



# An Architects Guide to Delivering Data Insights

Using the Microsoft Azure Data Platform

Paul Andrew | Group Manager & Analytics Architect





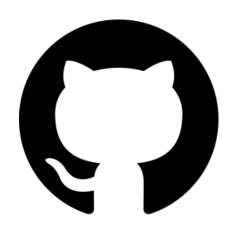












### https://github.com/mrpaulandrew

### ${\color{red}\textbf{Community}} \textbf{Events}$

Demo code, content and slides from various community events.

C++

{Event/Location}-{Month}-{Year}





### Question:

What is the answer to life, the universe and everything?

Answer:

42



Answer:







### Question:

What is big data?

### **Answer:**

It depends!



### **Answer:**

Any data that you cannot process in the time that you have/want using the technology you have.

- Buck Woody



### Goal





Paul's Magic Box -From the Hogwarts School of Witches & Wizardry



Data Sources Data Warehouse Data Insights

Data = Information = Knowledge = Power

### Goal





Clean Enrich Conform Translate Transform Curate Analyse Model Predict Master



Data Sources Data Warehouse

Data Insights



- Disaster recovery
- Transaction level restart ability

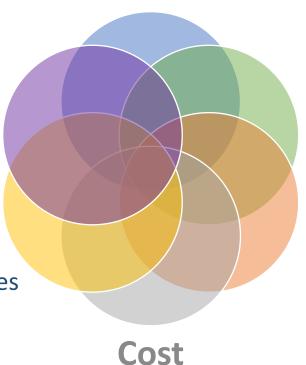
### Resilience

### **Rapid Delivery**

- Metadata driven
- Continuous deployments

### (Re)Usability

- Generic code libraries
- Dataset contracts



#### **Performance**

- Complex partitioning
- Large compute clusters

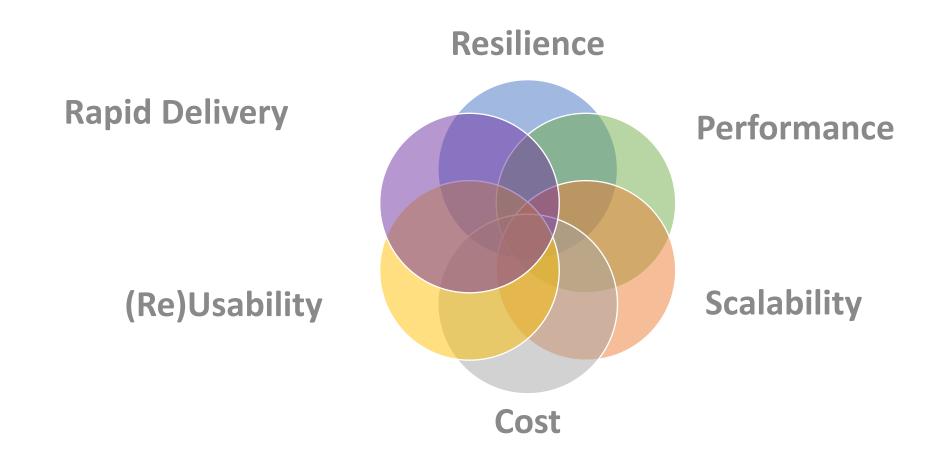
### **Scalability**

- Auto scaling microservices
- Event driven discrete batches

- Minimum resources used
- Dynamic resource management

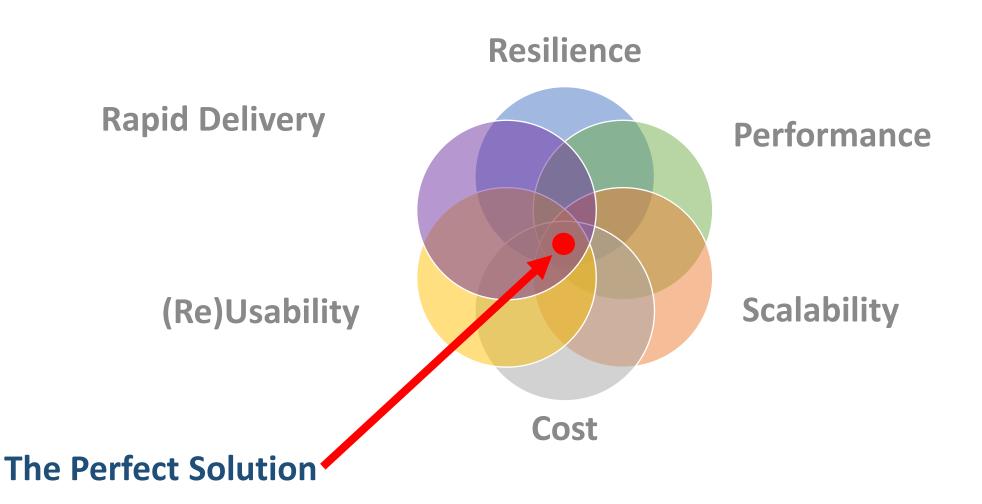






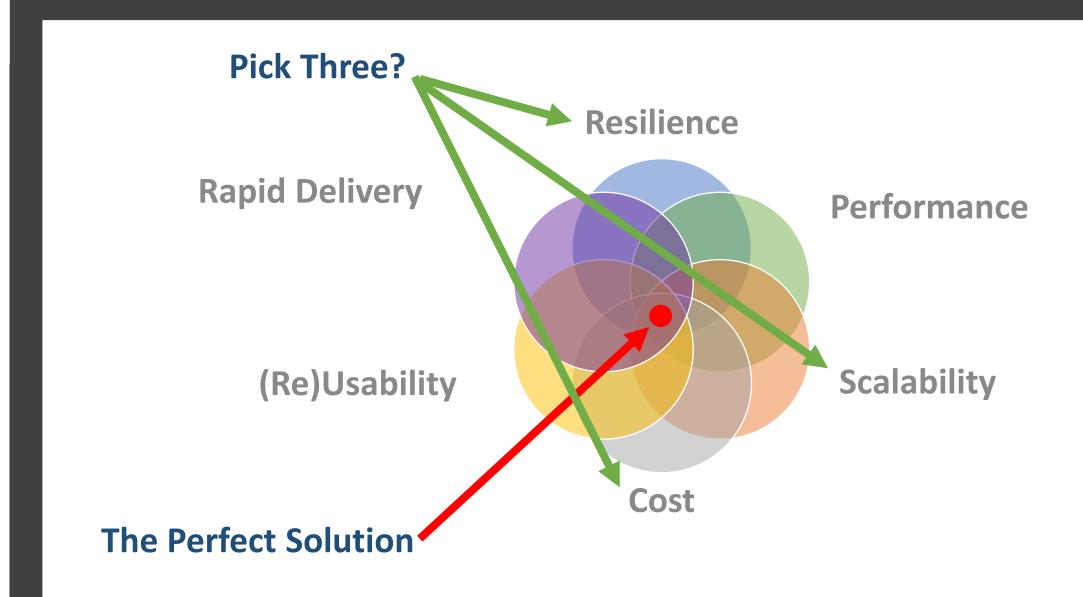




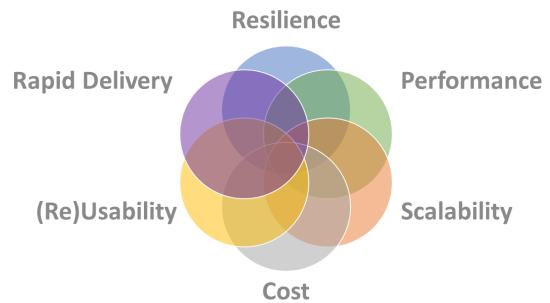


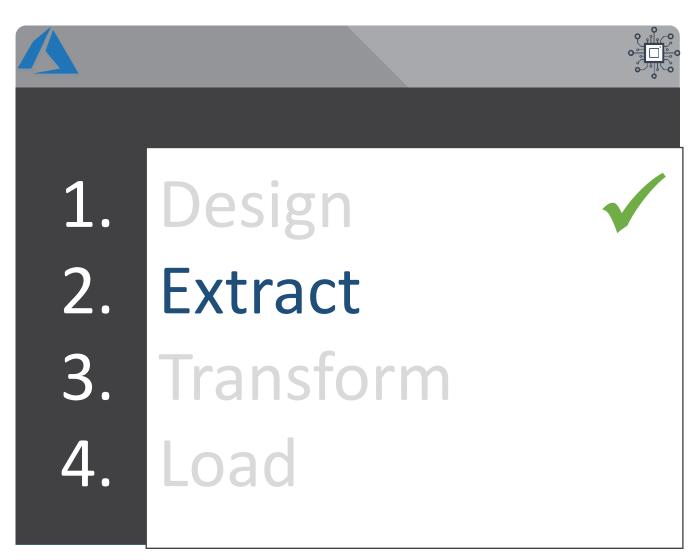


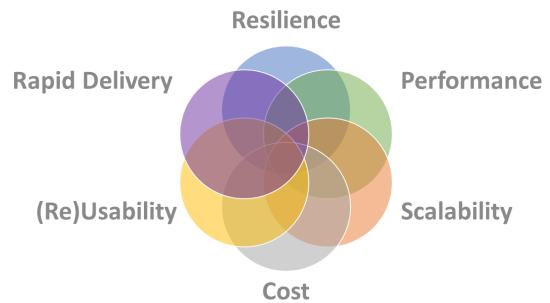














# Data Extraction & Ingestion







**Data Source** 



Push or Pull











Batch or Speed











Public or **Private Transfer** 







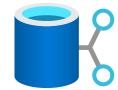




### Data Sensitivity











#### Data Volume













# Data Extraction & Ingestion – Spec v1

















Push or Pull











Batch or Speed











Public or Private Transfer







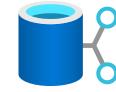




