

# Module 13 - *Bonus*

An Architects Recap



```
BEGIN --get ready
  SELECT
    *
  FROM
    [Training]
  WHERE
    [Module]
    BETWEEN 1 AND 12;
```

# Agenda



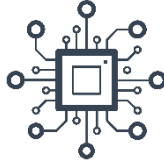
1. Design
2. Extract
3. Transform
4. Load

# Agenda



1. **Design**
2. Extract
3. Transform
4. Load

# Goal



Clean  
Enrich  
Conform  
Translate  
Transform  
Curate  
Analyse  
Model  
Predict  
Master



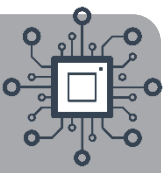
Data Sources



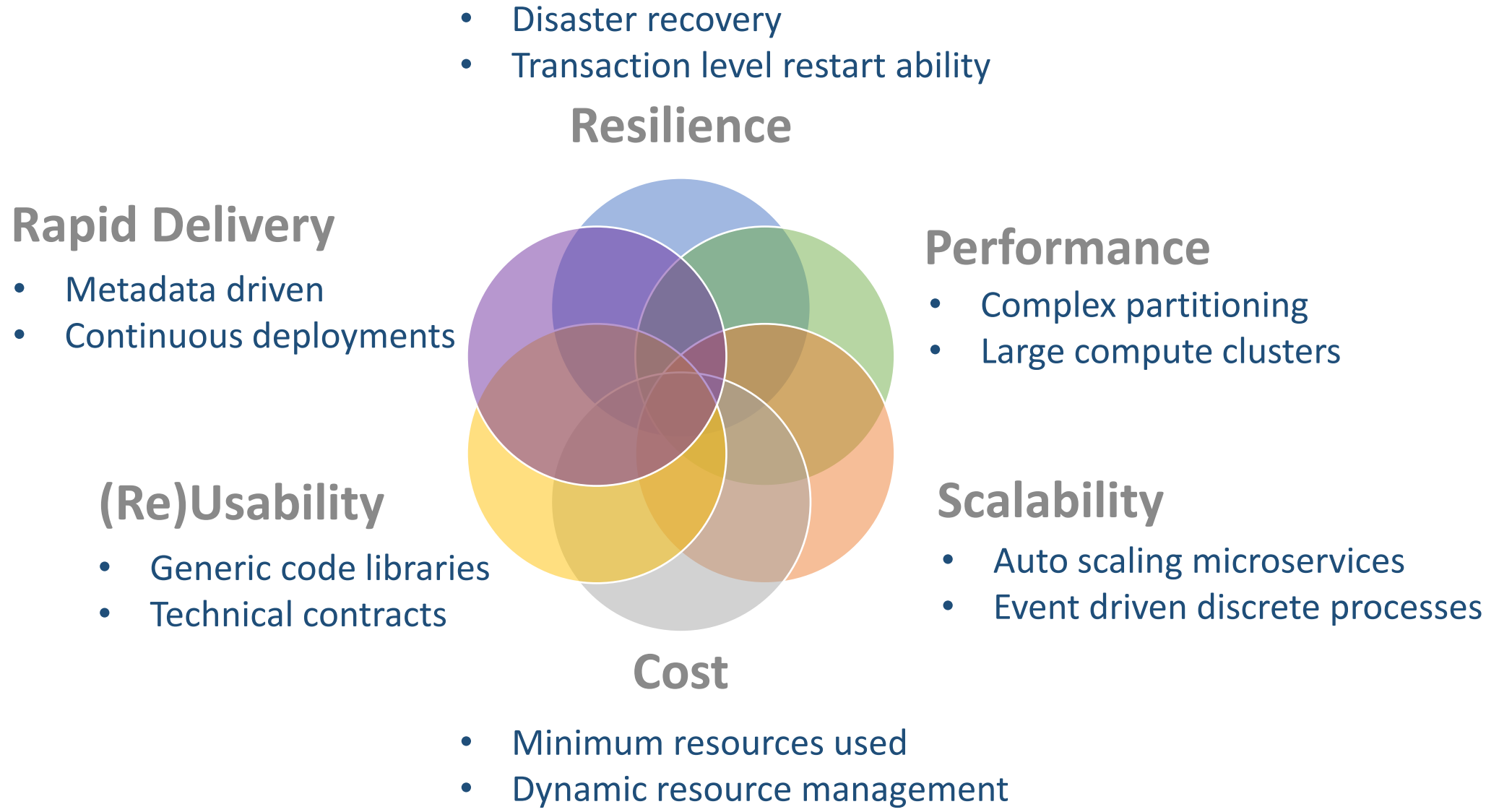
Data Warehouse

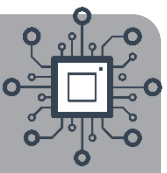


Data Insights

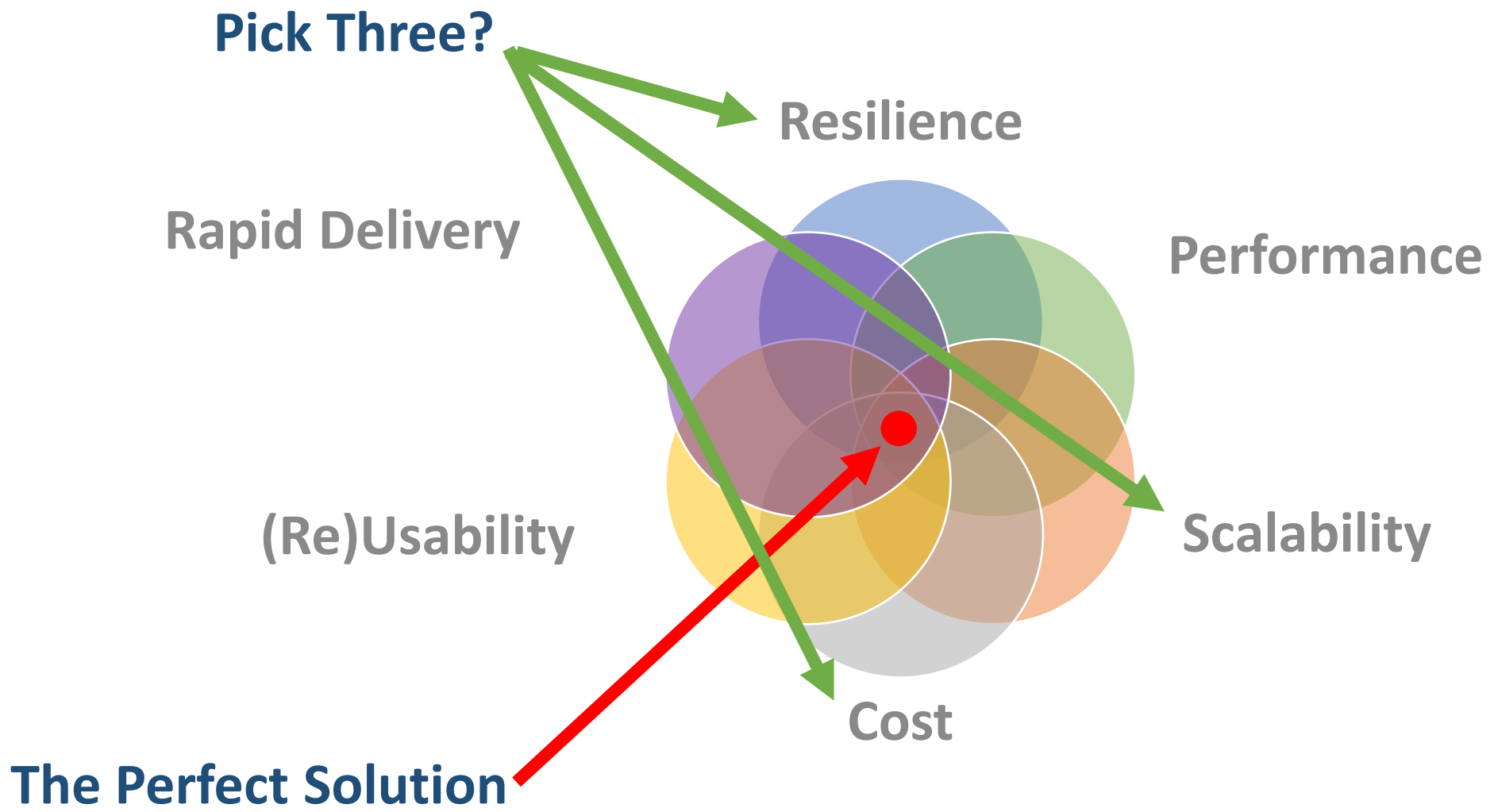


# What is your primary design focus?





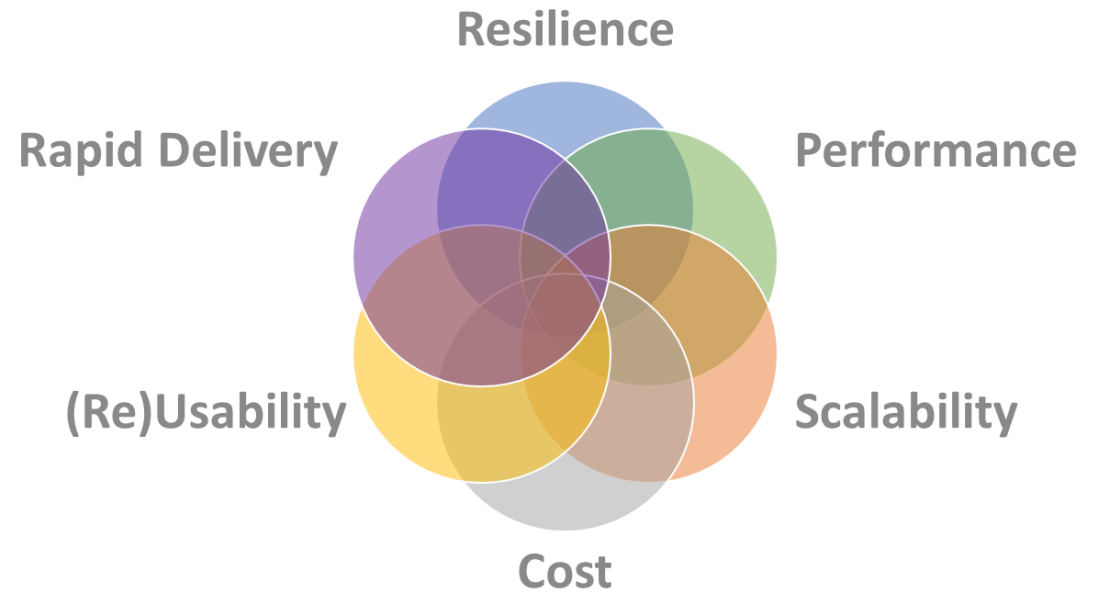
# What is your primary design focus?



