundwork for structured data management systems so users can efficiently store, retrieve, and manipulate information. As technology advanced, the need to accommodate semi-structured and unstructured data formats became apparent. The explosive proliferation of the internet, social media, and Internet of Things (IoT) devices has given rise to novel data paradigms, spurring the development of specialized storage solutions and advanced analysis techniques.

In today's data-driven world, where knowledge is power, a firm grasp of the core concepts of data is imperative for individuals and businesses alike. This chapter embarks on a captivating journey, exploring the foundational principles of managing, storing, and utilizing data. By mastering these core data concepts, you will attain a solid understanding of how data is structured, stored, and processed within the dynamic Azure landscape. Prepare to explore the vast landscape of data to harness its transformative potential.

Skills covered in this chapter:

- Skill 1.1 Describe ways to represent data
- Skill 1.2 Identify options for data storage
- Skill 1.3 Describe common data workloads
- Skill 1.4 Identify roles and responsibilities for data workloads

Skill 1.1: Describe ways to represent data

In the world of data, the way information is structured and organized plays a crucial role in its effective management and utilization. There are three fundamental ways you can represent data: structured, semi-structured, and unstructured. Each representation offers unique

characteristics and brings its own set of opportunities and challenges to the realm of data management.

By understanding these different ways to represent data, you and your organization can make informed decisions on how to store, process, and analyze your data effectively. Each data representation method offers distinct advantages and is suited to specific use cases. Whether it's the structured precision of relational databases, the flexibility of semi-structured data, or the untapped potential of unstructured data, embracing these representations empowers you to unlock the full value of your data assets.

This skill covers how to:

- Describe features of structured data
- Describe features of semi-structured data
- Describe features of unstructured data

Describe features of structured data

Structured data is a well-organized format that follows predefined schemas, providing efficient storage, retrieval, and analysis. It represents information in a tabular form with clear relation-ships between entities, making it highly suitable for relational databases. The structured nature of this data allows for easy sorting, searching, and querying using Structured Query Language (SQL). Examples of structured data include financial transaction records, customer profiles, and inventory management systems.

Imagine you have a structured data table representing sales transactions in a retail business. The table would have columns such as Transaction ID, Date, Customer Name, Product, Quality, and so on, as shown in Table 1-1. Each row would represent a specific sale, with corresponding values in each column. The structured format allows you and others in the business to efficiently track sales, analyze customer behavior, and generate insights for decision-making.

Transaction table						
Transaction ID	Date	Customer Name	Product	Quality	Unit Price	Total Amount
1	2023-05-10	John Smith	Smartphone	2	\$500	\$1,000
2	2023-05-11	Jane Doe	Laptop	1	\$1200	\$1,200
3	2023-05-11	Mark Johnson	Tablet	3	\$300	\$900

TABLE 1-1 Structured data

Structured data provides organizations with a consistent and organized way to store and manage critical information. It ensures data integrity and facilitates relational operations so you can perform complex queries, generate reports, and derive meaningful insights. Relational databases, such as Azure SQL Database, provide robust and scalable solutions for storing and processing structured data in a structured query language.

NEED MORE REVIEW? STRUCTURED DATA

You can learn more about structured data at *learn.microsoft.com/en-us/training/modules/* explore-core-data-concepts/2-data-formats.

Describe features of semi-structured data

Semi-structured data represents information that does not adhere to a rigid, predefined schema like structured data. It offers flexibility and accommodates varying formats, making it well-suited for capturing diverse attributes and evolving data structures. Unlike structured data, semi-structured data does not require fixed columns or tables. Instead, it uses formats such as JavaScript Object Notation (JSON) or eXtensible Markup Language (XML) to organize data in a hierarchical or nested structure.

Figure 1-1 shows an example of a semi-structured data document in JSON format that represents a social media post by the user JohnDoe123. The document contains fields for the author's username, the timestamp of the post, the content of the post, and the number of likes received. Additionally, there is an array of comments, with each comment containing the author's username, timestamp, and content. The flexible structure allows for optional fields or additional metadata, depending on the specific post, which means you can capture varying data attributes without the need to alter the underlying structure.

```
{
    "author": "JohnDoe123",
    "timestamp": "2023-05-16T10:30:00Z",
    "content": "Just had the most amazing hiking adventure! The views were breathtaking! 
    "Ikes": 102,
    "comments": [
    {
        "author": "JaneSmith456",
        "timestamp": "2023-05-16T11:15:00Z",
        "content": "Wow, those pictures are stunning! I need to visit that place someday.
        ",
        "author": "SamJones789",
        "timestamp": "2023-05-16T11:30:00Z",
        "content": "Tim so jealous! Hiking is my favorite activity. Can you share more details about the trail?"
    }
}
```

FIGURE 1-1 JSON representing semi-structured data

NOTE EXPLORING SEMI-STRUCTURED DATA REPRESENTATION

JSON is just one of several ways in which you can represent semi-structured data.

Semi-structured data is commonly encountered in various domains, including social media feeds, sensor data from IoT devices, and log files. With semi-structured data, businesses can capture and store diverse data sources that may have evolving schemas or complex relations.

To effectively manage and process semi-structured data, you can use specialized databases known as NoSQL (for "not only SQL" or "no SQL") databases. These databases, such as Azure Cosmos DB, provide scalable solutions for storing and querying semi-structured data. They offer flexibility and adaptability, making them suitable for handling diverse data formats and evolving schemas.

NEED MORE REVIEW? SEMI-STRUCTURED DATA

You can learn more about semi-structured data at *learn.microsoft.com/en-us/training/* modules/explore-core-data-concepts/2-data-formats.

Describe features of unstructured data

Unstructured data represents a vast and diverse category of information that lacks a predefined structure or format. It includes data in its rawest form, such as text documents, images, audio files, videos, and more. Unlike structured or semi-structured data, unstructured data does not fit neatly into tables or schemas, making it challenging to organize and analyze using traditional methods.

Data can encompass textual documents such as emails, news articles, or social media posts. It can also include images, such as photographs or scanned documents, audio recordings, and videos. Unstructured data may not have a consistent layout or specific attributes, making it difficult to extract insights using conventional data processing techniques.

In the example shown in Figure 1-2, each post represents an unstructured piece of data. The content varies from user to user, and there is no predefined structure or format governing the posts. Users can freely express their thoughts and emotions and use hashtags, mentions, or other forms of expression.

Effectively managing and deriving value from unstructured data requires specialized tools and techniques. Technologies such as natural language processing (NLP), image recognition, and audio transcription play a significant role in analyzing and extracting meaningful information from unstructured data sources.

In today's digital landscape, unstructured data is prevalent due to the exponential growth of internet, social media, and multimedia content. You can leverage unstructured data for sentiment analysis, customer feedback analysis, image recognition applications, and more. However, its sheer volume and lack of predefined structure pose significant challenges in terms of storage, processing, and analysis.



FIGURE 1-2 Unstructured data representation

To tackle these challenges, cloud-based storage solutions such as Azure Blob Storage provide scalable and cost-effective repositories for unstructured data. Advanced analytics platforms, such as Azure Cognitive Services, leverage machine learning algorithms to derive insights from unstructured data sources.

NEED MORE REVIEW? UNSTRUCTURED DATA

You can learn more about unstructured data at *learn.microsoft.com/en-us/training/modules/* explore-core-data-concepts/2-data-formats.

Skill 1.2: Identify options for data storage

In the modern era of data-driven decision-making, you face a multitude of choices when it comes to storing your valuable data. This skill dives into the realm of data storage options, providing you insights into the various formats and technologies available. By understanding the different options for data storage, you can make informed decisions that align with your organization's specific needs, ensuring efficient data management, scalability, and reliability. Let's embark on a journey to explore the landscape of data storage, where choices abound and the right storage solution can pave the way for you to unlock the true potential of your data assets.

