



Migrating from Hadoop to Data Lakehouse

Databricks Special Edition

by Stephanie Diamond

for
dummies[®]
A Wiley Brand

Migrating from Hadoop to Data Lakehouse For Dummies®, Databricks Special Edition

Published by
John Wiley & Sons, Inc.
111 River St.
Hoboken, NJ 07030- 5774
www.wiley.com

Copyright © 2022 by John Wiley & Sons, Inc.

No part of this publication may be reproduced, stored in a retrieval system or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning or otherwise, except as permitted under Sections 107 or 108 of the 1976 United States Copyright Act, without the prior written permission of the Publisher. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748- 6011, fax (201) 748- 6008, or online at <http://www.wiley.com/go/permissions>.

Trademarks: Wiley, For Dummies, the Dummies Man logo, The Dummies Way, Dummies.com, Making Everything Easier, and related trade dress are trademarks or registered trademarks of John Wiley & Sons, Inc. All other trademarks are the property of their respective owners. John Wiley & Sons, Inc., is not associated with any product or vendor mentioned in this book.

LIMIT OF LIABILITY/DISCLAIMER OF WARRANTY: WHILE THE PUBLISHER AND AUTHORS HAVE USED THEIR BEST EFFORTS IN PREPARING THIS WORK, THEY MAKE NO REPRESENTATIONS OR WARRANTIES WITH RESPECT TO THE ACCURACY OR COMPLETENESS OF THE CONTENTS OF THIS WORK AND SPECIFICALLY DISCLAIM ALL WARRANTIES, INCLUDING WITHOUT LIMITATION ANY IMPLIED WARRANTIES OF MERCHANTABILITY OR FITNESS FOR A PARTICULAR PURPOSE. NO WARRANTY MAY BE CREATED OR EXTENDED BY SALES REPRESENTATIVES, WRITTEN SALES MATERIALS OR PROMOTIONAL STATEMENTS FOR THIS WORK. THE FACT THAT AN ORGANIZATION, WEBSITE, OR PRODUCT IS REFERRED TO IN THIS WORK AS A CITATION AND/OR POTENTIAL SOURCE OF FURTHER INFORMATION DOES NOT MEAN THAT THE PUBLISHER AND AUTHORS ENDORSE THE INFORMATION OR SERVICES THE ORGANIZATION, WEBSITE, OR PRODUCT MAY PROVIDE OR RECOMMENDATIONS IT MAY MAKE. THIS WORK IS SOLD WITH THE UNDERSTANDING THAT THE PUBLISHER IS NOT ENGAGED IN RENDERING PROFESSIONAL SERVICES. THE ADVICE AND STRATEGIES CONTAINED HEREIN MAY NOT BE SUITABLE FOR YOUR SITUATION. YOU SHOULD CONSULT WITH A SPECIALIST WHERE APPROPRIATE. FURTHER, READERS SHOULD BE AWARE THAT WEBSITES LISTED IN THIS WORK MAY HAVE CHANGED OR DISAPPEARED BETWEEN WHEN THIS WORK WAS WRITTEN AND WHEN IT IS READ. NEITHER THE PUBLISHER NOR AUTHORS SHALL BE LIABLE FOR ANY LOSS OF PROFIT OR ANY OTHER COMMERCIAL DAMAGES, INCLUDING BUT NOT LIMITED TO SPECIAL, INCIDENTAL, CONSEQUENTIAL, OR OTHER DAMAGES.

For general information on our other products and services, or how to create a custom = * ý ŷ Ź book for your business or organization, please contact our Business Development Department in the U.S. at 877- 409- 4177, contact info@dummies.biz, or visit www.wiley.com/go/custompub. For information about licensing the = * ý ŷ Ź brand for products or services, contact BrandedRights&Licenses@Wiley.com.

ISBN: 978- 1- 119- 89420- 9 (pbk); ISBN: 978- 1- 119- 89421- 6 (ebk). Some blank pages in the print version may not be included in the ePDF version.

Publisher’s Acknowledgments

Some of the people who helped bring this book to market include the following:

Project Manager:
?|nna >qn_d ah) Hæcdpkj
Sr. Managing Editor: Rev Mengle
Managing Editor: Camille Graves

