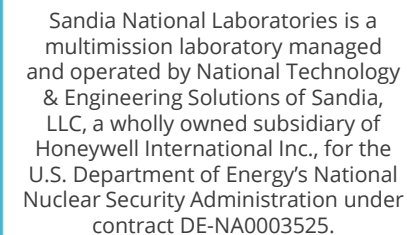




October 17, 2022

## Dr. Ian Brooks - Cloudera Government



# Presenters



## Dr. Patrick Carlson

- PhD graduate in Human-Computer Interaction (HCI) from Iowa State University
- Worked at Sandia National Labs as a Data Scientist since 2017
- Leads the Sandia Data Sciences Community of Practice (CoP)
- Interested in the technical aspects of Data Science such as data engineering, programming, analysis, modeling, machine learning, productionization, ml-ops, data-ops, and more



## Dr. Ian Brooks

- PhD graduate in Computer Science from the University of North Texas
- Focused on Big Data solutions since 2015
- Previous industry roles include Software Engineer, Data Architect, and Data Scientist
- At Cloudera, he is a Principal Solutions Engineer on the Public Sector team and a Machine Learning SME



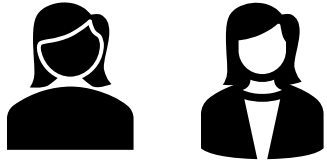
# Why Sandia Insights



- Reduce duplication and rework
- Faster turnaround time



## Data Scientists and Data Engineers



- Rapid access to data
- Robust tooling
- Common environment with libraries preinstalled
- HPC and model training

## Data Governance



- Data security
- Centralized access controls
- Central data catalog management

## Leadership



- Data-driven decision making

# Vision for Sandia Insights



- Empower Data Scientists and Analysts through easy access to vetted enterprise datasets
- Secure and document data through Data Governance and a data catalog
- Productionize and operationalize models, use data-ops and ml-ops practices
- Build out a pipeline and best practices for the entire data-lifecycle
- Focus on people and methods, not specific tools

# History of Sandia Insights



- 2015 – Tableau Server
- 2016 – Analytics for Sandia Knowledge (ASK)
  - Custom developed catalog
  - Tableau visualizations and reports
  - Vetted enterprise datasets via Data Central
  - Expertise finder and co-author network
- 2018 – CKAN data catalog
- 2018 – RStudio Connect (previously called Shiny Server)
- 2019 – Cloudera installation
- 2020 – Collibra data catalog

