

# Ingestão de Dados na Arquitetura Lakehouse na Oracle Cloud

**Thais Henrique**

**Andrea Rigoni**

**Trilha Inovação com dados em nuvem**





# Ingestão de Dados na Arquitetura Lakehouse na Oracle Cloud

**Thais Henrique**

**Andrea Rigoni**

**Trilha Inovação com dados em nuvem**



Este trabalho está licenciado sob uma Licença Creative Commons Atribuição-Compartilhalgual 4.0 Internacional. Para ver uma cópia desta licença, visite <http://creativecommons.org/licenses/by-sa/4.0/>.



# Agenda

1

Lakehouse

2

Data Lake

3

Ingestão de Dados

4

Governança de Dados

5

Data Warehouse

6

Hands-on Lab







# Thais Henrique



<https://www.linkedin.com/in/taishenrique/>



[thais.henrique@oracle.com](mailto:thais.henrique@oracle.com)







# Andrea Rigoni



<https://www.linkedin.com/in/rigoni/>



[andrea.carmo@oracle.com](mailto:andrea.carmo@oracle.com)



# A evolução do gerenciamento de análise de dados



# Apresentando o data lakehouse

## Data warehouse e data lake integrados

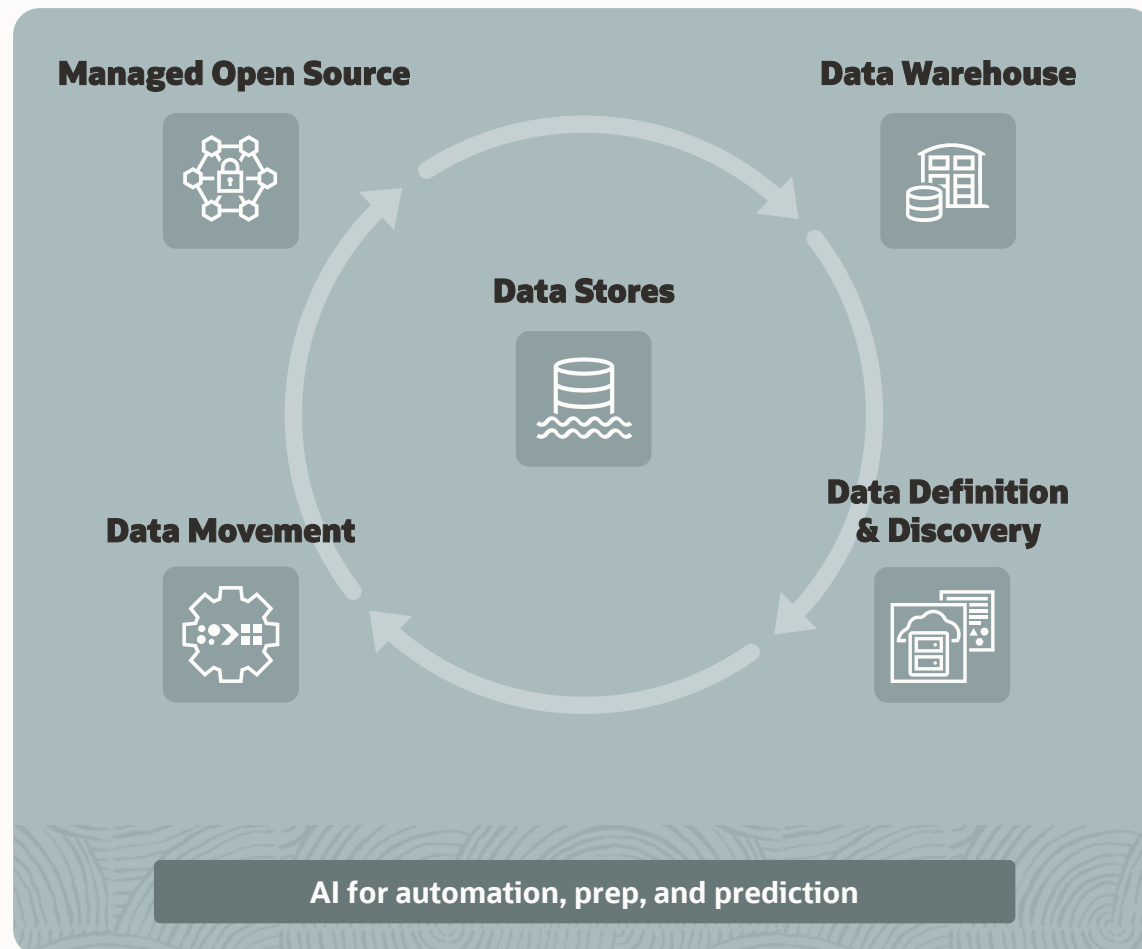
Permitindo análises integradas entre todos os dados

## Elimina data silos

Data irão se movimentar entre warehouse, lake e analysis tools conforme necessário

**Serviços nativos ou Open source** suportando ambos cenários e todo tipo de dado

## The data lakehouse



**Data Warehouse:** dados curados e de valor conhecido

**Data Lake:** dados brutos, ainda sendo explorados, antigos, ou de menor valor. Também utilizado como staging antes da carga para o data warehouse, archive, e repositório para treinamento de modelos de machine learning

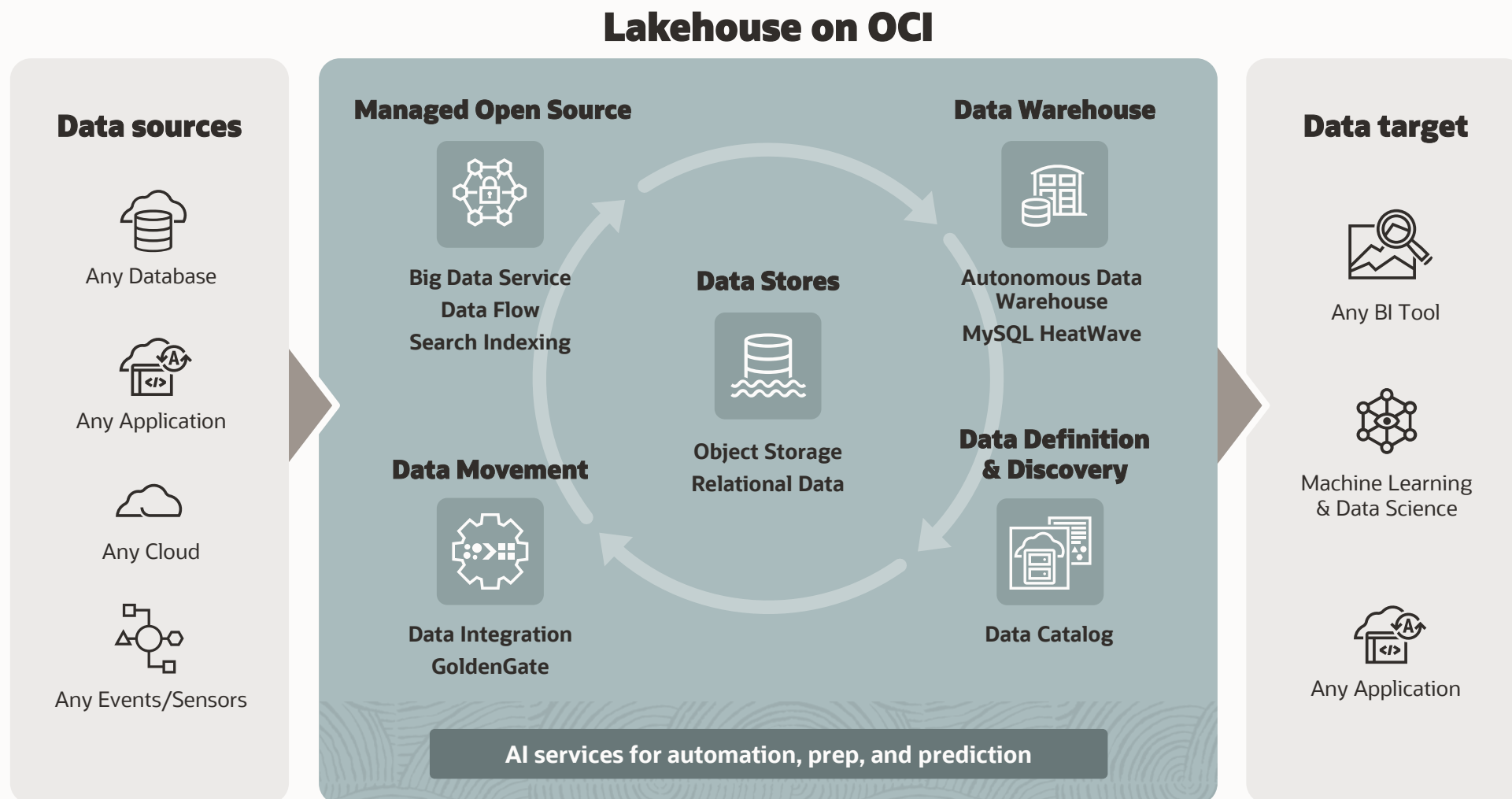
**Managed Open-source Services:** suporte as principais open-source tools para armazenamento e análise. Apache Hadoop, Apache Spark e Elasticsearch

**Data Integration:** dado no lakehouse poderá se mover entre lake, warehouse e open-source analytics, dependendo do caso de uso

