

Databricks Lakehouse Platform Cookbook

*100+ recipes for building a scalable
and secure Databricks Lakehouse*

Dr. Alan L. Dennis



www.bpbonline.com

Copyright © 2024 BPB Online

All rights reserved. No part of this book may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, without the prior written permission of the publisher, except in the case of brief quotations embedded in critical articles or reviews.

Every effort has been made in the preparation of this book to ensure the accuracy of the information presented. However, the information contained in this book is sold without warranty, either express or implied. Neither the author, nor BPB Online or its dealers and distributors, will be held liable for any damages caused or alleged to have been caused directly or indirectly by this book.

BPB Online has endeavored to provide trademark information about all of the companies and products mentioned in this book by the appropriate use of capitals. However, BPB Online cannot guarantee the accuracy of this information.

First published: 2024

Published by BPB Online

WeWork

119 Marylebone Road

London NW1 5PU

UK | UAE | INDIA | SINGAPORE

ISBN 978-93-55519-566

www.bpbonline.com

Dedicated to

My loving and supportive wife,

Kim

Foreword

There is no denying that data is the lifeblood of industry. Everyone understands that businesses that harness their data well, will thrive while the ones that do not, will fall by the wayside. Choosing the correct platform for your data estate is, perhaps, the most critical decision a business can make. The second most important factor is, of course, to hire the right team to build on this chosen platform.

2023 has shown us that GenAI is the future and I am convinced that there is only one data platform that is natively suited for the demands that GenAI will impose on the data estates of the future: the Lakehouse built on Databricks. Over the last 10 years, Databricks has transformed itself from a niche Spark-focused shop to a visionary organization building a holistic data platform that can support Analytics, Data Integration (a fancier term for ETL or ELT) and AI/ML: including the now white-hot GenAI use cases. While a lot of competing data platforms and data clouds make similar claims, there is only one platform, that is, Databricks: that has been doing cloud-native, multi-language data processing at scale: for 10 years now.

I consider myself fortunate that I discovered and fell in love with Apache Spark in 2014 and then got a chance to join Databricks in 2019. I am truly grateful that I got to pick the brains of some of the smartest minds in the universe during that time. The author of this book, Dr. Alan Dennis, is one such individual and it is my honor to count him as a friend and a business partner.

Whether you are a seasoned data professional or someone looking to unlock the potential of data in your organization, this book is your key to a data revolution. Databricks and the Lakehouse paradigm offer a fresh perspective on data management, and this book is your guide to mastering it. Alan's recipe-driven approach to teaching is perfect for the real world: it will enable you to deliver results quickly in the immediate term: and will help you connect the dots and build a strong foundation for self-learning over the longer term. The context-setting sections give you a quick history of features and approaches over the years. It will help you appreciate how the platform has evolved and most importantly, help you avoid old pitfalls and anti-patterns.

Apache Spark, Databricks and Lakehouse have transformed my life for the better: and it is my sincere hope and best wishes that you, the reader, have a similar fulfilling experience. Onward!

Subramanian Iyer

Principal, Speedboat Professional Services

Award-winning Architect and Certified Instructor on Databricks and Lakehouse

Ex-Brickster (2019-2023) and Spark fanboy since 2014

About the Author

Dr. Alan L. Dennis has been writing software for over 30 years. His experiences range from being one of the first employees at a startup to leading a team of over twenty developers. He has held titles such as Programmer, Architect, Chief Technical Officer, and member of technical staff. He has worked for many Fortune 50 companies, with a wide range of industry experience.

He holds a Doctorate in Computer Science with a concentration in Big Data Analytics, a Master's in Computer Science with specialization in Artificial Intelligence, and a Bachelor's of Business Administration with focus on Computer Information Systems. He teaches graduate classes at several universities and is a Databricks Certified Trainer.

About the Reviewers

- ❖ **Jay Kalathia** is an experienced Senior Software Engineer, with over two decades of experience in designing, developing, and optimizing cloud-based solutions on the Azure platform. He has a diverse background with proficiency in various programming languages such as Python, C#, JavaScript, and more. Jay is skilled in building infrastructure as code, developing CI/CD pipelines on Azure DevOps, and working on cloud-native solutions and tools including Azure, AWS, AKS, Kubernetes, and Terraform. Additionally, he has extensive experience with Azure Databricks and other cloud-based Big Data solutions. Jay is also a learner; from taking courses to completing certifications to stay up to date with latest technology trends.
- ❖ **Mahesh Das** is a Technology Evangelist and a Databricks Certified Data Science Professional. He is currently focused on the development of cloud-native solutions for ingesting and refining data from various sources using Azure Cloud, Databricks, Azure Data Factory, and Terraform scripts. Mahesh is actively engaged in projects implementing Large Language Models for diverse applications and remains a dedicated learner in the field of Machine Learning and Artificial Intelligence. With over 18 years of experience, he has contributed to numerous significant IT transformation projects for Fortune 500 clients across a range of industries, including Manufacturing, FMCG, oil and gas, Managed Print Services, and metals and mining. Mahesh possesses additional skills in AWS Cloud and SAP Master Data Governance, covering various functional domains such as sales and distribution, finance, and manufacturing. He is not only an avid reader of cutting-edge AI technologies but also serves as a technical reviewer for books within the same domain.

