



Kaiser Permanente EBC

A2.0 Refresh In Motion

February 20, 2020
Bellevue, Washington

EBC Agenda and Table of Contents

Time	Topic	Speaker(s)
8:30–9:00	Breakfast	
9:00–9:15	Attendees & Introductions	
9:15–9:45	A2.0 Program Overview and EBC Expectations	Vivian/Christopher Steve/Ganesh Manish/Dana
9:45–12:00	Azure Synapse Analytics	Rohan Kumar and Shane Risk
12:00-13:00	Lunch	
13:00-13:30	Technology Architecture	Kunal Jain
13:30-14:30	Azure Data Factory	Anand Subbaraj
14:30-15:30	HDInsight Roadmap	Ashish Thapliyal and Arindam Chatterjee
15:30-16:15	Azure Data Discovery & Governance (Babylon)	Jennifer Stevens
16:15-17:00	Azure Databricks	Yatharth Gupta
17:00	Closing and Next Step Planning	

Attendees and Introductions



KAISER
PERMANENTE®

Vivian Tan	VP, Strategic Information Management & Global Partnership
Christopher Guy	Executive Director, Analytics Information Management
Manish Vipani	VP, Cloud Application Services
Hovannes Daniels	VP, Chief Data Officer
Judy Sarles	Senior Director, Chief Data Office
Animesh Ghosh	Principal Architect, Cloud Application Platforms
Rahul Pendyal	Principal Architect, Cloud Application Platforms
Sam Gambarin	Executive Director, Cloud Application Platforms
Ganesh Thondikulam	Executive Director, Analytics Digital Foundation
Subha Parthasarathy	Senior Director, Analytics Digital Foundation
Ramesh Venkataramanujam	Principal Consultant, Analytics Digital Foundation
David Schaefer	Principal Architect, Finance Data Warehouse
Alex Rosenblum	Principal IT Engineer, Cloud Application Platforms
Brian Sikora	Executive Director, Decision Support
Erin Lamb	Director, Delivery System Analytics
Amarjit Hothi	Director, Analytics Service Delivery
Shannon Madsen	Director, Analytics Digital Foundation

	Microsoft
Rohan Kumar	CVP Azure Data
Dana Smith	Account Executive
Kunal Jain	Cloud Solution Architect
Mohammad Khan	Azure Specialist
Shane Risk	Azure Synapse Senior Program Manager
Ashish Thapliyal	Azure HDInsight Principal Program Manager
Arindam Chatterjee	Azure HDInsight Principal Program Manager
Alicia Li	Azure HDInsight, Principal Program Manager
Jennifer Stevens	Azure Data Catalog Product Manager
Anand Subbaraj	Azure Data Factory Product Manager
Igor Pavlovic	Senior Software Engineer
Yatharth Gupta	Azure Databricks Principal Program Manager

Planning Leads to Results

On January 17, Kaiser Permanente's A2.0 stakeholders and stewards partnered with Microsoft to examine the program's health and improve its prognosis.

This collaboration paved the way to explore the foundation's architecture, operation model, and tenant service framework.

Microsoft is enthused to present how its best of breed Enterprise solutions on Azure (as they are available today and their release schedule to general availability) will solve Kaiser Permanente's analytic requirements.

Partnership Update

Kaiser Permanente and Microsoft are working together.



Sharing Vision

Jointly architecting and designing technology solutions in line with Business/IT requirements.



Communication

KP and Microsoft Teams are communicating daily and with a focus on outcomes.



Smart Distribution

KP is managing requirements, tenants and operations while Microsoft is lending platform and industry expertise.



Journey

Project Planning for the near and long term are in motion with all agencies involved in identifying, estimating and executing.



Skilling

KP will continue to examine its skilling model and Microsoft is offering various roles for support (PFE, DSE, PM, SSE).



Shared Investment

Microsoft is investing in continuous outreach and offering a direct pulse of its Azure Product roadmap and direct access to its Product Managers.

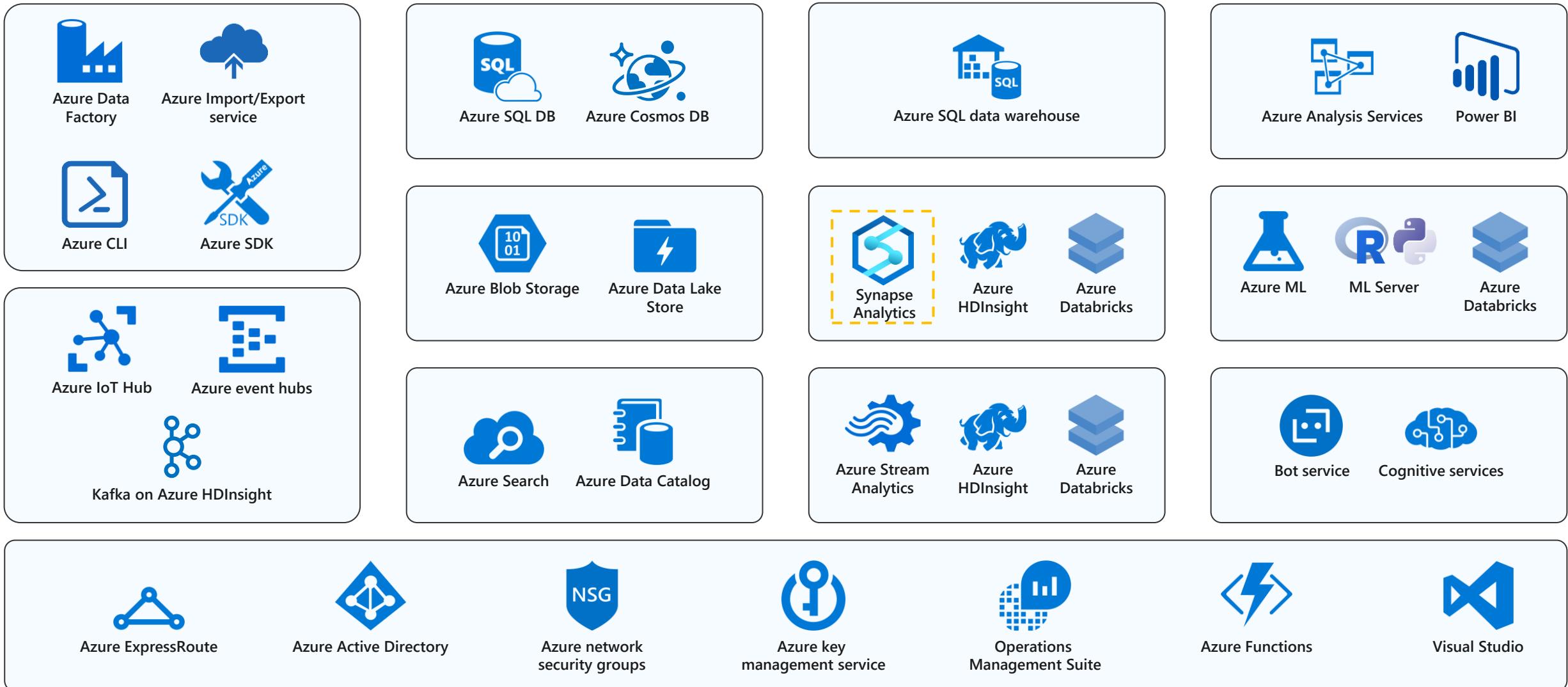


Azure Synapse Overview and Roadmap

Rohan Kumar, CVP Azure Data

Shane Risk, Azure Synapse Product Manager

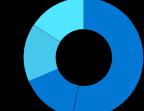
The Azure Data Landscape



Azure Synapse Analytics



Limitless Scale



Powerful Insights



Unified Experience



Unmatched Security

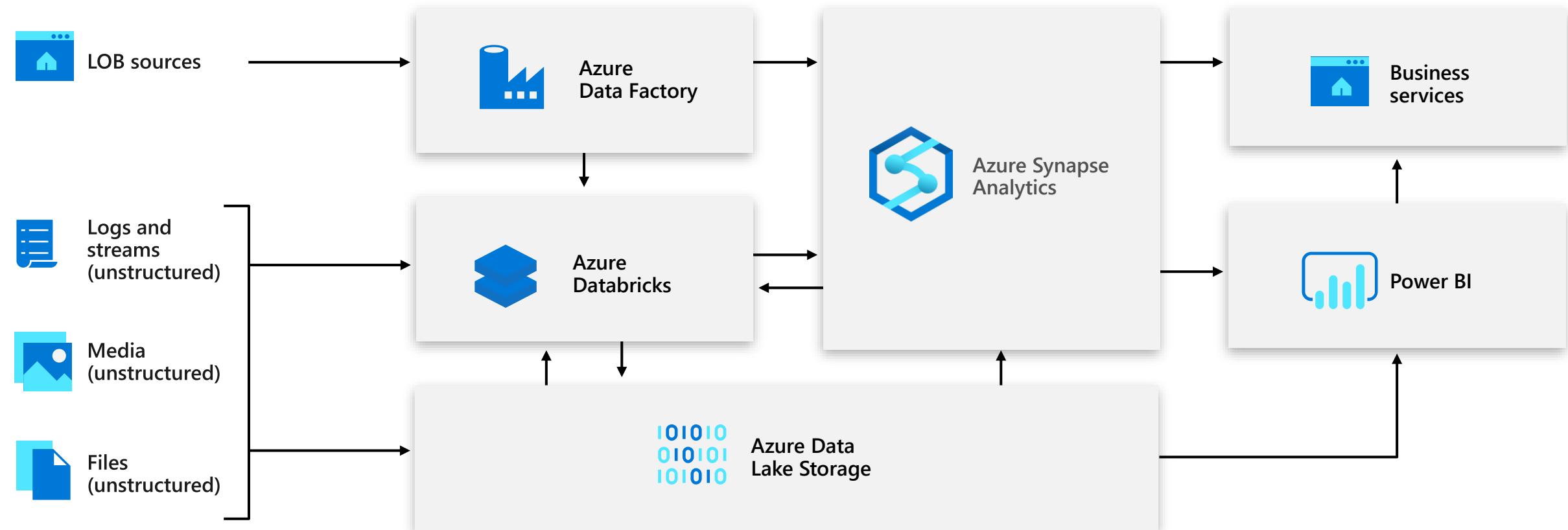
GENERALLY AVAILABLE

Provisioned Data Warehouse

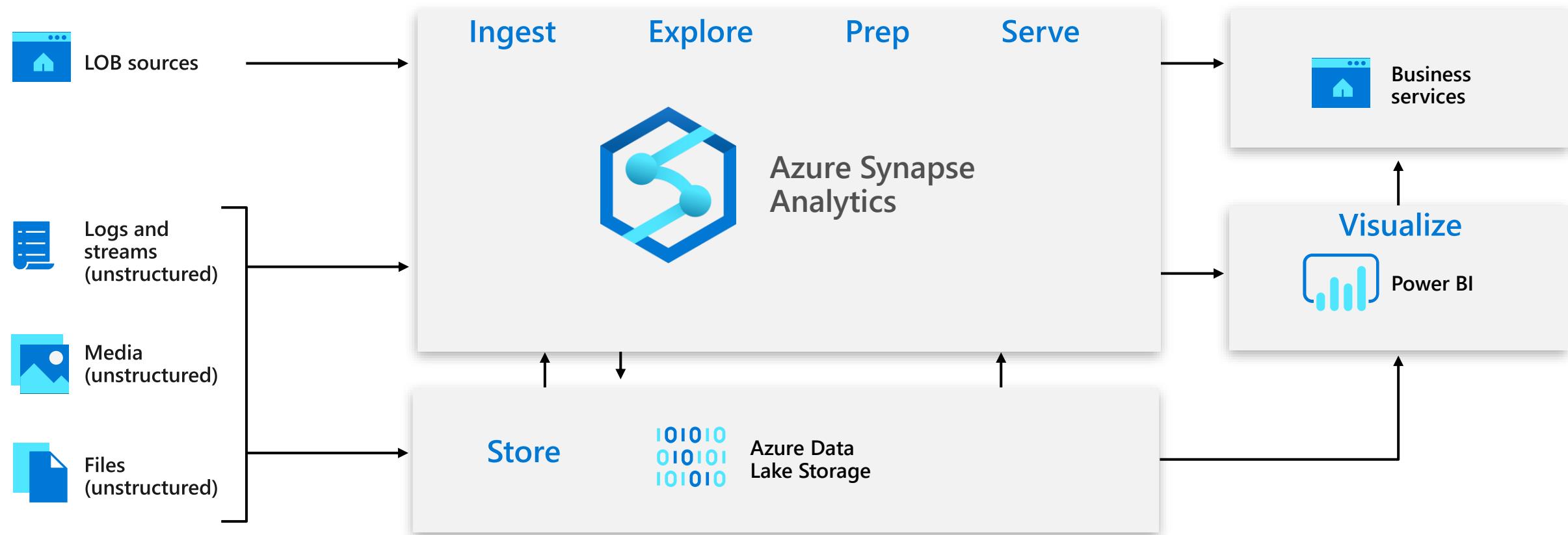
PREVIEW

On-demand Query as a Service

Typical Cloud Data Warehouse Architecture

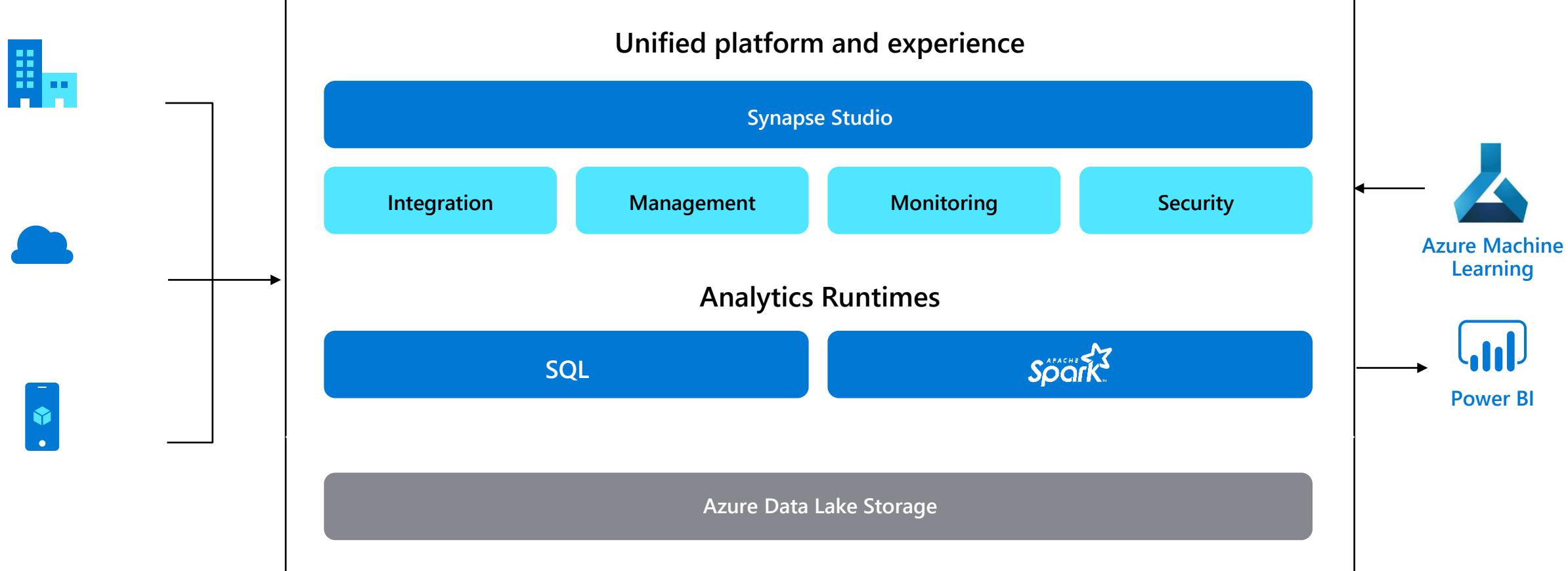


Cloud DW Architecture with Azure Synapse Analytics



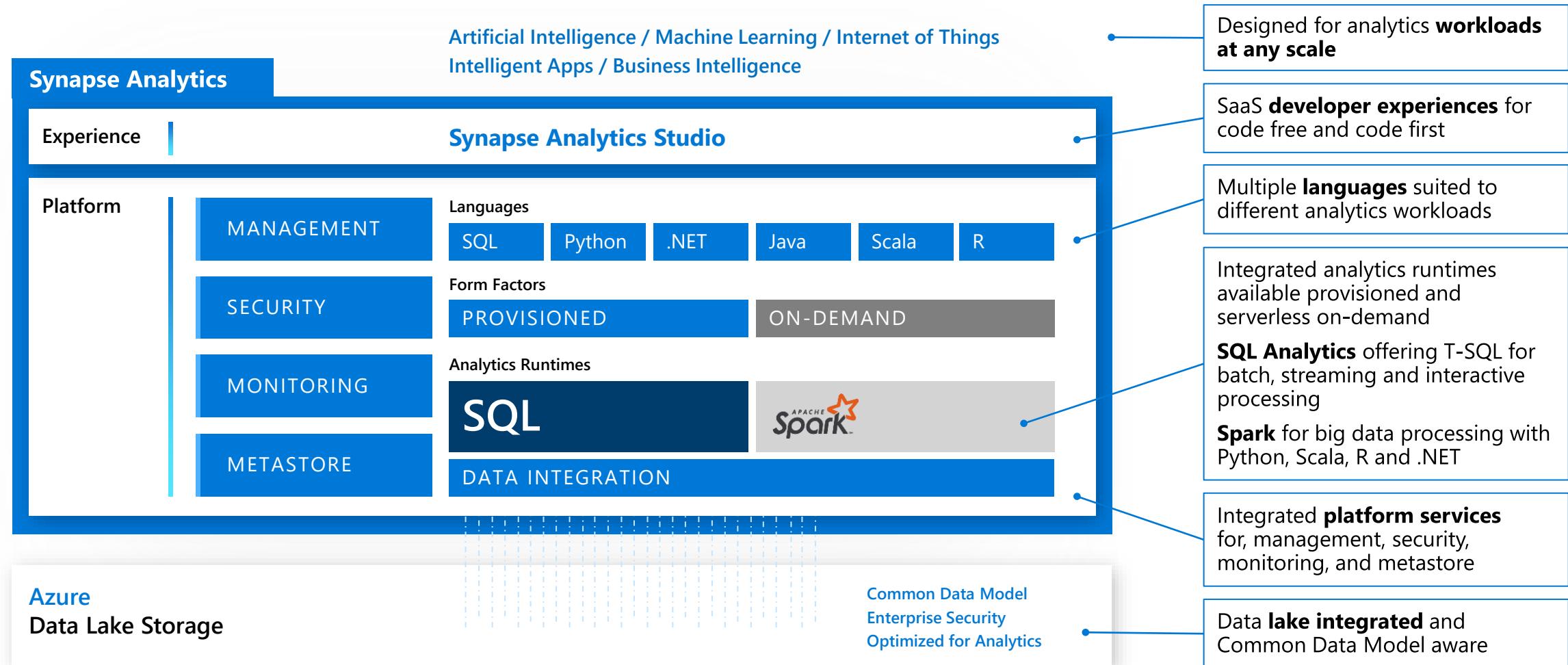
Azure Synapse Analytics

Limitless analytics service with unmatched time to insight

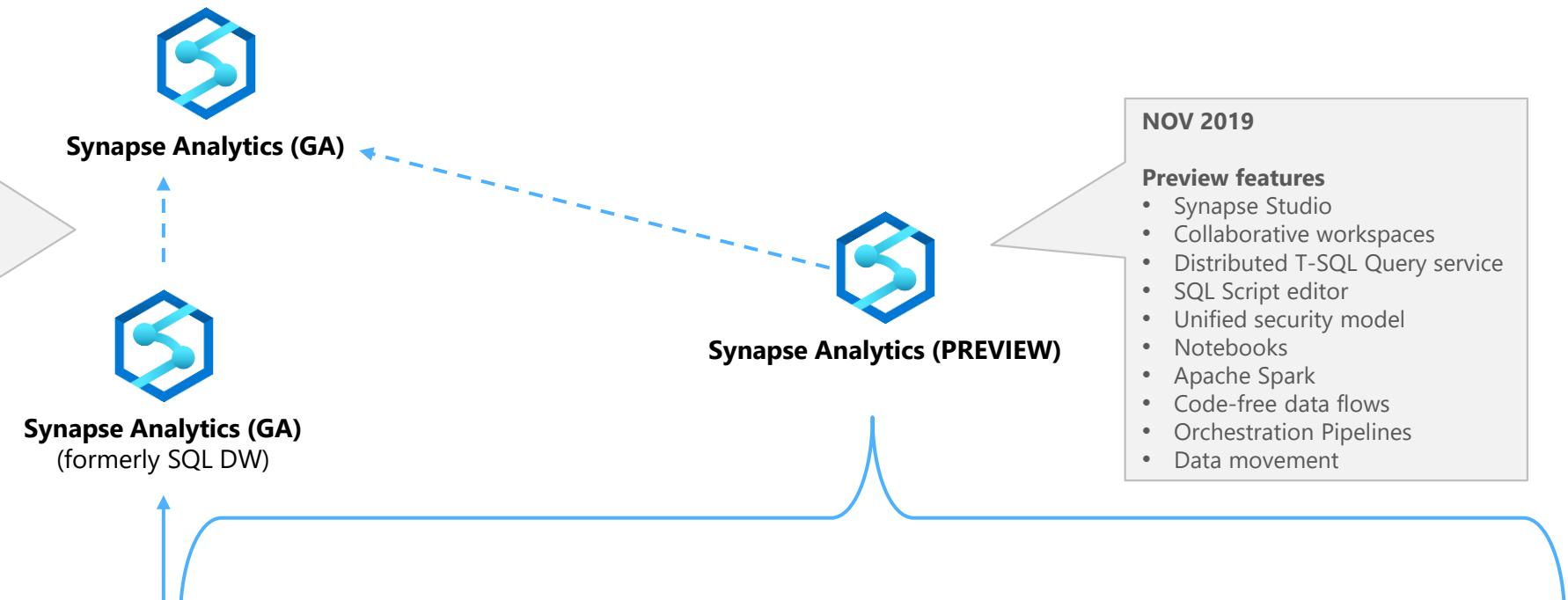
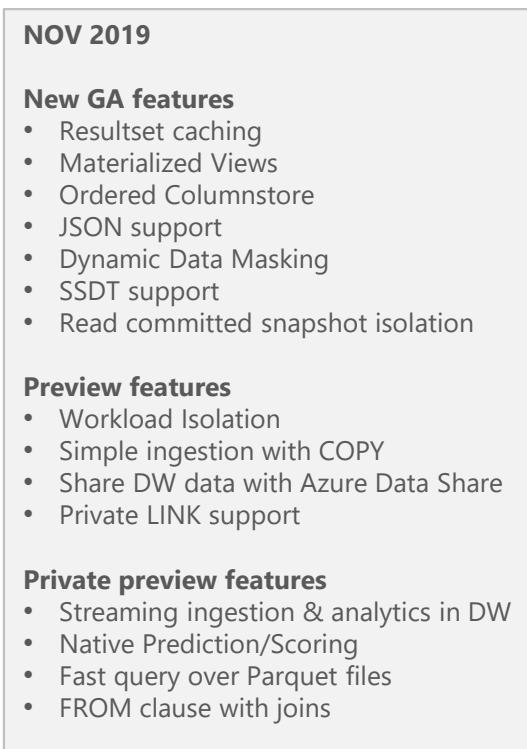


Azure Synapse Analytics

Limitless analytics service with unmatched time to insight



Azure Synapse Analytics



SQL ANALYTICS



APACHE SPARK



STUDIO



DATA INTEGRATION



General Availability (GA)
Azure Synapse Security

Comprehensive Security

Category	Feature	
Data Protection	Data in Transit	✓
	Data Encryption at Rest	✓
	Data Discovery and Classification	✓
Access Control	Object Level Security (Tables/Views)	✓
	Row Level Security	✓
	Column Level Security	✓
Authentication	Dynamic Data Masking	✓
	SQL Login	✓
	Azure Active Directory	✓
Network Security	Multi-Factor Authentication	✓
	Virtual Networks	✓
	Firewall	✓
Threat Protection	Azure ExpressRoute	✓
	Threat Detection	✓
	Auditing	✓
	Vulnerability Assessment	✓





Authentication

Azure Active Directory authentication

Overview

Manage user identities in one location.

Enable access to Azure SQL Data Warehouse and other Microsoft services with Azure Active Directory user identities and groups.

Benefits

Alternative to SQL Server authentication

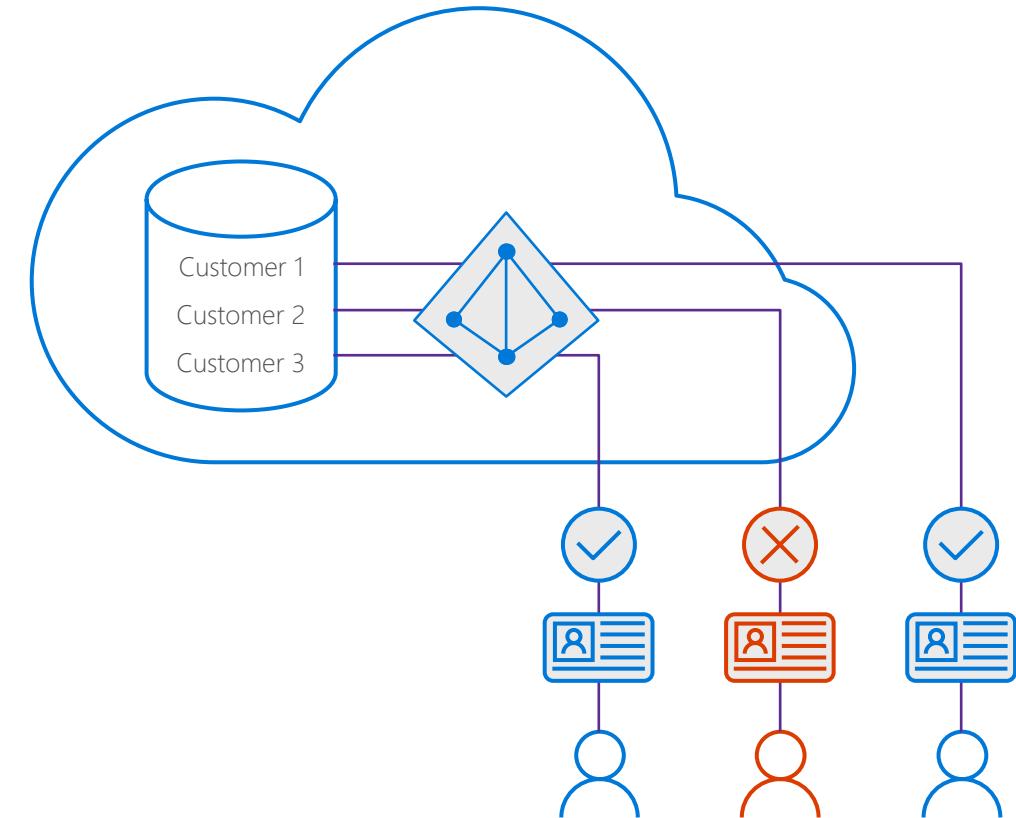
Limits proliferation of user identities across databases

Allows password rotation in a single place

Enables management of database permissions by using external Azure Active Directory groups

Eliminates the need to store passwords

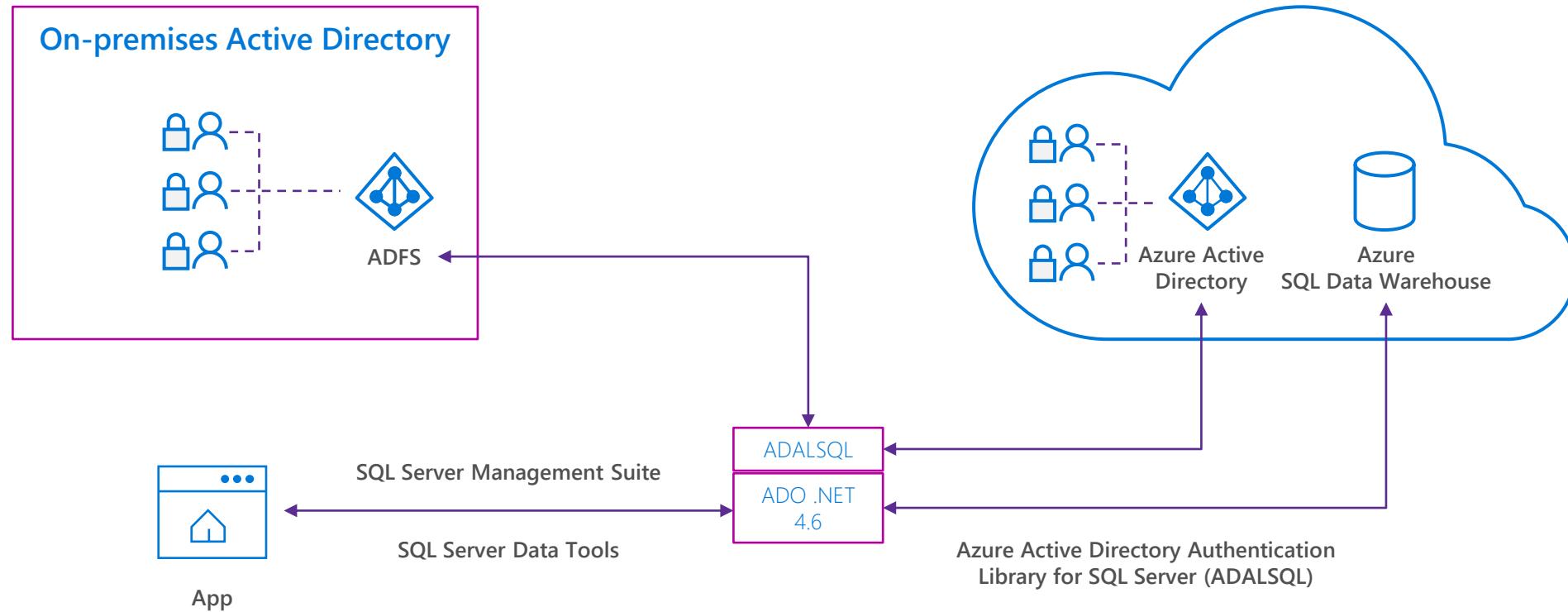
Azure SQL Data Warehouse



Azure Active Directory trust architecture



Azure Active Directory and Azure SQL Data Warehouse



Securing with firewalls

Overview

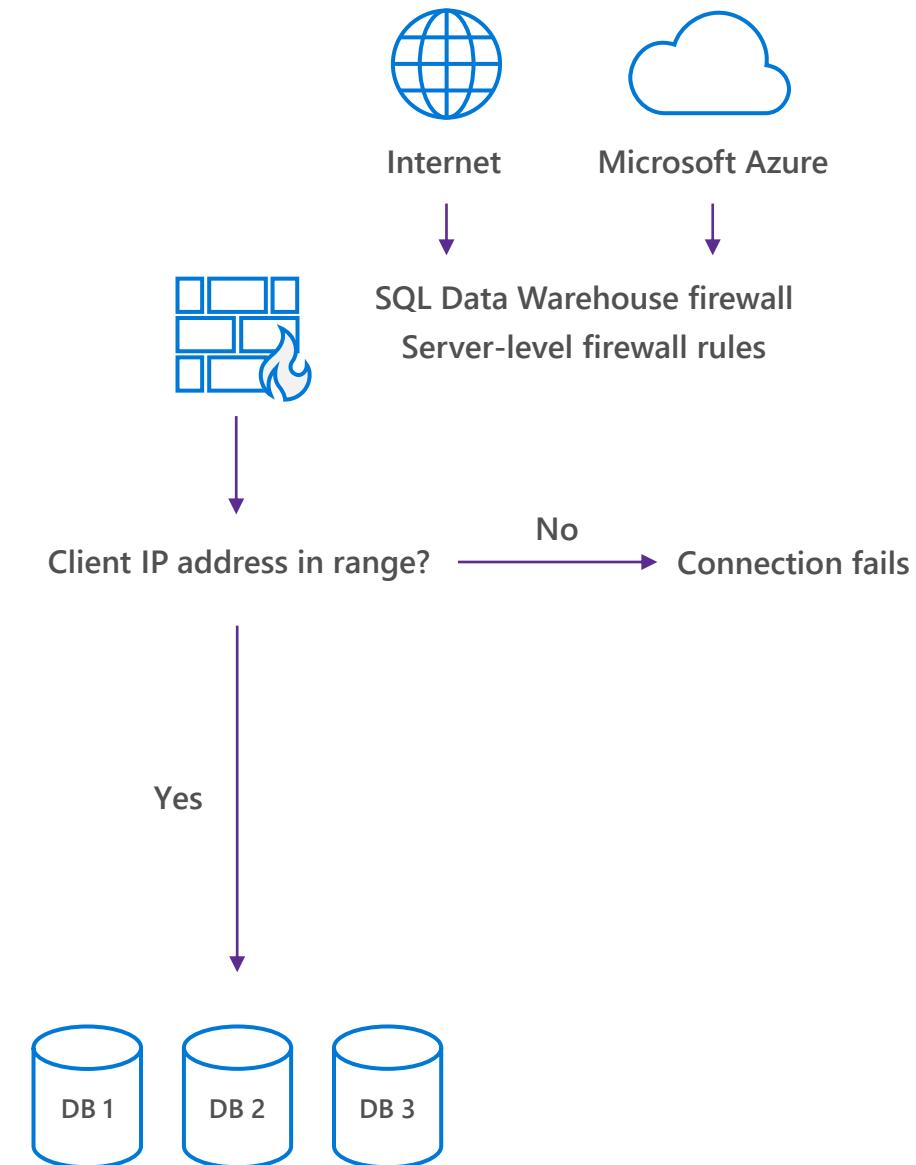
By default, all access to your Azure Synapse Analytics is blocked by the firewall.

