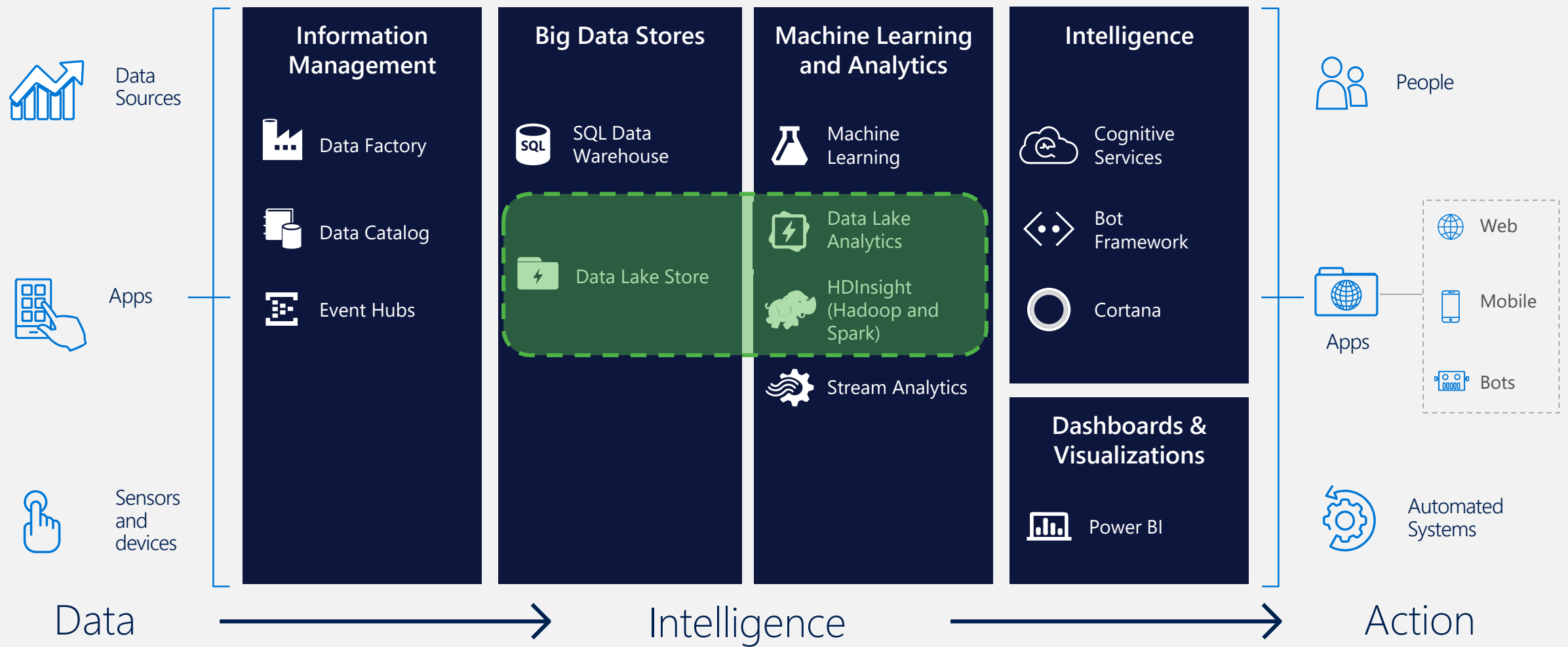# Azure Data Lake Store

*Manjunath Suryanarayana*
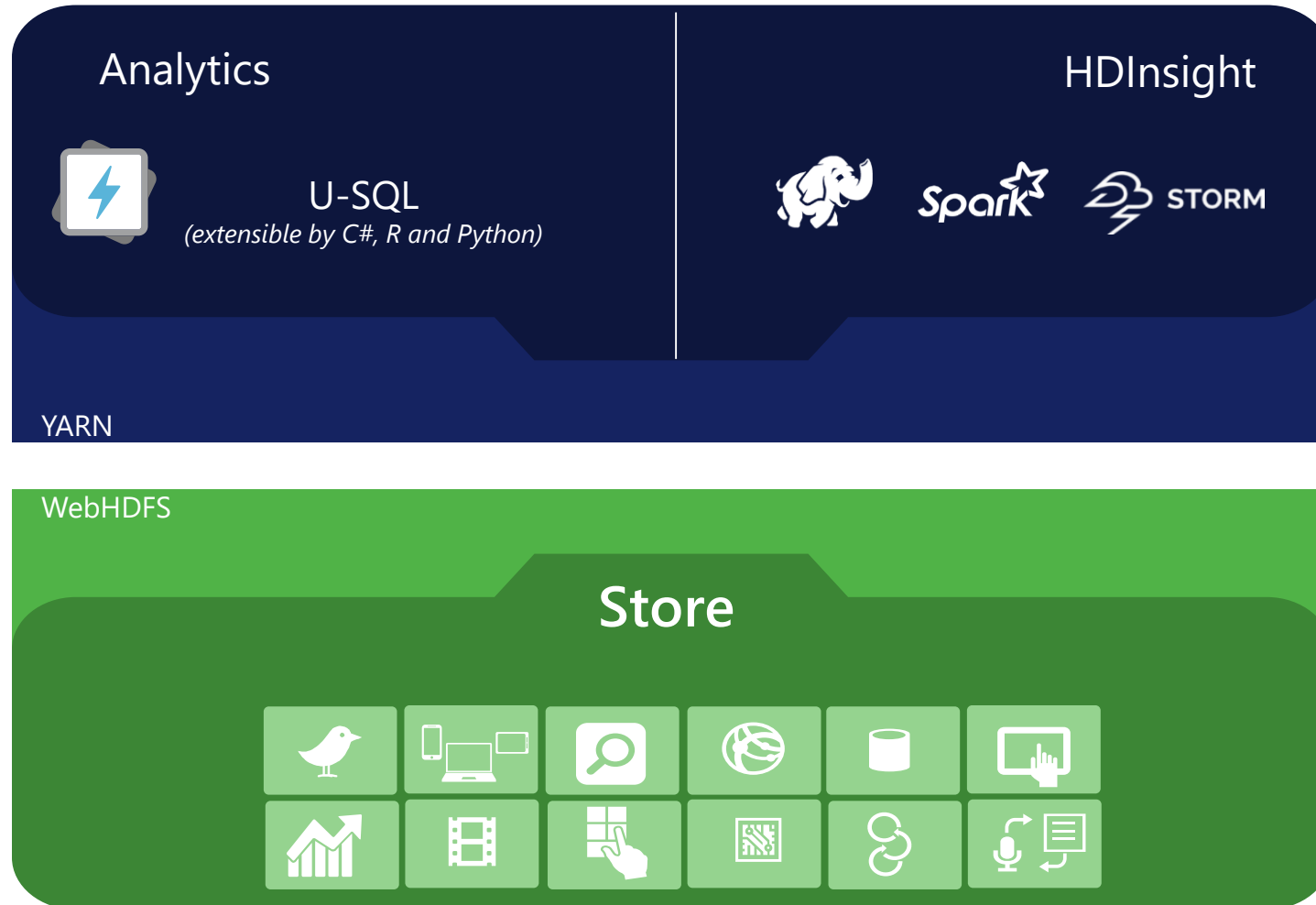
Microsoft

# Azure Data Lake
## as part of Cortana Intelligence Suite

**Data Sources**

**Apps**

**Sensors and devices**

### Information Management
- Data Factory
- Data Catalog
- Event Hubs

### Big Data Stores
- SQL Data Warehouse
- Data Lake Store

### Machine Learning and Analytics
- Machine Learning
- Data Lake Analytics