Acknowledgement

Many people have had a hand in this book. First, I would like to thank my parents for always supporting and encouraging me. They taught me to figure out how things work and, ideally, put them back together again afterward. My mother passed before this book could be completed. She will be missed, but her Heavenly Father called, and she went. I would also like to thank my wife Kim, for supporting and encouraging me and understanding when I disappear for weeks on end.

I would like to thank Jay, and Mahesh for their input to the book. They provided valuable feedback throughout the process. I would also like to thank Subramanian for his kind words in the foreword.

Lastly, I am incredibly thankful for you, gentle reader. I wrote this for you and I hope you find value in it. There is something for everyone in this book, let me know what you find.

Preface

It is commonly understood that valuable insights can be found in an organization's data. One way to extract that value is to construct a Data Lakehouse. This book helps you create a Lakehouse on the Databricks Platform. It is the culmination of decades of data processing design and implementation.

It is not easy to create a data ecosystem. There are many competing priorities and technical challenges. This book walks you through the process, providing hands-on examples. We organize those steps into recipes. This keeps the author from waxing on about theory and helps the reader find the information needed in a given situation. We cover the theory behind the approaches used and guide the reader to avoid common pitfalls.

We start with the basics, such as explaining what a Databricks Lakehouse is, why we need them, and what value it brings. We move on to applying the concepts in practice. Part of the reason for constructing a Data Lakehouse is to enable users to access its data. We then discuss the various personas that benefit from a Databricks Lakehouse.

While we start with the fundamentals, we rapidly move on to more advanced topics. A good understanding of SQL, Python, Spark, and cloud computing would benefit the reader but is not required.

Chapter 1: Introduction to Databricks Lakehouse – This chapter provides a brief history of Big Data, Spark, and Databricks. It introduces the reader to the community edition of Databricks as a starting point for using Databricks. We discuss why we construct a Lakehouse and present the overall architecture. We provide clear definitions for each of the layers of a Databricks Lakehouse. We discuss design considerations and compare Lakehouses to other data technologies.

Chapter 2: Setting-up a Databricks Workspace – This chapter presents the information necessary to provision and effectively use a Databricks environment. This includes examining core Databricks concepts, service tier selection, and cloud selection considerations. Deployment details are examined, including those with long-lasting implications. Access control and other configurations are discussed, along with the types of clusters and performance levels.

Chapter 3: Connecting to Storage – This chapter covers the approaches and tradeoffs to connect to storage. The Databricks File System is discussed in detail as it is an important element of the Lakehouse platform. The background of the file system is reviewed, and

various ways of connecting to storage are explored. The approaches to Lakehouse design are presented, with recommendations on how to organize a Lakehouse. Recommendations are provided regarding the documentation of allowed operations. Recipes containing various examples of connecting to Azure storage systems are provided.

Chapter 4: Creating Delta Tables – This chapter describes how to construct a Delta Lake, including a discussion of managed and external tables. Guidance is provided to help decide which type of table to create. Examples are provided of creating tables using SQL and the Spark API. Core concepts such as secret scopes are discussed, along with example of creating tables from AWS S3, GCP buckets, and Azure ADLS.

Chapter 5: Data Profiling and Modeling in the Lakehouse – This chapter examines two of the more important activities when constructing a Data Lakehouse. Various ways of performing profiling are examined, including Databricks' native Data Profile feature. Discussion of the Databricks Describe and Summary features are included, along with analysis at scale using `ydata_profiling`.

Chapter 6: Extracting from Source and Loading to Bronze – This chapter covers the first step in refining data. A discussion is presented regarding using the raw zone or skipping it and going from source to bronze. Several ways of incrementally ingesting data are presented, which is essential for a high-performance Databricks Lakehouse. These methods include self-managed watermarks, Auto Loader, Delta Live Tables, and streaming data.

Chapter 7: Transforming to Create Silver – This chapter continues the refinement journey, picking up data at the Bronze layer and moving it to Silver. Both incremental and full refinement are discussed. Several approaches to processing are discussed, including the importance of data quality rules and expectations. Common Silver-to-Silver operations are discussed, including denormalization, JSON exploding, and projection reshaping.

Chapter 8: Transforming to Create Gold for Business Purposes – This chapter continues the discussion of refining data, with the goal of answering business questions. Gold tables are built to answer a specific question. The sources for Gold tables are discussed, with implementations in PySpark and Delta Live Tables. As Gold tables are optimized for consumption, a brief discussion of support-related operations such as vacuuming and optimizing tables is present.

Chapter 9: Machine Learning and Data Science – Data scientists are common users of the Databricks Lakehouse. We examine using Machine Learning in Databricks, and the use of AutoML. Next, we discuss MLflow, and the importance it plays in deploying models to production. Lastly, we briefly discuss the Databricks feature store.