# Agenda



1. Design ✓
2. Extract
3. Transform
4. Load



# Agenda



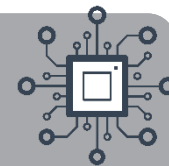
1. Design
2. **Extract**
3. Transform
4. Load







# Data Extraction & Ingestion



## Data Structure



## Data Source



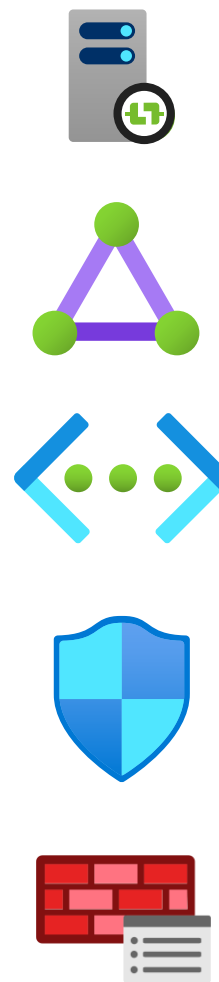
## Push or Pull



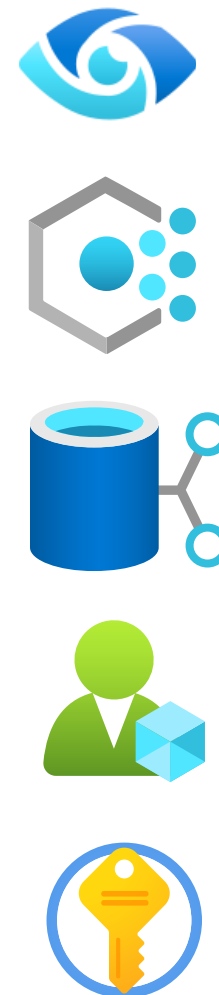
## Batch or Speed



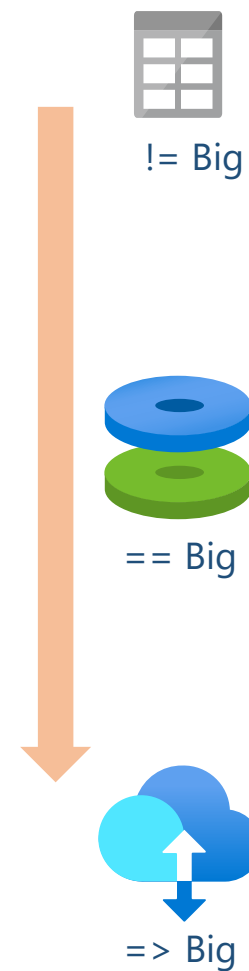
## Public or Private Transfer



## Data Sensitivity

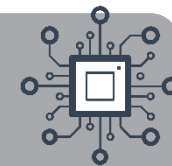


## Data Volume





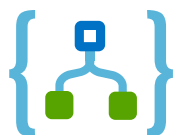
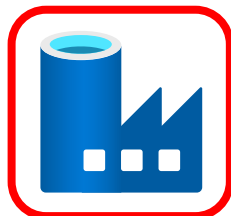
# Data Extraction & Ingestion – Spec v1



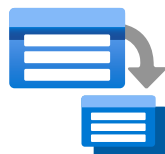
## Data Structure



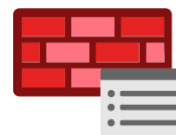
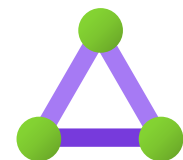
## Push or Pull



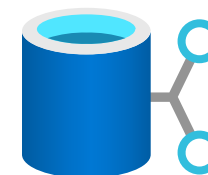
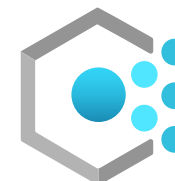
## Batch or Speed



## Public or Private Transfer



## Data Sensitivity



## Data Volume

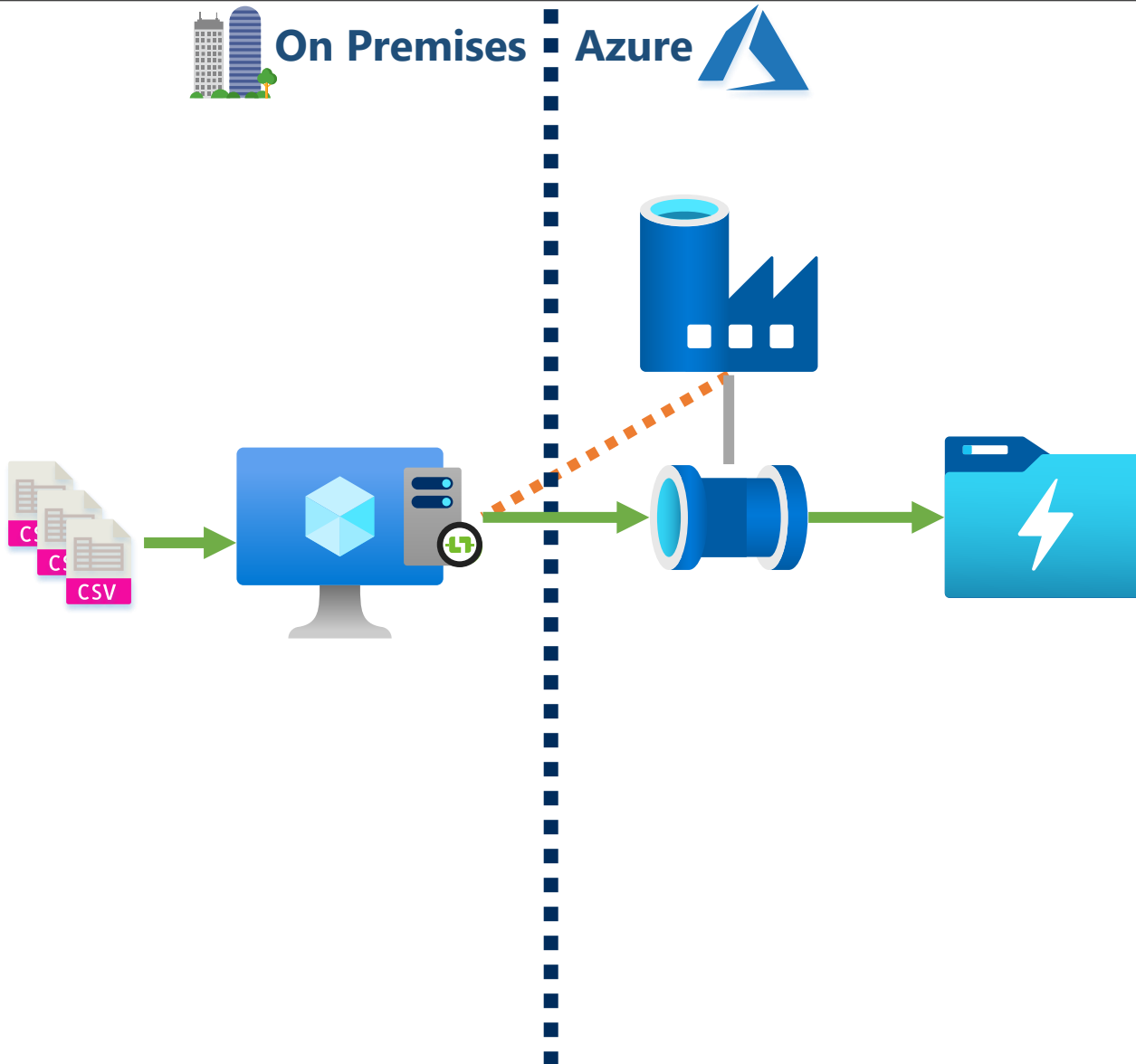


## Data Source





# Data Extraction & Ingestion – Solution 1

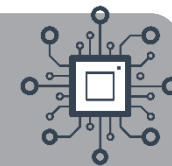


## Requirements:

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



# Data Extraction & Ingestion – Spec v2



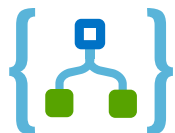
## Data Structure



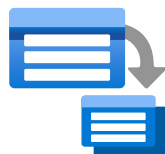
## Data Source



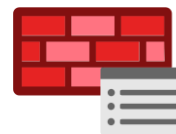
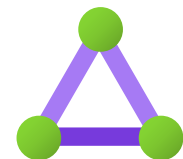
## Push or Pull



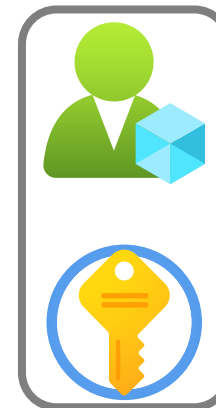
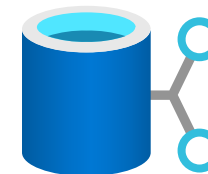
## Batch or Speed