Firewall also manages virtual network rules that are based on virtual network service endpoints.

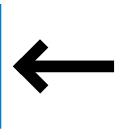
Rules

Allow specific or range of whitelisted IP addresses.

Allow Azure applications to connect.



Row-level security (RLS)



Overview

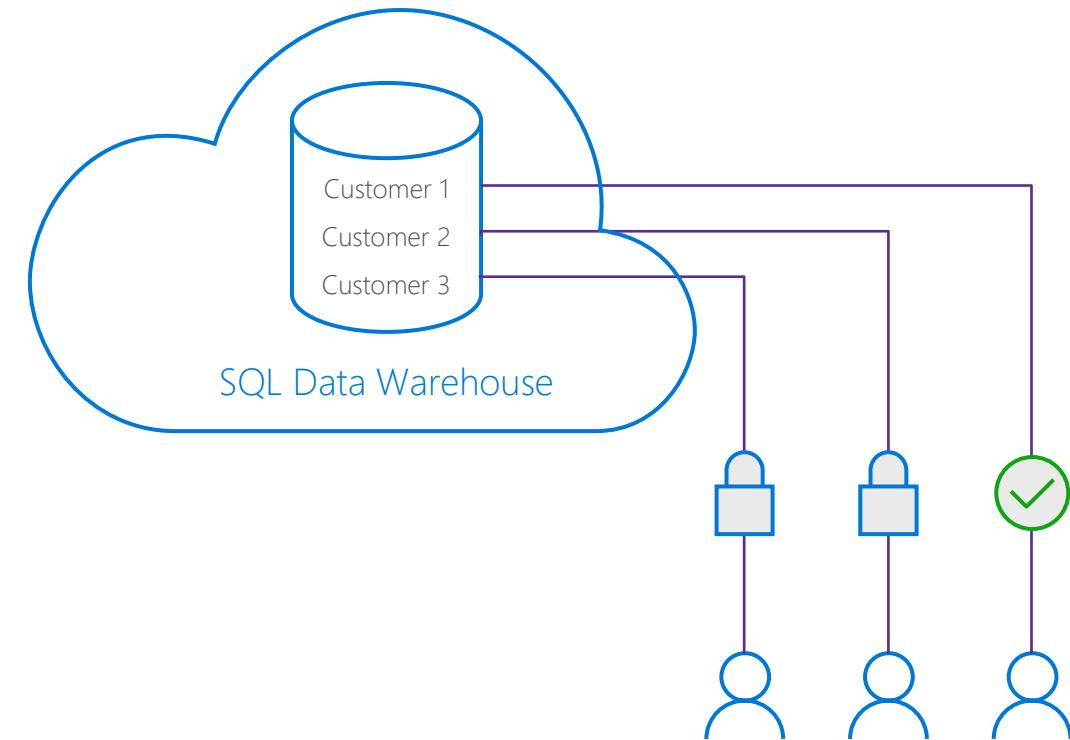
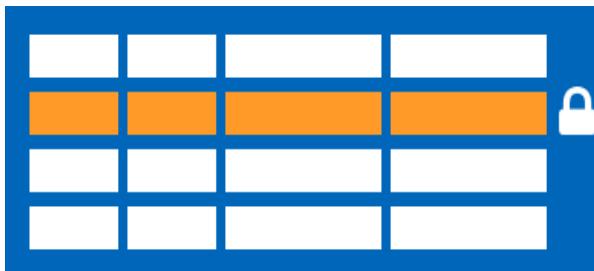
Fine grained access control of specific rows in a database table.

Help prevent unauthorized access when multiple users share the same tables.

Eliminates need to implement connection filtering in multi-tenant applications.

Administer via SQL Server Management Studio or SQL Server Data Tools.

Easily locate enforcement logic inside the database and schema bound to the table.



Row-level security



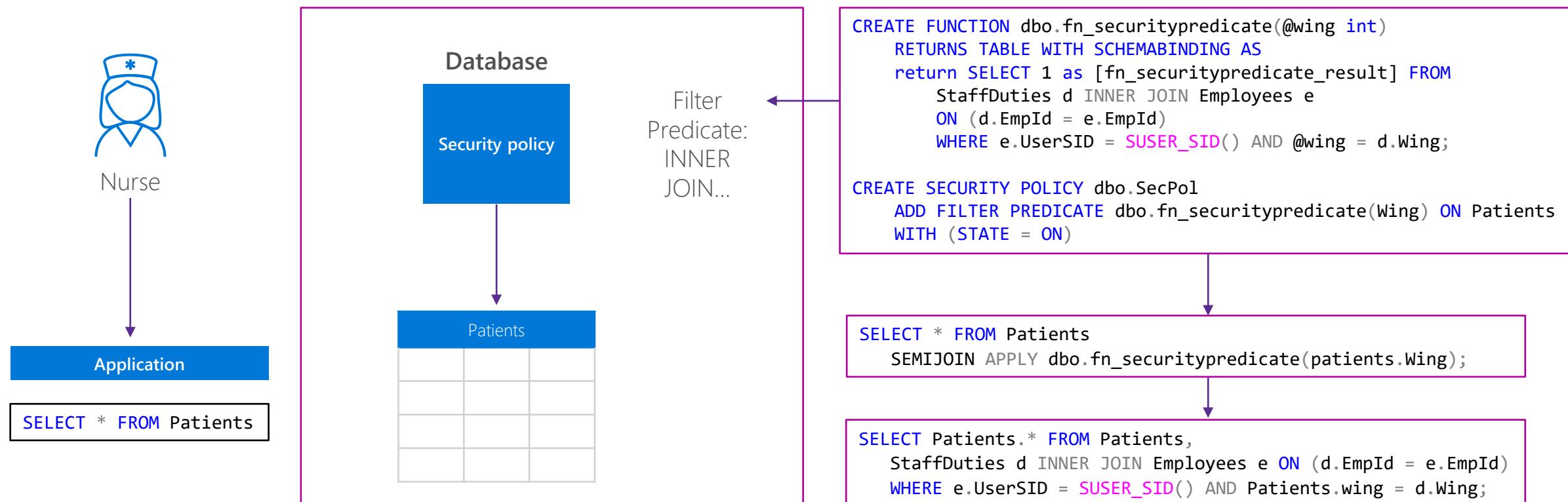
Access Control

Three steps:

1. Policy manager creates filter predicate and security policy in T-SQL, binding the predicate to the patients table.
2. App user (e.g., nurse) selects from Patients table.
3. Security policy transparently rewrites query to apply filter predicate.

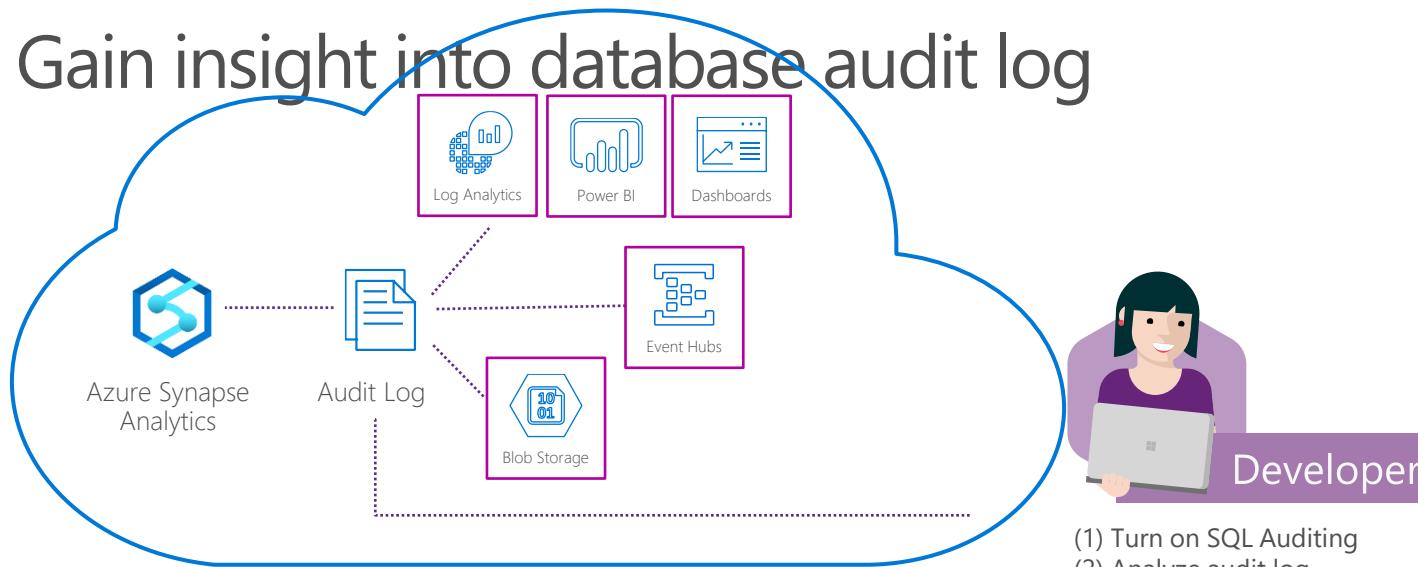


Policy manager



SQL auditing in Azure Log Analytics and Event Hubs

Gain insight into database audit log



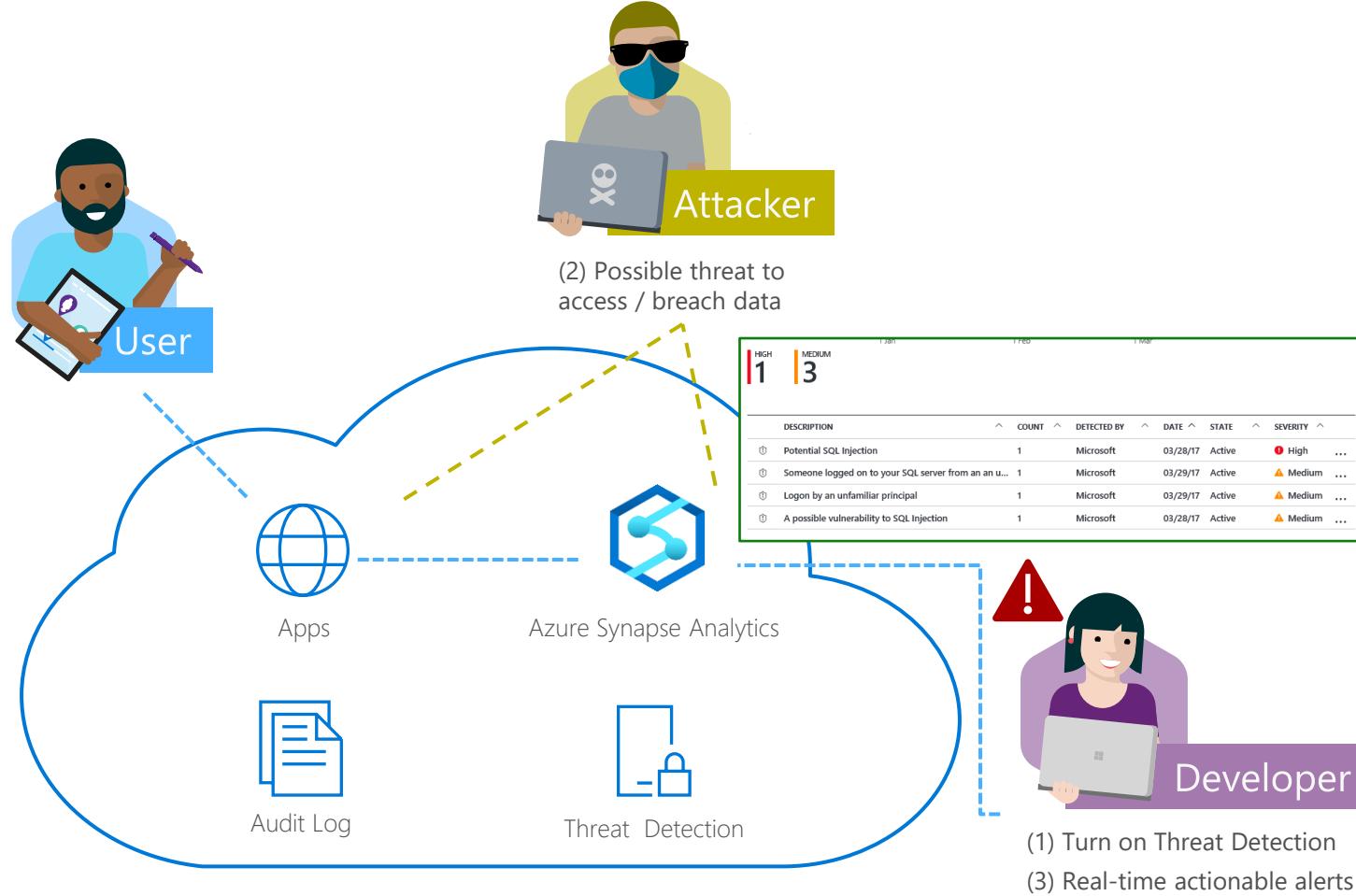
The screenshot shows a search results table with 62 entries. The columns are TimeGenerated, server_principal_name_s, statement_s, affected_rows_d, and SeverityLevel. The table lists various SQL statements executed by admin1 on 8/15/2018 at 12:00:22.521 AM, including SELECT queries and a permission grant statement.

TimeGenerated	server_principal_name_s	statement_s	affected_rows_d	SeverityLevel
8/15/2018 12:00:22.521 AM	admin1	exec sp_executesql N'SELECT tbl.name AS [Name] SCHEMA_NAME(tbl...	0	
8/15/2018 12:00:22.521 AM	admin1	exec sp_executesql N'SELECT ISNULL(HAS_PERMS_BY_NAME(QUOTEN...	1	
8/15/2018 12:00:22.521 AM	admin1	DECLARE @edition sysname; SET @edition = cast(SERVERPROPERTY(N'...	4	
8/15/2018 12:00:22.521 AM	admin1	exec sp_executesql N'SELECT CAST(@is_enabled AS bit) AS [isEnabled]...	0	
8/15/2018 12:00:22.521 AM	admin1	IF OBJECT_ID ('[sys].[database_query_store_options]') IS NOT NULL BE...		

- ✓ Configurable via audit policy
- ✓ SQL audit logs can reside in
 - Azure Storage account
 - Azure Log Analytics
 - Azure Event Hubs
- ✓ Rich set of tools for
 - Investigating security alerts
 - Tracking access to sensitive data

SQL threat detection

Detect and investigate anomalous database activity



- ✓ Detects potential SQL injection attacks
- ✓ Detects unusual access & data exfiltration activities
- ✓ Actionable alerts to investigate & remediate
- ✓ View alerts for your entire Azure tenant using Azure Security Center



SQL data classification

Discover, classify, protect and track access to sensitive data

The screenshot shows the Azure portal interface for 'Data discovery & classification (preview)'. On the left, a sidebar lists various database management options like Overview, Activity log, Tags, and Diagnose and solve problems. The main area has two donut charts: one for 'Classified columns' (10 / 109) and another for 'Information type distribution' (10 columns). Below these are dropdown filters for Schema, Table, Column, Information type, and Sensitivity label. A table lists columns from the 'ErrorLog' table with columns for UserName, Credentials, and Confidentiality level. A separate 'Settings - Information protection' window is open, showing a list of sensitivity labels: Public, General, Confidential, Confidential - GDPR, Highly confidential, and Highly confidential - GDPR. Each label has a description and a 'Configure' button.

- ✓ Automatic [discovery](#) of columns with sensitive data
- ✓ Add [persistent](#) sensitive data labels
- ✓ Audit and detect access to the sensitive data
- ✓ Manage [labels](#) for your entire Azure tenant using Azure Security Center

Industry-leading compliance



ISO 27001



SOC 1 Type 2



SOC 2 Type 2



PCI DSS Level 1

Cloud Controls
Matrix

ISO 27018

Content Delivery and
Security AssociationShared
AssessmentsFedRAMP JAB
P-ATOHIPAA /
HITECH

FIPS 140-2

21 CFR
Part 11

FERPA



DISA Level 2



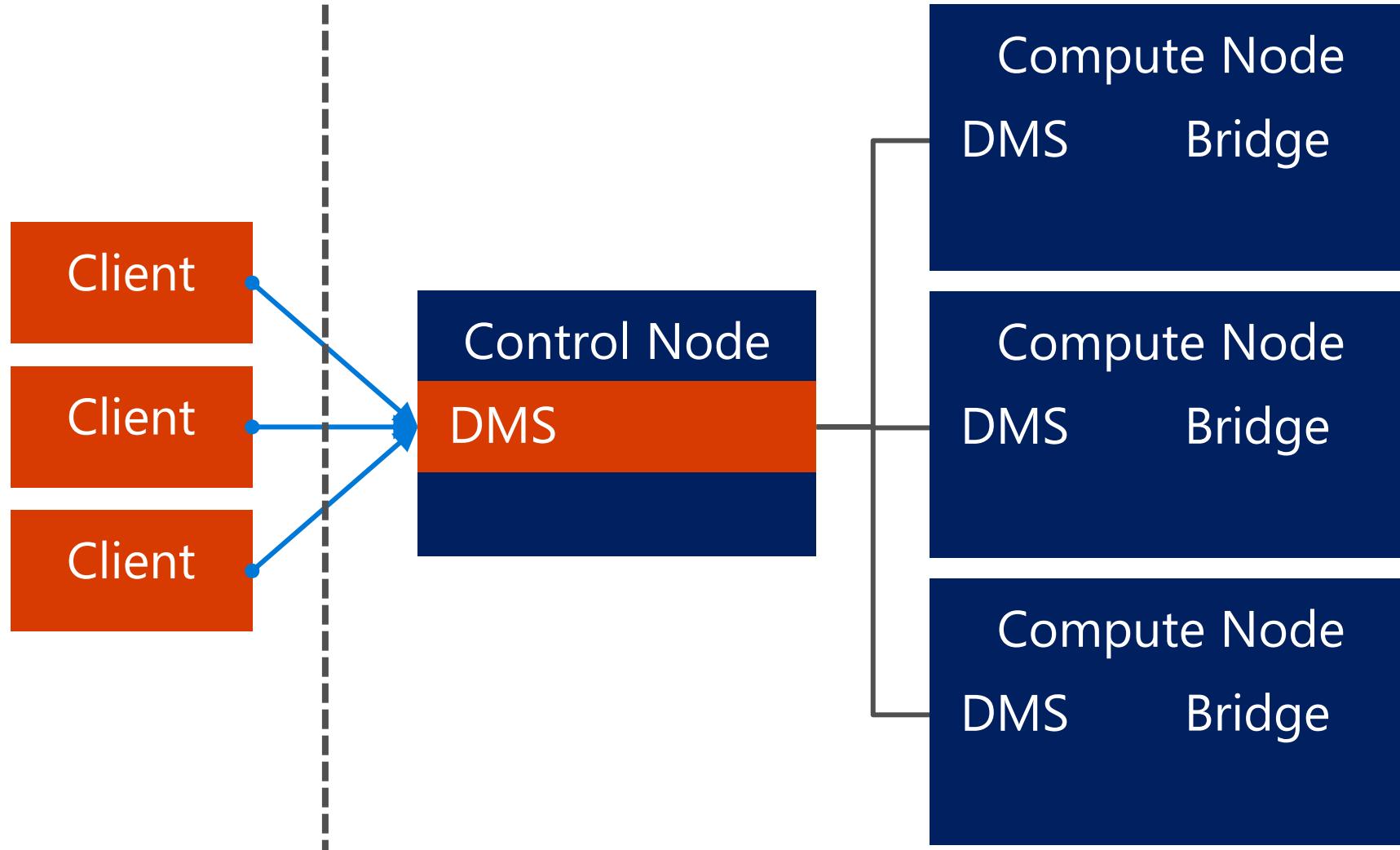
CJIS

IRS 1075
ITAR-readySection 508
VPATEuropean Union
Model ClausesEU Safe
HarborUnited
Kingdom
G-CloudChina Multi
Layer Protection
SchemeChina
GB 18030China
CCCPPFSingapore
MTCS Level 3Australian
Signals
DirectorateNew Zealand
GCIOJapan
Financial ServicesENISA
IAF

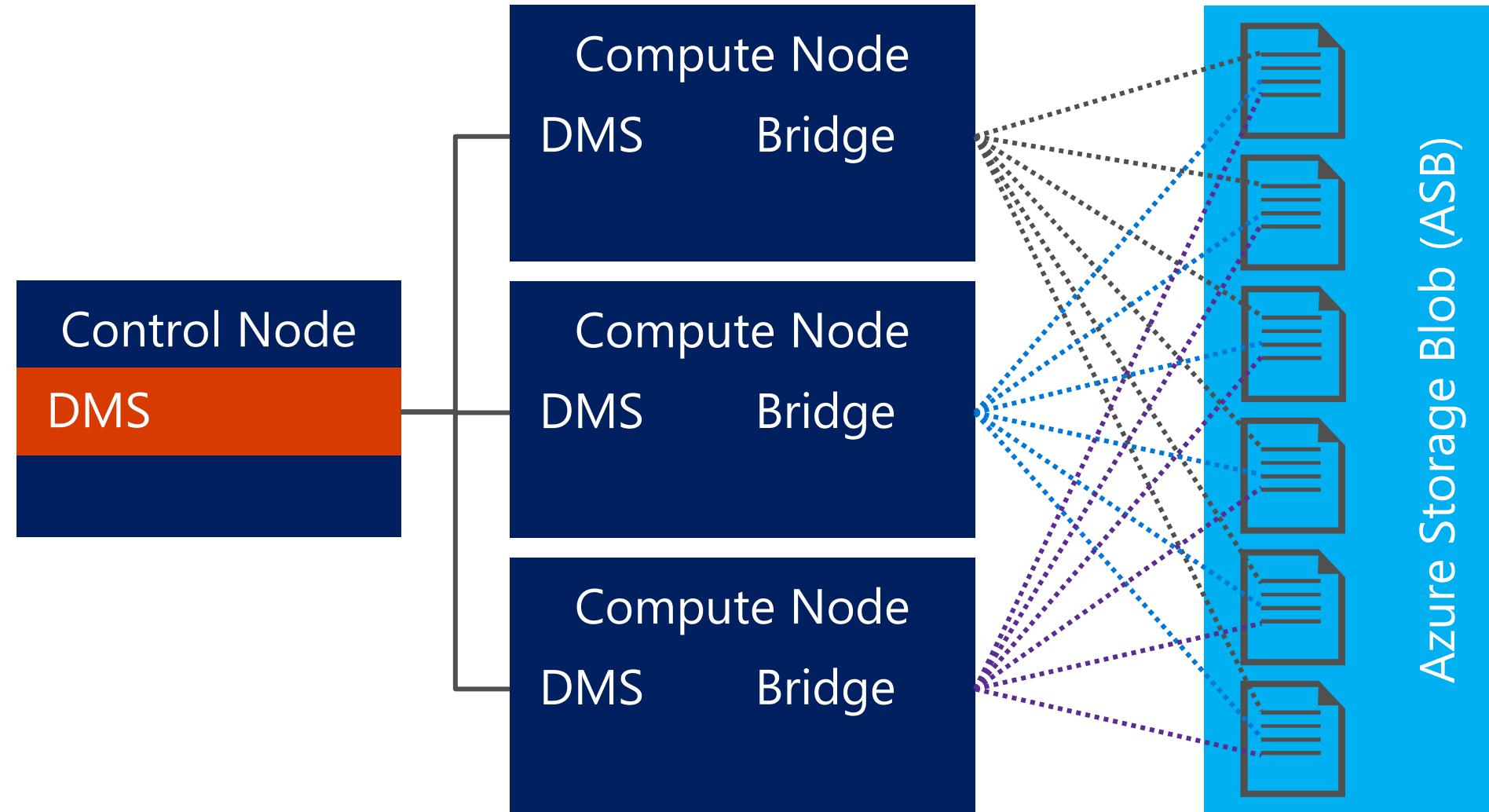


General Availability (GA) Azure Synapse Integration with Azure Services

Single gated client parallelized



Copy Into/Polybase parallel load to Azure Storage Blob



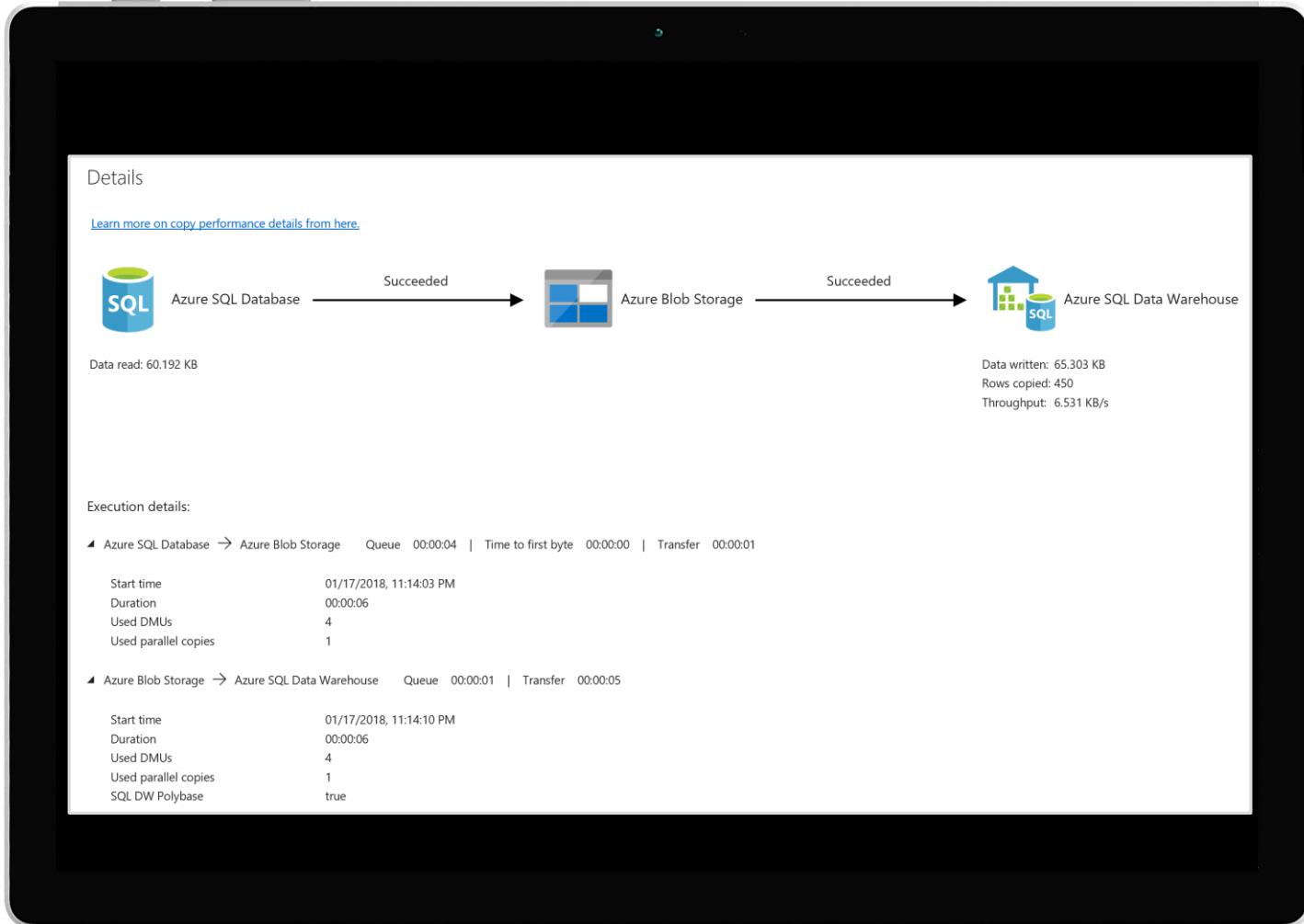
Azure Data Factory Copy activity

Overview

The Azure Data Factory Copy activity allows copying to and from Azure SQL Data Warehouse from any supported data store.

The Copy activity also supports retrieving data from a SQL source by using a SQL query or stored procedure. Authentication can be via:

- SQL Authentication
- Service principal token authentication
- Managed identity token authentication



Databricks – structured streaming

Overview

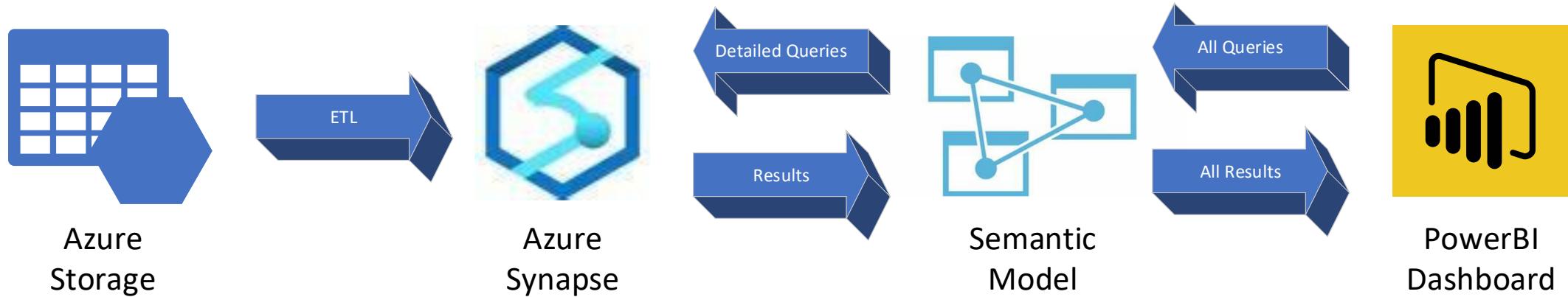
The Databricks SQL DW connector supports batch and structured streaming support for writing real-time data into Azure SQL Data Warehouse.

It uses Polybase and the Databricks structured streaming API to stream data from Kafka or Kinesis sources directly into SQL Data Warehouse at a user-configurable rate.