Chapter 10: SQL Analysis – SQL is one of the most widely known languages. We discuss the SQL Analysis features built into Databricks, including Databricks SQL. We show how to create and manage a SQL Warehouse. We discuss the usage of the SQL Editor and use it to write common queries. We create dashboards and alerts using those queries. We close with a discussion of cost and performance considerations.

Chapter 11: Graph Analysis – There are many ways to perform analysis; one way is to use mathematical graph algorithms. We discuss the nature of graphs and when using graph algorithms is appropriate. We discuss GraphX and GraphFrames, along with the operations they enable and associated algorithms. Lastly, we discuss reading data from Neo4J’s AuraDB from Databricks.

Chapter 12: Visualizations – There are many ways to present data; visualizations can be very powerful. We discuss visualization best practices and how to create a Databricks dashboard. We also discuss native visualization support within a Databricks notebook. We conclude the chapter by discussing the use of Power BI with Databricks.

Chapter 13: Governance – Without proper governance, a Databricks Lakehouse will not be successful. We discuss the role of data governance and the use of Databricks’ Unity Catalog. We walk through the installation and usage of Unity Catalog and review the major benefits. We discuss the steps to install and use Azure Purview in combination with Unity Catalog.

Chapter 14: Operations – This chapter covers the steps necessary to keep a Lakehouse working effectively, including source code management and orchestration. Preventive scheduled maintenance can help avoid unacceptable processing time and outages. We also discuss how to manage and maintain visibility of costs.

Chapter 15: Tips, Tricks, Troubleshooting, and Best Practices – This final chapter contains important elements that did not make it into other parts of the book. We revisit ingesting data, by ingesting relational data. Discuss performance optimizations such as using pools. We discuss how to orchestrate notebooks. Lastly, we conclude with a discussion of best practices and guiding principles.

Code Bundle and Coloured Images

Please follow the link to download the
Code Bundle and the *Coloured Images* of the book:

<https://rebrand.ly/llidt00>

The code bundle for the book is also hosted on GitHub at

<https://github.com/bpbpublications/Databricks-Lakehouse-Platform-Cookbook>.

In case there's an update to the code, it will be updated on the existing GitHub repository.

We have code bundles from our rich catalogue of books and videos available at **<https://github.com/bpbpublications>**. Check them out!

Errata

We take immense pride in our work at BPB Publications and follow best practices to ensure the accuracy of our content to provide with an indulging reading experience to our subscribers. Our readers are our mirrors, and we use their inputs to reflect and improve upon human errors, if any, that may have occurred during the publishing processes involved. To let us maintain the quality and help us reach out to any readers who might be having difficulties due to any unforeseen errors, please write to us at :

errata@bpbonline.com

Your support, suggestions and feedbacks are highly appreciated by the BPB Publications' Family.

Did you know that BPB offers eBook versions of every book published, with PDF and ePub files available? You can upgrade to the eBook version at www.bpbonline.com and as a print book customer, you are entitled to a discount on the eBook copy. Get in touch with us at :

business@bpbonline.com for more details.

At **www.bpbonline.com**, you can also read a collection of free technical articles, sign up for a range of free newsletters, and receive exclusive discounts and offers on BPB books and eBooks.

Piracy

If you come across any illegal copies of our works in any form on the internet, we would be grateful if you would provide us with the location address or website name. Please contact us at **business@bpbonline.com** with a link to the material.

If you are interested in becoming an author

If there is a topic that you have expertise in, and you are interested in either writing or contributing to a book, please visit **www.bpbonline.com**. We have worked with thousands of developers and tech professionals, just like you, to help them share their insights with the global tech community. You can make a general application, apply for a specific hot topic that we are recruiting an author for, or submit your own idea.

Reviews

Please leave a review. Once you have read and used this book, why not leave a review on the site that you purchased it from? Potential readers can then see and use your unbiased opinion to make purchase decisions. We at BPB can understand what you think about our products, and our authors can see your feedback on their book. Thank you!

For more information about BPB, please visit **www.bpbonline.com**.

Join our book's Discord space

Join the book's Discord Workspace for Latest updates, Offers, Tech happenings around the world, New Release and Sessions with the Authors:

<https://discord.bpbonline.com>



Table of Contents

1. Introduction to Databricks Lakehouse	1
Introduction.....	1
Structure.....	1
Objectives	2
Background	2
Brief history of Big Data, Spark, and Databricks	2
Databricks community edition.....	3
<i>Recipe 1: Signing up for the Databricks community edition.....</i>	<i>4</i>
<i>Recipe 2: Creating a notebook in the Databricks Community edition</i>	<i>5</i>
<i>Recipe 3: Changing a notebook's default language</i>	<i>7</i>
<i>Recipe 4: Create a table from CSV using SQL.....</i>	<i>7</i>
<i>Recipe 5: Query a table using SQL.....</i>	<i>8</i>
<i>Recipe 6: Examine a table's structure</i>	<i>9</i>
<i>Recipe 7: Use infer schema on CSV in SQL</i>	<i>9</i>
<i>Recipe 8: Compute mean in group by in SQL.....</i>	<i>10</i>
<i>Recipe 9: Importing a notebook</i>	<i>11</i>
<i>Recipe 10: Exporting a notebook in Databricks Community Edition.....</i>	<i>12</i>
Data Lakehouse value proposition	13
Lakehouse architecture.....	14
<i>Separation of computing and storage.....</i>	<i>15</i>
<i>Data lake.....</i>	<i>15</i>
<i>Delta Lake.....</i>	<i>15</i>
<i>Computational engine</i>	<i>16</i>
Design considerations.....	16
<i>Extraction and storage by system</i>	<i>16</i>
<i>Zones and their definitions.....</i>	<i>16</i>
<i>Source.....</i>	<i>17</i>
<i>Bronze</i>	<i>17</i>
<i>Silver</i>	<i>17</i>
<i>Gold.....</i>	<i>17</i>
Lakehouse compared to other data technologies	18

Extract load transform and extract transform load	18
Compared to traditional data lake approaches.....	18
Differences from Lambda architecture.....	19
Conclusion.....	19
Points to remember	19
2. Setting-up a Databricks Workspace	21
Introduction.....	21
Structure.....	21
Objectives	21
Core Databricks concepts	22
Databricks service tiers	22
Brief introduction of Databricks features	22
Machine Learning	22
Notebook access control.....	23
Databricks SQL and endpoints	23
Internet protocol addresses access control.....	23
Databricks pricing model	23
Pick your cloud.....	23
AWS	24
Azure	28
Google Cloud Platform.....	31
Deployment details	36
Public availability.....	36
Network size	36
Network peering.....	36
Initial configuration	37
Access control.....	37
Cluster types.....	38
All purpose.....	38
Job clusters	38
Cluster creation details.....	38
Single or multiple nodes.....	38
Access mode	38

Choosing performance level	39
Conclusion.....	40
3. Connecting to Storage	41
Introduction.....	41
Structure.....	41
Objectives	41
Databricks file system.....	42
Using mount points	42
Recipe 11: Using the DBFS file browser	42
Recipe 12: Using Databricks' web terminal.....	45
Recipe 13: Using Databricks Utilities' file system methods	46
The importance of DBFS	48
Lakehouse design.....	48
Source to Silver	48
Including raw.....	50
The document allowed operations crossing layers	50
Source to raw.....	51
Source to bronze	52
Raw to bronze.....	52
Bronze to silver.....	53
Silver to silver	53
Silver to gold and gold to gold.....	54
Recipes 14: Using the Lakehouse layer presentation.....	54
Azure.....	55
ADLS Gen2.....	55
Credential passthrough	55
Recipe 15: Creating a storage account for ADLS Gen2.....	56
Recipe 16: Creating a container and setting ACLs.....	57
Recipe 17: Using Passthrough authentication.....	60
Key vault and secret scope.....	62
Recipe 18: Link a key vault to a secret scope	62
Recipe 19: Displaying a redacted value.....	64
Blob storage.....	64

<i>Recipe 20: Account keys</i>	65
<i>Recipe 21: Service principle</i>	66
<i>Recipe 22: Shared access signatures</i>	72
Conclusion.....	75
4. Creating Delta Tables	77
Introduction.....	77
Structure.....	77
Objectives	77
Delta Lake.....	78
<i>Managed and unmanaged tables</i>	78
<i>Deciding table type</i>	78
<i>Schema and database</i>	79
Creating managed Delta tables	81
Ways to create	81
<i>Recipe 23: Upload data using Databricks workspace</i>	81
SQL	87
<i>Recipe 24: Reading the SQL language reference</i>	87
<i>Recipe 25: Creating a table with SQL</i>	90
<i>Recipe 26: Creating a table with SQL using AS</i>	91
Spark API.....	92
<i>Recipe 27: Creating a table using Spark API and random data</i>	92
<i>Recipe 28: Examining table history</i>	94
Managed tables details	96
<i>Recipe 29: Managed Delta table details</i>	97
<i>Recipe 30: Using Data Explorer to see table details</i>	99
Creating unmanaged tables.....	100
<i>Recipe 31: Using Databricks CLI to create a secret scope</i>	100
<i>Recipe 32: Accessing S3 from Databricks on AWS</i>	103
<i>Recipe 33: Creating an external Delta table in SQL on AWS</i>	104
<i>Recipe 34: Creating an external table in PySpark on AWS</i>	107
<i>Recipe 35: Creating an external Delta table in SQL on Azure</i>	109
<i>Recipe 36: Creating an external table with Python on Azure</i>	111
<i>Recipe 37: Accessing GCP buckets from Databricks</i>	112

