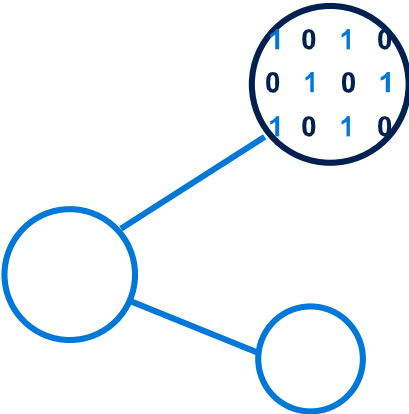


Big Data Ingestion and Storage

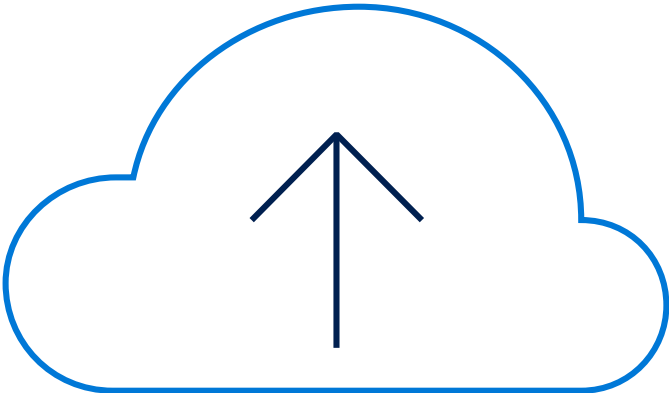
Darwin Schweitzer
Senior Program Manager



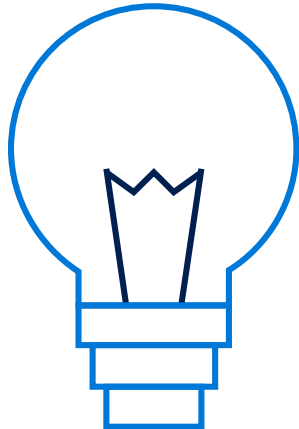
Business is being transformed by three trends