Source: <https://docs.azuredatabricks.net/spark/latest/data-sources/azure/sql-data-warehouse.html#streaming-support>

```
# Prepare streaming source; this could be Kafka,  
Kinesis, or a simple rate stream.  
df = spark.readStream \  
    .format("rate") \  
    .option("rowsPerSecond", "100000") \  
    .option("numPartitions", "16") \  
    .load()  
  
# Apply some transformations to the data then use  
# Structured Streaming API to continuously write the  
data to a table in SQL DW.  
df.writeStream \  
    .format("com.databricks.spark.sqldw") \  
    .option("url", <azure-sqldw-jdbc-url>) \  
    .option("tempDir",  
"wasbs://<containername>@<storageaccount>.blob.core.w  
indows.net/<directory>") \  
    .option("forwardSparkAzureStorageCredentials",  
"true") \  
    .option("dbTable", <table-name>) \  
    .option("checkpointLocation", "/tmp_location") \  
    .start()
```

PowerBI Integration



Analysis Services

- Summary data
- High Concurrency (via scale-out)
- Lower latency SLA (<3 sec)
- ~80% user queries

SQL DW

- Detailed data
- Low Concurrency
- Slightly longer duration acceptable (<30sec)
- ~20% user queries

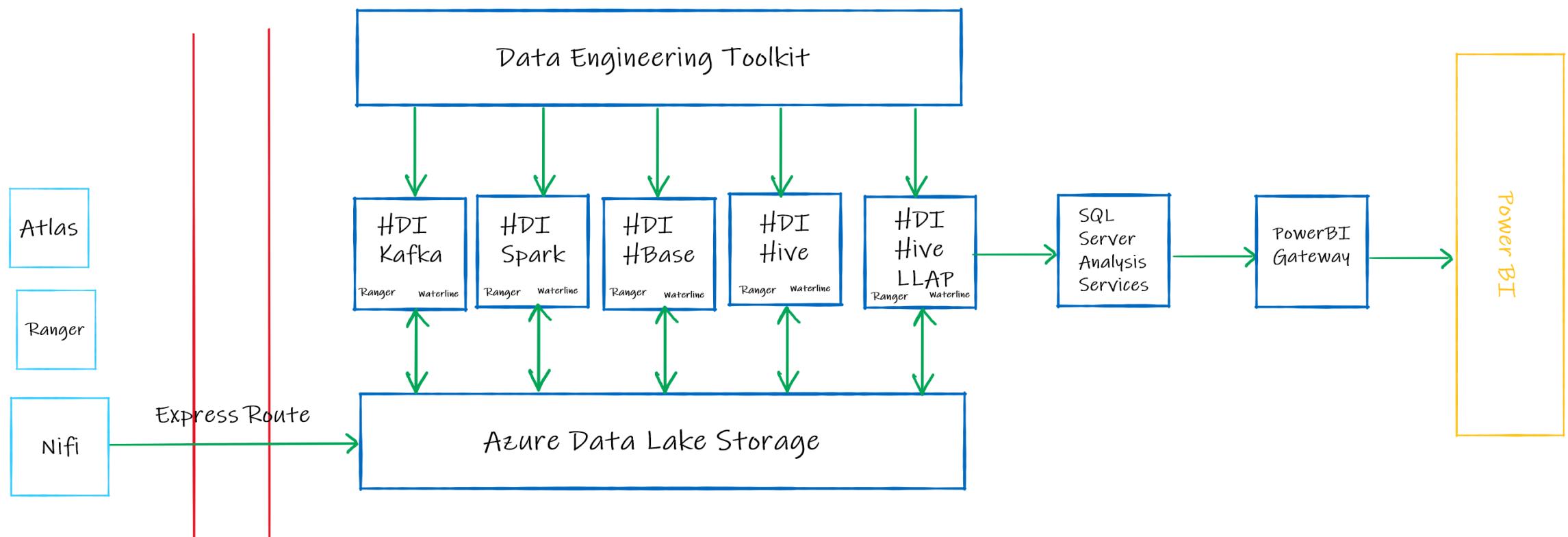


A2.0 Architecture Horizon Planning
Current, Bridge and Target Architecture
Kunal Jain, Sr. Cloud Solution Architect

Current Architecture

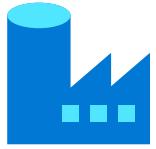
On-Premises

Azure



'Skate to where
the puck is
going, not
where it has
been'

-The Great One



Data Acquisition

This process is currently riddled with complexity and pluralism, with Watcher, Java and Spark code, and Kafka.

It is proving compute & memory intensive.

Azure Data Factory is our best of breed cloud based data integration platform.



Data Preparation

There is a heavy usage of HDInsight Edge Nodes, lacking in High Availability.

Azure Databricks provides superior data preparation in Azure Data Factory's auspice.

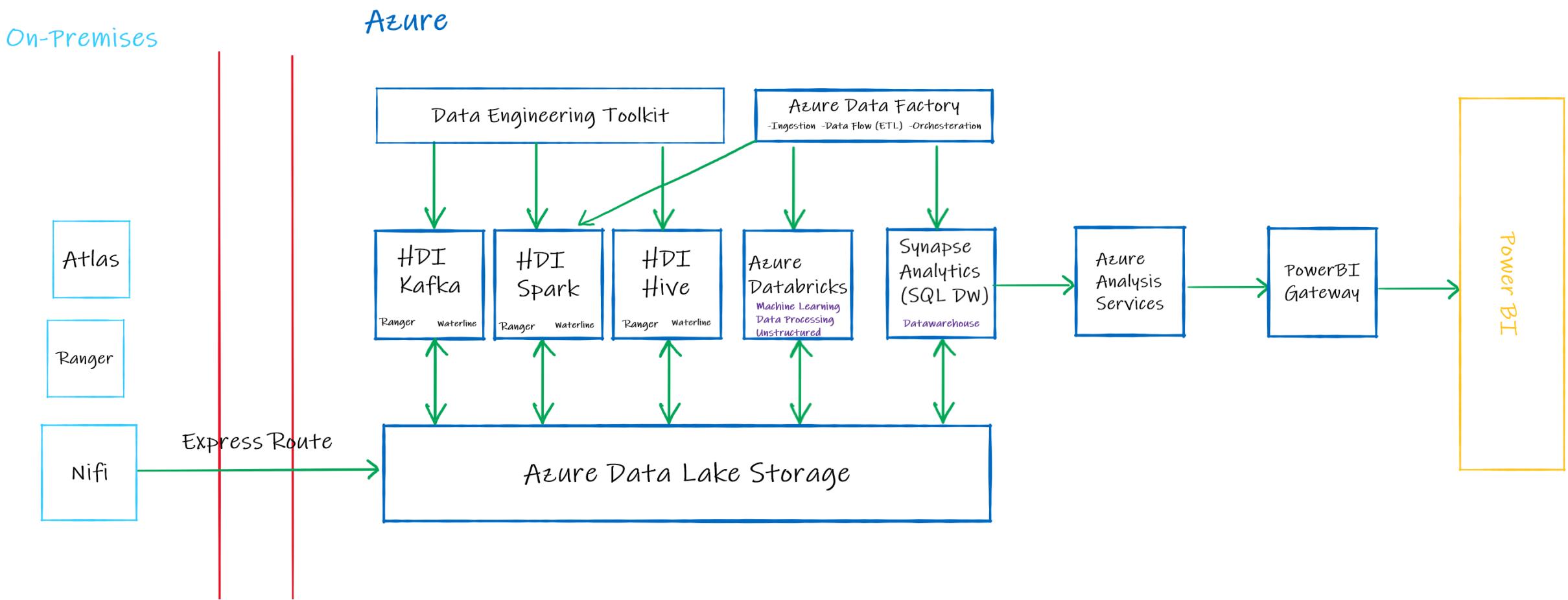


Tenant Processing

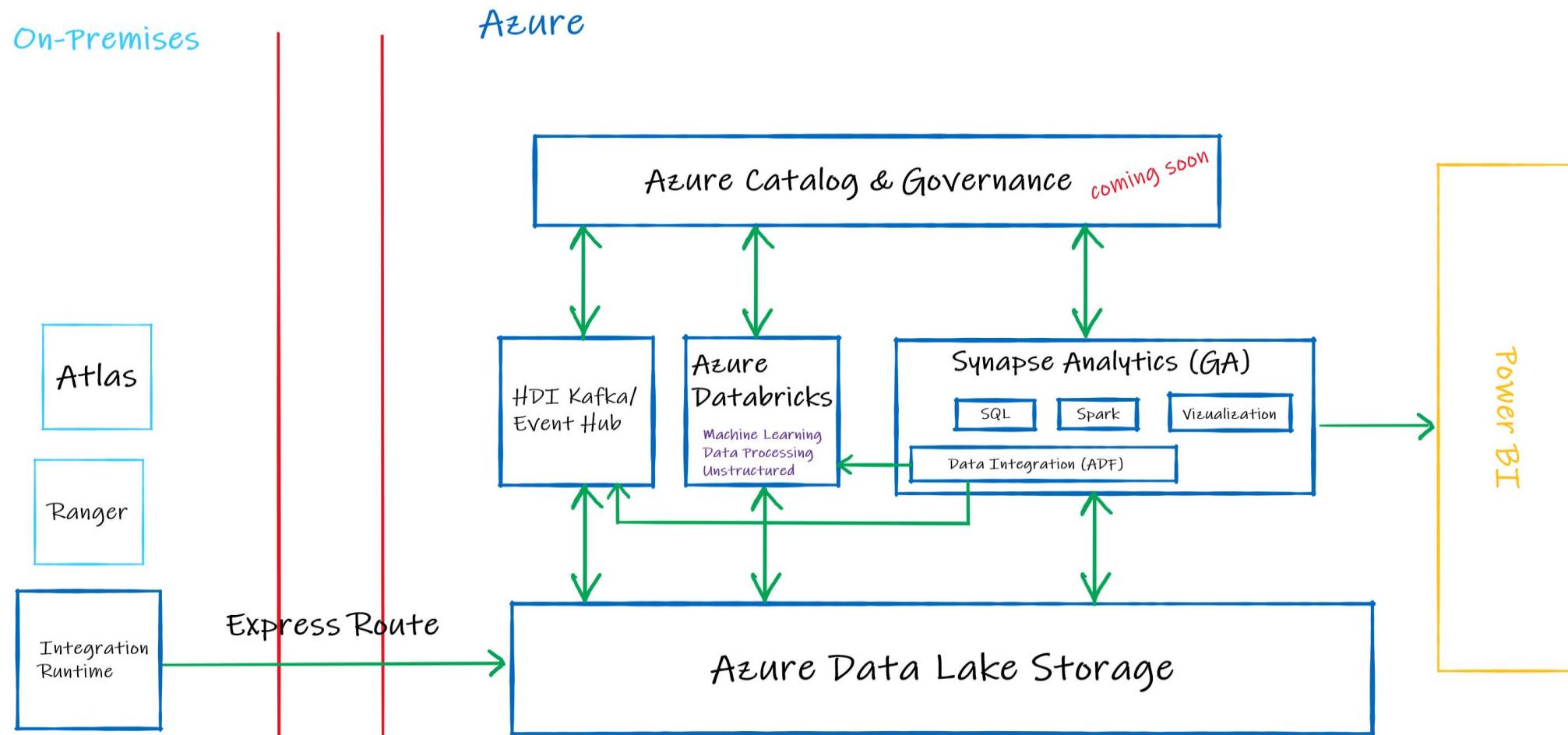
It is challenging to build Data Warehouses on HIVE and extract reports.

Azure Synapse Analytics solves specifically for this use case – a limitless analytics service and enterprise data warehouse.

Phase 1: Interim Architecture (MSFT GA Services)



Phase 2: Long-term Architecture (MSFT 2H 2020 tentative)





The Azure Data Solutions Big Picture

Microsoft Release Schedule and Kaiser Permanente GA Planning

Azure Data Services – Microsoft Release Roadmap



Microsoft GA – Q4 2019

Azure Synapse



- ✓ Materialized Views
- ✓ Ordered Columnstore
- ✓ JSON support
- ✓ Dynamic Data Masking
- ✓ SSDT support
- ✓ Read committed snapshot isolation
- ✓ Resultset caching

Azure Data Factory



- ✓ Manageability enhancements for self-hosted IR
- ✓ Availability Zone support
- ✓ VNET support
- ✓ Region expansion
- ✓ Hybrid, serverless, scalable data integration in Synapse
- ✓ ADF Dataflow on Synapse Spark
- ✓ ADF Synapse VNET integration
- ✓ Mapping Data Flows UX improvements
- ✓ CDM/ODI integration
- ✓ Mapping Data Flows Embeddability, Error Handling, Hierarchy Support
- ✓ Mapping Data Flows - Continuous ETL
- ✓ Wrangling Data Flows GA

Microsoft Preview

- ✓ Synapse Studio
- ✓ Collaborative workspaces
- ✓ Distributed T-SQL Query service
- ✓ SQL Script editor
- ✓ Unified security model
- ✓ Notebooks
- ✓ Apache Spark
- ✓ Code-free data flows
- ✓ Orchestration Pipelines
- ✓ Data movement

- ✓ Workload Isolation
- ✓ Simple ingestion with COPY
- ✓ Share DW data with Azure Data Share
- ✓ Private LINK support
- ✓ Streaming ingestion & analytics in DW
- ✓ Native Prediction/Scoring
- ✓ Fast query over Parquet files
- ✓ FROM clause with joins

Azure HDInsights



- ✓ ORC (r/w) for Data Flow
- ✓ Excel (r), XML (r), Snowflake (r/w), Delta lake (r/w) for Copy & Data Flow
- ✓ SharePoint List (r) for Copy
- ✓ Oracle CDC
- ✓ Copy connector extensibility [Public Preview]
- ✓ Data consistency check, skip error files
- ✓ Editable metadata-driven approach for 100's of sources & 10's of 1000's of table objects

- Code-Free Transformation**
- ✓ Mapping: Known Schema to Known Schema Mapping (GA)
- ✓ Wrangling: Self-Service Data Exploration (Preview)

- Trustworthy Computing**
- ✓ Virtual Networks support at scale (GA)
- ✓ Private IP Support (GA)
- ✓ Availability Zone Support (GA)
- Data Movement**
- ✓ Change Data Capture (Preview)
- ✓ Extensibility and New Connectors (Preview/GA)

- ✓ Kafka REST proxy (GA)
- ✓ Spark 2.4.4, Spark 3.0 (Preview), HBase 2.1, Kafka 2.3
- ✓ Spark .Net GA
- ✓ Spark Synapse DW connector, Spark ADX Connector

- Orchestration & Monitoring**
- ✓ Pipeline SLA management and violation debugging
- ✓ Late data handling enhancements
- ✓ Gantt and operational lineage views

- Data Catalog and Metadata Management**
- ✓ Publish ADF Datasets for discovery and reuse (Preview)
- ✓ Publish Data Flows and Pipelines for Data Lineage (Preview)

- ✓ Premium File Share support for HBase WAL and Kafka Logs
- ✓ Hive/Spark HBase/Spark Integration Improvements
- ✓ Auto-scale for HBase & LLAP (GA)
- ✓ Availability Zone support

Kaiser Permanente Adoption Feasibility Roadmap



Kaiser Permanente GA – June 2020

Azure Synapse



- ✓ Materialized Views
- ✓ Ordered Columnstore
- ✓ JSON support
- ✓ Dynamic Data Masking
- ✓ SSDT support
- ✓ Read committed snapshot isolation
- ✓ Resultset caching

Azure Data Factory



- ✓ Manageability enhancements for self-hosted IR
- ✓ Availability Zone support
- ✓ VNET support
- ✓ Region expansion
- ✓ Hybrid, serverless, scalable data integration in Synapse
- ✓ ADF Dataflow on Synapse Spark
- ✓ ADF Synapse VNET integration
- ✓ Mapping Data Flows UX improvements
- ✓ CDM/ODI integration
- ✓ Mapping Data Flows Embeddability, Error Handling, Hierarchy Support
- ✓ Mapping Data Flows - Continuous ETL
- ✓ Wrangling Data Flows GA

Azure HDInsights



Kaiser Permanente GA

(Microsoft GA +3-6 Months)

- ✓ Synapse Studio
- ✓ Collaborative workspaces
- ✓ Distributed T-SQL Query service
- ✓ SQL Script editor
- ✓ Unified security model
- ✓ Notebooks
- ✓ Apache Spark
- ✓ Code-free data flows
- ✓ Orchestration Pipelines
- ✓ Data movement

- ✓ Workload Isolation
- ✓ Simple ingestion with COPY
- ✓ Share DW data with Azure Data Share
- ✓ Private LINK support
- ✓ Streaming ingestion & analytics in DW
- ✓ Native Prediction/Scoring
- ✓ Fast query over Parquet files
- ✓ FROM clause with joins

Code-Free Transformation

- ✓ Mapping: Known Schema to Known Schema Mapping (GA)
- ✓ Wrangling: Self-Service Data Exploration (Preview)

Orchestration & Monitoring

- ✓ Pipeline SLA management and violation debugging
- ✓ Late data handling enhancements
- ✓ Gantt and operational lineage views

Trustworthy Computing

- ✓ Virtual Networks support at scale (GA)
- ✓ Private IP Support (GA)
- ✓ Availability Zone Support (GA)

Data Catalog and Metadata Management

- ✓ Publish ADF Datasets for discovery and reuse (Preview)
- ✓ Publish Data Flows and Pipelines for Data Lineage (Preview)

Data Movement

- ✓ Change Data Capture (Preview)
- ✓ Extensibility and New Connectors (Preview/GA)

Kafka REST proxy (GA)

- ✓ Spark 2.4.4, Spark 3.0 (Preview), HBase 2.1, Kafka 2.3
- ✓ Spark .Net GA
- ✓ Spark Synapse DW connector, Spark ADX Connector

Premium File Share support for HBase WAL and Kafka Logs

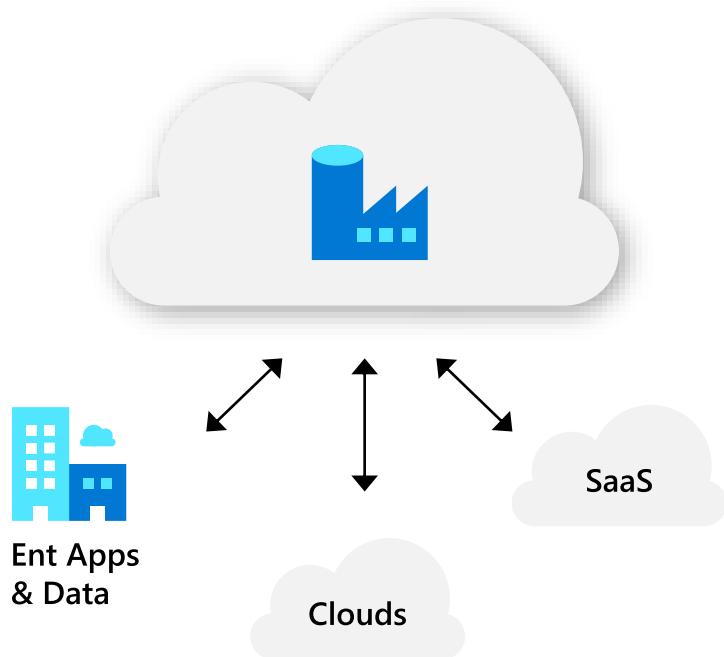
- ✓ Hive/Spark HBase/Spark Integration Improvements
- ✓ Auto-scale for HBase & LLAP (GA)
- ✓ Availability Zone support



Azure Data Factory
Anand Subbaraj, Principal Program Manager

Azure Data Factory

Data Integration Service: Serverless, Scalable, Hybrid



Data Movement and Transformation @Scale

Cloud & Hybrid w/ 90+ connectors provided
Up to 2 GB/s, ETL/ELT in the cloud

Hybrid Pipeline Model

Seamlessly span: on premise, Azure, other clouds & SaaS
Run on-demand, scheduled, data-availability or on event

Author & Monitor

Programmability w/ multi-language SDK
Visual Tools

SSIS Package Execution

Lift existing SQL Server ETL to Azure
Use existing tools (SSMS, SSDT)

Azure Data Factory Statistics

>10 MILLION
data factories created in
the past 12 months

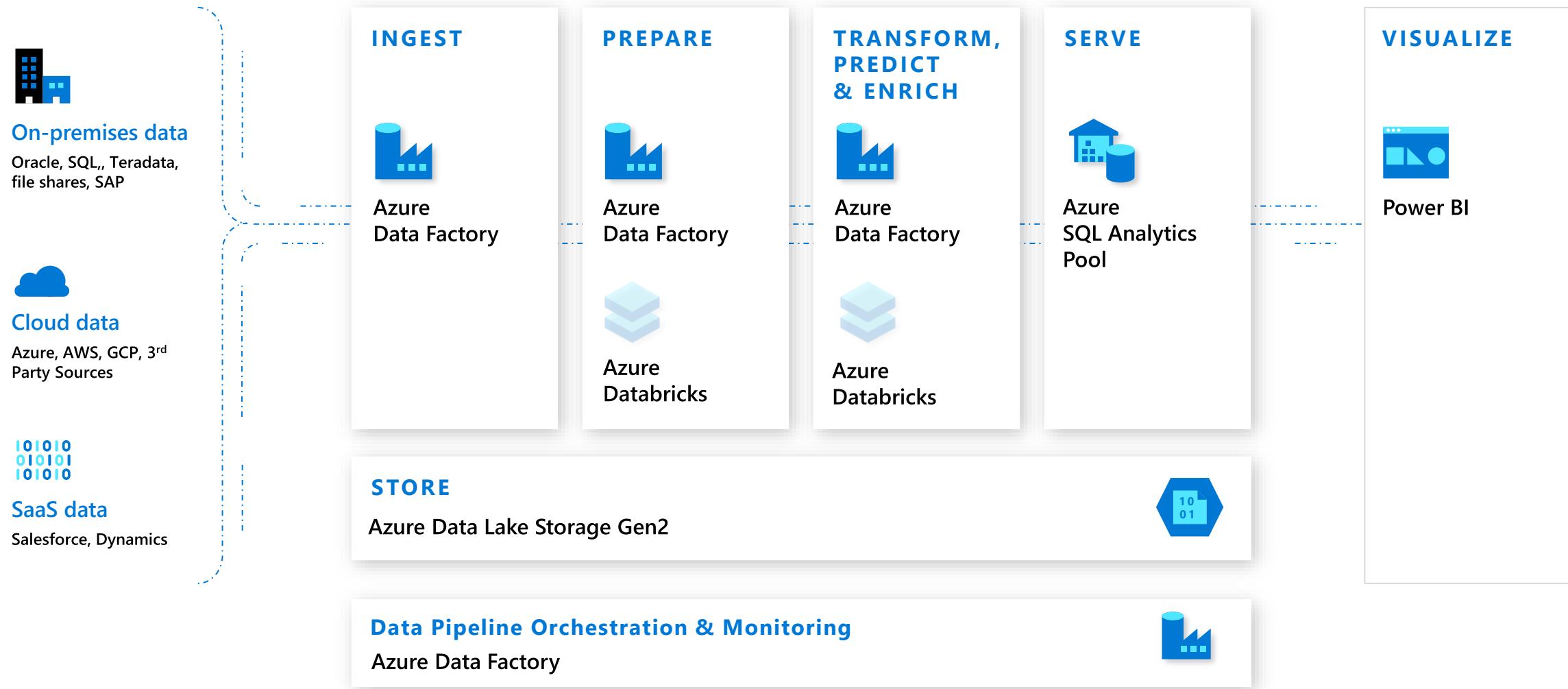
>1 BILLION
activity runs per month

30+ PBs
of data moved monthly

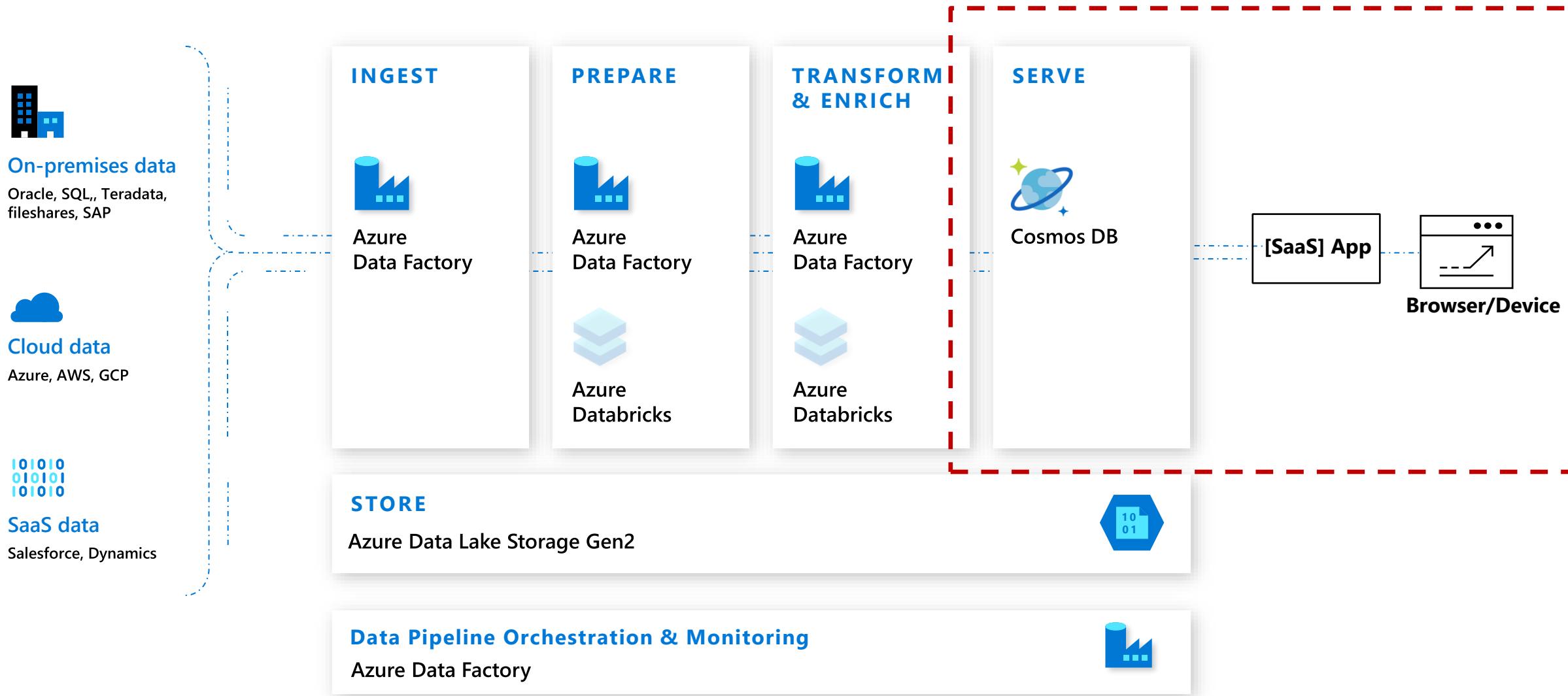
Azure Data Factory

Modern Data Engineering Scenarios

Modern Data Warehouse (MDW)

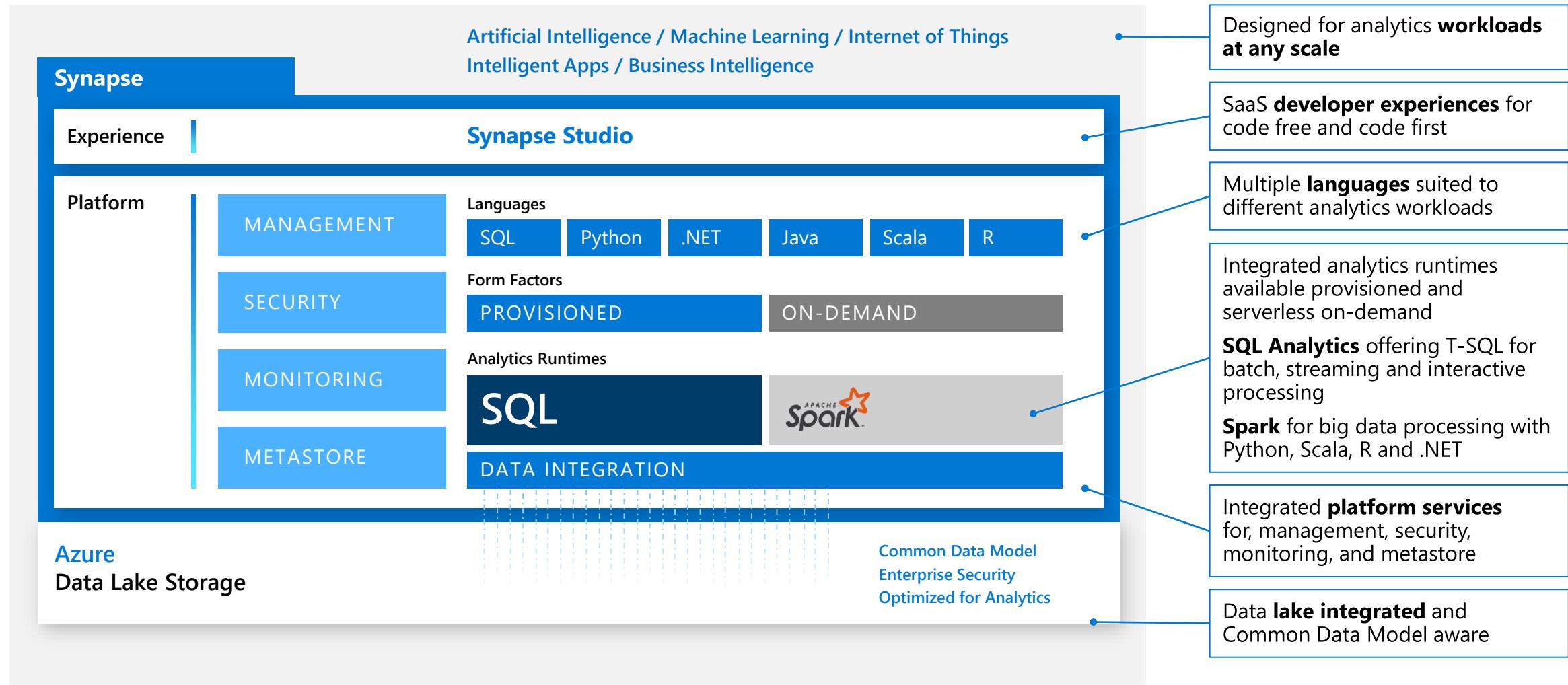


Analytics for data-driven apps



Azure Synapse Analytics

Integrated data platform for BI, AI and continuous intelligence



Data Movement



Scalable

per job elasticity
Up to 4 GB/s

Simple

Visually author or via code (Python, .Net, etc)
Serverless, no infrastructure to manage

Access all your data

90+ connectors provided and growing (cloud,
on premises, SaaS)
Data Movement as a Service: 20 points of
presence world wide
Self-hostable Integration Runtime for hybrid
movement

Access all your data - 90+ Built-in Connectors & Growing

Azure (15)	Database & DW (26)		File Storage (6)	File Formats (6)	NoSQL (3)	Services & Apps (28)		Generic (4)			
Blob Storage	Amazon Redshift	Oracle	Amazon S3	AVRO	Cassandra	Amazon MWS	Oracle Service Cloud	HTTP			
Cosmos DB – SQL API	DB2	Phoenix	File System	Binary	Couchbase	CDS for Apps	PayPal	OData			
Cosmos DB – MongoDB API	Drill	PostgreSQL	FTP	Delimited Text	MongoDB	Concur	QuickBooks	ODBC			
ADLS Gen1	Google BigQuery	Presto	Google Cloud Storage	JSON		Dynamics 365	Salesforce	REST			
ADLS Gen2	Greenplum	SAP BW Open Hub	HDFS	ORC		Dynamics AX	SF Service Cloud				
Data Explorer	HBase	SAP BW MDX	SFTP	Parquet		Dynamics CRM	SF Marketing Cloud				
Database for MariaDB	Hive	SAP HANA				Google AdWords	SAP C4C				
Database for MySQL	Impala	SAP Table				HubSpot	SAP ECC				
Database for PostgreSQL	Informix	Spark				Jira	ServiceNow				
File Storage	MariaDB	SQL Server				Magento	Shopify				
SQL Database	Microsoft Access	Sybase				Marketo	Square				
SQL Database MI	MySQL	Teradata				Office 365	Web Table				
SQL Data Warehouse	Netezza	Vertica				Oracle Eloqua	Xero				
Search Index						Oracle Responsys	Zoho				
Table Storage											
					Support read & write						
					Support read only						

Data Movement Scenarios

Ingest data using ADF to bootstrap your analytics workload

KEY SCENARIO

WHY ADF

Data migration for data lake & EDW

1. Big data workload migration from AWS S3, on-prem Hadoop File System, etc
2. EDW migration from Oracle Exadata, Netezza, Teradata, AWS Redshift, etc

- **Tuned for perf & scale:** PBs for data lake migration, tens of TB for EDW migration
- **Cost effective:** serverless, PAYG
- Support for **initial snapshot & incremental catch-up**

Data ingestion for cloud ETL

1. Load as-is from variety of data stores
2. Stage for data prep and rich transformation
3. Publish to DW for reporting or OLTP store for app consumption

- **Rich built-in connectors:** file stores, RDBMS, NoSQL.
- **Hybrid connectivity:** on-prem, other public clouds, VNet/VPC
- **Enterprise grade security:** AAD auth, AKV integration
- **Developer productivity:** code-free authoring, CICD
- **Single-pane-of-class monitoring** & Azure Monitor integration

Customer Wins

Scenario: Data migration from Amazon S3 to Azure Storage. Migrate existing analytics workload from AWS into Azure, over time modernize workload using Azure data services.

CASE #1

Requirements:

- Minimal migration duration
- Max load throughput

Result:

- **2PB data from S3 → Blob in <11 days**
- Initial load: 1.9 PB, avg. throughput **2.1 GB/s**
- Incremental: 221 TB, avg. throughput **3.6 GB/s**

CASE #2

Requirements:

- Data traverse through private channel instead of public internet;
- Optimize network egress charges out of AWS by going through AWS Direct Connect

Result:

- **1PB from S3 → Blob over 10Gbps private link**
- Avg. throughput **787 MB/s**

Customer Wins (continued)

CASE #3

Scenario:

- Standardize on an enterprise-wide analytics solution following typical MDW pattern
- **Lake hydration from on-prem Netezza**

Requirements:

- Designated migration window between 18:00-08:00 from Monday to Friday
- Limit overhead on Netezza: max 8 concurrent DB connections

Result:

- **25TB on-prem Netezza to ADLS**
- Total migration duration: ~3 weeks

CASE #4

Scenario:

- Build a modern multi-platform DW that **brings together silos of data across BUs into a central data lake.**

Requirements:

- Better perf and scale with elasticity to deal with big data.
- Fit-for-purpose tools for transformation and analytics that supports a wide spectrum of personas (data engineer, data scientists, data integrator).