### Data Sensitivity









### Data Volume



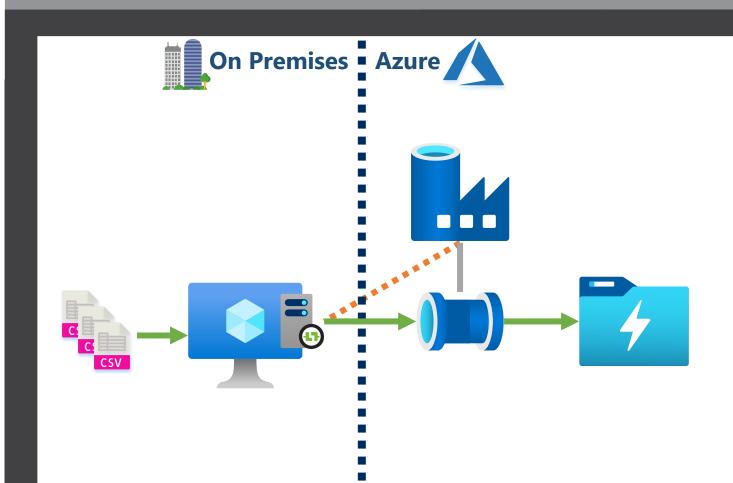






# Data Extraction & Ingestion – Solution 1





#### Requirements:

- Flat files
- From local storage
- Pulled from source
- Batch load
- Public connections
- No PII data
- Small data volumes



CSV

# Data Extraction & Ingestion – Spec v2

















### Push or Pull











Batch or Speed











Public or Private Transfer







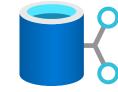




### Data Sensitivity

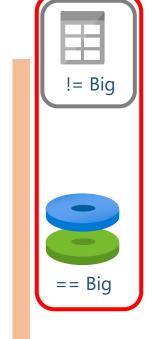








### Data Volume

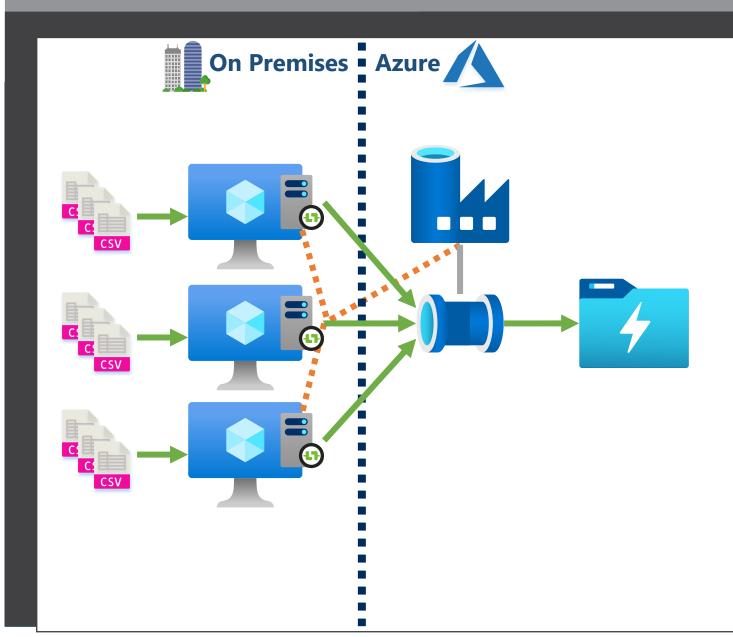






## Data Extraction & Ingestion – Solution 2





#### Requirements:

- Flat files
- From local storage
- Pulled from source
- Batch load
- Public connections
- No PII data
- <u>Large</u> data volumes



# Data Extraction & Ingestion – Spec v3

















Push or Pull











Batch or Speed



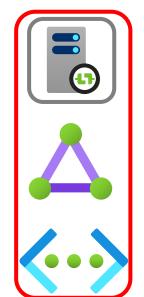








Public or Private Transfer







### Data Sensitivity

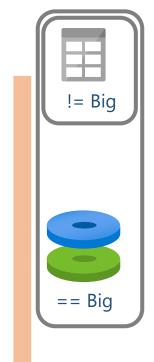








#### Data Volume

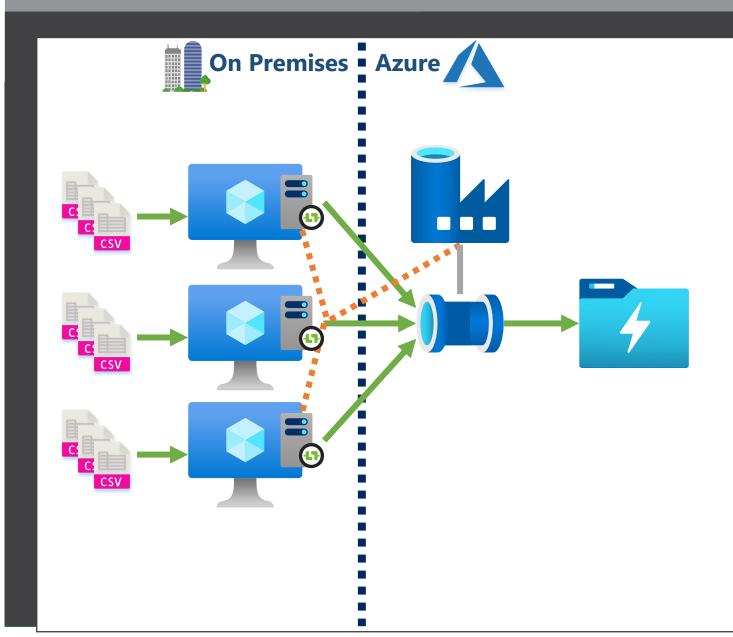






## Data Extraction & Ingestion – Solution 3





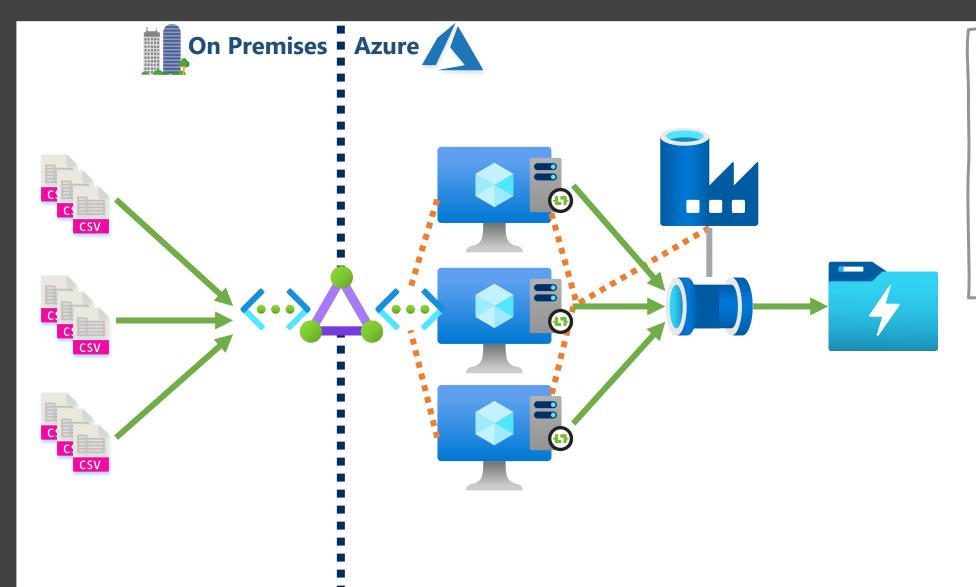
#### Requirements:

- Flat files
- From local storage
- Pulled from source
- Batch load
- Private connections
- No PII data
- Large data volumes



## Data Extraction & Ingestion – Solution 3





#### Requirements:

- Flat files
- From local storage
- Pulled from source
- Batch load
- <u>Private</u> connections
- No PII data
- Large data volumes



# Data Extraction & Ingestion – Spec v4







Data Source



Push or Pull











Batch or Speed



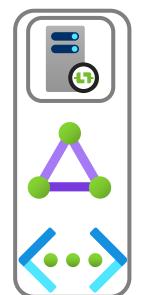








Public or Private Transfer



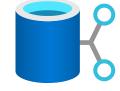




### Data Sensitivity

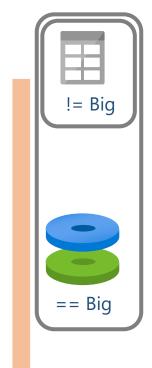








#### Data Volume

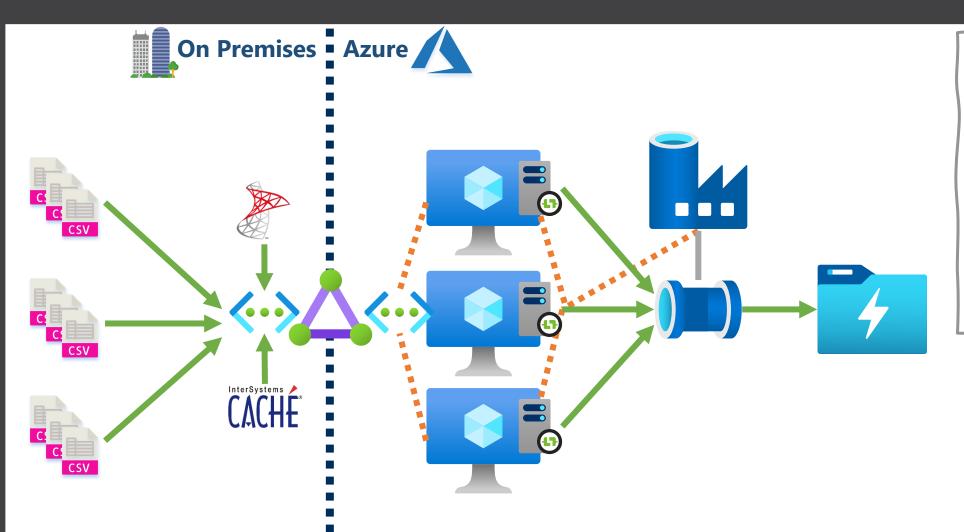






## Data Extraction & Ingestion – Solution 4





#### Requirements:

- Flat files
- From local storage& database tables
- Pulled from source
- Batch load
- Private connections
- No PII data
- Large data volumes