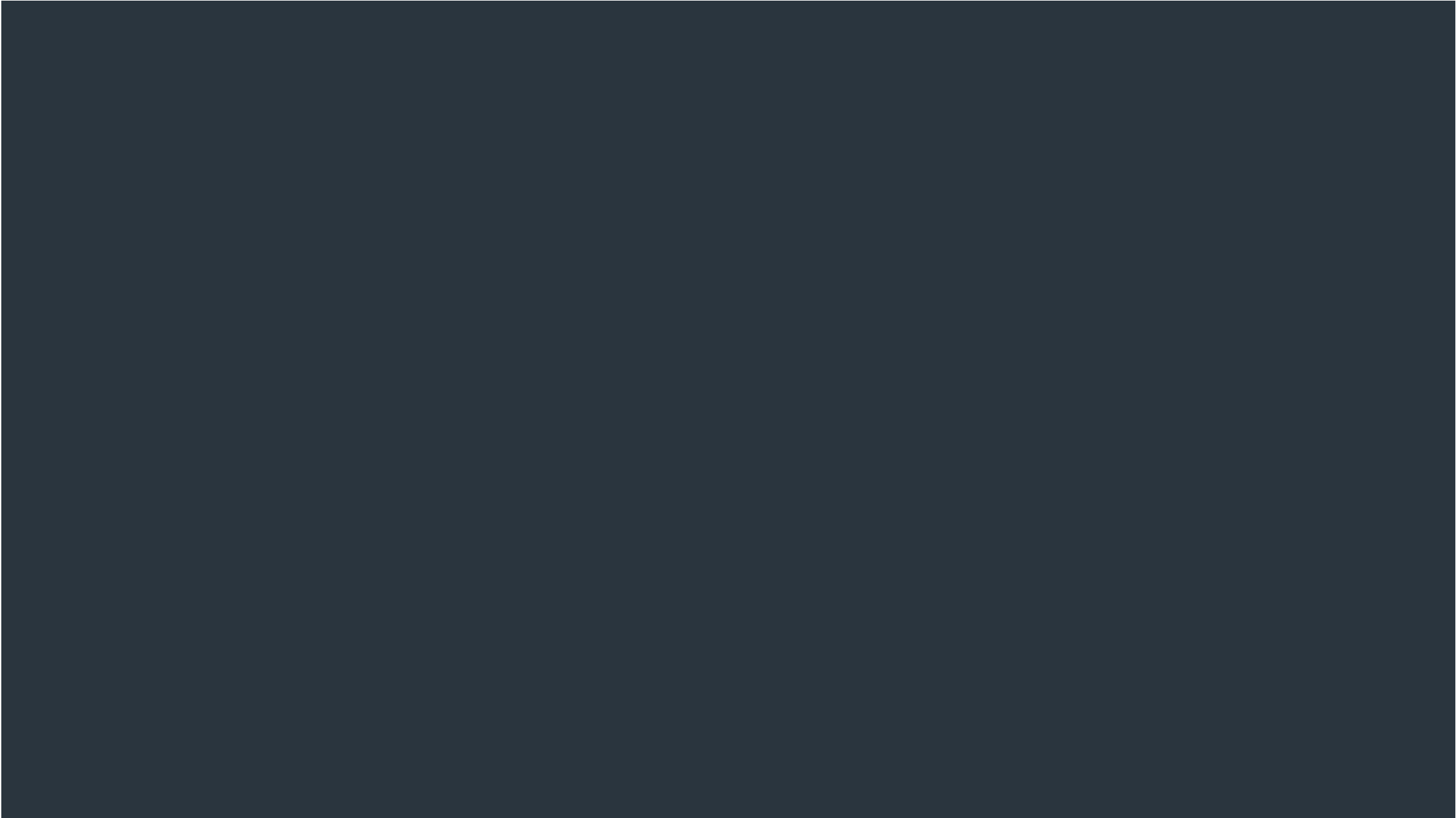


 **Studio** Connect

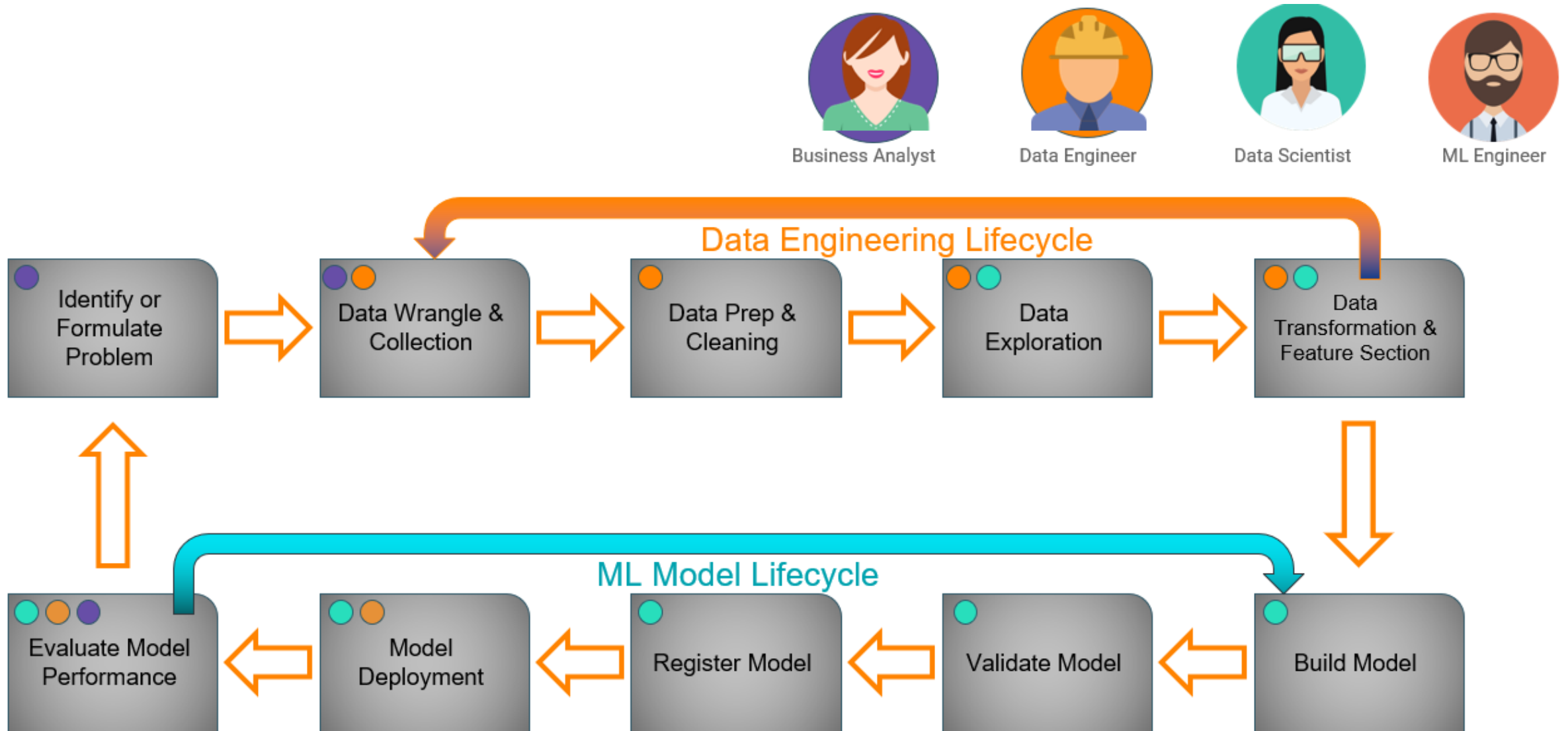
**CLOUDERA**



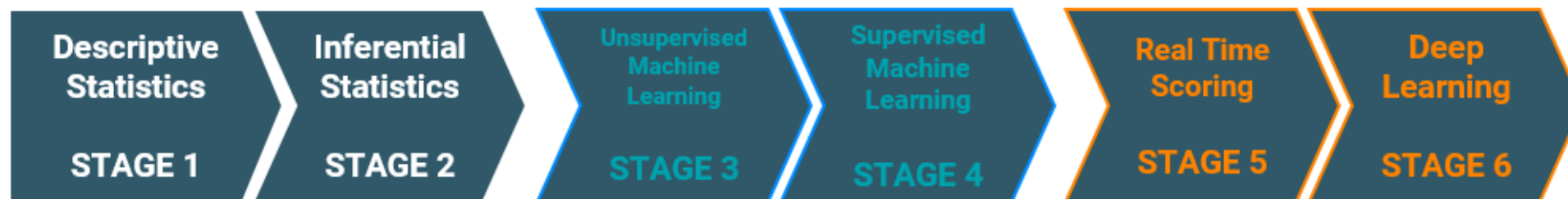
# Current State



# Advanced Analytics Data Lifecycle

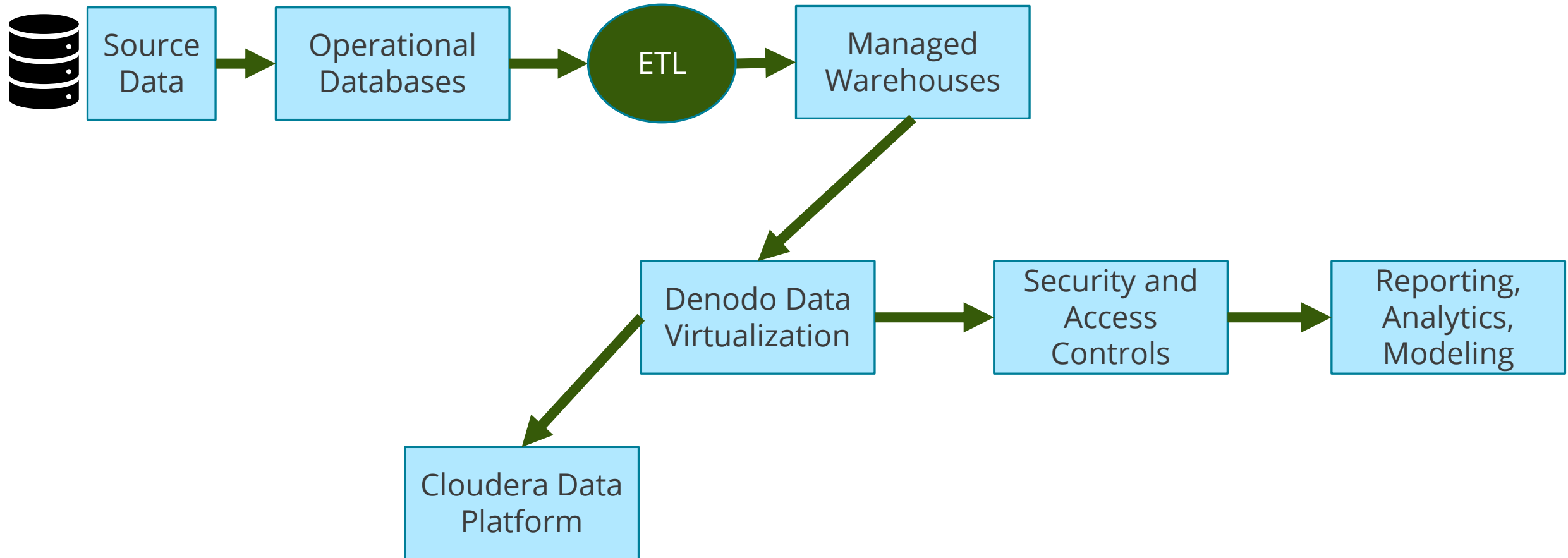


# Progression of Enterprise Data Science Practices



	Phase 1	Phase 2	Phase 3
Goals	Gain Insight Into Data	Enable Proactive Approach	Enable Sub Second Decision Making
Deliverables	Reporting / Charts	ML Models	DL Models and Real Time Scoring Engine
Tools	SQL, CLI, and BI	Notebooks, IDE and ML Libraries	Model Management, Streaming Data Management, RT Scoring Engines, DL Libraries, GPUs

# Current Data Pipeline



# Challenges



- Existing legacy ETL process
- Batch ETL process is well suited for slowly changing datasets or tables (once a week updates)
- Batch ETL process ties up system for many hours, prevents usage of data for reporting and analytics
- Existing process is rigid and prevents agile development