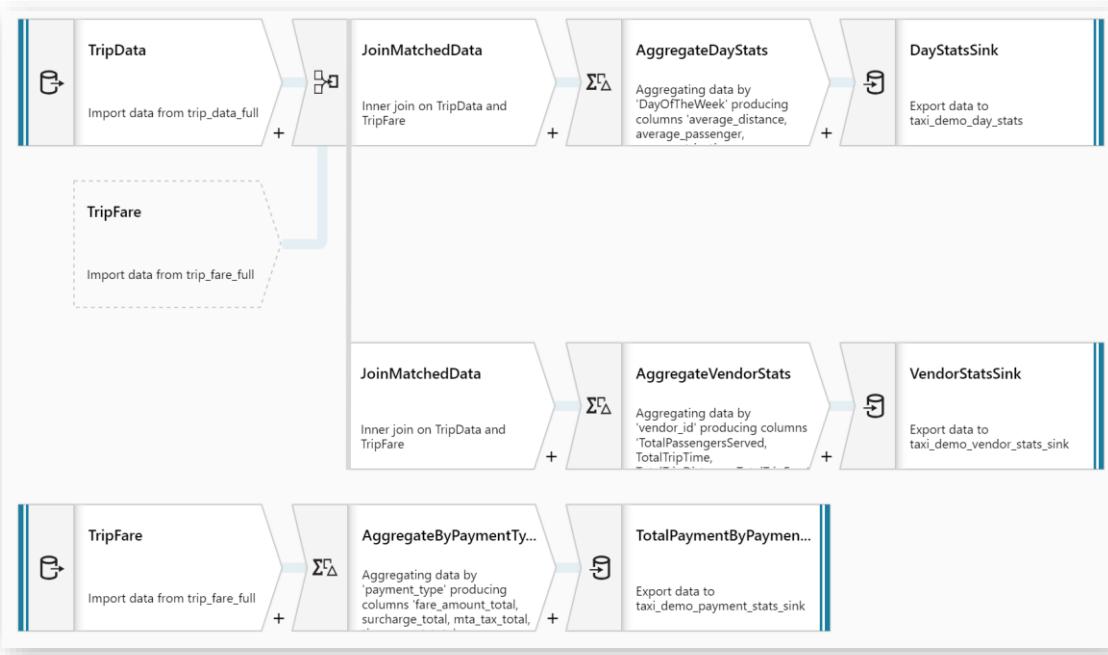
Result:

- Ingest data from **various sources** – SFTP, SAP, Amazon S3, Google BigQuery, SQL Server, Oracle, Teradata, File, APIs, etc.

No Code Data Transformation At Scale

Focus on building business logic and data transformation

- Data cleansing, transformation, aggregation, conversion, etc.
- Cloud scale via Spark execution
- Easily build resilient data flows



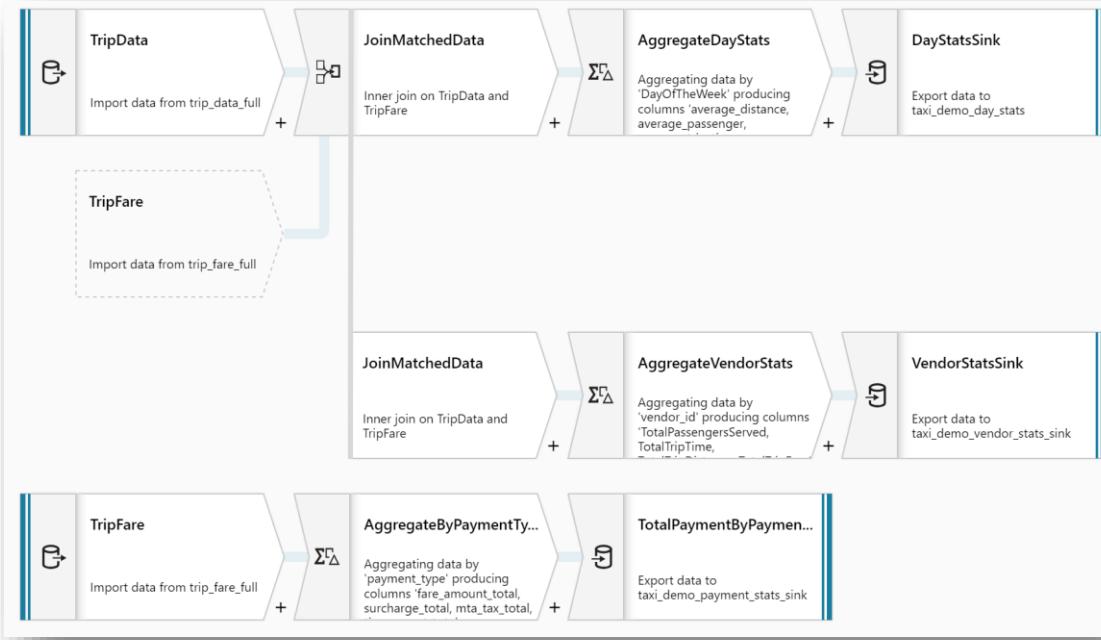
... not

```
File MovieRecommendationE2EDemo.txt X
1 HOI Cluster Details:
2 Adfd1.azurehdinsight.net
3 Admin
4 Adf@123456
5 Storage:
6 adfd1storage
7 /anyPw#6G1j7T8lBMMm1So/YgdJyGT4d+S1]Ar+sMs7b7gb954706gjOM1oksZ19uXso740xZo8xIKhdWQ==
8
9 Cluster Remote Login Details:
10 Adf
11 India1234
12
13 HiveQuery:
14 DROP TABLE IF EXISTS MovieRatings;
15 CREATE EXTERNAL TABLE MovieRatings
16 (
17   UserID int,
18   MovieID int,
19   Rating int,
20   Timestamp string
21 ) ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' LINES TERMINATED BY '\n' STORED AS TEXTFILE LOCATION '$(hiveconf:MovieRatings)';
22
23 DROP TABLE IF EXISTS MovieTitles;
24 CREATE EXTERNAL TABLE MovieTitles
25 (
26   MovieID int,
27   MovieName string
28 ) ROW FORMAT DELIMITED FIELDS TERMINATED BY ',' LINES TERMINATED BY '\n' STORED AS TEXTFILE LOCATION '$(hiveconf:MovieTitles)';
```

Generally Available

MAPPING DATAFLOW

Code-free data transformation @scale



PUBLIC PREVIEW

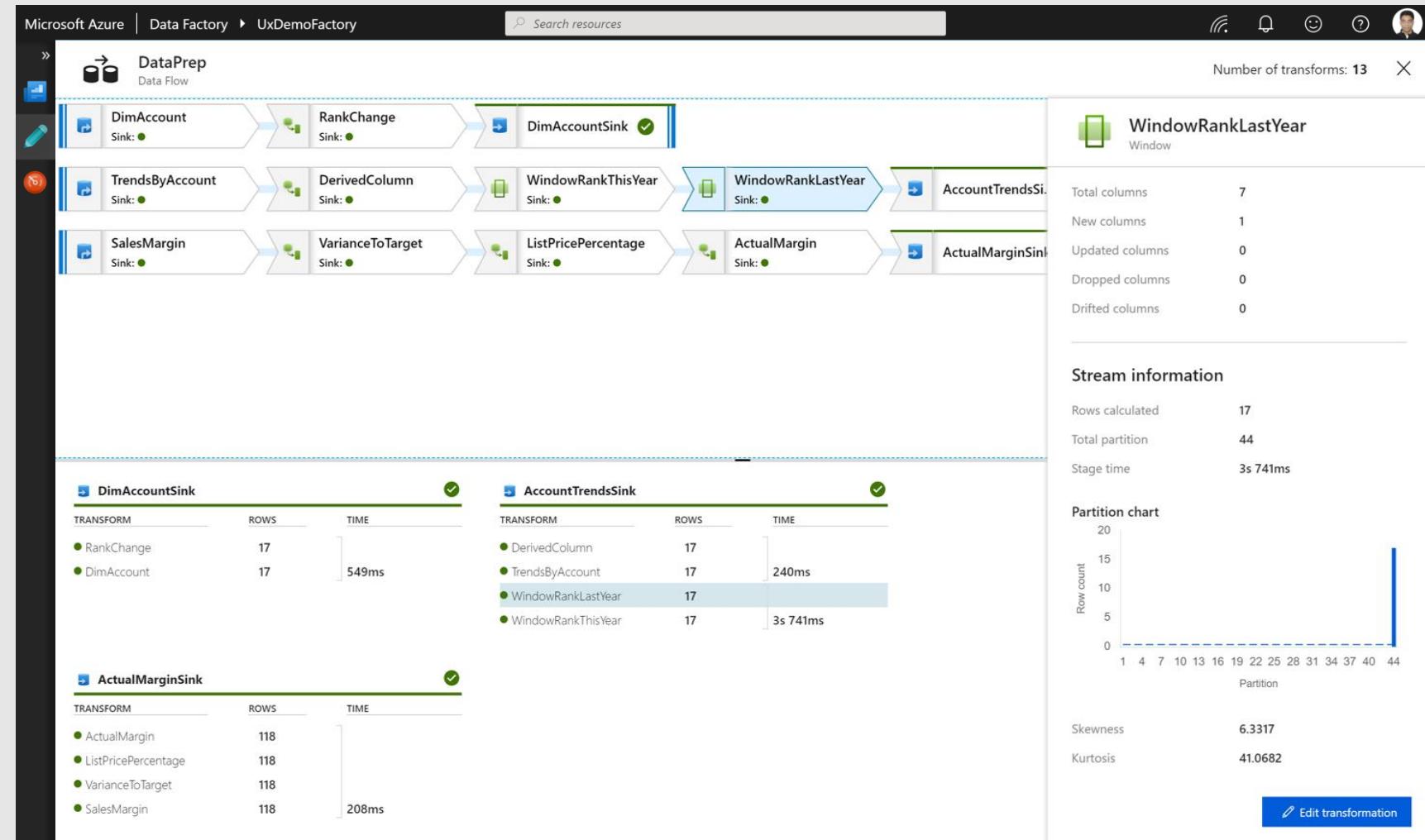
WRANGLING DATAFLOW

Code-free data preparation @scale

The screenshot shows the Azure Data Factory Wrangling Dataflow interface. On the left, there's a sidebar with "ADFRessource [2]" and "WrangleUserQuery". The main area displays a table titled "Table.RemoveColumns(#'Renamed CustomerID', '(CustomerID')"). The table has columns: #C CustId, #C FirstName, #C LastName, #C City, #C ZIP, #C Email, #C State, and #C BasePay. The data consists of 24 rows of Harry Potter characters. On the right, there's a panel titled "Name: WrangleUserQuery" with "Applied steps" listed: "Source", "Renamed columns", "Merged queries", "Renamed CustomerID", and "Removed columns".

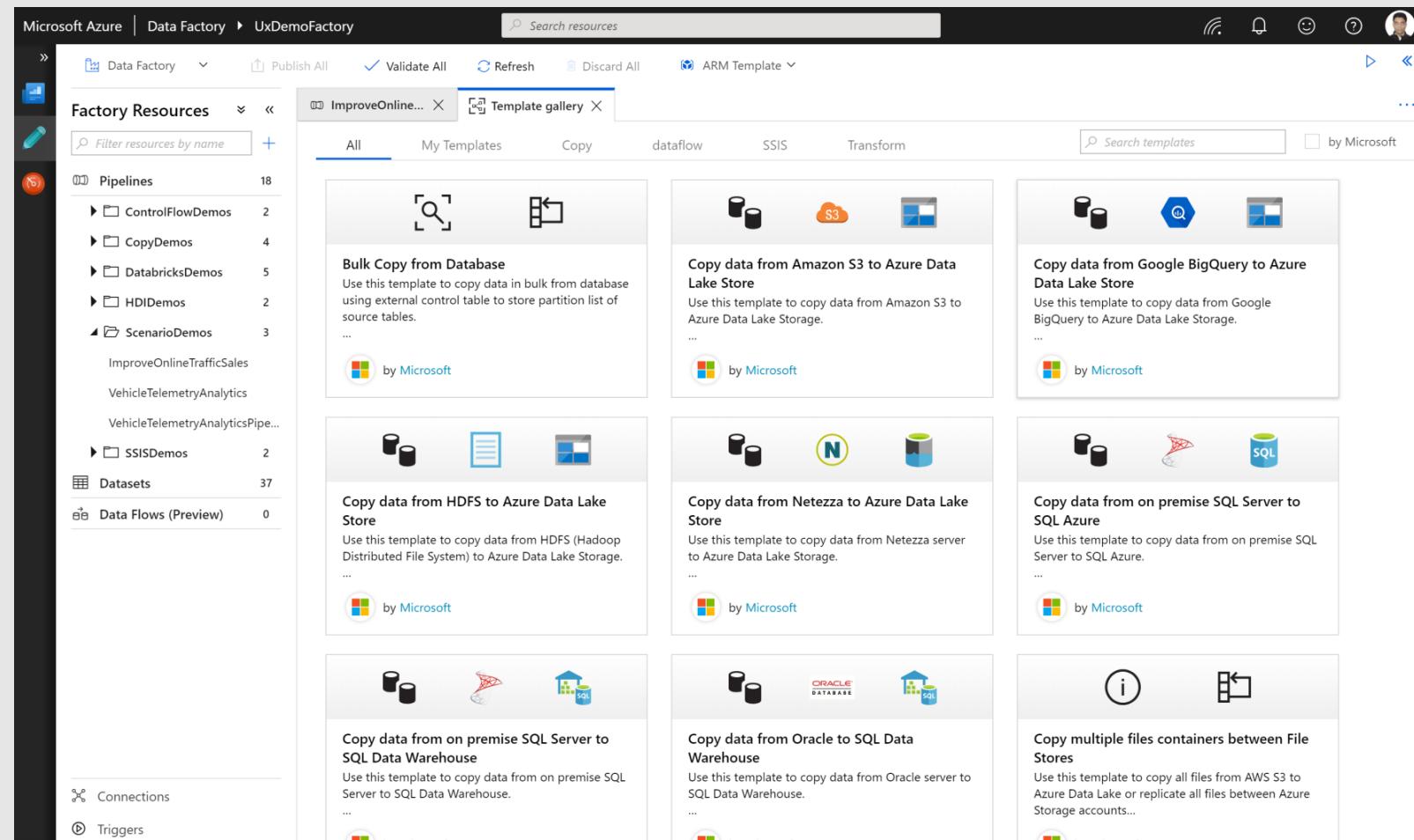
Best in class monitoring and management

- Monitor Pipeline and Activity Runs
- Rich language to query Runs
- Operational lineage between parent-child pipelines
- Azure Monitor Integration
 - Diagnostics logging
 - Metrics & Alerts
 - Events
- Restate Pipeline and Activities

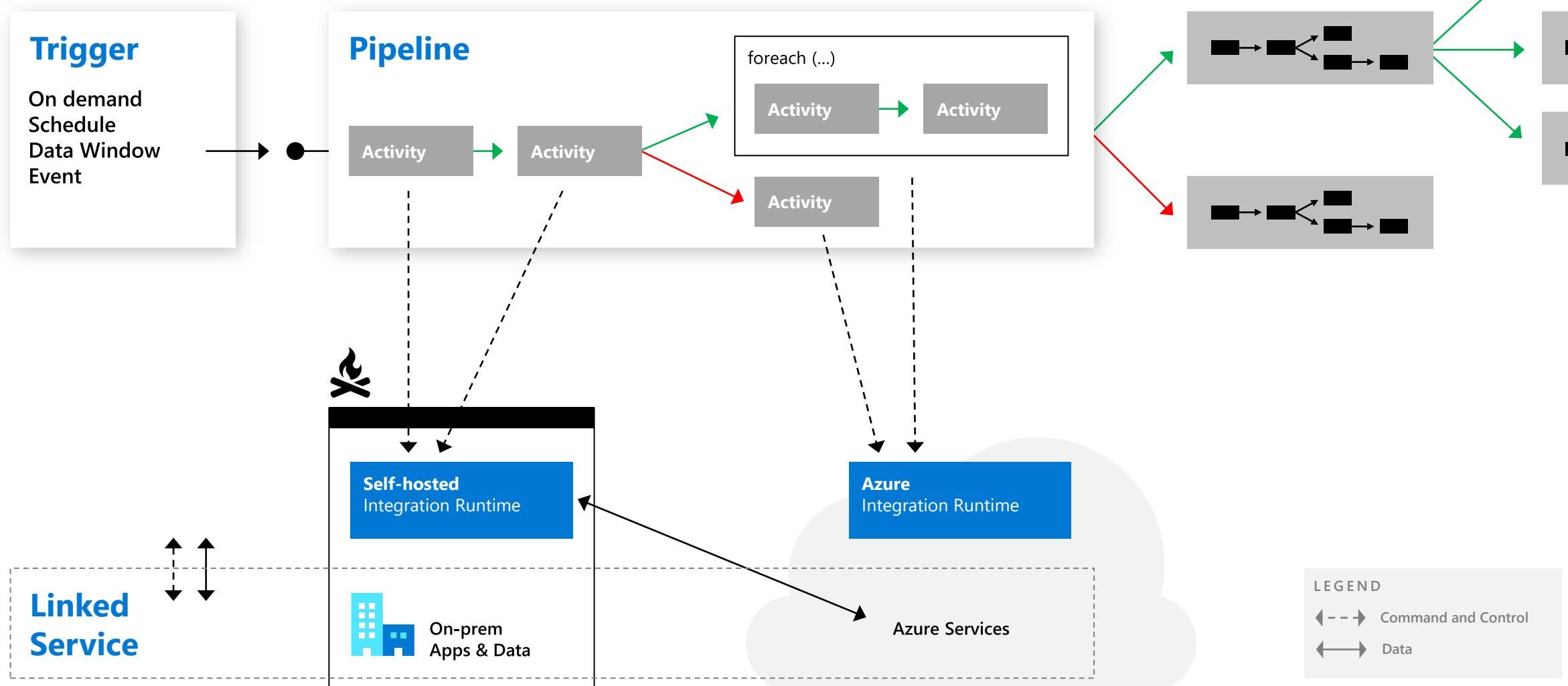


Use Templates to quickly get started with ADF

- Quickly get started with building data integration solutions
- Avoid building same workflows repeatedly.
Simply instantiate a template
- Improve developer productivity along with reducing development time for repeat processes

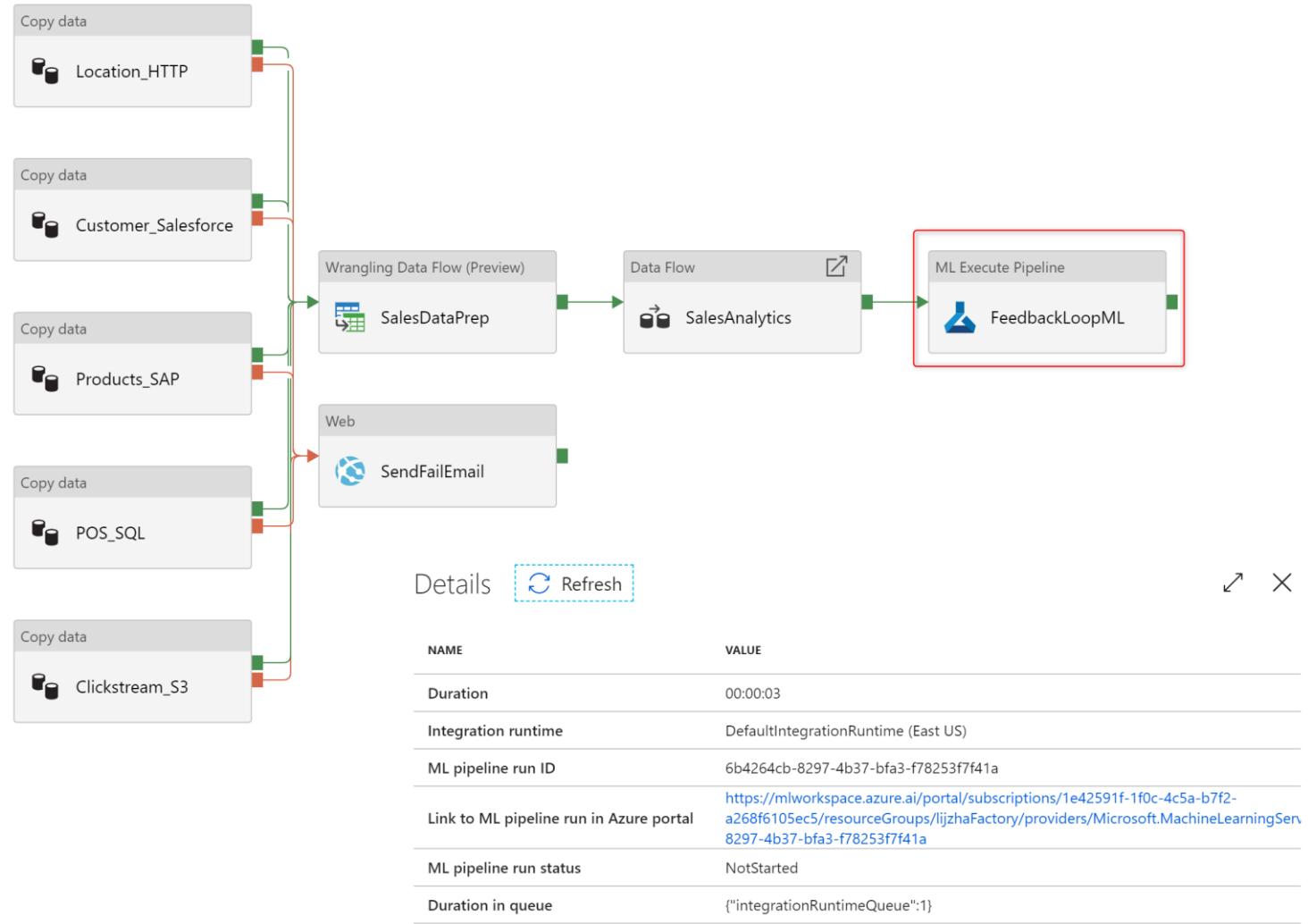


Orchestration @ Scale

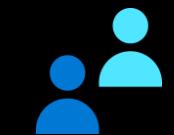


Run AML service pipelines as a step in your ADF pipelines

- Use AML service to quickly train, deploy machine learning pipelines at scale. Execute deployed AML service pipelines as part of your ETL/ELT ADF workflows.
- Ingest data from multiple sources (90 plus connectors) and run AML service pipelines that do training, retraining and batch prediction scenarios
- Deep linking to monitor Azure ML service pipeline runs from ADF UX

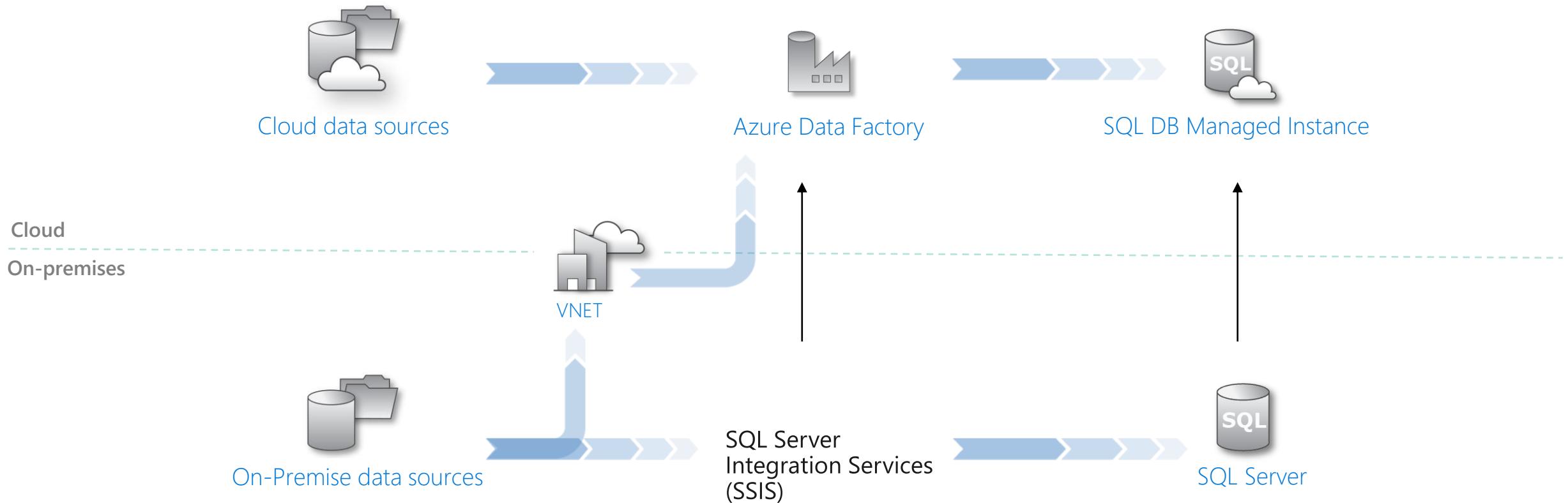


DEMO



Lift & Shift

Lift Existing SSIS + SQL Projects to Cloud



Partners



SentryOne®



KingswaySoft



Additional Information: <https://aka.ms/ssisazureupdates>

CY19 H1 Updates

Code-Free Transformation

- Mapping Data Flow (Preview)
- Wrangling Data Flow (Preview)
- Managed Spark

Data Movement

- New Connectors (85+ connectors)
- Error Handling – Skip Fault Rows, Resume from checkpoint
- Performance – 2X improvement in cloud-to-cloud copy performance

Developer Productivity

- Continuous Integration / Deployment
- Git Integration
- Pipeline Templates

SQL Integration Services

- Custom Setup + 3rd Party Licensing/Extensibility/Ecosystem
- First-Class SSIS Activities in ADF Pipelines
- Enterprise Edition

CY19 H2 Deliverables

Code-Free Transformation

- Mapping: Known Schema to Known Schema Mapping (GA)
- Wrangling: Self-Service Data Exploration (Preview)

Trustworthy Computing

- Virtual Networks support at scale (GA)
- Private IP Support (GA)
- Availability Zone Support (GA)

Data Movement

- Change Data Capture (Preview)
- Extensibility and New Connectors (Preview/GA)

Orchestration & Monitoring

- Pipeline SLA management and violation debugging
- Late data handling enhancements
- Gannt and operational lineage views

Data Catalog and Metadata Management

- Publish ADF Datasets for discovery and reuse (Preview)
- Publish Data Flows and Pipelines for Data Lineage (Preview)

ADF Roadmap: H1 CY20

Item	Description
Enterprise Readiness & fundamentals	<ul style="list-style-type: none">Manageability enhancements for self-hosted IRAvailability Zone supportVNET supportRegion expansion
Synapse Analytics	<ul style="list-style-type: none">Hybrid, serverless, scalable data integration in SynapseADF Dataflow on Synapse SparkADF Synapse VNET integration
Code-free ETL	<ul style="list-style-type: none">Mapping Data Flows UX improvementsCDM/ODI integrationMapping Data Flows Embeddability, Error Handling, Hierarchy SupportMapping Data Flows - Continuous ETLWrangling Data Flows GA
Code-based ETL	<ul style="list-style-type: none">Databricks activities support MSIScript activity against data stores for transformation: Snowflake, SQL
Data Movement	<p><u>Connectivity</u></p> <ul style="list-style-type: none">ORC (r/w) for Data FlowExcel (r), XML (r), Snowflake (r/w), Delta lake (r/w) for Copy & Data FlowSharePoint List (r) for CopyOracle CDCCopy connector extensibility [Public Preview] <p><u>Data Movement Reliability</u></p> <ul style="list-style-type: none">Data consistency check, skip error filesEditable metadata-driven approach for 100's of sources & 10's of 1000's of table objects
Orchestration	<ul style="list-style-type: none">Trigger dependenciesGlobal ParametersLate data handling



Azure HDInsight

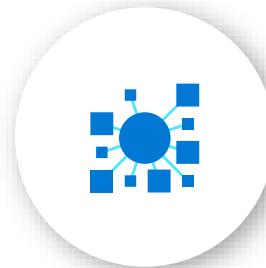


Open Source, Customizable & Extensible



- 100% Apache Open Source
- Data ingestion (Kafka) to interactive exploration (LLAP, Spark SQL)
- No fuss migration of workloads from on-prem or other clouds
- Curated application platform for wide variety of use cases

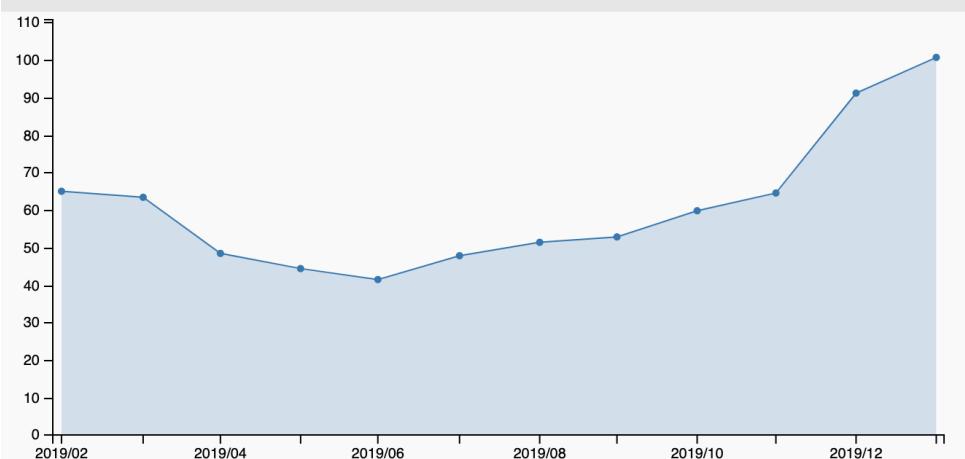
Enterprise ready, Managed & Secure



- 99.9% availability SLA
- Cluster Health Monitoring with Ambari and Azure Log Analytics
- Cost control through intelligent Autoscale
- Azure AD, Kerberos authentication and Role-based access control
- Strong network isolation through VNETs and service endpoint support



HDInsight is focused on meeting customers where they are



KP Growth (# of Cluster running in HDInsight)

Greetings Dana and Microsoft team,

I know that on most occasions, communication from me seems to be escalations. Sorry about that 😊.

This communication is about acknowledging and celebrating the accomplishments we have achieved TOGETHER. All the hard work and focused attention from the Joint KP/Microsoft team is paying off.

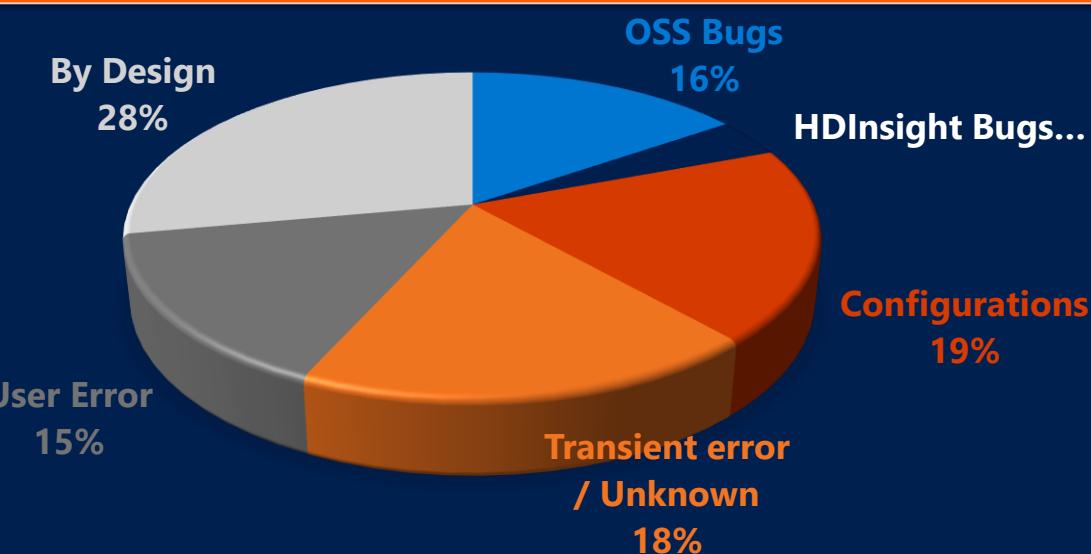
Here are some key highlights of progress to date.

1. **Service Enablement:** 11 MS Azure services have been hardened and enabled; Azure certified for PHI.
2. **Automation:** 25+ ARM Templates created and integrated to DevOps pipeline
3. **Data:** One time data loads for Clarity, Onelink, PDW, CDW and Member month amounting to 530 terabyte has been onboarded. Incremental data loads in progress.
4. **Data Management:** Tag based security with Database, row, column level security enabled
5. **Capabilities:** 18 real-time events from Epic enabled from Epic in production; predictive models deployed.
6. **Business Value:** 4 use cases tenants onboarded to production; Insight Driven – OWL App, Metrics Mart, PDR Suite, Risk Adjustment, including 400+ users.
7. **Performance:** Performance benchmark completed with comparable query response times between onpremise Engineered DataWarehouse and HDI which had 1/10 cores; concurrent query testing underway.

