# Microsoft Ignite
Australia 2017

# Orchestrating Big Data Pipelines with Azure Data Factory

DA332
Lace Lofranco

Microsoft

# Survey

Microsoft

# Session Objective

Learn to design, build and manage big data orchestration pipelines using Azure Data Factory

# Agenda

## Design
Big data pipelines
Lamda Architecture
Data Factory Concepts

## Build
Data Movement
Data Transformation

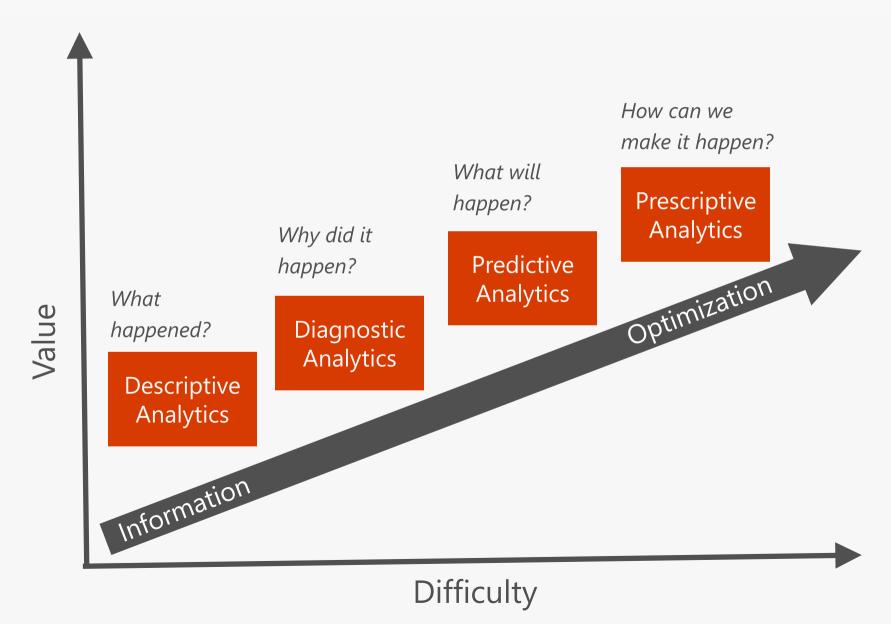## Manage
Monitor pipeline health
Developer tools
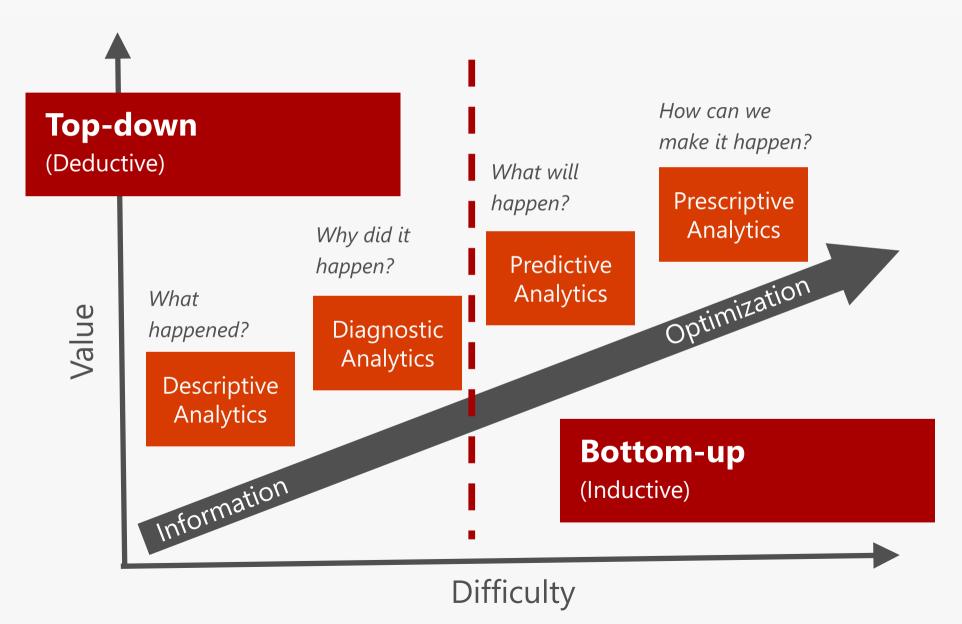
## Compare
Data Factory vs Oozie