**Data Catalog:** mantém uma visão completa de todo o dado disponível, possibilitando descobertas e governança.

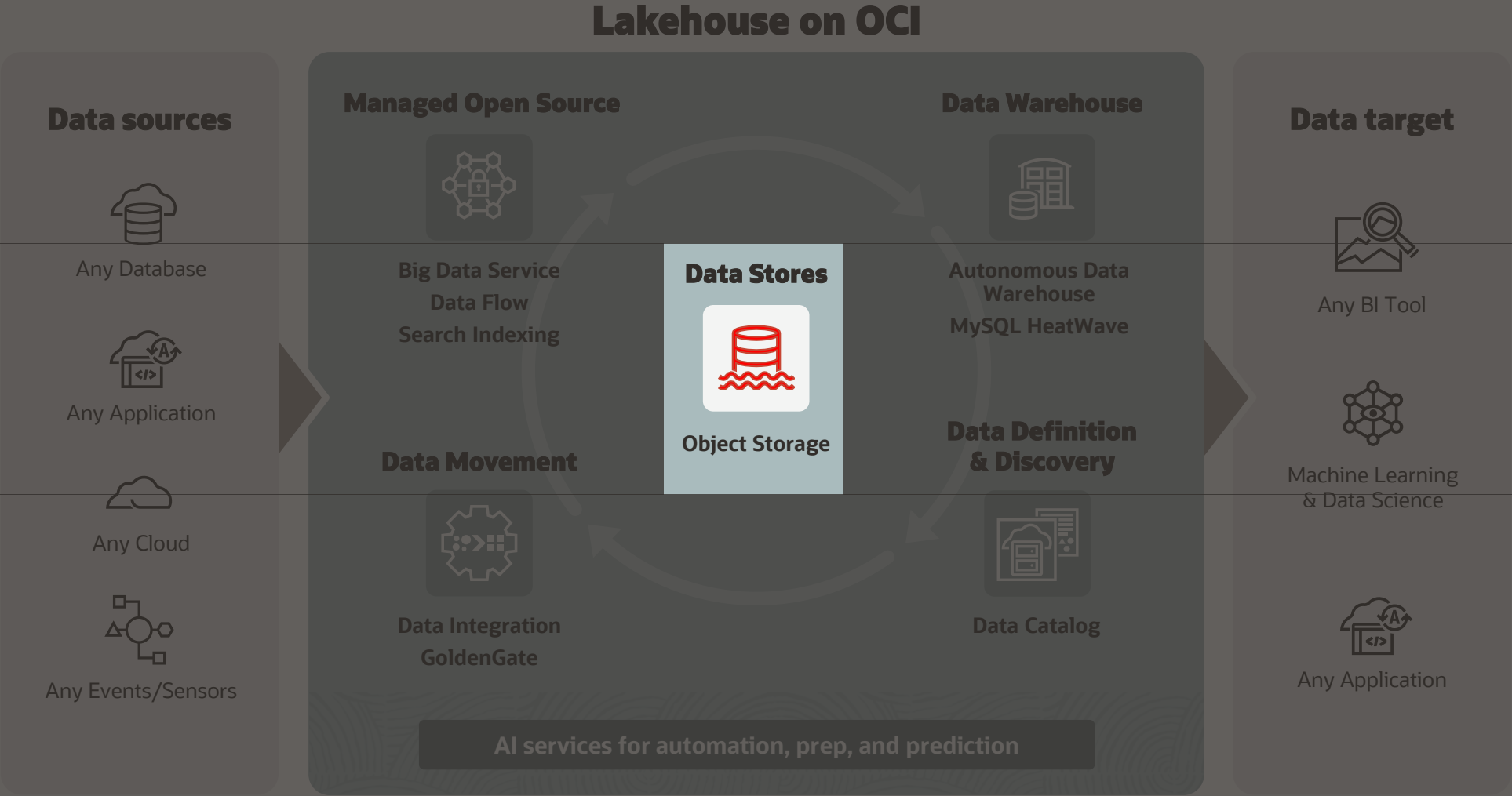


# Construa um Lakehouse em OCI





# Construa um Lakehouse em OCI





Object Storage

# Armazenamento de Dados

## OCI Object Storage

### Agnóstico

Repositório capaz de armazenar qualquer tipo de dado

**Controle**  
Versionamento de objetos

### Rápido

Estrutura baseada discos SSD NVMe

**Flexível**  
Defina a temperatura dos objetos

### Big Data

Construído para atuar como Data Lake



Object Storage

**Integrável**  
Conecte outros serviços presentes na OCI

### Elástico

Cresça até escala de Petabytes

**Seguro**  
Defina políticas de acesso e consumo aos dados

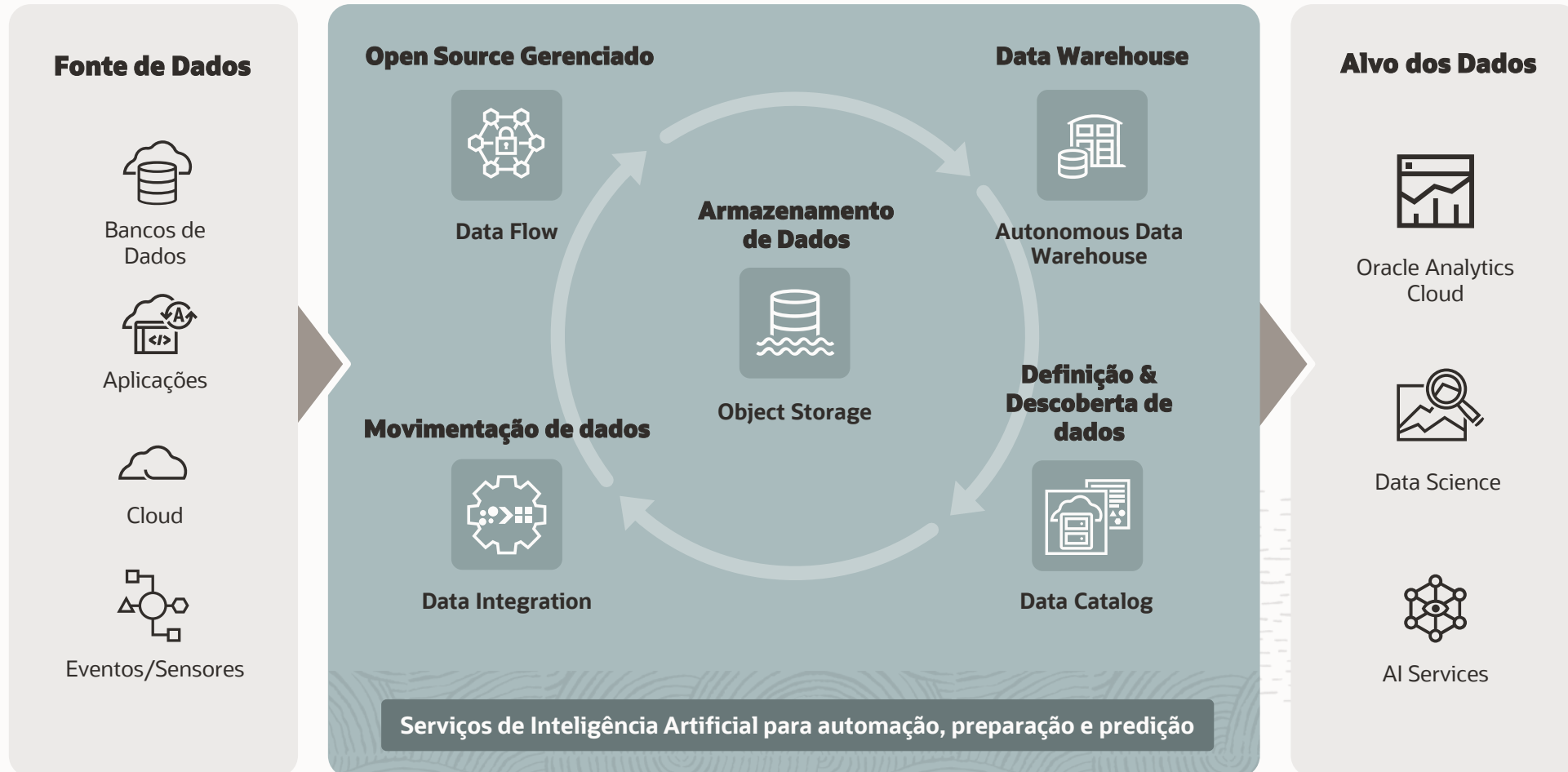




# Oracle Lakehouse

Arquitetura de dados moderna