Acquisitions Editor: =odlau ?k au
Sr. Client Account Manager:
Matt Cox

Table of Contents

INTRODUCTION 1

 About This Book 1

 Icons Used in This Book..... 2

 Beyond the Book..... 2

CHAPTER 1: Understanding the Shortfalls of Legacy 8UUs@_Yg..... 3

 Reviewing the History of Data Lakes..... 3

 H.Y.VYbY hg..... 4

 The challenges..... 4

 Recognizing the Limitations of Hadoop..... 5

 The Divided Worlds of Data 5

 Keeping Up with Changing Regulatory Requirements 6

CHAPTER 2: Prioritizing Your Data and AI Strategy 7

 Examining Top Strategic Goals 7

 A cj]b[Zcfk UfX'cb 'hY'5 'a Uh f]mWfj Y..... 8

 Reducing risks 8

 Controlling costs 9

 Focusing on Creating a Data Culture 9

 Improving data literacy 10

 Collaborative tooling 10

 Capitalizing on the Vs of Data..... 11

 Focusing on Business Value..... 11

 Reviewing Options in Future Technology..... 12

 ; YH]b['6YHf' bg][\hgZfca '8UH'UbXd g]b['5 #A @UhGWY..... 12

CHAPTER 3: '8Uk b]b['cZh Y'@_Y\ci gY..... 13

 : i hi fY!dfcc b['M:i f'Cf[Ub]nUhcb..... 13

 Evolving from the EDW to the Data Lake to the Lakehouse 14

 Reviewing What a Lakehouse Is 15

 Looking at Legacy Data Lakes..... 16

CHAPTER 4: '6YbY H]b[Zfca '@_Y\ci gY'A][fUh]cb 17

 HfUbgZfca]b['M:i f'Cf[Ub]nUhcb..... 18

 8UHVF]Wg@_Y\ci gY'7i g'ca Yf'Gi WggsGcf]Yg..... 19

 CBC Radio-Canada..... 19

 H&M Group..... 20

CHAPTER 5: **FYj JYk]b['K \nhtc'A][fUHY'hc'h Y'@U_Y\ci gY** 21

FYj JYk]b['Ub '5[]Y' hYfUHj Y '5ddfcUW 'hcA][fUH]cb 21

A cj]b['<UXccd '5fW]hYVh fY'hcsh Ys7'ci X 23

Planning the Migration Journey 23

Migrating from Hadoop to the Lakehouse 25

 Figuring out what you're administering 25

 Migrating your data 25

 Processing your data 26

FYUd]b['h Y'VYbY 'hg'cZgYWf]hmUbX\$[cj YfbUbW 26

 Considering your SQL and BI workloads 26

CHAPTER 6: **HYb 'Dc]bhg'hc'9bUV'YUbX\$GM'Y'M;i f'8UHJ'**

UbX\$5 'GhfUH[m 27

Introduction

The cost, complexity, and viability of existing Hadoop platforms have failed to deliver on business value due to the lack of data science capabilities, the high cost of operations, the inability to scale, the lack of agility, and poor performance. As a result, enterprises are looking to migrate their existing Hadoop platforms to the data lakehouse.

The *data lakehouse* is a cloud-native platform for data management that provides a powerful engine for data processing and simple tools for developers, analysts, data scientists, and business users in an intuitive user interface (UI). It enables you to build, deploy, scale quickly, and manage analytical applications in minutes instead of hours or days. The data lakehouse is an open data architecture that combines the best of data warehouses and data lakes on one platform.

About This Book

Welcome to *Migrating from Hadoop to Data Lakehouse For Dummies*, Databricks Special Edition. This book looks at why and how to migrate from Hadoop to the lakehouse to prepare your organization to meet the future. It also uncovers what your leaders need to know to meet new challenges in data management. This book covers several topics:

- » The history of data lakes
- » Determining where your company resides on the data and
- » FYJ JYk Jb['k \UhU'XUH'U_Y\ci gY'Jg
- » 5W7f Jb['VYbY hg'cZa J[fUhJb['Zfca 'cUXccd'hc'hY'XUH' lakehouse
- » What to consider when you migrate to the lakehouse
- » HYb'ghYdg'hc'YbUV'Y'UbX'gWY'nci f'XUH'UbX'5'ghfUH[m

Icons Used in This Book

Prdrkqcdkqppdœ^kkg(`e araj pe_kj o]m qœa` pk dœdœdpœ l kn-
tant information. Here's what they mean:



TIP

The Tip icon adds information to help you manage processes faster and easier.



REMEMBER

The Remember icon points out content to remember when searching your memory bank.



WARNING

The Warning icon alerts you to information that you should be aware of that can be harmful to you or your company.



TECHNICAL
STUFF

Sometimes I give you a few tidbits of research or facts beyond the basics. If you want to know the technical details, watch out for this icon.

Beyond the Book

This book can help you discover more about migrating from Hadoop to the data lakehouse, but if you want resources beyond s d] pœdœ ^kkg k arœ(_da_g kqppla tkllks ej c hœ goœ

- » databricks.com/discover/lakehouse: See why your 'Y[UWrXUHuK UFY\ci gY Vlb hgi ddcfhHAY UXj UbWX bYYXg' nœi \Uj YhcXUm'
- » databricks.com/product/data-lakehouse: See how the lakehouse platform can support all your data, analytics, and 5 i gY WgYgcb Uglã d'YœcdYbœa i 'hWœi X d'Uhcfa "
- » databricks.com/product/Databricks-sql . 8]gWœj Yf\ck ' 8UHUVf]WgGE @U'ck g'nœi 'œ'GE @k cf_cUXg'cb 'HAY`U_Y- \ci gY UFVW]HVM fY"
- » databricks.com/try.Hfœci hHAY 8UHUVf]Wg'@U_Y\ci gY"
- » databricks.com/p/ebook/migration-guide-hadoop-to-databricks. 5'ghd!VmgHd hYVb]W' [i]XY'cb\ck 'œ' a][fUHY Zœa 'cUXccd 'œ'8UHUVf]Wg'