Microsoft Ignite

Microsoft Ignite
Australia 2017

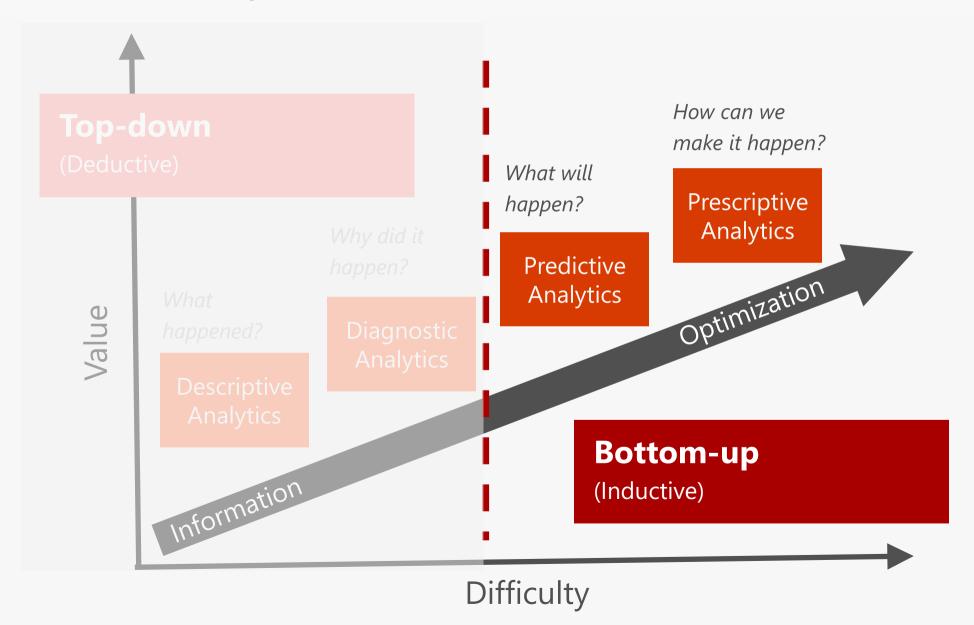# Design: Big Data Pipelines

Microsoft
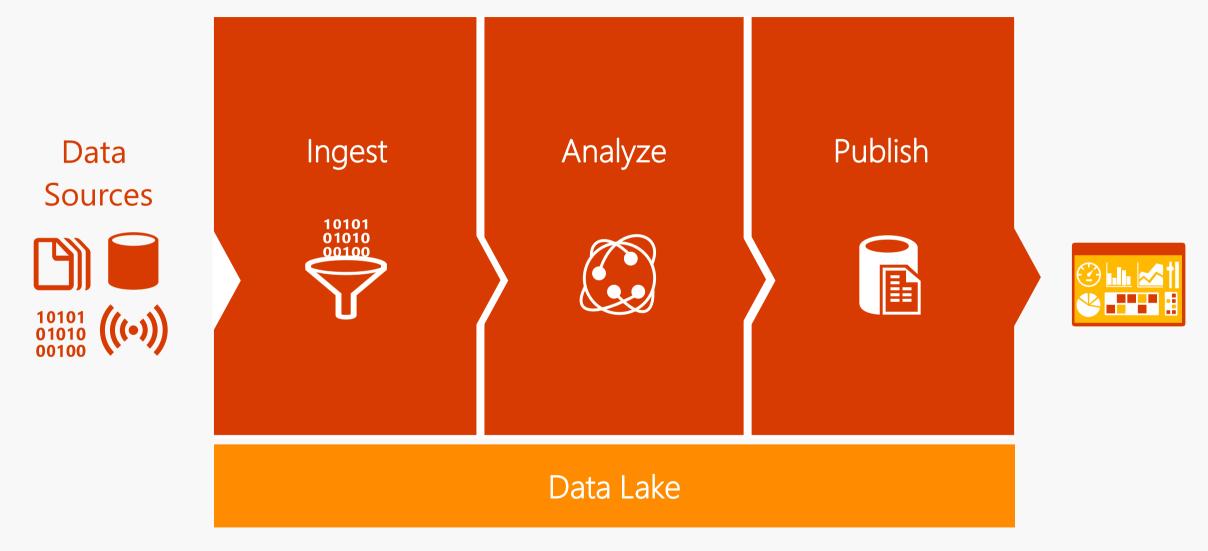
# Types of analytics

# Types of analytics

# Types of analytics

# Big Data Pipelines Examples
*Velocity, Variety, Volume*

- Optimizing Ad revenue
- Demand Forecasting
- Predictive Maintenance
- Supply chain optimization
- Portfolio optimization