- HDInsight (Hadoop and Spark)
- Stream Analytics

### Intelligence
- Cognitive Services
- Bot Framework
- Cortana

### Dashboards & Visualizations
- Power BI

**People**

**Apps**
- Web
- Mobile
- Bots

**Automated Systems**

Data → Intelligence → Action

Microsoft

# Azure Data Lake



**Analytics**

U-SQL
*(extensible by C#, R and Python)*

**HDInsight**

Spark   STORM
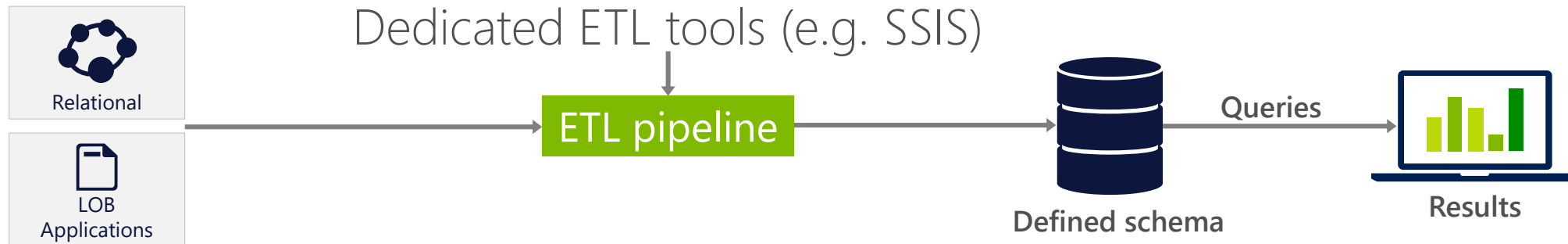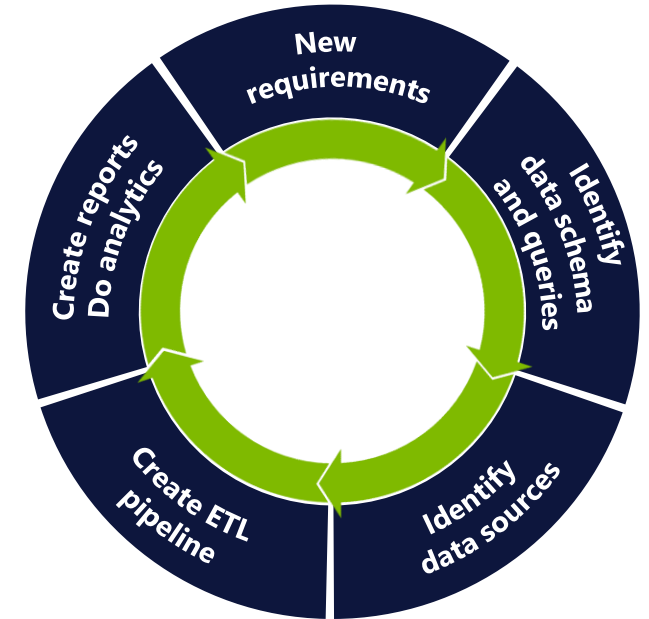
YARN