<i>Recipe 38: Creating an external Delta table in SQL on GCP</i>	114
<i>Recipe 39: Creating an external Delta table in Python on GCP</i>	116
Conclusion.....	118
5. Data Profiling and Modeling in the Lakehouse	119
Introduction	119
Structure	119
Objectives	119
Data profiling.....	120
<i>Recipe 40: Using Azure Data Factory to ingest raw</i>	120
<i>Recipe 41: Reorganize files</i>	129
<i>Recipe 42: Creating tables from a directory programmatically</i>	130
<i>Recipe 43: Data profiling using Databricks native functionality</i>	131
<i>Recipe 44: Listing row counts for all</i>	133
<i>Recipe 45: Using DBUtils summarize</i>	134
<i>Recipe 46: Using a DataFrames describe and summary methods</i>	136
<i>Recipe 47: Descriptive data analysis with Pandas profiling</i>	139
Data modeling	142
Common modeling approaches.....	142
Entity-relationship data modelling.....	143
Star schema	143
Snowflake schema.....	144
Standardized data models	145
Retrieval optimized models.....	145
Design approach	146
Conclusion	147
6. Extracting from Source and Loading to Bronze	149
Introduction	149
Structure	149
Objectives	150
To raw or not to raw	150
Using change data feed.....	150
Overview of change data feed	150
<i>Recipe 48: Creating a table with change data feed on</i>	151

Recipe 49: Using Python to enable CDF.....	151
Recipe 50: Ensure CDF is enabled for all tables.....	154
Loading files using self-managed watermarks.....	155
Incremental ingestion example.....	155
Recipes 51: Using incremental load of files	157
Recipes 52: Convert Event Hub data to JSON	167
Recipes 53: Full load of files	170
Loading files using Auto Loader	172
Auto Loader overview.....	172
Recipe 54: Incremental ingestion of files Avro using Auto Loader in Python.....	172
Recipe 55: Incremental ingestion of CSV files using Auto Loader in Python	174
Loading files using Delta Live Tables	177
Delta Live Tables overview	177
Recipe 56: Using the DLT SQL API to ingest JSON	178
Recipe 57: Incremental ingestion using DLT using Python API	187
Recipe 58: Full ingestion using DLT using SQL API.....	188
Recipe 59: Full ingestion using DLT using Python API.....	189
Loading streaming data	190
Recipe 60: Parameterizing pipelines	190
Recipe 61: Stream processing with DLT Python API.....	191
Recipe 62: Using Spark structured streaming	197
Conclusion	202
7. Transforming to Create Silver.....	203
Introduction	203
Structure	203
Objectives	204
Bronze to silver	204
Incremental refinement.....	204
Recipe 63: Incremental refinement using Delta Live Tables	204
Recipe 64: Incremental refinement using PySpark	206
Full refinement	212
Recipe 65: Full update refinement using Delta Live Tables	212
Recipe 66: Full refinement using PySpark.....	214

Data quality rules.....	216
Recipe 67: Using expectations in DLT with SQL.....	216
Recipe 68: Using expectations in DLT with PySpark.....	218
Silver to silver	219
Reshaping projection	220
Recipe 69: Projection reshaping using Python	220
Recipe 70: Projection reshaping using Delta Live Tables.....	221
Splitting tables	222
Recipe 71: Splitting table into multiple in PySpark	222
Recipe 72: Splitting table into multiple in Delta Live Tables	223
Enrichment.....	225
Recipe 73: Creating lookup data from telemetry	225
Recipe 74: Combining tables using DLT.....	227
Conclusion	228
8. Transforming to Create Gold for Business Purposes	229
Introduction	229
Structure	229
Objectives	230
Silver to gold.....	230
Aggregation	230
Recipe 75: Aggregation in Delta Live Tables	230
Dimensional tables using PySpark	232
Recipe 76: Creating a time dimension.....	232
Recipe 77: Creating a dimension from telemetry	233
Recipe 78: Creating a fact table from telemetry.....	235
Dimensional tables in Delta Live Tables	237
Recipe 79: Dimensional models with Delta Live Table	237
Using Common Data Models with Delta Live Tables.....	240
Microsoft Common Data Model	240
Gold to gold	240
Table optimization for consumption	240
Optimize	242
Recipe 80: Manually optimize a table	242

<i>Vacuum</i>	242
<i>Recipe 81: Vacuum a Delta table</i>	243
Conclusion	244
9. Machine Learning and Data Science	245
Introduction	245
Structure	245
Objectives	246
Machine Learning in Databricks.....	246
Using AutoML	246
<i>Recipe 82: Creating an ML cluster</i>	247
<i>Recipe 83: Importing data with the Databricks web page</i>	248
<i>Recipe 84: Creating and running an AutoML experiment</i>	249
Setting up and using MLflow.....	253
<i>Recipe 85: Setting up an MLflow experiment</i>	253
<i>Recipe 86: Using MLflow for non-ML workflows</i>	255
Deploying models to production.....	258
<i>Recipe 87: Registering a model</i>	258
<i>Recipe 88: Using a model for inference</i>	260
Using Databricks feature store	264
<i>Recipe 89: Importing an HTML notebook</i>	264
<i>Recipe 90: Basic interaction with Databricks Feature Store</i>	265
Conclusion	268
10. SQL Analysis	269
Introduction	269
Structure	269
Objectives	270
Databricks SQL.....	270
Creating and managing a SQL Warehouse.....	270
<i>Recipe 91: Creating a SQL Warehouse</i>	271
<i>Recipe 92: Connect to a SQL Warehouse from a Python Jupyter Notebook</i>	274
Using the SQL Editor.....	275
Writing queries	277
<i>Common interview queries</i>	278