# Data Extraction & Ingestion – Spec v5







Data Source



Push or Pull











Batch or Speed



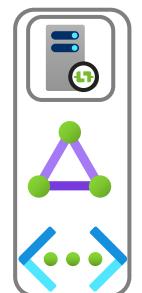








Public or Private Transfer



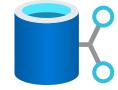




### Data Sensitivity

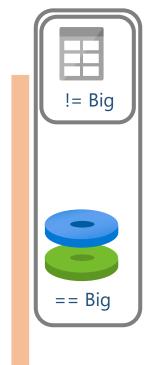








### Data Volume

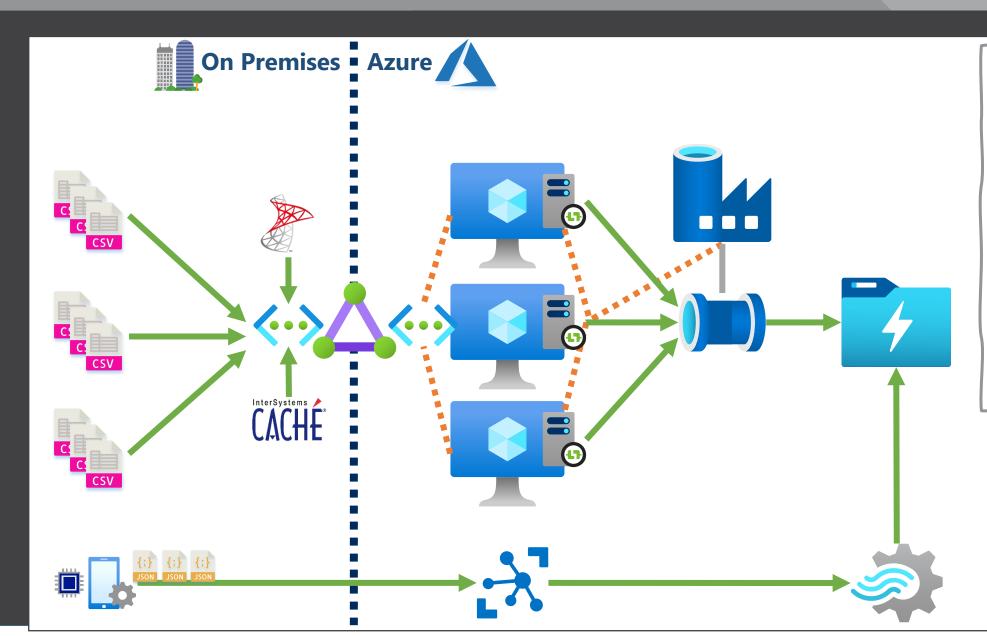






### Data Extraction & Ingestion – Solution 5





#### Requirements:

- Flat files & JSON
- From local storage& database tables
- Pulled from source& pushed
- Batch load & streamed
- Private connections
- No PII data
- Large data volumes



## Data Extraction & Ingestion – Spec v6







Data Source



Push or Pull











Batch or Speed



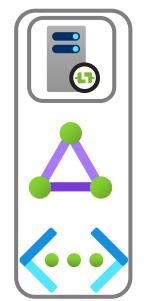




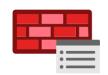




Public or Private Transfer



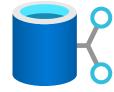




### Data Sensitivity

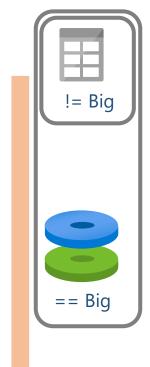








#### Data Volume

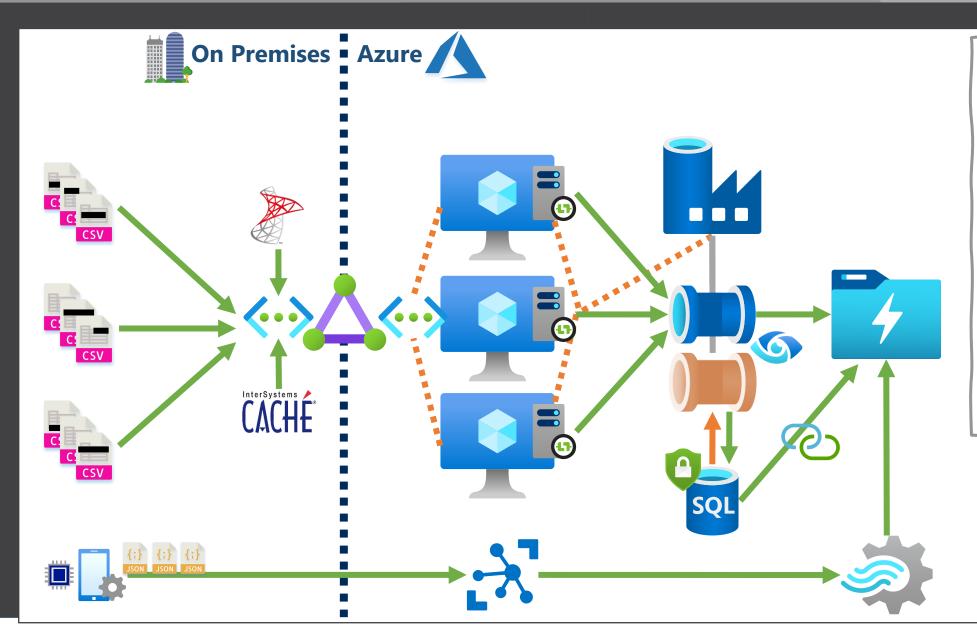






## Data Extraction & Ingestion – Solution 6





#### Requirements:

- Flat files & JSON
- From local storage& database tables
- Pulled from source& pushed
- Batch load & streamed
- Private connections
- Both PII & none
  PII data
- Large data volumes