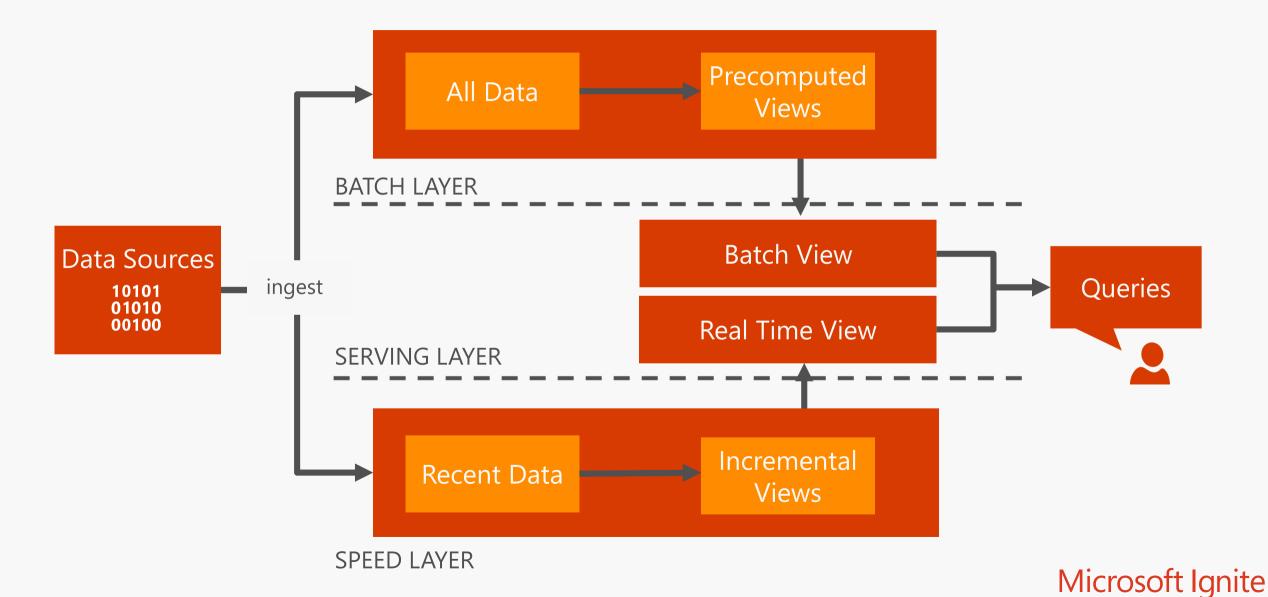Microsoft Ignite

# Big Data Pipelines

**Data Sources**

**Ingest**

**Analyze**

**Publish**
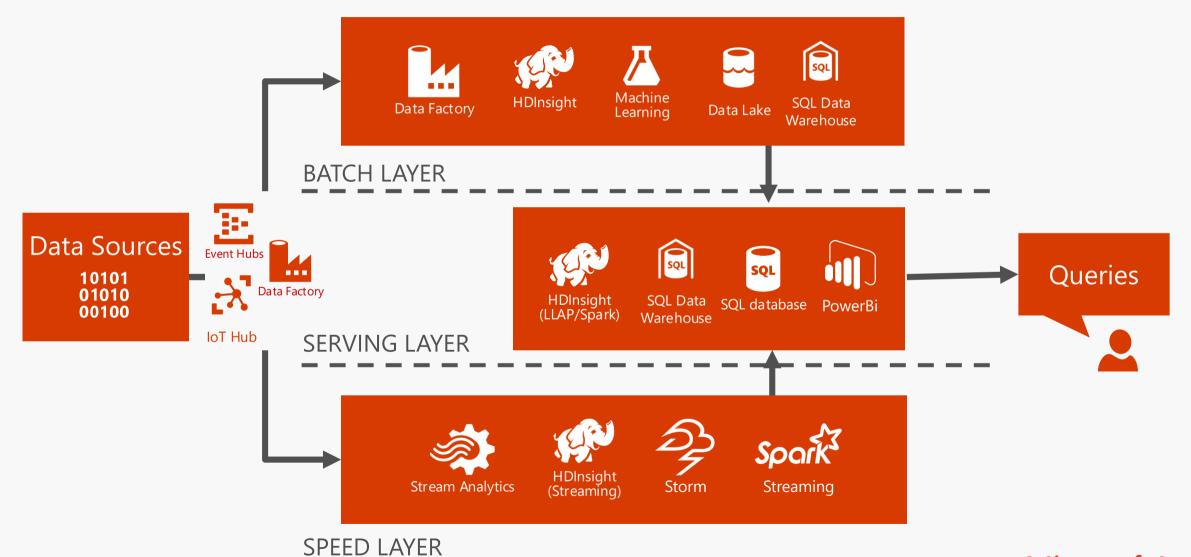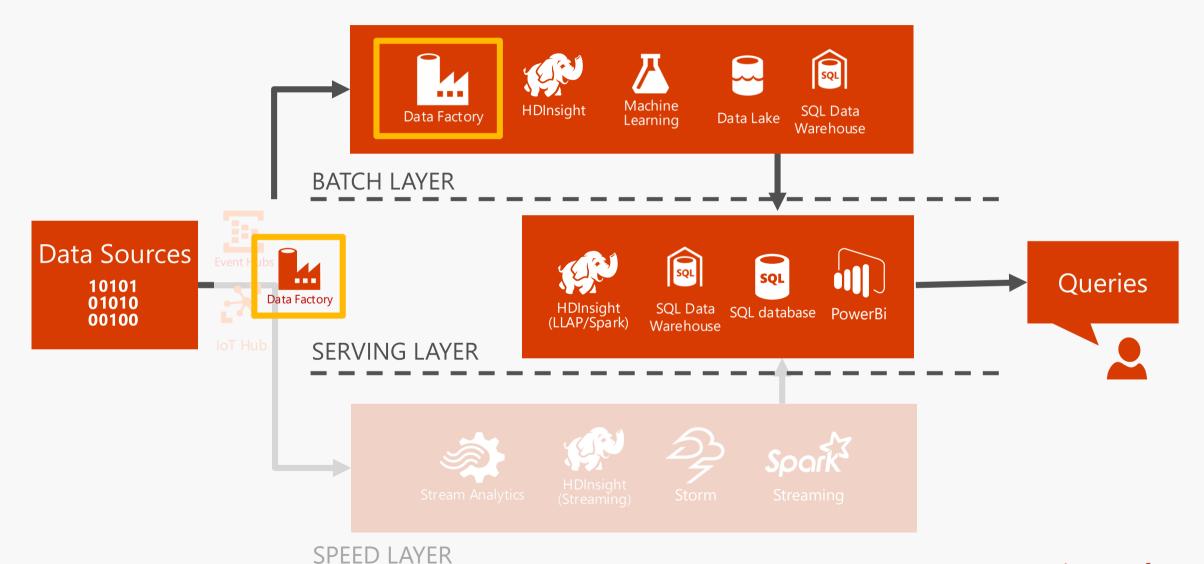
**Data Lake**

# Lamda Architecture

Data architecture for designing Big Data applications

Three layers: Batch, Speed, Serving

Popularized by Nathan Marz

# Lamda Architecture

**Data Sources**
10101
01010
00100

ingest

## BATCH LAYER

All Data → Precomputed Views

Batch View

Real Time View

Queries

## SERVING LAYER

Recent Data → Incremental Views

## SPEED LAYER

# Lamda Architecture

# Lamda Architecture



BATCH LAYER