WebHDFS

**Store**

Microsoft

# Why data lakes?

# Traditional business analytics process

1. Start with end-user requirements to identify desired reports and analysis

2. Define corresponding database schema and queries

3. Identify the required data sources

4. Create a Extract-Transform-Load (ETL) pipeline to extract required data (curation) and transform it to target schema ('*schema-on-write*')

5. Create reports. Analyze data



Dedicated ETL tools (e.g. SSIS)

Relational

LOB Applications

ETL pipeline

Defined schema

Queries

Results

**All data not immediately required is discarded or archived**

Microsoft

# New big data thinking: All data has value

- All data has potential value
- Data hoarding
- No defined schema—stored in native format
- Schema is imposed and transformations are done at query time *(schema-on-read)*.
- Apps and users interpret the data as they see fit

| Gather data from all sources | Store indefinitely | Analyze | See results |

Iterate

Microsoft

# Data Lake Store: Technical Requirements

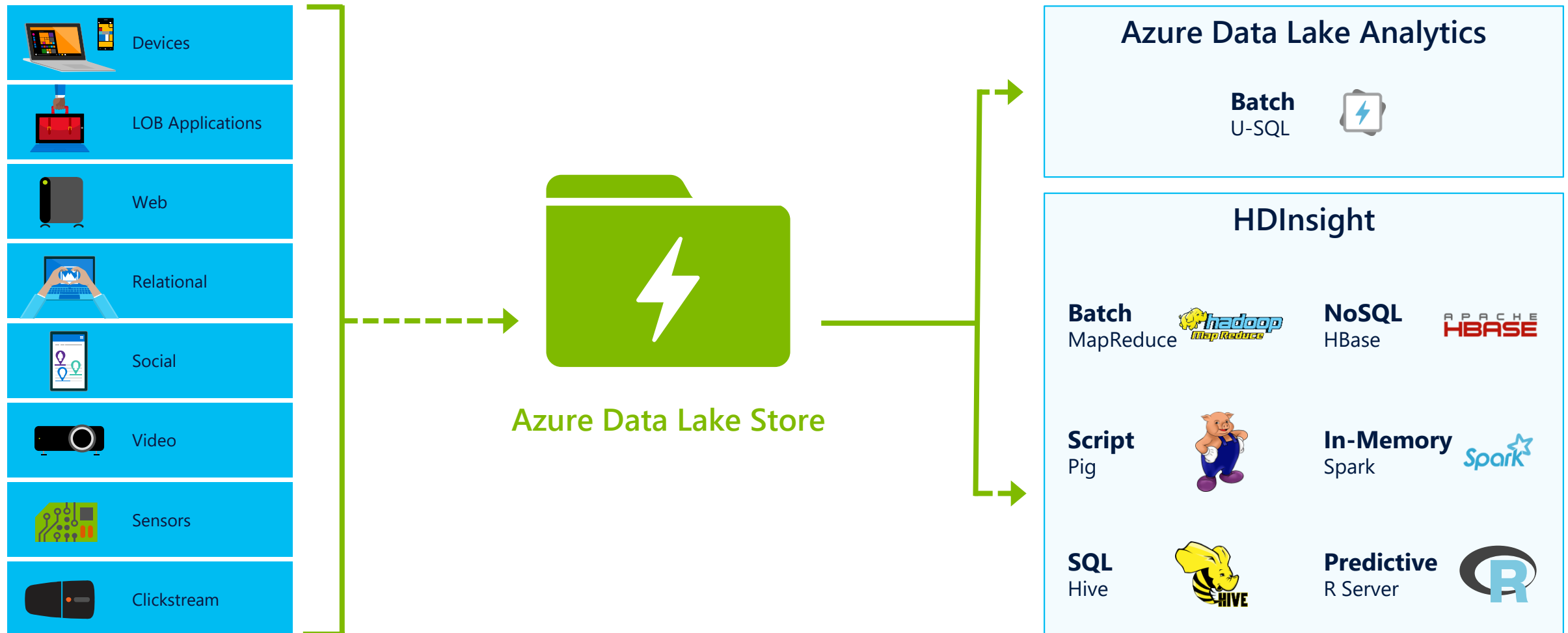| | | |
|---|---|---|
| | **Secure** | Must be highly secure to prevent unauthorized access (especially as all data is in one place). |
| | **Scalable** | Must be highly scalable. When storing all data indefinitely, data volumes can quickly add up |
| | **Reliable** | Must be highly available and reliable (no permanent loss of data). |
| | **Throughput** | Must have high throughput for massively parallel processing via frameworks such as Hadoop and Spark |
| | **Details** | Must be able to store data with all details; aggregation may lead to loss of details. |
| | **Native format** | Must permit data to be stored in its 'native format' to track lineage & for data provenance. |
| | **All sources** | Must be able ingest data from a variety of sources-LOB/ERP, Logs, Devices, Social NWs etc. |
| | **Multiple analytic frameworks** | Must support multiple analytic frameworks—Batch, Real-time, Streaming, ML etc.<br>No one analytic framework can work for all data and all types of analysis. |