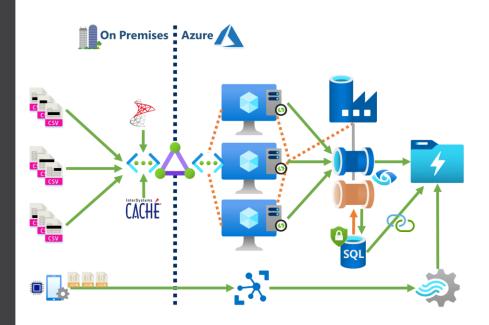
## Overall Architecture

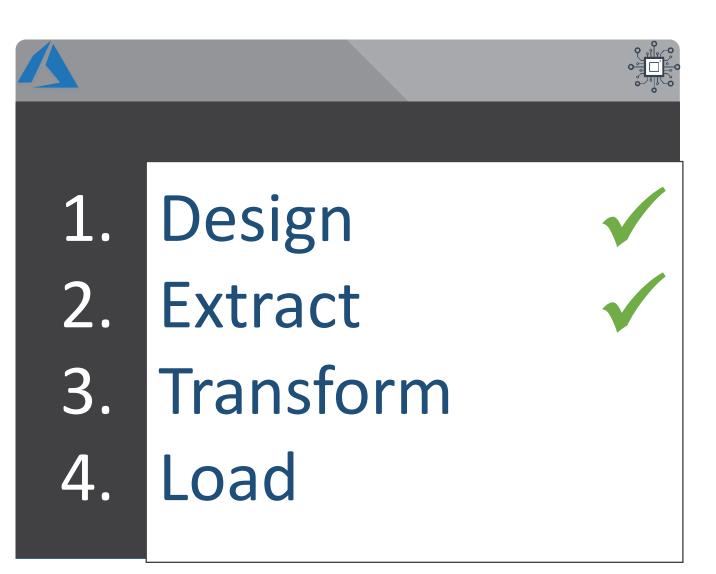


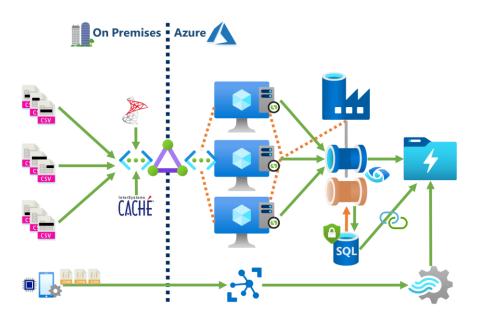
# Extract

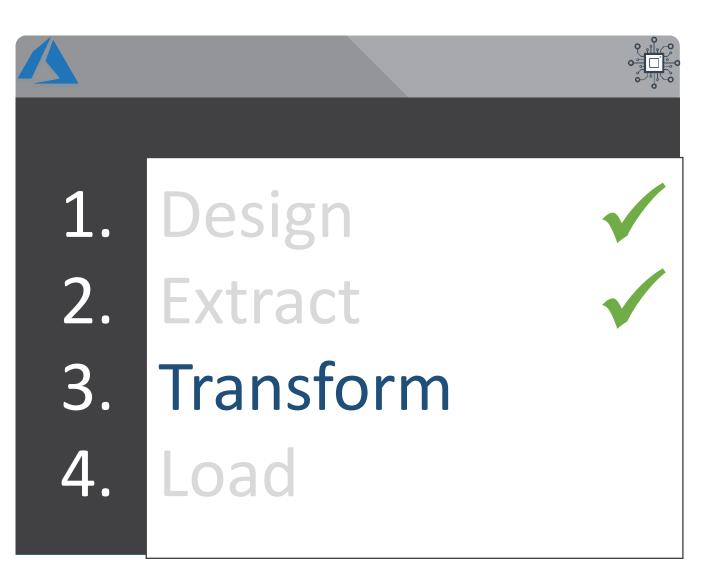
# Transform

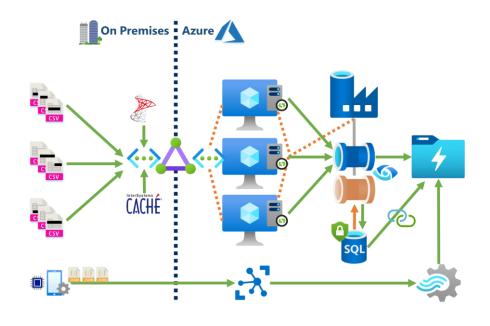
Load

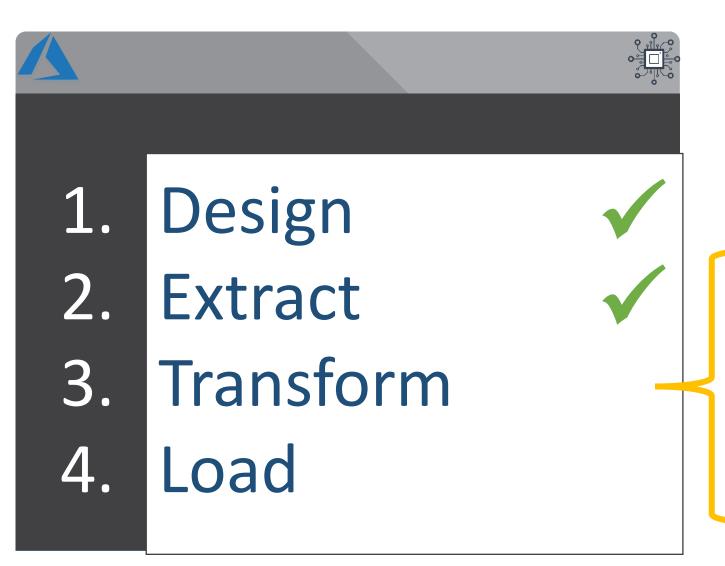








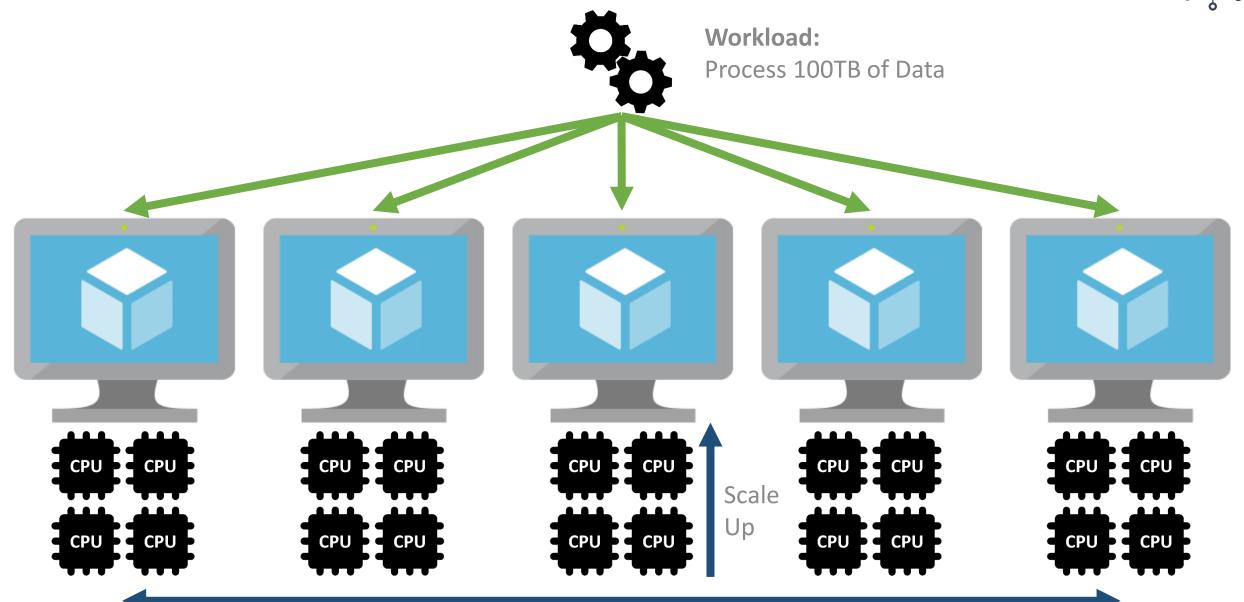




Compute
Storage, Structure
& Data Format

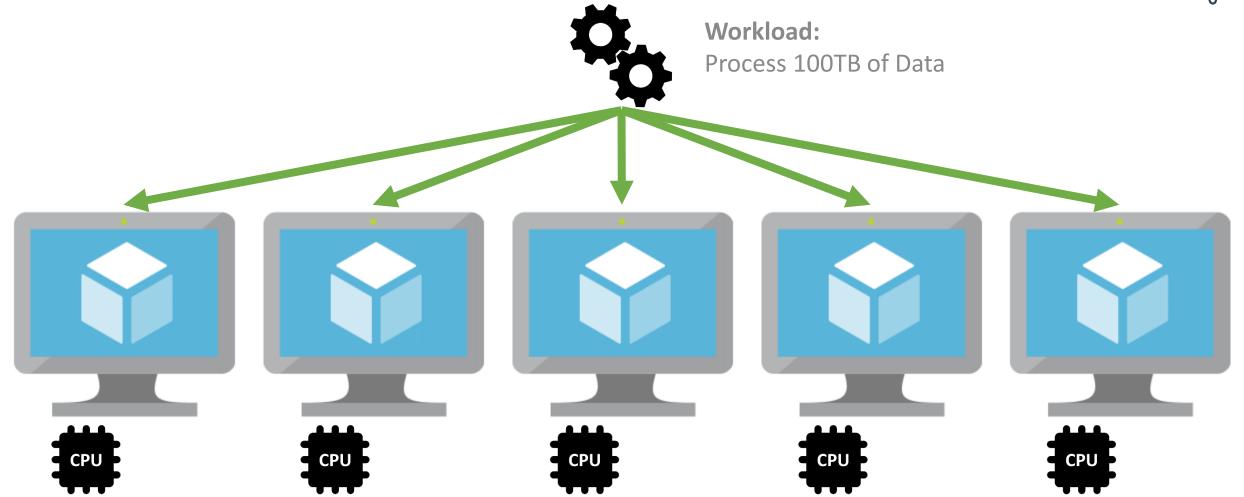
### Scaling Up and/or Scaling Out





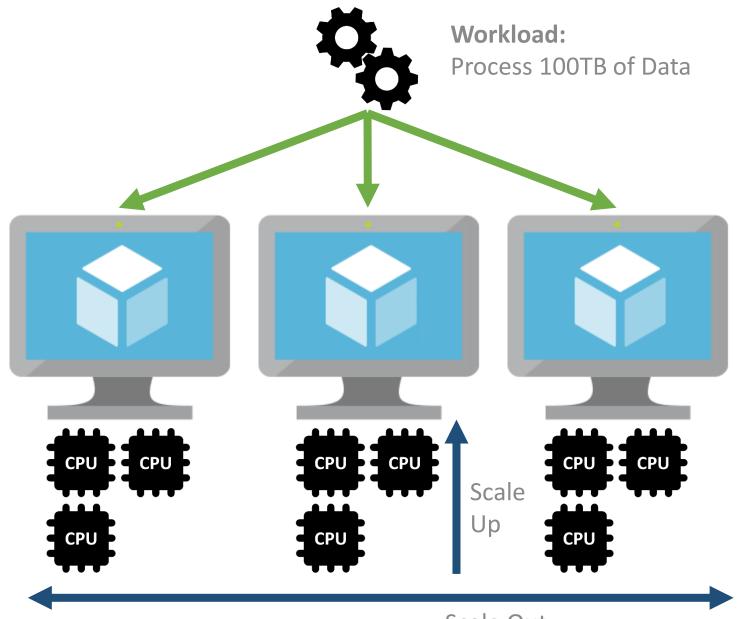
# Scaling Up and/or Scaling Out





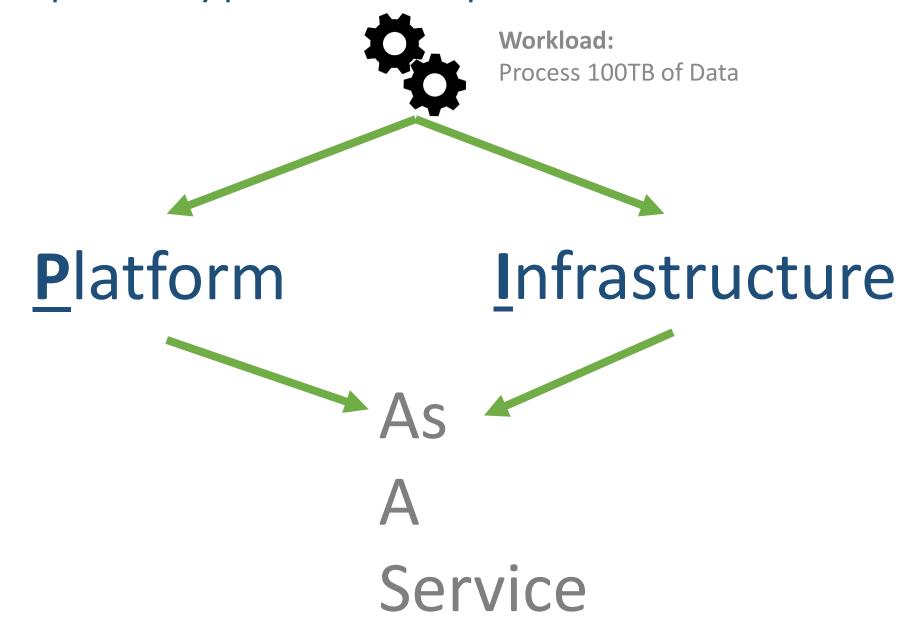
## Scaling Up and/or Scaling Out





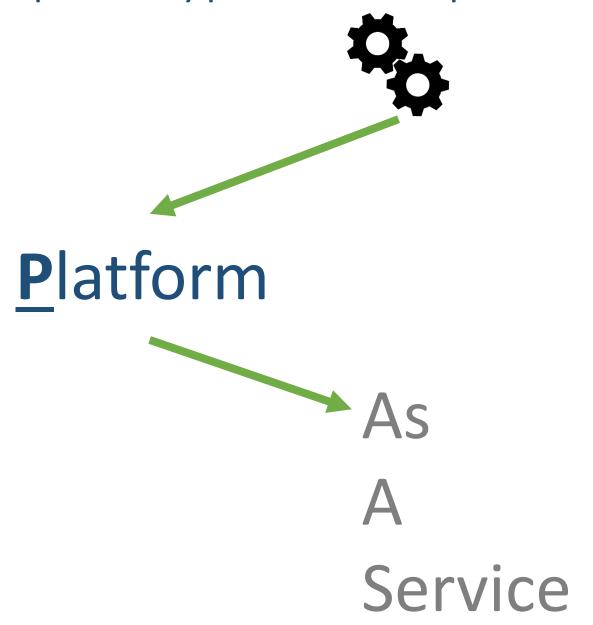
## What Compute Type of Compute?





#### What Compute Type of Compute?







#### Data Transformation – Compute



Data Lake **Analytics** 