While we have much work ahead of us to accelerate the adoption, I wanted to take a moment to acknowledge the tremendous progress made to date and **THANK YOU** for the partnership.

Ganesh Thondikulam,
Executive Director, Analytics Digital Foundation

Root Cause	Root Cause Drill Down	2020-02	2020-01	2019-12	2019-11	2019-10	2019-09	2019-08	Total
Azure platform issues	Azure VM issues	0	0	1	0	0	0	0	1
Bug	HDInsight	3	2	4	1	1	1	0	12
Bug	Hadoop -OSS	0	0	1	0	0	0	0	1
Bug	OMS Agent	3	0	0	0	0	0	0	3
By Design	HDInsight	3	6	2	3	1	1	2	18
By Design	Hadoop - OSS	1	3	2	2	0	0	1	9
By Design	Other	0	0	1	0	0	0	1	2
Configuration	HDInsight	3	3	1	0	2	0	2	11
Configuration	Hadoop - OSS	0	0	1	0	0	0	4	5
Configuration	Other	1	1	2	0	0	0	0	4
HDInsight Custom configuration		0	0	2	0	0	0	1	3
Lack of documentation		0	0	1	0	0	0	0	1
Scenario not supported		1	0	0	0	0	0	0	1
Transient error / Unknown		0	5	3	2	2	1	6	19
User Authentication and authorization issues		1	1	0	0	0	0	0	2
User Environment/Platforms		0	0	0	0	0	1	0	1
User Error		7	1	2	1	1	2	2	16



KP/HDInsight Issues by category

HDInsight strategy & roadmap



Azure HDInsight =

Open Source Frameworks

+

Platform

Key Drivers:

- Stability Reliability of Mission Critical Production applications
- Market Churn & Dependency on Hortonworks "HDP"
- Contributing to the OSS community
- Azure/Cloud Specific improvements in OSS
- Onboarding new OSS frameworks

Key Drivers:

- Containers & Microservices
- Cluster Upgrade without downtime
- Version Management
- Simpler Security Model
- Cost Control/ Effective use of resources
- Cluster Scale up Scale Down latencies
- Configuration Management

HDInsight & Apache OSS Frameworks



HDInsight /w HDP



HDInsight /w "Microsoft Distribution"

Customer Promise# Maintain 100% Compatibility

Microsoft distribution of popular Apache Analytics frameworks

Apache analytics projects built, delivered and supported completely by Microsoft

Significant test Investments to improve customer reliability

Leveraged by multiple Microsoft products (HDInsight, SQL BDC, Cosmos, Synapse)

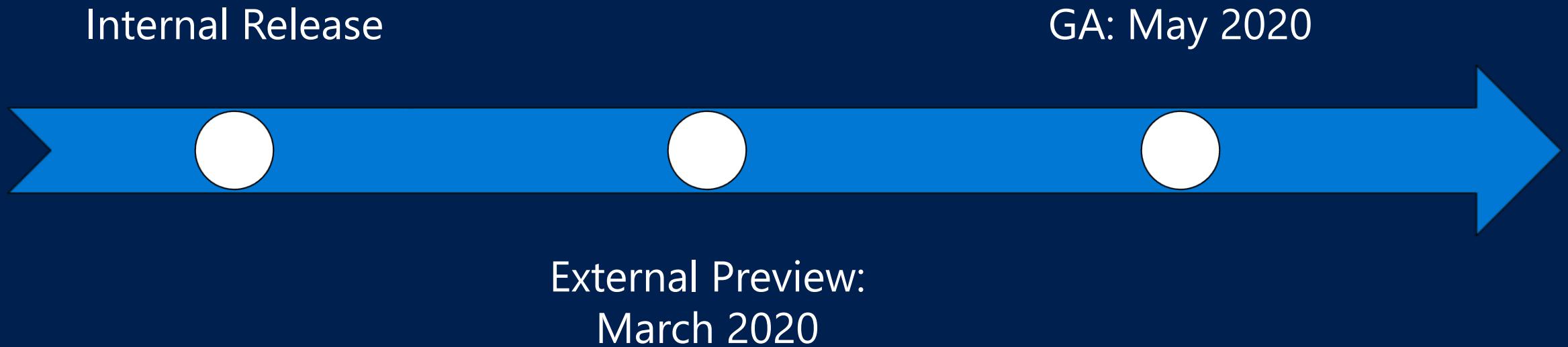
Apache projects **enhanced with Microsoft's** years of experience with Big Data analytics

Innovations by Microsoft offered back to the community



Frameworks	HDInsight 4.0	HDInsight 3.6
Apache Hadoop and YARN	3.1.1	2.7.3
Apache Tez	0.9.1	0.7.0
Apache Pig	0.16.0	0.16.0
Apache Hive	3.1.0	2.1.0, 1.2.1
Apache Ranger	1.1.0	7.0
Apache HBase	2.1	1.1.2
Apache Sqoop	1.4.7	1.4.6
Apache Oozie	4.3.1	4.2.0
Apache Zookeeper	3.4.6	3.4.6
Apache Storm		1.1.0
Apache Mahout		0.9.0
Apache Phoenix	5	4.7.0
Apache Spark	2.4	2.1, 2.2, 2.3.1
Apache Livy	0.5	0.3, 0.4
Apache Kafka	2.1	1.0, 1.1, 1.1.1
Apache Ambari	2.7.0	2.6.0
Apache Zeppelin	0.8.0	0.7.3

Microsoft distribution of popular Apache Analytics frameworks

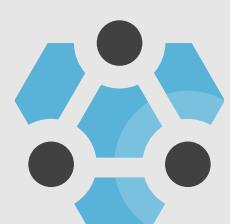


We are ready to onboard your test subscriptions

HDInsight Platform

HDInsight on Kubernetes (Planned Features)

- **Multi-Cluster Management**
- **Cluster Cloning**
- **Configuration Management**
- **Zero Downtime Upgrades**
- **Auto Deletion**
- **Version Management**
- **Private Networks**
- **Workload Monitoring**



H2 FY20 Outlook

OPEN SOURCE ANALYTICS

HDInsight w/ (Microsoft Open Source Analytics Distribution)

Kafka REST proxy (GA)
Spark 2.4.4, Spark 3.0 (Preview), HBase 2.1, Kafka 2.3
Spark .Net GA
Spark Synapse DW connector, Spark ADX Connector
Premium File Share support for HBase WAL and Kafka Logs
Hive/Spark HBase/Spark Integration Improvements

PLATFORM IMPROVEMENTS

Auto-scale for HBase & LLAP (GA)

Availability Zone support
HBase: Encryption at REST (Write Ahead Logs)
Support **ADLS Gen 2** (Premium)

ENTERPRISE GRADE SECURITY

TLS Enforcement

HDInsight Identity Broker for **MFA** (GA)
Deploy HDInsight clusters in **Private VNETs**
Disk Encryption, Encryption of intra-node traffic & shuffle data

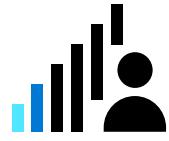


Azure Data Discovery & Governance

Codename: *Babylon*

Jennifer Stevens, Azure Catalog Program Manager

Why Your Data Estate Needs Governance



Data Consumers &
Data Producers

What data do I have?
What is in the data?

Where did the data originate?
Can I make a change?

Who else used the data?
Can I trust the data?



Data Officers

What business sensitive or
PII data exists in my org?
How can I understand risk
level across my data?

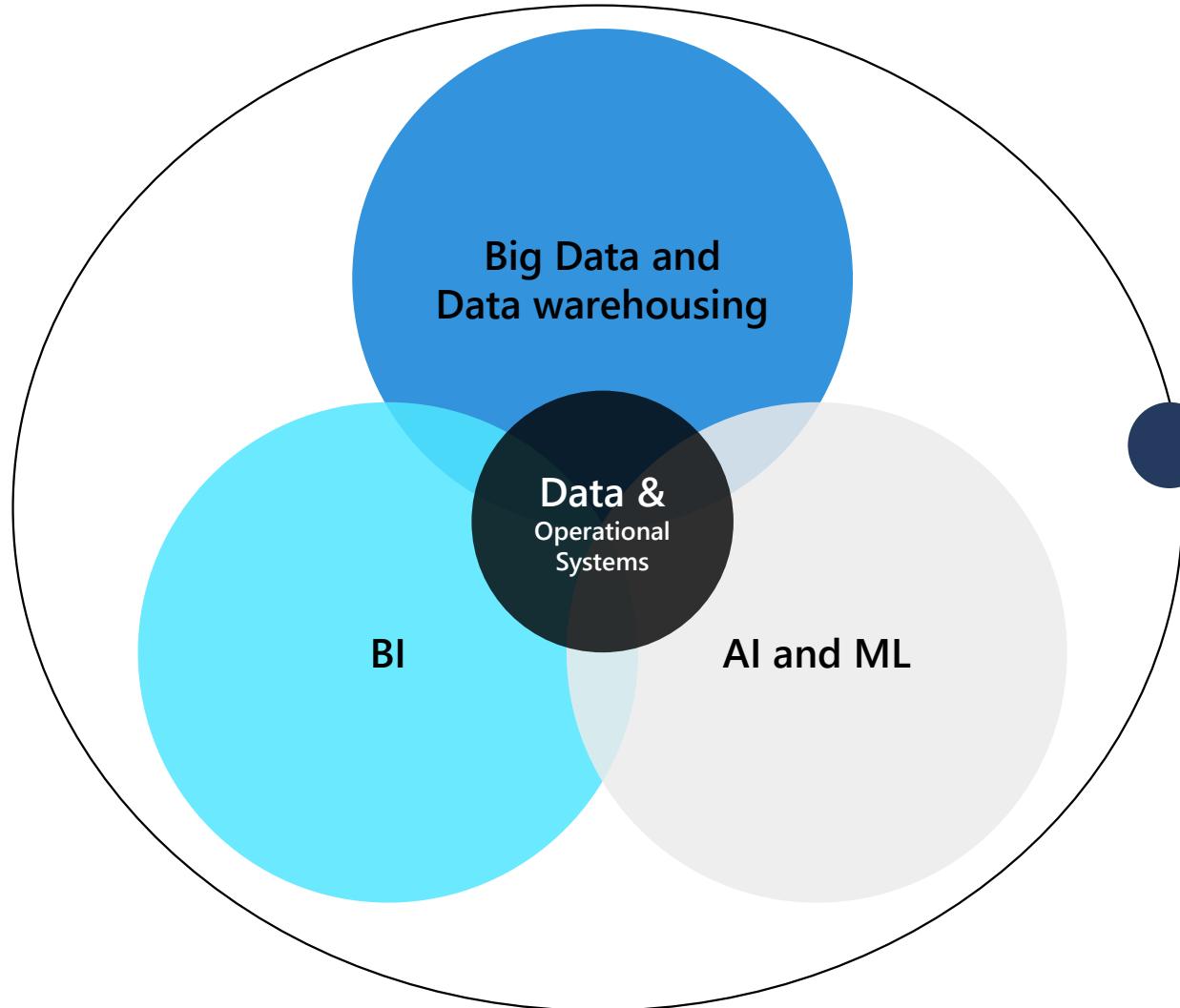
Who is using the data and
for what purpose?

Is our data compliant with
regulations and corporate
policies?

Microsoft
Empowers You
to Answer These Questions



A single pane of glass for data estate governance



A single pane of glass to...
Inventory data
Discover & curate
Assess data compliance & protection
Set & manage data policy

Data Governance and Protection



Map

Automated data estate scanning & classification tells you **what data you have**, and where it is



Discover

Data Catalog to find data, explore metadata & lineage to understand if it is fit for your purpose

Technical users see technical metadata

Business users explore and search using concepts familiar to them



Insights

Intelligent “data estate guardians” keep watch over your data estate alerting you to key issues and helping you assess help & risk

Assess regulatory risk, ensure **sensitive data is protected** and **unearth data breaches** as soon as they occur



Protect

Data **policies** that enforce **right-use** of all data with no code

Ensure data is accessible to the right users at the right time for a well-known purpose.

Ensure data location, movement and sharing complies to enterprise & regulatory rules



Self-Service

Self service without compromise. Users can access the data they need — and only what they’re allowed to see

Map & Discover Capabilities - Private Preview

Connectors

- SQL Server, SQL DB/DW
- Azure Data Lake Gen1 and Gen2 plus blobs
- Cosmos DB
- Azure Files

Parsers

- JSON, CSV/TSV/SSV, XML, ORC, Parquet & Avro

Classification Rules

- Support for detecting common personal data types, e.g., name, address, email, phone number, credit card, geolocation, etc.
- Custom classification rule support

Platform

- Atlas APIs for scanning ingest, search, and metadata management
- Azure Data Factory copy and data flow-based lineage
- Resource deletion detection
- Resource Sets - A logical representation of partitions under common folder and file naming patterns in a data store.

User Experiences

- Admin and UI based experience for configuring, scanning, and classification of data
- Core catalog experience with ability to search, understand, and annotate assets
- Business glossary
- Glossary import from CSV
- Pre-built roles

Map & Discover - What's Coming Next

Scanning and Connectors

- Improvement to UI based scanning experience including on-premises SQL Server
- Power BI scanning and lineage support
- Resource Sets Improvements
- Targeted scanning and filtering at sub-resource levels, e.g., scan only one SQL table, scan whole SQL server and exclude one SQL table

Classification Rules

- Focus on accuracy and increase the number of PII-based classifications.

Platform

- ACLs
- Unification across policy, protect and insights

User Experiences

- Lineage view improvements
- Updated portal experience including browse for data
- Business rules to apply labels
- Data steward editing experience
- Insights reports for the data steward

Roadmap

Private Preview

- Private Preview Milestone 1
17 customers
Delivered Feb 2019

Public Preview

- Expect to announce in early 2020* for Map, Catalog and Insights

GA

- Date will be determined based on customer feedback
Not expected until late 2020

- Private Preview Expansion
400+ customers
Onboarding customers upon request.

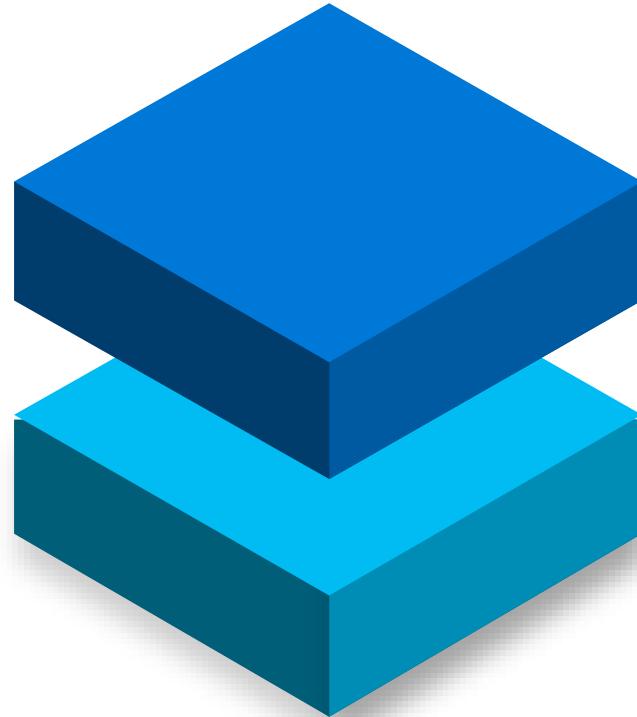
* This date is not publicly announced and subject to change.



Azure Databricks

Yatharth Gupta, Azure Databricks Program Manager

Azure Databricks



Fast, Easy, and Collaborative





Azure Databricks – Introduction

Fast, easy, and collaborative Apache Spark™-based analytics platform



Increase productivity



Build on a secure, trusted cloud



Scale without limits



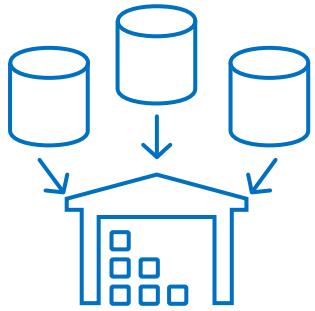
Built with your needs in mind

- Role-based access controls
- Effortless autoscaling
- Live collaboration
- Enterprise-grade SLAs
- Best-in-class notebooks
- Simple job scheduling



Seamlessly integrated with the Azure Portfolio

Our customers have three common objectives



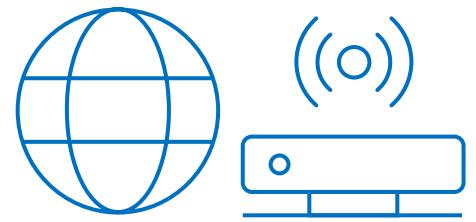
"We want to extend to un tapped sources"

Modern Data Warehouse



"We want to use ML and AI to get deeper insights from our data"

Advanced Analytics



"We want to get insights from our devices in real-time"

Real-time Analytics



Azure Databricks – Use Cases

Retail

CONSUMER ENGAGEMENT



Real-time Pricing Optimization

- Demand-Elasticity
- Personal Pricing Schemes
- Promotion events
- Multi-channel engagement

Healthcare

SENSOR DATA



IoT DEVICE ANALYTICS

- Aggregation of streaming events
- Predictive Maintenance
- Anomaly Detection

Financial

RISK AND REVENUE MANAGEMENT



Risk and Fraud, Threat Detection

- Real-time anomaly detection
- Card Monitoring and Fraud Detection
- Risk Aggregation

Advertising

RECOMMENDATION ENGINE



Next Best and Personalized Offers

- Right product, promotion, at right time
- Real time Ad bidding platform
- Personalized Ad Targeting

Oil/Gas & Energy

GRID OPS, ASSET OPTIMIZATION



Industrial IoT

- Preventive Maintenance
- Smart Grids and Microgrids
- Asset performance as a Service
- UAV image analysis

Media Entertainment

CONSUMER ENGAGEMENT ANALYSIS



Sentiment Analysis

- Demand-Elasticity
- Social Network Analysis
- Promotion events
- Multi-channel Attribution

Security

ACTIONABLE THREAT INTELLIGENCE



Security Intelligence

- Real-time firewall, network, and auth log correlation
- Anomaly detection
- Security context, enrichment
- Security Orchestration

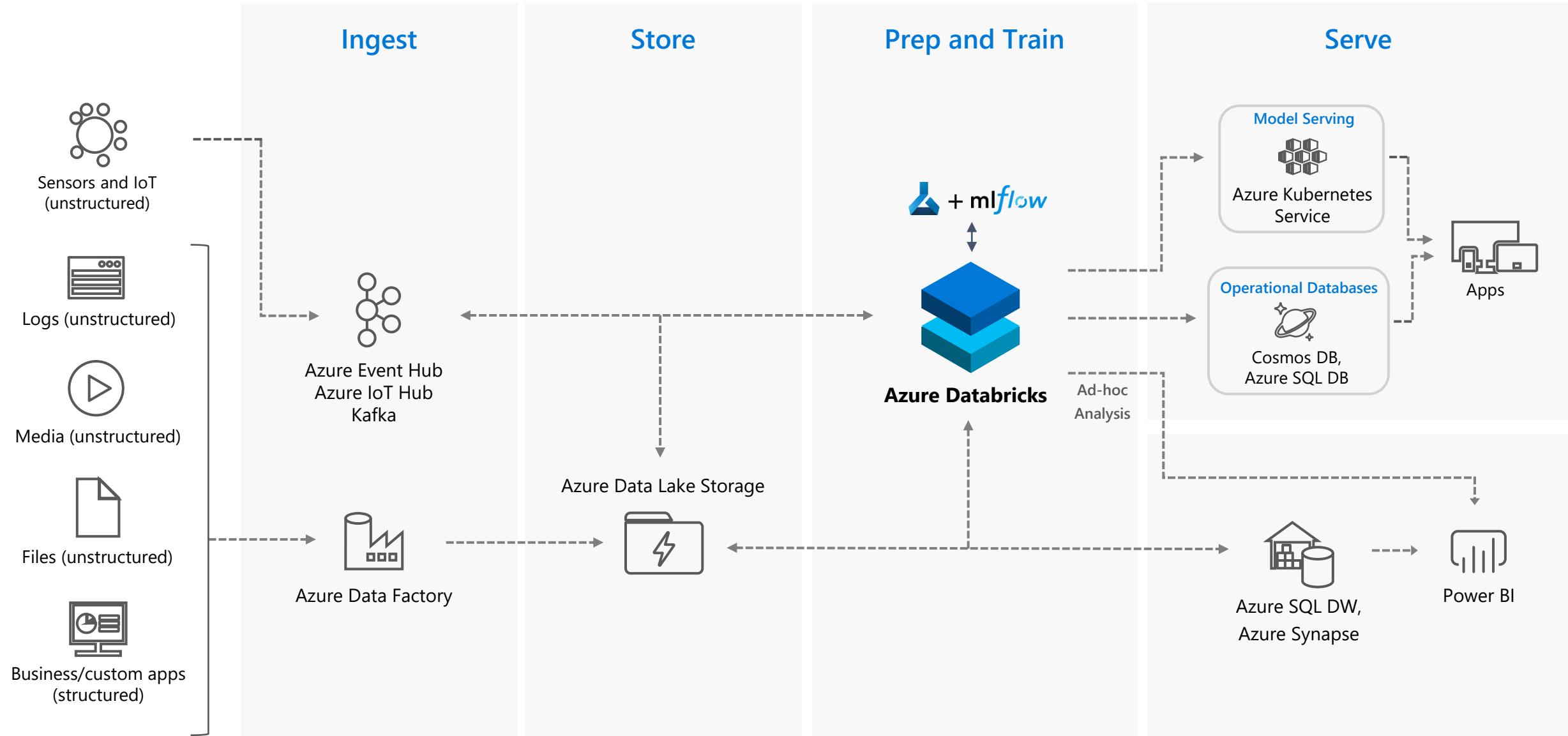
And Much More!



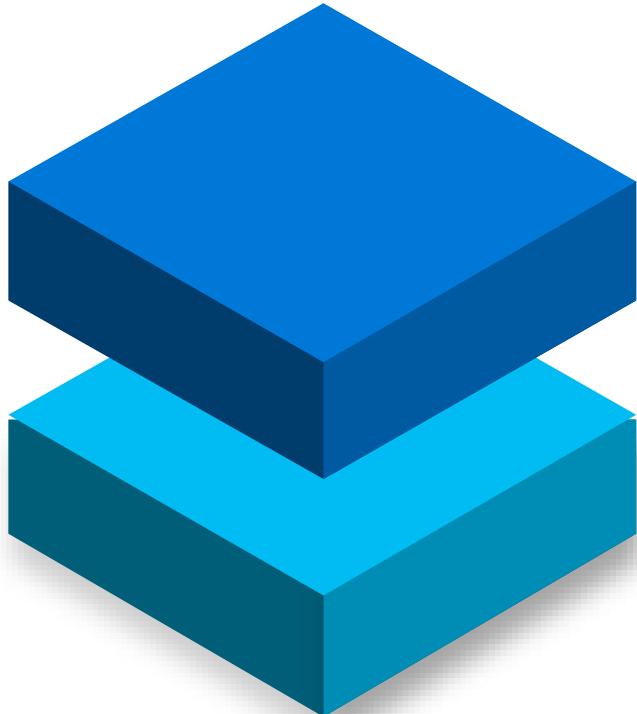
Azure Databricks – How does it fit in?



Azure Databricks – Architecture Example

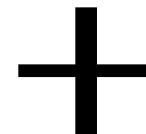


Azure Databricks



Fast, Easy, and Collaborative

DELTA LAKE



Bring Reliability to your Data Lake

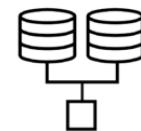
DELTA LAKE - Reliable Data Lakes at Scale on Azure



Data Versioning



ACID Transactions



Optimized Layouts



Fast Streaming



Efficient Upserts



Schema Enforcement

ACID Transaction Guarantees

- Atomic, Consistent, Isolated, Durable

Schema Enforcement

- Control schema evolution for reliable Data Lakes

Efficient Upserts

- *MERGE, DELETE, UPDATE*

Small file compaction w/ no interrupt to availability

- *OPTIMIZE and VACUUM*

Z-Order partitioning w/ up to 100x perf

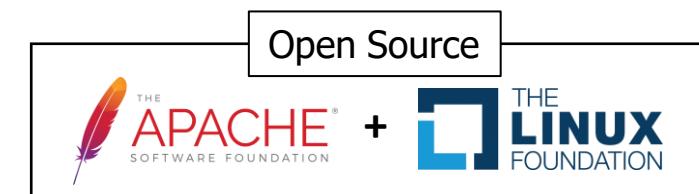
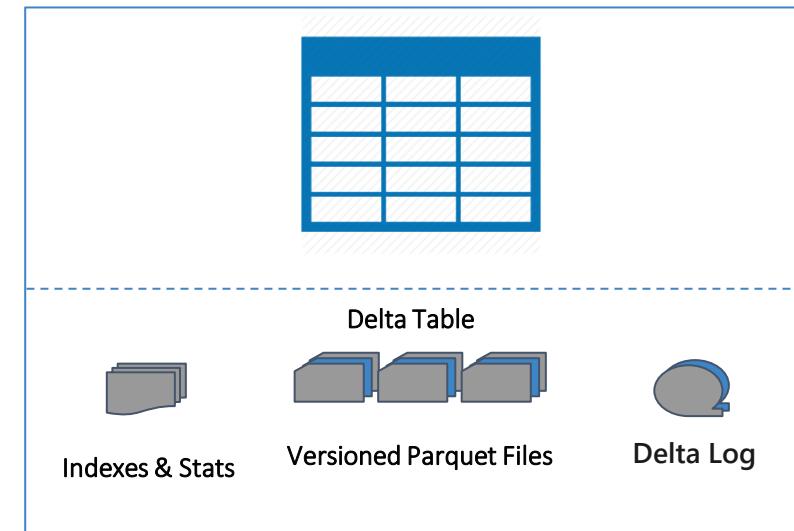
- New multidimensional partitioning enables data skipping

Time Travel

- Audit history, Pipeline Debugging, Data Reproducibility

Delta Table =

Parquet + Transaction Log + Indexes/Stats





DELTA LAKE - Easy to use

BEFORE

```
CREATE TABLE ...  
USING parquet  
...
```

or

```
dataframe  
.write  
.format("parquet")  
.save("/data")
```

AFTER

```
CREATE TABLE ...  
USING delta  
...
```

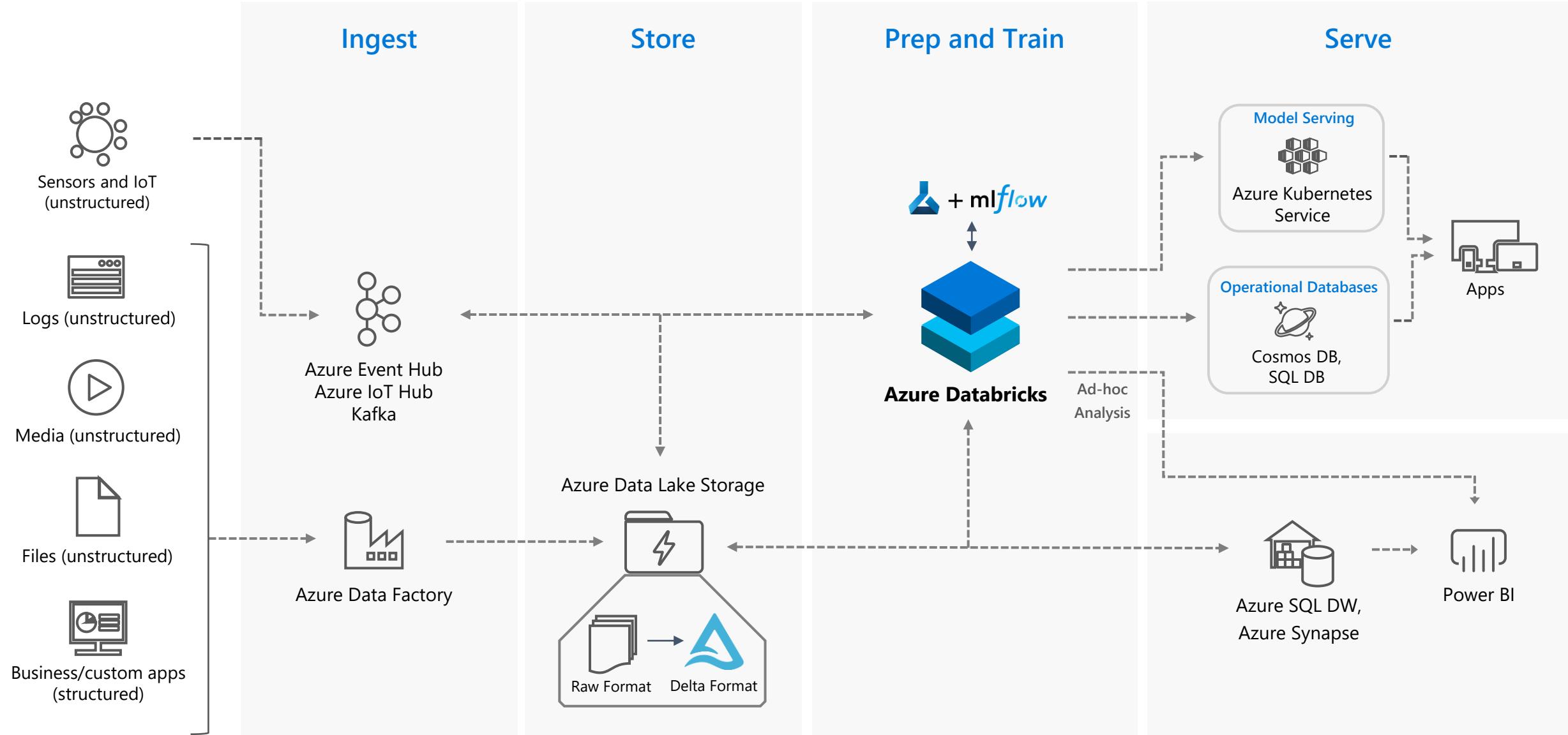
or

```
dataframe  
.write  
.format("delta")  
.save("/data")
```



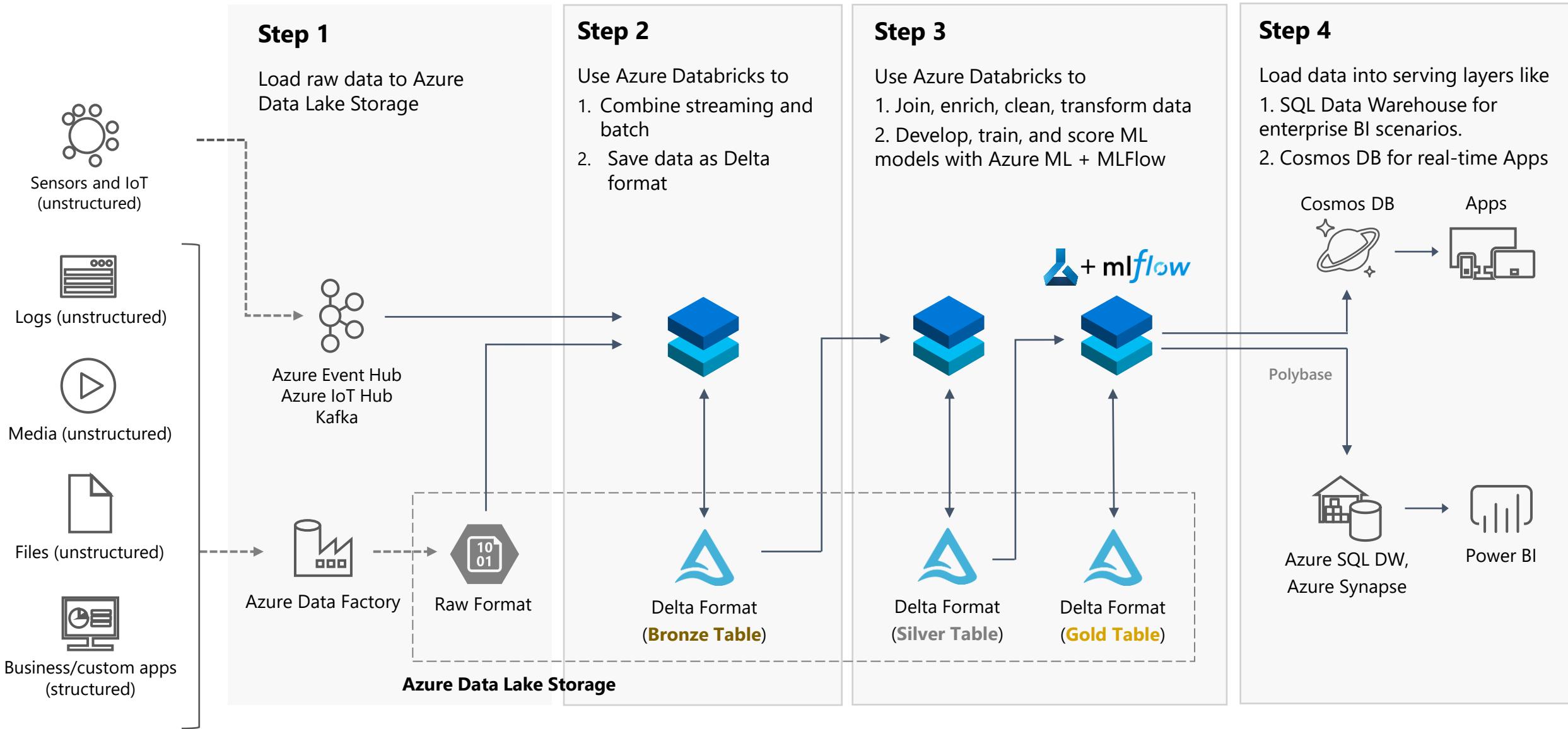


Azure Databricks – Architecture Example





Azure Databricks – Delta Lake at Scale on Azure





DELTA LAKE – Perf Tuning





Azure Databricks – Delta, Performance Tuning Best Practices

Partitioning

Partition on low cardinality columns. Try to aim for ~1GB per partition

File Compaction - *OPTIMIZE*

Run your *OPTIMIZE* command in a separate cluster (F or Fsv2 Series recommended)

ZORDER

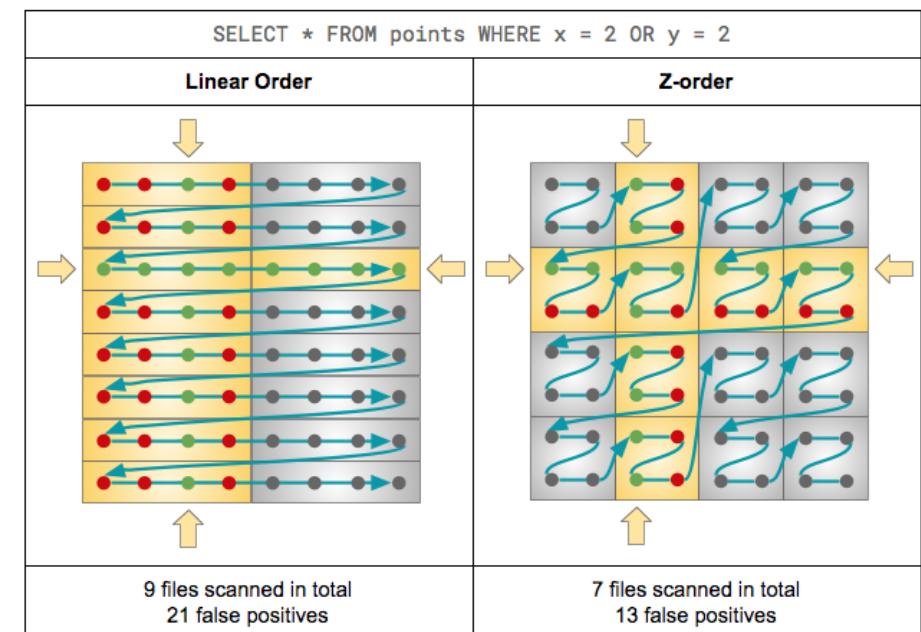
ZORDER on HIGH cardinality columns. This enables data skipping.