Big Data

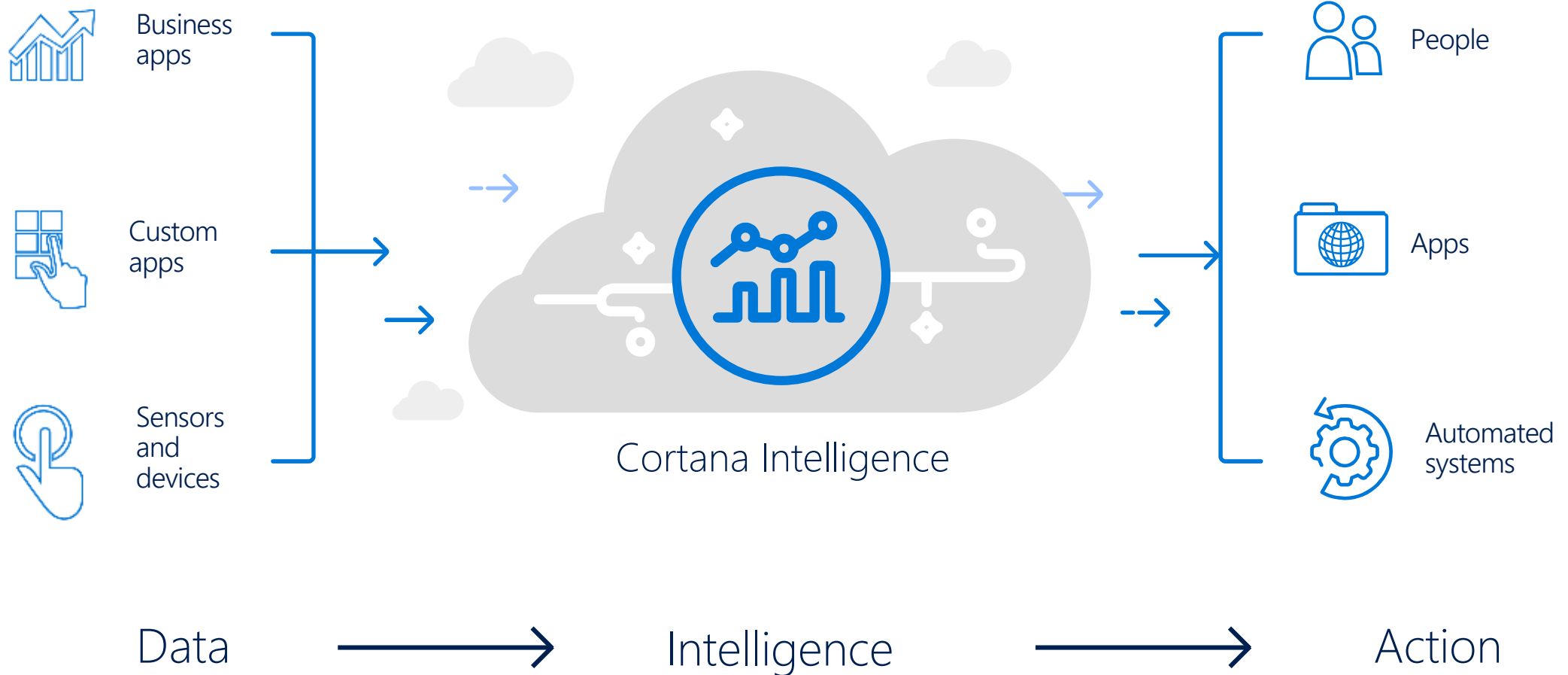


Cloud

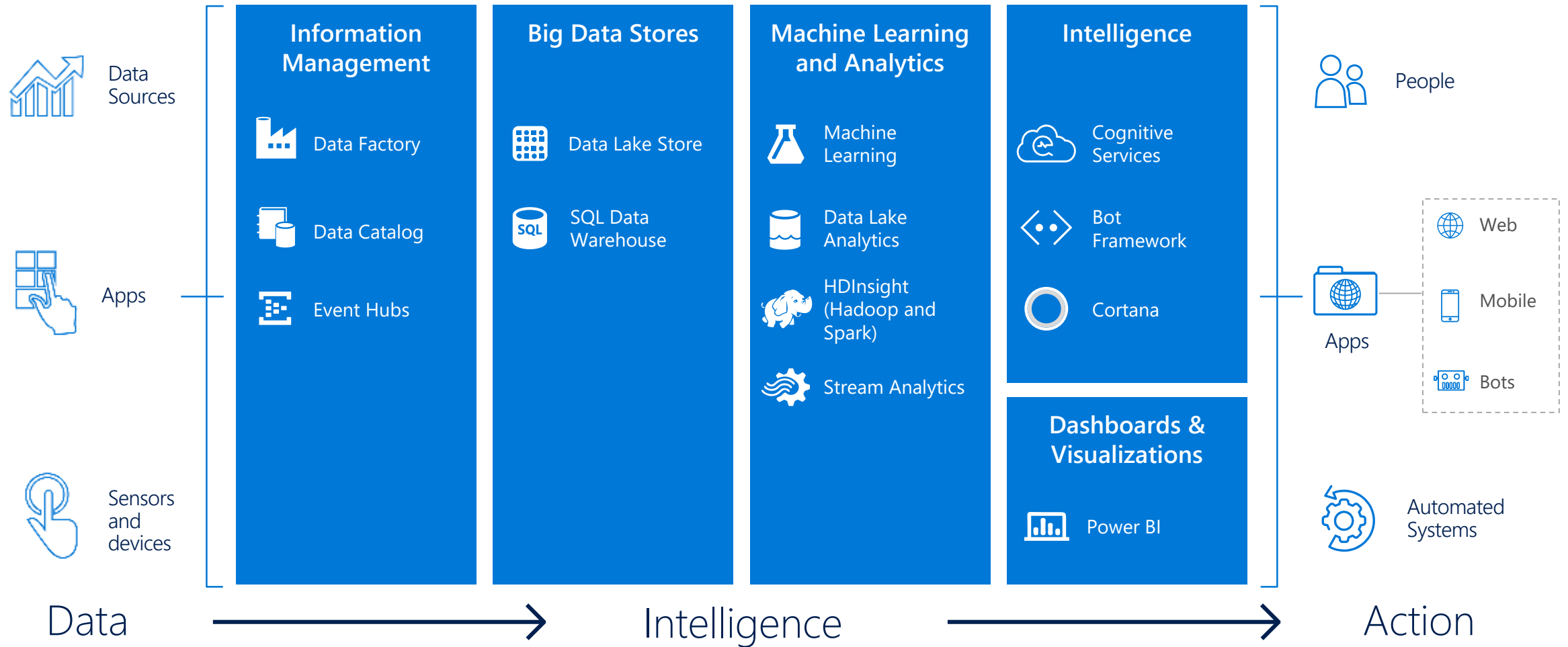


Intelligence

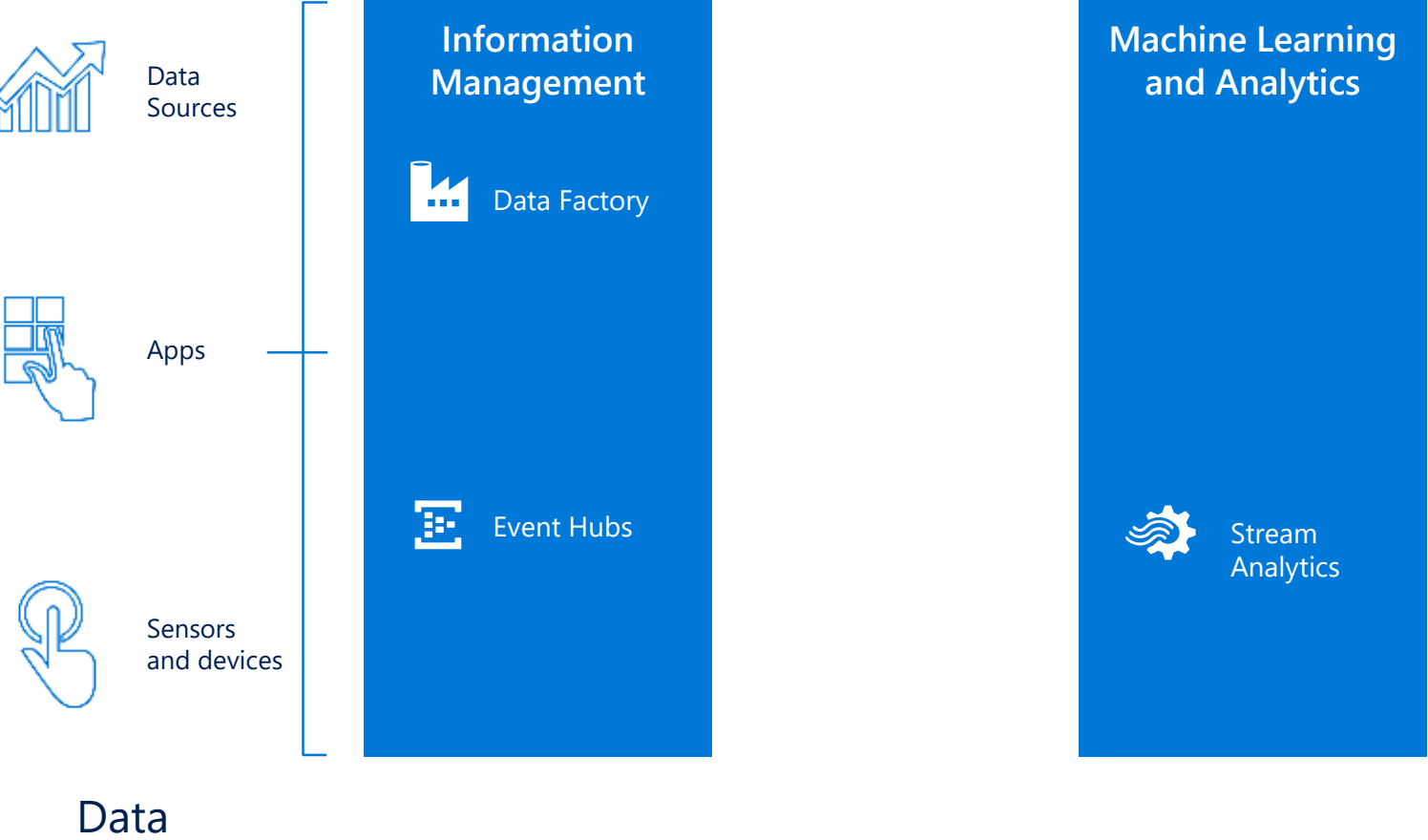
Stay ahead of the curve with Cortana Intelligence Suite



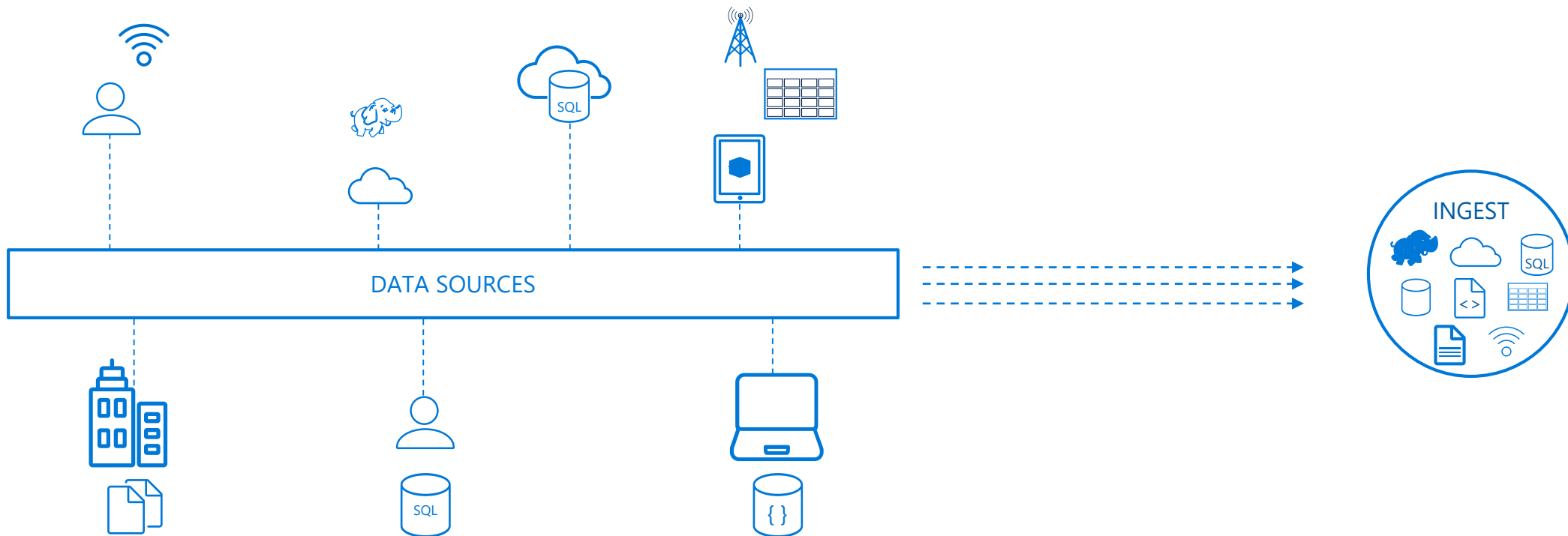
Easily turn data into intelligent action



Big Data Ingestion

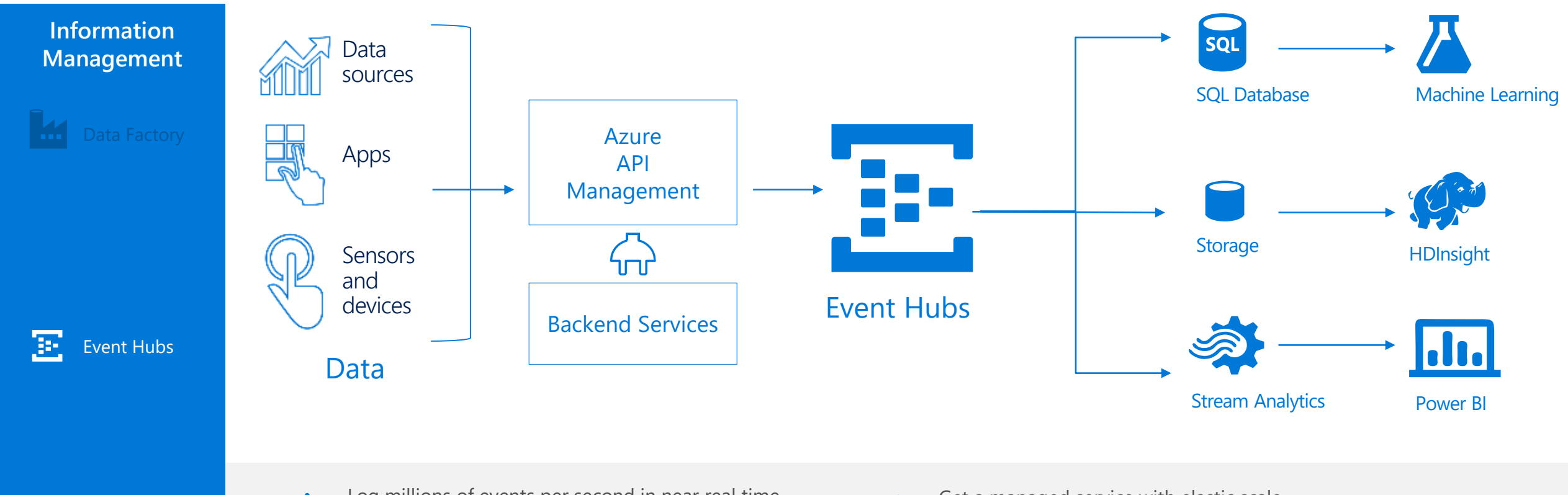


Compose and orchestrate data services at scale



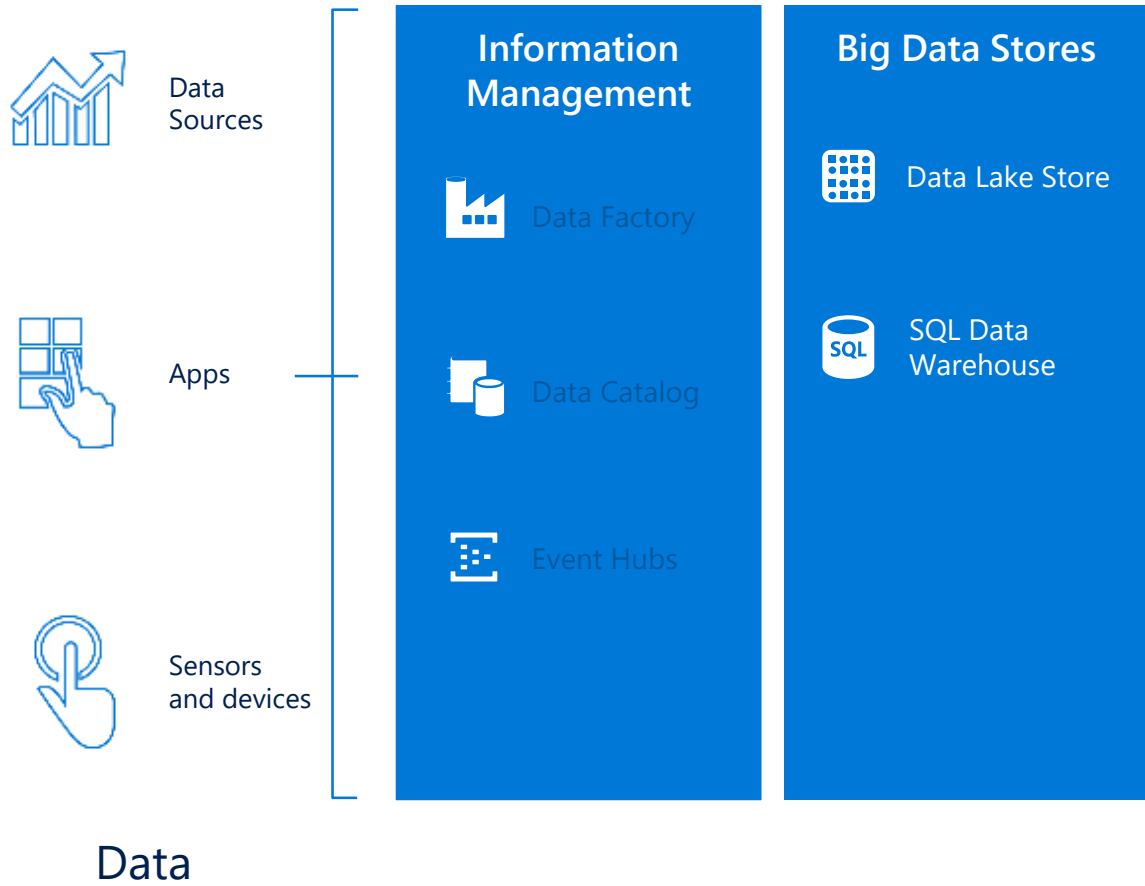
- Create, schedule, orchestrate, and manage data pipelines
- Visualize data lineage
- Connect to on-premises and cloud data sources
- Monitor data pipeline health
- Automate cloud resource management
- Move relational data for Hadoop processing
- Transform with Hive, pig, or custom code