Microsoft

# Azure Data Lake Store
## Overview

# Big Data analytics workloads

A highly scalable, distributed, parallel file system in the cloud
specifically designed to work with a variety of big data analytics workloads

| Devices |
| LOB Applications |
| Web |
| Relational |
| Social |
| Video |
| Sensors |
| Clickstream |

**Azure Data Lake Store**

## Azure Data Lake Analytics

**Batch**
U-SQL

## HDInsight

**Batch**
MapReduce

**NoSQL**
HBase

**Script**
Pig

**In-Memory**
Spark

**SQL**
Hive

**Predictive**
R Server

Scale, performance, reliability

Microsoft

# Azure Data Lake Store: no scale limits

Azure Data Lake Store integrates with
Azure Active Directory (AAD) for:

- ⚡ Amount of data stored
- ⚡ How long data can be stored
- ⚡ Number of files
- ⚡ Size of the individual files
- ⚡ Ingestion throughput

**Seamlessly scales
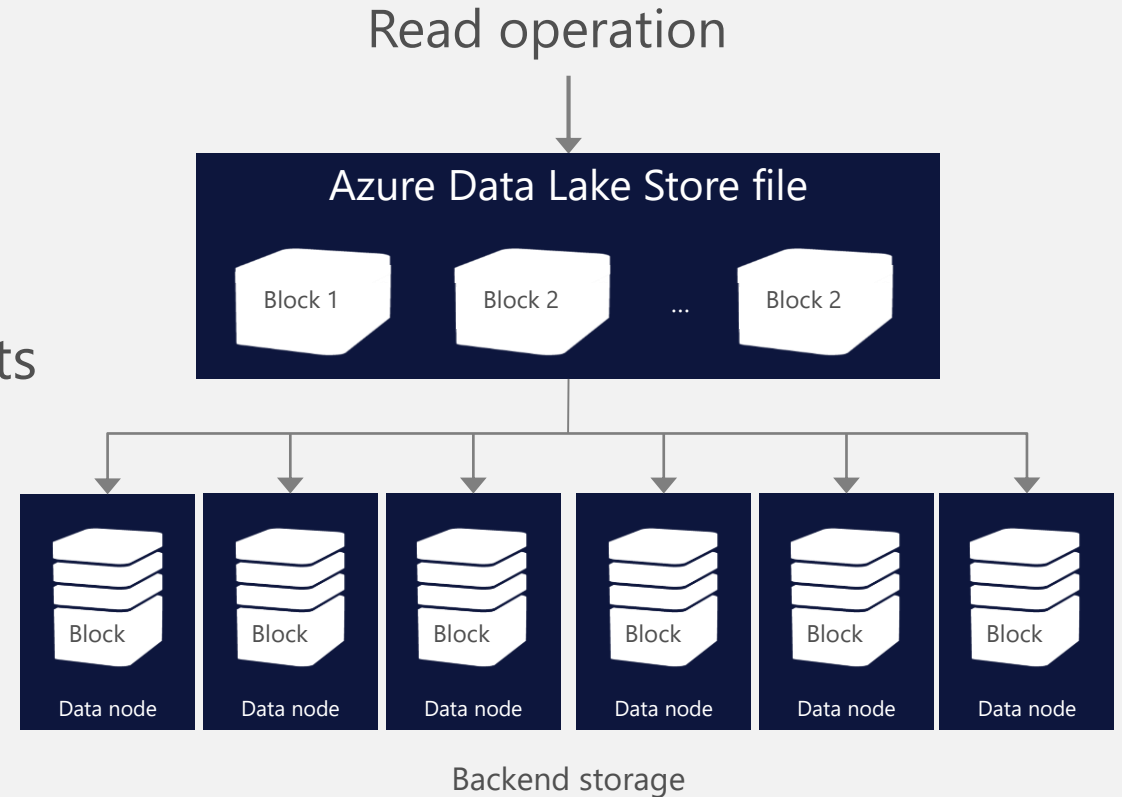from a few KBs
to several PBs**

Microsoft

# ADL Store Unlimited Scale – How it works

⚡ Each file in ADL Store is sliced into blocks

⚡ Blocks are distributed across multiple data nodes in the backend storage system

⚡ With sufficient number of backend storage data nodes, files of any size can be stored

⚡ Backend storage runs in the Azure cloud which has virtually unlimited resources