# Future State

# Productionization of Data Science and Collaboration



- ml-ops – productionizing machine learning applications
- Cloudera Data Science Workbench (CDSW) – containerized environment for Data Scientists, collaborative development
- Future MLFlow support for experiment tracking and hyper-parameter tuning

```

quakes.r
10 ## Setting Up Spark Context
11 sc <- spark_connect(master = "yarn",
12                      spark.shuffle.service.enabled = "True",
13                      spark.dynamicAllocation.enabled = "True",
14                      spark.sql.broadcastTimeout = "1200",
15                      app_name = "Quakes!")
16 sc
17
18 ## Load Training and Test Data Frames
19 #train <- spark_read_csv(sc, name="train", path = "hdfs://tmp/quake/train/train.csv", header
20 #train <- spark_read_csv(sc, name="train", path = "hdfs://tmp/quake/train/train_sample.csv",
21
22 #test <- spark_read_csv(sc, name="test", path = "hdfs://tmp/quake/test/*.csv", infer_schem
23 #test <- spark_read_csv(sc, name="test", path = "hdfs://tmp/quake/test/kag_000307.csv", infer
24
25 summary(train)
26 summary(test)
27
28 ## Show Tables
29 src_tbls(sc)
30 sdf_schema(test)
31 sdf_schema(train)
32
33 ## Display Dataframe Values
34 summarise(train, acoustic_data, time_to_failure)
35
36 featureList <- list('acoustic_data')
37
38 ## Linear Regression Model
39 linearReg_pipeline <- ml_pipeline(sc) %>%
40   ft_vector_assembler(
41     input_cols = featureList,
42     output_col = "feature_vector"
43   ) %>%
44   ft_standard_scaler(
45     input_col = "feature_vector",
46     output_col = "scaled_features"
47   ) %>%
48   ft_formula(formula = "time_to_failure ~ acoustic_data",
49              prediction_col = "predicted_failure_time",
50              label_col = "time_to_failure"
51             ) %>%
52   ml_linear_regression(features_col = "scaled_features",
53                       label_col = "time_to_failure",
54                       standardization = TRUE)
55
56 linearReg_pipeline
57
58 #fitted_pipeline <- ml_fit(linearReg_pipeline, train)
  