**HDInsight** 

Relational Database

Synapse – **SQL** Pools or **Spark Pools** 

**Databricks** 

**Batch Service** 

Data Explorer















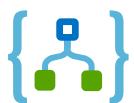
**Automation** 



**Functions** 

Power BI Data Flows

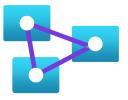




Data Factory **Data Flows** 









Cosmos







#### Data Transformation – Compute



Data Lake Analytics

HDInsight

Relational Database Synapse – SQL Pools or Spark Pools

Databricks

Batch Service

Data Explorer















Automation











Analysis Services



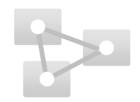




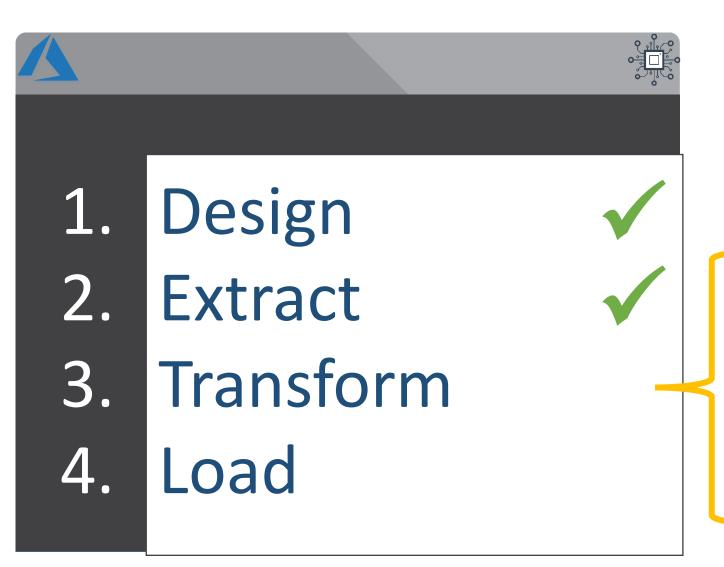








# Agenda



Compute 

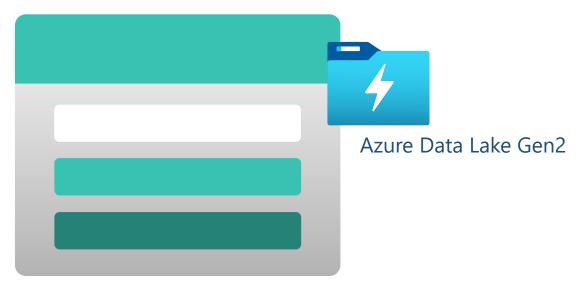
Storage, Structure

& Data Format





Azure Storage Account



Hadoop Distributed File System (HDFS)