Recipe 93: Show the contents of a table	278
Recipe 94: Select with filtered ordered limited result	279
Recipe 95: Aggregation of records	280
Recipe 96: Using grouping to find duplicate records	281
Recipe 97: Generating synthetic data	282
Recipe 98: Calculate rollups	284
Recipe 99: Types of joins	285
Inner	285
Left and right outer joins	286
Full outer join	287
Cross join	287
Creating dashboards	288
Recipe 100: Creating a quick dashboard	288
Recipe 101: Schedule dashboard refresh	291
Setting alerts	293
Recipe 102: Create a query for an alert	293
Recipe 103: Create an alert	293
Cost and performance considerations	298
Conclusion	300
11. Graph Analysis	301
Introduction	301
Structure	301
Objectives	302
What is a graph	302
When to use graph operations	303
GraphX	303
GraphFrames	304
Recipe 104: Creating a GraphFrame	304
Recipe 105: Using example graphs	306
Graph operations and algorithms	308
Recipe 106: Breadth-first search	308
Recipe 107: PageRank	309
Recipe 108: Shortest path	309

Recipe 109: Connected components	310
Recipe 110: Strongly connected components	311
Recipe 111: Label Propagation Algorithm	312
Recipe 112: Motif finding	314
Neo4J and Databricks	315
Recipe 113: Using AuraDB	315
Recipe 114: Reading Neo4J's AuraDB from Databricks	318
Conclusion	321
12. Visualizations	323
Introduction	323
Structure	323
Objectives	323
Visualization best practices	324
Visually appealing	324
Keep it simple	324
Explain unfamiliar graph types	324
Follow conventions	324
Tell a story	325
Databricks dashboards	325
Recipe 115: Importing sample dashboards	325
Recipe 116: Data preparation for a new dashboard	329
Recipe 117: Creating a dashboard	332
Visualizations in Databricks notebooks	341
Recipe 118: Using visualizations in notebooks	341
Power BI	344
Recipe 119: Connecting Power BI to Databricks	344
Conclusion	348
13. Governance	349
Introduction	349
Structure	349
Objectives	349
Role of data governance	350
Using Unity Catalog	350

<i>Recipe 120: Configuring Unity Catalog in Azure</i>	350
<i>Creating storage</i>	350
<i>Create a managed identity</i>	350
<i>Create Access Connector for Azure Databricks</i>	351
<i>Grant managed identity access</i>	355
<i>Creating a metastore</i>	356
<i>Unity Catalog object model</i>	360
<i>Recipe 121: Creating a new catalog</i>	361
<i>Recipe 122: Uploading data</i>	362
<i>Recipe 123: Creating a table</i>	363
Installing and using Purview	364
<i>Recipe 124: Installing Purview</i>	364
<i>Recipe 125: Connecting Purview to Databricks</i>	365
<i>Recipe 126: Scanning a Databricks workspace</i>	367
<i>Recipe 127: Browsing the Data Catalog</i>	371
Conclusion	372
14. Operations	373
Introduction	373
Structure	373
Objectives	373
Source code management and orchestration	374
<i>Recipe 128: Use GitHub with Databricks</i>	374
<i>Recipe 129: Create workflows to orchestrate processing</i>	379
<i>Recipe 130: Saving a Job JSON</i>	385
<i>Recipe 131: Use Airflow to coordinate processing</i>	387
Scheduled and ongoing maintenance	389
<i>Recipe 132: Repairing damaged tables</i>	389
<i>Recipe 133: Vacuum unneeded data</i>	393
<i>Recipe 134: Optimize Delta tables</i>	395
Cost management	395
<i>Recipe 135: Use cluster policies</i>	395
<i>Recipe 136: Using tags to monitor costs</i>	397
Conclusion	398