Data storage is a critical aspect of any data-driven organization, encompassing the selection of suitable formats and technologies to store and manage data effectively. This section delves into data storage options, providing you with a comprehensive understanding of the choices available. From traditional file formats to modern database systems, the landscape of data storage is diverse and evolving.

Contextually, the rapid growth of digital data in recent years has necessitated the development of scalable and efficient data storage solutions. Traditional data storage approaches, such as file-based storage systems, have given way to more advanced technologies designed to handle the ever-increasing volume, velocity, and variety of data. You need to consider factors such as data access speed, scalability, security, and cost-effectiveness when selecting the appropriate data storage solution.

By exploring the options for data storage, including common formats for data files and different types of databases, you and others within your organization can gain insights into the benefits and use cases of each option. This knowledge empowers you to make informed decisions about data storage that align with your organization's specific requirements, ensuring data availability, reliability, and accessibility in a rapidly evolving data landscape.

This skill covers how to:

- Describe common formats for data files
- Describe types of databases

Describe common formats for data files

In this vast world of data storage, selecting the appropriate file format is essential for efficient data management and interoperability. Exploring the common formats of data files sheds light on their characteristics and use cases. By understanding these formats, you can make informed decisions about storing and exchanging your data to ensure seamless integration and accessibility across different systems and platforms.

Delimited file format

Delimited file formats are widely used for storing and exchanging tabular data, where fields are separated by specific delimiter characters. These formats offer you simplicity, versatility, and compatibility with various system and applications. In a delimited file, each record is represented as a line, and individual fields within the record are separated by a delimiter character, such as a comma, tab, or pipe. Let's explore the characteristics and benefits of delimited file formats.

Figure 1-3 shows an example of the comma-separated values (CSV) file format, where each line represents an employee record. The fields of each record, such as name, age, job title, and location, are separated by commas. Each comma (or delimiter) acts as a marker to distinguish one field from the other.

John Doe,35,Software Engineer,New York Jane Smith,28,Data Analyst,San Francisco Mark Johnson,42,Project Manager,Chicago

FIGURE 1-3 CSV file format

The first line of a CSV file typically serves as a header, specifying the name of the fields. Subsequent lines contain the actual data, with each field representing a specific attribute of the field. With delimited files, you can easily import/export the data into various software applications and seamlessly exchange data between different systems.

Delimited file formats offer you several advantages. They are human-readable and widely supported, making them accessible across different platforms and programming languages. Because of the simplicity of the format, users can easily process and manipulate data using various tools. Delimited files are also lightweight and space-efficient, as they do not require complex data structures or encoding schemas.

These formats are commonly used for data interchange, data migration, and integration between different systems. They provide a standardized and straightforward approach to represent tabular data. While CSV is the most prevalent delimited format, other variations such as tab-separated values (TSV) and pipe-separated values (PSV) offer alternative delimiters for specific use cases.

NEED MORE REVIEW? DELIMITED TEXT FILES

You can learn more about delimited text files at *learn.microsoft.com/en-us/training/modules/* explore-core-data-concepts/3-file-storage.

JavaScript Object Notation file format

JSON is a widely used file format for representing and exchanging structured data. It provides a lightweight, human-readable, and platform-independent format that is easy to parse and generate. JSON files use a hierarchical structure with key-value pairs to represent complex data structures. Let's explore the characteristics and benefits of the JSON file format.

Figure 1-4 shows a JSON file that stores information about employees in a company. The file consists of key-value pairs, with each pair representing a specific attribute of an employee. The structure allows for nesting and therefore can represent complex data relationships. An employee's record might include attributes such as name, position, department, salary, and more. Nested objects could contain additional information such as address or skills.

```
comment": "Employee 1",
"employees": [
"id": 1,
"name": "John Doe",
"position": "Software Engineer",
"department": "Engineering",
"salary": 75000,
"join date": "2022-01-15"
},
{
" comment": "Employee 2",
"id": 2,
"name": "Jane Smith",
"position": "Marketing Manager",
"department": "Marketing",
"salary": 85000,
"join_date": "2021-06-10"
},
{
    "_comment": "Employee 3",
"id": 3,
"name": "Michael Johnson",
"position": "Sales Representative",
"department": "Sales",
"salary": 60000,
"join_date": "2022-03-01"
1
}
```

FIGURE 1-4 JSON file format representation

NOTE UNLOCKING THE POWER OF JSON

JSON files offer several advantages. They are easily readable by both humans and machines, facilitating data understanding and manipulation. The format is supported by a wide range of programming language and frameworks, making it highly interoperable. JSON's hierarchical structure can represent complex data structures, making it suitable for a variety of use cases.

JSON files are commonly used for data interchange between web-based systems and APIs. They are the preferred format for transmitting structured data over HTTP requests and responses. Additionally, JSON is frequently used for configuration files, log files, and data storage in NoSQL databases because of its flexibility and ease of use.

NEED MORE REVIEW? JAVASCRIPT OBJECT NOTATION

You can learn more about the JSON file format at *learn.microsoft.com/en-us/training/* modules/explore-core-data-concepts/3-file-storage.

XML data format

XML is a versatile file format used for storing and exchanging structured data. It provides a hierarchical structure using tags to define elements and using attributes to describe element properties. XML files are human-readable, platform-independent, and widely supported, making them suitable for a variety of applications.

Take, for instance, an XML file that stores information about books in a library, as shown in Figure 1-5. Each book is represented as an XML element, with nested elements representing different attributes such as title, author, genre, and more. XML allows for flexible nesting and customization and therefore can represent complex data structures and relationships.

library>
<book></book>
<title>The Great Gatsby</title>
<author>F. Scott Fitzgerald</author>
<genre>Classic</genre>
<pre><publication_year>1925</publication_year></pre>
<isbn>9780743273565</isbn>
<book></book>
<title>To Kill a Mockingbird</title>
<author>Harper Lee</author>
<genre>Classic</genre>
<pre><publication_year>1960</publication_year></pre>
<isbn>9780060935467</isbn>
<book></book>
<title>The Catcher in the Rye</title>
<author>J.D. Salinger</author>
<genre>Coming-of-age</genre>
<pre><publication_year>1951</publication_year></pre>
<isbn>9780316769488</isbn>

FIGURE 1-5 The XML file format

XML files offer several advantages. They are easily readable and understandable by both humans and machines. The hierarchical structure of XML facilitates the representation of structured and semi-structured data, making it suitable for diverse use cases. XML files are also self-descriptive, as the tags and attributes provide meaningful metadata about the data they represent.

XML is widely supported by programming languages, databases, and web technologies, making it highly interoperable. It is commonly used for exchanging data, writing configuration files, and processing complex data structures in various domains such as web services, content management systems, and scientific research.

NEED MORE REVIEW? EXTENSIBLE MARKUP LANGUAGE

You can learn more about XML at *learn.microsoft.com/en-us/training/modules/* explore-core-data-concepts/3-file-storage.

Parquet file format

Parquet is a columnar storage file format designed for efficient data processing and analytics in big data environments. It provides a highly optimized and compressed representation of structured data, making it well-suited for handling large datasets. Parquet files organize data by columns rather than rows and therefore support advanced column pruning and predicate push-down optimizations.

Figure 1-6 shows a Parquet file that stores sales data for an e-commerce business. Instead of storing all the attributes of a sale in row-based format, Parquet stores each column of data separately. This columnar organization allows for efficient compression and encoding techniques specific to each column, reducing storage requirements and improving query performance. The example shown contains these columns: OrderID, CustomerID, ProductID, Quantity, Price, and Timestamp.