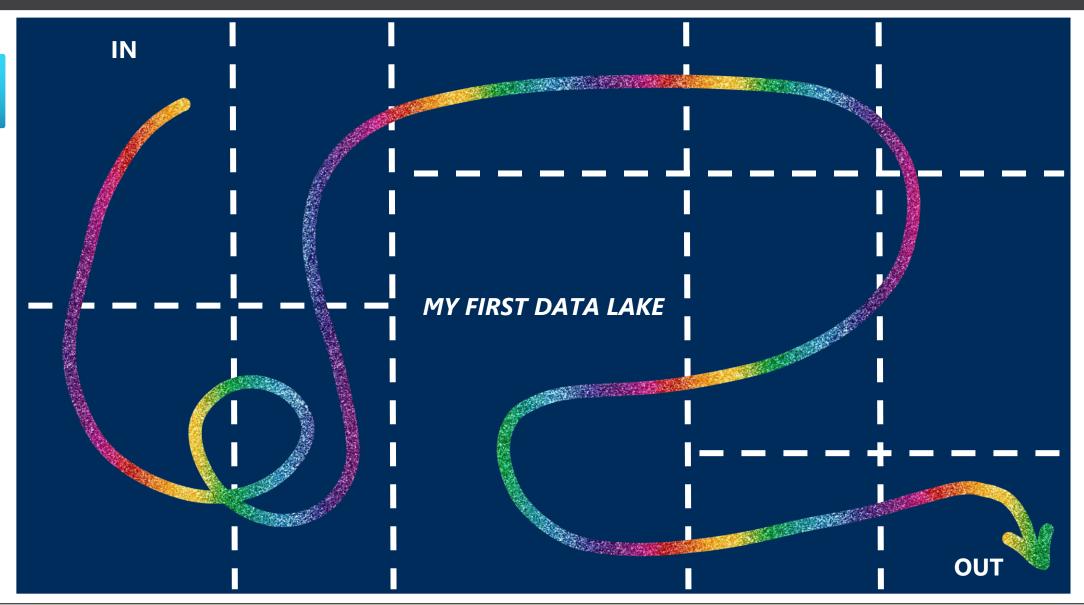






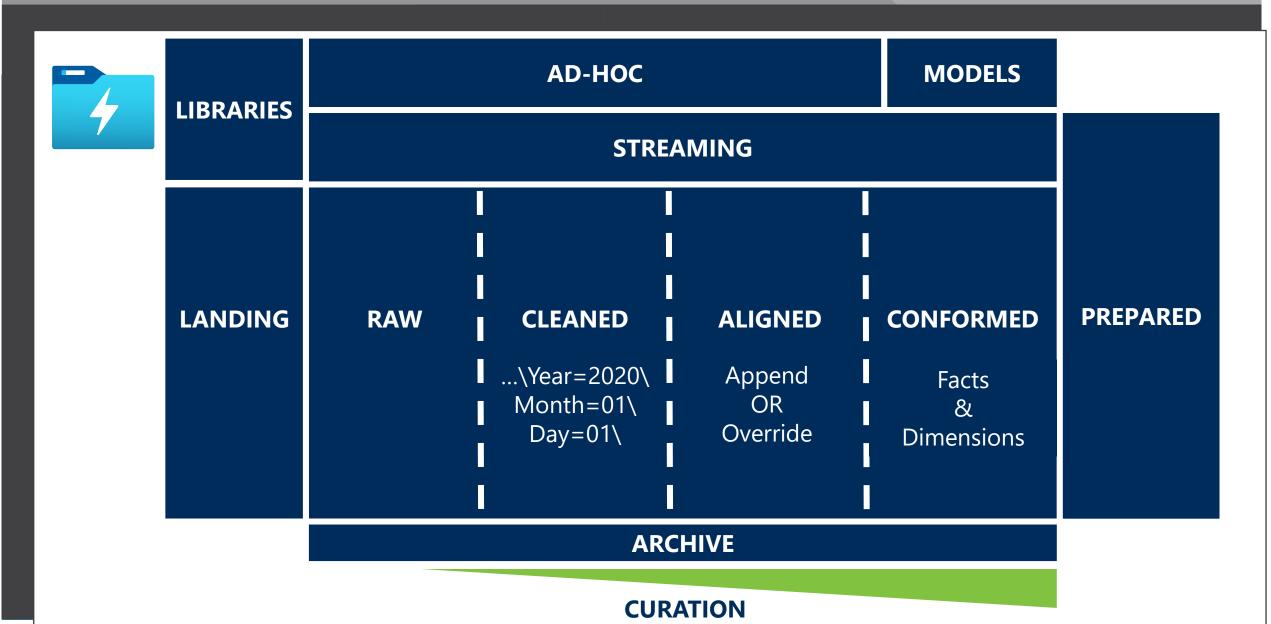






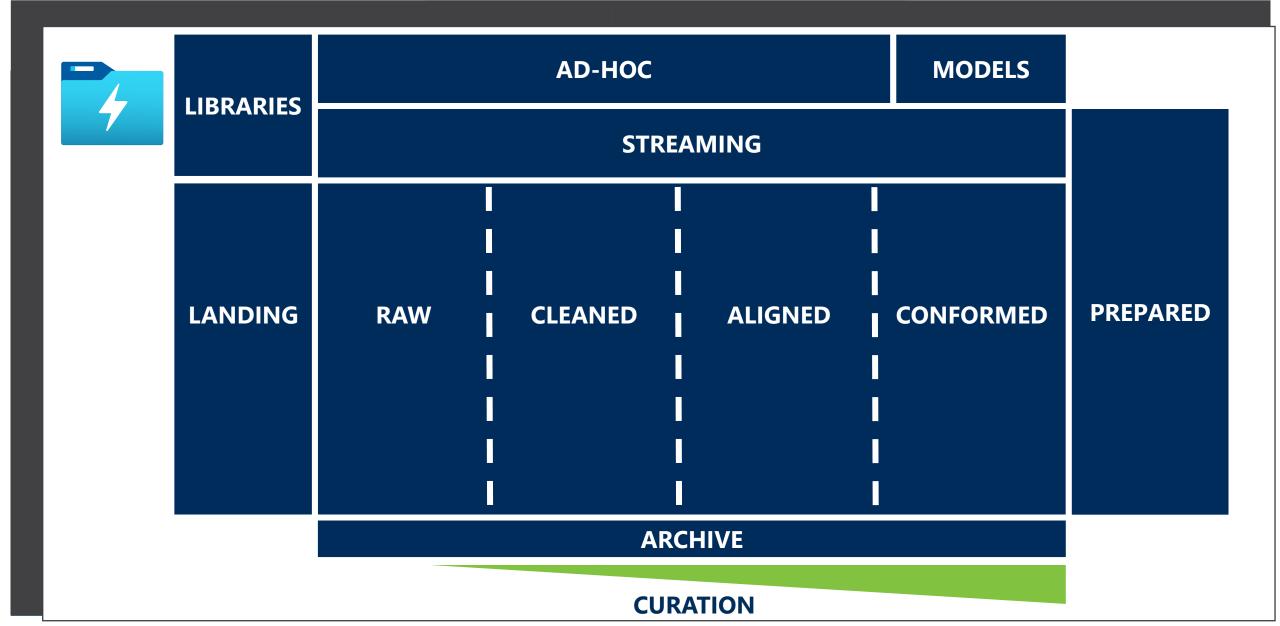






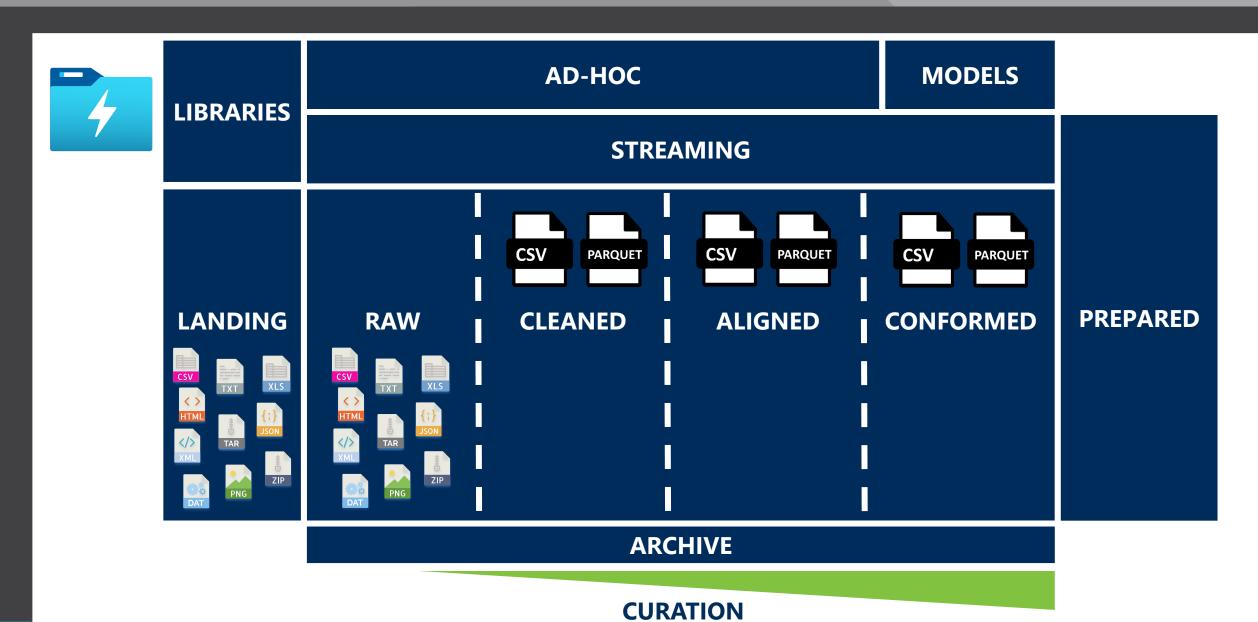






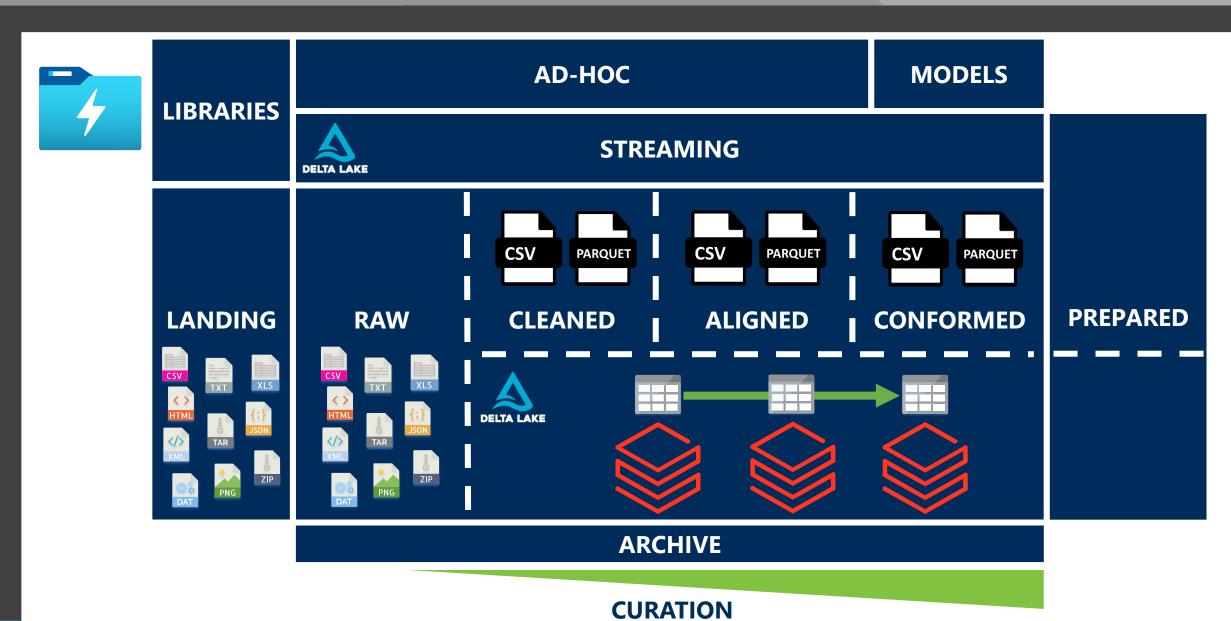




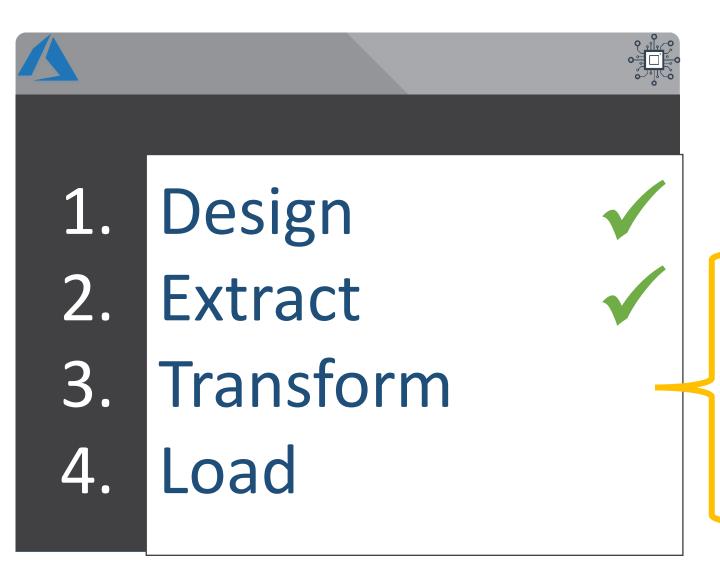








# Agenda



Compute 

Storage, Structure

& Data Format

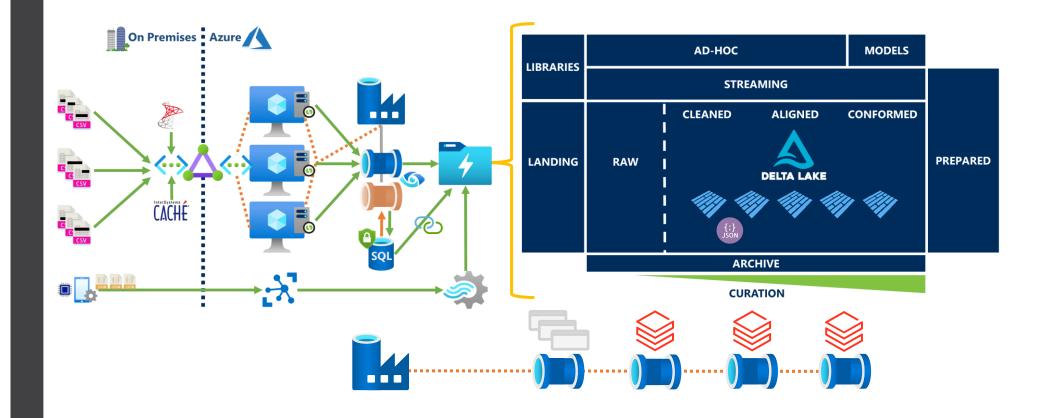




## Extract

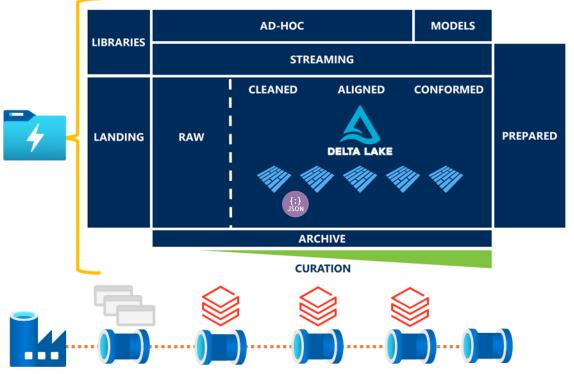