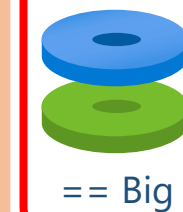
## Public or Private Transfer



## Data Sensitivity

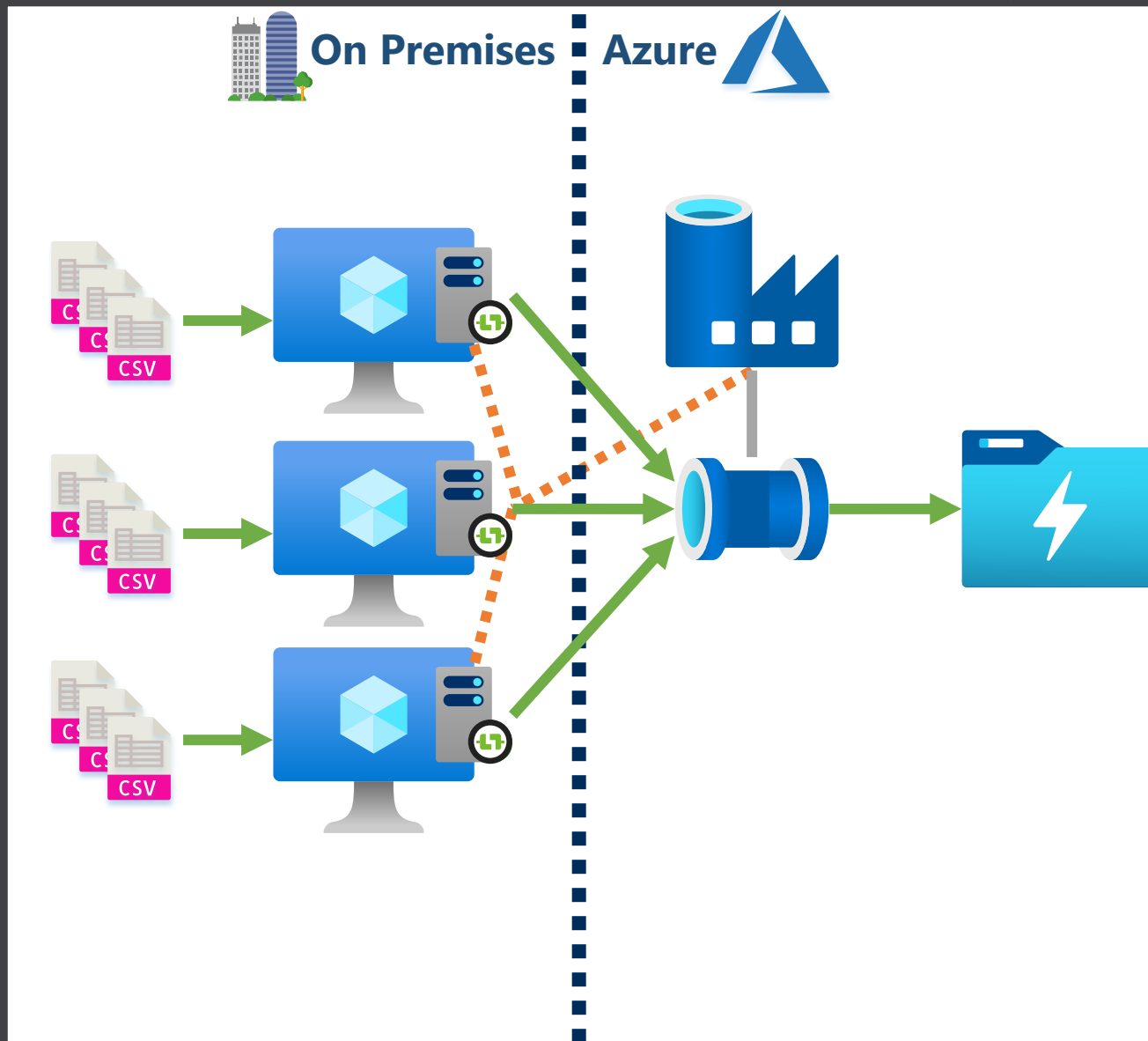
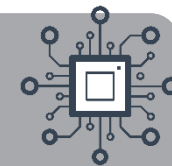


## Data Volume





# Data Extraction & Ingestion – Solution 2

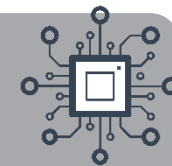


## Requirements:

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



# Data Extraction & Ingestion – Spec v3



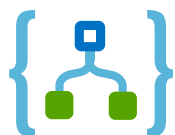
## Data Structure



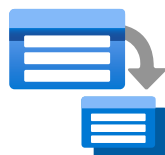
## Data Source



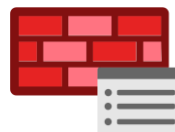
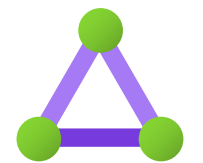
## Push or Pull



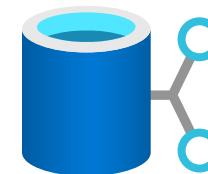
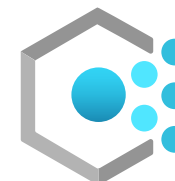
## Batch or Speed



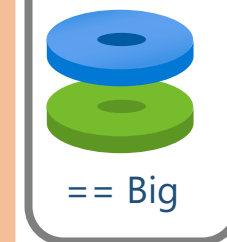
## Public or Private Transfer



## Data Sensitivity

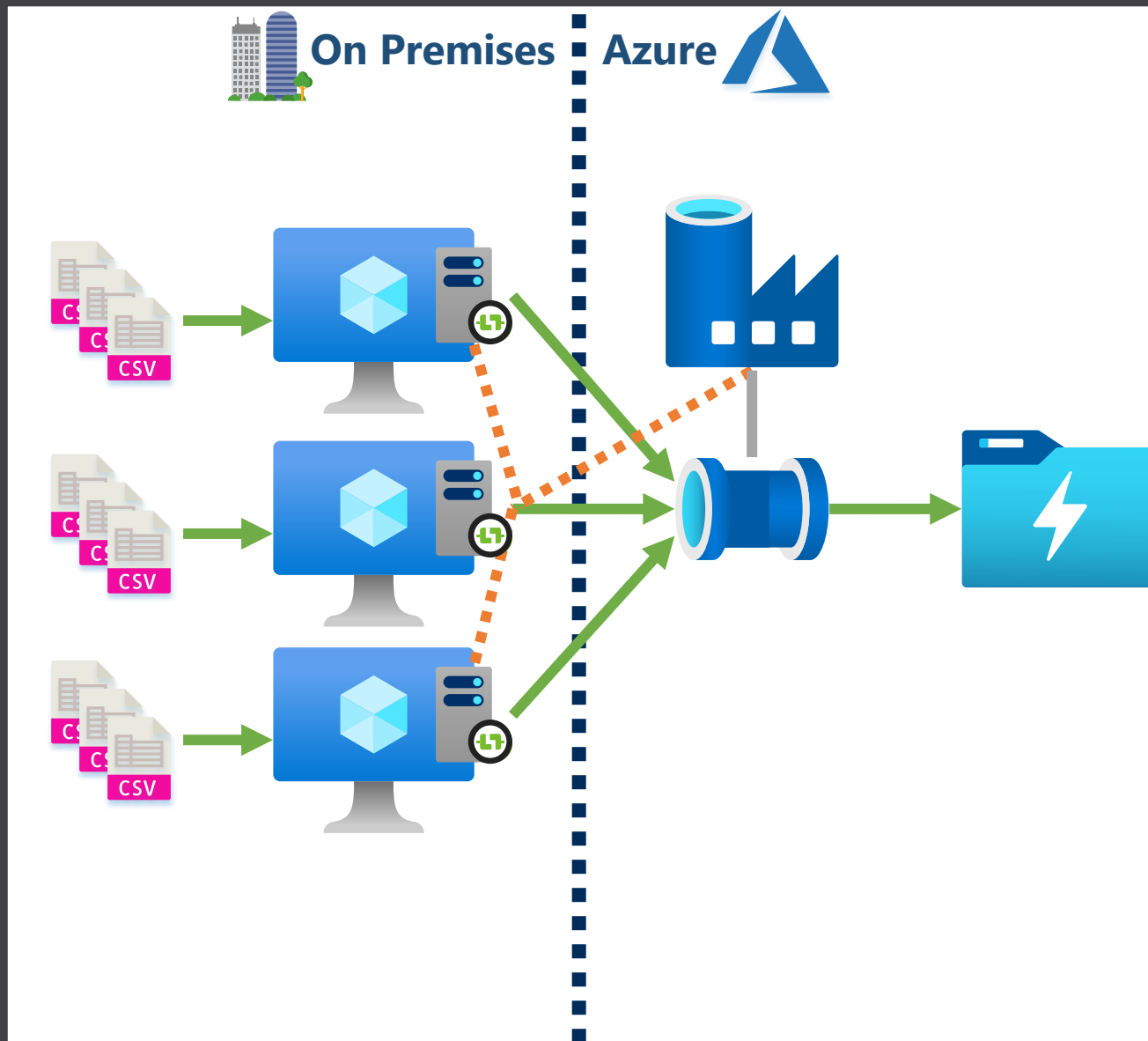
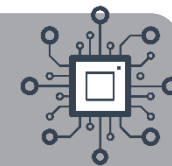


## Data Volume





# Data Extraction & Ingestion – Solution 3

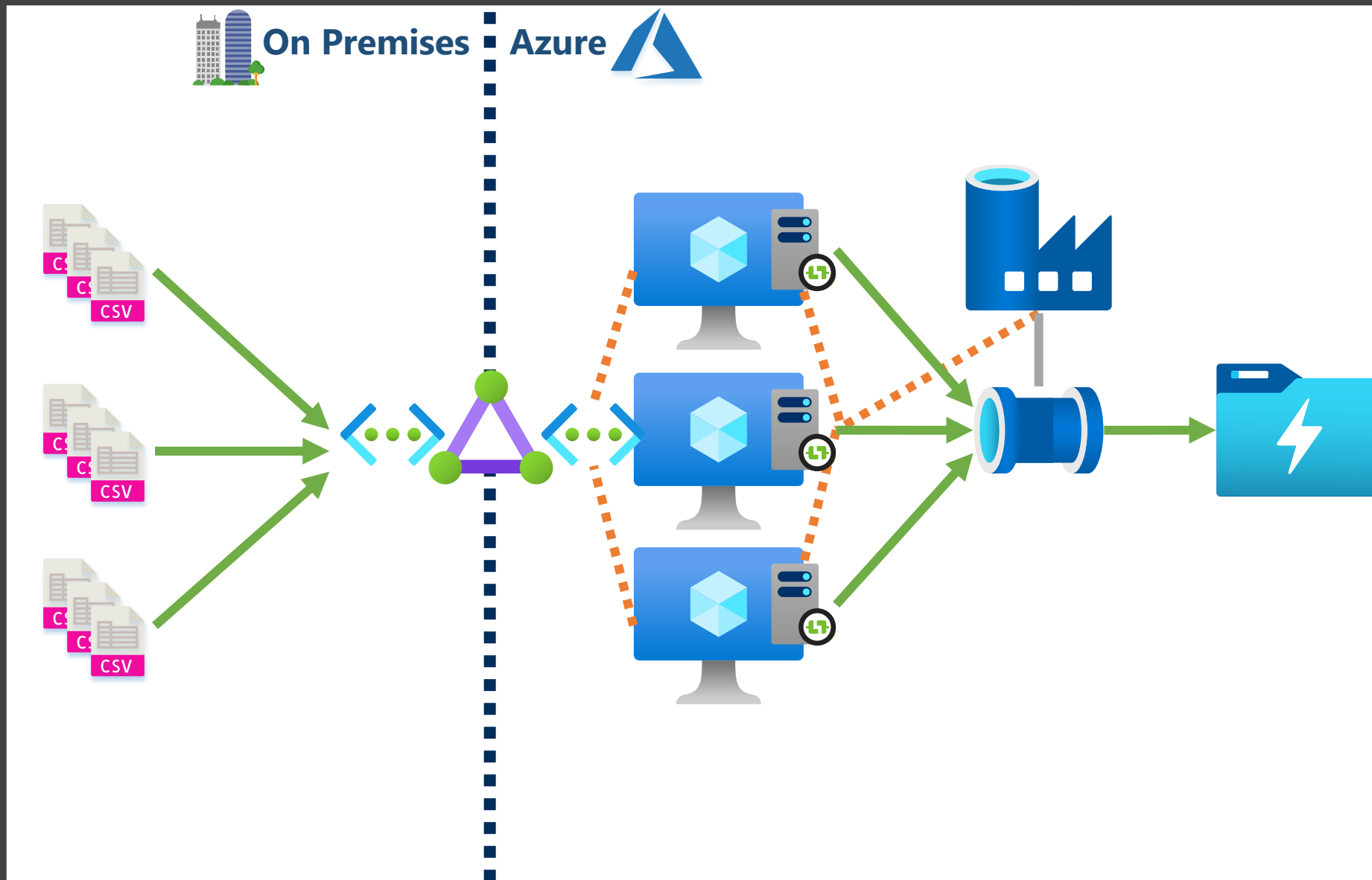
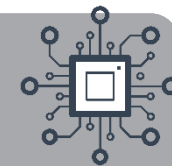


