# Advanced Analytics & AI

**Semih Eskin**
Principal Technologist

Data Lakehouse

# Storage Disaggregation

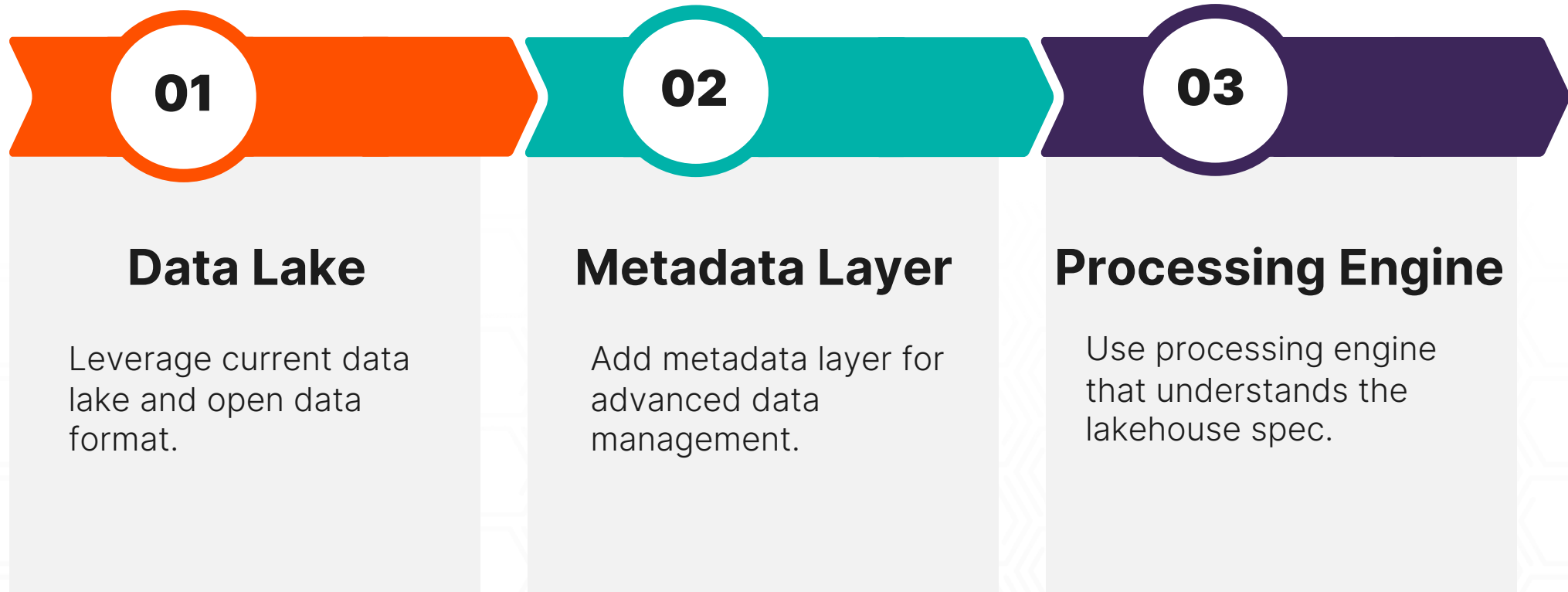**DIRECT ATTACHED STORAGE**

**DISAGGREGATED STORAGE AND COMPUTE**

**INDEPENDENTLY** Scale Compute & Storage
**REDUCE** Infrastructure Expense
**INCREASED** Compute Performance
**DECREASE** Software Spend
**FASTER** Recovery Time
**CONSISTENT** Application Performance

**PURE**STORAGE®

# Three Key Components in a Data Lakehouse

A simplified view of how to implement an open data lakehouse

## 01
### Data Lake

Leverage current data lake and open data format.

## 02
### Metadata Layer

Add metadata layer for advanced data management.

## 03
### Processing Engine

Use processing engine that understands the lakehouse spec.