## Lakehouse em OCI

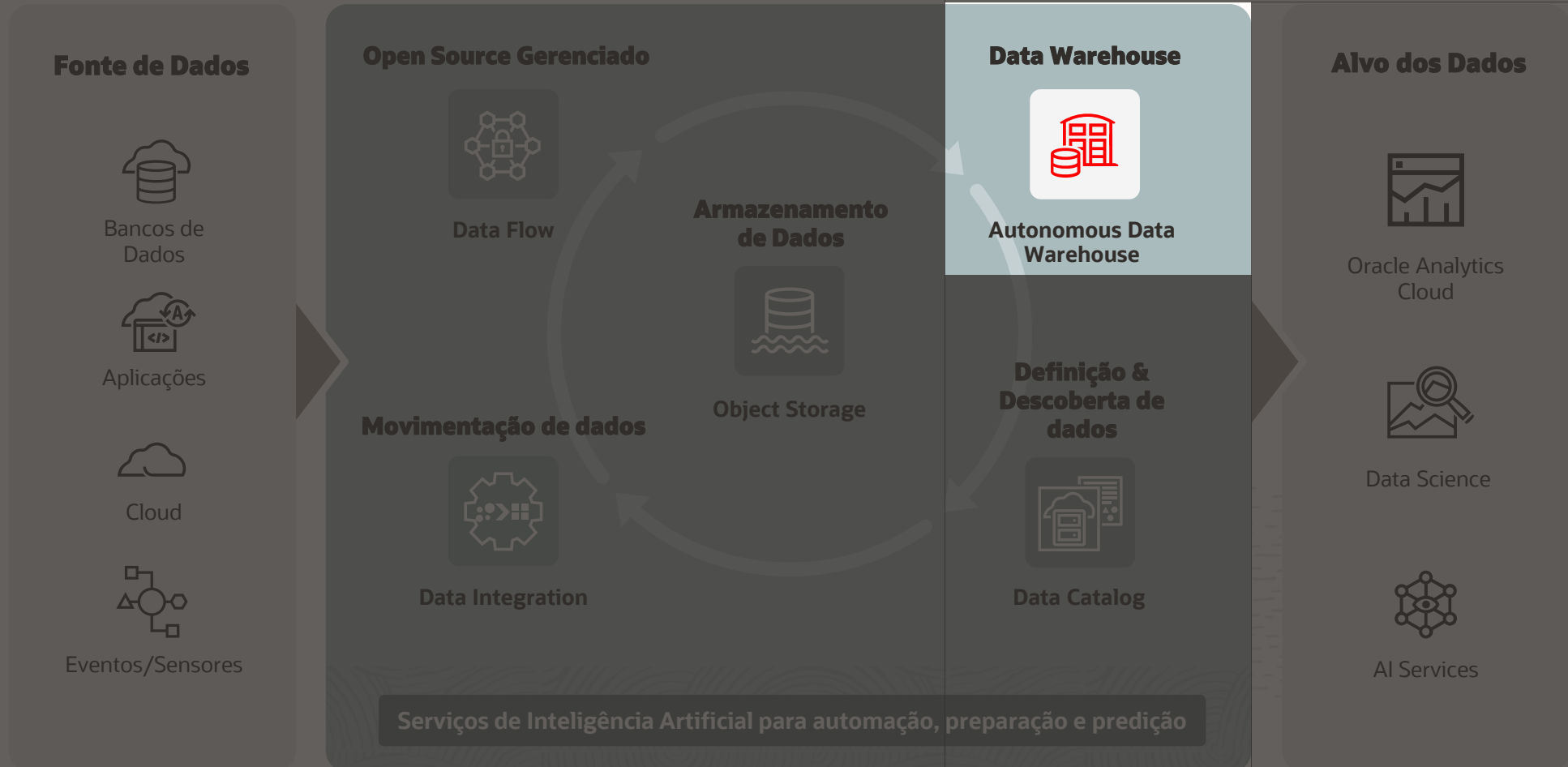


# Oracle Lakehouse

Arquitetura de dados moderna

## Data Warehouse

### Lakehouse em OCI





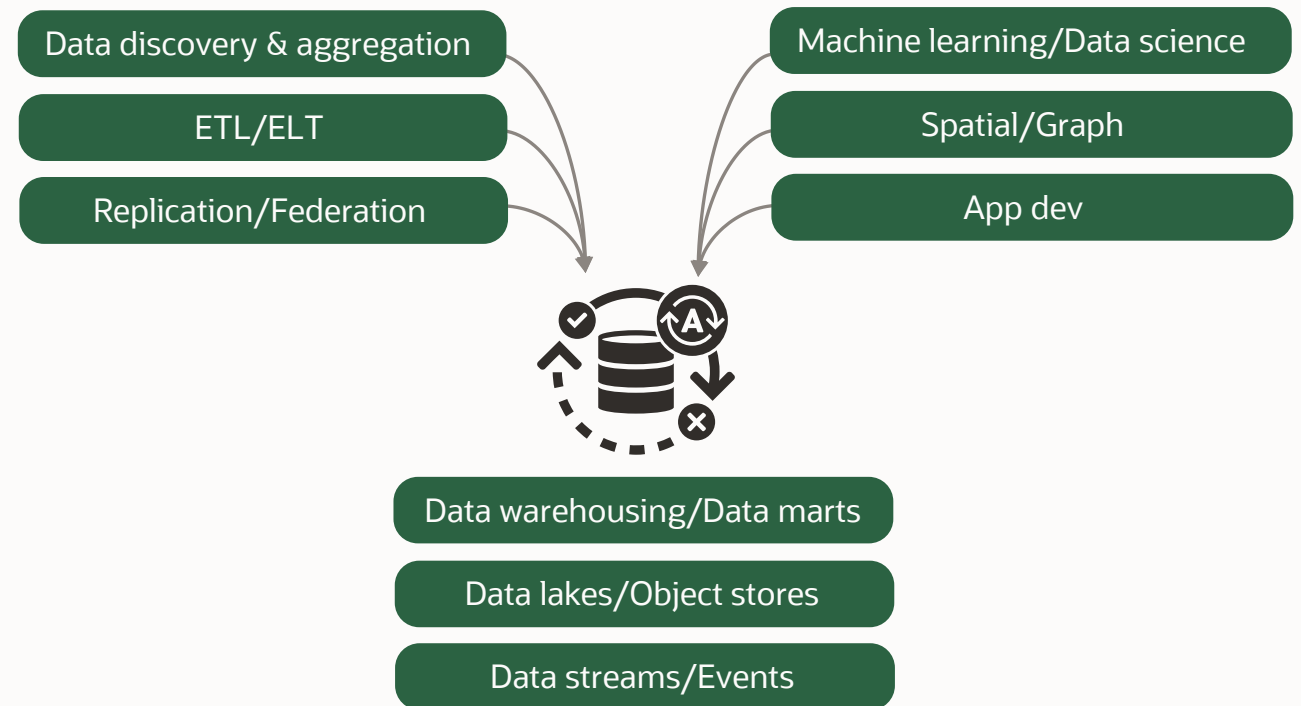
# OCI Autonomous Data Warehouse



Autonomous Data  
Warehouse

- **Gerenciamento automatizado de data warehouse** – sem administração manual
- **Uma solução completa com análises integradas** – mas também suporte para análises de terceiros
- **Implementação fácil e rápida** – carregue todos os dados e obtenha uma variedade de insights
- **Resposta rápida** – independentemente do tamanho dos dados, tipo de análise ou número de usuários simultâneos
- **Proteção abrangente de dados e privacidade** – sem brechas de segurança ou necessidade de serviços adicionais
- **Escolha** provisionar em nuvem pública ou em seu DC