## Requirements:

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



# Data Extraction & Ingestion – Solution 3

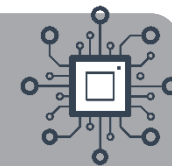


- Requirements:**
- Flat files
  - From local storage
  - Pulled from source
  - Batch load
  - Private connections
  - No PII data
  - Large data volumes





# Data Extraction & Ingestion – Spec v4



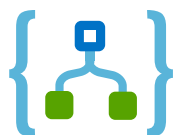
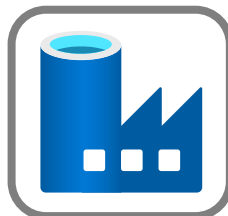
## Data Structure



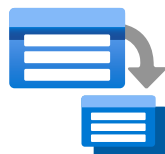
## Data Source



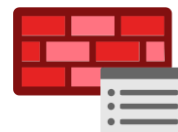
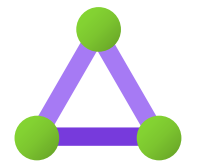
## Push or Pull



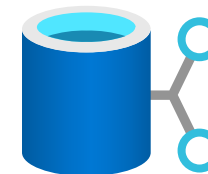
## Batch or Speed



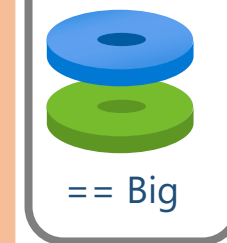
## Public or Private Transfer



## Data Sensitivity

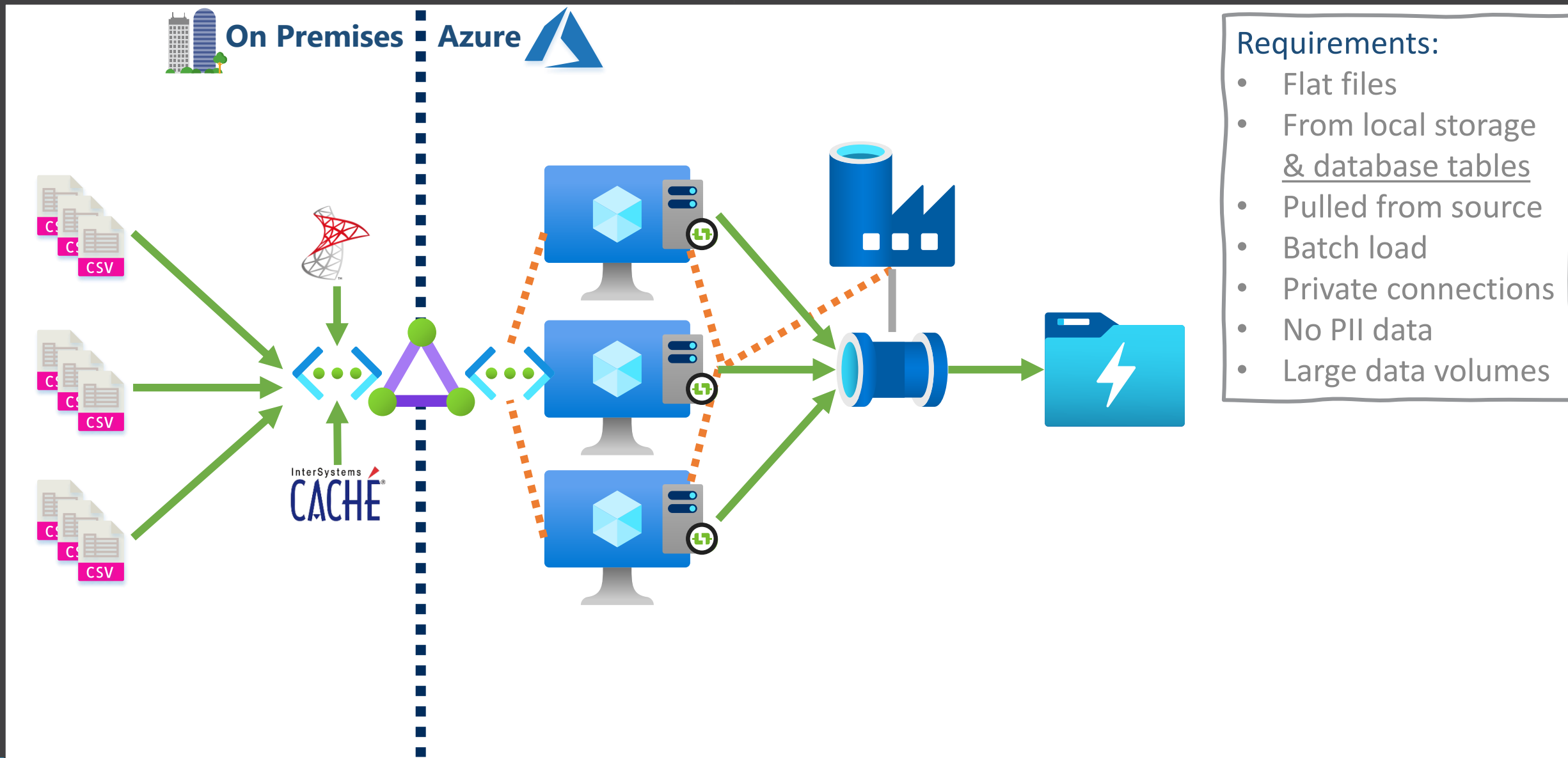
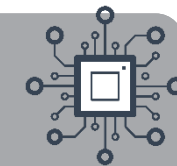


## Data Volume





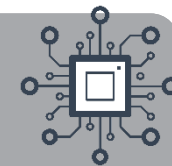
# Data Extraction & Ingestion – Solution 4



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



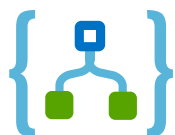
# Data Extraction & Ingestion – Spec v5



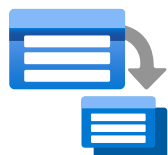
## Data Structure



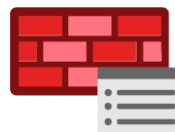
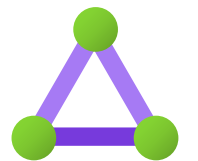
## Push or Pull



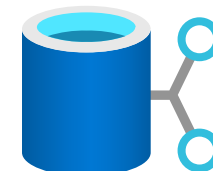
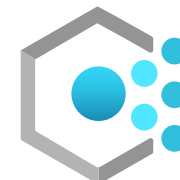
## Batch or Speed



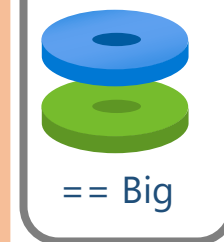
## Public or Private Transfer



## Data Sensitivity



## Data Volume

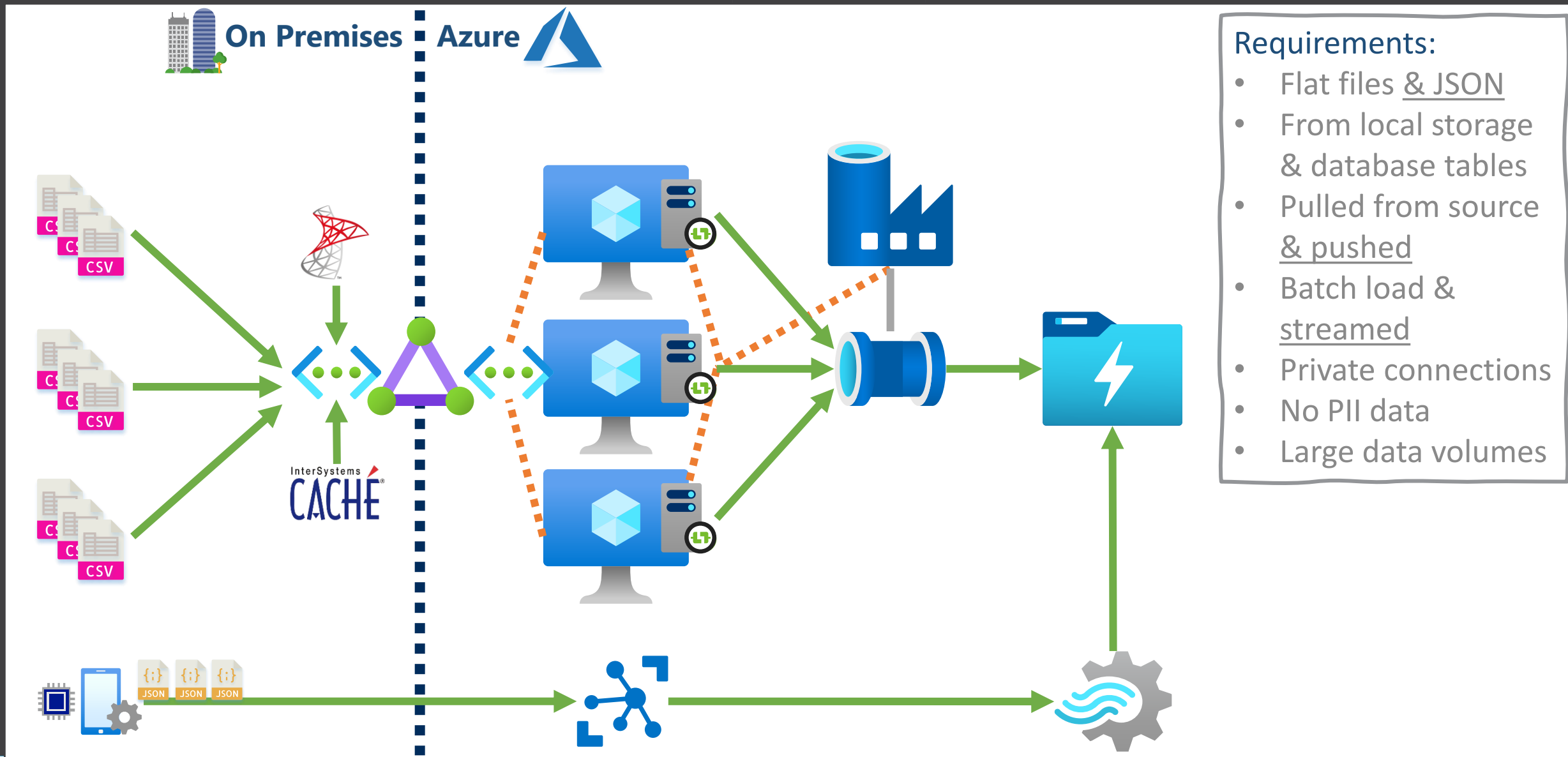
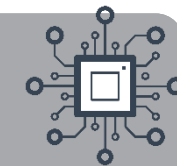