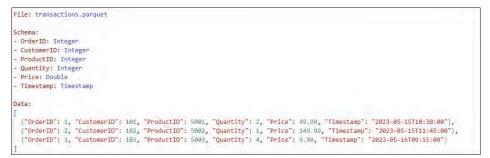


FIGURE 1-6 The Parquet file format

Parquet files offer several advantages. The columnar storage format reduces disk I/O and improves query performance, especially when queries involve only a subset of columns. Parquet leverages advanced compression techniques, such as run-length encoding (RLE) and dictionary encoding, to further reduce the storage footprint while maintaining data integrity. The format also supports nested and complex data types and therefore can represent hierarchical structures.

Parquet is widely adopted in big data processing frameworks such as Apache Hadoop and Apache Spark, as it accelerates analytics workloads by efficiently reading and processing only the required columns. Its compatibility with various data processing tools makes it a preferred choice for high-performance data analytics and data warehousing scenarios.

NEED MORE REVIEW? PARQUET FILE FORMAT

You can learn more about the Parquet file format at *learn.microsoft.com/en-us/training/* modules/explore-core-data-concepts/3-file-storage.

Avro file format

When it comes to storing and exchanging structured data efficiently, the Avro file format stands out as a compact and dynamic solution. Embracing a binary encoding approach, Avro files offer high-performance data processing and storage capabilities. What sets Avro apart is its ability to provide a self-describing schema alongside the data, facilitating schema evolution and dynamic typing.

Figure 1-7 shows an example of the Avro schema defining the structure of the user profile data. The type field specifies that it is a record, the name is User Profile, and the fields are defined in the fields array. In this case, the fields include the user's name (string), age (int), location (string), and interests (an array of strings).

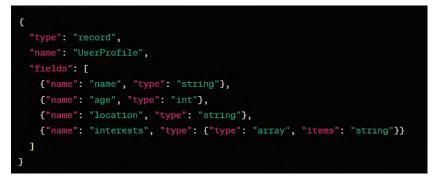


FIGURE 1-7 The Avro file format

Within this file, both the data and its accompanying schema are encapsulated. This inclusion empowers flexibility, allowing the schema to evolve over time without compromising compatibility with existing data. The schema is a rich and diverse representation of structured information.

Avro files offer remarkable advantages. Through their binary encoding, these files achieve exceptional compactness and efficiency, making them ideal for managing and transmitting large volumes of data. The self-describing nature of Avro allows seamless schema evolution, facilitating the adaptation to changing data requirements. Additionally, Avro supports complex data structures, including nested and hierarchical relationships, and therefore can represent intricate information models.

With its outstanding performance, flexibility, and schema evolution capabilities, the Avro file format has garnered widespread adoption within big data processing frameworks such as Apache Hadoop and Apache Spark. Its ability to efficiently handle high-volume data processing, analytics, and seamless interoperability between systems positions Avro as a powerful and dynamic choice.

NOTE EMPOWERING DATA STORAGE

When exploring the features of the Avro file format, it's important to note the following:

- Avro supports popular compression codecs such as Snappy, Deflate, and Bzip2.
- You can define an Avro schema using JSON.

NEED MORE REVIEW? AVRO FILE FORMAT

You can learn more about the Avro file format at *learn.microsoft.com/en-us/training/modules/ explore-core-data-concepts/3-file-storage*.

ORC file format

The optimized row columnar (ORC) file format is designed to optimize performance and storage efficiency for big analytics. It leverages columnar storage and advance compression techniques to deliver high-speed data processing and a reduced storage footprint.

Let's say you have an ORC file that stores sales data for an e-commerce company. Instead of storing data in a row-by-row format, ORC organizes the data into columns, as shown in Figure 1-8.

- Each column contains data of a specific type (e.g., integers, strings, dates).
- Data within each column is stored together, which allows for efficient compression and encoding.
- Rows are typically stored in *stripes* for further compression and optimization.

ORC files offer several advantages. The columnar storage approach reduces disk I/O by reading only the necessary columns during data processing, resulting in faster query performance. Additionally, ORC utilizes advanced compression techniques, such as RLE, dictionary encoding, and bloom filters, to minimize storage requirements while maintaining data integrity. These optimizations make ORC a highly efficient format for big data analytics and data warehousing.

The OCR file format is widely adopted in big data ecosystems, including Apache Hadoop and Apache Hive. Its compatibility with various data processing tools and frameworks makes it a popular choice for optimizing data processing pipelines. By leveraging the benefits of ORC, you can achieve significant performance gains and storage savings when dealing with largescale data analytics.

++
Column1 (e.g., CustomerID)
11
12345
56789
10111
1
++
Column2 (e.g., ProductID)
789 456
1 123
++
++
Column3 (e.g., Quantity)
5
2
10
1
++
++
Column4 (e.g., OrderDate)
11
2023-01-15
2023-02-10
2023-03-05
1
++

FIGURE 1-8 The ORC file format

NEED MORE REVIEW? OPTIMIZED ROW COLUMNAR FORMAT

You can learn more about the ORC format at *learn.microsoft.com/en-us/training/modules/* explore-core-data-concepts/3-file-storage.

Describe types of databases

Effective data management is essential in today's world, where information fuels innovation and drives business success. Choosing the right database type plays a pivotal role in storing, retrieving, and analyzing data with efficiency and precision. Let's journey through the realm of database types, where you will explore the diverse options available and their specific use cases. By gaining a deep understanding of these database types, you can make informed decisions that align with your unique data requirements and unlock the true potential of your data assets. Get ready to embark on an exploration of the vast landscape of databases, where choices abound and where the right selection can pave the way for streamlined data management and impactful analytics.

The rapid growth of digital data in recent years has led to the development of various database types, each tailored to address specific needs and data characteristics. In the context

of the DP-900 exam, it is essential to have a comprehensive understanding of the characteristics, strengths, and use cases of different database types. Whether it's the structured precision of relational databases, the flexibility of NoSQL databases, or the cloud-native capabilities of Azure databases, the right database choice sets the foundation for successful data management and empowers organizations to harness the full potential of their data assets.

Relational database

Relational databases are a cornerstone of modern business data management, providing a structured and efficient approach to organizing and manipulating data. From customer information to financial transactions, relational databases offer a robust foundation for businesses to store, retrieve, and analyze their critical data.

Relational databases organize data into tables, where each table consists of rows and columns. This tabular structure allows businesses to represent entities and their relationships in a clear and organized manner. For example, consider a customer database table with columns such as Customer ID, Name, Email, and Phone Number. Each row represents a specific customer, and the columns store their corresponding attributes. With this structured representation, users can retrieve data efficiently and write complex queries using SQL.

Relational databases offer several advantages for businesses. First, they ensure data consistency and integrity by enforcing relationships between tables through keys and constraints. This reliability means users get accurate reporting, regulatory compliance, and trustworthy decision-making. Second, relational databases support transactional processing, allowing businesses to perform reliable and secure operations on their data, such as creating, updating, and deleting records. This is crucial for maintaining data accuracy and auditability.

Furthermore, relational databases provide powerful query capabilities so that businesses can extract valuable insights from their data. SQL queries can be crafted to join multiple tables, filter databases on specific criteria, aggregate information, and perform complex calculations. This ability to efficiently retrieve and analyze data empowers businesses to make data-driven decisions and gain a competitive edge.

Relational databases have been widely adopted across industries and are supported by a variety of database management systems (DBMSs) such as Microsoft SQL Server, Oracle Database, and MySQL. These systems provide robust features for data storage, indexing performance optimization, and security.