# The Data Lakehouse Architecture

**Unites the strength of data lakes and data warehouse**

- Supports ETL, SQL, and ML workloads
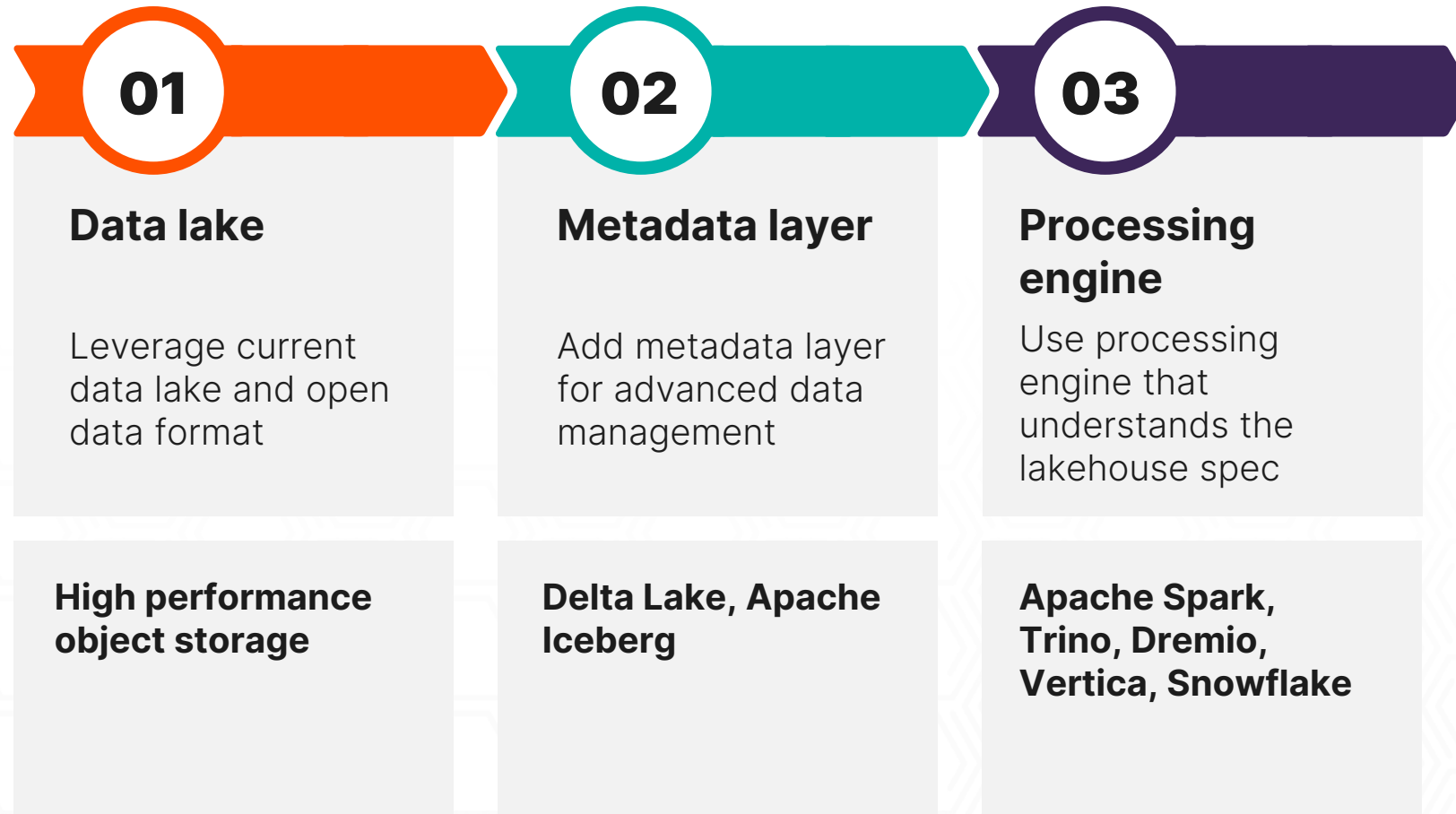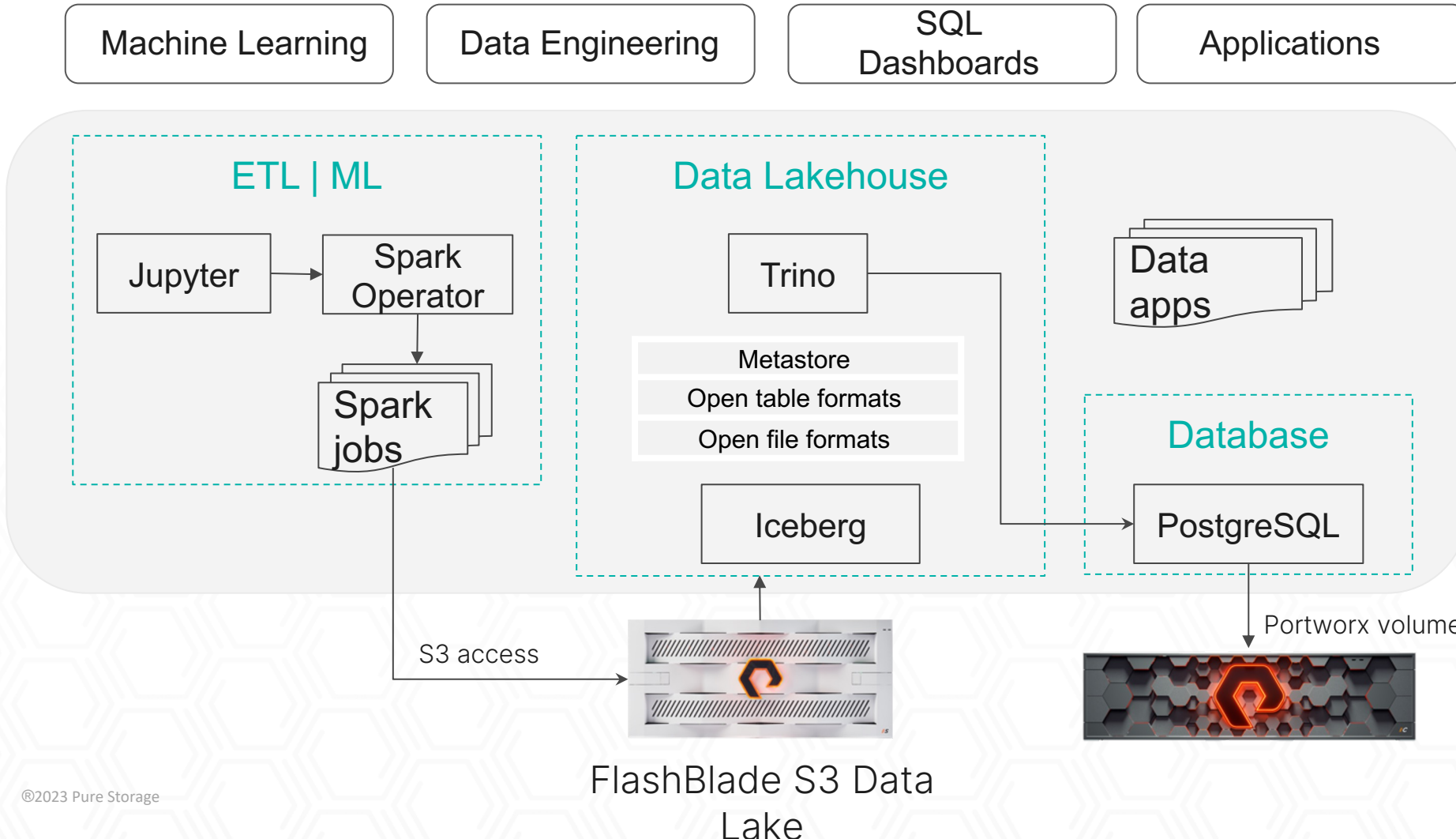- In a single system