- » Looking at the history of data lakes
- » Migrating away from Hadoop
- » Dealing with changing regulations

Chapter 1

Understanding the Shortfalls of Legacy Data Lakes

Traditional data lakes such as Hadoop (and even their cloud variants) are coming to the end of their life cycles. As a result, data lakes are primed for the next phase in their evolution. Businesses have realized that to use their data and analytics as the premier assets they are, migrating to the data lakehouse is the next logical step from the Hadoop data lake. This chapter looks at the evolution of the data lake (including Hadoop) and how it's poised to play a role in the data lakehouse architecture.

Reviewing the History of Data Lakes

Before the creation of data lakes, companies relied on enterprise data warehouses (EDWs) to store data. This storage worked for a time. However, as the volume and types of data grew exponentially, aggregating databases in EDWs produced a collection of data silos that were not well suited for many business purposes.



WARNING

The main drawback of this approach was that they couldn't be queried across business unit boundaries to deliver enterprise-level actionable insights. As a result, the data lake was created around the year 2011 to address this problem.



REMEMBER

The *data lake* is a repository for raw unstructured data and is usually Hadoop was deployed to lower costs and scale out the resources using commodity hardware, albeit using tightly coupled storage and compute resources. Now Apache Spark runs data lakes in the cloud. Data lakes were cheap, and they could store all kinds of data based on open-standard formats. That meant that no bottlenecks existed between the data lake and its external sources.

History

EDWs were formalized data models that aggregated the data at an enterprise level for certain reports and dashboards. They usually didn't have the raw, granular data that business teams needed for self-service analytics and exploratory and advanced machine learning. Data lakes have the capacity to scale in storage and compute to house all the data — structured and unstructured — for the whole enterprise. As a result, the data lake was created.



TIP

Data lakes were created to address the limitations of EDWs.

- » All data was up to date, and it was easy to add new sources of data.
- » You didn't need to maintain multiple copies of the database.
- » Large-scale data cleansing and transformations were possible.
- » You could run ad hoc queries against the entire data set.
- » You could easily extract data from the data lakes and send it to other locations.
- » It supported open-source ML libraries.

The challenges

Data lakes faced several challenges.

- » No support for atomicity, consistency, isolation, and durability (ACID) transactions

- » No enforcement of data quality or governance
- » Failed jobs and missed data
- » Poor business intelligence (BI) support
- » Poor performance



WARNING

Another problem with data lakes was that they often became a source of shortcomings such as inadequate data governance and the lack of attention to data quality, data lakes were sometimes referred to as *data swamps*. Read more about how data lakes evolved to support the new data lakehouse in Chapter 3.

Recognizing the Limitations of Hadoop



WARNING

Hadoop solved a data problem. However, as time went on and organizations faced new challenges, several limitations of Hadoop came to the fore. They included

- » **Wasted hardware capacity:** Over-capacity is a given in on-premises implementations so you can scale up to your peak time needs, but the result is that much of that capacity sits idle but continues to add to the operational and maintenance costs.
- » **DevOps burden:** You can assume roughly four to eight full-time employees (FTEs) for every 100 nodes. (This is based on the experience of Databricks' customers.)
- » **Increased power costs:** Expect to pay as much as \$800 per server annually based on consumption and cooling. That's \$80,000 per year for a 100-node Hadoop cluster.
- » **Software version upgrades:** These upgrades are often mandated to ensure that the support contract is retained. Those projects take months at a time, deliver little new functionality, and take up precious bandwidth.

The Divided Worlds of Data

Cloud data lakes, by design, require customers to use EDWs for data integration and governance. This is a challenge because data lakes are designed to be a single source of truth, while EDWs are designed to be a single source of truth for a specific business process.

management, and governance of data. Data handling is the key to why the lakehouse is so valuable. The value of using the lakehouse architecture is that you can

- » Use data lakes as a landing zone for all your data for SQL, BI, data science, and ML use cases.
- » Easily meet regulatory requirements to mask data that contains private information before it enters your data lake.
- » Secure your data lake with role-and-review-based access controls.
- » Catalog the data in your data lake.
- » Get the same price performance of an EDW with the scale of a data lake.

Keeping Up with Changing Regulatory Requirements

One of the crucial things that organizations must do is ensure that they comply with changing regulations. Regulatory change d[o g _r]æ` 1, , 1 a_n aj pøj _a pda . . , 4 clk^] h j lj _d h_rææ(and it's boosted regulatory costs. Two key regulations include the General Data Protection Regulations (GDPR) and the California Consumer Privacy Act (CCPA).