Ingest events from websites, apps and devices at cloud scale



- Log millions of events per second in near real time
- Connect devices using flexible authorization and throttling
- Use time-based event buffering
- Get a managed service with elastic scale
- Get a managed service with elastic scale
- Reach a broad set of platforms using native client libraries
- Pluggable adapters for other cloud services

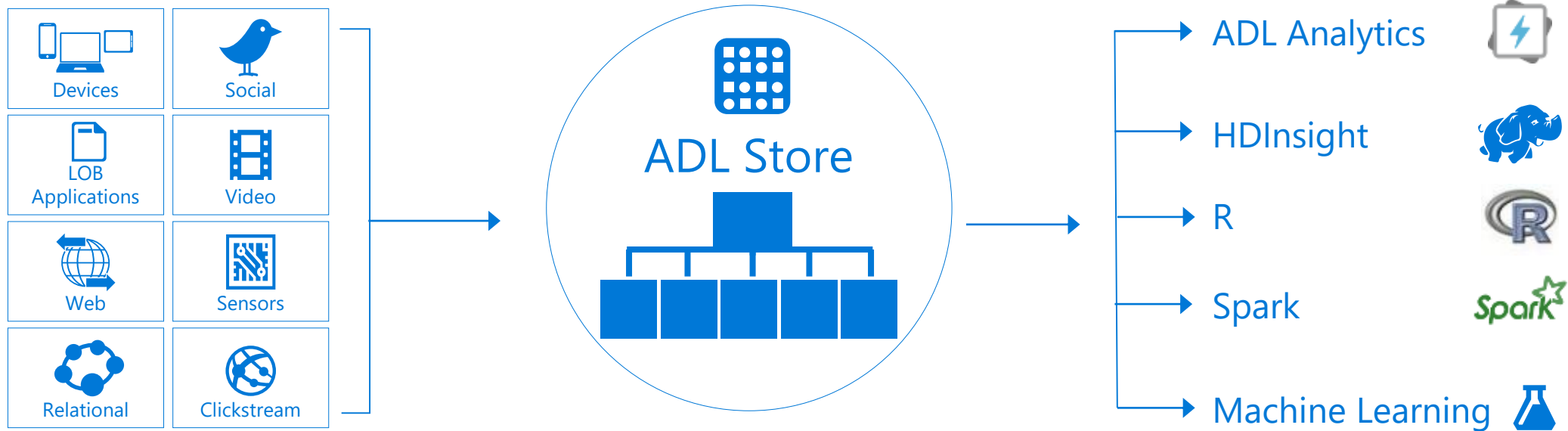
Big Data Stores



A hyper-scale repository for big data analytics workloads

Big Data Stores

- Data Lake Store
- SQL Data Warehouse



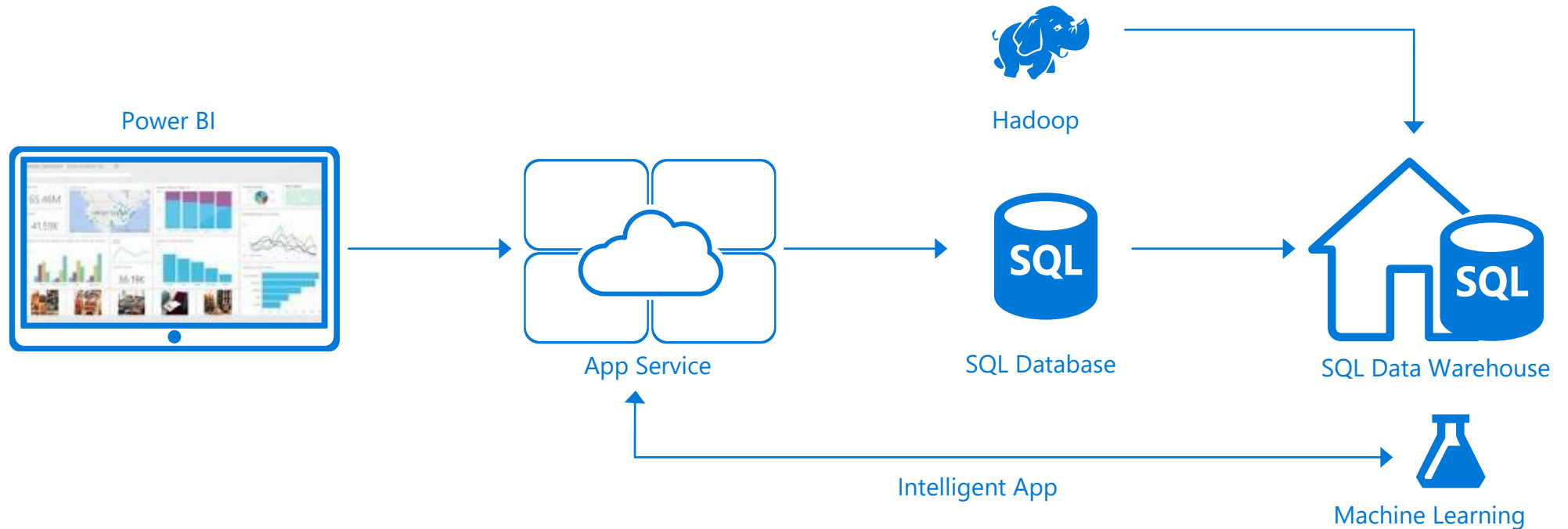
- A Hadoop Distributed File System for the cloud
- No fixed limits on file size
- No fixed limits on account size
- Unstructured and structured data in their native format
- Massive throughput to increase analytic performance
- High durability, availability, and reliability
- Azure Active Directory access control

Elastic data warehouse as a service with enterprise-class features

Big Data Stores

Data Lake Store

SQL Data Warehouse



- Petabyte scale with massively parallel processing
- Independent scaling of compute and storage—in seconds
- Transact-SQL queries across relational and non-relational data

- Full enterprise-class SQL Server experience
- Works seamlessly with Power BI, Machine Learning, HDInsight, and Data Factory

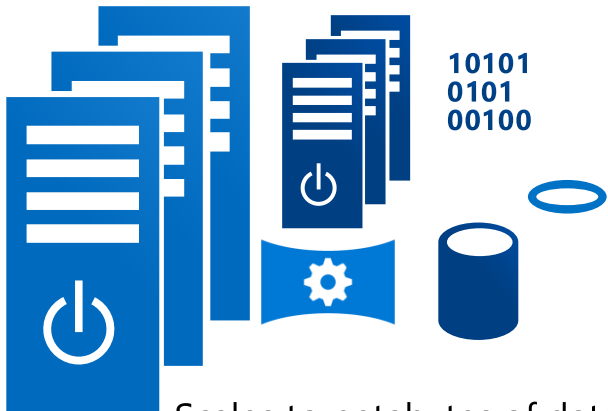
Introducing Azure SQL DW Service

A relational **data warehouse-as-a-service**, fully managed by Microsoft.

Industries first **elastic** cloud data warehouse with proven SQL Server capabilities.

Support your **smallest to your largest** data storage needs.