NEED MORE REVIEW? RELATIONAL DATABASES

You can learn more about relational databases at *learn.microsoft.com/en-us/training/* modules/explore-core-data-concepts/4-databases.

Non-relational databases

Non-relational databases, also known as NoSQL databases, have emerged as powerful alternatives to traditional relationship databases for businesses. These databases provide a flexible and scalable approach to storing and managing a vast amount of unstructured and semi-structured data. Let's explore the key aspects and benefits of non-relational databases from a business perspective.

Non-relational databases offer a variety of data models, including document databases, key-value stores, columnar databases, and graph databases. Each data model is designed to address specific data needs and use cases, providing businesses with the flexibility to adapt to diverse data structures.

Consider a business that operates an e-commerce platform. A document database could be used to store customer profiles, where each customer's data is represented as a document containing attributes such as Name, Email, Address, and Purchase History. This flexibility document structure allows for easily storing and retrieving customer information even while accommodating variations in data formats.

Non-relational databases provide several advantages for businesses. First, they offer horizontal scalability, allowing businesses to handle massive data volumes and high-velocity data streams. These databases are built to distribute data across multiple servers and therefore offer seamless scalability as data requirements grow. This scalability supports businesses in handling dynamic workloads and accommodating evolving data demands.

Second, non-relational databases excel in handling unstructured and semi-structured data. They offer schema-less designs, allowing for flexibility in data modeling and accommodating changing data structures over time. This flexibility is crucial in scenarios where data structures are not predefined or where data formats vary greatly.

Furthermore, non-relational databases are highly performant, which results in fast data retrieval and processing. They often employ distributed computing techniques and optimized data storage mechanisms to deliver real-time data access and analytics capabilities. This performance advantage contributes to efficient decision-making and empowers businesses to derive meaningful insight from their data.

Non-relational databases have gained significant traction in modern business environments and are supported by popular systems such as Azure Cosmo DB, MongoDB, and Cassandra, just to name a few. These systems provide robust features for data storage, scalability, and distributed computing.

In a nutshell, non-relational databases offer businesses the flexibility, scalability, and performance required to handle diverse and rapidly growing data. With their ability to handle unstructured data and adapt to changing requirements, non-relational databases empower businesses to unlock the full potential of their data assets.

NEED MORE REVIEW? NON-RELATIONAL DATABASES

You can learn more about non-relational databases at *learn.microsoft.com/en-us/training/* modules/explore-core-data-concepts/4-databases.

Cloud databases

Cloud databases have revolutionized the way businesses store, manage, and analyze their data. These databases leverage the power of cloud computing, offering scalability, flexibility, and ease of management.

With cloud databases, businesses can store and access their data in the cloud, eliminating the need for on-premises infrastructure and maintenance. By leveraging cloud services, businesses can focus on their core competencies while benefiting from the scalability and agility provided by cloud databases.

An example of a cloud database service is Microsoft Azure SQL Database. It offers a fully managed, scalable, and secure database platform so that businesses can rapidly provision and scale their databases based on demand. With Azure SQL Database, businesses can focus on their applications and data, leaving the infrastructure management to the cloud provider.

Cloud databases provide several advantages for businesses. First, they offer scalability, allowing businesses to scale their databases up or down based on workload demands. This flexibility ensures optimal performance and cost efficiency, as resources can be allocated as needed. Businesses can easily handle spikes in user activity or accommodate data growth without significant infrastructure investments.

Second, cloud databases enhance agility and collaboration. They provide seamless access to data from anywhere, which means geographically distributed teams can work together effectively. Cloud-based collaboration and integration tools further streamline workflows and facilitate real-time decision-making, boosting business productivity and efficiency.

Cloud databases also have robust security measures and compliance capabilities. Cloud providers invest heavily in security infrastructure, implementing stringent access controls, encryption, and continuous monitoring to protect business data. Compliance certifications such as International Organization for Standardization (ISO) 2700, the Health Insurance Portability and Accountability Act (HIPAA), and General Data Protection Regulation (GDPR) demonstrate adherence to industry standards and regulatory requirements, providing peace of mind for businesses.

Furthermore, with cloud databases, businesses can leverage advanced analytics and machine learning capabilities. Cloud providers offer integrated analytics services, such as Azure Synapse Analytics, so businesses can derive actionable insights from their data. These services provide powerful tools for data exploration, visualization, and predictive analytics, driving data-driven decision-making and enhancing business competitiveness.

Cloud databases are supported by leading cloud providers, including Microsoft Azure, Amazon Web Services (AWS), and Google Cloud Platform (GCP). These providers offer a wide range of database options, including relational, NoSQL, and specialized databases, tailored to meet diverse business needs. These cloud databases empower businesses by providing scalable, agile, and secure data management solutions. With their ability to scale on demand, foster collaboration, ensure data security, and provide advanced analytics, cloud databases are integral to driving business growth and innovation in the digital era.

Skill 1.3: Describe common data workloads

In the dynamic world of data management, understanding common data workloads is essential for data professionals seeking to harness the transformative potential of data. This skill explores the realm of common data workloads, providing insights into different types of data processing scenarios and their specific requirements. By gaining a deep understanding of these workloads, individuals can effectively design and implement data solutions that align with business needs and drive meaningful insights. Let's take a look at some common data workloads and unlock the work power of data.

In today's data-driven landscape, organizations encounter two primary types of data workloads: transactional workloads and analytical workloads. Transactional workloads focus on the efficient and reliable processing of business transactions, such as capturing customer orders, processing financial transactions, or updating inventory levels. These workloads require strong data consistency, durability, and atomicity/consistency/isolation/durability (ACID) properties to ensure data integrity and reliability.

On the other hand, analytical workloads revolve around deriving insights and knowledge from data support decision-making and strategic planning. Analytical workloads involve complex queries, aggregations, data transformations, and statistical analysis to uncover patterns, trends, and correlations within the data. These workloads typically require scalable processing power, efficient data retrieval, and advanced analytics capabilities to unlock valuable insights and drive informed decisions.

As data volumes continue to grow exponentially and organizations increasingly rely on data-driven insights, understanding and effectively managing these common data workloads become paramount. By comprehending the distinct requirements and characteristics of transactional and analytical workloads, individuals can design appropriate data architectures, select suitable database systems, and implement robust data processing solutions to meet business objectives.

This book delves into the intricacies of these common data workloads, ensuring that data professionals process the knowledge and skills necessary to navigate the dynamic world of data management. By grasping the nuances of transactional and analytical workloads, individuals can contribute to the design and implementation of efficient data solutions, paving the way for business success in an increasingly data-centric era.

This skill covers how to:

- Describe features of transactional workloads
- Describe features of analytical workloads

Describe features of transactional workloads

Transactional workloads play a critical role in ensuring the smooth operation of businesses and maintaining data integrity. These workloads encompass activities such as capturing customer orders, processing financial transactions, and updating inventory levels.

Transactional workloads are designed to handle business operations that involve data modifications, ensuring the accuracy, consistency, and reliability of data. Let's consider an e-commerce platform that processes customer orders. Each customer order represents a transaction that requires capturing the order details, updating inventory levels, and recording the financial transaction. These transactions must be executed reliably and in an atomic manner, meaning they should either complete successfully or be rolled back entirely if an error occurs.

Transactional workloads offer several advantages for business. First, they ensure data consistency and integrity. The ACID properties guide transactional processing, ensuring that data remains in a consistent state even in the event of failure or concurrent access. This integrity is crucial for financial systems, inventory management, and other critical business functions.