## Data Source



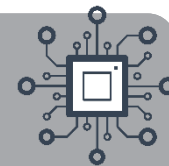


# Data Extraction & Ingestion – Solution 5

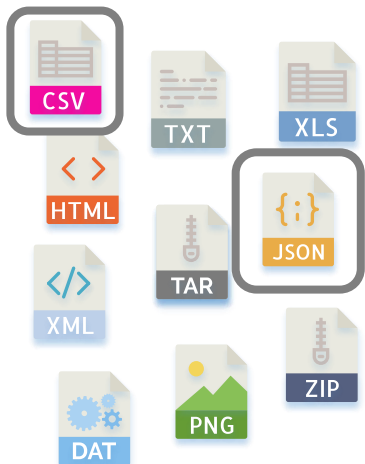




# Data Extraction & Ingestion – Spec v6



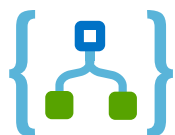
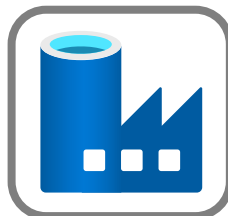
## Data Structure



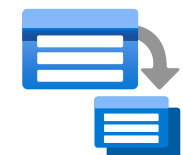
## Data Source



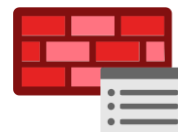
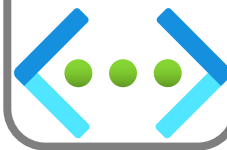
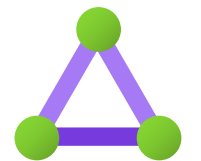
## Push or Pull



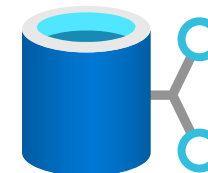
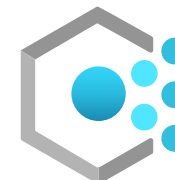
## Batch or Speed



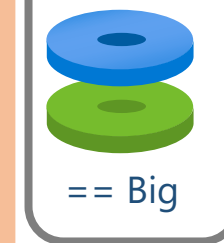
## Public or Private Transfer



## Data Sensitivity

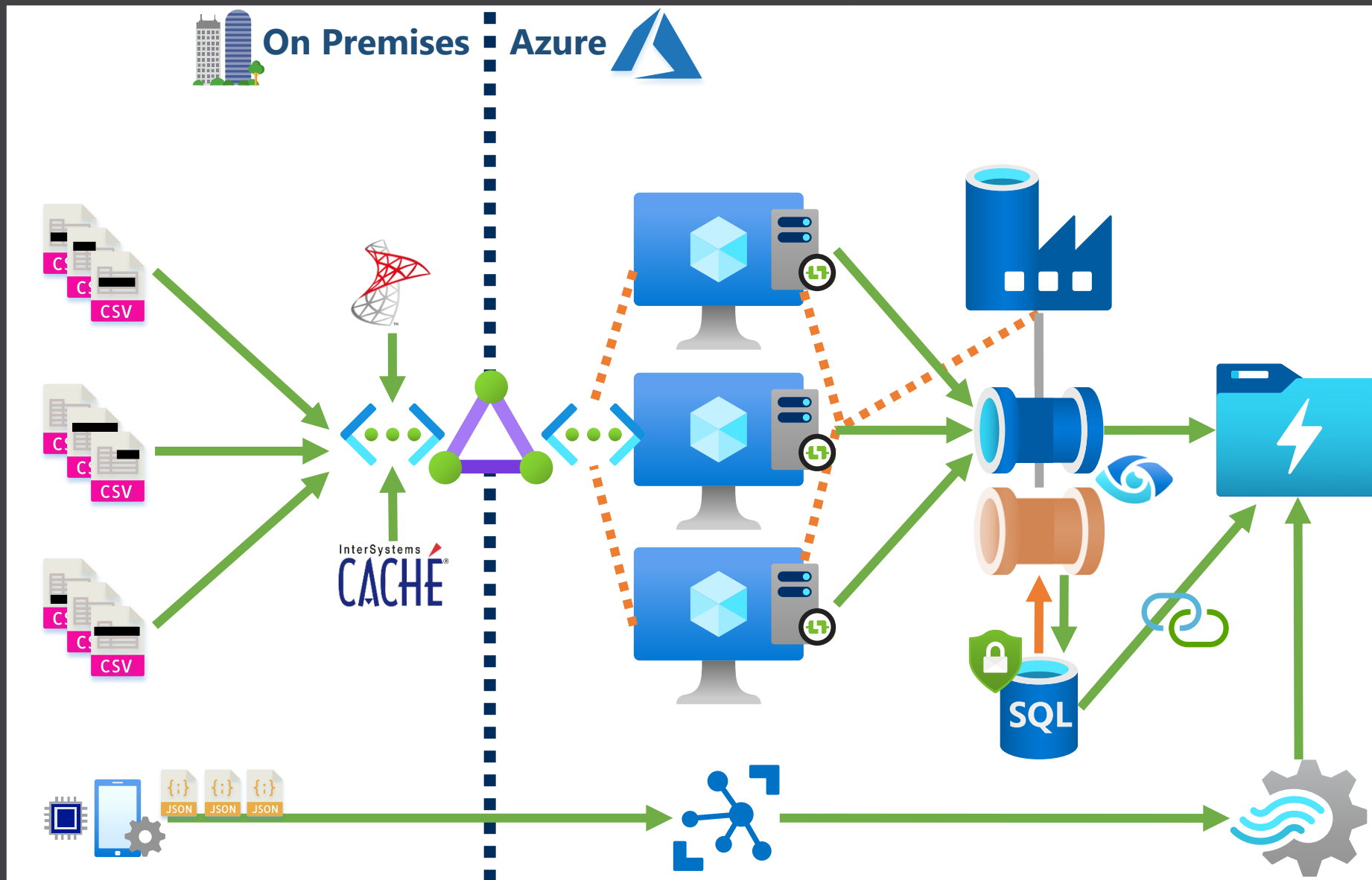
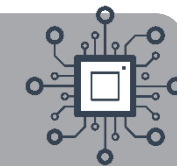