TECHNICAL
STUFF

Delta Lake on Databricks manages GDPR and CCPA compliance for your data lake. Delta Lake adds a transactional layer that provides structured data management on top of your data lake. As a result, it can dramatically simplify and speed up your ability to locate and remove personal information (also known as *personal data*) in response to consumer GDPR or CCPA requests. For more information, see <https://docs.databricks.com/product/delta-lake-on-databricks>.

- » docs.databricks.com/product/delta-lake-on-databricks
- » docs.databricks.com/security/privacy/gdpr-delta.html

- » Looking at top strategic goals
- » Developing a data culture
- » Using the Vs of data
- » Increasing your business impact
- » Looking at tech options
- » Getting useful insights

7\UdhYfs2

Prioritizing Your Data and AI Strategy

In today's competitive environment, it's not enough to have the right architecture to support your organization's data. You also need a comprehensive strategy that serves all the essential components of your organization. This strategy should include leveraging people, business goals, and technology. It's the key to long-term business success. Ultimately, the technology should be an enabler of the strategy and not the other way around.

Examining Top Strategic Goals

Data and technology executives who want to nail their data and AI strategy should focus on the following top strategic goals:

- » Moving forward on the AI maturity curve
- » Reducing risks
- » Controlling costs

You look at each goal in this section.

Moving forward on the AI a Uri f]msWfj Y

Pk i aappda b]p]ra]j` pk i kra]k]c]`a j a` i]p]rpu _q]ra(_ki l]j aej aa` pk]`kl p]`]p]]oap]j` =Ei g` oap* Pda ck]h is to go from being descriptive to prescriptive. To this end, Databricks created a maturity model that your organization can qoa pk qj` aro]j` pda _q]raj pcp] pa kbukqnf]kq] au pk`]p]]j` =E i]p]rpu* Pda i k` ahao]o k]l]ks o6

- » **Explore:** At this stage, your organization is beginning to explore Big Data and AI and understand the possibilities and potential of a few starter projects and experiments.
- » **Experiment:** Organizations at this stage are building the basic capabilities and foundations to begin to explore a more expansive data and AI strategy, but they're lacking vision, long-term objectives, or leadership buy-in.
- » **Formalize:** At this stage, data and AI are budding into a Xf] Yf]c]j U'i Y'Zcf Vi g]bYgg'i gYfg'U][bYX'hc'gdYVY Wdfc'YVtg' and initiatives, as the core tenets of data and AI are integrated into corporate strategy.
- » **Optimize:** Data and AI are core drivers of value across the cf[Ub]nUh]cbs` g]fi V]i fYX'UbX'Wb]fU'hc'h.Y V]f]dcfUHY' strategy with a scalable architecture that meets business needs and buy-in from across the organization.
- » **Transform:** At this stage, data and AI are at the heart of the V]f]dcfUHY'g]fUHY[m]UbX'Ub']bj U'i UV'Y'X] YfYb]h]Urc'f'UbX' driver of competitive advantage.

Do you recognize your organization at one of these stages? If so, do you know how to move forward to achieve a higher maturity? You can schedule a custom business value assessment at databricks.com/p/business-value-assessment-databricks.

Reducing risks

=j k]dan op] p]ce_ ck] h]kn h]]`aro kb krc]j e] p]k] o eo pk m` q_a several potential risks such as weak data management, failed IT projects, missing out on innovation due to the lack of advanced analytics platforms, and the ever-present threat of cyberattacks. These threats make it imperative to have a consistent way

to store, process, manage, and secure data. However, this goal is
i j` a i kra _ki l hat ^u pda bklks ej c j aa` o6

- » Adhering to the evolving privacy regulations landscape, such as the General Data Protection Regulation (GDPR) and the California Consumer Privacy Act (CCPA)
- » Adjusting data management to contend with new privacy directives such as those from Google and Apple
- » Figuring out how to take advantage of new data sources that can supersede typical behavioral, demographic, and engagement data

Controlling costs

Leaders must always contend with the need to control costs. Data lakes can get expensive fast as the amount of data it manages grows. On top of it, there are overheads from data center equipment, database administration operations and maintenance, and many locked-in vendor agreements.

Pda oklupkj pk pde l rk^hai \$ _ki i k`] pjc _qraj p`] p] j` =Eej qej praofao ei l hai aj pjc] _lkq`] n dqa ppa p] po alj qe] j` at e^ha pk]`] l ppk pda _d] j c ej c ^qos aooj aa` o] j` d] ockk` price performance in the cloud as data size increases.



Utilizing simpler architectures also produces more agility for data and technology executives to integrate and have actionable insights without delays or IT intervention.