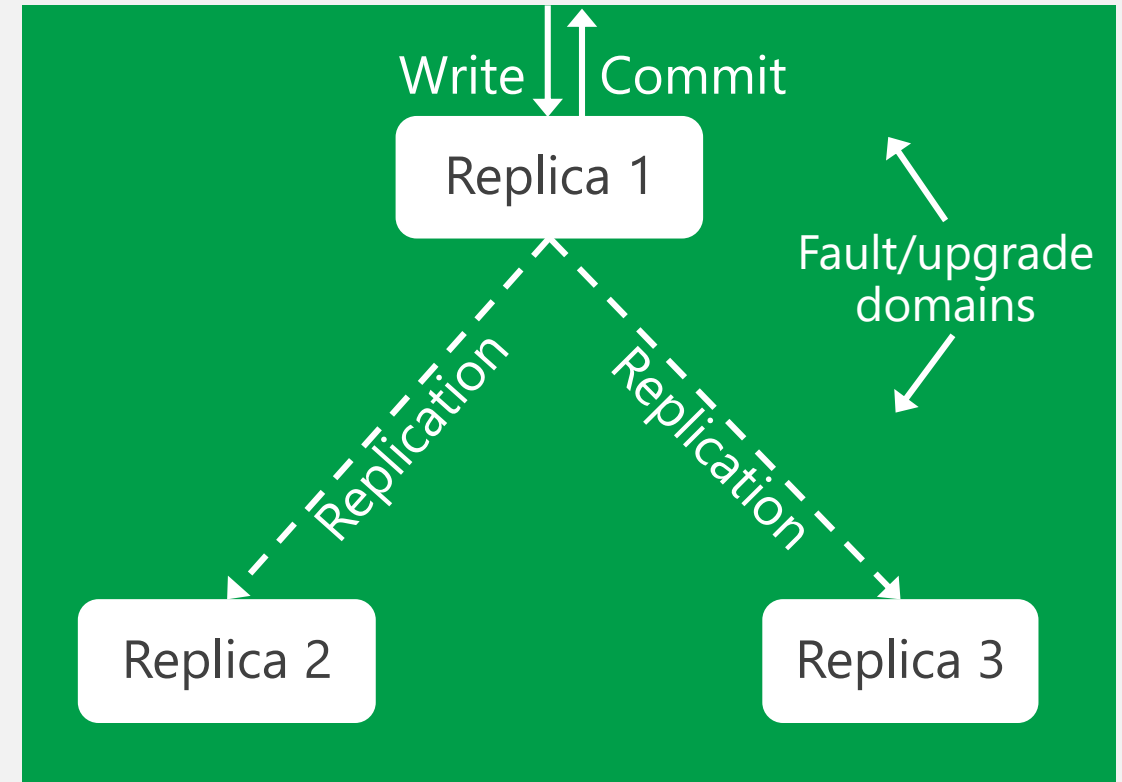⚡ Metadata is stored about each file
No limit to metadata either.

Azure Data Lake Store file

| Block 1 | Block 2 | ... | Block 2 |

| Block | Block | Block | Block | Block | Block |
| Data node | Data node | Data node | Data node | Data node | Data node |

Backend Storage

Microsoft

# ADL Store offers massive throughput

⚡ Through read parallelism ADL Store provides massive throughput

⚡ Each read operation on a ADL Store file results in multiple read operations executed in parallel against the backend storage data nodes

Read operation

**Azure Data Lake Store file**

Block 1    Block 2    ...    Block 2

Block    Block    Block    Block    Block    Block

Data node    Data node    Data node    Data node    Data node    Data node

Backend storage

Microsoft

# ADL Store: high availability and reliability

- Azure maintains 3 replicas of each data object per region across three fault and upgrade domains

- Each create or append operation on a replica is replicated to other two

- Writes are committed to application only after all replicas are successfully updated

- Read operations can go against any replica

Write Commit

Replica 1

Fault/upgrade domains

Replication Replication

Replica 2 Replica 3

Data is never lost or unavailable even under failures

Microsoft

14

# The building blocks

Ingestion, processing, egress, visualization, and orchestration tools

Microsoft

# Big Data Flow

**Business apps**

**Custom apps**

**Sensors and devices**

**Ingestion**
Bulk, Event Ingestion

**Processing**
Preparation, Analytics

**Azure Data Lake Store**

**Visualization**

**Discovery**

People

# Ingestion tools – Getting started

## Data on your desktop

### Azure Portal

- Easy to use
- Good for small amount of data
- Analyzing data using Portal

### PowerShell

- Upload file and folders
- Control parallelism
- Control format of upload
- Need to use other services

### ADL Tools for Visual Studio

- Integrated experience
- Drag-and-drop
- Programmatic Analytics

### CLI

- Linux, Mac
- Most features of PowerShell

## Data located in other stores

### Azure Data Factory

- Copy Wizard for intuitive one-time copy from multiple sources

### AdlCopy

- Copy data easily from Azure Storage at least cost

### OSS tools on HDI

- Distcp, Sqoop
- If analyzing data using HDInsight

# Azure Data Factory

Compose, orchestrate & monitor data services at scale

- ⚡ Fully managed service
- ⚡ Any data on-premises or in the cloud
- ⚡ Single pane of glass management
- ⚡ Global service infrastructure
- ⚡ Cost Effective