Second, transactional workloads support concurrency control and isolation in a multiuser environment, where multiple transactions can occur simultaneously. Transactional processing mechanisms ensure that transactions are executed independently and do not interfere with each other, maintaining data integrity and preventing conflicts.

Furthermore, transactional workloads facilitate data durability and reliability. Transactional systems employ techniques such as write-ahead logging and database recovery mechanisms to ensure that committed transactions persist even in the face of system failures. This durability ensures that critical business operations can be restored and recovered without data loss.

Transactional workloads are supported by various database systems such as relational databases, where ACID properties are typically enforced. These systems provide transaction management features that guarantee data consistency, durability, and isolation. We can say transactional workloads are essential for maintaining accurate data, supporting reliable business operations, and ensuring data integrity. By executing operations in an atomic and consistent manner, businesses can confidently process customer orders, handle financial transactions, and manage inventory levels, fostering trust and reliability in their operations.

NEED MORE REVIEW? TRANSACTIONAL WORKLOADS

You can learn more about transactional workloads at *learn.microsoft.com/en-us/training/modules/explore-core-data-concepts/5-transactional-data-processing*.

Describe features of analytical workloads

Analytical workloads play a pivotal role in extracting valuable insights and patterns from data to support informed decision-making and strategic planning within businesses. These work-loads involve complex data analysis, aggregations, and transformations to uncover meaningful information.

Analytical workloads encompass a range of activities, such as data exploration, statistical analysis, data mining, and predictive modeling. The process begins by identifying relevant data sources and extracting the required data. The advanced analytics techniques, such as data visualization, machine learning, and statistical algorithms, are applied to gain insights and patterns from the data. The results are interpreted and translated into actionable business intelligence, which results in data-driven decision-making.

Analytical workloads serve different data personas within organizations.

- Data analysts: Data analysts leverage analytical workloads to explore and analyze data, uncovering trends, correlations, and patterns that provide valuable insights. They use statistical techniques and data visualization tools to communicate their findings effectively to stakeholders, resulting in evidence-based decision-making.
- Data scientists: Data scientists go beyond analyzing data and utilize advanced analytical methods to develop predictive models, machine learning algorithms, and datadriven solutions. They leverage analytical workloads to build models that forecast future trends, identify opportunities, and optimize business processes.
- Business executives: Business executives rely on analytical workloads to gain highlevel insight and make strategic decisions. They rely on reports, dashboards, and visualizations generated by analytical workloads to monitor key performance indicators, track business metrics, and assess the effectiveness of strategies.
- Data engineers: Data engineers support analytical workloads by designing and implementing the data infrastructure necessary for data analysis. They ensure that data is ingested, processed, and made available in a format that facilitates efficient analysis. They collaborate with data analysts and scientists to ensure data quality and reliability.

Analytical workloads are supported by various technologies and tools, including data platforms, machine learning frameworks, and business intelligence tools. These solutions provide capabilities for data exploration, modeling, visualization, and advanced analytics.

NEED MORE REVIEW? ANALYTICAL WORKLOAD

You can learn more about analytical workloads at *learn.microsoft.com/en-us/training/* modules/explore-core-data-concepts/6-analytical-processing.

Skill 1.4: Identify roles and responsibilities for data workloads

This section focuses on the critical aspect of identifying roles and responsibilities for data workloads. In the world of data management, different professionals contribute their expertise to ensure the efficient handling, processing, and utilization of data. Understanding these roles and responsibilities is vital for organizations to effectively manage and leverage their data assets.

In today's data-driven landscape, organizations rely on dedicated professionals to fulfill specific roles related to data management. This skill highlights the significance of recognizing and assigning the appropriate roles within data workloads. By identifying the individuals responsible for specific tasks, organizations can streamline their data operations, promote collaboration, and optimize the overall data management process.

Assigning roles and responsibilities for data workloads ensures that the right expertise is applied to each aspect of data management. Database and administrators, data engineers, and data analysts play pivotal roles in supporting data workloads, each with their unique skill sets and responsibilities.

Identifying these roles helps establish clear line of responsibility and accountability. By understanding and assigning these roles, organizations can foster collaboration and coordination among professionals involved in data workloads. This alignment promotes effective data management, offers smooth data workflows, and maximizes the value derived from data assets.

This exam skill emphasizes the importance of recognizing these roles and responsibilities in the broader context of data workloads. By understanding the significance of each role and its contribution to successful data management, individuals can grasp the collaborative efforts required to leverage data effectively. Let's take a closer look at each of these data roles and their responsibilities.

This skill covers how to:

- Describe responsibilities for database administrators
- Describe responsibilities for data engineers
- Describe responsibilities for data analysts

Describe responsibilities for database administrators

As a database administrator (DBA), your role is crucial in the management and maintenance of databases, ensuring their smooth operation, integrity, and performance. You are the guardian of data within your organization, responsible for various tasks that contribute to efficiently storing, retrieving, and securing data.

You are involved in the entire life cycle of databases, starting from the initial design and creation to ongoing maintenance and optimization. You work closely with stakeholders to understand data requirements and design database structures that optimize performance and scalability. You determine data models, create database schemas, and define relationships between tables.

Ensuring data security is a critical aspect of your role. You implement access controls, user authentication, and encryption mechanisms to protect sensitive data from unauthorized access or malicious activities. You establish backup and recovery procedures to safeguard against data loss, ensuring the continuity of business operations.

Monitoring databases and optimizing performance are essential responsibilities. You constantly monitor database performance, identifying and resolving bottlenecks to enhance system efficiency. You analyze query performance, tune database configurations, and optimize indexing strategies to improve overall performance and ensure timely data retrieval.

Your expertise also extends to backup and recovery. You design and implement robust backup and recovery strategies to protect data from system failure, human errors, or disasters. You schedule regular backups, perform restoration tests, and maintain disaster recovery plans to ensure data availability and minimize downtime.

Keeping databases up to date is another aspect of your role. You oversee database upgrades and apply patches, ensuring that the database systems are equipped with the latest features, bug fixes, and security updates. You perform compatibility tests and ensure seamless transitions to new versions or releases.

Your role as a DBA is instrumental in maintaining data integrity, ensuring system availability, and supporting business continuity. Your expertise ensures that databases operate efficiently, adhere to industry standards, and meet regulatory requirements. With your skills and knowledge, you contribute to the data-driven systems within your organization functioning smoothly.

NEED MORE REVIEW? DATABASE ADMINISTRATORS

You can learn more about data administrators at *learn.microsoft.com/en-us/training/modules/* explore-roles-responsibilities-world-of-data/2-explore-job-roles.

Describe responsibilities for data engineers

As a data engineer, your role is vital in designing, constructing, and maintaining the data infrastructure and pipeline to promote efficient data processing and analysis. You play a crucial part in the data management process, ensuring that data flows seamlessly across systems and remains accessible for analysis.

Your primary responsibility is to design and construct the data infrastructure necessary for effective data management. You collaborate with stakeholders to understand their data requirements, identify relevant data sources, and determine the best approach to data integration. You develop data pipelines, ensuring the smooth and reliable flow of data from the source systems to the target destinations.

You are involved in data ingestion, where you extract data from various sources such as databases, files, or APIs. You transform and cleanse the data to ensure its quality and consistency, making it suitable for downstream analysis. This may involve tasks such as data extraction, data validation, data cleansing, and data enrichment.

In addition to data ingestion, you are responsible for data transformation and integration. You apply data processing techniques to convert raw data into a usable format, ensuring it aligns with the required data model's schema. This may involve tasks such as data aggregations, data filtering, data normalization, and data enrichment. Data engineering also involves developing data processing workflows. You design and implement efficient workflows that orchestrate the movement and transformation of data, ensuring optimal performance and reliability. This may include using workflow management tools or frameworks to schedule and monitor data processing tasks.