Elastic scale & performance



Scales to petabytes of data

Massively Parallel Processing

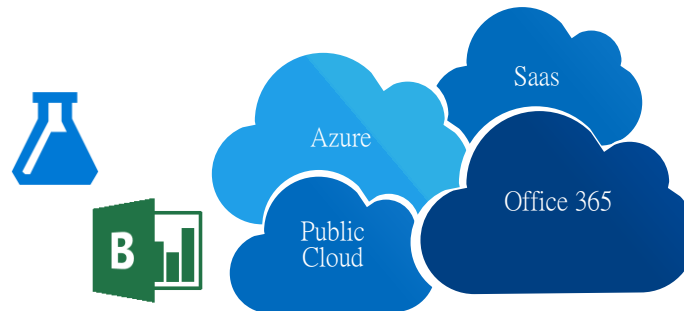
Instant-on compute scales in seconds

Query Relational / Non-Relational

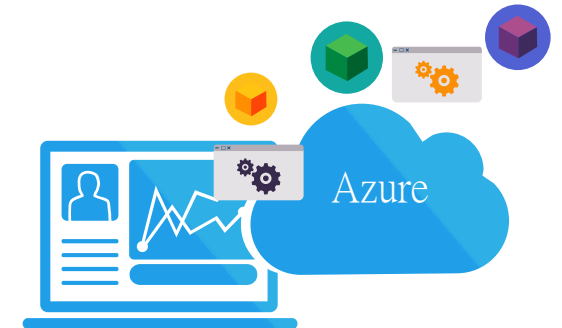
Powered by the Cloud

Get started in minutes

Integrated with Azure ML, PowerBI & ADF



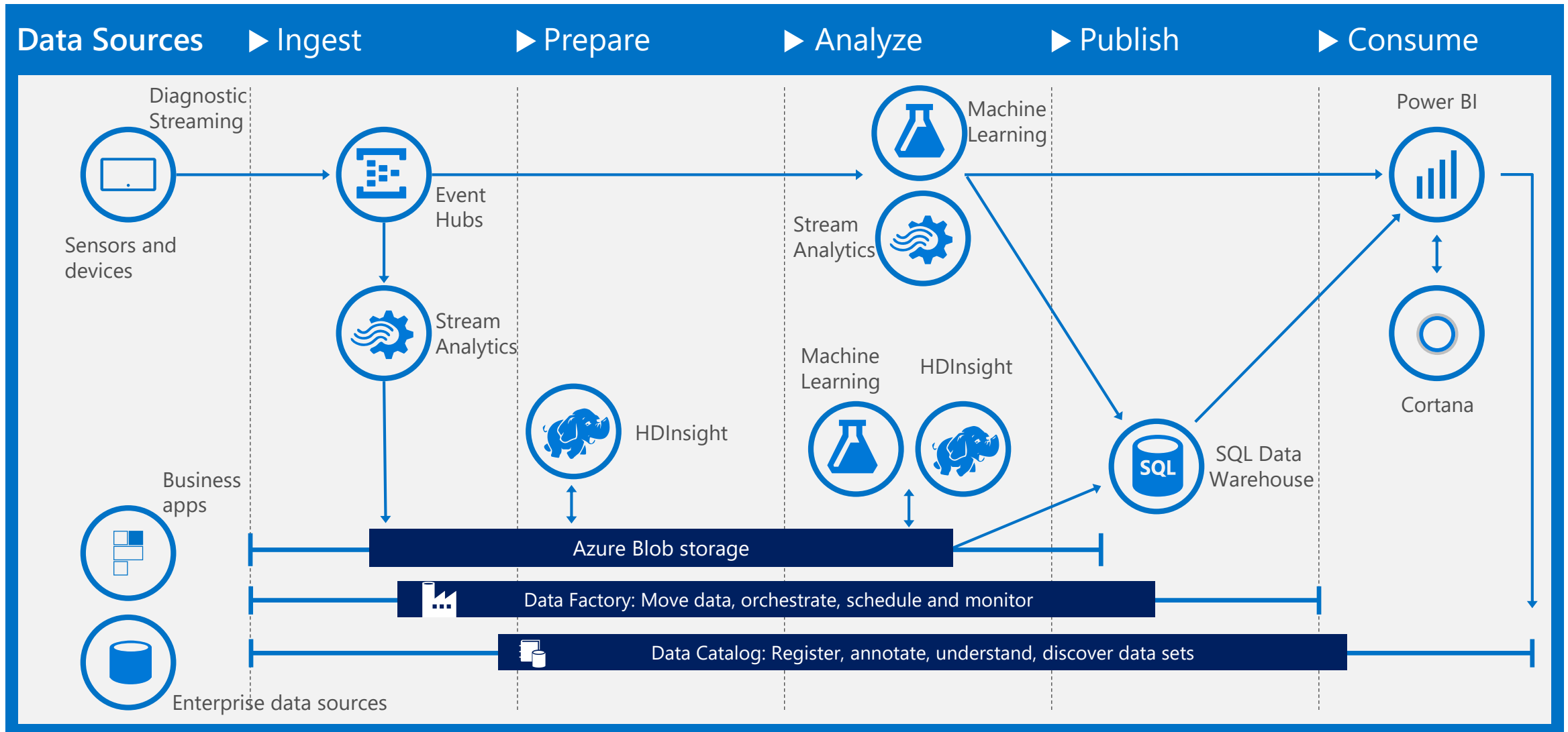
Market Leading Price & Performance



Simple billing compute & storage

Pay for what you need, when you need it with dynamic pause

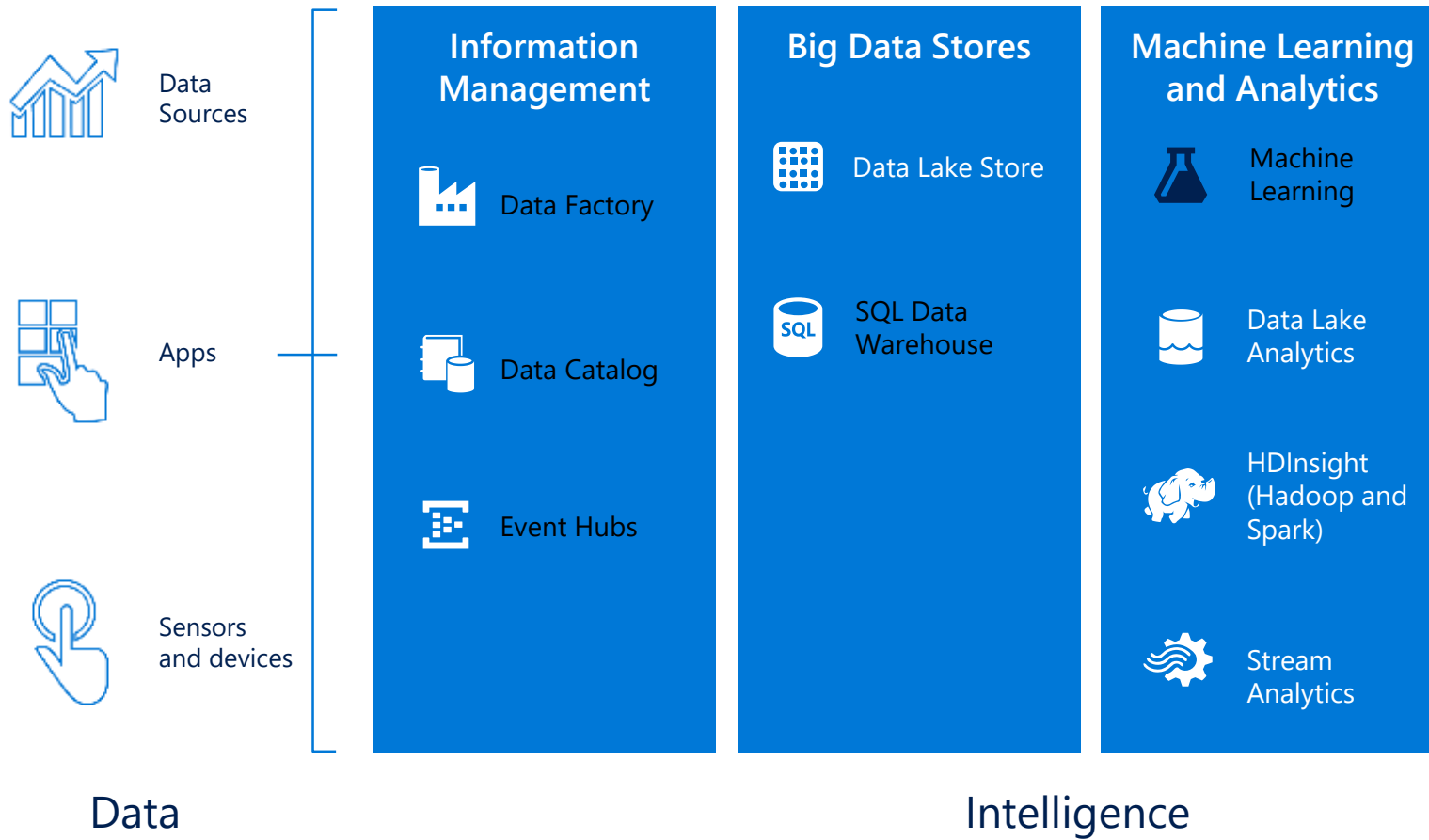
Example of Cortana Intelligence Suite in action



Demo Azure SQL Data Warehouse



Machine Learning and Analytics



Big data analytics made easy



Data Lake Analytics



- Analyze data of any kind and size
- Develop faster, debug and optimize smarter
- Interactively explore patterns in your data
- No learning curve—use U-SQL, Spark, Hive, HBase and Storm

- Managed and supported with an enterprise-grade SLA
- Dynamically scales to match your business priorities
- Enterprise-grade security with Azure Active Directory
- Built on YARN, designed for the cloud

Machine Learning and Analytics



Machine Learning



Data Lake Analytics







HDInsight
(Hadoop and Spark)

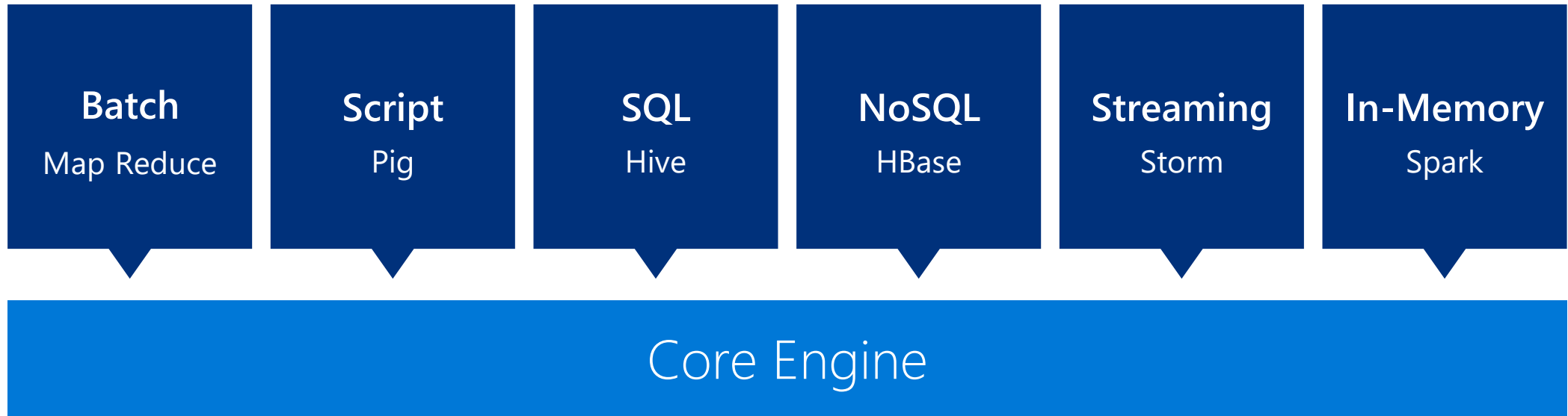


Stream Analytics

Comprehensive set of managed Apache big data projects

Machine Learning and Analytics

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



- Scale to petabytes on demand
- Process unstructured and semi-structured data
- Develop in Java, .NET, and more
- Skip buying and maintaining hardware
- Deploy in Windows or Linux
- Spin up an Apache Hadoop cluster in minutes
- Visualize your Hadoop data in Excel
- Easily integrate on-premises Hadoop clusters