Hadoop on Azure

Data Lake Analytics

Custom Code

Stored Procedures

Machine Learning

VM

No SQL

ADL

1110
1010
1010

1110
1110
1010
1010

Trusted data

BI & analytics

# Azure Data Factory

Connects ADL Store out-of-the-box to all your stores

| Category | Data store | Supported as source | Supported as sink |
|---|---|:---:|:---:|
| Azure | **Azure Data Lake Store** | ● | ● |
| | Azure Blob storage | ● | ● |
| | Azure SQL Database | ● | ● |
| | Azure SQL Data Warehouse | ● | ● |
| | Azure Table storage | ● | ● |
| | Azure DocumentDB | ● | ● |
| Databases | SQL Server* | ● | ● |
| | Oracle* | ● | ● |
| | MySQL* | ● | |
| | DB2* | ● | |
| | Teradata* | ● | |
| File | HDFS* | ● | |
| | Others | ● | |

# Visualizing data



Azure Data Factory

SQL Data Warehouse

Azure Data Lake Store

HDInsight

Power BI dashboard

Power BI desktop

Excel*

Jupyter Data Science Notebooks

* Post General Availability

# Customizing using SDKs/APIs

**Your application**

| ADL PowerShell | ADL XPlat CLI | | |
|---|---|---|---|
| ADL .NET SDK | ADL Java SDK | ADL Node.js SDK | ADL Python* |

Only on Windows OS

**Azure and ADL Store REST APIs**

* At General Availability

# Building pipelines - Management and orchestration

## Out-of-the-box tools

**Azure Data Factory**

First-class support

**Azure Stream Analytics**

Seamlessly stream data

**OSS tools**

Supports OSS tools

**PowerShell**

Management with Workflow & Script Runbooks

## Custom tools

**ADL Store SDK**

Available in multiple languages

**REST APIs**

For unsupported languages and platforms

Security

# Security features

| | |
|---|---|
| **Identity Management & Authentication** | Azure Active Directory |
| **Access Control & Authorization** | Azure RBAC for Account Management<br><br>File & Folder level POSIX ACLs |
| **Auditing** | Azure Diagnostic Audit Logs |
| **Data Protection & Encryption** | Encryption on the wire using HTTPS<br><br>Transparent Service side encryption using service & customer managed keys |

# ADL Store Security: AAD integration

⚡ Multi-factor authentication based on OAuth2.0

⚡ Integration with on-premises AD for federated authentication

⚡ Role-based access control

⚡ Privileged account management

⚡ Application usage monitoring and rich auditing

⚡ Security monitoring and alerting

⚡ Fine-grained ACLs for AD identities

Microsoft

# Leveraging Azure Active Directory



OAuth token → Data Lake Store → Graph APIs

OAuth token → Data Lake Store → Graph APIs

Basic auth (HTTPS) → HDInsight → Kerberos, LDAPS

Azure active directory tenant

Azure directory domain services instance

1. Create ADDS instance in separate VNET

2. Add users to AAD Tenant

3. Add users to ADLA RBAC roles

4. Add users to ADLS RBAC roles & file system ACLs

5. Join HDInsight cluster to ADDS instance

# ADL Store security: Role-based access

⚡ Each file and directory is associated with an owner and a group

⚡ Files or directories have separate permissions (read(r), write(w), execute(x)) for owners, members of the group, and for all other users

⚡ Fine-grained access control lists (ACLs) rules can be specified for specific named users or named groups

# Granular control of file and folder access
## POSIX-Style ACLs with full compatibility with HDFS/WebHDFS

- Generate default ACLs for files and folders
- Customize for fine-tuned control
- Access ACLs control how a user can access to the file or folder
- Default ACLs used to construct the Access ACL of new children
- Default ACLs copied to the Default ACL of new child folders

**Child File**

**Folder**

Access ACLS

Access ACLS

Default ACLS

**New Child Folder**

**New Child File**

Default ACLS

Access ACLS

Access ACLS

# IP address ACLs

- Access rights based on IP range
- Applies to traffic from inside or outside Azure
- Cannot be used to filter VNETs

**Azure Data Lake Store**

15.23.1.5

64.34.55.130

**64.34.55.130 – 64.34.55.135**

**IP range whitelist**

# Encryption of data at rest*

⚡ Provides transparent server-side encryption

⚡ Choice made at account creation to enable encryption

⚡ Service managed keys or user managed keys

**Azure Data Lake Store**

**Azure Key Vault**

* In Private Preview

# Audit logs for data access

⚡ Logs are available in JSON format

⚡ Sample U-SQL scripts are available on GitHub to-read logs

⚡ Enhancement to logs will continue through GA



**Azure Data Lake Analytics**

**Azure Data Lake Store**

**Azure Blob Store**

[T1] Alice, Write
[T2] Bob, Read

Bob

Alice

# ADL Store
## Hadoop integration

# ADL Store is HDFS-compatible

With a WebHDFS endpoint Azure Data Lake Store is a Hadoop-compatible file system that integrates seamlessly with Azure HDInsight

| Map reduce | Hive query | HBase transactions | Spark queries |
|---|---|---|---|