An example of your role as a data engineer could be working on a project to develop a real-time analytics platform for a financial institution. You would be responsible for designing and implementing the data infrastructure, ingesting real-time transaction data from multiple sources, transforming and aggregating the data, and making it available for real-time analysis and reporting.

Your expertise in data engineering contributes to the overall success of data-driven initiatives within your organization. By building robust data pipelines, ensuring data quality and reliability, and implementing efficient data processing, you facilitate effective data analysis and drive actionable insights.

NEED MORE REVIEW? DATA ENGINEER

You can learn more about data engineers at *learn.microsoft.com/en-us/training/modules/* explore-roles-responsibilities-world-of-data/2-explore-job-roles.

Describe responsibilities for data analysts

As a data analyst, your role is crucial in uncovering valuable insight and patterns within data to support informed decision-making within your organization. You play a pivotal role in analyz-ing, interpreting, and visualizing the data to derive meaningful information that drives business strategies. Let's explore your comprehensive role as a data analyst from a business perspective.

Your primary responsibility is to explore and analyze data to uncover trends, correlations, and patterns that provide valuable insights. You work with various data sources, ranging from structured databases to unstructured text files, and use statistical techniques and analytical tools to extract meaningful information.

Data exploration is an essential part of your role. You dive deep into the data, examining its structure, quality, and relationships. You identify relevant variable and metrics to analyze, ensuring that the data is appropriate for the questions or problems at hand.

Once you have gathered and cleaned the data, you apply statistical analysis techniques to identify patterns and relationships. You may perform tasks such as descriptive statistics, hypothesis testing, regression analysis, or clustering to extract insights from the data. These analyses help you uncover trends, anomalies, and relationships that can guide decision-making.

Data visualization is another crucial aspect of your role. You use visual tools and techniques to present data in a clear and concise manner. By creating charts, graphs, and dashboards, you transform complex datasets into easily understandable visual representations. These visualizations help stakeholders to grasp insights quickly and make informed decisions.

An example of your role as a data analyst is analyzing customer behavior data for an e-commerce company. You would examine the data to understand customer preferences, identify purchasing patterns, and segment customers based on their buying behaviors. These insights would then inform marketing strategies, product recommendations, and customer retention efforts.

Your expertise as a data analyst contributes to evidence-based decision-making within your organization. By analyzing and interpreting data, you provide insights that support strategic planning, optimize operations, and drive business growth.

NEED MORE REVIEW? DATA ANALYSTS

You can learn more about data analysts at *learn.microsoft.com/en-us/training/modules/* explore-roles-responsibilities-world-of-data/2-explore-job-roles.

EXAM TIP

When preparing for the exam, you should focus on understanding the core data concepts and their practical application. Familiarize yourself with different data representation formats, storage options, and common data workloads. Pay attention to the roles and responsibilities of database administrators, data engineers, and data analysts. Additionally, practice relating these concepts to real-world scenarios to reinforce your understanding. Being able to apply your knowledge to practical situations will help you excel on the exam and in realworld data management scenarios.

Chapter summary

- Data concepts
 - Data can be represented as structured, semi-structured, or unstructured.
 - There are various ways to store data, which include common file formats and different types of databases.
 - It's essential to understand the difference between transactional and analytical data workloads.
- Roles and responsibilities
 - Database administrators are responsible for managing and maintaining databases.
 - Data engineers play a pivotal role in building and maintaining the data infrastructure.
 - Data analysts primarily focus on extracting valuable insights and analyzing data to inform decisions.

- Key takeaways
 - The representation of data can vary widely depending on its structure and format.
 - Options for data storage range from various file formats to databases.
 - Transactional analytical workloads process unique and distinct characteristics.
 - Several roles work in tandem to ensure the effective management and utilization of data.

Thought experiment

In this thought experiment, demonstrate your skills and knowledge of the topics covered in this chapter. You can find answers to this thought experiment in the next section.

- 1. Your company is planning to implement a new database system to handle its growing customer data. Who among the following roles is primarily responsible for ensuring data security, backup, and system performance?
 - A. Data engineer
 - B. Database administrator
 - C. Data analyst
 - D. None of the above
- 2. Your organization is dealing with a massive influx of data from various sources and needs to design efficient data pipelines. Which role is responsible for designing and building these data pipelines?
 - A. Data engineer
 - B. Database administrator
 - C. Data analyst
 - D. All of the above
- 3. You need to analyze customer purchase data to identify trends and create reports for your management team. Which role is best suited for this task?
 - A. Data engineer
 - B. Database administrator
 - C. Data analyst
 - D. All of the above
- 4. Your organization is facing performance issues with its database system, resulting in slow query execution. Which role would you consult to optimize the database performance?
 - A. Data engineer
 - B. Database administrator

- C. Data analyst
- D. All of the above
- 5. Your company wants to implement real-time analytics to track user activities on its website. Which role would be instrumental in setting up the infrastructure for real-time data processing?
 - A. Data engineer
 - B. Database administrator
 - C. Data analyst
 - D. None of the above
- 6. Your organization needs to ensure data is stored securely and can be recovered in the case of disasters. Which role is responsible for setting up data backups and recovery procedures?
 - A. Data engineer
 - B. Database administrator
 - C. Data analyst
 - D. All of the above
- 7. You are tasked with designing a database schema for a new project. Which role is typically responsible for this task?
 - A. Data engineer
 - B. Database administrator
 - C. Data analyst
 - D. All of the above
- 8. Your company is planning to implement a data warehouse for historical data analysis. Which role would be involved in selecting and configuring data warehousing services?
 - A. Data engineer
 - B. Database administrator
 - C. Data analyst
 - D. All of the above
- 9. Your organization needs to extract insights from unstructured text data. Which role is most suitable for this task?
 - A. Data engineer
 - B. Database administrator
 - C. Data analyst
 - D. All of the above

- 10. Your organization is expanding its use of cloud-based data services. Which role would likely be involved in selecting and implementing these cloud services?
 - A. Data engineer
 - B. Database administrator
 - C. Data analyst
 - D. All of the above
- 11. Your organization needs to store large volumes of structured data efficiently. Which file format is most suitable for this requirement?
 - A. JSON
 - B. XML
 - C. Parquet
 - D. ORC
- 12. Your team is tasked with handling semi-structured data such as log files from various sources. Which file format is designed for storing and querying semi-structured data?
 - A. JSON
 - B. XML
 - C. Parquet
 - D. Delimited text
- 13. Your company deals with extensive financial data and needs to ensure data consistency and eliminate data redundancy. What type of database is best suited for this requirement?
 - A. Relational database
 - B. Non-relational database
 - C. Graph database
 - D. Columnar database
- 14. Your organization operates in a highly regulated industry and must maintain transactional data integrity. Which type of database is essential for ensuring atomicity, consistency, isolation, and durability (ACID)?
 - A. Relational database
 - B. Non-relational database
 - C. Columnar database
 - D. Document database
- 15. Your company needs to store historical data for business analytics. Which data warehousing service in Azure is suitable for this requirement?
 - A. Azure SQL Database
 - B. Azure Cosmos DB

- C. Azure Synapse Analytics
- D. Azure Data Lake Storage

Thought experiment answers

This section contains the answers for the thought experiment. Each answer explains why the answer choice is correct.

1. **B** Database administrator

Explanation: Database administrators (DBAs) are essential for maintaining the integrity, security, and performance of the database system. They oversee access controls, implement encryption, schedule regular backups, and optimize queries to ensure the database functions reliably and efficiently.