## Autonomous Data Warehouse (ADW)

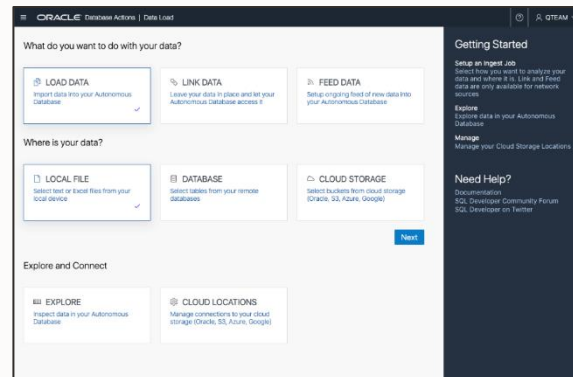


# OCI Autonomous Data Warehouse



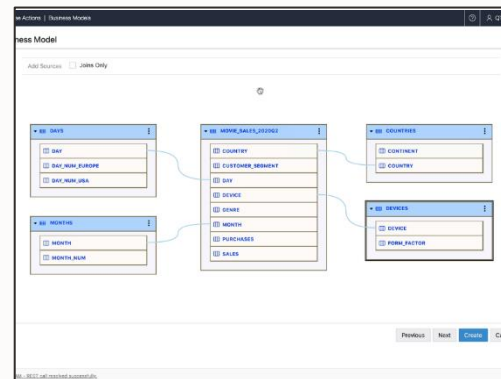
Autonomous Data  
Warehouse

## Load



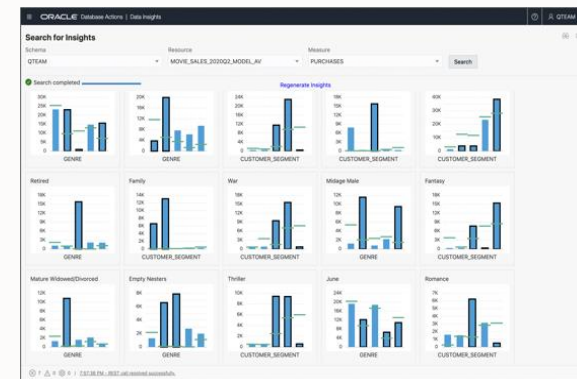
Carga Drag & drop  
simples

## Business Model



Crie Modelos de  
Negócios  
Automaticamente

## Data Insights



Descubra padrões e  
Anomalias  
escondidos  
Automaticamente





# OCI Autonomous Data Warehouse



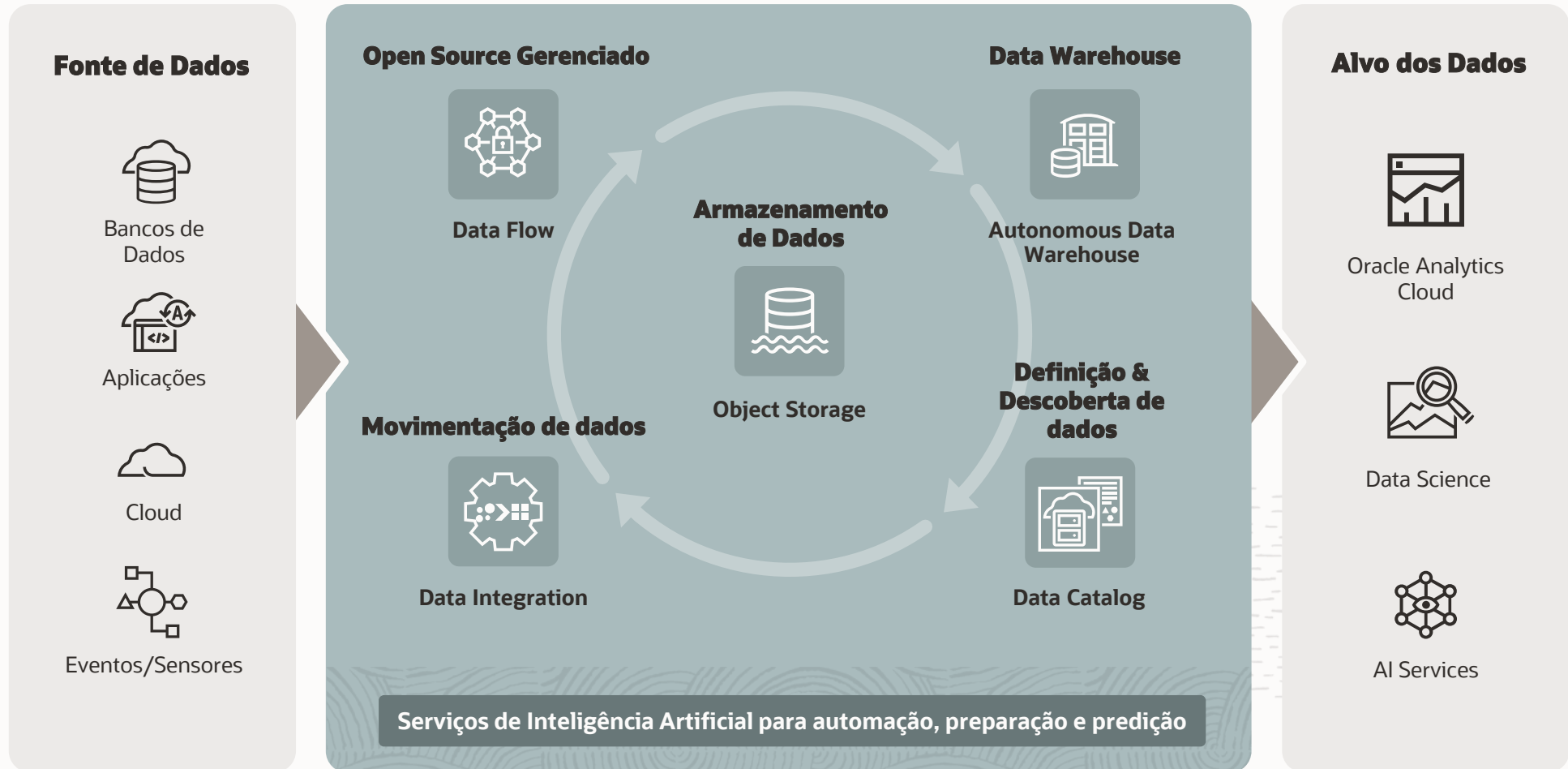
Autonomous Data  
Warehouse



# Oracle Lakehouse

Arquitetura de dados moderna

## Lakehouse em OCI

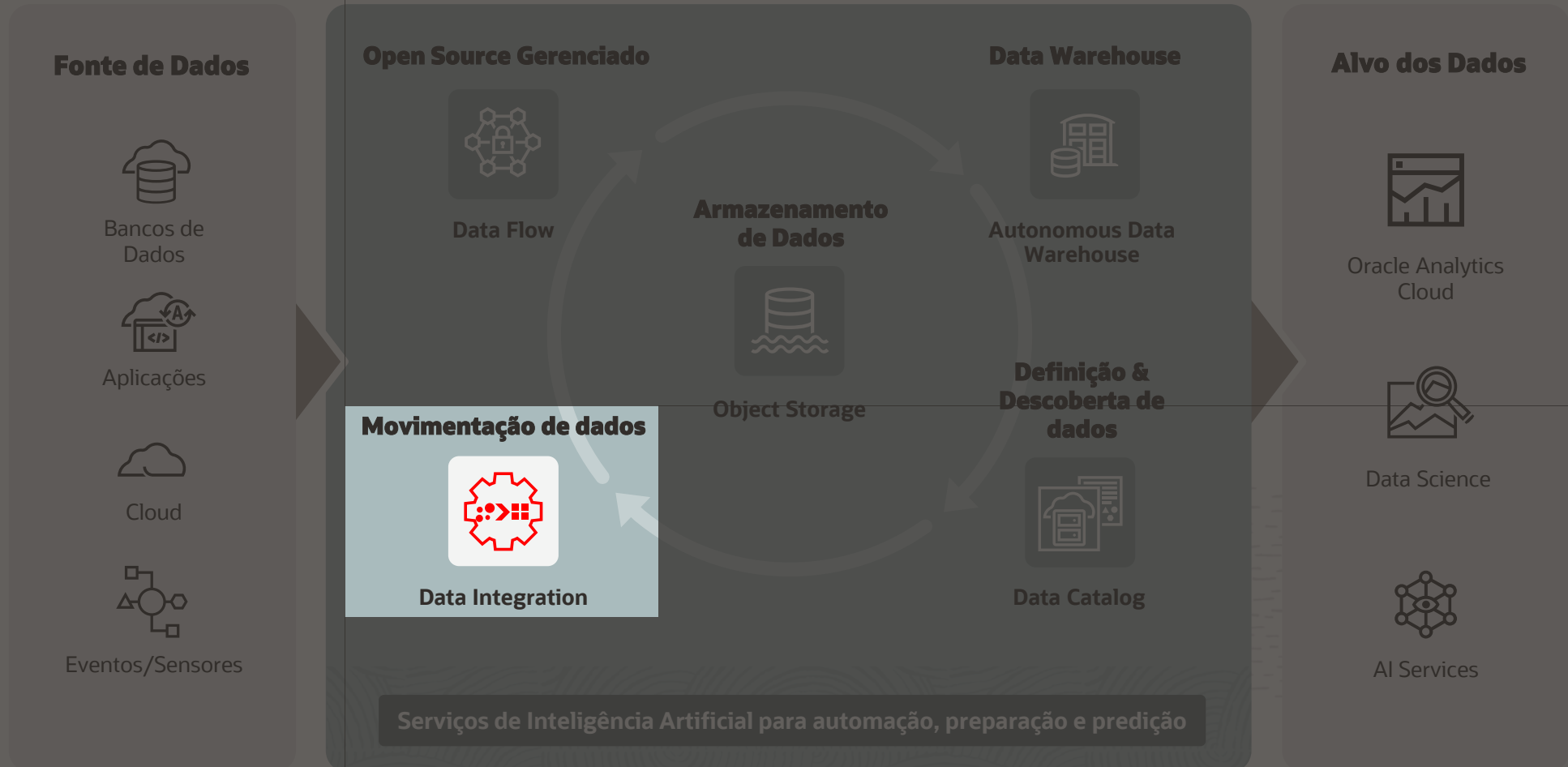


# Oracle Lakehouse

Arquitetura de dados moderna

## Ingestão de Dados

### Lakehouse em OCI



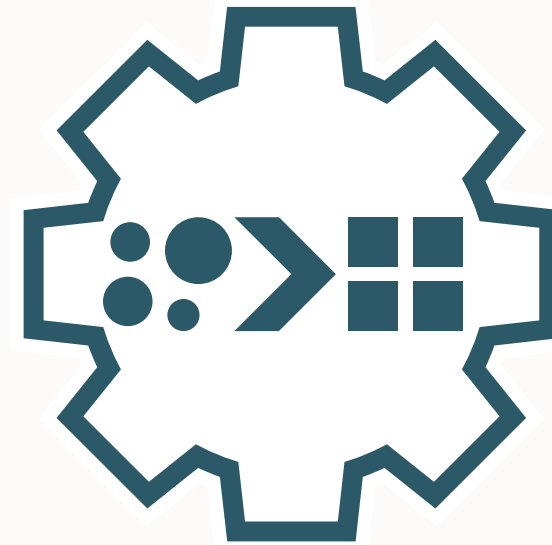




# Oracle Cloud Infrastructure (OCI) Data Integration

**O serviço de ETL nativo em nuvem e serveless para integração, transformação e movimentação de dados dentro do ecossistema OCI**