**Any HDFS application**

**Azure HDInsight**

**Hadoop WebHDFS client**

**WebHDFS-compatible REST API**

**Azure Data Lake Store**

Microsoft

# ADL Store: ingress and egress

# ADL Store: Ingress



Data can be ingested into Azure Data Lake Store from a variety of sources

Azure SQL DB

Azure SQL DW

Azure tables
Table Storage

On-premises databases

Azure Data Factory

ADL Store

Azure Stream Analytics

Azure Event Hubs

ADL built-in copy service
Azure Data Factory
Hadoop DistCp

Azure Storage Blobs

.NET SDK
CLI
Azure Portal
Azure PowerShell

Custom programs

Microsoft

# ADL Store: Egress

Data can be exported from Azure Data Lake Store into numerous targets/sinks

Azure SQL DB

Azure SQL DW

Azure Tables

Table Storage

On-premises databases

Azure Data Factory

Hadoop DistCp

Azure Data Factory

Apache Sqoop

ADL Store

.NET SDK
CLI
Azure Portal
Azure PowerShell

Azure Storage Blobs

Custom programs

Microsoft

36

# ADL Store: Azure Portal integration

# Creating a new ADL Store

# ADL Store: Properties

# Viewing Users and their Roles & Privileges

# Adding Users

Owner – Lets you manage everything, including access to resources

Contributor – Lets you manage everything, except access to resources

Reader – Lets you view everything, but not make changes

User Access Administrator – Lets you manage user access to Azure resources

Microsoft

41

# File Upload

Azure Portal lets you upload files directly to ADL Store

# File Preview



- Input and output files can be previewed directly in the portal without having to download them.
- The preview shows the first few rows.
- Column numbers are automatically assigned
- Understands CSV and TSV formats.

# App Development – Languages and Tools
## Azure Data Lake Store supports multiple languages for application development

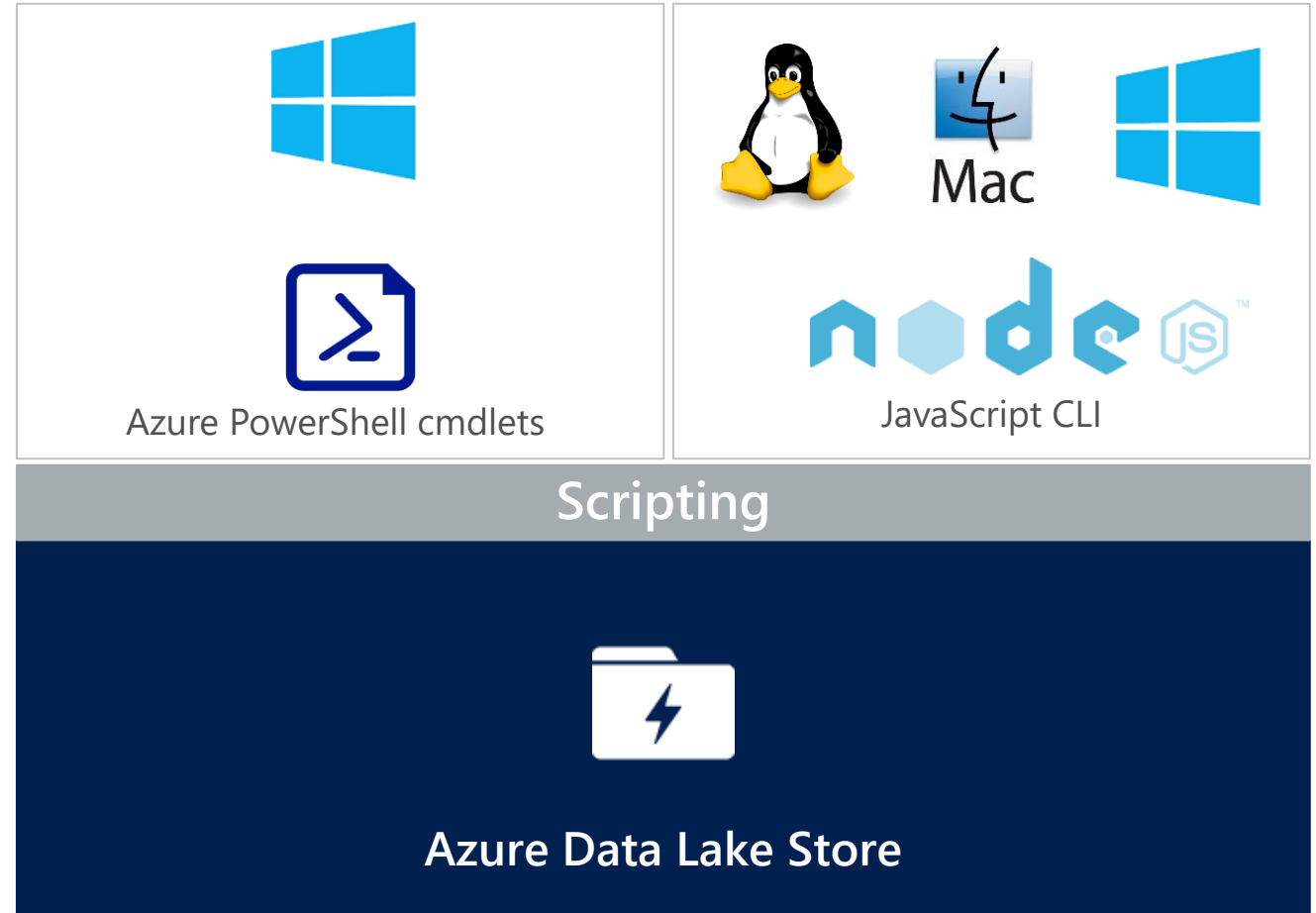| | | |
|---|---|---|
| Java Developers | → | WebHDFS |
| C++ Developers | → | LibWebHDFS |
| .NET Developers | → | Azure .NET SDK |
| Other languages | → | x-plat SDK |

Visual Studio

Note: If you are using Hadoop (Map Reduce programs or Hive or HBase) or Spark, then you will not be programming directly to the Azure Data Lake Store as they all will transparently access Azure Data Lake Store under the covers.