Z-Ordering is a [technique](#) to colocate related information in the same set of files. This co-locality is automatically used by Delta data-skipping algorithms to dramatically reduce the amount of data that needs to be read. To Z-Order data, you specify the columns to order on in the `ZORDER BY` clause:

```
OPTIMIZE events  
WHERE date >= current_timestamp() - INTERVAL 1 day  
ZORDER BY (eventType)
```

Copy

You can specify multiple columns for `ZORDER BY` as a comma-separated list. However, the effectiveness of the locality drops with each additional column.





DELTA LAKE – Time Travel



DELTA LAKE - Time Travel



```
SELECT count(*) FROM my_table TIMESTAMP AS OF "2019-01-01"  
SELECT count(*) FROM my_table TIMESTAMP AS OF date_sub(current_date(), 1)  
SELECT count(*) FROM my_table TIMESTAMP AS OF "2019-01-01 01:30:00.000"
```

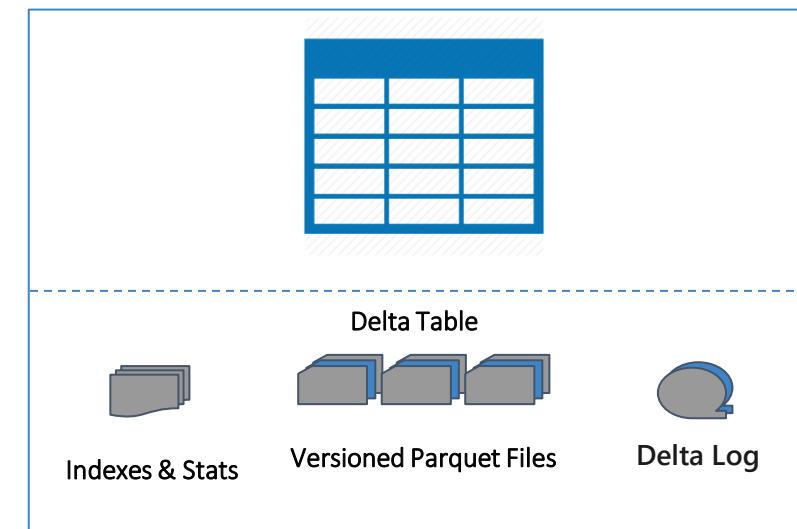
```
SELECT count(*) FROM my_table VERSION AS OF 5238  
SELECT count(*) FROM my_table@v5238  
SELECT count(*) FROM delta.`/path/to/my/table@v5238`
```

Applications Include:

1. Audit History
2. Data Reproducibility
3. Data Pipeline Debugging
4. Immediate Rollbacks

Delta Table =

Parquet + Transaction Log + Indexes/Stats



DELTA LAKE - Time Travel, Audit History

Audit Data Changes

- History of all operations are recorded for audit history
- Audit operation types, userIds, clusterIds, notebookIds, timestamps and versions

DESCRIBE HISTORY

Provides provenance information, including the operation, user, and so on, for each write to a table. This information is not recorded by versions of Databricks Runtime below 4.1 and tables created using these versions will show this information as `null`. Table history is retained for 30 days.

```
display(spark.sql("DESCRIBE HISTORY events"))
```

▶ (2) Spark Jobs

version	timestamp	userId	userName	operation	operationParameters	job	notebook	clusterId	readVersion	isolationLevel
2	2019-01-29 00:38:19			OPTIMIZE	▶ {"predicate": "[]", "zOrderBy": "[]", "batchId": "0"}	null	▶ {"notebookId": "2433269420249641"}		1	SnapshotIsolation
1	2019-01-29 00:38:10			WRITE	▶ {"mode": "Append", "partitionBy": "["date\"]"}	null	▶ {"notebookId": "2433269420249641"}		0	WriteSerializable
0	2019-01-29 00:37:58			WRITE	▶ {"mode": "Overwrite", "partitionBy": "["date\"]"}	null	▶ {"notebookId": "2433269420249641"}		null	WriteSerializable

DELTA LAKE - Time Travel, Data Reproducibility

Data reproducibility

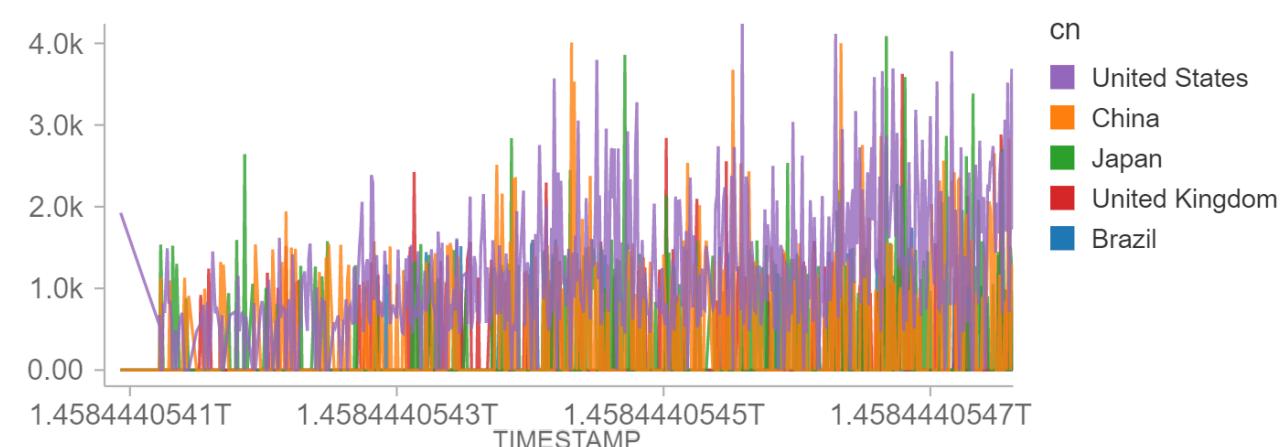
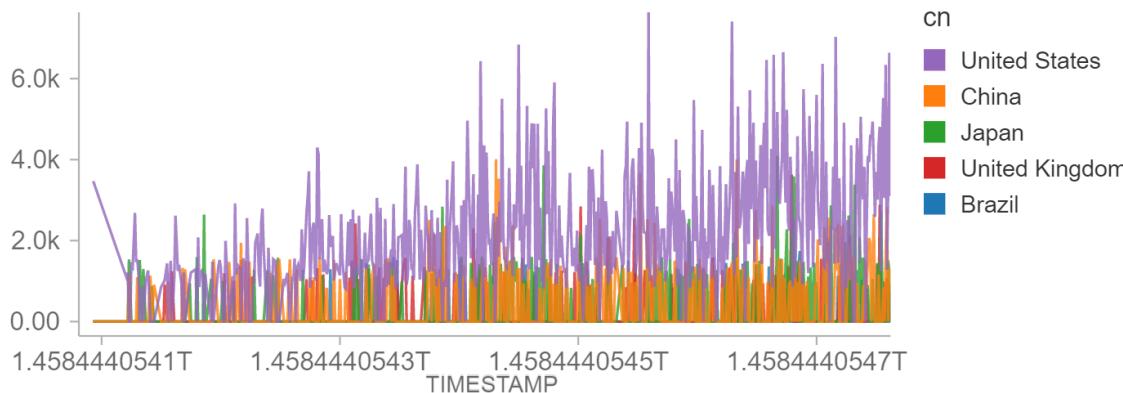
Reproduce query results and reports

- Go back to the exact same data that was used to train an ML model version in the past.

```
SELECT count(*) FROM my_table AS OF "2019-01-01"
```

```
1 %sql
2 SELECT c02_level, deviceRaw.cn, timestamp
3 FROM deviceRaw
4 WHERE cn IN ("United States", "United Kingdom", "China", "Japan", "Brazil")
```

```
1 %sql
2 SELECT * FROM deviceRaw TIMESTAMP AS OF '2019-10-31T05:11:54.000+0000'
3 WHERE cn IN ("United States", "United Kingdom", "China", "Japan", "Brazil")
```





DELTA LAKE - Time Travel, Rollbacks

Rollbacks

Time travel also makes it easy to do rollbacks in case of bad writes. For example, if your GDPR pipeline job had a bug that accidentally deleted user information, you can easily fix the pipeline:

```
INSERT INTO my_table
SELECT * FROM my_table TIMESTAMP AS OF date_sub(current_date(), 1)
WHERE userId = 111
```

You can also fix incorrect updates as follows:

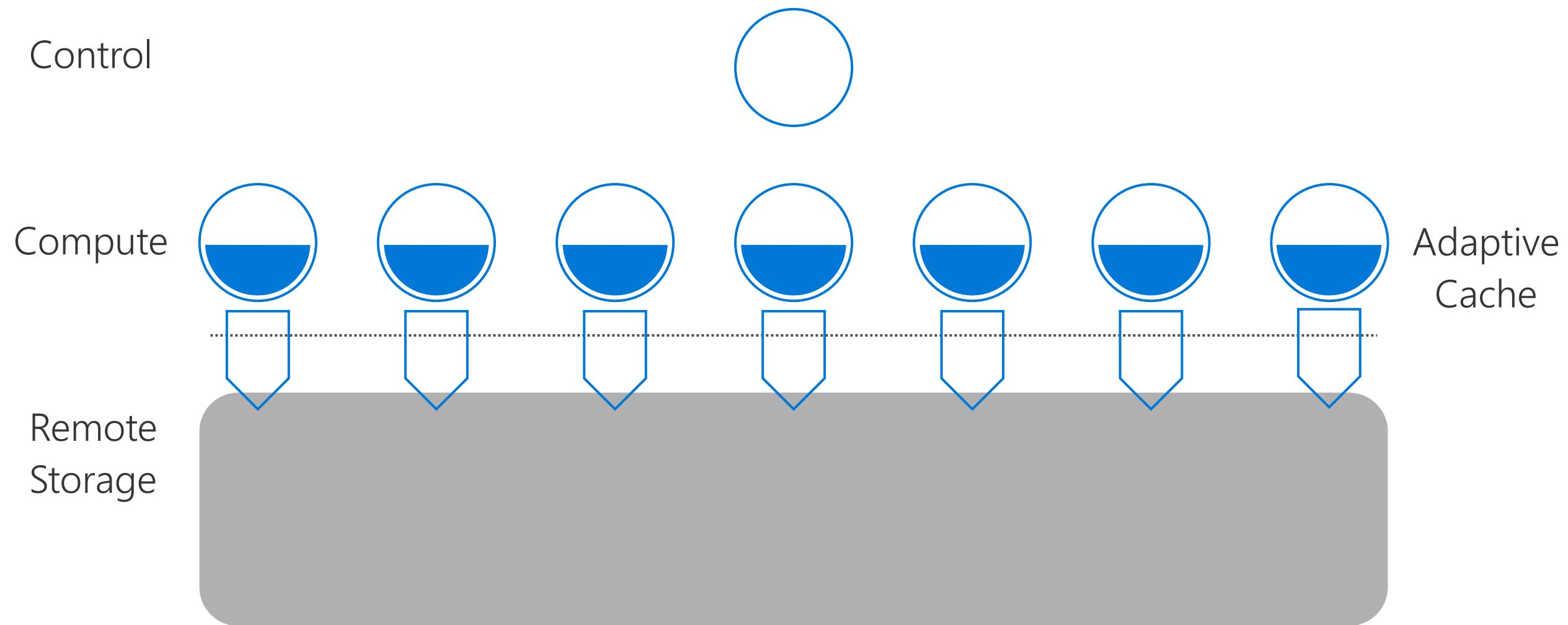
```
MERGE INTO my_table target
USING my_table TIMESTAMP AS OF date_sub(current_date(), 1) source
ON source.userId = target.userId
WHEN MATCHED THEN UPDATE SET *
```



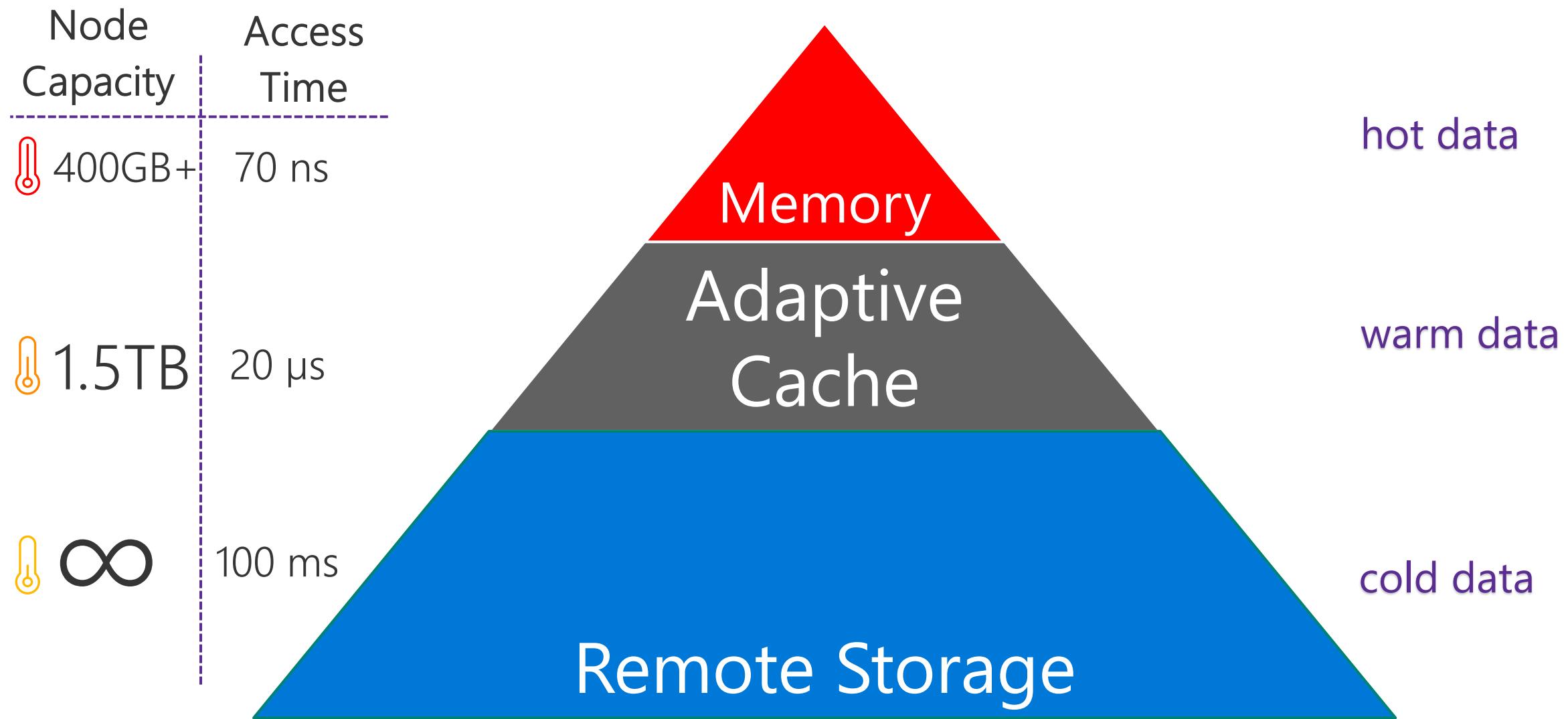
Microsoft

Appendix and Supporting Materials

Provisioned SQL Architecture



Automated Tiering Of Storage Layers



Tables – Indexes

Clustered Columnstore index (Default Primary)

Highest level of data compression

Best overall query performance

Support for ordered Columnstore segments

Clustered index (Primary)

Performant for looking up a single to few rows

Heap (Primary)

Faster loading and landing temporary data

Best for small lookup tables

Nonclustered indexes (Secondary)

Enable ordering of multiple columns in a table

Allows multiple nonclustered on a single table

Can be created on any of the above primary indexes

More performant lookup queries

-- Create table with index

```
CREATE TABLE orderTable
```

```
(
```

```
    OrderId INT NOT NULL,
```

```
    Date DATE NOT NULL,
```

```
    Name VARCHAR(2),
```

```
    Country VARCHAR(2)
```

```
)
```

```
WITH
```

```
(
```

```
    CLUSTERED COLUMNSTORE INDEX |
```

```
    HEAP |
```

```
    CLUSTERED INDEX (OrderId)
```

```
);
```

-- Add non-clustered index to table

```
CREATE INDEX NameIndex ON orderTable (Name);
```

Tables – DW Indexes Illustrated

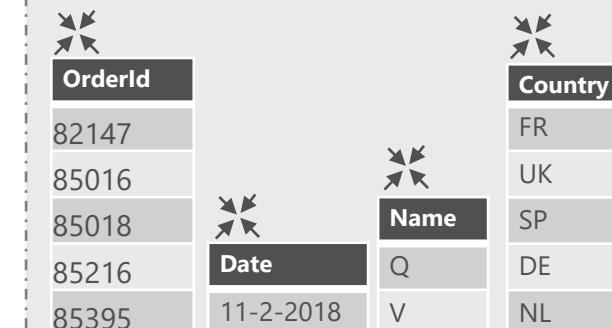
Logical table structure

OrderId	Date	Name	Country
85016	11-2-2018	V	UK
85018	11-2-2018	Q	SP
85216	11-2-2018	Q	DE
85395	11-2-2018	V	NL
82147	11-2-2018	Q	FR
86881	11-2-2018	D	UK
93080	11-3-2018	R	UK
94156	11-3-2018	S	FR
96250	11-3-2018	Q	NL
98799	11-3-2018	R	NL
98015	11-3-2018	T	UK
98310	11-3-2018	D	DE
98979	11-3-2018	Z	DE
98137	11-3-2018	T	FR
...

Clustered columnstore index

(OrderId)

Rowgroup1
Min (OrderId): 82147 | Max (OrderId): 85395



Delta Rowstore

OrderId	Date	Name	Country
98137	11-3-2018	T	FR
98310	11-3-2018	D	DE
98799	11-3-2018	R	NL
98979	11-3-2018	Z	DE

- Data stored in compressed columnstore segments after being sliced into groups of rows (rowgroups/micro-partitions) for maximum compression
- Rows are stored in the delta rowstore until the number of rows is large enough to be compressed into a columnstore

Clustered/Non-clustered rowstore index

(OrderId)

OrderId	PagId
82147	1001
98137	1002

OrderId	PagId
82147	1005
85395	1006

OrderId	Date	Name	Country
82147	11-2-2018	Q	FR
85016	11-2-2018	V	UK
85018	11-2-2018	Q	SP
85395	11-3-2018	R	NL

OrderId	PagId
98137	1007
98979	1008

OrderId	Date	Name	Country
98137	11-3-2018	T	FR
98310	11-3-2018	D	DE
98799	11-3-2018	R	NL
98979	11-3-2018	Z	DE

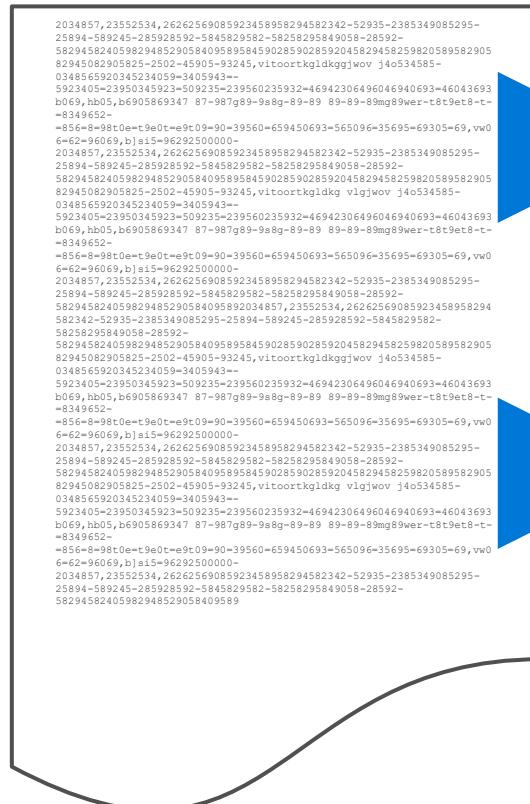
- Data is stored in a B-tree index structure for performant lookup queries for particular rows.
- Clustered rowstore index: The leaf nodes in the structure store the data values in a row (as pictured above)
- Non-clustered (secondary) rowstore index: The leaf nodes store pointers to the data values, not the values themselves

Column store taxonomy

Data

Row Group

Segments Column store



Tables – Distributions

Round-robin distributed

Distributes table rows evenly across all distributions at random.

Hash distributed

Distributes table rows across the Compute nodes by using a deterministic hash function to assign each row to one distribution.

Replicated

Full copy of table accessible on each Compute node.

```
CREATE TABLE dbo.OrderTable
(
    OrderId INT NOT NULL,
    Date    NOT NULL,
    Name    VARCHAR(2),
    Country VARCHAR(2)
)
WITH
(
    CLUSTERED COLUMNSTORE INDEX,
    DISTRIBUTION = HASH([OrderId]) |
                    ROUND ROBIN |
                    REPLICATED
);
```

Tables – Partitions

Overview

Table partitions divide data into smaller groups

In most cases, partitions are created on a date column

Supported on all table types

RANGE RIGHT – Used for time partitions

RANGE LEFT – Used for number partitions

Benefits

Improves efficiency and performance of loading and querying by limiting the scope to subset of data.

Offers significant query performance enhancements where filtering on the partition key can eliminate unnecessary scans and eliminate IO.

```
CREATE TABLE partitionedOrderTable
(
    OrderId INT NOT NULL,
    Date DATE NOT NULL,
    Name VARCHAR(2),
    Country VARCHAR(2)
)
WITH
(
    CLUSTERED COLUMNSTORE INDEX,
    DISTRIBUTION = HASH([OrderId]),
    PARTITION (
        [Date] RANGE RIGHT FOR VALUES (
            '2000-01-01', '2001-01-01', '2002-01-01',
            '2003-01-01', '2004-01-01', '2005-01-01'
        )
    )
);
```

Logical table structure

OrderId	Date	Name	Country
85016	11-2-2018	V	UK
85018	11-2-2018	Q	SP
85216	11-2-2018	Q	DE
85395	11-2-2018	V	NL
82147	11-2-2018	Q	FR
86881	11-2-2018	D	UK
93080	11-3-2018	R	UK
94156	11-3-2018	S	FR
96250	11-3-2018	Q	NL
98799	11-3-2018	R	NL
98015	11-3-2018	T	UK
98310	11-3-2018	D	DE
98979	11-3-2018	Z	DE
98137	11-3-2018	T	FR
...

Physical data distribution

(Hash distribution (OrderId), Date partitions)

Distribution1

(OrderId 80,000 – 100,000)

11-2-2018 partition

OrderId	Date	Name	Country
85016	11-2-2018	V	UK
85018	11-2-2018	Q	SP
85216	11-2-2018	Q	DE
85395	11-2-2018	V	NL
82147	11-2-2018	Q	FR
86881	11-2-2018	D	UK
...

11-3-2018 partition

OrderId	Date	Name	Country
93080	11-3-2018	R	UK
94156	11-3-2018	S	FR
96250	11-3-2018	Q	NL
98799	11-3-2018	R	NL
98015	11-3-2018	T	UK
98310	11-3-2018	D	DE
98979	11-3-2018	Z	DE
98137	11-3-2018	T	FR
...

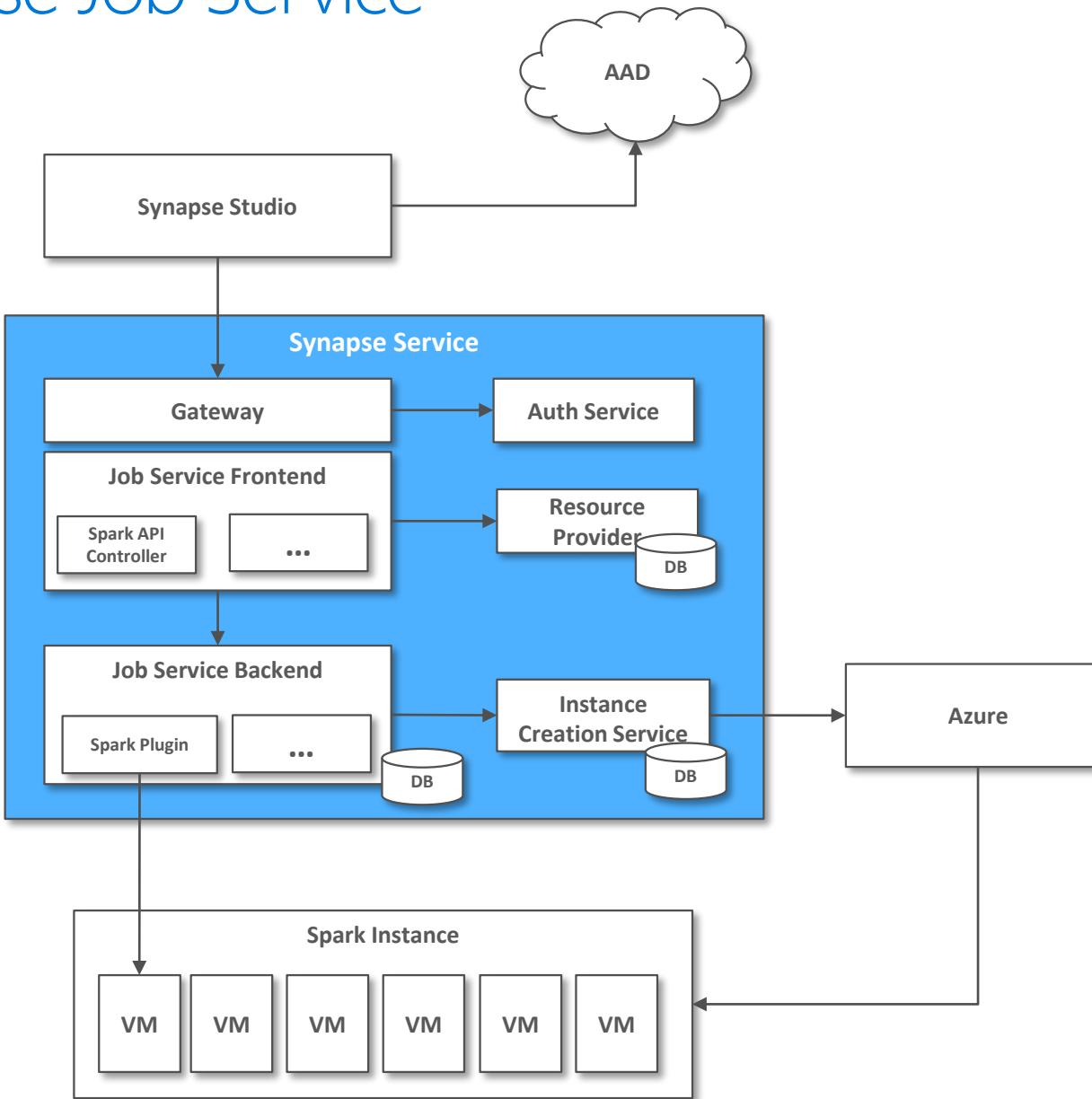
Tables – DW
Distributions &
Partitions Illustrated

...

x 60 distributions (shards)

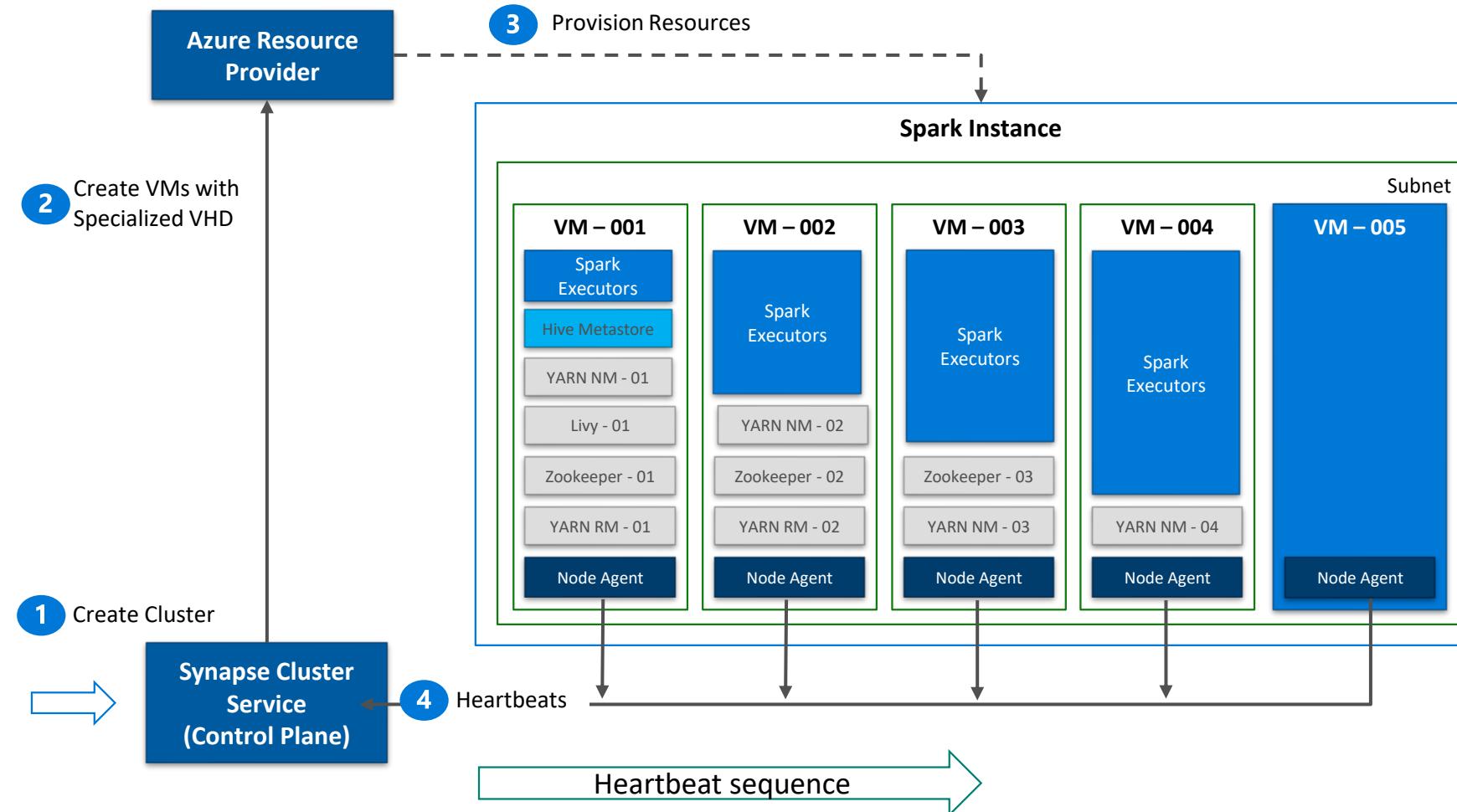
- Each shard is partitioned with the same date partitions
- A minimum of 1 million rows per distribution and partition is needed for optimal compression and performance of clustered Columnstore tables

Synapse Job Service



- User creates Synapse Workspace and Spark pool and launches Synapse Studio.
- User attaches Notebook to Spark pool and enters one or more Spark statements (code blocks).
- The Notebook client gets user token from AAD and sends a Spark session create request to Synapse Gateway.
- Synapse Gateway authenticates the request and validates authorizations on the Workspace and Spark pool and forwards it to the Spark (Livy) controller hosted in Synapse Job Service frontend.
- The Job Service frontend forwards the request to Job Service backend that creates two jobs – one for creating the cluster and the other for creating the Spark session.
- The Job service backend contacts Synapse Resource Provider to obtain Workspace and Spark pool details and delegates the cluster creation request to Synapse Instance Service.
- Once the instance is created, the Job Service backend forwards the Spark session creation request to the Livy endpoint in the cluster.
- Once the Spark session is created the Notebook client sends Spark statements to the Job Service frontend.
- Job Service frontend obtains the actual Livy endpoint for the cluster created for the particular user from the backend and sends the statement directly to Livy for execution.