Focusing on Creating a Data Culture

Taking advantage of the value of data to organizations has created the need for companies to establish and maintain a data culture.
= rk^qop`] p] _qlppra aj oqrao p] pukqnkr:] j ev] p] s dhi aappra future. You can spot a company that has a strong data culture in ps k gau s] uo6

- » You see the entire organization making daily informed business decisions by using available, relevant data.
- » Data is given greater weight than experience, intuition, or tenure. The data and insights provide the proof.

Improving data literacy

Collaborative tooling

DEMOCRATIZING YOUR DATA

- **Embedding data scientists directly into business units:** This process enables them to interact directly with data users.
- **Putting access to the analytics in the hands of the users:** This action permits users to draw insights themselves.
- **Providing senior leaders with access to visual tools with simple interfaces:** This provision allows your senior leaders to get data insights as needed.

Capitalizing on the Vs of Data

Coping with the wave of data that hits companies every day is impossible to manage without the right data management



TIP

» **Volume:** The size of the data

» **Veracity:** The extent to which you can rely on the data

» **Velocity:** The speed at which the data is consumed

The assumption is that it should be analyzed as close to receipt as possible.

» **Variety:** The range of data types and sources

» **Value:** The usefulness of data in decision making (not all data is equal)



TIP

Copying any data. This process allows you to include all properties of unstructured data sets from social media posts and metadata to catalog images in their analysis and models.

Focusing on Business Value

Now that so many more unique forms of data are available, businesses realize that their legacy platforms can't scale and meet the increasing demands for better data analytics.

Data and technology executives want to use that data to get a lower-cost approach that improves the user experience and increases collaboration across data personas. This goal moves them away from complex and expensive on-premises architectures.

Reviewing Options in Future Technology

Lakehouse architecture allows you to future-proof your technology with open standards and the use of multi-cloud. Many drivers exist for a multi-cloud approach, but one of the biggest is economic leverage.



REMEMBER

Multi-cloud architecture allows you to future-proof your technology with open standards and the use of multi-cloud. Many drivers exist for a multi-cloud approach, but one of the biggest is economic leverage. You should adopt multi-cloud architecture for a variety of reasons.

- » Data is naturally multi-cloud.
- » It improves negotiation leverage.
- » You get the best technology for the workload.
- » Cost is optimized.

Getting Better Insights from Data

Using Machine Learning to Get Better Insights from Data

Machine learning (ML) helps organizations make better decisions by analyzing data. ML can be used to predict future trends, identify patterns, and optimize processes. ML can help organizations make better decisions by analyzing data in today's marketplace.

- » : i h fY!dfcc b['nei f'cf[Ub]nU]cb
- » Understanding the evolution to the lakehouse
- » Looking at lakehouse architecture and data lakes

7\UdhYfs3

Dawning of the Lakehouse

Lessons learned from working with enterprise data warehouses (EDWs) and data lakes have paved the way for the lakehouse's modern cloud-based data architecture. It combines both the best properties and capabilities to provide a far i kra l ks artq[h]j` at'e'la`]p l h]p[kri p]l j l koe'la g p]a past. This chapter looks at the evolution of the modern cloud-based lakehouse and the need to future-proof your organization.

: i h fY!dfcc b['M:i f'Cf[Ub]nU]cb

One of the critical requirements of enterprise leaders is that they successfully prepare their organizations to meet the future. Relying on antiquated processes to manage data could overwhelm your organization, put you at a competitive disadvantage, and s] opa _rpe]hdqi lj lj` j lj _d h_]l q] l t



TECHNICAL
STUFF

According to Statista, the total amount of data created, captured, copied, and consumed worldwide is forecast to increase by 152.5 percent from 2020 to 2024 to 149 Zettabytes.

Reviewing What a Lakehouse Is

The lakehouse takes the greatest elements of data warehouses and data lakes and combines them into a single platform that gives you the best of both worlds. Operating a lakehouse architecture is the foundation that enables you to

- » Manage all data use cases on one single source of truth for all your data.
- » 6Y'a cFYfYgdcBgJj Y UbX' bX'bYk jbgj[\hg ZJghYf"
- » 5`ck Yj YfncbYtc`cc_UhH.Y'gUa Yj Yfgjcb`cZH.Y'XUHU"
- » Simplify existing architectures and security by reducing silos and the number of systems and tools that you need to manage.
- » Have the ability to consolidate and tie your data marts and 98K gk Jh' cH.Yf'i bgf'i Vh fYX XUHU Zcf YbfjWka YbhUbX' create innovative data products.
- » Perform extract, transform, load (ETL) operations on the XUHUk Jh Jb`H.Y'XUHU`U_Y\ci gY"



REMEMBER

Lakehouse architecture is

- » **Simple:** Unify your data, analytics, and AI on one platform.
- » CdYb.'l bJznci f'XUHU'YVtgngh'a 'k Jh' cdYb'g'UbXUfXg'UbX' Zcfa Uhgzk \JW'dFYj Ybhj YbXcf`cW!\jB"
- » **Multi-cloud:** Maintain consistent management, security, and governance experience across all clouds and enable your H'Ua g'hc ZcW'g'cb'di Hjb['U`nci f'XUHU'hc'k cf`hc'XjgVtj Yf' bYk jbgj[\hg"

=o Jj at J i l ha(pda @) p ^ne go H gadkqoa l H pkri eo odks j g
Figure 3- 2.