15. Tips, Tricks, Troubleshooting, and Best Practices	399
Introduction	399
Structure	399
Objectives	400
Ingesting relational data with Databricks	400
<i>Recipe 137: Loading data from MySQL.....</i>	<i>400</i>
<i>Recipe 138: Extending a Python class and reading using</i> <i>Databricks runtime format</i>	<i>402</i>
<i>Recipe 139: Caching DataFrames</i>	<i>403</i>
<i>Recipe 140: Loading data from MySQL using workers</i>	<i>405</i>
Performance optimization	407
<i>Using Databricks event log</i>	<i>407</i>
<i>Exploring the Spark UI jobs tab</i>	<i>408</i>
<i>Using the Spark UI SQL/DataFrame tab</i>	<i>411</i>
<i>Recipe 141: Using pools to improve performance</i>	<i>413</i>
Programmatic deployment and interaction	415
<i>Recipe 142: Creating a workspace with ARM Template</i>	<i>415</i>
<i>Recipe 143: Using the Databricks API.....</i>	<i>419</i>
Reading a Kafka stream	422
<i>Recipe 144: Creating a Kafka cluster.....</i>	<i>422</i>
<i>Recipe 145: Using confluent cloud.....</i>	<i>425</i>
Notebook orchestration	429
<i>Recipe 146: Running a notebook with parameters.....</i>	<i>429</i>
<i>Recipe 147: Conditional execution of notebooks.....</i>	<i>431</i>
Best practices.....	432
<i>Organize data assets by source until silver.....</i>	<i>432</i>
<i>Use automation as much as possibly.....</i>	<i>433</i>
<i>Use version control.....</i>	<i>433</i>
<i>Keep each step of the process simple</i>	<i>433</i>
<i>Do not be afraid to change.....</i>	<i>434</i>
Conclusion	434
Index	435-442

CHAPTER 1

Introduction to Databricks Lakehouse

Introduction

Welcome to our journey of learning and mastering the Databricks Lakehouse Platform. This is a hands-on book. While we will cover each topic's theoretical and technical foundations, you will have code to help you learn how to build a Lakehouse and succeed using Databricks.

Structure

In this chapter, we will cover the following topics:

- Background
- Brief history of Big Data, Spark, and Databricks
- Databricks community edition
- Data Lakehouse value proposition
- Lakehouse architecture
- Design considerations
- Lakehouse compared to other data technologies

Objectives

This chapter introduces nomenclature commonly used when discussing the Databricks Lakehouse Platform. By the end of the chapter, you should be able to describe a typical Lakehouse configuration and understand the architectural components and the value proposition of the lakehouse architecture.

Background

It is often important to understand a phenomenon's history that influenced its creation, and Data Lakehouse is no exception. We start with a brief history of Big Data, Spark, and Databricks. We briefly discuss the Databricks community edition and perform our first analysis in Databricks. We close this section by discussing the value proposition that drives the adoption of Data lakehouse, particularly Databricks Lakehouse.

Brief history of Big Data, Spark, and Databricks

When looking at how things came to be, we often discuss supporting and challenging forces. In the case of Big Data, several forces were driving its adoption. One key supporting force was the shift of the Internet from companies and government entities sharing information with their customers to users of platforms creating content in social media. Companies also learned that online sales had many advantages over traditional outlets, including lower operating costs. This shift generated vast amounts of data that previously was minimal. Online retailers learned that examining those log files could give insights into their customers that previously was not possible. The desire to process this information, which was too large to process with traditional file-processing approaches, led to the creation of a new set of technologies. Big Data was used to label these distributed, software-based fault-tolerant algorithms and technologies.

One of the early success stories of Big Data was Hadoop. Hadoop is a collection of open-source projects related to processing large, fast, or variant data. An early processing approach in Hadoop was called MapReduce. MapReduce was a framework that simplified the process of creating distributed solutions. Before Big Data frameworks like MapReduce, software developers coordinated activities between various workers attempting to

work together to solve a problem. Often, one or more of those workers would become unavailable. The software developer's job was to determine how to address this and many other challenges. With MapReduce, a developer was tasked with writing a few functions called by a framework to simplify the process. While MapReduce was a significant advancement, it was limited by its original design and purpose. MapReduce was focused on processing large or numerous files. Due to this design goal, it failed to support iteration and relied heavily on disk drives.

Spark was developed to address many of these challenges. Spark is a computational solution that relies on other technologies for storage. It also favors processing data in memory, resulting in significant performance improvements over MapReduce. Spark also enabled iteration during processing. These advancements lead to its rapid adoption and increasing popularity. Many of the creators of Spark formed Databricks in 2013. Databricks is a cloud company supporting the major cloud vendors. In 2017 Azure Databricks was announced. This partnership was notable because of the high integration between Azure and Databricks.

A data lakehouse is an architecture that combines the best elements of data lakes to address data warehousing needs. It is an open standards-based set of technologies. A key distinction of data lakehouse from data lakes is that it uses a schema and **Atomic, Consistent, Isolated, and Durable (ACID)** transactions. Lakehouses allow updates to a record, while data lakes treat data as immutable. In Databricks, Spark is the computational engine supporting all lakehouse processing, and Delta Lake is the storage format used to enable ACID transactions and schemas. Delta Lake is based on the Parquet format, with transaction logs in **JavaScript Object Notation (JSON)** to journal interactions with data. A Delta Lake exists on top of data lakes and cloud storage containers.