# Transform

Load



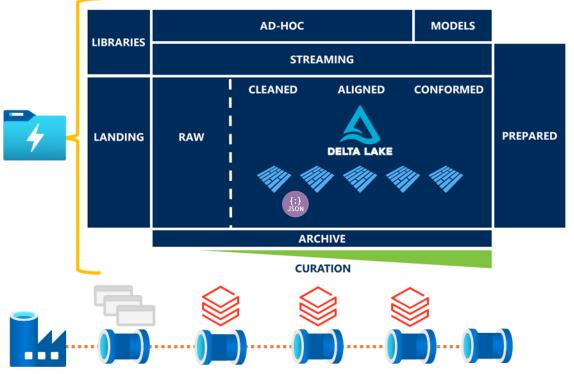
# Agenda





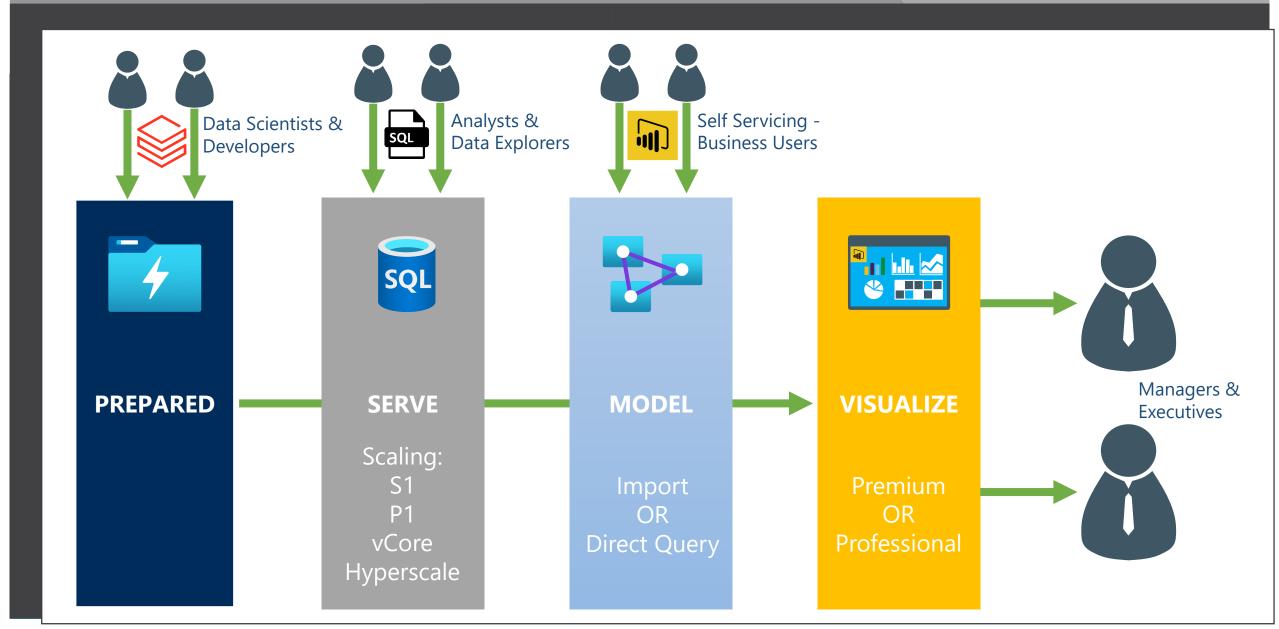
# Agenda





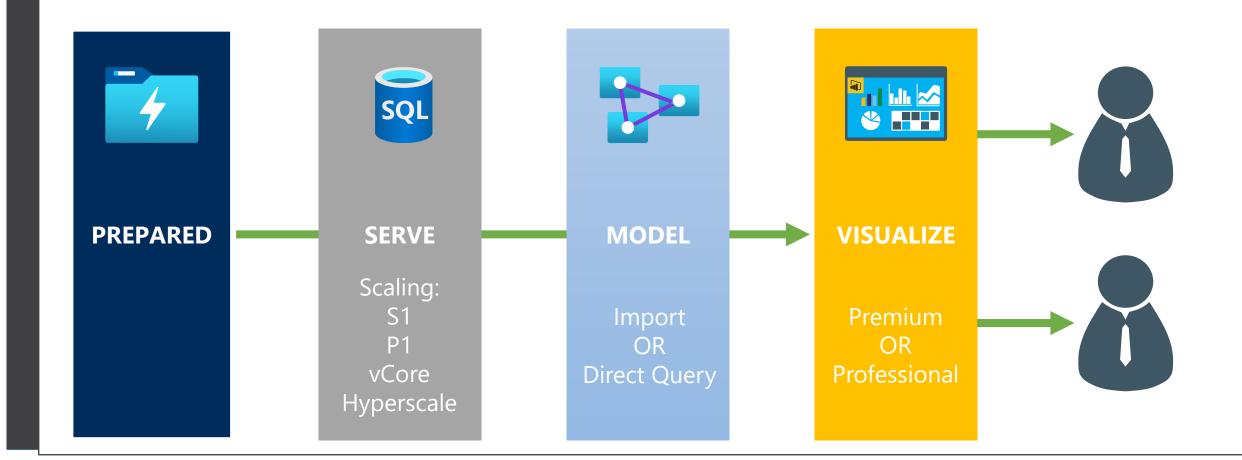






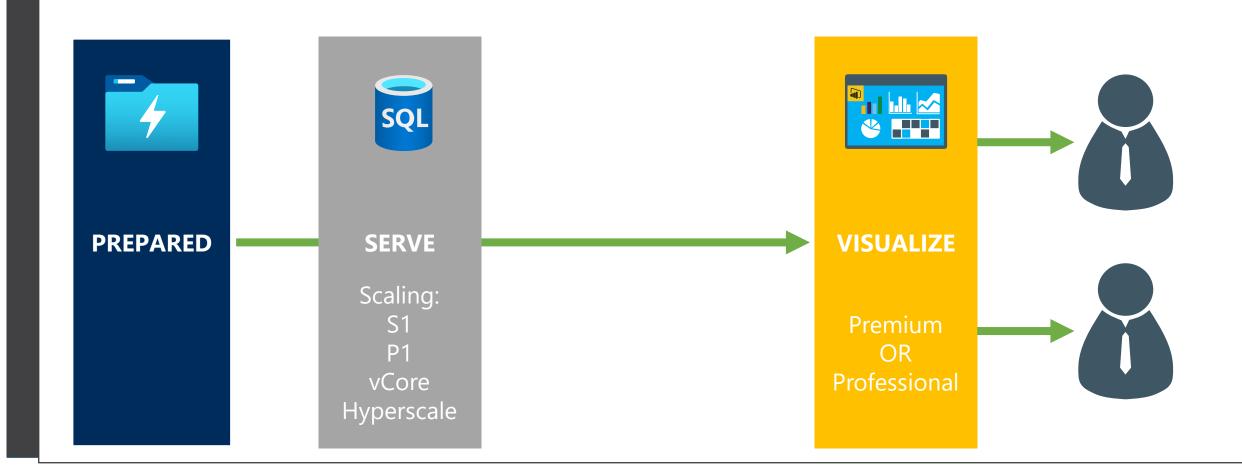






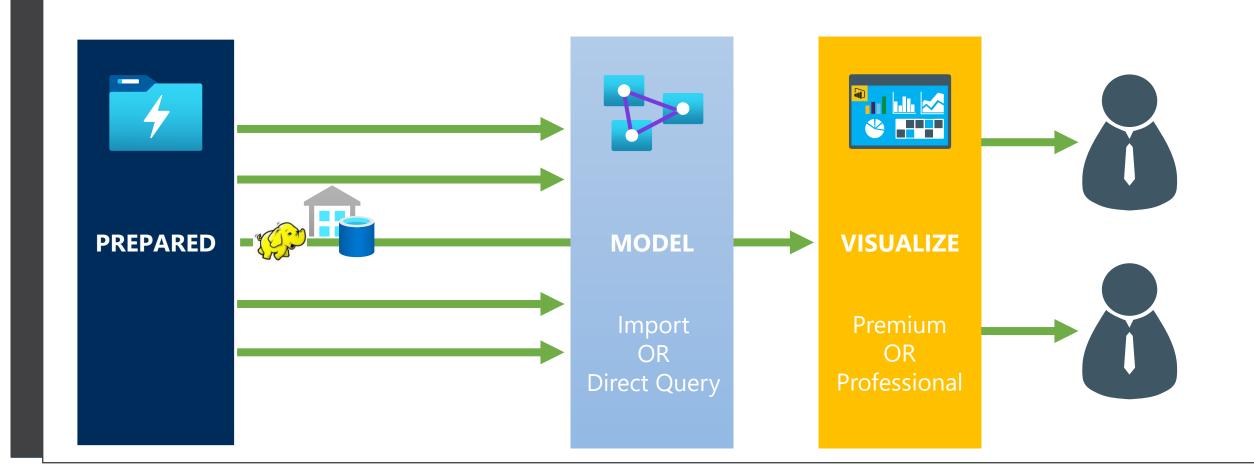






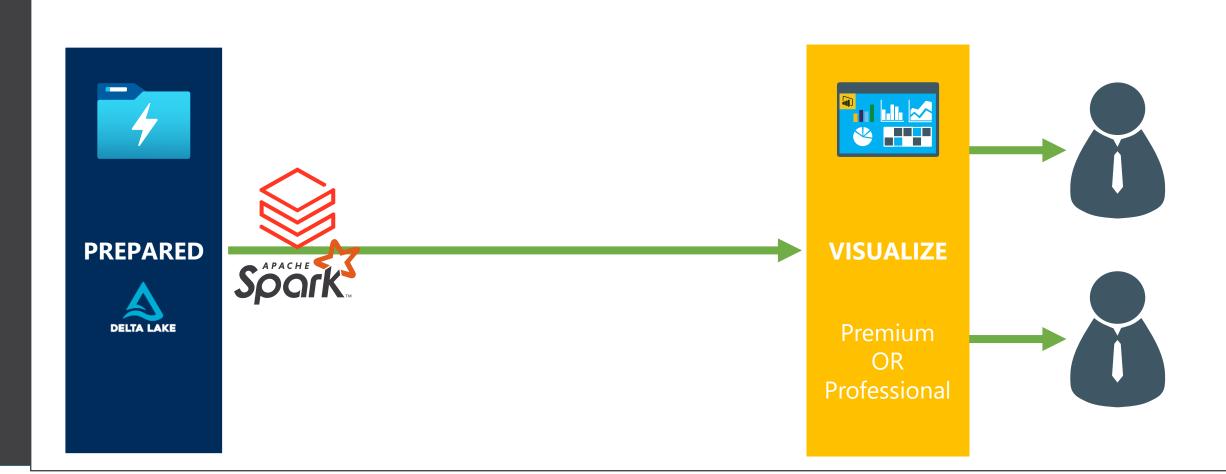












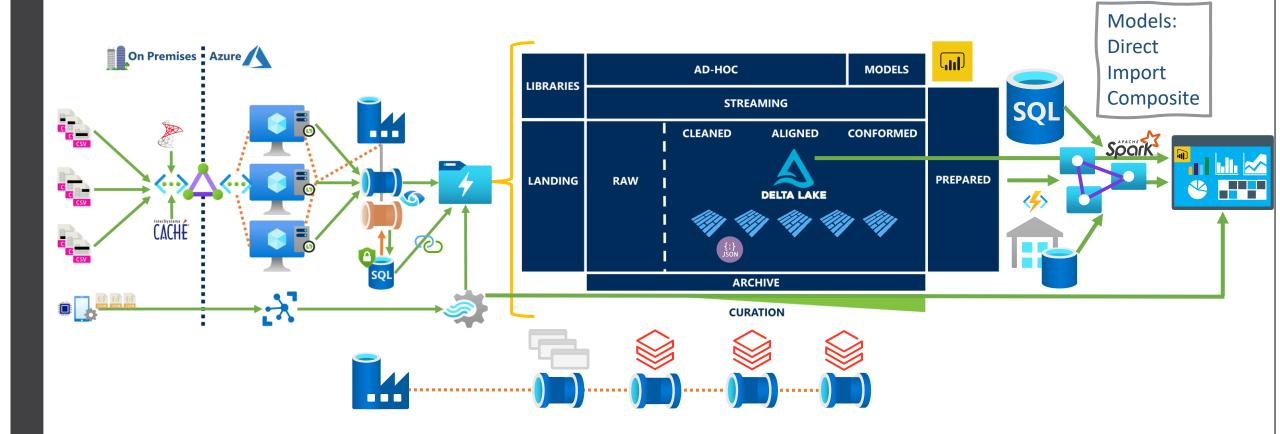




#### **Extract**

# Transform

# Load



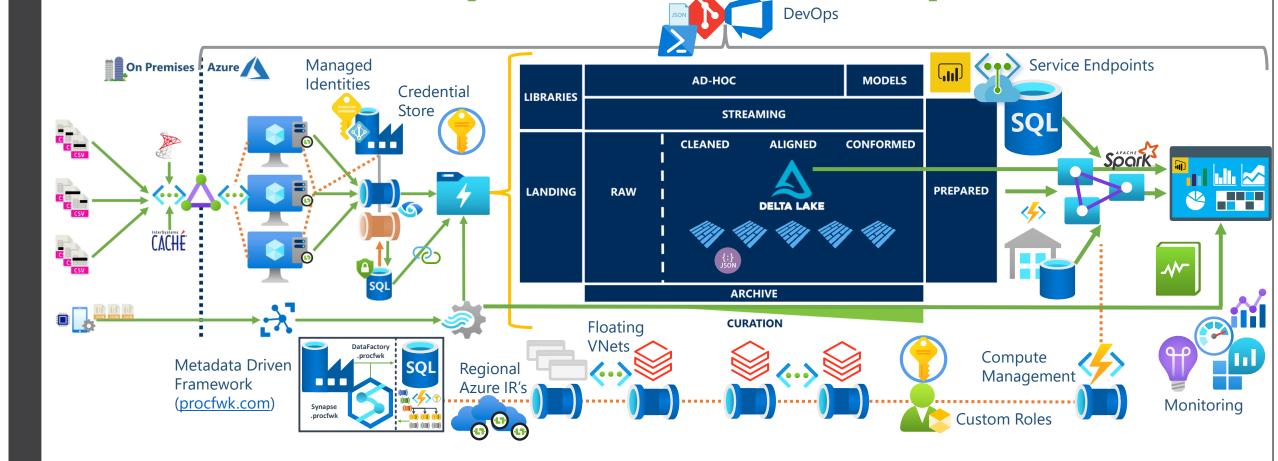