Azure HDInsight running Linux

Choice of Windows or Linux clusters

Managed & supported by Microsoft

Re-use common tools, documentation, samples from Hadoop/Linux ecosystem

Add Hadoop projects that were authored on Linux to HDInsight

Easier transition from on-premises to cloud



Strata Announcements

<https://blogs.technet.microsoft.com/machinelearning/2016/03/29/microsoft-makes-big-data-analytics-easier-in-the-cloud/>

Strata+Hadoop WORLD

MAKE DATA WORK
MARCH 28-29, 2016: TRAINING
MARCH 29-31, 2016: CONFERENCE
SAN JOSE, CA

SCHEDULE | SPEAKERS | EXPO | EVENTS | VENUE | ABOUT | RESOURCES

Strata + Hadoop World 2016 Speakers

New speakers are added regularly. Please check back to see the latest updates agenda.

All Speakers

Sessions

Tutorials

2-day Training

Search Speakers

Joseph S



Joseph Sirosh (Microsoft), @josephsirosh

Connected brains Keynote

Joseph Sirosh is the corporate vice president of the Data Group, le machine-learning products, as well as a talented team of engineer: developing tools and services to transform data at scale into action

Machine Learning Blog

Microsoft Makes Big Data and Analytics Easier in the Cloud

★★★★★

March 29, 2016 by [ML Blog Team](#) // [0 Comments](#)

[f](#) 0 [t](#) 0 [in](#) 0

This post is by Joseph Sirosh, Corporate Vice President of the Data Group at Microsoft.

This week I'm joining thousands of people attending [Strata + Hadoop World](#) in San Jose to explore the technology and business of big data and data science. As part of our participation in the conference, we are announcing several important investments to continue delivering on our commitment to make big data processing and analytics simpler and more accessible:

- **Advanced analytics at scale with R Server for HDInsight and the latest version of Spark for HDInsight are now available in preview:** Customers can leverage their existing R skills and reuse current code to run at scale. R Server for HDInsight offers popular scalable R algorithms and the ability to parallelize any existing R function. We are also releasing the latest version of Spark for HDInsight, which can deliver **7x** performance over MapReduce for most analytics. These capabilities give our customers the ability to train and run advanced analytics and ML models on larger datasets, and much faster than previously possible in the cloud.
- **Out-of-the-box application integration, providing easier access to popular big data apps:** Customers can now discover and deploy popular big data applications with HDInsight without any code or scripting required. Leading solutions such as [Datameer Cloud](#) offer code-free data preparation, [AtScale](#) has cloud-based OLAP BI on Hadoop, and an ecosystem of other big data applications can now be deployed alongside HDInsight.
- **Azure Data Catalog**, previously announced as a [public preview](#) will be [generally available tomorrow](#): Data Catalog is an enterprise



Microsoft Azure Data Lake

Analytics Service

HDInsight

U-SQL spark STORM

YARN

HDFS

Store



Demo Azure Data Lake



If you would like Azure Data Lake Preview Access

Name	Azure Email Account	Azure SubscriptionID
Darwin Schweitzer	datacommoner@outlook.com	bcb1d5d2-e6ea-492d-b9c7-xxxxxxxxxxxx

Send email to darsch@microsoft.com

With your:

To use in HOL or during
the Hackathon

Name

Azure Email Account

Azure SubscriptionID

https://caqs.azure.net

The screenshot shows a web browser window displaying the Cortana Intelligence Quick Start Gallery. The browser address bar shows `caqs.azure.net/#gallery`. The gallery interface includes a sidebar with navigation options: HOME, GALLERY, and a search icon. The main content area displays several pattern cards:

- Microsoft Data Science Virtual Machine** (READY): Features an image of a flask with green liquid and bubbles. Description: "This pattern provisions a Data Science Virtual Machine".
- Enterprise Data Warehouse with Hive** (DRAFT): Includes a detailed diagram of a data pipeline. It shows data flowing from "On Premise Data Sources" through "Ingest" (Azure Blob Storage), "Sanitize" (Azure HDInsight), and "Data Warehouse (Data Mart)" (Azure SQL Data Warehouse) to "Visualization" (Power BI, Databricks, etc.).
- Energy Demand Forecasting** (DRAFT): Shows a flow from "External Data" through "Azure Services" (Azure WebJob, Azure Stream Analytics, Azure Data Factory, Azure SQL) to "Power BI" and "AML Model".
- Multi-step Automation** (FOR INTERNAL USE): Features the GitHub logo and code symbols (</> and >).
- NY Rides** (READY): Shows a group of people in a yellow taxi cab.
- AML with On-Premises SQL Server** (READY): Includes a diagram showing data flow from an "On Premise SQL Server" through "Azure Data Factory" and "Azure Storage Blob Input" to "Azure Data Factory" and "Azure Data Factory Pipeline".

For tomorrows session Power BI with Big Data Stores

Homework

- Go to <https://caqs.azure.net/#gallery/datasciencevm>
- Sign In with your Azure Subscription account
- Accept the Terms of Use for your Azure Subscription
[Configure Programmatic Deployment](#)
- Click the Continue button to provision the Data Science VM
- Fill in parameters and click Create
- Connect to the VM and login

Microsoft Data Science Virtual Machine

Gopi Kumar · published on 02/08/2016

This pattern will provision the **Microsoft Data Science Virtual Machine** on Azure via CAQS.

This Linked Resource pattern will help you get started with the following:

- Provision a Data Science VM as a CAQS Project Linked Resource
- Help you find your Data Science VMs without going to the Azure Portal
- Use the Data Science VM with one or more other CAQS Design Pattern Projects
- If are wondering about things you can do with the DSVM read this [How-To Guide to the Data Science Virtual Machine](#)

STOP before you proceed Want to deploy programmatically? Get Started by accepting the Terms of Use for your Azure Subscription [Configure Programmatic Deployment](#)

Ingredients

- Azure Storage
- Azure IaaS Virtual Machine

[Continue](#) (Sign-in is required)



Create new project

Project name

(Project name must be between 3 and 11 characters, start with a lowercase letter, and contain only lowercase letters and numbers.)

Subscription

Azure ML Build-Demo

(0f74ada4-f6bb-41d0-850e-e9ceb738df2a)

Locations

Southeast Asia

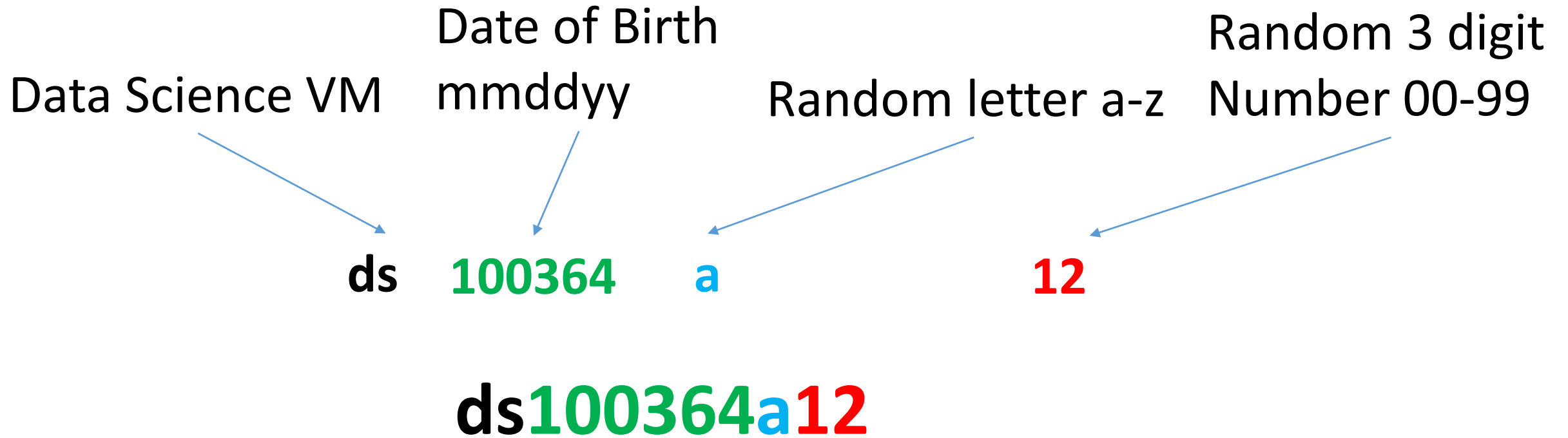
Description (optional)

Microsoft Data Science Virtual Machine

[Cancel](#) [Create](#)

CAQS Project Naming recommendation

Pattern Id(first two digits) your DOB (next 6 digits mmddyy) Random letter a-z , Random 2 digit number between 00-99



Create new project ✕

Project name

 ✓

(Project name must be between 3 and 11 characters, start with a lowercase letter, and contain only lowercase letters and numbers.)