Data Factory · HDInsight · Machine Learning · Data Lake · SQL Data Warehouse

Data Sources
10101
01010
00100

Event Hubs

Data Factory

IoT Hub

SERVING LAYER

HDInsight (LLAP/Spark) · SQL Data Warehouse · SQL database · PowerBi

Queries

SPEED LAYER

Stream Analytics · HDInsight (Streaming) · Storm · Spark Streaming

Microsoft Ignite

# Melbourne Foot Traffic

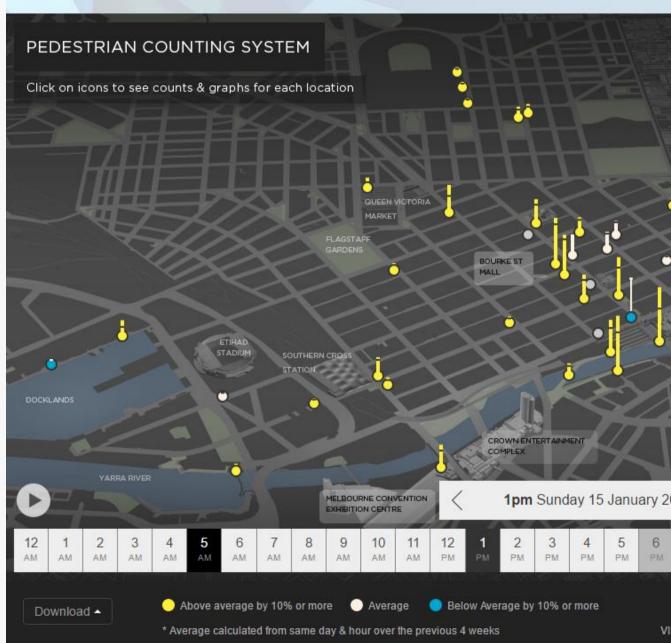Predicting **tram load** based on **foot traffic**

Model training data from yearly tram load survey conducted by Public Transport Victoria
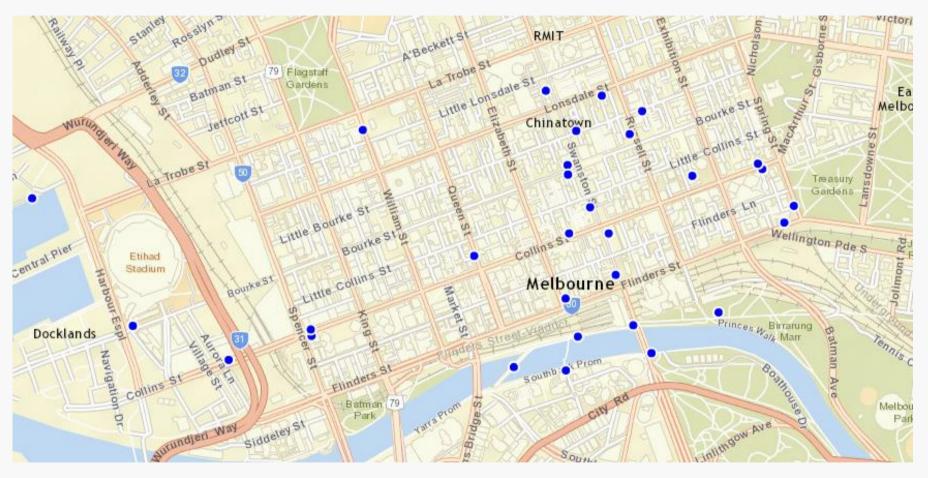
Data sources:

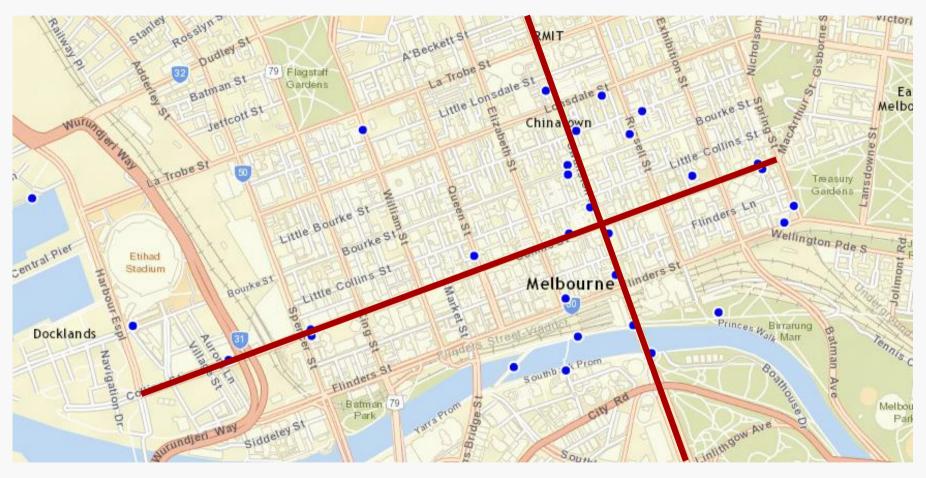https://data.melbourne.vic.gov.au

https://www.ptv.vic.gov.au

# Sensor locations



Source: https://data.melbourne.vic.gov.au

Microsoft Ignite

# Sensor locations



Source: https://data.melbourne.vic.gov.au

Microsoft Ignite

# Demo Solution Architecture

# Demo Solution Architecture*

HDInsight

Machine Learning

data.melbourne.vic.gov.au

HTTP

Batch

Storage blob

Data Lake

SQL DW

PowerBI

Data Lake Analytics

Data Factory

*Designed for the purpose of demoing Data Factory activities

Microsoft Ignite

# Data Factory Concepts

Data Store

# Data Factory Concepts

# Data Factory Concepts

# Data Factory Concepts

# Data Factory Concepts

# Data Factory Concepts

Data Store

Linked Service

Dataset

Activity

Dataset

Pipeline

Linked Service

Data Store

# Scheduling and Execution



Monthly — Source Dataset

Jan
Feb
Mar

Run 1
Run 2
Run 3

Activity

Monthly — Sink Dataset

Jan
Feb
Mar

# Scheduling and Execution



Monthly — Source Dataset
Jan
Feb
Mar

Run 1
Run 2
Run 3

Activity

Monthly — Sink Dataset
Jan
Feb
Mar

Hourly — Source Dataset
1
...
24

Run 1

Activity

Daily — Sink Dataset
Mon
Tues
Wed

Microsoft Ignite

# Demo: Copy Activity



data.melbourne.vic.gov.au

HTTP

Batch

Storage blob

Data Lake

HDInsight

Machine Learning

SQL DW

PowerBI

Data Lake Analytics

Data Factory

# Data Factory Design Considerations

- Ideal for time series data
- Design repeatable activity windows
- Handle 'late' / out-of-order runs
- Try to finalize pipeline schedules in advance
- No first class support for event driven pipelines and control flow activities
  - But with workarounds ☺

# Agenda

| Design | Build | Manage | Compare |
|--------|-------|--------|---------|
| Lamda Architecture<br>Intro to Data Factory<br>Considerations | Data Movement<br>Data Transformation | Monitor pipeline health<br>Developer tools | Data Factory vs Oozie |

Microsoft Ignite

Microsoft Ignite
Australia 2017

# Build: Data Movement

Microsoft

# Data Movement Activity: Copy

## Cloud to Cloud

| Source | WAN | Azure Data Factory | | | | WAN | Sink |
|--------|-----|-------------------|--|--|--|-----|------|
| | | Serialization-Deserialization | Compression | Column Mapping | ... | | |

## On-premises to Cloud

| Source | LAN/WAN | Data Management Gateway | | | | LAN/WAN | Sink |
|--------|---------|------------------------|--|--|--|---------|------|
| | | Serialization-Deserialization | Compression | Column Mapping | ... | | |

Microsoft Ignite

# Data Management Gateway

Client agent installed on-premises environment to copy data between cloud and on-premises data stores

# Supported Data Stores

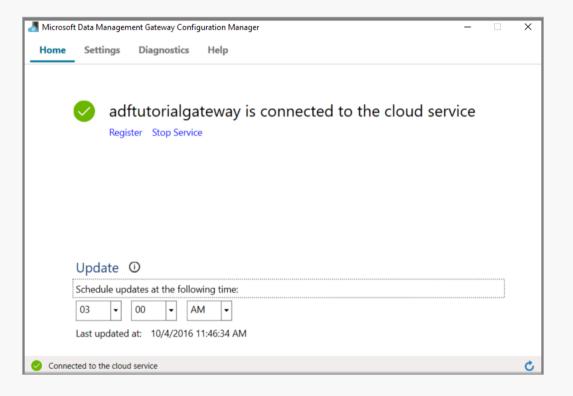| Category | Data store | Source | Sink |
|---|---|:---:|:---:|
| Azure | Azure Blob storage | ✓ | ✓ |
| | Azure Data Lake Store | ✓ | ✓ |
| | Azure SQL Database | ✓ | ✓ |
| | Azure SQL Data Warehouse | ✓ | ✓ |
| | Azure Table storage | ✓ | ✓ |
| | Azure DocumentDB | ✓ | ✓ |
| | Azure Search Index | ✓ | ✓ |
| Databases | SQL Server* | ✓ | ✓ |
| | Oracle* | ✓ | ✓ |
| | MySQL* | ✓ | |
| | DB2* | ✓ | |
| | Teradata* | ✓ | |
| | PostgreSQL* | ✓ | |
| | Sybase* | ✓ | |
| | Cassandra* | ✓ | |
| | MongoDB* | ✓ | |
| | Amazon Redshift | ✓ | |