Microsoft

# Developing scripting applications

Provides native Windows and cross-platform (Mac, Linux) scripting experience

Scripting operations include
- ⚡ Create new directories
- ⚡ Listing the contents of a directory
- ⚡ Upload files to directory
- ⚡ Delete files/directories
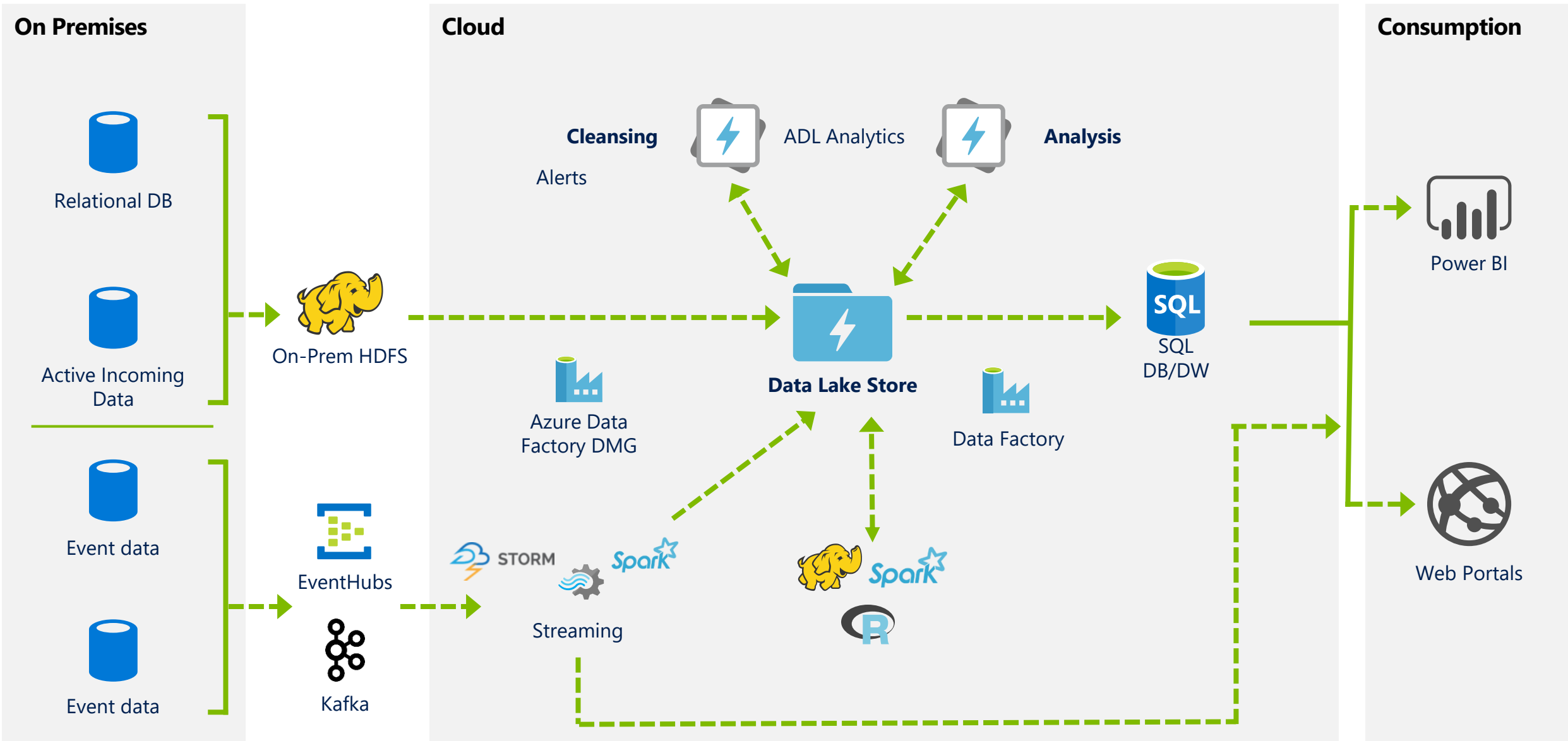- ⚡ Rename files/directories
- ⚡ ...

**Azure PowerShell cmdlets**

**JavaScript CLI**

**Scripting**

**Azure Data Lake Store**

Microsoft

Implementation
Common customer patterns

Microsoft

# Lambda architecture

# ADL Store

Costs

# Costs breakdown by stage

| | |
|---|---|
| **Ingestion** | Number of write transactions |
| **Storage** | Data stored per month |
| **Processing** | Number of read transactions<br>Number of write transactions |
| **Egress** | Number of read transactions |

Get all the advantages of ADL Store with cost concepts you are familiar with

# Get started today!

For more information visit:
http://azure.com/datalake