- » How legacy architectures can't keep up with modern organizations
- » How real-time data, predictive insights, and Hadoop's rising cost prevent teams from driving high-impact business outcomes

4

Legacy on-premise analytics architectures aren't keeping up with the needs of modern organizations

Things like the inability to use real-time data, getting predictive insights, and Hadoop's rising cost prevent teams from driving high-impact business outcomes

Legacy on-premise analytics architectures aren't keeping up with the needs of modern organizations. Things like the inability to use real-time data, getting predictive insights, and Hadoop's rising cost prevent teams from driving high-impact business outcomes.

As organizations look to do more with their data, empower data scientists, and reduce infrastructure maintenance and data management costs, the world of data and AI needs a Hadoop alternative. Organizations worldwide have realized that it's no longer a matter of *if* migration is required to stay competitive and innovate but a matter of *when*.

Cloud data lakes and data warehouses give you some customer success stories.

How Data Leaders Are Preparing for the Future

When making changes that will future-proof their organization, data leaders are highly concerned about the following:

- » The opportunity cost of not migrating
- » The cost of migrating away from their existing solution
- » The learning curve that data teams will encounter and the availability of skills for building solutions on the new platform
- » The breadth of the ecosystem built around the solution

They believe that four truths should guide their decisions:

- » **A data maturity curve** To gain better insights from their data, leaders want to move up the maturity curve so they can take full advantage of a data lakehouse. "Data maturity is a curve that shows the progression of data management from silos to a unified data lakehouse. Leaders want to move up the curve to take full advantage of the data lakehouse." "Data maturity is a curve that shows the progression of data management from silos to a unified data lakehouse. Leaders want to move up the curve to take full advantage of the data lakehouse."
- » **Open source is preferred** Leaders are uneasy about locking themselves in with any vendor and want to use open source whenever possible. Using open source and open formats also helps with more readily available skilled resources compared to using proprietary technologies and formats.
- » **Multi-cloud strategy** Most companies are moving to the cloud to strengthen their data capabilities. However, in adopting a multi-cloud strategy, it's important to have the ability to run your analytics and AI workloads across all clouds.
- » **Eliminate complexity** Leaders want to eliminate complexity, minimize silos, and avoid multiple copies of the same data. They want one governance, security, and lineage model for all their data across all clouds.

Regarding these truths, consider how the lakehouse architecture meets each of these needs:

- » **A data maturity curve** The lakehouse uses ML and AI from the ground up.
- » **Open source is preferred** The lakehouse uses open formats and standards to provide greater data portability and avoid vendor lock-in.

- » D'Ub Zcf hY'a i 'h!Wci X" The lakehouse leverages low-cost cloud object stores to store *all* enterprise data across the major cloud providers.
- » I gY'gla d'YXUHUfW]hVMI fY" The lakehouse supports all use-case platforms, including data engineering, data warehousing, real-time streaming, data science, and ML across the major clouds.

S daj ukq i æn] pã pk pda h gadkqoa(ukq ` anæra oaran] h^aj a _ej h outcomes. These include the ability to

- » Continue leveraging low-cost cloud object stores.
- » Fi b X] YfybhXUHUk cf _cUXgZca '6 'hc'5 'WghY YWij Y'mi UbX Y VYbhm'
- » Utilize open formats and standards. You won't experience vendor lock-in.
- » Future-proof your investment in modern cloud data architecture, including the ability to leverage multi-clouds.

8UHUf]Wg@U_Y\ci gY'7i ghca Yf' G WWggGhcf]Yg

If you're curious how other companies have migrated from Hadoop to the lakehouse, this section provides you with two com-
l]j æopd] pra] læv` æej e _lj poq_aos daj pdaui æn] pã pk pda Databricks Lakehouse platform.

767'FUX]c!7UbUXU

CBC Radio-Canada's mission is to enlighten and entertain its diverse audience. However, its legacy Hadoop infrastructure prevented the company from using its diverse data to personalize its customer interactions in the way it wanted to. Personalization was a driving goal. CBC Radio-Canada turned to Databricks for dah]j ` læran] ca` pda a _æj _æo kb pda _hkq`]j ` pda l ks ank b a lakehouse architecture to give data teams direct access to all of the company's data.

Databricks SQL empowered CBC Radio- Canada to quickly run SQL queries to derive insights into its digital audiences and behaviors. The reduction in time- to- insight was lowered by 50 percent. The company is now able to provide more visibility into its digital audiences and develop strategies and services that will boost engagement and retention.