```



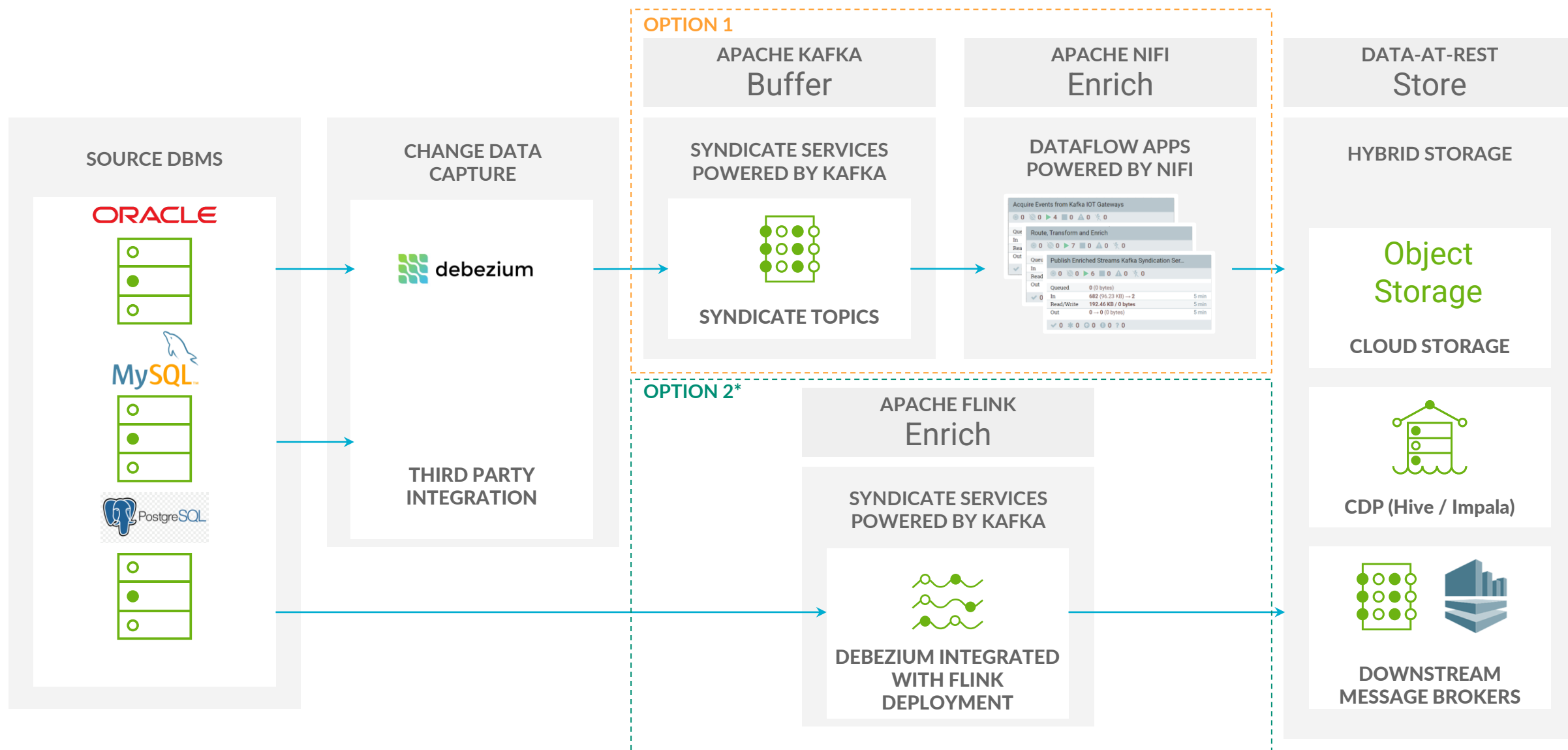
# Data Governance and Cataloging



- Document more datasets with descriptions and metadata
- Leverage Denodo as singular integration point for structured data and access controls
- Tagging data in Denodo and Cloudera (Apache Atlas)
- Test attribute-based access controls in Denodo and Cloudera (Apache Ranger)
- Working on labels and sensitivity categories to address Controlled Unclassified Information (CUI) requirements



# Change Data Capture (CDC)

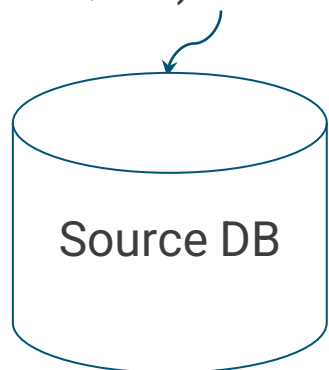


# Integrated Change Data Capture



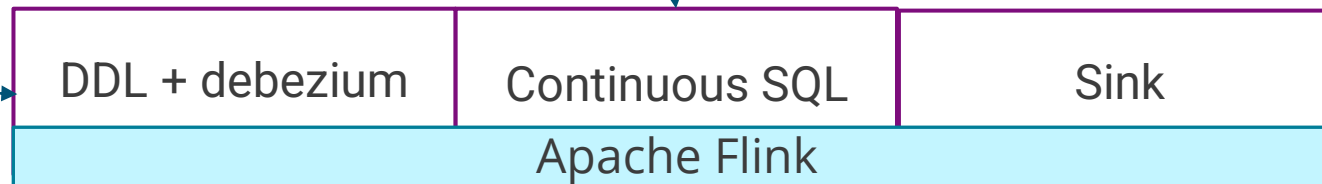
**Key Takeaway; Run CDC right from Flink**  
**note: no Kafka between Flink and DB, no connectors, etc.**

Capture source changes in native change log (WAL, Logminer, etc)



SQL queries on CDC stream directly.