| Category | Data store | Source | Sink |
|---|---|:---:|:---:|
| File | File System* | ✓ | ✓ |
| | HDFS* | ✓ | |
| | Amazon S3 | ✓ | |
| | FTP | ✓ | |
| Others | Salesforce | ✓ | |
| | Generic ODBC* | ✓ | |
| | Generic OData | ✓ | |
| | Web Table (table from HTML) | ✓ | |
| | GE Historian* | ✓ | |

Data stores with * can be on-premises or on Azure IaaS, and require you to install [Data Management Gateway](#) on an on-premises/Azure IaaS machine.

Microsoft Ignite

# Copy Wizard



Microsoft Ignite

# Considerations

- Copy service typically runs at region closest to sink
- Performance & Tuning
  - Use Parallel Copy if you have several small files
  - Increase Cloud data movement units (DMUs)
  - Staged copy and Compression
  - Column mapping and binary copy

Microsoft Ignite
Australia 2017

# Build: Data Transformation

Microsoft

# Data Transformation Activities

- Stored Procedure Activity
- HDInsight Activities
- Data Analytics U-SQL Activity
- Machine Learning Activities
- .NET Custom Activity

# Stored Procedure Activity

## Azure SQL Database

## Azure SQL Datawarehouse

Use Polybase for loading large datasets

## SQL Server Database (on-premises or IAAS)

Needs Data Management Gateway

# Polybase and CTAS

Polybase to access non-relational sources

'CREATE TABLE AS' is fully parallelized operation to create and load tables from a SELECT statements

Super-charge SELECT INTO statement

Microsoft Ignite

Microsoft Ignite
Australia 2017

# Demo: Stored Proc Activity

Loading data to SDW via Polybase

Microsoft

# Demo: Stored Procedure Activity



data.melbourne.vic.gov.au

HTTP

Batch

Storage blob

HDInsight

Machine Learning

Data Lake

SQL DW

PowerBI

Data Lake Analytics

Data Factory

# HDInsight Activities

**On-demand** or **Bring Your Own**

| Hive<br>SQL-on-hadoop | Pig<br>ETL / Scripting | | Spark<br>In-memory processing |
|---|---|---|---|
| Tez<br>Execution engine | | MapReduce<br>Java / Classic MR | |

YARN (Resource Manager)

HDFS

# Hive Activity Advanced Properties

| Property | Hadoop config |
| --- | --- |
| coreConfiguration | core-site.xml |
| hBaseConfiguration | hbase-site.xml |
| hdfsConfiguration | hdfs-site.xml |
| hiveConfiguration | hive-site.xml |
| mapReduceConfiguration | mapred-site.xml |
| oozieConfiguration | oozie-site.xml |
| stormConfiguration | storm-site.xml |
| yarnConfiguration | yarn-site.xml |

# Demo: Hive Activity



data.melbourne
.vic.gov.au

HTTP

Batch

Storage blob

HDInsight

Machine
Learning

Data Lake

Data Lake
Analytics

SQL DW

PowerBI

Data Factory

# U-SQL Activity

**Data Lake Analytics** is a distributed analytics service w/ federated access across Azure data stores

Comes with U-SQL
Big data query language that combines SQL and C#

Pay per job (no cluster needed)

# Demo: U-SQL Activity



data.melbourne.vic.gov.au

HTTP

Batch

Storage blob

HDInsight

Machine Learning

Data Lake

Data Lake Analytics

SQL DW

PowerBI

Data Factory

# AzureML Activities

1. ## Batch Execution Activity
   Scoring datasets
   Retraining models

2. ## Update Resource Activity
   Updating ML models

# Batch Execution Activity

1. ## Use Web service inputs/outputs
   Pass datasets as inputs / outputs of web service
   Need to specify single file name as inputs

2. ## Use AzureML import/export modules
   Pass parameters (globalParameters) to specify reader/writer properties
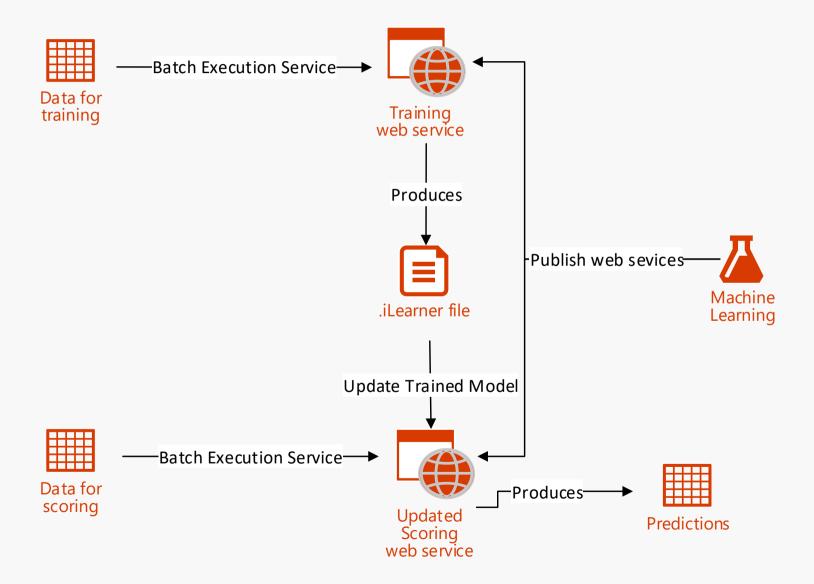   Can pass a directory as input data to readers