</ A ; fci d

H&M Group is well- known as a major disruptor and innovator in the fashion and retail industry. However, its on- premise Hadoop system crippled its ability to ingest and analyze data generated by millions of customers needed to power predictive models. To solve this problem, H&M Group moved to the Databricks Lakehouse Platform to simplify infrastructure management, enable performant data pipelines at scale, and simplify the ML life cycle. This allowed the group to make data- driven decisions that accelerated business growth.

The results were dramatic:

- » A 70 percent reduction in operational costs
- » 6YHfWcgghUa Vt`UVcfUjcb UbX'UfYXi Wjcb Jb h.Y' number of components needed to go in production with easy setup and management
- » Faster time-to-insight

- » Taking an agile approach to migration
- » Detailing the Hadoop migration

7\UdhYfs5

Reviewing Why to Migrate to the Lakehouse

Migrating to a new architecture can be a complex process. If you've decided to migrate from your legacy Hadoop, think through several essential considerations before architectural choices helps you make well-informed decisions about how best to proceed with your modernization initiatives. This chapter looks at the value of taking an agile, iterative approach to migration and suggests how to plan and execute the migration journey.

Reviewing an Agile Iterative Approach

How you migrate your data to the lakehouse is a critical decision. Regardless of the path you follow, balance is required. Databricks recommends that you do all your migrations in a phased agile manner:

- » **Make the decision.** Migrating from Hadoop to a modern cloud data platform can be daunting, but staying with



TECHNICAL
STUFF

existing solutions can be even worse. The technical debt you

» **Use a balanced approach to lift and shift and modernization.**

Lift and shift the code as well as modernize in one iteration. Use lift and shift with automated code converters and immediately modernize to optimal Databricks patterns.

Here *lift and shift* refers to the movement of an application design and code from one environment to another without making massive changes. But don't wait to modernize and

» **Learn what worked and didn't and iterate.** Add on

additional use cases and workloads as you go.

» **Show success in shorter sprints and adapt.** This way you accelerate the value delivered and demonstrate immediate success to the stakeholders. The learnings and feedback also help you improve the next iteration of migration.



WARNING

Just simple lift and shift is rarely the answer; if you just lift and

» You won't get to fully utilize the lakehouse.

» You lose out on any innovation you may discover if you just move everything as is.

» You won't have the chance to improve your tech strategy to optimize for cost savings, agility, and scale.

Both lift and shift and total re- engineering have pros and cons:

» **Lift and shift:** On the pro side, it's faster and more critical to do if you have an upcoming license renewal deadline. The downside is that you may not take the opportunity to

» **Total re-engineering:** On the other hand, it can take a long time to complete and comes at a high upfront cost.

Moving Hadoop Architecture

hcsH Ys7`ci X

I krēj c D]` kkl ꞑk ꞑda _lkq` eo]j at_ahaj ꞑ rop qpal ēj ukqn migration journey. However, it doesn't unlock your data like moving to the lakehouse does. Consider what the impact of mov- ēj c ukqn`] ꞑ ꞑk ꞑda h gadkqoa] kn` o ukq* Ukq _]j

- » Deliver data faster to business users for better and more timely business decisions.
- » Deliver consistent data from a shared data lake that's properly governed to ensure the entire organization is k cf_]b['c 'hY'gUa Y'XUHU"
- » Achieve greater scalability to take on the largest AI and ML i gY'WgYg'UbX'YbUV'Y'a cfY'Y YWij Y'UbUrhVg'UbX'5 'Vm bX]b['bYk 'a Uf_YhgZ]bWYUg]b['fYj Ybi YzfYXi V]b['V'g'g'UbX' lowering risk.
- » Make larger data sets available to business decision makers and give them the ability to visualize your entire data lake while _YYd]b['V'g'g'ck 'h'fci [\ 'U'dUhtZc'f!k \Uhtnci !i gY'UfW'jhVmi fY"

Planning the Migration Journey

When considering a migration journey, carefully plan each step along the way. Therefore, the journey can be depicted as a set of steps. Here's how each step works:

1. Ask internal questions.

This step covers the discovery phase. The key to this step is to answer two questions: Where am I now, and where do I need to go? Make sure that you collect questionnaires from all your XUHU'hUa gZ'W]YZ]bZcfa Uh'cb'c 'Wf'g'UbX'ch'Yf'fY'Yj Ubh' stakeholders. Be prepared for a lot of new learning and gY'Z'X]gW'j YfmUg'hUa g'h'gh'UbX'j U'XUH'Uggj a dh'cbg"

2. Make a migration assessment.

8i f]b['h'Y'UggYgga Ybhd\UgYžnci 'k Ubh'ic'fY' bY'UbX' evaluate the solutions on the table. Take an inventory of all workloads and prioritize the use cases.

Migrating from Hadoop to the Lakehouse

No matter what architecture you're working with, migrations are j arana] ou* @] p] ^ne go oqccaop ukq _kj oe an ra] na] o kb tk_qo when migrating from Hadoop to minimize adverse impact, ensure ^qoqj aoo _kj p] qeju(] j ` i] j] ca_kop a a_prahi*Pdkoa] na] o] na covered in this section.