- Interface gráfica, code-free
- Preparação e criação de perfis de dados interativos
- Proteção de Schema evolution (Drift)
- Powered by Spark ETL ou E-LT Push-Down



## OCI DI



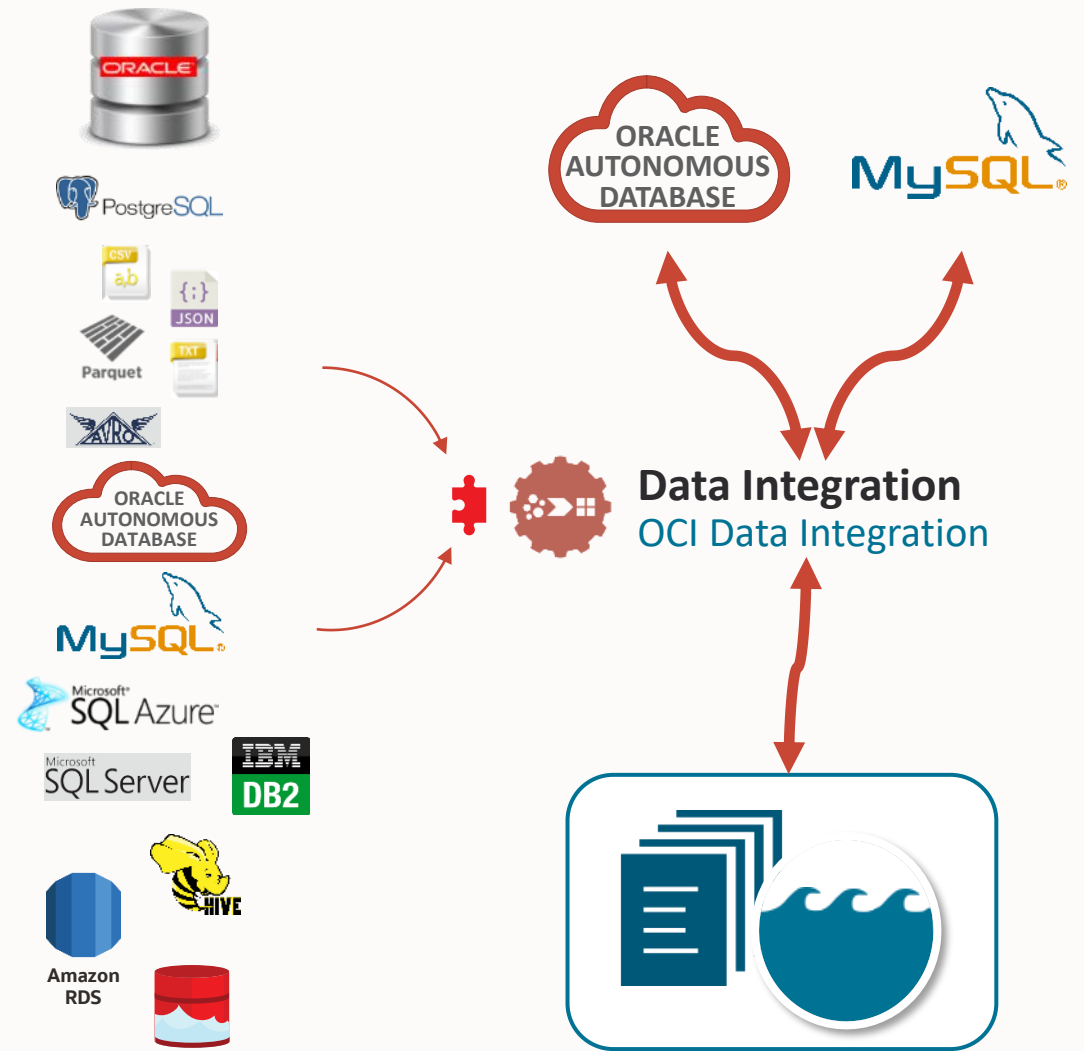
# Ampla Conectividade

Optimized for Oracle Cloud

- Acesso seguro – público e privado – inclusive para on-premises
- Conexão para:
  - Oracle Autonomous Database (ADW/ATP)
  - Oracle Database & Exadata DB Systems
  - Oracle Object Storage: CSV, JSON, Parquet, Avro
  - Oracle Fusion
  - MySQL / PostgreSQL / Apache Hive
  - Microsoft SQL Server & Azure SQL Database & Azure Synapse Analytics
  - Amazon RDS (MySQL, Oracle, Microsoft SQL Server)
  - IBM DB2

## Benefícios

- A melhor conectividade para Oracle Cloud
- Expansão de um conjunto de adaptadores nativos fáceis de usar



Data Integration





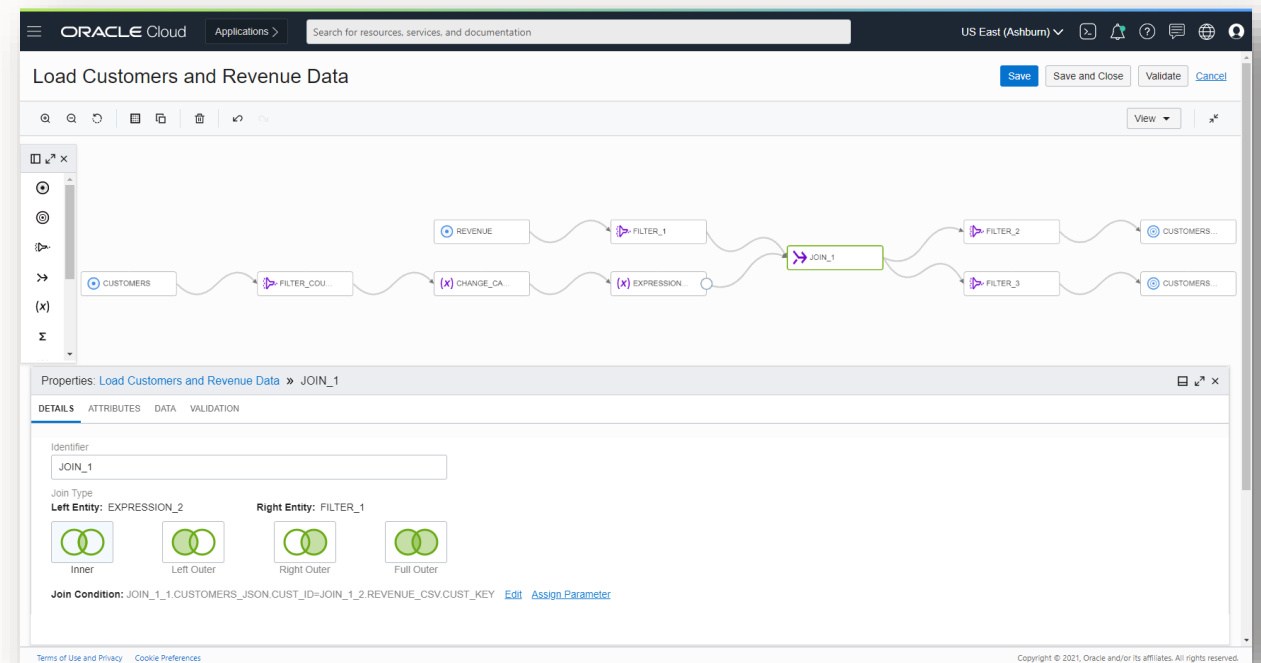
# Maximize a produtividade do desenvolvedor

## No Code Data Flow Design

- Poderoso editor gráfico para criar fluxos de dados (data flows)
- Preveja dados com Data Xplorer
- Fluxos de dados parametrizáveis para máxima flexibilidade

## Benefícios

- Permite o foco em inovação
- Simplifica os processos de ETL / Manutenção
- Permite integrações e transformações poderosas