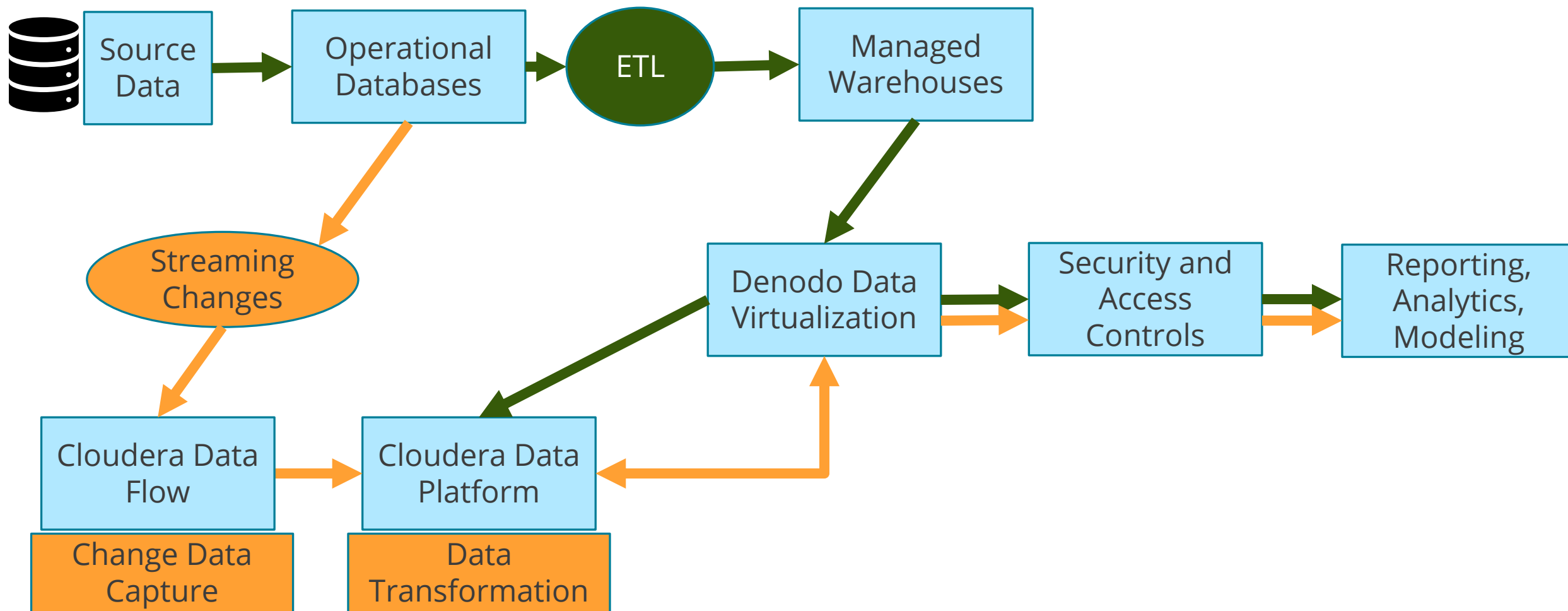
Route output to any sink, including Kafka or even MV's



Debezium as a library running *\*inside\** Flink, capturing changes as a direct source

- Oracle, PostgreSQL, MySQL, DB2, MSSQL
- Logminer (open) and Xstream (licensed) interfaces for Oracle
- Uses DDL Interface directly, one DDL configuration per CDC stream.
- Can specify query on capture

# Sandia Automated Data Pipeline



# Demo or Video of CDC



# Impact of Automated Data Pipelines



- Rapid access to current data enables reporting and analytics access to quickly changing datasets or tables
- Providing robust tooling to augment and not interfere with existing data pipelines such as ETL
- More agile development of future data pipelines
- Ensure end-to-end security and governance throughout entire pipeline
- Provide leadership with real-time dashboards and decision-making opportunities

# Lessons Learned and Best Practices



Challenges	Actions
Culture of “ownership” of data	Implement Data Governance
Pockets of Data Scientists around the lab	Data Sciences Community of Practice (CoP)
Lack of leadership support and direction	Get a Chief Data Office (CDO) or Deputy Chief Data Officer
Lack of documentation around data	Implement a data catalog
Data architecture is old	Understand architectural paradigms: ETL, ELT, data warehouse, data lake, lambda and kappa architectures, etc.

## PRIORITIZE:

People  
Processes  
Architectures  
Culture  
Results



Tools  
Software  
One-Off Projects

# Comments or Questions?



[pcarloso@sandia.gov](mailto:pcarloso@sandia.gov)



<https://carlsonp.github.io>

<https://www.linkedin.com/in/patrick-carlson-13753628>

[ibrooks@cloudera.com](mailto:ibrooks@cloudera.com)



<https://github.com/Brookslan>

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