Subscription

 ▼

(bcb1d5d2-e6ea-492d-b9c7-f461c2a94a92)

Locations

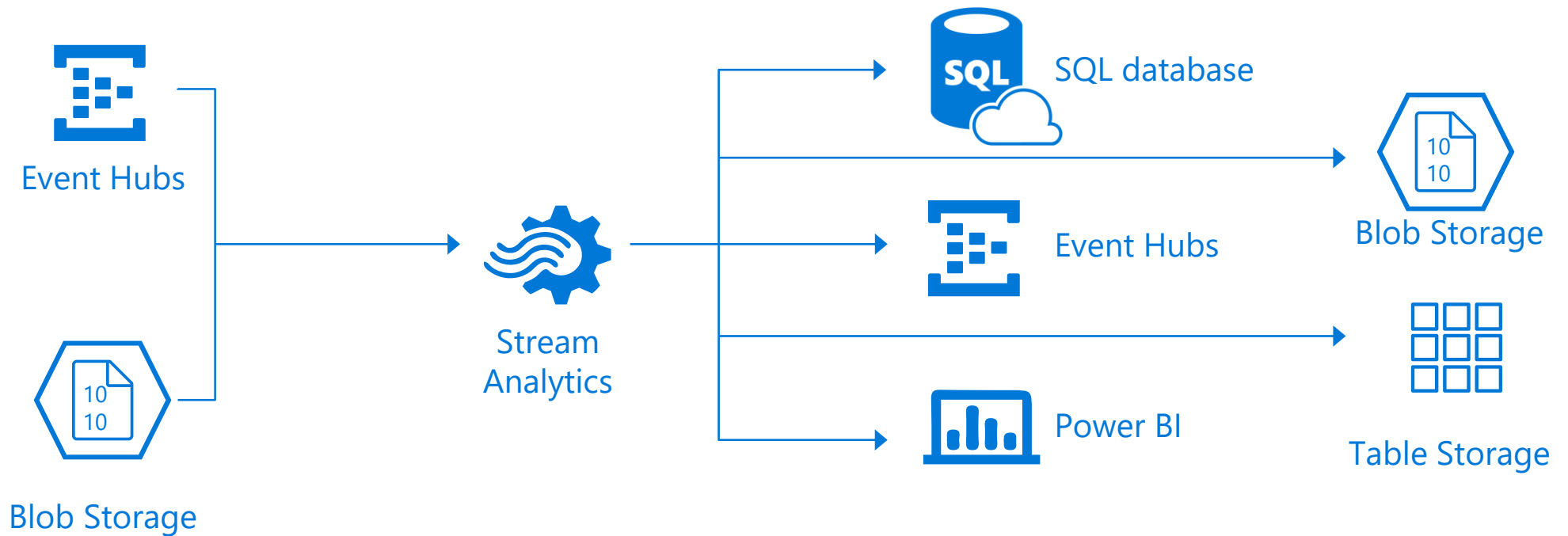
 ▼

Description (optional)

Real-time stream processing in the cloud

Darwin Schweitzer
@DataSnowman

Stream Analytics



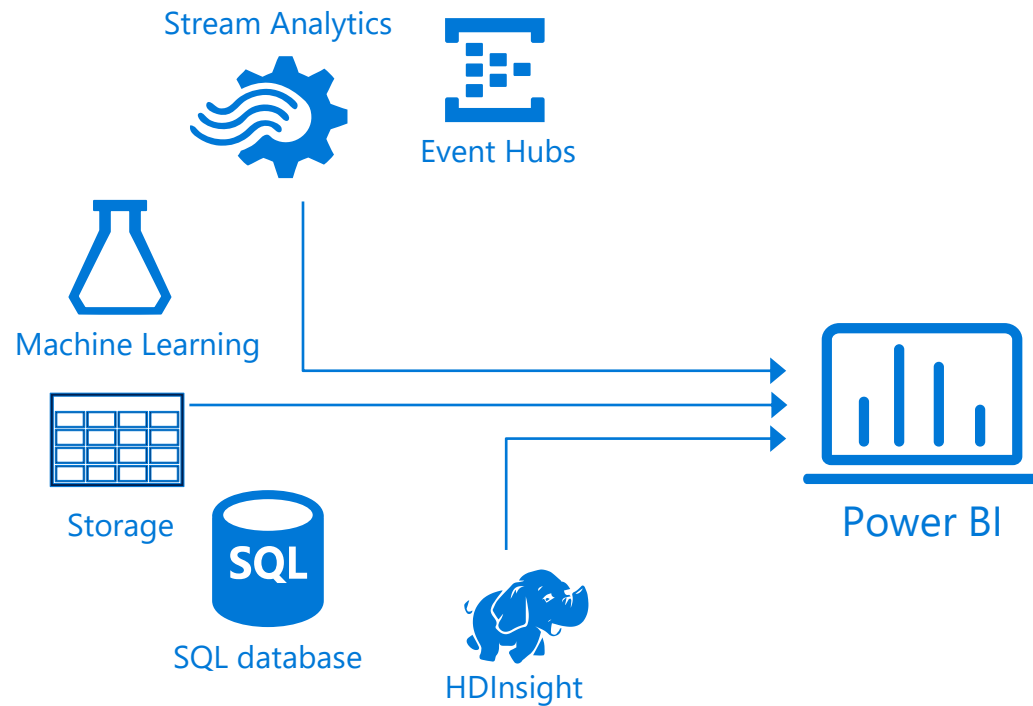
- Perform real-time analytics for your Internet of Things solutions
- Stream millions of events per second
- Get mission-critical reliability and performance with predictable results

- Create real-time dashboards and alerts over data from devices and applications
- Correlate across multiple streams of data
- Use familiar SQL-based language for rapid development

Keep a pulse on your business with live, interactive dashboards

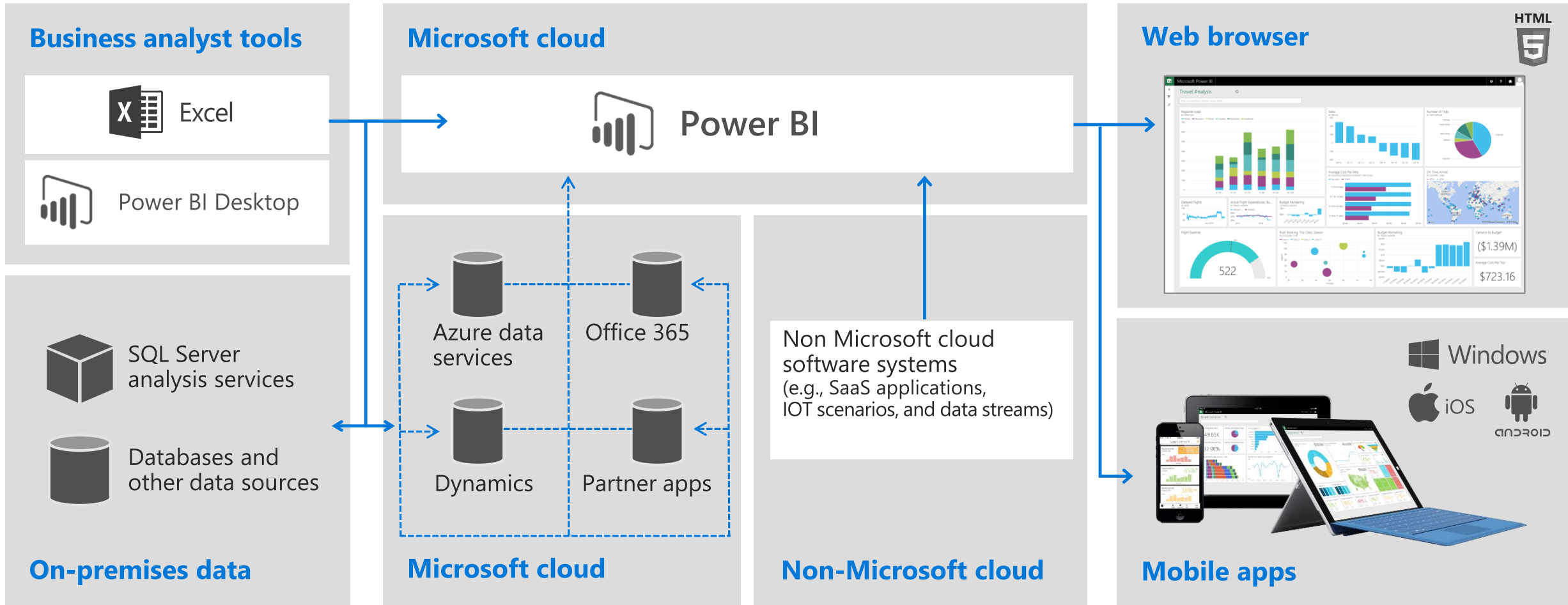
Dashboards & Visualizations

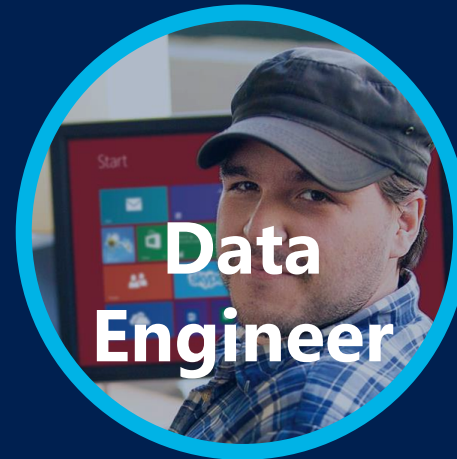
 Power BI



- Analytics for everyone, even non-data experts
- Your whole business on one dashboard
- Create stunning, interactive reports
- Drive consistent analysis across your organization
- Embed visuals in your applications
- Get real-time alerts when things change