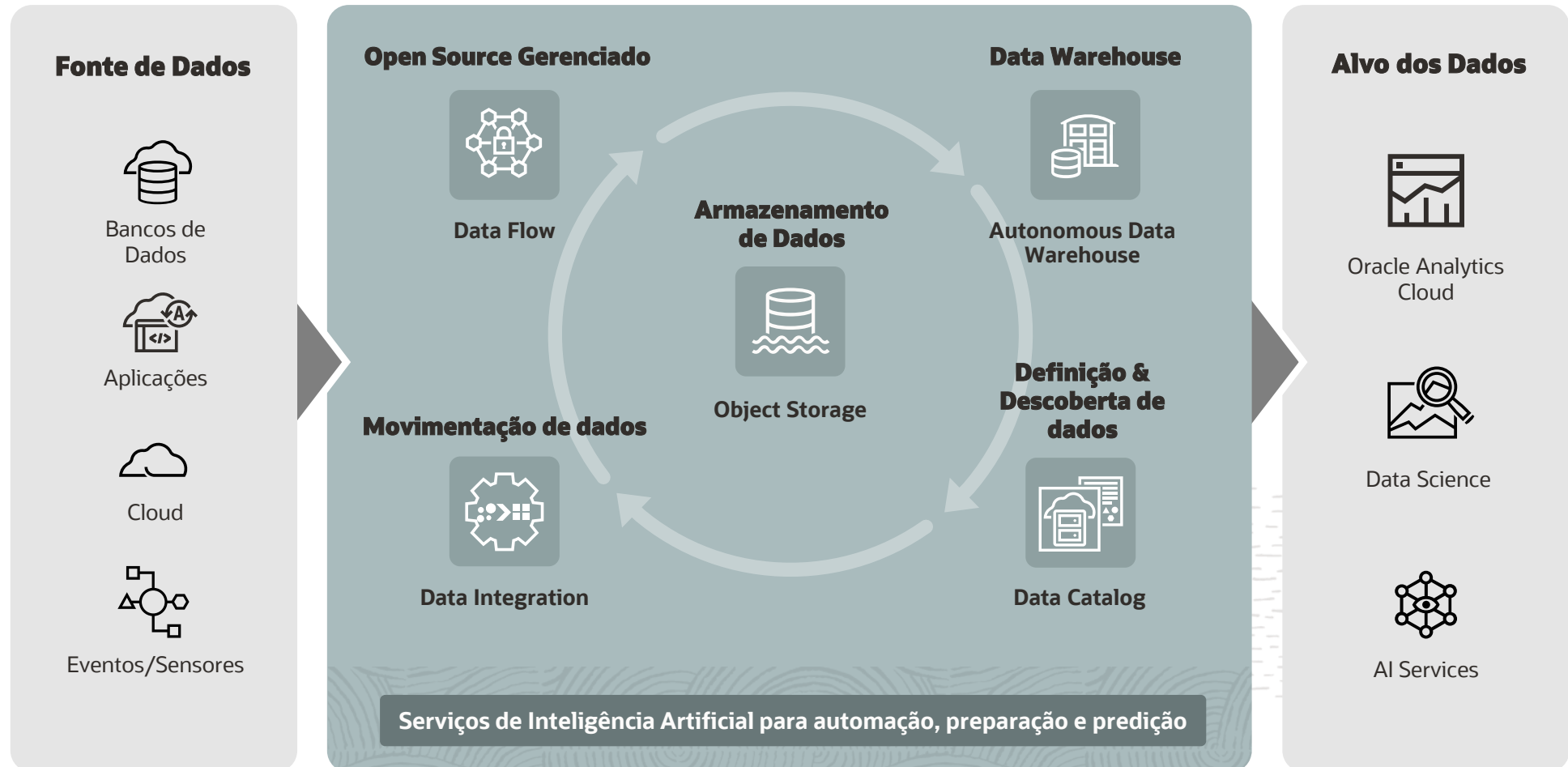




# Oracle Lakehouse

## Arquitetura de dados moderna

### Lakehouse em OCI

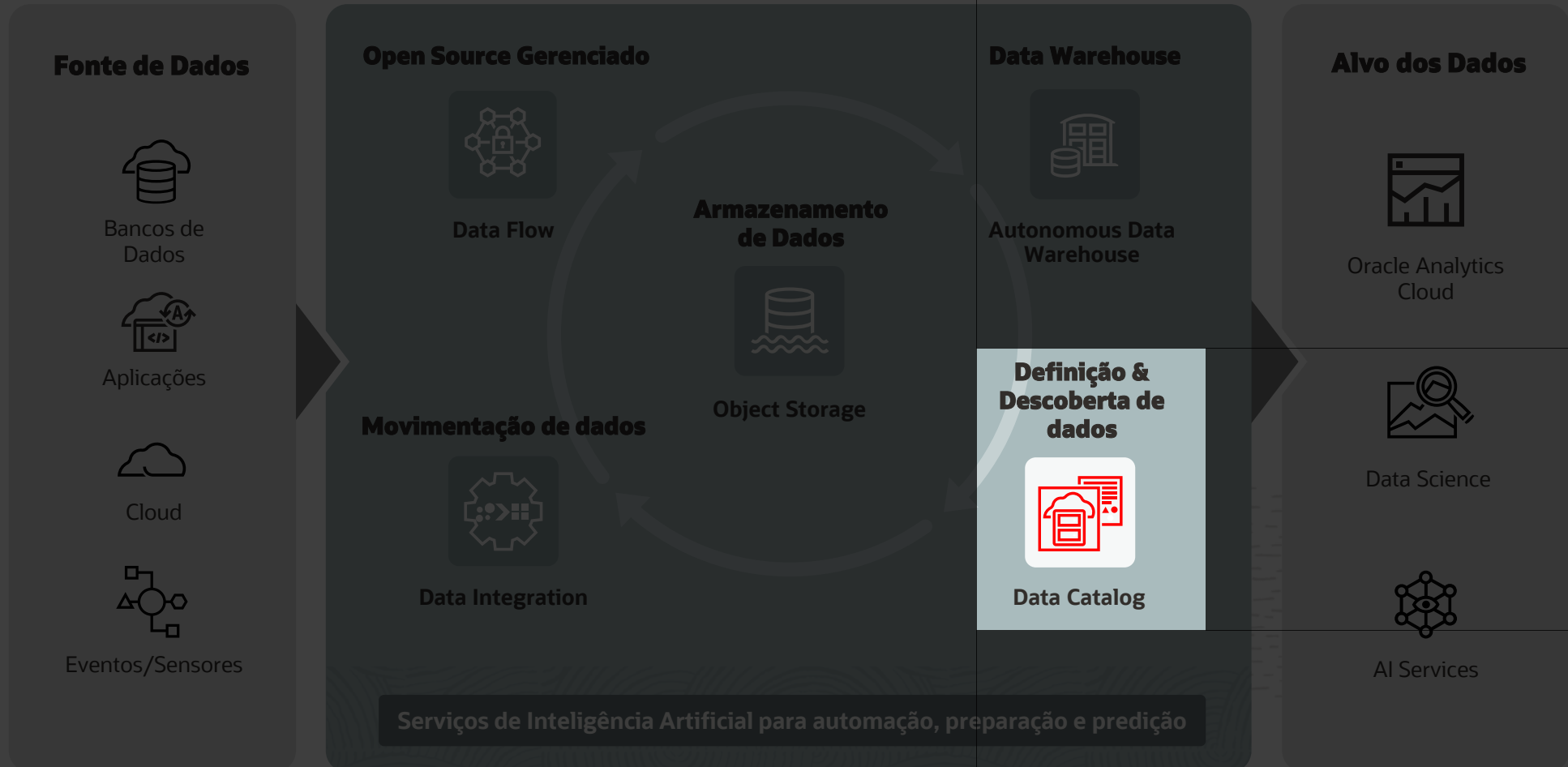


# Oracle Lakehouse

Arquitetura de dados moderna

## Governança de Dados

### Lakehouse em OCI



# Por que organizações precisam de um Catálogo de Dados na Nuvem Oracle?



Data Catalog



Analista de Dados



Cientista de Dados



Engenheiros de Dados



Data Stewards

## Dificuldade de encontrar o dado correto para Analytics

- Falta de visão holística sobre os dados
- Difícil de solucionar problemas de dados
- Poucas informações sobre os dados

## Dificuldade de entender dados em um Data Lake

- Sem dicionário de dados
- Definição manual de schema
- Incapacidade de compartilhar modelos de dados entre aplicativos

## Necessidade de melhora na Governança de Dados

- Falta de visão sobre donos dos dados
- Falta de conceitos de negócios comuns
- Sem colaboração para resolver problemas de dados
- Proliferação de dados sensíveis







# Overview de Capacidades

## Coleta de Metadados

- Coleta automatizada de metadados técnicos
- Suporta fontes de dados na nuvem e on-premise
  - OCI Object Storage, Autonomous Database
  - Oracle DB, MySQL, Hive, Kafka
  - MS SQL Server, Azure SQL DB, IBM DB2, PostgreSQL
- Entidades lógicas para agrupar arquivos relacionados em data lakes

## Curadoria de Metadados

- Glossários de Negócio com Termos e Categorias
- Propriedades de enriquecimento personalizadas definidas pelo usuário
- Tags de formato livre para anotações
- Vincule ativos a termos comerciais, tags

## Busque e Filtre

- Ambiente colaborativo para consumidores de dados
- Pesquisa com base em nomes técnicos, termos comerciais, tags e propriedades personalizadas
- Busque dados por hierarquias do sistema
- Ver detalhes técnicos e de contexto de negócios

## Otimizado para a nuvem Oracle

- Seguro, escalável, serverless, nativo da nuvem
- REST APIs e SDKs em Java, Python, Ruby, Go para integração
- Políticas baseadas em IAM para gerenciamento de controle de acesso