BI   Reports   Data Science   Machine Learning

Metadata, Caching, and Indexing Layer

ETL

Data Lake

Structured, Semi-structured & Unstructured Data

# The Data Lakehouse Architecture

**Unites the strength of data lakes and data warehouse**

- Supports ETL, SQL, and ML workloads
- In a single system

# Recap - Three Key Components in a Data Lakehouse

Build an open data lakehouse with high performance object storage.

| 01 | 02 | 03 |
|---|---|---|
| **Data lake** | **Metadata layer** | **Processing engine** |
| Leverage current data lake and open data format | Add metadata layer for advanced data management | Use processing engine that understands the lakehouse spec |
| **High performance object storage** | **Delta Lake, Apache Iceberg** | **Apache Spark, Trino, Dremio, Vertica, Snowflake** |

# Open Data Lakehouse in Action

An example architecture for simple, fast and open data lakehouse with **Pure Stora**

Machine Learning

Data Engineering

SQL Dashboards

Applications

## ETL | ML

Jupyter → Spark Operator → Spark jobs

## Data Lakehouse

Trino

Metastore
Open table formats
Open file formats

Data apps

Iceberg

### Database

PostgreSQL

- No lock-in
- Inexpensive
- Fast data exploration
- Simple
- Cloud ready

S3 access

FlashBlade S3 Data Lake

Portworx volume

# AI

# ...AI is more than just the Model

*"Hidden Technical Debt in Machine Learning Systems", Google NIPS 2015*
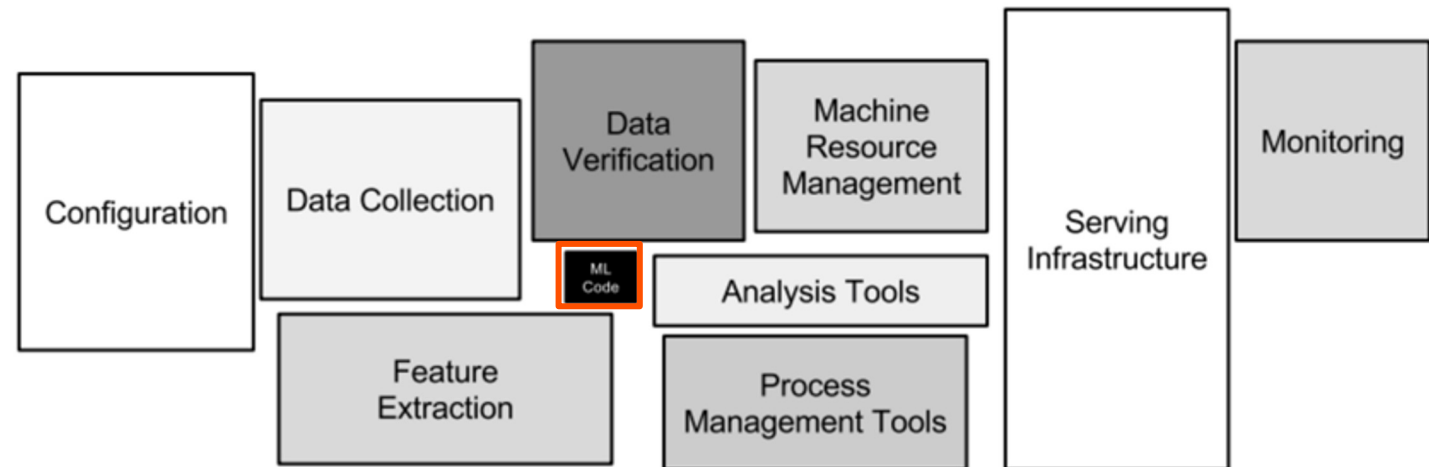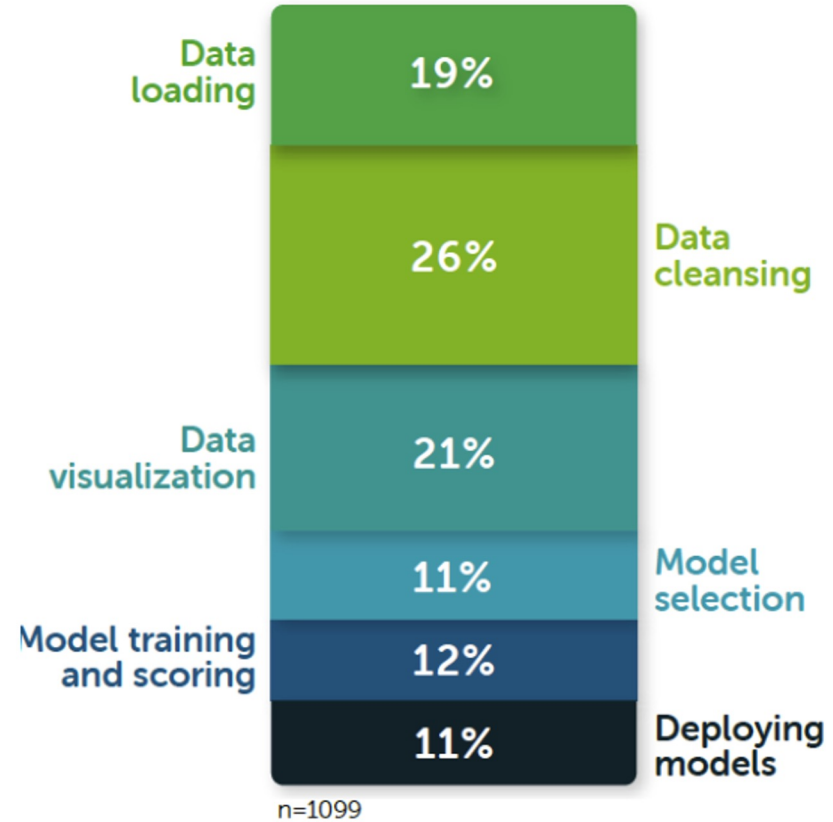


Figure 1: Only a small fraction of real-world ML systems is composed of the ML code, as shown by the small black box in the middle. The required surrounding infrastructure is vast and complex.

# So why should AI care about Pure Storage

## Time to science



| | |
|---|---|
| Data loading | 19% |
| | 26% Data cleansing |
| Data visualization | 21% |
| | 11% Model selection |
| Model training and scoring | 12% |
| | 11% Deploying models |

n=1099

How data scientists spend their time (Image courtesy Anaconda "2020 State of Data Science: Moving From Hype Toward Maturity.")