2. A Data engineer

Explanation: Data engineers play a pivotal role in constructing data pipelines that retrieve data from various sources, transform it to meet specific requirements, and load it into data storage systems. Their expertise lies in maintaining data integrity, applying data transformation logic, and streamlining the flow of data for analysis. Data engineers ensure that the data is efficiently prepared for analytical processes so that businesses can derive meaningful insights from a diverse array of data sources.

3. C Data analyst

Explanation: Data analysts possess the analytical expertise to examine data, apply statistical methods, and generate reports. They can identify patterns in customer behavior, extract actionable insights, and present findings to support data-driven decision-making.

4. **B** Database administrator

Explanation: Database administrators specialize in fine-tuning database performance. They identify and resolve bottlenecks, optimize SQL queries, configure indexing, and manage system resources to ensure responsive database operations.

5. A Data engineer

Explanation: Data engineers are responsible for building the infrastructure needed for real-time data processing. They design data pipelines so users can efficiently capture, process, and analyze streaming data.

6. B Database administrator

Explanation: Database administrators are the primary custodians of data security and recovery. They establish robust backup and disaster recovery plans to safeguard data integrity and availability.

7. **D** All of the above

Explanation: Database design is often a collaborative effort involving data engineers, who design the schema structure; DBAs, who ensure it meets performance and security requirements; and data analysts, who provide input on data access needs.

8. A Data engineer

Explanation: Data engineers are responsible for selecting, configuring, and maintaining data warehousing services. They ensure data is stored and structured optimally for analytical purposes.

9. C Data analyst

Explanation: Data analysts are skilled in natural language processing (NLP) and text analysis. They can extract valuable insights from unstructured text data, such as customer reviews or social media comments.

10. **D** All of the above

Explanation: Cloud adoption involves multiple roles. Data engineers choose appropriate cloud-based data storage and processing solutions. DBAs ensure database integration and security in the cloud, and data analysts use cloud tools for analysis and reporting collaboration among all roles, which is crucial for a successful transition to cloud-based data services.

11. C Parquet

Explanation: Parquet is a columnar storage file format that excels at storing large volumes of structured data efficiently, making it an ideal choice for such scenarios.

12. A JSON

Explanation: JavaScript Object Notation (JSON) is a flexible and widely used format for storing and querying semi-structured data, making it suitable for log files and other similar data.

13. A Relational database

Explanation: Relational databases are known for their data consistency, structured schema, and ability to eliminate data redundancy, making them well-suited for financial data.

14. A Relational database

Explanation: Relational databases are known for their support of ACID properties, making them critical in scenarios where transactional data integrity is paramount.

15. C Azure Synapse Analytics

Explanation: Azure Synapse Analytics is a powerful data warehousing service designed for historical data storage and advanced analytics.

Index

Numerics

1NF (first normal form), 32 2NF (second normal form), 32 3NF (third normal form), 33 4NF (fourth normal form), 33 5NF (fifth normal form), 33

A

access tiers, Azure Blob storage, 80-81 ACID (atomicity/consistency/isolation/durability), 17, 18 ADX (Azure Data Explorer), 124 ALTER statement, 36 analytical workload, 17, 18-19, 110 columnar data stores, 111-112 NoSQL data store, 111 relational data stores, 110–111 time-series data stores, 112-113 analytics deep, 106-107 large-scale, 103, 104 data ingestion and processing, 107-110 data transformation, 106 data warehousing architecture, 104-105 practical scenario, 105-106 real-time data, 103, 117-118 API, Azure Cosmos DB, 90 Azure Table, 92 Cassandra, 91-92 Core (SQL), 90 Gremlin, 92 MongoDB, 91 append blob, 80 architecture Azure Data Lake Storage Gen2, 82-83 Azure File storage, 84-85 Azure SQL Managed Instance, 51 large-scale data warehousing, 104-105 Archive access tier, 81

Azure Analysis Services, 106 Azure Blob storage, 5, 74 access tiers, 80-81 container, 79 storage account, 75-79 key features, 75-78 pricing tiers, 79 types of blobs, 80 Azure Cosmos DB, 87-88, 106 API/s, 90 Azure Table, 92 Cassandra, 91-92 Core (SQL), 90 Gremlin, 92 MongoDB, 91 use cases, 88-89 Azure Data Factory, 106, 115 Azure Data Lake Storage Gen2, 81 architecture, 82-83 key enhancements, 81-82 Azure Database for MariaDB, 58-59 business benefits, 59-60 use cases, 60 Azure Database for MySQL, 56-57 business benefits, 57-58 compliance offerings, 57 use cases, 58 Azure Database for PostgreSQL, 60-61 business benefits, 61 use cases, 62 Azure Databricks, 114 Azure File storage, 83-84 architecture, 84-85 integration with Azure File Sync, 85 key features, 84 security and access control, 85-86 use cases, 86 Azure File Sync, integration with Azure File storage, 85

Avro, 11–12

Azure HDInsight

Azure HDInsight, 114 Azure SQL Database, 45, 47, 106 business benefits, 48 use cases, 47-48 Azure SQL Edge, 46, 49-50 Azure SQL family, comparing services, 46-47 Azure SQL Managed Instance, 45, 48-49, 50-51 architectural components, 51 example scenario, 52-53 key features, 51–52 use cases, 52 Azure SQL Server on Virtual Machines, 45, 50 Azure Stream Analytics, 122–124 Azure Synapse Analytics, 46, 49, 106, 113 Azure Table API, 92 Azure Table storage, 86-87

B

backup and recovery, 21 batch data, 118-119 BCNF (Boyce-Codd normal form), 33 **BEGIN TRANSACTION statement**, 41 big data, 10, 12 blob storage, 74 access tiers, 80-81 container, 79 storage account key features, 75-78 pricing tiers, 79 types of blobs, 80 block blob, 80 Boyce, Raymond F., 34 building blocks, Azure Table storage, 86-87 business benefits Azure Database for MariaDB, 59-60 Azure Database for MySQL, 57–58 Azure Database for PostgreSQL, 61 Azure SQL Database, 48 business executive, 19

C

calculated column, 134 Cassandra API, 91–92 Chamberlin, Donald D., 34 cloud database, 15–16 provider, 16 Codd, Edgar F., 1 columnar data stores, 111–112 COMMIT statement, 41 containers, blob storage, 79 Cool access tier, 81 Core (SQL) API, 90 CREATE INDEX statement, 43 CREATE statement, 36 CSV (comma-separated values) file format, 6–7

D

data, 1 analyst, 19, 22-23 batch, 118-119 categorization, 135-136 durability, 18 engineer, 19, 21-22 file format, 6 Avro, 11-12 delimited, 6-7 JSON, 7-8 ORC, 12-13 Parquet, 10-11 XML, 9-10 ingestion, 21, 107-110 integrity, 21, 31 pipeline, 108-109 scientist, 19 semi-structured, 3-4 storage, 5-6 streaming, 119-120 structured, 2-3, 30 unstructured, 4-5 variety, 108 velocity, 108 views, 42 visualization, 22, 127 volume, 108 warehousing Azure Data Factory, 115 Azure Databricks, 114 Azure HDInsight, 114