## Data Volume





# 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



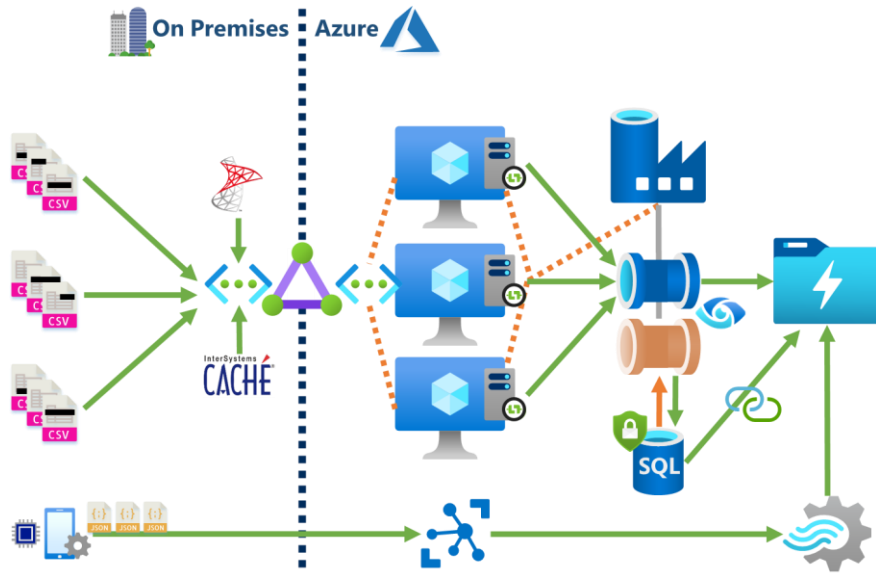
# Overall Architecture



## Extract

## Transform

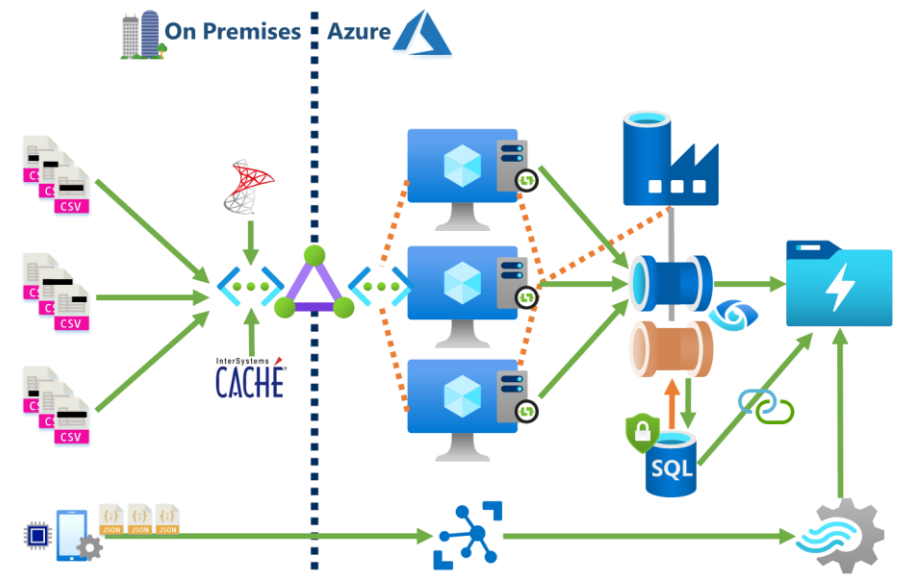
## Load



# Agenda



1. Design ✓
2. Extract ✓
3. Transform
4. Load

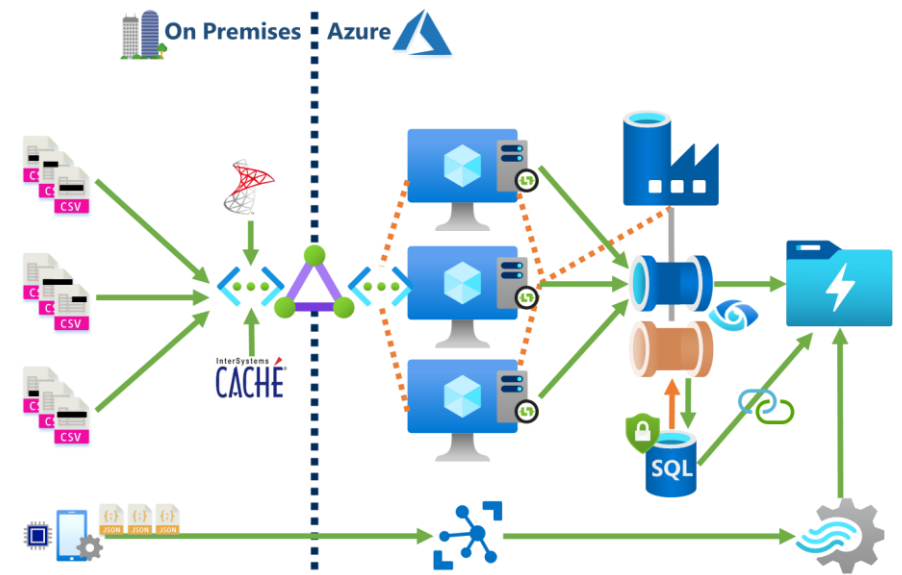




# Agenda



1. Design ✓
2. Extract ✓
3. **Transform**
4. Load

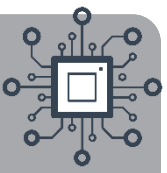


# Agenda



- 1. Design ✓
- 2. Extract ✓
- 3. Transform
- 4. Load

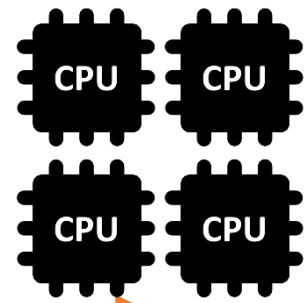
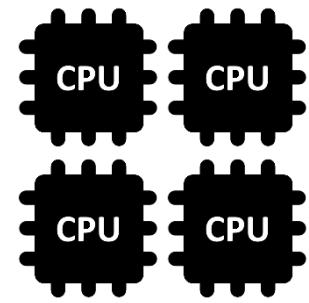
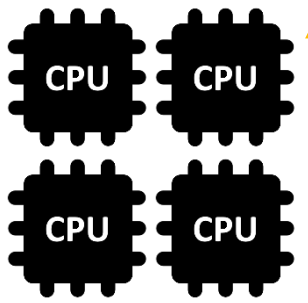
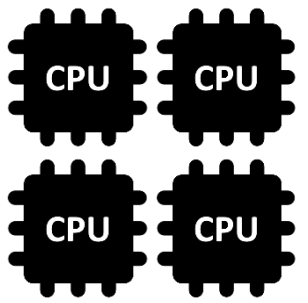
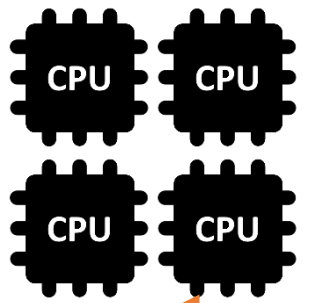
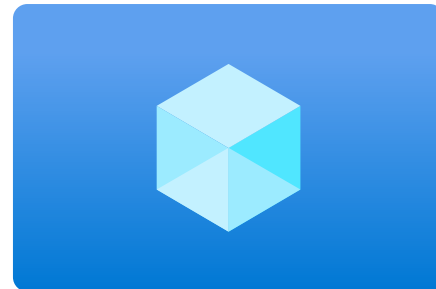
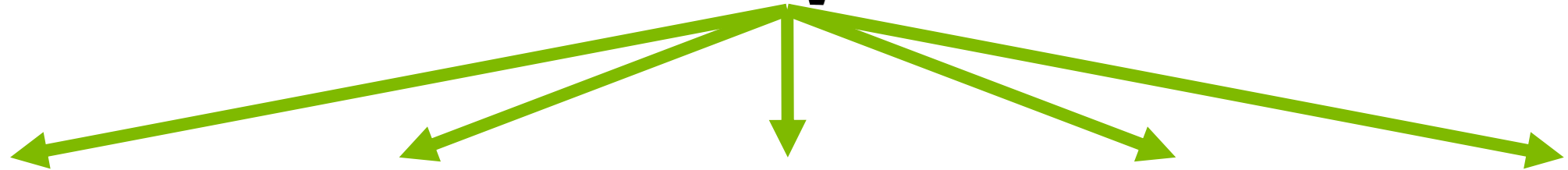
Compute  
Storage, Structure  
& Data Format



# Scaling Up and/or Scaling Out



Workload:  
Process 100TB of Data



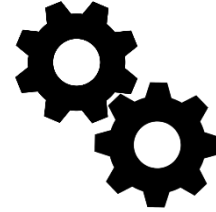
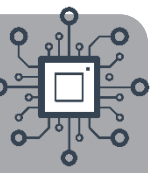
Scale Out

Scale Up





# What Compute Type of Compute?



Workload:  
Process 100TB of Data

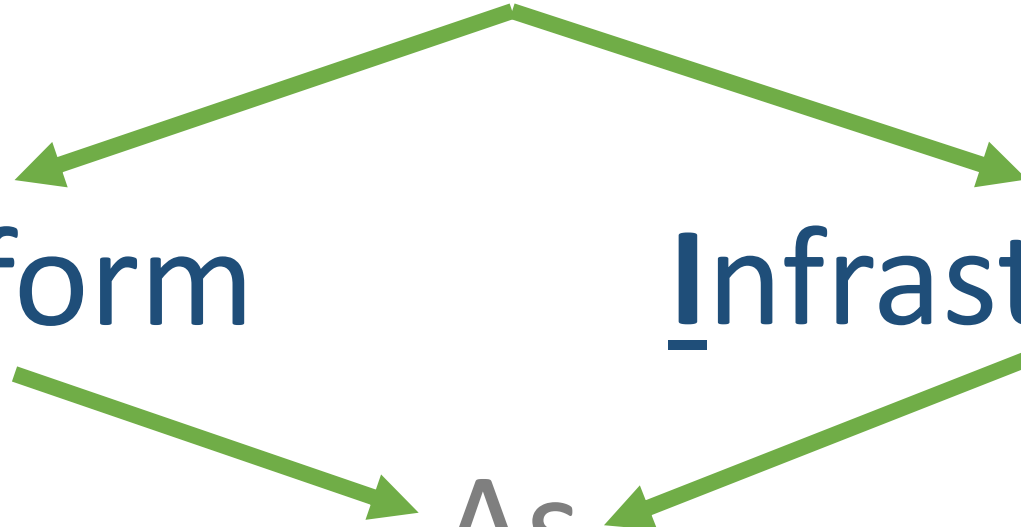
Platform

Infrastructure

As

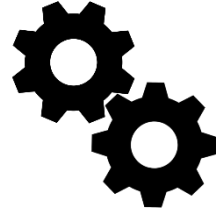
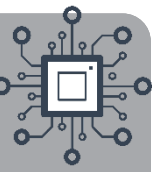
A

Service





# What Compute Type of Compute?



Workload:

Process 100TB of Data

Platform

As

A

Service

IaaS

PaaS

Applications

Applications

Data

Data

Runtime

Runtime

Middleware

Middleware

Operating System

Operating System

Virtualization

Virtualization

Servers

Servers

Storage

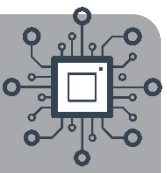
Storage

Networking

Networking



# Data Transformation – Compute



Data Lake Analytics



HDInsight



Relational Database



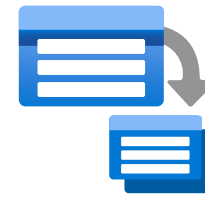
Synapse – SQL Pools or Spark Pools



Databricks



Batch Service



Data Explorer



Automation



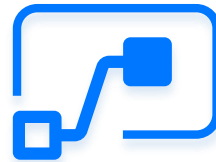
Cosmos



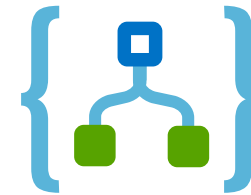
Functions



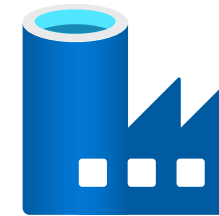
Power BI Data Flows



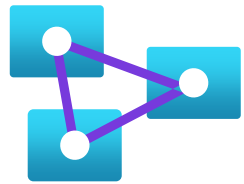
Logic Apps



Data Factory Data Flows

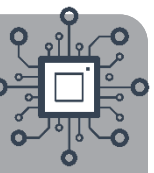


Analysis Services





# Data Transformation – Compute



Data Lake Analytics



HDInsight



Relational Database



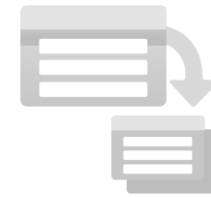
Synapse – SQL Pools or Spark Pools



Databricks



Batch Service



Data Explorer



Automation



Cosmos



Functions



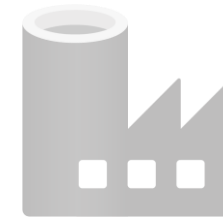
Power BI Data Flows



Logic Apps



Data Factory Data Flows

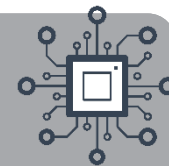


Analysis Services





# Data Transformation – Compute



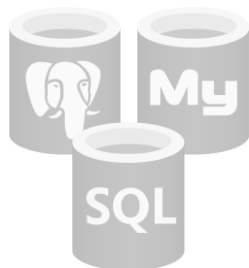
Data Lake Analytics



HDInsight



Relational Database



Batch Service



Data Explorer



Automation



Cosmos



Functions



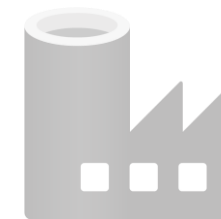
Power BI Data Flows



Logic Apps



Data Factory Data Flows



Analysis Services





# Agenda

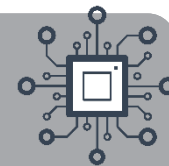


- 1. Design ✓
- 2. Extract ✓
- 3. Transform
- 4. Load

Compute ✓  
Storage, Structure  
& Data Format



# Data Transformation – Storage & Format



Azure Storage Account



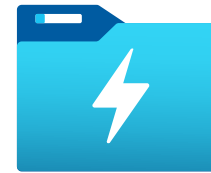
Azure Data Lake Gen2

Hadoop Distributed File System ( HDFS )