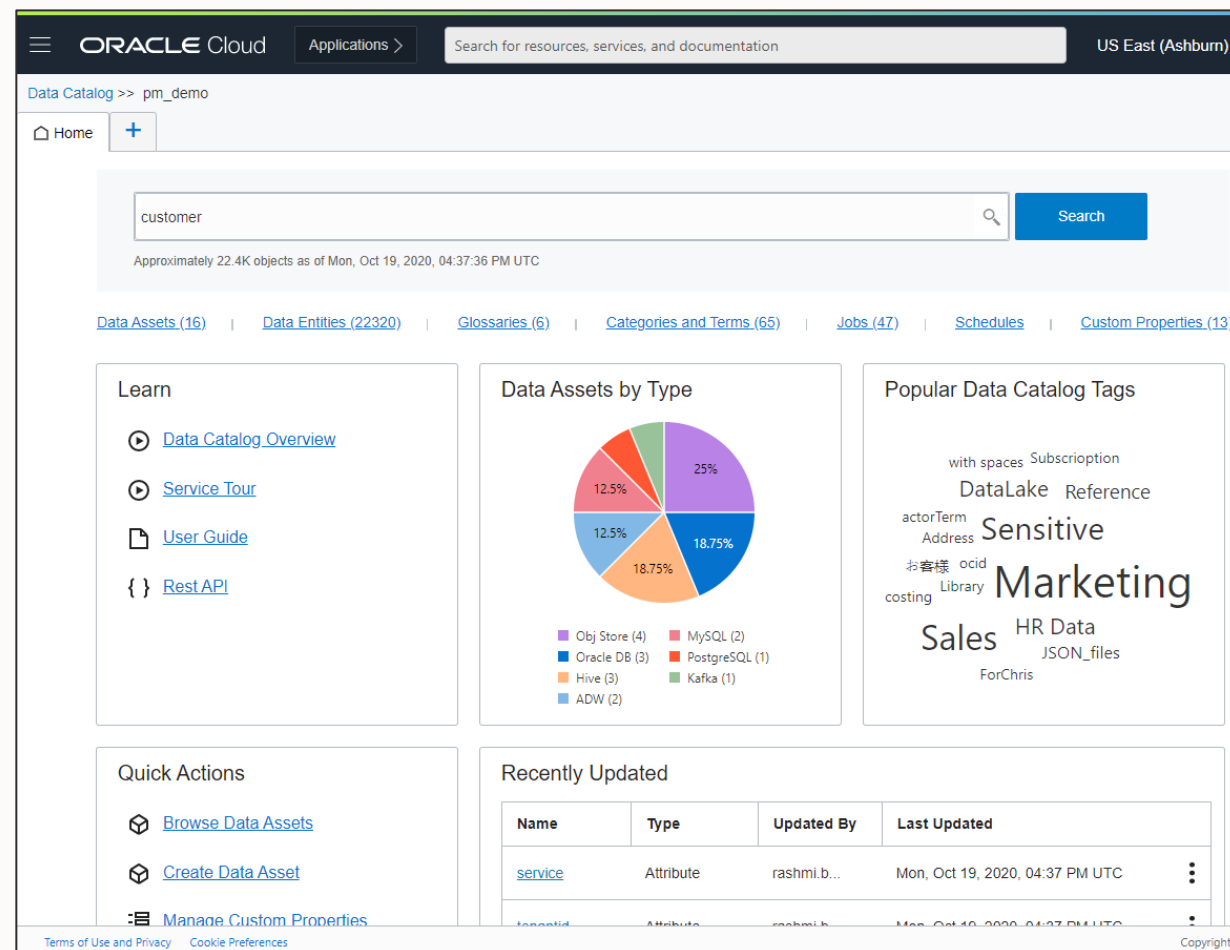
\* Capacidades atuais disponíveis sem custo





# Ambiente Colaborativo

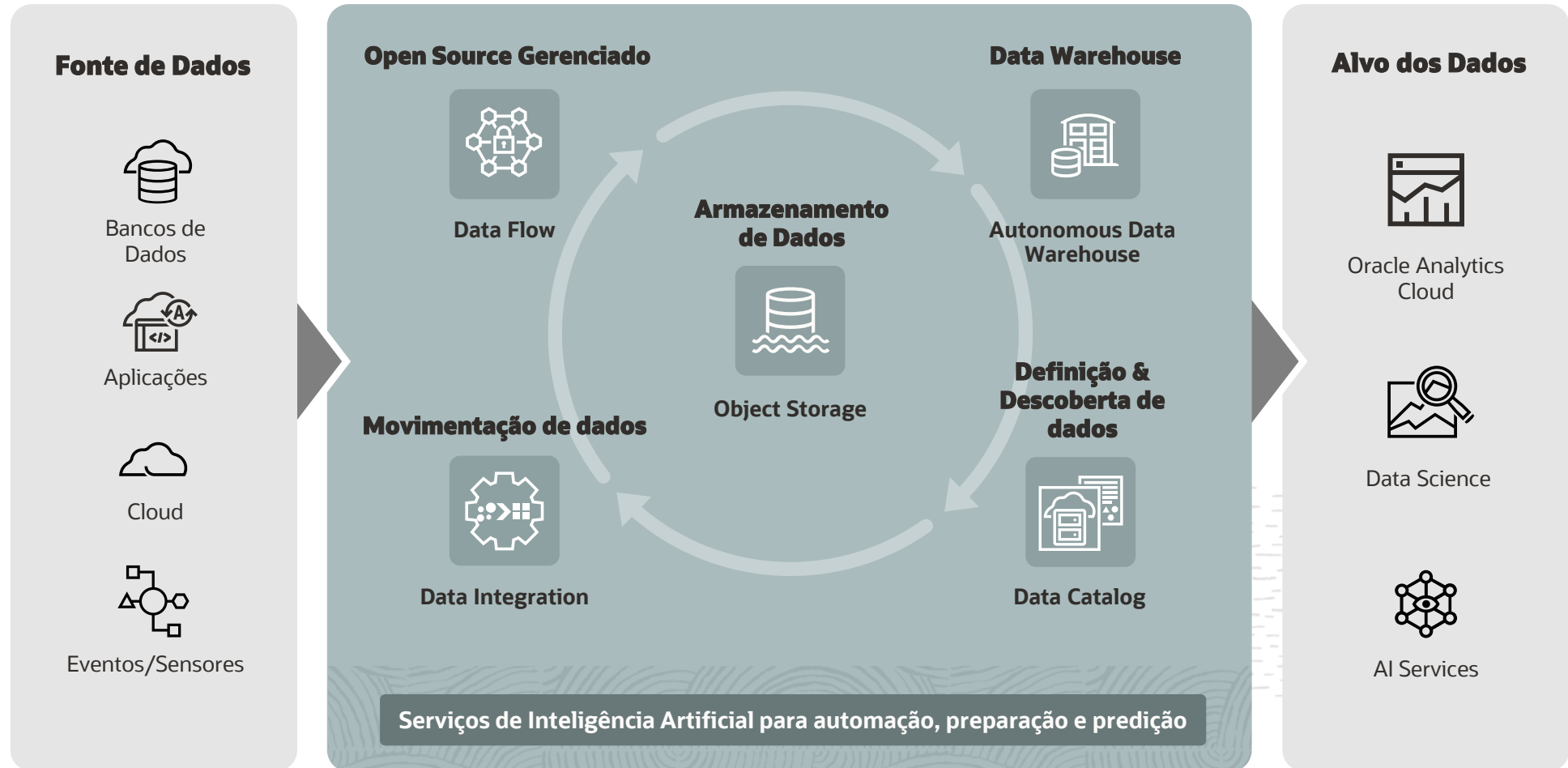
- Página inicial com atalhos e estatísticas operacionais
- Ações rápidas para gerenciar ativos de dados, trabalhos, propriedades personalizadas, padrões de nome de arquivo, etc.
- Utilize as opções de Tags populares e Objetos atualizados recentemente para acesso rápido