Databricks community edition

Databricks understands that learning technology is essential for its adoption. Databricks offers a community edition of its platform to enable learning and smaller workloads. The community edition offers limited functionality compared to the enterprise-class versions available on AWS, Azure, and **Google Cloud Platform (GCP)**. The community edition has several restrictions, including little computational power and lacks automation capabilities via an API. To learn more about the Databricks community edition, go to <https://docs.databricks.com/getting-started/community-edition.html>.

Recipe 1: Signing up for the Databricks community edition

To sign up for the Databricks community edition, go to <https://www.databricks.com/try-databricks> and fill out the form, as shown in *Figure 1.1*. You will be asked for your name, email, company, and job title:

The screenshot shows a web browser window with the URL `databricks.com/try-databricks#account`. The page features the Databricks logo and the heading "Try Databricks free". Below this, it states: "Test-drive the full Databricks platform free for 14 days on your choice of AWS, Microsoft Azure or Google Cloud." Three bullet points highlight key features:

- ✓ Simplify data ingestion and automate ETL. Ingest data from hundreds of sources. Use a simple declarative approach to build data pipelines.
- ✓ Collaborate in your preferred language. Code in Python, R, Scala and SQL with coauthoring, automatic versioning, Git integrations and RBAC.
- ✓ 12x better price/performance than cloud data warehouses. See why over 7,000 customers worldwide rely on Databricks for all their workloads from BI to AI.

On the right side, there is a "Create your Databricks account" form (1/2). The form includes the following fields:

- First name:
- Last Name:
- Email:
- Company:
- Title:
- Phone (Optional):
- Country:

Below the form, there is a disclaimer: "By submitting, I agree to the processing of my personal data by Databricks in accordance with our [Privacy Policy](#). I understand I can [update my preferences](#) at any time." At the bottom of the form is a red "Continue" button.

Figure 1.1: Sign-up for Databricks Community Edition

After clicking **Continue**, you are presented with a page asking you to choose your cloud provider, as shown in *Figure 1.2*. Under the section that refers to not having a cloud account, there is a link titled **Get Started** with Community Edition. It is relatively small and easy to miss, but it is how to sign up for the free community edition.

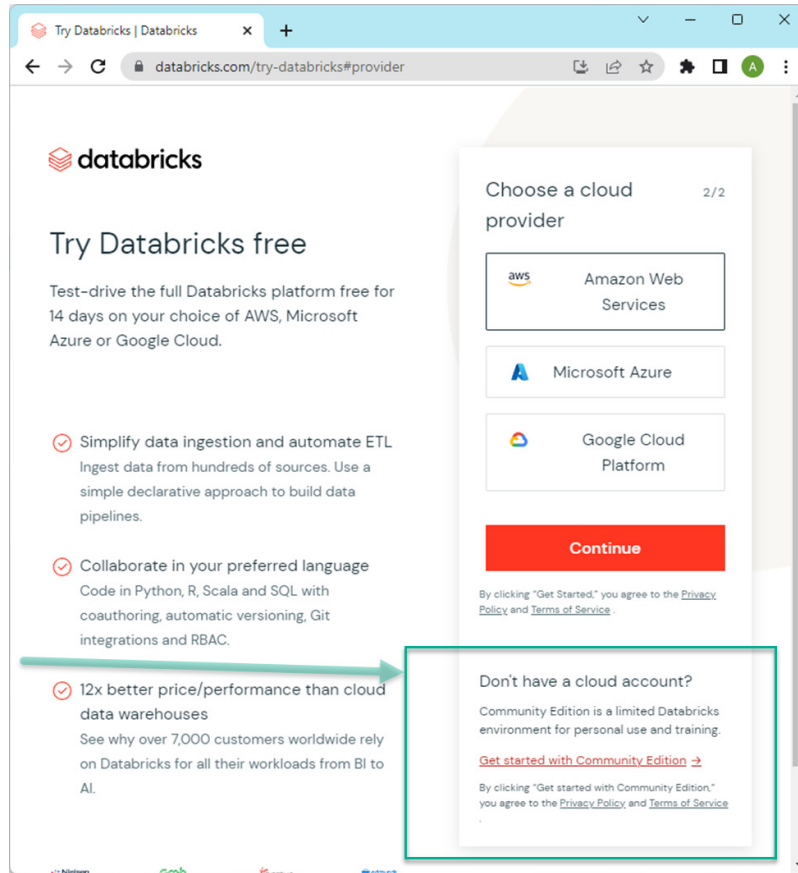


Figure 1.2: Select getting started with Databricks Community Edition

After clicking the link, you will likely be asked to prove you are a human by solving a simple puzzle. After solving the puzzle, you will be redirected to a page asking you to confirm your email address. Check your email and click the link in the body to confirm receipt of the email message. You are then asked to provide a password for logging into the tool. After supplying a password, you will be redirected to the Databricks community edition home page.

Recipe 2: Creating a notebook in the Databricks Community edition

The areas of the community edition Databricks workspace are similar to that of the enterprise-class cloud-hosted versions, as shown in *Figure 1.3*. The layout is organized