Azure Synapse Analytics, 113 Microsoft Fabric, 115-117 workload, 17 analytical, 17, 18-19 roles and responsibilities, 19-20 transactional, 17, 18 database, 13-14 administrator, 20-21 Azure Cosmos DB, 87-88 APIs, 90 Azure Table API, 92 Cassandra API, 91-92 Core (SQL) API, 90 Gremlin API, 92 MongoDB API, 91 use cases, 88-89 cloud, 15-16 index, 42-43 non-relational, 14-15 NoSQL, 4 objects, 30, 42 relational, 1, 3, 14, 30, 31-32 keys, 31 normalization, 31, 32-34 relationships, 31 tables, 14, 30-31 table, 42 DBA (database administrator), 20-21 DCL (Data Control Language), 39 DENY statement, 40 GRANT statement, 39-40 **REVOKE statement**, 40 DDL (Data Definition Language) ALTER statement, 36 CREATE statement, 36 DROP statement, 36-37 **TRUNCATE statement**, 37 deep analytics, 106–107 **DELETE statement**, 38 delimited file format, 6–7 DENY statement, 40 dialects, SQL (Structured Query Language), 35 DMBS (database management system), 14 DML (Data Manipulation Language), 37 **DELETE statement**, 38 **INSERT** statement, 38 MERGE statement, 39 SELECT statement, 38 **UPDATE statement**, 38

DP-900 Microsoft Azure Data Fundamentals exam, updates, 149–153 DROP statement, 36–37

E-F

file format, 6 Avro, 11–12 delimited, 6–7 JSON, 7–8 ORC, 12–13 Parquet, 10–11 XML, 9–10

G

GRANT statement, 39–40 Gremlin API, 92 GRS (geo-redundant storage), 109

Η

hierarchies, Power BI, 132–134 hopping window, 123 Hot access tier, 81

I-J

IaaS (Infrastructure-as-a-Service), SQL Server on Azure VMs, 49 IBM, 34 index, database, 42–43 INSERT statement, 38 insight, 19 IoT Hub, 106

JSON (JavaScript Object Notation), 3, 7-8

Κ

key features Azure File storage, 84 Azure SQL Managed Instance, 51–52 SQL Server on Azure VMs, 54 storage account, 75–78 keys, 31, 130

L

large-scale analytics, 103, 104 data ingestion and processing, 107–110 data transformation, 106 data warehousing architecture, 104–105 practical scenario, 105–106

Μ

MariaDB, 56. See also Azure Database for MariaDB measure, 134 MERGE statement, 39 Microsoft Azure SQL Database, 16 Microsoft Fabric, 115–117 Microsoft Power BI. See Power BI MongoDB API, 91 MySQL, 35, 56. See also Azure Database for MySQL

Ν

non-relational database, 14–15 normalization, 30, 31, 32–34 challenges of, 34 when to use, 34 NoSQL, 73. *See also* unstructured data data store, 111 database, 4, 14–15

О-Р

ORC (optimized row columnar) file format, 12-13

page blob, 80 Parquet, 10–11 pipeline, data, 108–109 PL/SQL (Procedural Language/SQL), 35 PostgreSQL, 35, 56. *See also* Azure Database for PostgreSQL Power BI, 103, 107 capabilities, 128–129 data categorization, 135–136 data models, relationships, 129–134 data visualization, 127 hierarchies, 132–134 measures and calculated columns, 134–135 Quick Measures, 136–138 schema, 131–132 pricing tiers, storage account, 79 provider, cloud, 16 PSV (pipe-separated values) file format, 7

Q-R

query, SQL, 14 Quick Measures, 136–138

RDBMS (relational database management system), 56 MariaDB, 56 MySQL, 56 PostgreSQL, 56 real-time data analytics, 103, 117-118 Microsoft cloud services ADX (Azure Data Explorer), 124 Azure Stream Analytics, 122–124 Spark Structured Streaming, 125–126 relational data stores, 110-111 relational database, 1, 3, 14, 30, 31-32 Azure, 44-45 keys, 31 normalization, 31, 32-34 challenges of, 34 when to use, 34 relationships, 31 tables, 14, 30-31 relationships, 31, 130-131 responsibilities data analyst, 22–23 data engineer, 21-22 DBA (database administrator), 20-21 **REVOKE statement**, 40 **ROLLBACK statement**, 41

S

SAVEPOINT statement, 41 schema Avro, 11 snowflake, 131–132 star, 131 security Azure File storage, 85–86 cloud database, 16 SELECT statement, 38 semi-structured data, 3-4 sliding window, 123 snowflake schema, 131-132 Spark Structured Streaming, 125–126 SQL (Structured Query Language), 2, 30, 31, 34-35 dialects, 35 query, 14 statement/s DCL DENY, 40 DCL GRANT, 39-40 DCL REVOKE, 40 DDL (Data Definition Language), 35-37 DDL ALTER, 36 DDL CREATE, 36 DDL DROP, 36-37 DDL TRUNCATE, 37 DML (Data Manipulation Language), 37 DML DELETE, 38 DML INSERT, 38 DML MERGE, 39 DML SELECT, 38 DML UPDATE, 38 TCL BEGIN TRANSACTION, 41 TCL COMMIT, 41 TCL ROLLBACK, 41 TCL SAVEPOINT, 41 SQL Server on Azure VMs, 49, 53 example scenario, 55 key aspects, 53-54 key features, 54 use cases, 54-55 star schema, 131 statement/s CREATE INDEX, 43 DCL (Data Control Language), 39 DENY, 40 GRANT, 39-40 **REVOKE**, 40 DDL (Data Definition Language) ALTER, 36 CREATE, 36 DROP, 36-37 RENAME, 37 TRUNCATE, 37 DML (Data Manipulation Language), 37 DELETE, 38 INSERT, 38

MERGE, 39 SELECT, 38 UPDATE, 38 stored procedure, 43-44 TCL (Transaction Control Language), 40 **BEGIN TRANSACTION, 41** COMMIT, 41 ROLLBACK, 41 SAVEPOINT, 41 storage, 5-6, 74 account, 79 Azure Data Lake Storage Gen2, 81 architecture, 82-83 key enhancements, 81-82 Azure File, 83-84 architecture, 84-85 integration with Azure File Sync, 85 key features, 84 security and access control, 85-86 use cases, 86 Azure Table, 86 building blocks, 86-87 blob, 74 access tiers, 80-81 container, 79 storage account, 75-79 types of blobs, 80 geo-redundant, 109 stored procedure, 43-44 streaming data, 119-120 anomaly detection, 123 time windowing, 123 structured data, 2-3, 30

Т

table/s, 14, 30–31, 42. See also Azure Table storage keys, 130
normalization, 33–34
relationships, 130–131
schema, 131–132
TCL (Transaction Control Language), 40
BEGIN TRANSACTION statement, 41
COMMIT statement, 41
ROLLBACK statement, 41
SAVEPOINT statement, 41
TCO (total cost of ownership), Azure SQL Database, 48

time windowing

time windowing, 123 time-series data stores, 112–113 transactional workload, 17, 18 triggers, 44, 108 TRUNCATE statement, 37 T-SQL (Transact-SQL), 35 TSV (tab-separated values) file format, 7 tumbling window, 123

U

unstructured data, 4–5 UPDATE statement, 38 updates, DP-900 Microsoft Azure Data Fundamentals exam, 149–153 use cases, 62 Azure Cosmos DB, 88–89 Azure Database for MariaDB, 60 Azure Database for MySQL, 58 Azure File storage, 86 Azure SQL Database, 47–48 Azure SQL Managed Instance, 52 SQL Server on Azure VMs, 54–55

V

variety, data, 108 velocity, data, 108 view, 42 volume, data, 108

W

warehousing, large-scale data, 104–105 workload, 17 analytical, 17, 18–19, 110 columnar data stores, 111–112 NoSQL data store, 111 relational data stores, 110–111 time-series data stores, 112–113 transactional, 17, 18

X-Y-Z

XML (eXtensible Markup Language), 3, 9-10