# Oracle Lakehouse

Arquitetura de dados moderna

## Lakehouse em OCI





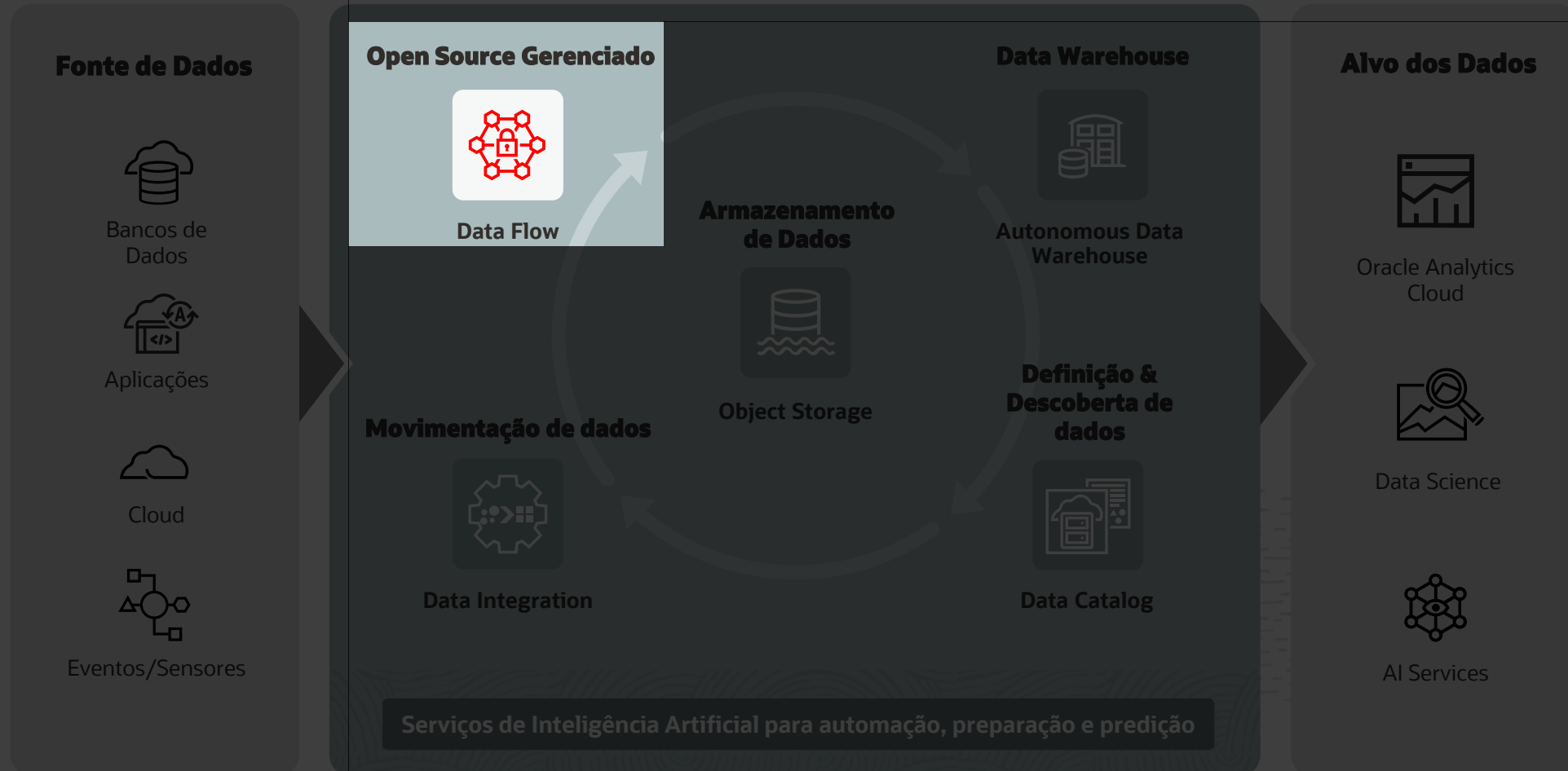
# Oracle Lakehouse

Arquitetura de dados moderna

## Open Source Gerenciado

- Processamento
- Aplicações efemêras
- Ingestão de Dados

### Lakehouse em OCI



# OCI Data Flow



Data Flow

- Data Flow permite entregar **Spark-based** Big Data ou ML applications de forma **rápida**.
- Mantenha o foco em suas aplicações e **esqueça a Infraestrutura**.
- Pague somente pelo que for utilizar e quando utilizar.



# OCI Data Flow



Data Flow

## Spark on Demand

Implante clusters Spark em minutos sem nada para manter



## Operação simplificada

Uis para monitoramento, alertas e diagnósticos.



## Enterprise Security

RBAC e Impersonation para controle total.



## Baixo Custo

Não existe custo pelo serviço, Pague apenas pela Infraestrutura utilizada.



# Data Flow torna a operação mais simples.



## Operations

- Sort and filter to quickly find the biggest jobs.
- Find out who is consuming the most.
- Stop jobs that are running too long.

Data Flow

Applications

**Runs**

List Scope

COMPARTMENT

Marketing

paasdevssstest (root)/Marketing

Tag Filters

add | clear

no tag filters applied

Filters

STATE

Succeeded

LANGUAGE

All

CREATED START DATE

Aug 1, 2019 00:00 UTC

CREATED END DATE

Oct 31, 2019 00:00 UTC

Runs in Marketing Compartment

Name	Language	State	Owner	Created	Duration	Total oCPU	Data Read	Data Written	
<a href="#">PSR_AirlineETL</a>	SQL	Succeeded	user@example.com	Fri, Aug 16, 2019, 21:46:38 UTC	15m	22	151 GB	6 GB	⋮
<a href="#">PSR_Scala_test</a>	Scala	Succeeded	user@example.com	Thu, Aug 15, 2019, 18:20:59 UTC	1h 49m 7s	22	62 GB	92 MB	⋮
<a href="#">Simple Test App</a>	Java	Succeeded	user@example.com	Mon, Oct 7, 2019, 21:30:33 UTC	1m 40s	2	68 MB	31 MB	⋮
<a href="#">PSRTESTING_DJJWKEXEPCESGA</a>	Java	Succeeded	user@example.com	Thu, Oct 17, 2019, 09:30:17 UTC	1m 14s	4	68 MB	31 MB	⋮
<a href="#">PSRTESTING_BRNOFFVEEQONTL</a>	Java	Succeeded	user@example.com	Thu, Oct 17, 2019, 09:00:02 UTC	1m 11s	4	68 MB	31 MB	⋮
<a href="#">Simple Test App1</a>	Java	Succeeded	user@example.com	Thu, Oct 17, 2019, 07:53:25 UTC	1m 14s	4	68 MB	31 MB	⋮
<a href="#">PSRTESTING_TXDCZCCQOISEKO</a>	Java	Succeeded	user@example.com	Thu, Oct 17, 2019, 08:49:06 UTC	1m 13s	4	68 MB	31 MB	⋮
<a href="#">PSRTESTING_GUYKXIINPUXPJZ</a>	Java	Succeeded	user@example.com	Thu, Oct 17, 2019, 09:11:32 UTC	1m 12s	4	68 MB	31 MB	⋮
<a href="#">Simple Test App</a>	Java	Succeeded	user@example.com	Mon, Oct 7, 2019, 22:38:33 UTC	1m 41s	2	68 MB	31 MB	⋮
<a href="#">Simple Test App1</a>	Java	Succeeded	user@example.com	Thu, Oct 17, 2019, 07:14:59 UTC	1m 14s	4	68 MB	31 MB	⋮

Showing 10 Items < Page 1 >

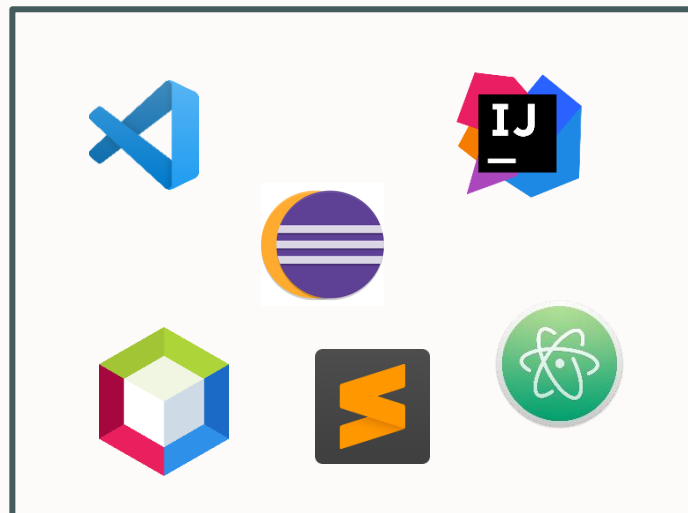




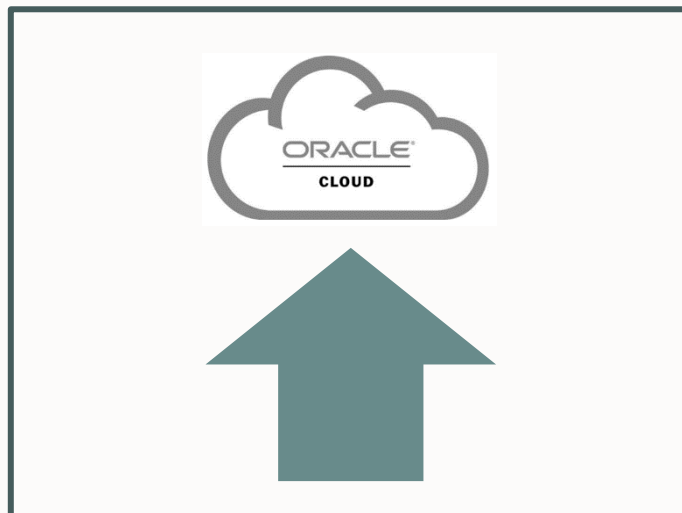
# OCI Data Flow



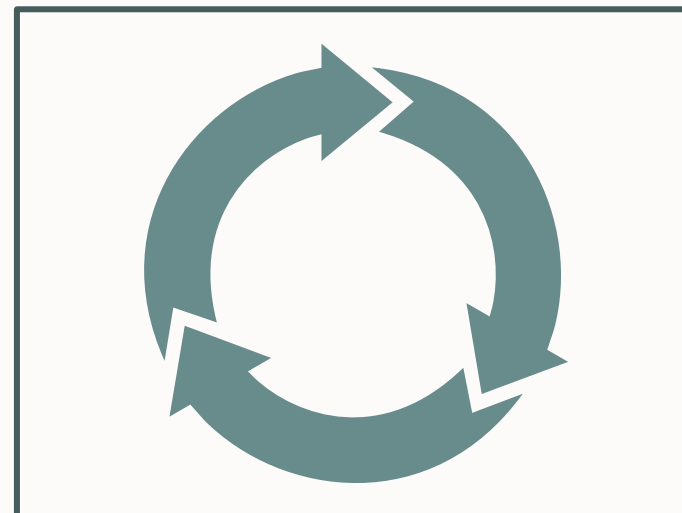
Data Flow



Desenvolva e teste aplicativos Spark localmente usando o IDE de sua escolha.



Carregue aplicativos como estão, sem modificações.



Execute na nuvem em qualquer escala.





# OCI Data Flow

## Convertendo CSV para Parquet

- Validação de Pré-requisitos
- Edição do script Python
- Criação da Data Flow Application
- Execução da aplicação
- Validação do resultado

Hands-on – Lab Data Flow



**Data Flow**