Synapse Spark Instances



1. Synapse Job Service sends request to Cluster Service for creating BBC clusters per the description in the associated Spark pool.
2. Cluster Service sends request to Azure using Azure SDK to create VMs (required plus additional) with specialized VHD.
3. The specialized VHD contains bits for all the services that are required by the Cluster type (for e.g. Spark) with prefetch instrumentation.
4. Once VM boots up, the Node Agent sends heartbeat to Cluster Service for getting node configuration.
5. The nodes are initialized and assigned roles based on their first heartbeat.
6. Extra nodes get deleted on first heartbeat.
7. After Cluster Service considers the cluster ready, it returns the Livy endpoint to the Job Service.

Creating a Spark pool (1 of 2)

Provision Spark Pool through Azure Portal with default settings or per requirements

Basic Settings – Minimum details required from user

Home > Synapse workspaces > euang-synapse-nov-ws - Apache Spark pools > Create Apache Spark pool

Create Apache Spark pool

Basics * Additional settings * Tags Summary

Create a Synapse Analytics Apache Spark pool with your preferred configurations. Complete the Basics tab then go to Review + create to provision with smart defaults, or visit each tab to customize.

Apache Spark pool details

Name your Apache Spark pool and choose its initial settings.

Apache Spark pool name *

Node size family

Node size *

Autoscale * ⓘ

Number of nodes *

Only required field from user

Default Settings

Enter Apache Spark pool name

MemoryOptimized

Medium (8 vCPU / 64 GB)

Enabled Disabled

3 40

Creating a Spark pool (2 of 2) - optional

Additional Settings offer optional settings to customize Spark pool

Customize component versions, auto-pause

Import libraries by providing text file containing library name and version

Home > Synapse workspaces > euang-synapse-nov-ws - Apache Spark pools > Create Apache Spark pool

Create Apache Spark pool

Basics * Additional settings * Tags Summary

Customize additional configuration parameters including autoscale and component versions.

Auto-pause

Enter required settings for this Apache Spark pool, including setting auto-pause and picking versions.

Auto-pause * ⓘ Enabled Disabled

Number of minutes idle * 15

Component versions

Select the Apache Spark version for your Apache Spark pool.

Apache Spark *	2.4
Python	3.6.1
Scala	2.11.12
Java	1.8.0_222
.NET Core	3.0
.NET for Apache Spark	0.6.0
Delta Lake	0.4.0

Packages

Upload environment configuration file ("PIP freeze" output).

File upload Select a file Upload

Library Management - Python

Overview

Customers can add new python libraries at Spark pool level

Benefits

Input requirements.txt in simple pip freeze format

Add new libraries to your cluster

Update versions of existing libraries on your cluster

Libraries will get installed for your Spark pool during cluster creation

Ability to specify different requirements file for different pools within the same workspace

Constraints

The library version must exist on PyPI repository

Version downgrade of an existing library not allowed

In the Portal

Specify the new requirements while creating Spark Pool in Additional Settings blade

Microsoft Azure (Preview) Restore default configuration Report a bug Search resources, services, and data

Home > nushuklasynapsewestus2 > Create Apache Spark pool

Create Apache Spark pool

Enter required settings for this Apache Spark pool, including setting auto-pause and picking versions.

Auto-pause * Enabled Disabled

Number of minutes idle *

Component versions

Select the Apache Spark version for your Apache Spark pool.

Component	Version
Apache Spark *	2.4
Python	3.6.1
Scala	2.11.12
Java	1.8.0_222
.NET Core	3.0
.NET for Apache Spark	0.6.0
Delta Lake	0.4.0

Packages

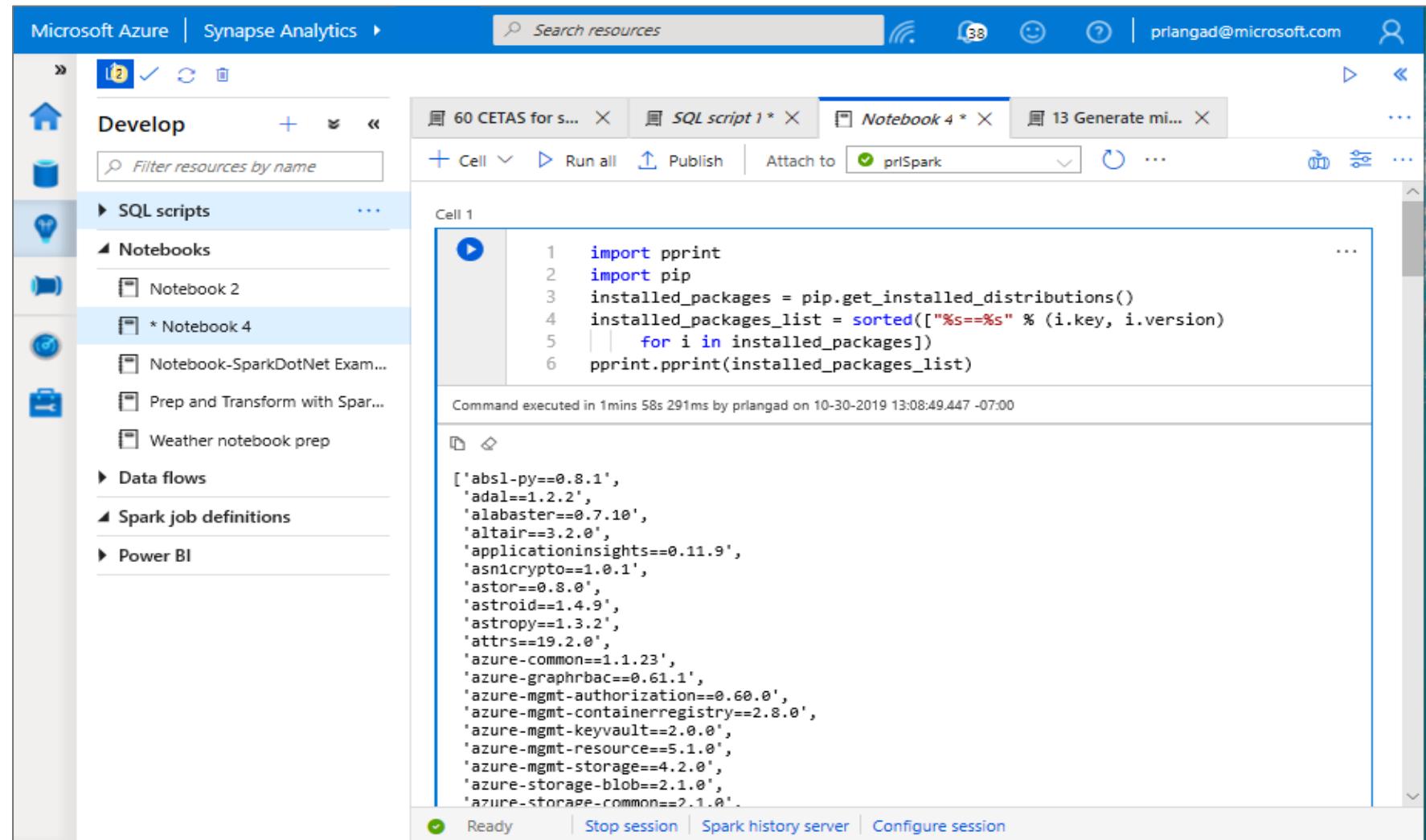
Upload environment configuration file ("PIP freeze" output).

File upload Upload

Review + create < Previous Next: Tags >

Library Management - Python

Get list of installed libraries with version information



The screenshot shows the Microsoft Azure Synapse Analytics workspace interface. On the left, the 'Develop' sidebar is open, showing a list of resources: SQL scripts, Notebooks, Data flows, Spark job definitions, and Power BI. Under 'Notebooks', 'Notebook 4' is selected. In the main area, there are four tabs at the top: '60 CETAS for s...', 'SQL script 1 * X', 'Notebook 4 * X' (which is active), and '13 Generate mi... X'. Below the tabs, a toolbar includes 'Cell', 'Run all', 'Publish', 'Attach to', and session controls. A session named 'priSpark' is attached. A code cell labeled 'Cell 1' contains the following Python script:

```
1 import pprint
2 import pip
3 installed_packages = pip.get_installed_distributions()
4 installed_packages_list = sorted([ "%s==%s" % (i.key, i.version)
5         for i in installed_packages])
6 pprint.pprint(installed_packages_list)
```

The output of the cell shows a long list of installed Python packages and their versions. Some examples from the list include:

- 'absl-py==0.8.1'
- 'adal==1.2.2'
- 'alabaster==0.7.10'
- 'altair==3.2.0'
- 'applicationinsights==0.11.9'
- 'asn1crypto==1.0.1'
- 'astor==0.8.0'
- 'astroid==1.4.9'
- 'astropy==1.3.2'
- 'attrs==19.2.0'
- 'azure-common==1.1.23'
- 'azure-graphrbac==0.61.1'
- 'azure-mgmt-authorization==0.60.0'
- 'azure-mgmt-containerregistry==2.8.0'
- 'azure-mgmt-keyvault==2.0.0'
- 'azure-mgmt-resource==5.1.0'
- 'azure-mgmt-storage==4.2.0'
- 'azure-storage-blob==2.1.0'
- 'azure-storage-common==2.1.0'

At the bottom of the cell, it says 'Command executed in 1mins 58s 291ms by priLangad on 10-30-2019 13:08:49.447 -07:00'. The status bar at the bottom of the workspace shows 'Ready'.

Microsoft Machine Learning for Apache Spark

v1.0-rc

Microsoft's Open Source
Contributions to Apache Spark



Distributed
Machine Learning



Fast Model
Deployment



Microservice
Orchestration



Multilingual Binding
Generation

www.aka.ms/spark

 [Azure/mmlspark](https://github.com/Azure/mmlspark)



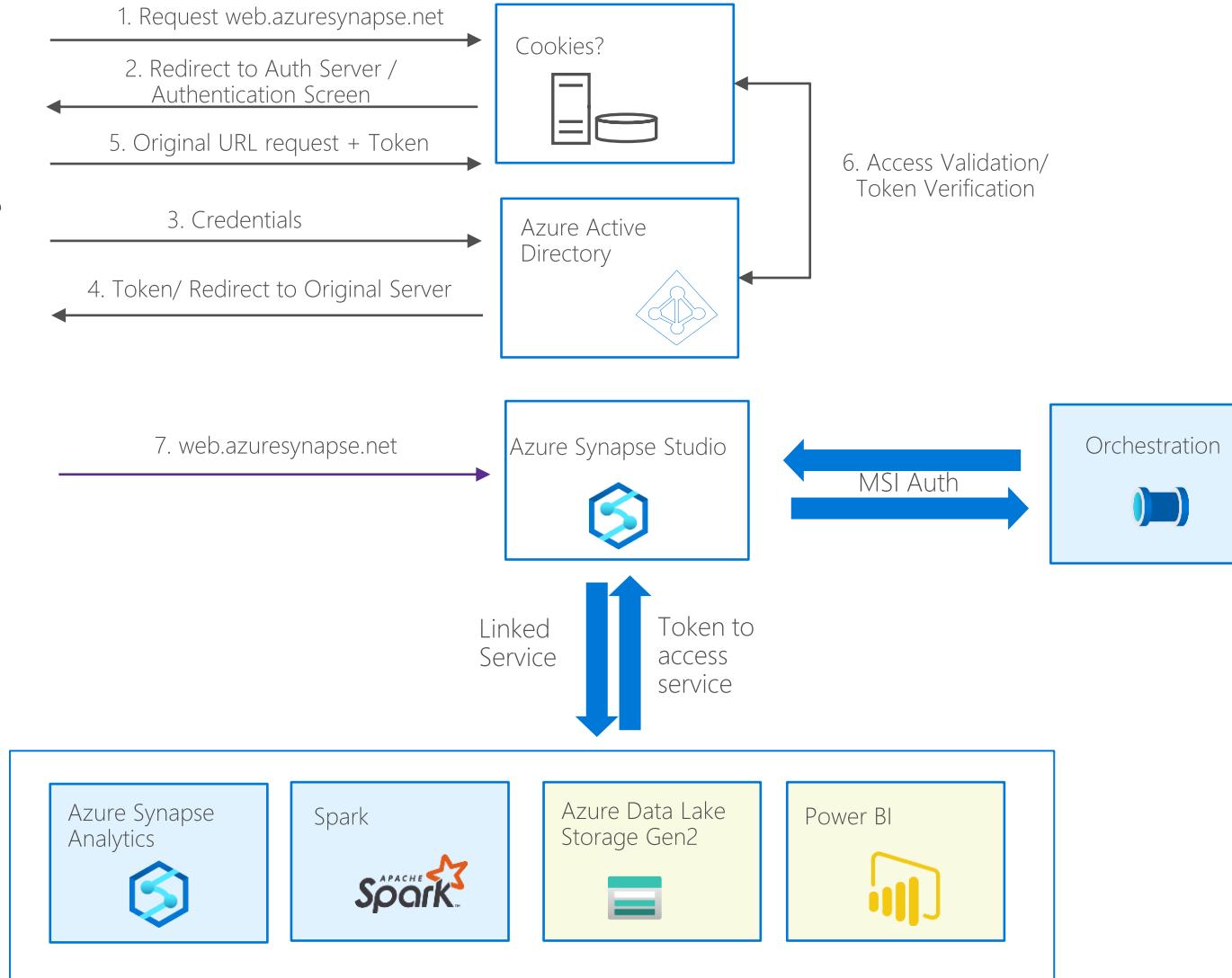
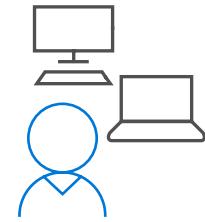
Azure Synapse Analytics Security

Enterprise-grade security



Single Sign-On

Synapse Foundation Components
 Synapse Linked Services



Industry-leading compliance



ISO 27001



SOC 1 Type 2



SOC 2 Type 2



PCI DSS Level 1



Cloud Controls Matrix



ISO 27018



Content Delivery and Security Association



Shared Assessments



FedRAMP JAB P-ATO



HIPAA / HITECH



FIPS 140-2



21 CFR Part 11



FERPA



DISA Level 2



CJIS



IRS 1075 / ITAR-ready



European Union Model Clauses



EU Safe Harbor



United Kingdom G-Cloud



China Multi Layer Protection Scheme



China GB 18030



China CCCPPF



Singapore MTCS Level 3



Australian Signals Directorate



New Zealand GCIO



Japan Financial Services



ENISA IAF



SQL On-Demand

SQL On-Demand

Overview

An interactive query service that provides T-SQL queries over high scale data in Azure Storage.

Benefits

Serverless

No infrastructure

Pay only for query execution

No ETL

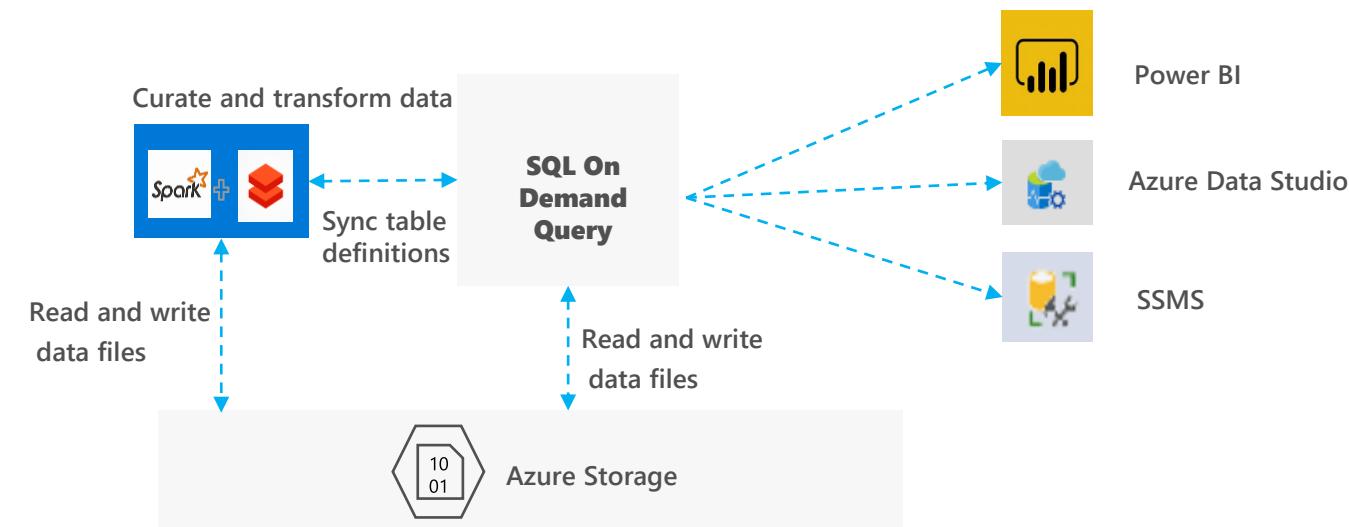
Offers security

Data integration with Databricks, HDInsight

T-SQL syntax to query data

Supports data in various formats (Parquet, CSV, JSON)

Support for BI ecosystem



SQL On-Demand – Querying on storage

Microsoft Azure | Synapse Analytics > prlangadws2

Search resources

Publish all 2 Validate all Refresh Discard all

Data HolidayDataPi... Load Data to S... Pipeline 1 Data flow 1 SQL script 1 Copy Open Da... nytlc

Upload Download New Folder Select All Rename Manage Access Properties Delete Refresh

Storage accounts prlangaddemosa (Primary)

- filesystem
- holidaydatacontainer
- isdweatherdatacontainer
- nytlc**
- prlangaddemosa
- tmpcontainer
- wwimporters

New SQL script and open in new tab

part-00133-120938564719836543-aea5b543-5e83-4a7d-8d31-69f72c50b05d-15253-1.c000.snappy.parquet

New SQL script New notebook Copy ABFS path Manage Access... Rename... Download Delete Properties...

LAST MODIFIED 10/25/2019, 2:20:23 PM

Microsoft Azure | Synapse Analytics > prlangadws2

Search resources

Publish all 3 Validate all Refresh Discard all

Data HolidayDataPi... Load Data to S... Pipeline 1 Data flow 1 SQL script 1 Copy Open Da... nytlc SQL script 2

Run Publish Query plan Connect to SQL Analytics on-demand Use database master

```

1 SELECT
2   TOP 100 *
3   FROM
4     OPENROWSET(
5       BULK 'https://prlangaddemosa.dfs.core.windows.net/nytlc/yellow/puYear=2015/puMonth=3/part-00133-tid-210938564719836543-aea5b543-5e83-4a7d-8d31-69f72c50b05d-15253-1.c000.snappy.parquet'
6       FORMAT='PARQUET'
7     ) AS nyc;
8

```

Results Messages

View Table Chart Export results

VENDORID	TPEPICKUPDATETIME	TPEPDROPOFFDATETIME	PASSENGERCOUNT	TRIPDISTANCE	PULOCATIONID	DOLOCATIONID	STARTLON	STARTLAT	ENDLON	ENDLAT
2	2015-02-28T23:5...	2015-03-01T00:0...	6	1.63	NULL	NULL	-74.000846862793	40.7306938171387	-73.	
1	2015-03-28T19:2...	2015-03-28T19:2...	1	2.2	NULL	NULL	-73.977653503418	40.7631607055664	-73.	
2	2015-02-28T23:5...	2015-03-01T00:1...	5	3.23	NULL	NULL	-73.96012878417...	40.7621574401855	-73.	
1	2015-03-28T19:2...	2015-03-28T19:3...	1	2.1	NULL	NULL	-73.98143005371...	40.7815055847168	-74.	
2	2015-02-28T23:5...	2015-03-01T00:1...	1	3.52	NULL	NULL	-73.98373413085...	40.7497062683105	-74.	
?	2015-02-28T00:0...	2015-02-28T00:0...	5	...	NULL	NULL	-73.9814077707	40.748055847168	-73.	

00:01:00 Query executed successfully.

SQL On-Demand – Querying Spark Tables

- SQL On-Demand shares metastore with Spark
- Spark databases automatically created in SQL On-Demand

Spark Code

```
new_rows = [('CA',22, 45000),("WA",35,65000),  
 ,("WA",50,85000)]  
  
demo_df = spark.createDataFrame(new_rows, ['state', 'age',  
 'salary'])  
  
demo_df.write.saveAsTable('demo_df')
```

SQL OD Code

```
Select * From demo_df
```

state	age	salary
CA	22	45000
WA	35	65000
WA	50	85000

SQL On-Demand – Creating views

Overview

Create views using SQL On-Demand queries

Benefits

Works same as standard views

```
USE [mydbname]
GO

IF EXISTS(select * FROM sys.views where name = 'populationView')
DROP VIEW populationView
GO

CREATE VIEW populationView AS
SELECT *
FROM OPENROWSET(
    BULK 'https://XXX.blob.core.windows.net/csv/population/population.csv',
    FORMAT = 'CSV',
    FIELDTERMINATOR = ',',
    ROWTERMINATOR = '\n'
)
WITH (
    [country_code] VARCHAR (5) COLLATE Latin1_General_BIN2,
    [country_name] VARCHAR (100) COLLATE Latin1_General_BIN2,
    [year] smallint,
    [population] bigint
) AS [r]
```

```
SELECT
    country_name, population
FROM populationView
WHERE
    [year] = 2019
ORDER BY
    [population] DESC
```

	country_name	population
1	China	1389618778
2	India	1311559204
3	United States	331883986
4	Indonesia	264935824
5	Pakistan	210797836
6	Brazil	210301591
7	Nigeria	208679114
8	Bangladesh	161062905
9	Russia	141944641
10	Mexico	127318112

SQL On-Demand – Creating views

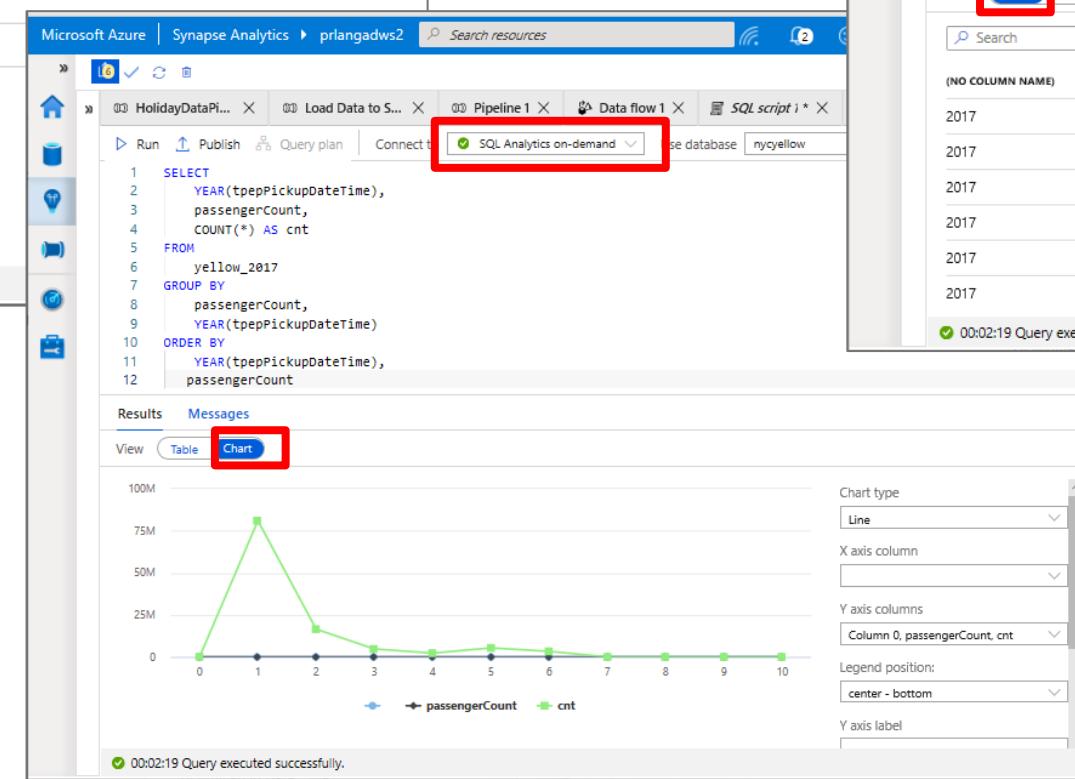
```

-- type your sql script here, we now have intellisense
CREATE VIEW yellow_2017 AS
SELECT *
FROM
OPENROWSET(
    BULK 'https://prlangaddemosa.dfs.core.windows.net/nyctlc/yellow/puYear=2017/*/*',
    FORMAT='PARQUET'
) AS nyc

```

Results Messages

00:00:17 Query executed successfully.



(NO COLUMN NAME)	PASSENGERCOUNT	CNT
2017	0	166086
2017	1	81034075
2017	2	16545571
2017	3	4748869
2017	4	2257813
2017	5	5407319

00:02:19 Query executed successfully.

Create External Table As Select

Overview

Creates an external table and then exports results of the Select statement. These operations will import data into the database for the duration of the query

Steps:

1. Create Master Key
2. Create Credentials
3. Create External Data Source
4. Create External Data Format
5. Create External Table

```
-- Create a database master key if one does not already exist
CREATE MASTER KEY ENCRYPTION BY PASSWORD = 'S0me!Info'
;

-- Create a database scoped credential with Azure storage account key as the secret.
CREATE DATABASE SCOPED CREDENTIAL AzureStorageCredential
WITH
    IDENTITY = '<my_account>'
, SECRET   = '<azure_storage_account_key>'
;
-- Create an external data source with CREDENTIAL option.
CREATE EXTERNAL DATA SOURCE MyAzureStorage
WITH
(
    LOCATION  = 'wasbs://daily@logs.blob.core.windows.net/'
, CREDENTIAL = AzureStorageCredential
, TYPE      = HADOOP
)
-- Create an external file format
CREATE EXTERNAL FILE FORMAT MyAzureCSVFormat
WITH (FORMAT_TYPE = DELIMITEDTEXT,
      FORMAT_OPTIONS(
          FIELD_TERMINATOR = ',',
          FIRST_ROW = 2))
--Create an external table
CREATE EXTERNAL TABLE dbo.FactInternetSalesNew
WITH(
    LOCATION = '/files/Customer',
    DATA_SOURCE = MyAzureStorage,
    FILE_FORMAT = MyAzureCSVFormat
)
AS SELECT T1.* FROM dbo.FactInternetSales T1 JOIN dbo.DimCustomer T2
ON ( T1.CustomerKey = T2.CustomerKey )
OPTION ( HASH JOIN );
```



Azure Synapse Spark

Azure Synapse Apache Spark - Summary



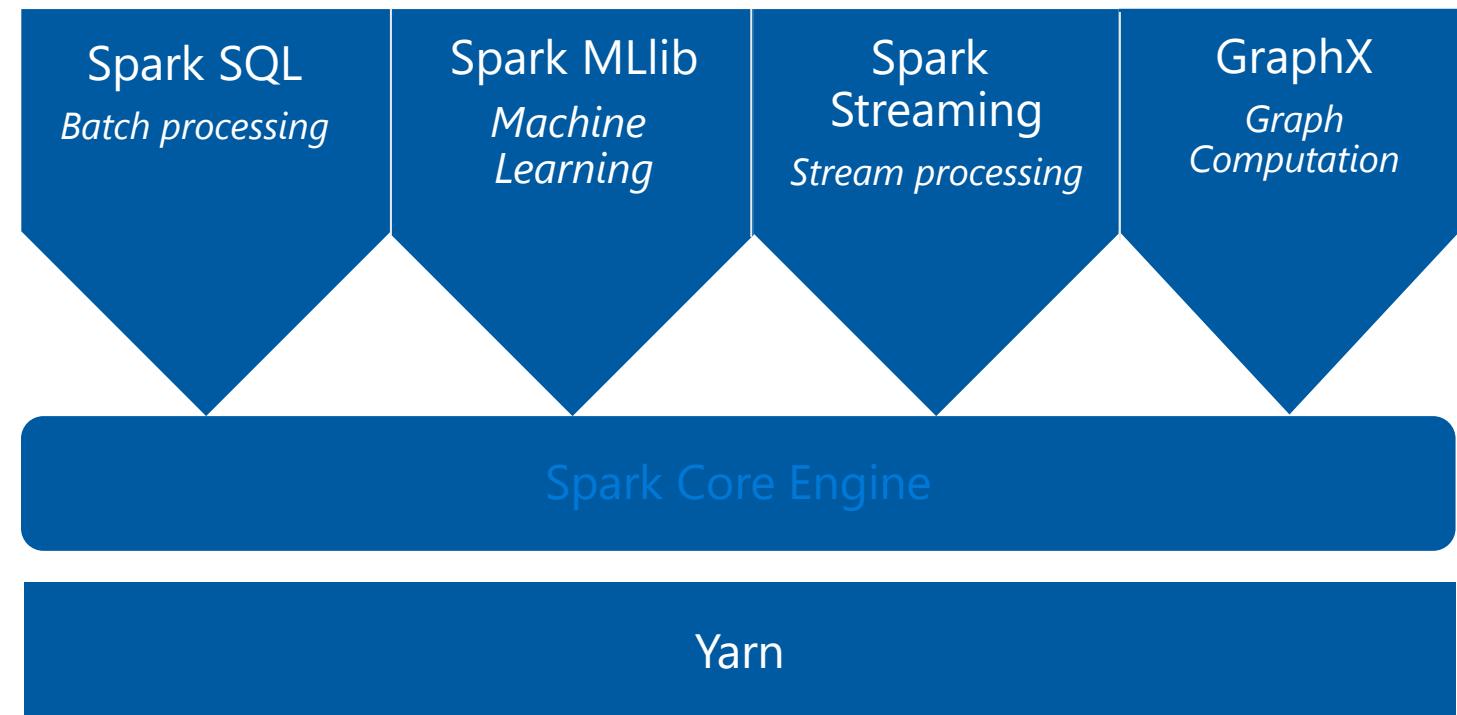
- Apache Spark 2.4 derivation
 - Linux Foundation Delta Lake 0.4 support
 - .Net Core 3.0 support
 - Python 3.6 + Anacondas support
- Tightly coupled to other Azure Synapse services
 - Integrated security and sign on
 - Integrated Metadata
 - Integrated and simplified provisioning
 - Integrated UX including Jupyter based notebooks
 - Fast load of SQL Analytics pools
- Core scenarios
 - Data Prep/Data Engineering/ETL
 - Machine Learning via Spark ML and Azure ML integration
 - Extensible through library management
- Efficient resource utilization
 - Fast Start
 - Auto scale (up and down)
 - Auto pause
 - Min cluster size of 3 nodes
- Multi Language Support
 - .Net (C#), PySpark, Scala, Spark SQL, Java

Apache Spark

A unified, open source, parallel, data processing framework for Big Data Analytics

Spark Unifies:

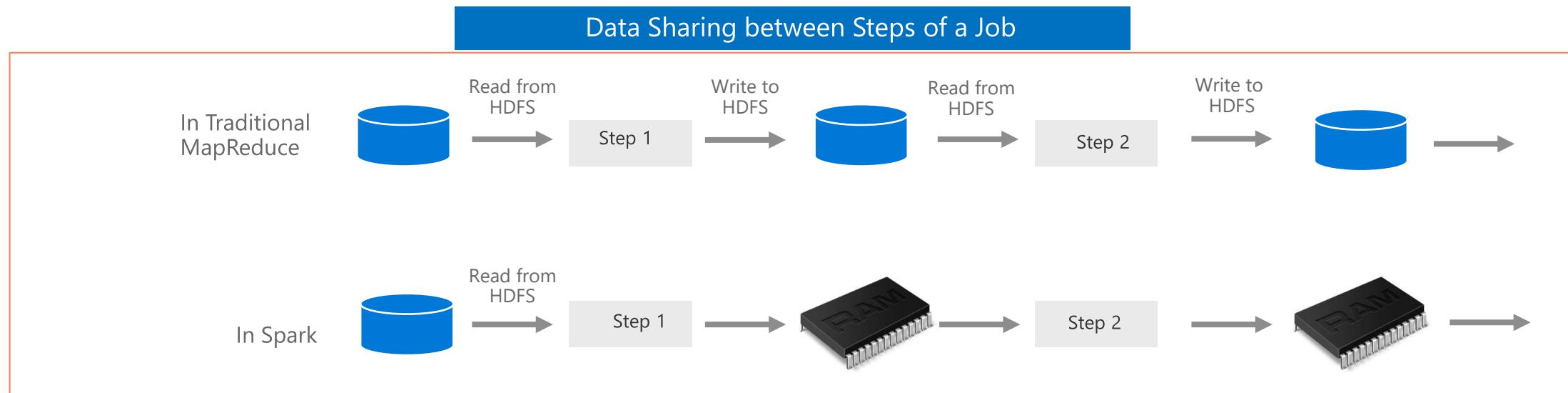
- Batch Processing
- Interactive SQL
- Real-time processing
- Machine Learning
- Deep Learning
- Graph Processing



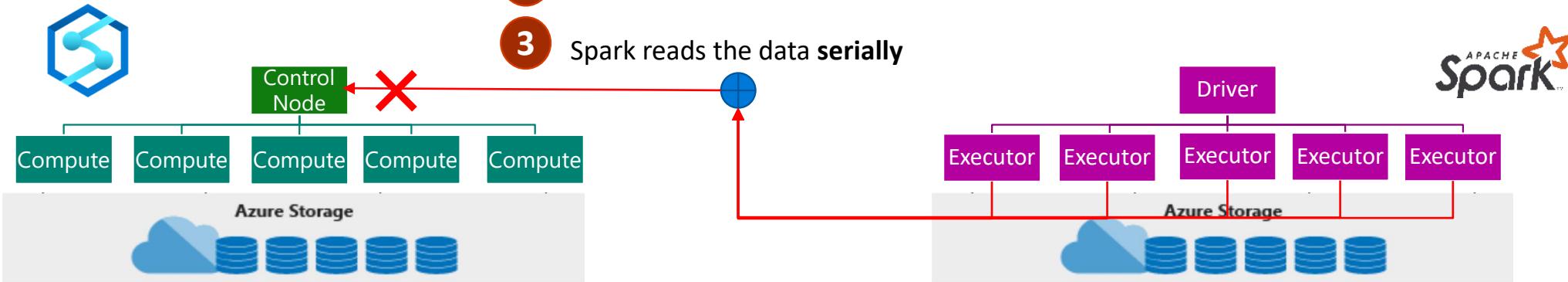
<http://spark.apache.org>

What makes Spark fast?