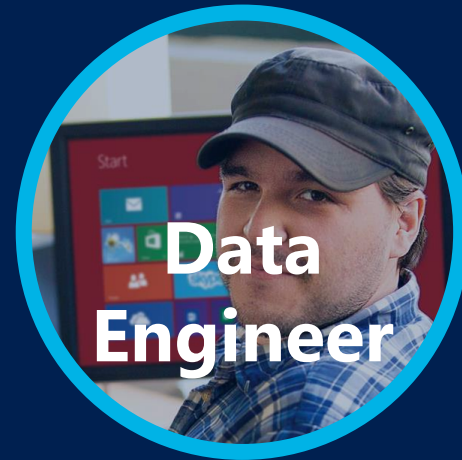
Power BI – Business Users and Data Analysts



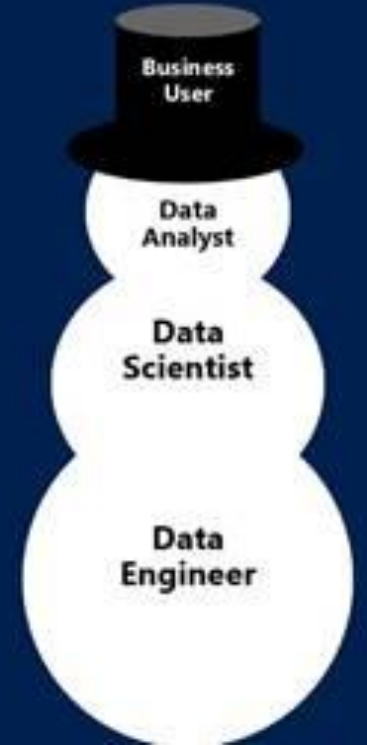


Intuitive
Accessible

- Simple
- Visual



DataSnowman



[DASHBOARD](#)[MONITOR](#)[INPUTS](#)[FUNCTIONS **PREVIEW**](#)[QUERY](#)[OUTPUTS](#)[SCALE](#)[CONFIGURE](#)

Need help with your query? Check out some of the most common Stream Analytics query patterns [here](#).

query

```
1 WITH TipStream As (SELECT System.Timestamp AS BookingTime, PickupLocation, DestinationLocation,
2 count(*) as NumBookings, max(TipAmount) as MaxTipAmount, min(TipAmount) as MinTipAmount,
3 avg(TipAmount) as AvgTipAmount
4 FROM bookingstream TIMESTAMP BY CreatedTime
5 GROUP BY PickupLocation, DestinationLocation, TumblingWindow(second,30) )
6 SELECT BookingTime, PickupLocation, DestinationLocation, NumBookings, AvgTipAmount
7 INTO driverstip FROM TipStream
8 SELECT * INTO pbipassenger FROM [passengerstream] SELECT * INTO pbidriver FROM [driverstream]
```



Azure Stream Analytics: easy to write SQL queries over streaming data

Azure Portal and Visual Studio - Data Engineers

Start Page - Microsoft Visual Studio

File Edit View Debug Team Data Lake Tools Test Analyze Window Help

Server Explorer

One or more subscriptions are not supported by Server Explorer. [Open Cloud Explorer to view all subscriptions.](#)

- Azure (darsch@microsoft.com - 3 subscriptions)
 - App Service
 - Cloud Services
 - Data Factory
 - Data Lake Analytics
 - HDInsight
 - Mobile Services
 - Notification Hubs
 - Service Bus
 - SQL Databases**
 - Storage
 - Virtual Machines
- Data Connections
 - atest.master.dbo
 - atest.pbicp.dbo
 - bkic5twjyh.TollDataDB.dbo
 - gf0ag27whj.tweets.dbo
- Servers
 - darwinsx1c

Start Page - Microsoft Visual Studio

File Edit View Debug Team Data Lake Tools Test Analyze Window Help

SQL Server Object Explorer

- SQL Server
 - (localdb)\MSSQLLocalDB (SQL Server 12.0.2000 - REDMOND\darsch)
 - (localdb)\ProjectsV12 (SQL Server 12.0.2000 - REDMOND\darsch)
 - afinal0127.database.windows.net (SQL Server 12.0.2000 - darwin, nyrides)
 - asqldw0303dwsrv.database.windows.net (SQL Server 10.0.2531.0 - username, asqldw0303db)
 - Databases
 - System Databases
 - master
 - asqldw0303db
 - Tables
 - System Tables
 - External Tables
 - dbo.bad_coordinate_count
 - dbo.fare_aggregate_by_paytype
 - dbo.learningbycount_feature
 - dbo.nyctaxi_fare
 - dbo.nyctaxi_trip
 - dbo.tipped_distribution
 - dbo.tip_class_distribution
 - dbo.trip_distribution_by_med**
 - dbo.trip_distribution_by_med_hack
 - Views
 - Programmability
 - Security
 - External Resources
 - Security
 - atest.database.windows.net (SQL Server 12.0.2000 - darwin)
 - atest.database.windows.net (SQL Server 12.0.2000 - darwin, pbicp)
 - tchurn02261sql.database.windows.net (SQL Server 11.0.9231 - custchurnadmin, customerchurn)
 - w52etupmdm.database.windows.net (SQL Server 10.0.2531.0 - darwin, snowman)
 - Projects



- New
- All resources
- Recent
- SQL databases
- Data Lake Store
- Data factories
- HDInsight Clusters
- Subscriptions
- Storage accounts
- Storage accounts (classic)
- Data Lake Analytics

mwinklehivefastpath
HDInsight Cluster

Settings Dashboard Remote Desktop... Scale Cluster Delete

Essentials

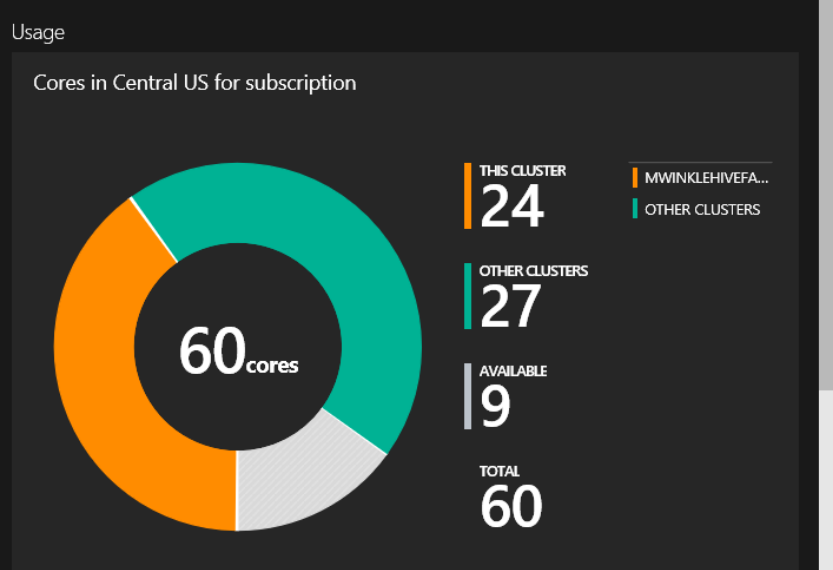
Resource group: **mwinklestrata**
Status: **Running**
Location: **Central US**
Subscription name: **BDHadoopTeamPMTTestDemo2**
Subscription id: **ff78024c-d5a1-48ae-88eb-d61f0f60f8ff**

URL: **mwinklehivefastpath.azurehdinsight.net**
Cluster Type: **Hadoop on Windows**
Head Node, Worker Nodes: **D12 (x2), D12 (x4)**
Learn more: [Documentation](#)
Getting Started: [Quickstart](#)

[All settings](#)

Quick Links

Cluster Dashboard ... Documentation **Scale Cluster**



Cluster nodes

Scale Cluster
mwinklehivefastpath

Save

Number of Worker nodes:

Worker Nodes Pricing Tier: **D12 (4 nodes)**

Head Node Pricing Tier: **D12 (2 nodes)**

WORKER NODES	0.81 x 4 = 3.22
HEAD NODES	0.81 x 2 = 1.61
TOTAL COST	4.84

USD/HOUR (ESTIMATED)

Using 51 of 60 total cores in Central US

This estimate does not include subscription discounts or storage costs.

Questions? [Contact billing support](#)
To learn more, visit our [pricing page](#).



- New
- All resources
- Recent
- SQL databases
- Data Lake Store
- Data factories
- HDInsight Clusters
- Subscriptions
- Storage accounts
- Storage accounts (classic)
- Data Lake Analytics

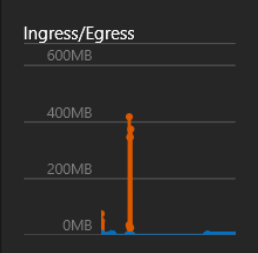
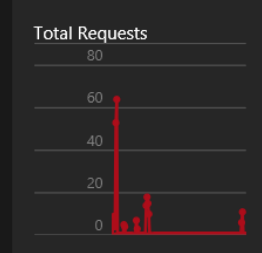
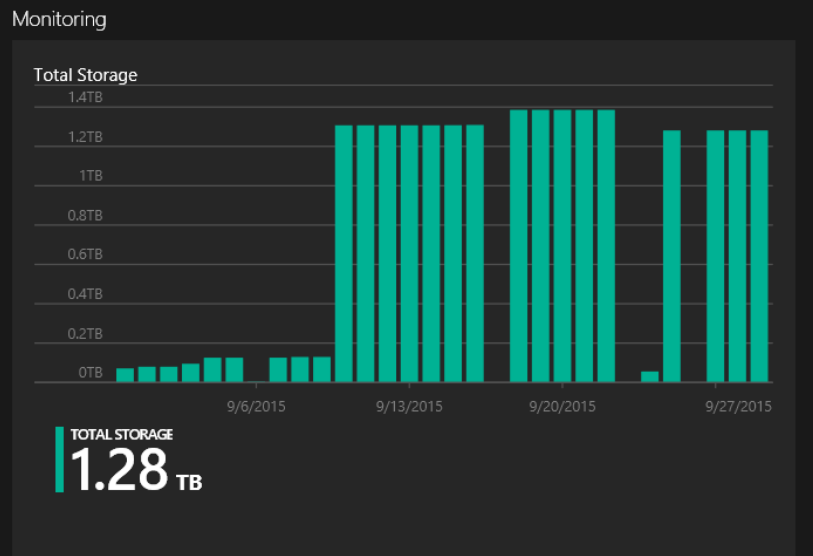
mwinkleadl
Data Lake Store