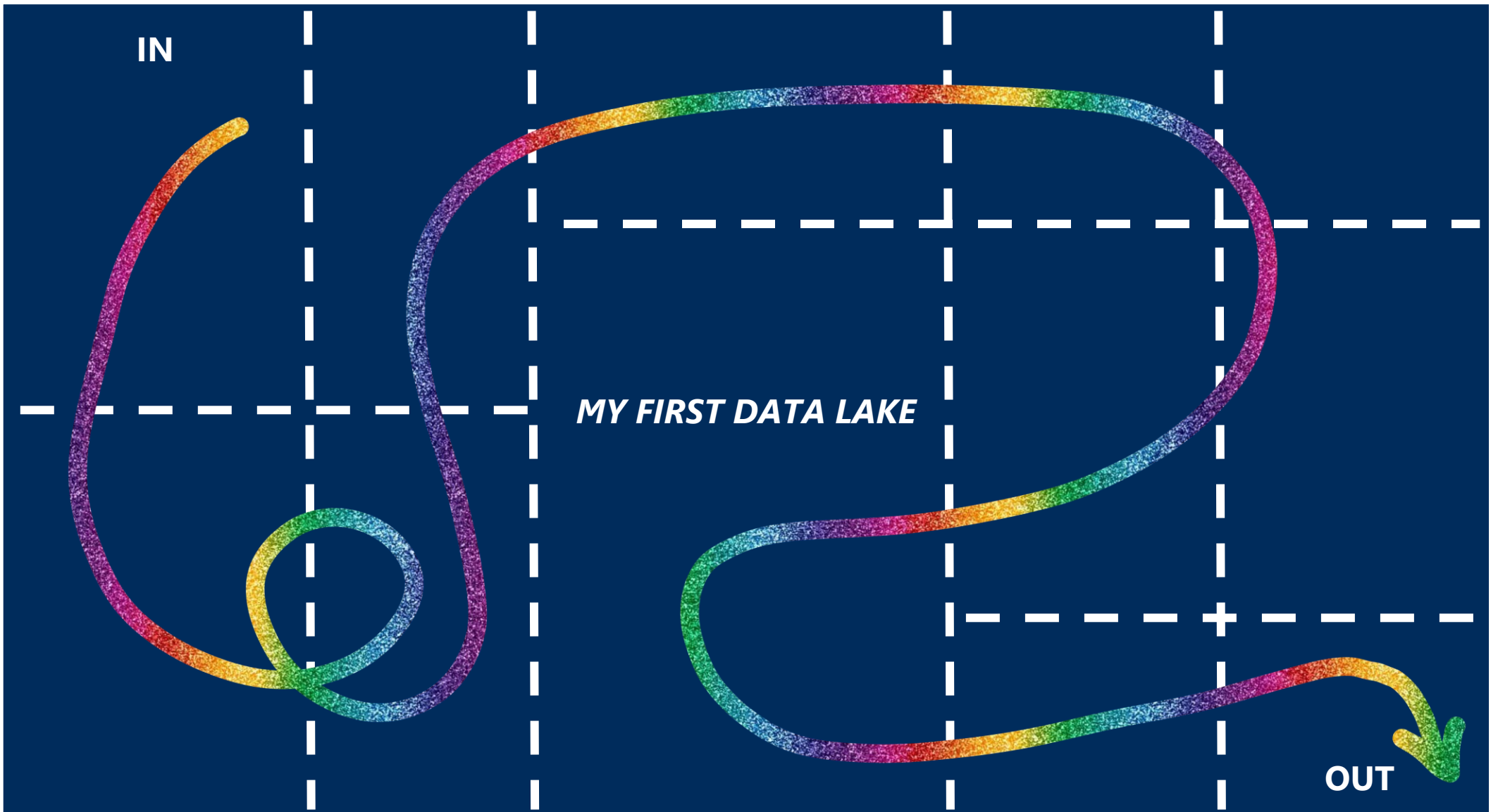
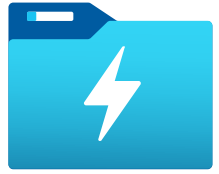
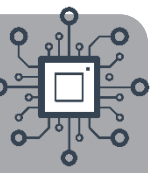


# Data Transformation – Storage & Format



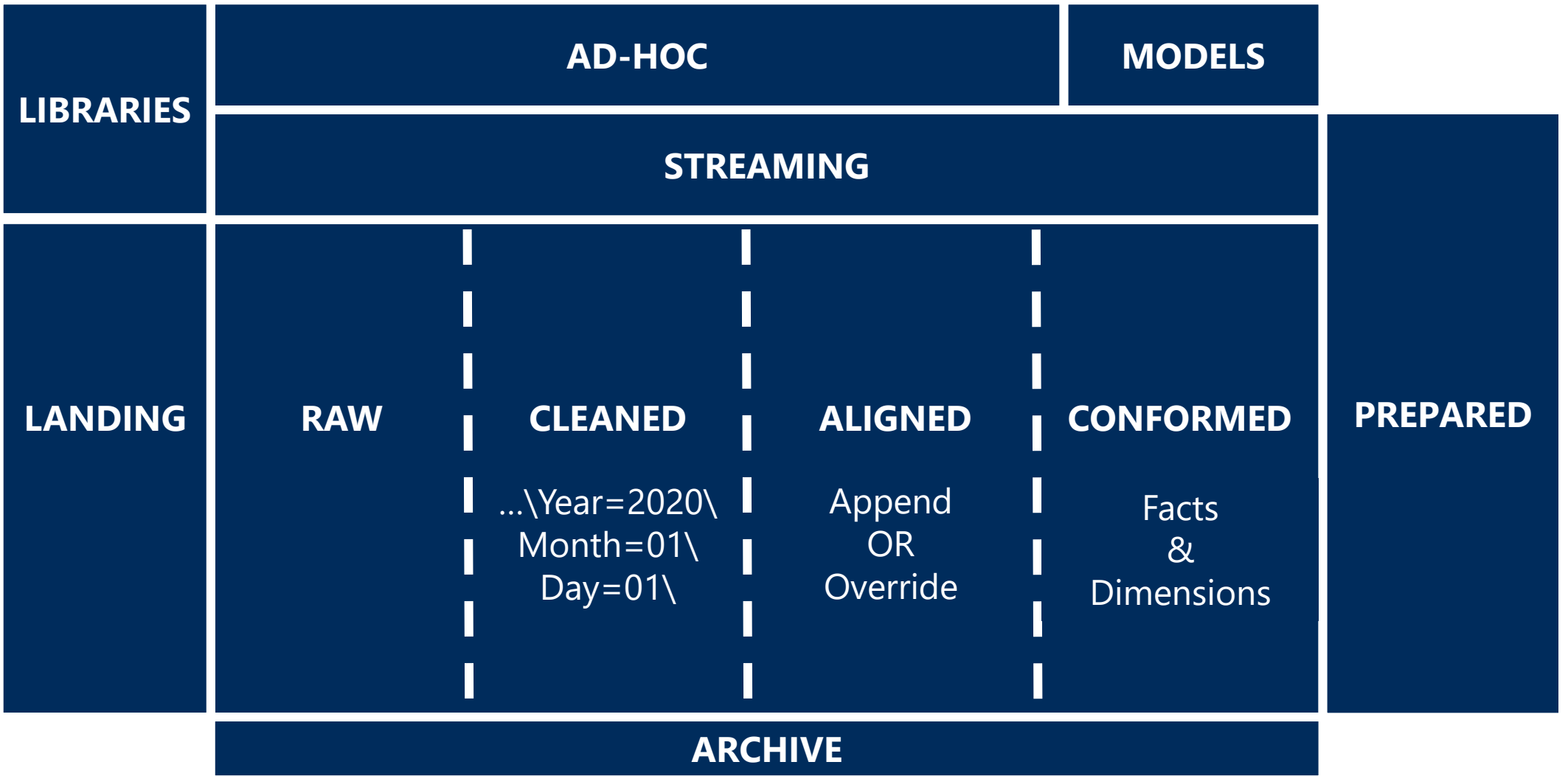
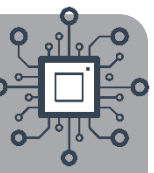