REMEMBER

Don't forget to bring key business stakeholders along on this journey. It's as much a technology decision as it is a business ^ a_eokj * Ukq j aa` ukqn ^qoqj aoo op] gadkhi an ^rkqcdp ej pk pda journey and its end state to have a successful migration.



TECHNICAL
STUFF

Bkn] pa_dj e] h`aal `era kj pda ra i eon] p] kj] na] o_kraa` ej this section, visit databricks.com/blog/2021/08/06/5-key-steps-to-successfully-migrate-from-hadoop-to-the-lakehouse-architecture.html.

Figuring out what you're administering

D]` kkl aj rækj i aj p d] ra ^aaj _kj j a` ^u oqj cha _ki l qpa environments. As migrations happen, consider how to take advantage of a multi- compute environment with customization to size, run time, and more — all while gaining greater control on who has access to what data.

Migrating your data

Consider a balanced approach to your data migration to ensure business continuity while enabling new use cases. Dual ingestion allows for organizations to feed new data to your preferred cloud storage, and historical data can be migrated in parallel or at a later stage. Also, with data migration, you need to decide between a push- oriented data migration and a pull data migration. A *push-oriented data migration* is when the source pushes data to the target system, and a *pull data migration* is when the target systems pull data from the source system.

Ukq _dkkoa ^aps aaj pda ps k ^u` apari ej ej c dks _repe] hpda _qn- raj ps krglk] ` o] na pk pda ^qoqj aoo] j ` dpda ^qoqj aoo] j] kni downtime with the migration. Push migrations tend to be easier and more automated, whereas pull migrations (maybe needed for a subset of your data) can be more complex and intrusive as you invite third parties to extract the data.

Processing your data

Data processing is transferred over to a data lakehouse across the full stack. From data engineering, ETL, batch processing, to SQL and ML, the data lakehouse architecture easily allows organizations continuity with their data processing needs for like-to-like migration to a modern cloud platform.

FYUd]b['h Y VYbY hgjcZgYWf]hmi UbX\$ cj YfbUbW

With authentication and governance capabilities to meet your organization's requirements, you can migrate your data to a modern cloud platform.

- » **Authentication:** Single sign on (SSO) with SAML 2.0 supported corporate directory
- » **Authorization:** Fine-grained access control to secure your data
- » **Metadata management:** Integration with your preferred metadata management tool
- » **Quick and easy curation of data for your organization**
With enhanced discovery capabilities, and know your data lineage.

Considering your SQL and BI workloads

Don't just transfer your Hadoop workloads; consider your SQL and BI workloads as well. Easily transfer over users to Databricks SQL with familiar experiences, tools, and third-party visualization integrations.

- » Identifying business value
- » Building successful teams
- » Democratizing access to quality data

6

Ten Points to Enable Your Data and AI Strategy

To enable and scale your data and AI strategy with a modern data stack. To achieve that goal, Databricks suggests you keep in mind the following:

- » **Move to production and drive adoption.** Establish a clear set of metrics to measure adoption and track the net promoter score (NPS) so the user experience continues to improve over time.
- » **Establish goals and your business value.** Most organizations establish goals for their data, analytics, and AI journey across business outcomes, people, and technology.
- » **Identify and prioritize use cases.** Use cases should avoid the trappings of “cool” data science and machine learning (ML) projects and focus on delivering business value.

- » **Build successful data teams.** Hiring, training, and leadership. Building a successful data team requires a mix of technical skills, business acumen, and leadership. The team should be able to work together to solve complex problems and deliver value to the business.
- » **Deploy a modern cloud data stack.** The capabilities, legacy and disparate data architectures inspired the next generation of data architectures. The modern cloud data stack refers to as the *data lakehouse*.
- » **Improve data governance and compliance.** We recommend adopting data policies and practices that help the business to realize value through centralizing storage, metadata, and governance on a single platform. This enables you to ensure data quality and consistency across the organization. Tools your organization uses.
- » **Democratize access to quality data.** Legacy data platforms lack the ability to democratize access to quality data. The amount of quality data available than on the sophistication or complexity of the model or algorithm.
- » **Increase the productivity of your workforce.** What tools should you provide to the user community so it can be the best of both worlds? The user experience, move beyond best-of-breed tools, unify the platform and personas, and make sure the tools are simple and self-service.
- » **Make informed build versus buy decisions.** When going through modernization initiatives, your engineers may face a choice between building or buying. The build versus buy decision is a complex one. Also consider if the technology can evolve quickly as your data needs change.
- » **Allocate, monitor, and optimize costs.** Any decision to modernize your data platform requires you to reduce complexity and the costs of managing legacy platforms to ensure the new platform is cost-effective.

WILEY END USER LICENSE AGREEMENT

Go to www.wiley.com/go/eula to access Wiley's ebook EULA.