Microsoft Ignite

# Retraining your model

# Retraining your model



Batch Execution Activity

Data for training → Batch Execution Service → Training web service

Produces → .iLearner file

Publish web sevices → Machine Learning

Update Trained Model

Data for scoring → Batch Execution Service → Updated Scoring web service → Produces → Predictions

# Retraining your model



**Update Resource Activity**

Data for
training

Batch Execution Service →

Training
web service

Produces

.iLearner file

Publish web sevices

Machine
Learning

Data for
scoring

Batch Execution Service →

Update Trained Model

Updated
Scoring
web service

Produces →

Predictions

Microsoft Ignite

# Demo: AzureML Batch Scoring Activity

# .NET Custom Activity

Use if data source/sinks not supported by ADF
Compute: **Azure Batch** or **HDInsight**

# .NET Custom Activity

```csharp
public class ApiDownloadActivity : IDotNetActivity
{
        public IDictionary<string, string> Execute(
                IEnumerable<LinkedService> linkedServices,
                IEnumerable<Dataset> datasets,
                Activity activity,
                IActivityLogger logger)
        {
            //your code

        }
}
```
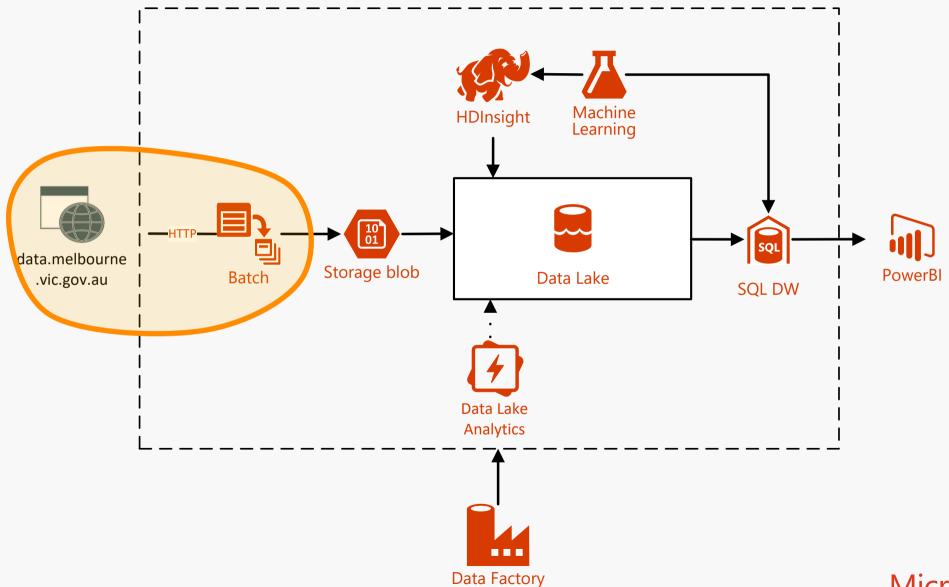
Microsoft Ignite
Australia 2017

# Demo: .NET Custom Activity

Calling an API with Data Factory

Microsoft

# Demo: Custom Activity on Batch



data.melbourne
.vic.gov.au

HTTP

Batch

Storage blob

HDInsight

Machine Learning

Data Lake

SQL DW

PowerBI

Data Lake Analytics

Data Factory

# .NET Custom Activity on Batch



Microsoft Ignite

Microsoft Ignite
Australia 2017

# PowerBI Dashboard

Predicting tram load from pedestrian traffic

Microsoft

# Melbourne Pedestrian Foot Traffic



Microsoft Ignite

# Tram Load Predictions

location
- ☐ Collins St West End (Southern Cross Station)
- ☑ Swanston St / Flinders St (Federation Sq)