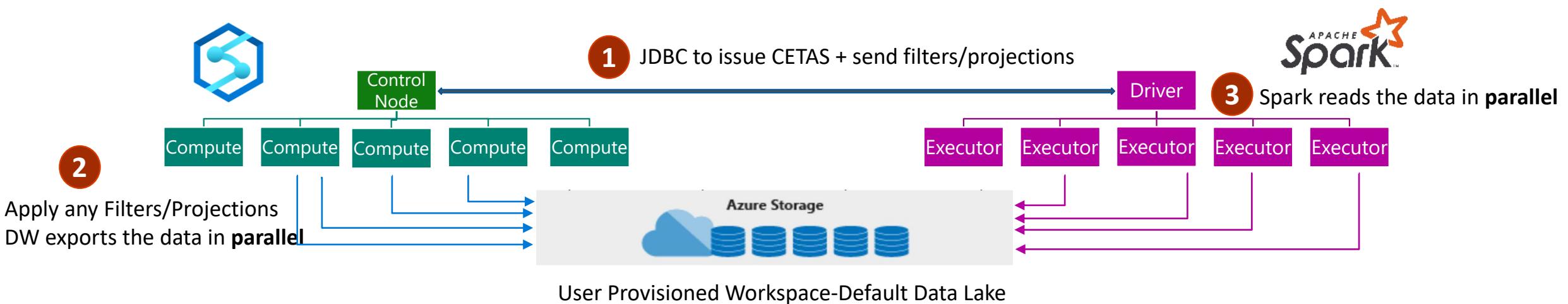
- **In-memory cluster computing:** Spark provides primitives for *in-memory* cluster computing. A Spark job can *load and cache* data into memory and query it repeatedly (iteratively) much quicker than disk-based systems.
- **Scala Integration:** Spark integrates into the Scala programming language, letting you manipulate distributed datasets like local collections. No need to structure everything as map and reduce operations
- **Faster Data-sharing:** Data-sharing between operations is faster as data is in-memory:
 - In (traditional) Hadoop data is shared through HDFS which is expensive. HDFS maintains three replicas.
 - Spark stores data in-memory *without any replication*.



Existing Approach: JDBC



New Approach: JDBC and Polybase



Code-Behind Experience

Existing Approach

```
val jdbcUsername = "<SQL DB ADMIN USER>"  
val jdbcPwd = "<SQL DB ADMIN PWD>"  
val jdbcHostname = "servername.database.windows.net"  
val jdbcPort = 1433  
val jdbcDatabase = "<AZURE SQL DB NAME>"  
  
val jdbc_url =  
  s"jdbc:sqlserver://${jdbcHostname}:${jdbcPort};database=${jdbcDatabase};"  
  encrypt=true;trustServerCertificate=false;hostNameInCertificate=*.databas  
e.windows.net;loginTimeout=60;"  
  
val connectionProperties = new Properties()  
  
connectionProperties.put("user", s"${jdbcUsername}")  
connectionProperties.put("password", s"${jdbcPwd}")  
  
val sqlTableDf = spark.read.jdbc(jdbc_url, "dbo.Tbl1", connectionProperties)
```

New Approach

```
// Construct a Spark DataFrame from SQL Pool  
var df = spark.read.sqlAnalytics("sql1.dbo.Tbl1")  
  
// Write the Spark DataFrame into SQL Pool  
df.write.sqlAnalytics("sql1.dbo.Tbl2")
```

HDInsight Appendix

AutoScale for Apache Spark, Hive, HBase, LLAP

& MR

Automatically scale cluster size up and down:

1 Load-based: define min and max

2 Schedule-based: create custom schedule
e.g. 25 nodes @ 9 AM & 3 nodes @ 10 PM

SETUP AUTOSCALE

1

Basics Storage Security + networking Configuration + pricing Review + create

Configure cluster performance and pricing. [Learn more](#)

Node configuration

Configure your cluster's size and performance, and view estimated cost information.

The cost estimate represented in the table does not include subscription discounts or costs related to storage, networking, or data transfer.

This configuration will use 40 to 88 of 1440 available cores in the East US region. [View cores usage](#)

Node type	Node size	Number of nodes	estimated cost/hour
Head node	D12 v2 (4 Cores, 28 GB RAM), 0.37 USD...	2	0.75 USD
Worker node	D13 v2 (8 Cores, 56 GB RAM), 0.75 USD...	4	

Enable autoscale (preview) [Learn more](#)

Autoscale type Load-based Schedule-based **Min** 4 **Max** 10 2.99 to 7.48 USD

Load-based autoscale will scale the number of worker nodes used based on the cluster's activity.

2

Autoscale configuration

Configure the schedule-based autoscale settings for your cluster.

Time zone [\(UTC-08:00\) Pacific Time \(US & Canada\)](#)

Conditions

Based on the conditions below, your cluster will scale to the targeted number of worker nodes.

Each condition applies to one or more days of the week, and multiple conditions cannot share the same day of the week.

+Add condition

Days	Time	Number of nodes
MON,TUE,WED,THURS,FRI	09:00	20
	18:00	6
	22:00	3

contosocluster1234 - Cluster size

HDInsight cluster

Search (Ctrl+I) Save Revert changes Feedback

The cost estimate represented in the table does not include subscription discounts or costs related to storage, networking, or data transfer.

This configuration will use 40 of 64 available cores in the East US region. [View core usage](#)

NODE TYPE	NODE SIZE	# OF NODES	ESTIMATED COST/HOUR
Head node	D12 V2 (4 cores, 28 GB RAM) - 0.37 USD/hour	2	0.74 USD
Worker node	D3 V2 (4 cores, 14 GB RAM) - 0.30 USD/hour	2	0.60 USD

Enable autoscale

Autoscale will allow the number of worker nodes used to adjust based on the cluster's activity.
Minimum # of nodes 2 Maximum # of nodes 6 Estimated cost/hour 0.60 to 1.80 USD [Configure autoscale settings](#)

Cluster size history

100 80 60 40 20 0

7/2 7/3 7/4 7/5 7/6 7/7 7/8 7/9 7/10 7/11 7/12 7/13 7/14 7/15 7/16 7/17 7/18 7/19 7/20 7/21 7/22 7/23 7/24 7/25 7/26 7/27 7/28 7/29 7/30 7/31 7/32 7/33 7/34 7/35 7/36 7/37 7/38 7/39 7/40 7/41 7/42 7/43 7/44 7/45 7/46 7/47 7/48 7/49 7/50 7/51 7/52 7/53 7/54 7/55 7/56 7/57 7/58 7/59 7/60 7/61 7/62 7/63 7/64 7/65 7/66 7/67 7/68 7/69 7/70 7/71 7/72 7/73 7/74 7/75 7/76 7/77 7/78 7/79 7/80 7/81 7/82 7/83 7/84 7/85 7/86 7/87 7/88 7/89 7/90 7/91 7/92 7/93 7/94 7/95 7/96 7/97 7/98 7/99 7/100

Worker nodes 20 Min # of nodes 10 Max # of nodes 80 View in Azure Metrics

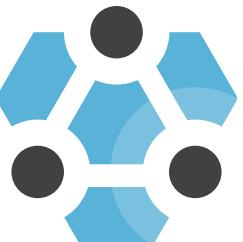
3

4

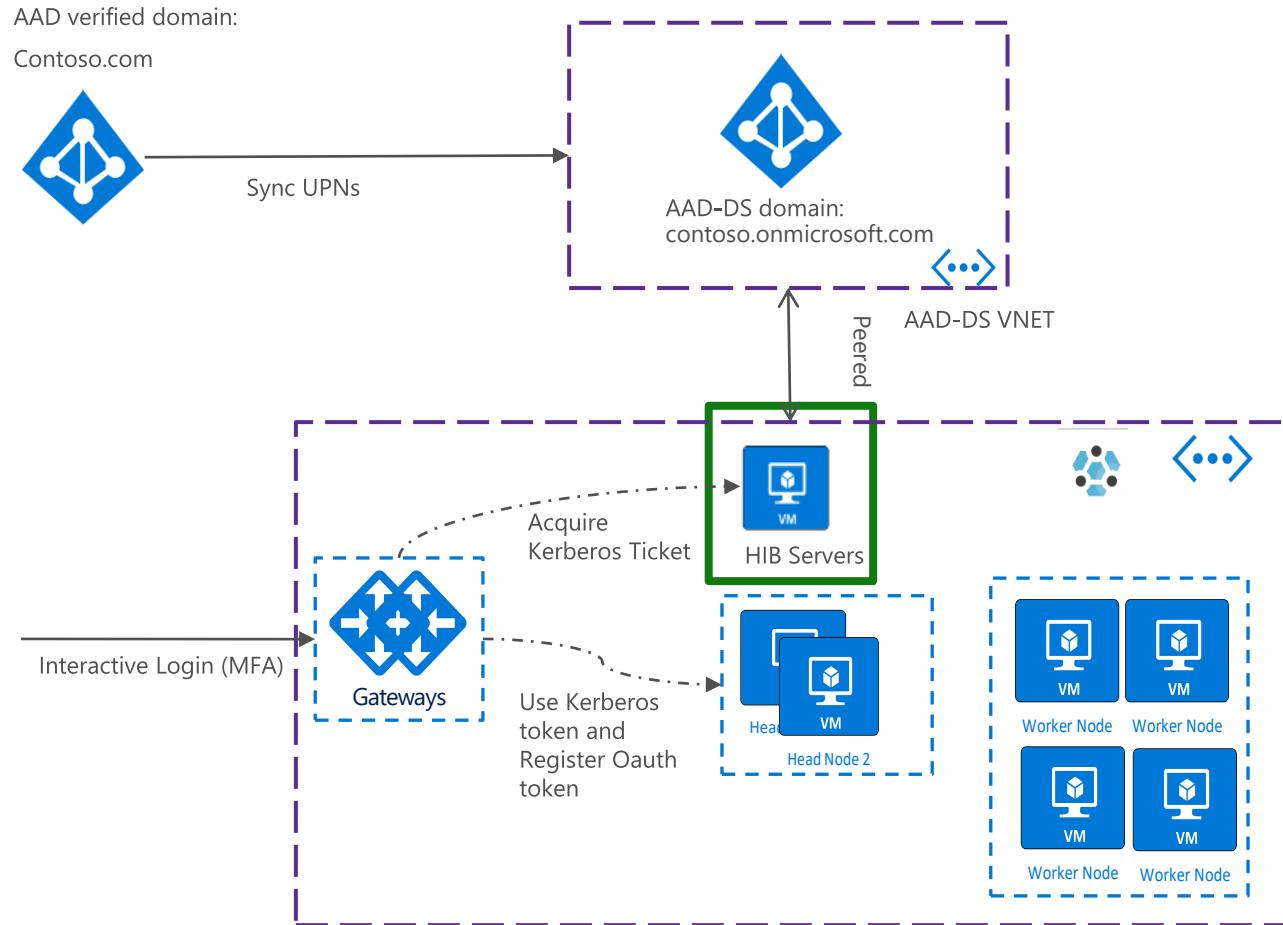
4

Pay for only what you need
Easily Monitor scaling history

MONITOR AUTOSCALE



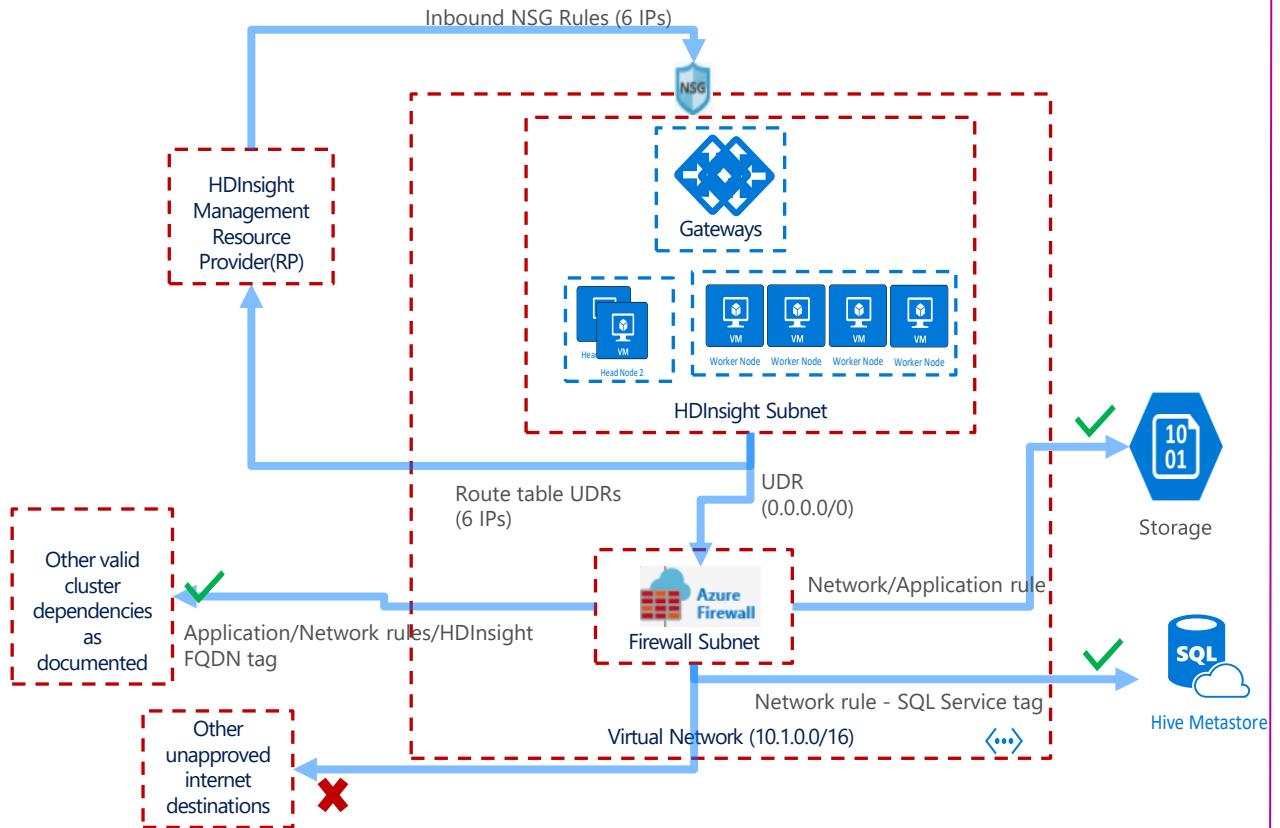
HDInsight ID Broker (HIB) – Preview



- Login to cluster using **OAuth and MFA**
- **Single Sign-On** experience
- **No Password Hash Sync** to AAD required for gateway access.
- Tools login using MFA (**IntelliJ HDInsight plugin**)



Restricting Outbound Traffic with Azure Firewall (GA)



Prevent large data transfers to unauthorized destinations

Simple **HDInsight FQDN tag** in Azure Firewall

See **denied access logs** in Azure Monitor

The screenshot shows the 'Add application rule collection' interface in the Azure Firewall portal. The 'FQDN tags' section is highlighted, showing the 'HDInsight' tag selected.

Left sidebar: Microsoft.AzureFirewall-20191025165441 - Overview > Firewall > Rules > Add application rule collection

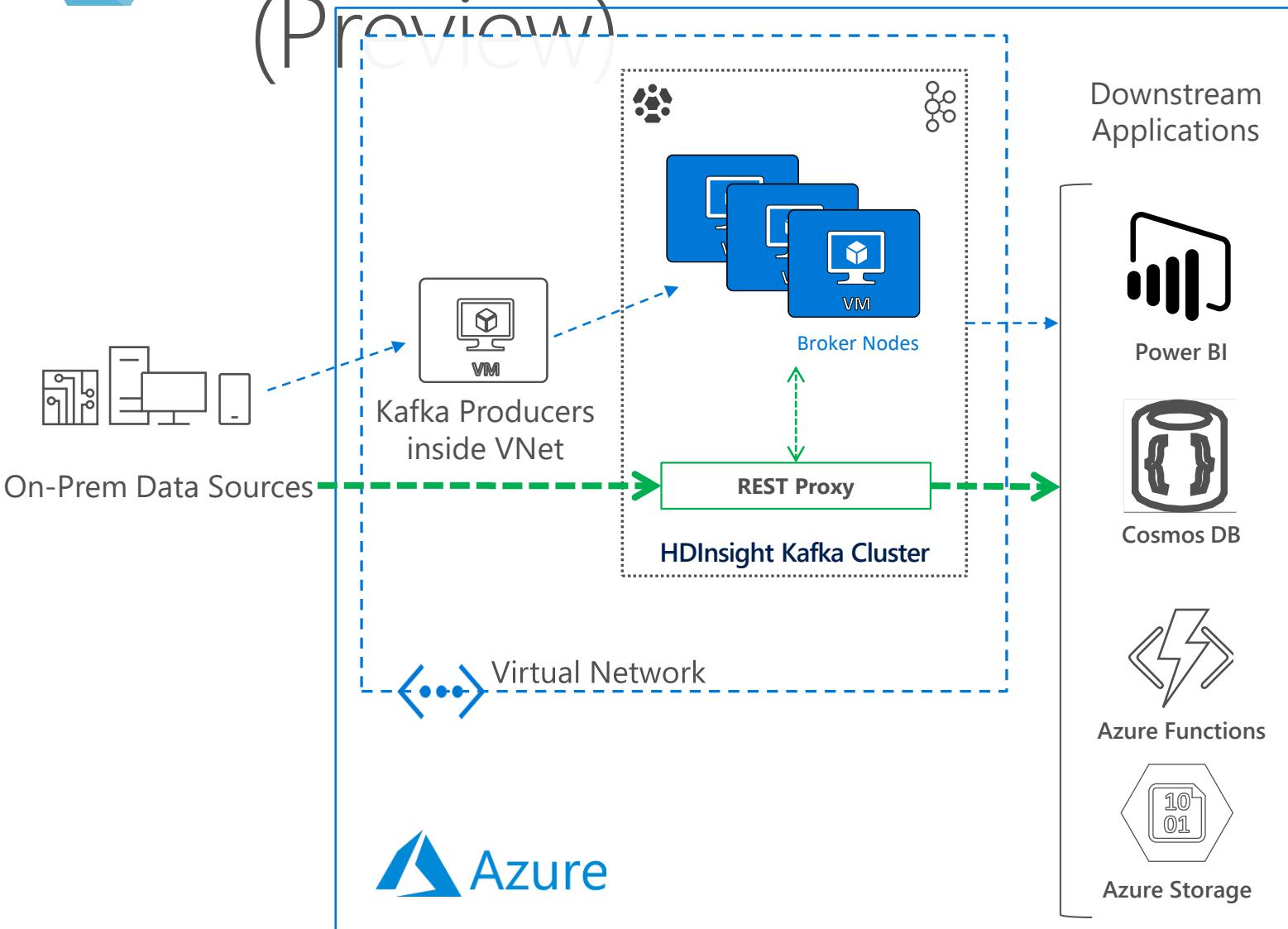
Right pane:

- Name:** [empty]
- Priority:** allowed numeric values between 100-65000
- Action:** Allow
- Rules:** [empty]
- FQDN tags:**
 - Source Addresses:** [empty]
 - Target FQDNs:** [empty]
 - Protocol/Port:** [empty]
- Notes:** FQDN tags may require additional configuration. Learn more.
- Bottom:** Add button



REST Proxy with Kafka on HDIInsight

(Preview)



One-click deploy of **Highly Available** REST proxy with your Kafka cluster

Connect with your Kafka cluster from anywhere

Secured by AAD authorization and OAuth protocol

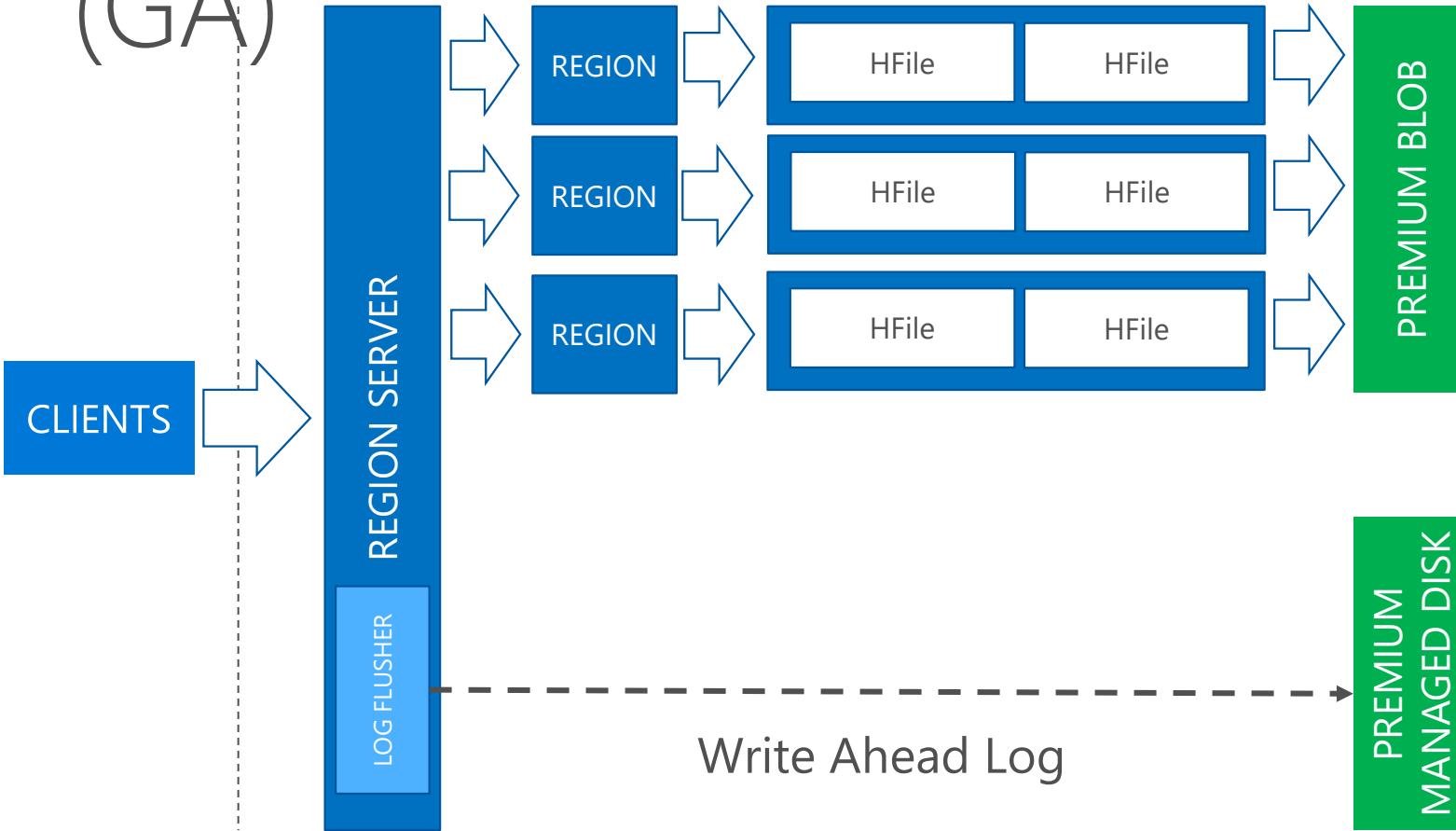
Metadata APIs supported – Get topics, partitions, brokers

Admin APIs supported – Create topics

Data APIs supported – Produce and consume records

HBase Enhanced Read & Accelerated Writes

(GA)



Enhanced Large Reads (YCSB)

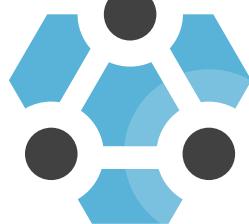
- ✓ Up to 25X greater throughput
- ✓ Up to 11X lower latency

Read/Write to Premium Blob

Accelerated Writes (YCSB)

- ✓ Up to 8X greater throughput
- ✓ Up to 3X lower latency

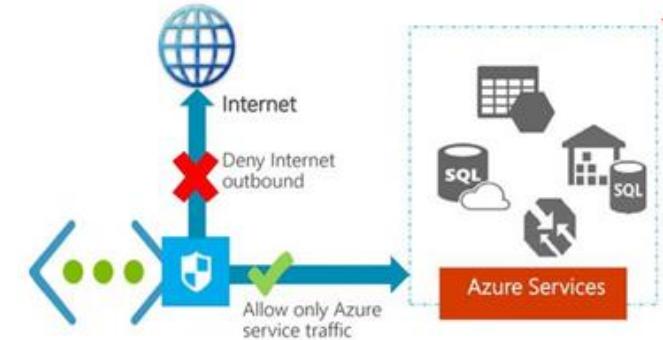
Write the WAL entries to Premium Managed Disk



Service Tag in NSGs

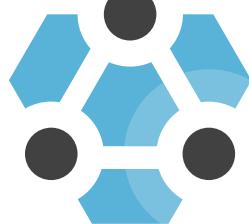
Service tags provide an alternative method for allowing inbound traffic from specific IP addresses

- Customers --restrict network access to just the Azure services you use
- HDInsight --maintenance of IP addresses for each tag provided by Azure



Network Security Group (NSG)				
Action	Name	Source	Destination	Port
Allow	AllowStorage	VirtualNetwork	Storage	Any
Allow	AllowSQL	VirtualNetwork	Sql.EastUS	Any
Deny	DenyAllOutBound	Any	Any	Any

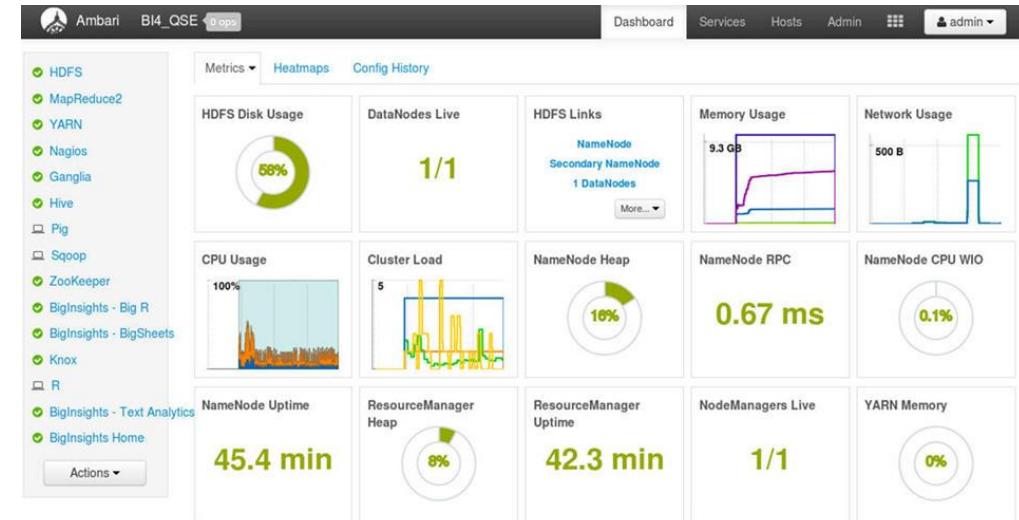
Options	Example	Note
Global service tag	HDInsight global	Open your virtual network to all of the IP Addresses that the HDInsight service is using to monitor clusters across all regions
Regional service tag	HDInsight.NorthCentralUS	open your virtual network to only the IP Addresses that HDInsight is using in that specific region. Some regions may need to whitelist multiple tags due to HDInsight has global RP running for that region.



Bring your own Ambari Database

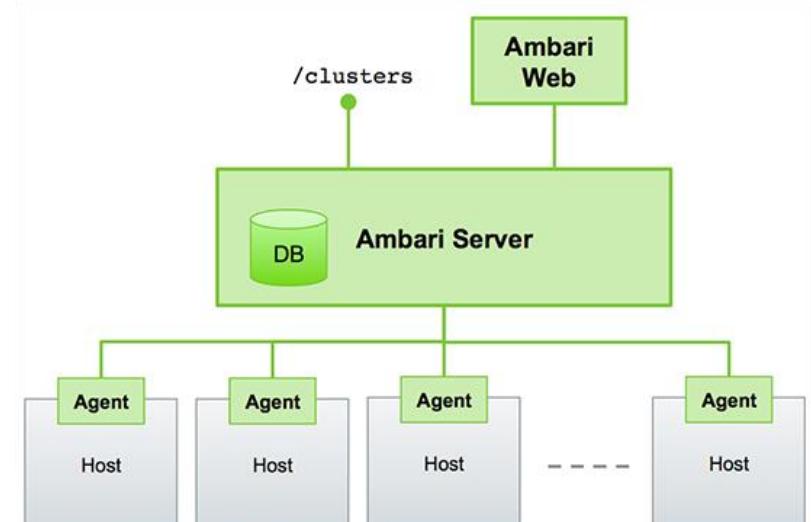
Ambari Database options in HDInsight:

1. HDInsight provisions a default S0 DB for customers
2. For critical workload or heavy usage, customers can bring their own Ambari DB to meet their performance and growth needs



Requirements:

1. Multiple clusters cannot use the same Ambari DB
2. The database must be empty
3. The IP addresses (from HDInsight service) need to be allowed in the SQL Server.



Other noteworthy improvements

- . BYOK: Customer managed key support for disks
- Audit logs [Who is accessing HDInsight?]
- Cluster Creation New UX
- Spark 2.4 & Kafka 2.1
- Azure CLI
- HDInsight 4.0: Hive/Spark Metastore integration for external tables
- F series support