Settings Delete Data Explorer

Essentials

Resource group ajjaz-sandbox	Pricing tier Pay-As-You-Go
Status Running	Location East US 2
URL https://mwinkleadl.azuredatalake.net	Subscription BDHadoopTeamPMTTestDemo2
WebHDFS URI swebhdfs://mwinkleadl.azuredatalake.net	Subscription ID ff78024c-d5a1-48ae-88eb-d61f0f60f8ff

[All settings](#) →



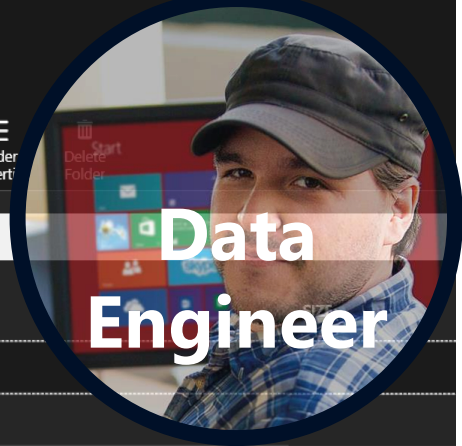
Data Explorer
mwinkleadl

- mwinkleadl
 - catalog
 - ngrams
 - Outputs
 - sampledata
 - Samples
 - system
 - tmp
 - users

mwinkleadl
Data Lake Store

New Folder Upload Access Rename Folder Folder Properties

- mwinkleadl
- NAME
- catalog
 - ngrams
 - Outputs
 - sampledata
 - Samples
 - system
 - tmp
 - users



mwinkle
Data Lake Analytics

Settings | New Job | Add Data Source | View All Jobs | Data Explorer | Delete

Essentials

Resource group: **ajjaz-sandbox**
 Status: **Running**
 Location: **East US 2**
 Subscription Id: **ff78024c-d5a1-48ae-88eb-d61f0f60f8ff**
 Subscription: **BDHadoopTeamPMTTestDemo2**

Pricing tier: **Pay-As-You-Go**
 Default Data Lake Store: **mwinkleadl**
 Learn: [Explore sample jobs](#)
 Getting Started: [Explore interactive tutorials](#)

Job Management

TOTAL **307** | SUCCEEDED **109** | FAILED **184** | CANCELLED **14**

Usage

Compute Hours

Legend: SUCCEEDED (green), FAILED (red), CANCELLED (blue)

New U-SQL Job

Submit Job | Data Explorer | Import Local File

* Job Name: Priority: Parallelism:

```

3 @allRequests = EXTRACT UTCDate string
4                       , ActivityId string
5                       , Account string
6                       , operation string
7                       , HttpStatus string
8                       , Latency long
9                       FROM @"users/mwinkle/output/telemetryExplore/jobRequests3.csv" USING
Extractors.Csv();
10
11 @allRequests = SELECT *, Microsoft.Analytics.Internal.Telemetry.Helpers.Dates.GetBeginningOfWeek
(UTCDate).ToShortDateString() AS BusinessWeek FROM @allRequests;
12
13 @summary = SELECT DISTINCT BusinessWeek
14           , operation
15           , PERCENTILE_CONT(0.25) WITHIN GROUP (ORDER BY Latency) OVER (PARTITION BY
BusinessWeek, operation) AS [25thPercentileLatency]
16           , PERCENTILE_CONT(0.5) WITHIN GROUP (ORDER BY Latency) OVER (PARTITION BY
BusinessWeek, operation) AS MedianLatency
17           , PERCENTILE_CONT(0.75) WITHIN GROUP (ORDER BY Latency) OVER (PARTITION BY
BusinessWeek, operation) AS [75thPercentileLatency]
18           , PERCENTILE_CONT(0.90) WITHIN GROUP (ORDER BY Latency) OVER (PARTITION BY
BusinessWeek, operation) AS [90thPercentileLatency]
19           , PERCENTILE_CONT(0.95) WITHIN GROUP (ORDER BY Latency) OVER (PARTITION BY
BusinessWeek, operation) AS [95thPercentileLatency]
20           , PERCENTILE_CONT(0.99) WITHIN GROUP (ORDER BY Latency) OVER (PARTITION BY
UTCDate, operation) AS [99thPercentileLatency]
21           FROM @allRequests;
22
23 @summaryCountMinMaxAvg = SELECT BusinessWeek
24           , operation
25           , COUNT(DISTINCT Account) AS Accounts
26           , COUNT(DISTINCT ActivityId) AS Operations
27           , MIN(Latency) AS MinLatency
28           , MAX(Latency) AS MaxLatency
29           , AVG(Latency) AS AvgLatency

```



Azure Data Lake

Currently in Public Preview

A **cloud scale HDFS store** designed for parallel processing workloads

Accessible to all HDFS compliant analytics applications and tools

No limits to scale

- Petabyte files, unlimited account size
- High throughput to increase analytics performance
- Low latency ingestion with immediate access to data

Intelligent data storage

- Multiple tiers of storage
- Dynamic tiering based workload performance requirements
- Data lifecycle management to manage costs

Enterprise grade security

- AAD based access control of files and folders
- Encryption of data at rest



AzureBigAnalyticsDataLakeSample
Data factory

Delete

Essentials

Resource group: ADF
Subscription name: MDP_492270
Subscription id: 1e42591f-1f0c-4c5a-b7f2-a268f6105ec5

Location: WestUS
Provisioning state: Succeeded

Summary

- Author and deploy
- Diagram** (No permission to resource group: ADF)

Contents

Datasets	Pipelines	Linked services
4	3	4
With errors: 2	ComputeEventsByRegion...	Data Stores: 4
	EgressBlobToDataLakePi...	Data Gateways: 0

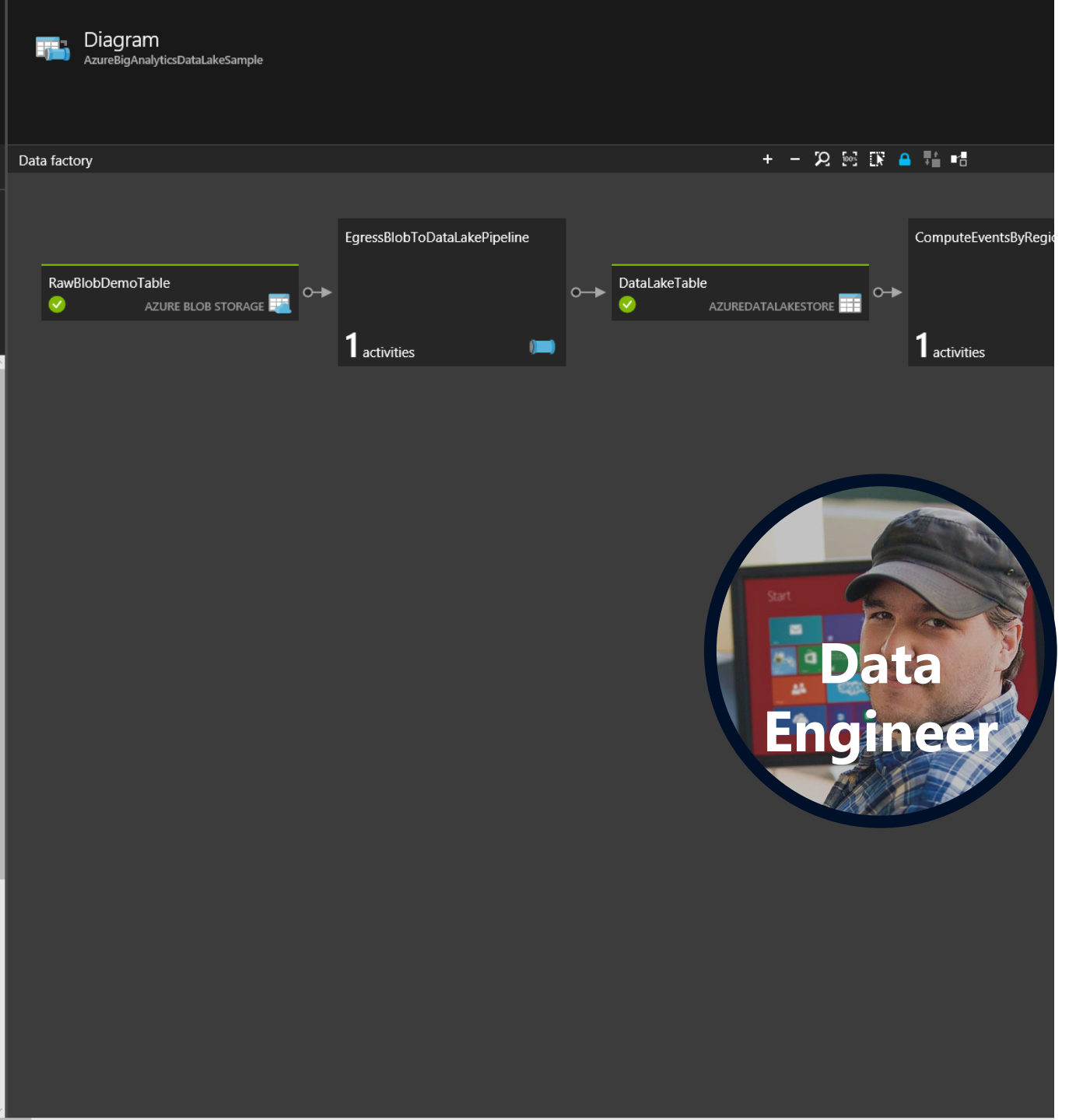
Operations

Events: AZUREBIGANALYTICSDATALAKESAMPLE

Alert rules: AZUREBIGANALYTICSDATALA...

Monitoring

Data factory metrics



Questions or
Comments?



Thank You!
darsch@microsoft.com