#### Extract

# Transform

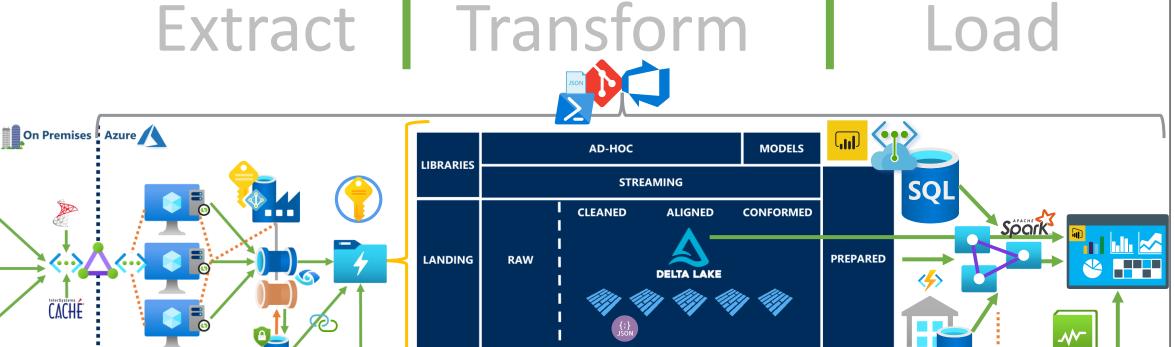
Load





procfwk.com





ARCHIVE

**CURATION** 

Q: Should we build our data platform solution like this?... A: It depends!





# Thank you for listening...

#### Paul Andrew





Blog: mrpaulandrew.com

YouTube: c/mrpaulandrew

Email: paul@mrpaulandrew.com

Twitter: @mrpaulandrew

LinkedIn: In/mrpaulandrew

GitHub: github.com/mrpaulandrew