**0.58**
Average of scored_prob

hourly_counts and scored_prob by hour
- hourly_counts  ● scored_prob

scored_prob by Day Name

Microsoft Ignite

# Agenda

| Design | Build | Manage | Compare |
|--------|-------|--------|---------|

Big data pipelines
Lamda Architecture
Data Factory Concepts

Data Movement
Data Transformation

Monitor pipeline health
Developer tools

Data Factory vs Oozie

Microsoft Ignite

# Monitor and Manage Dashboard

Pause and resuming pipelines

Creating alerts

Re-running failed pipelines

# Demo: Monitor and Manage

Microsoft

# Developer tools

Azure Portal

VS Data Factory Extension

PowerShell

Azure Resource Manager Templates

C# SDK

# Agenda

| Design | Build | Manage | Compare |
|---|---|---|---|

Big data pipelines

Lamda Architecture

Data Factory Concepts

Data Movement

Data Transformation

Monitor pipeline health

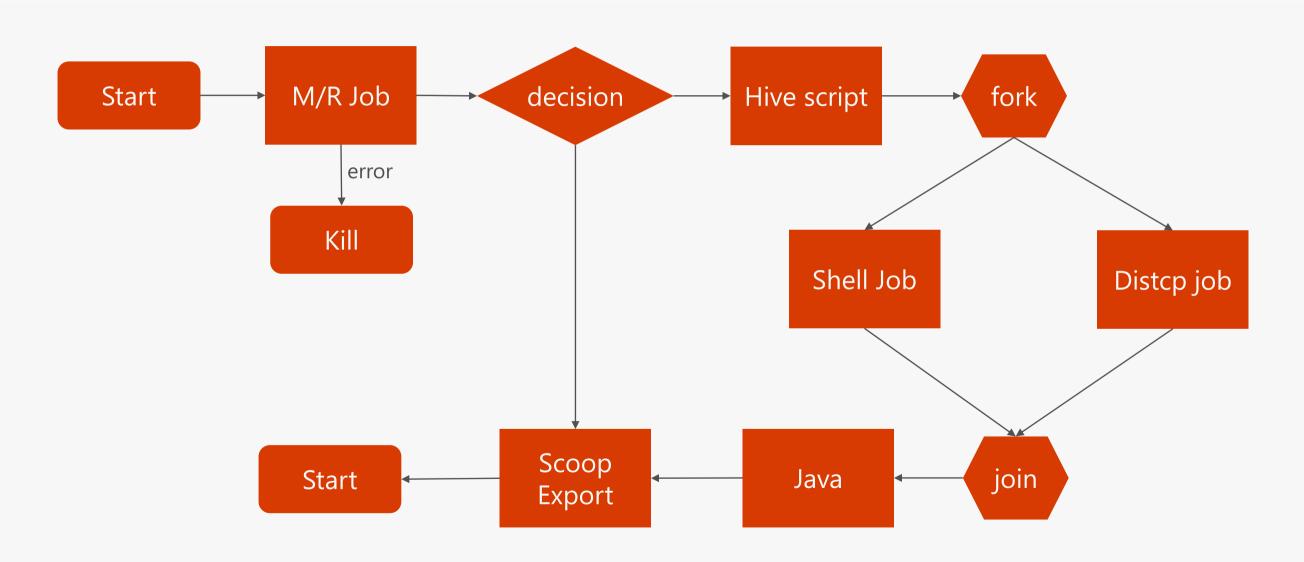Developer tools

Data Factory vs Oozie

# What is Oozie?

Oozie is a workflow scheduler system to manage Apache Hadoop jobs

De-facto workflow scheduler of the Hadoop stack

Out-of-the-box integration with Hadoop jobs

Java map-reduce, Streaming map-reduce, Pig, Hive, Sqoop and shell

# Sample Oozie workflow

Start → M/R Job → decision → Hive script → fork

M/R Job → (error) → Kill

decision → Scoop Export

fork → Shell Job

fork → Distcp job

Shell Job → join

Distcp job → join

join → Java → Scoop Export → Start

# Data Factory vs Oozie

| | Data Factory | Oozie |
|---|---|---|
| **Sources and Sinks** | Azure<br>Hadoop Stack<br>Numerous 3rd party | Hadoop Stack<br>Relational through Sqoop<br>Shell, Distcp |
| **Hybrid pipelines** | Data Management Gateway | None |
| **Tooling** | Visual Studio, Portal, PS | Hue, Eclipse plugin, 3rd party |
| **Extensibility** | C# Custom Activity | Java Action Node |
| **Control flow** | Limited | Fork and Join, Decision Control, Kill (on error) |
| **Event-based** | No first class support. Some workarounds | Dataset polling |
| **Performance** | Designed for big data workloads | Designed for big data workloads |
| **Maturity** | Preview in Oct 2014, GA in August 2015 | Developed since 2008, Open sourced since 2010 |

# Overall Summary

- Data Factory

...can orchestrate data pipelines at scale

...has tight integration with Azure's big data PAAS offerings and with variety of 3$^{rd}$ party source systems

...offers data Movement and Data Transformation activities

...first choice orchestrator for **Azure** services

# Session evaluation

Complete your session evaluation on MyIgnite
for your chance to **WIN** one of many daily prizes.

(image of prizes tbc)

# Continue your Ignite learning path

Visit Channel 9 to access a wide range of Microsoft training
and event recordings https://channel9.msdn.com/

Head to the TechNet Eval Centre to download trials of the latest
Microsoft products http://Microsoft.com/en-us/evalcenter/

Visit Microsoft Virtual Academy for free online training visit
https://www.microsoftvirtualacademy.com

Microsoft Ignite

# Thank you

Chat with me in the Speaker Lounge
Find me @ linkedin.com/in/lacelofranco/

Microsoft