# Data Transformation – Storage & Format





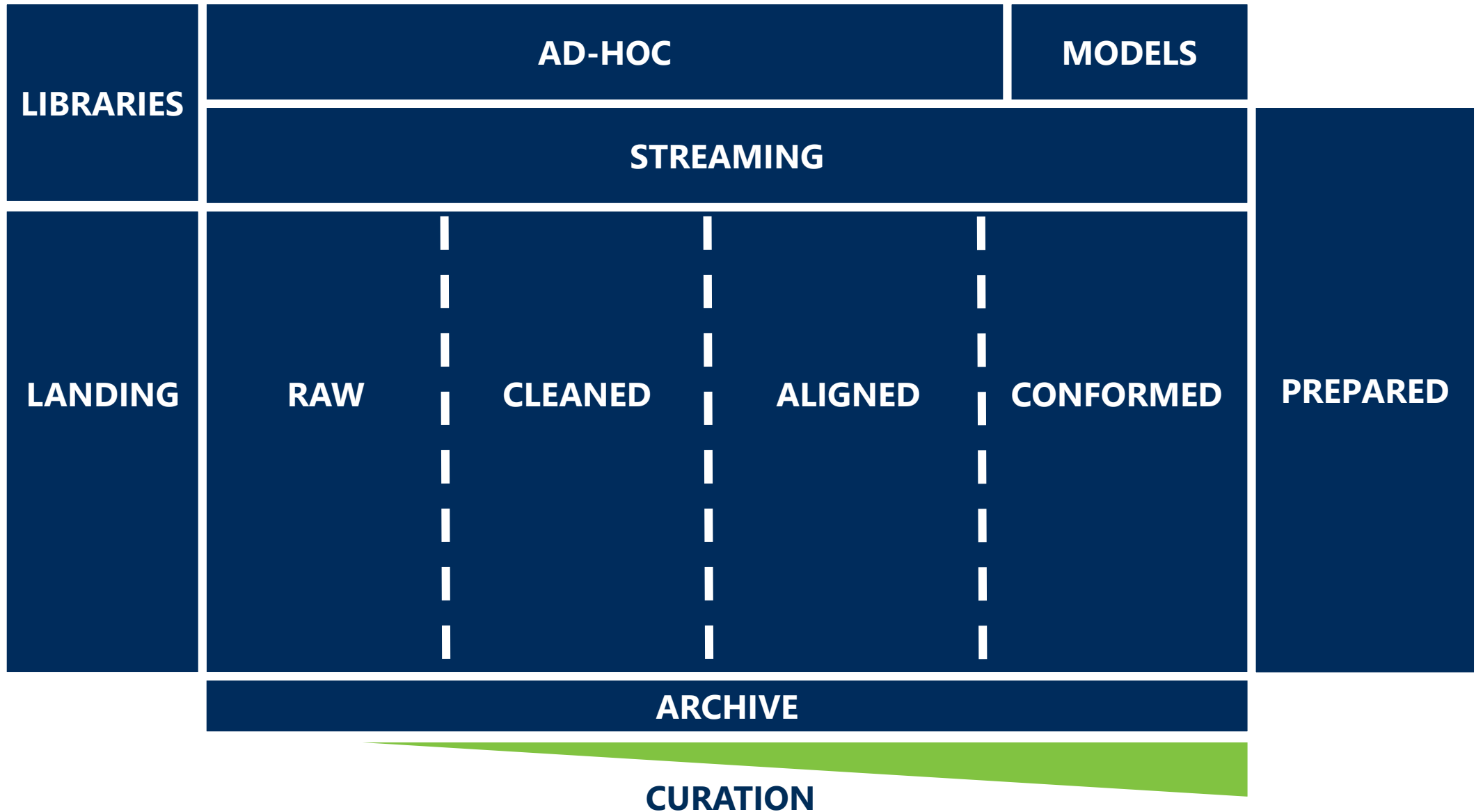
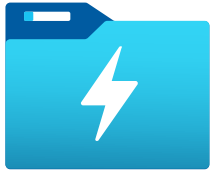
# Data Transformation – Storage & Format



**CURATION**

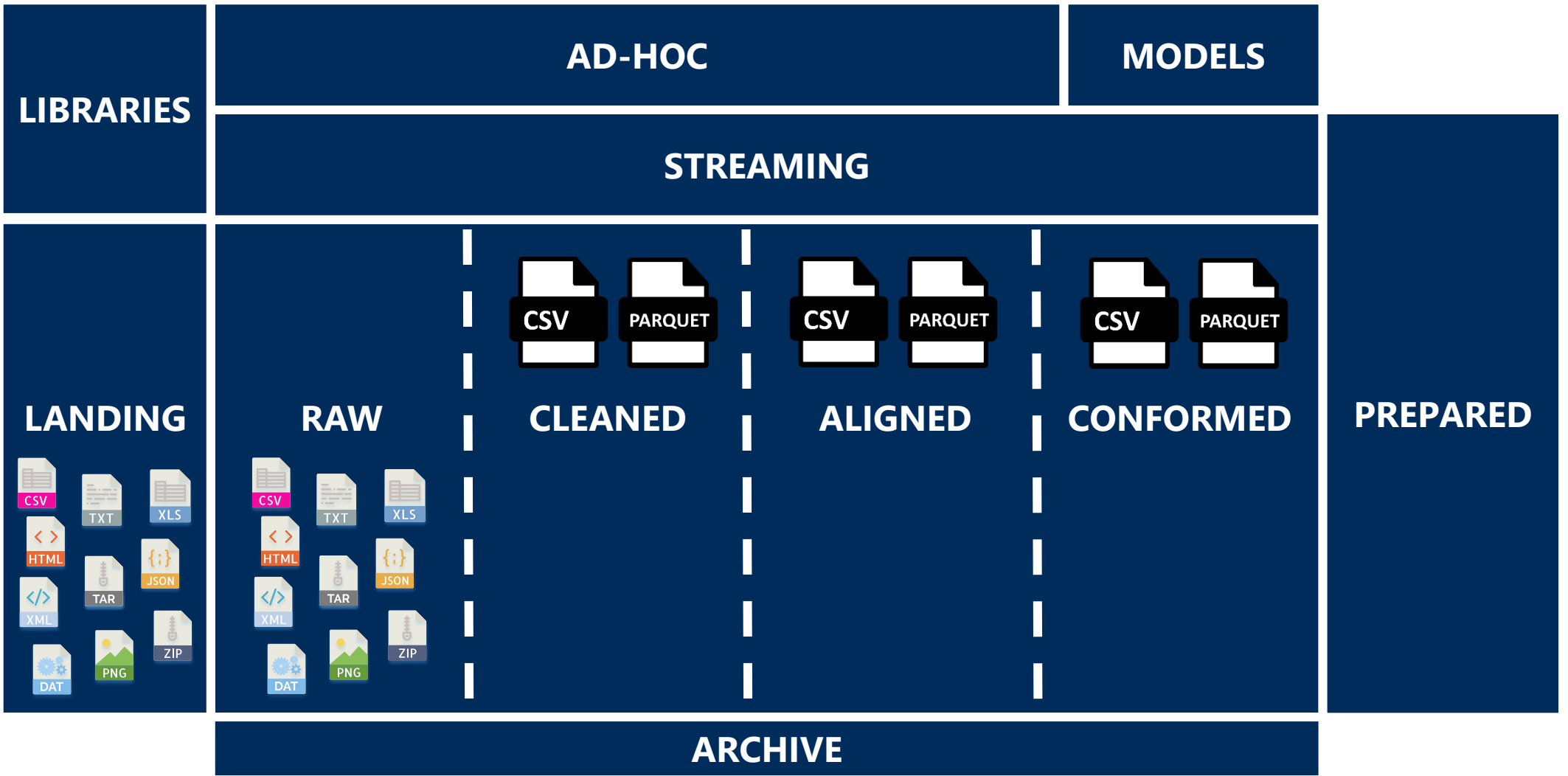
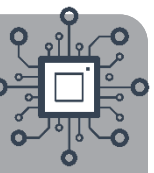


# Data Transformation – Storage & Format





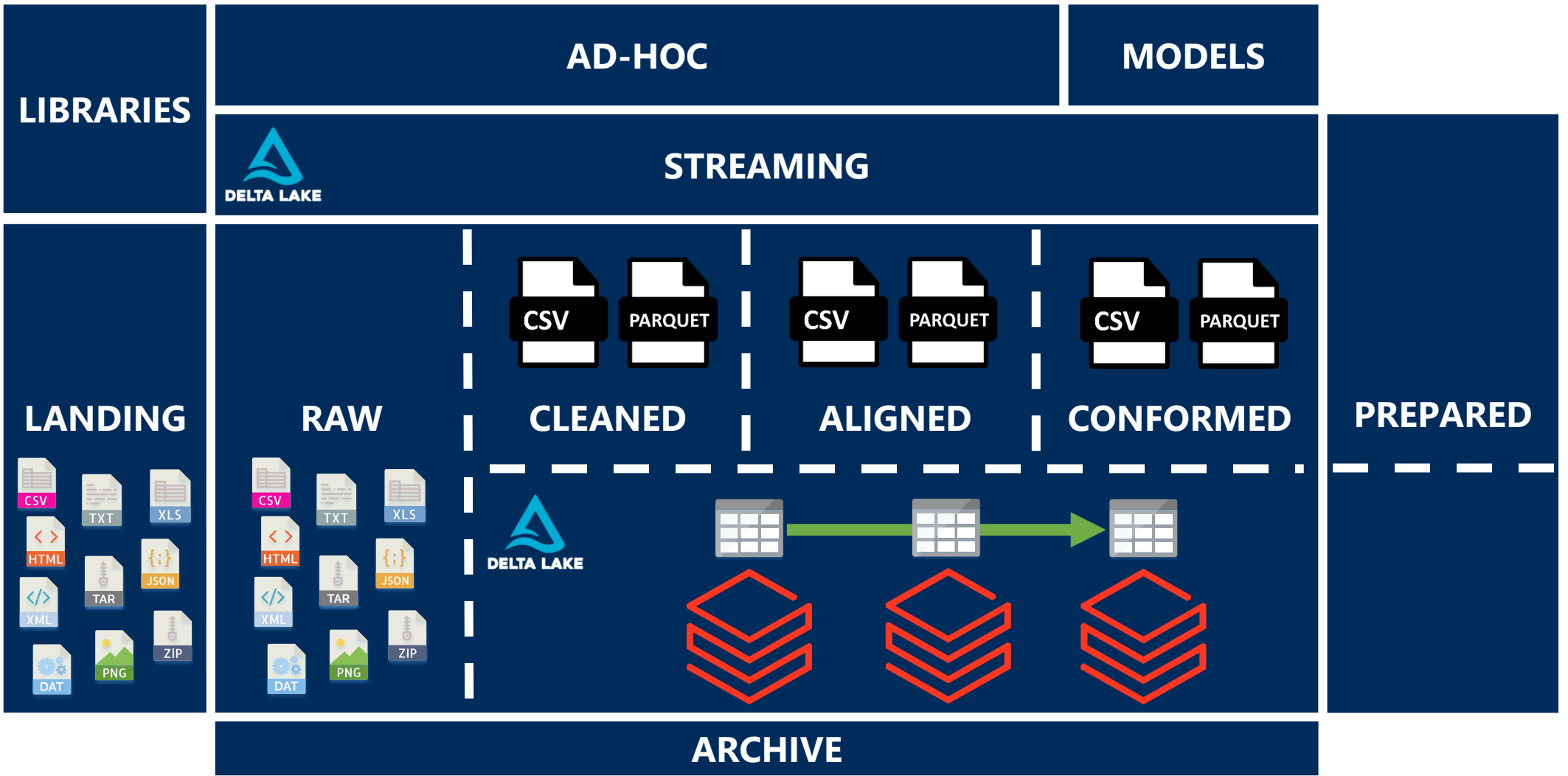
# Data Transformation – Storage & Format



**CURATION**



# Data Transformation – Storage & Format





# Agenda

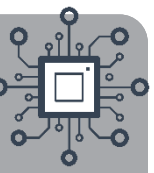


- 1. Design ✓
- 2. Extract ✓
- 3. Transform
- 4. Load

Compute ✓  
Storage, Structure  
& Data Format ✓



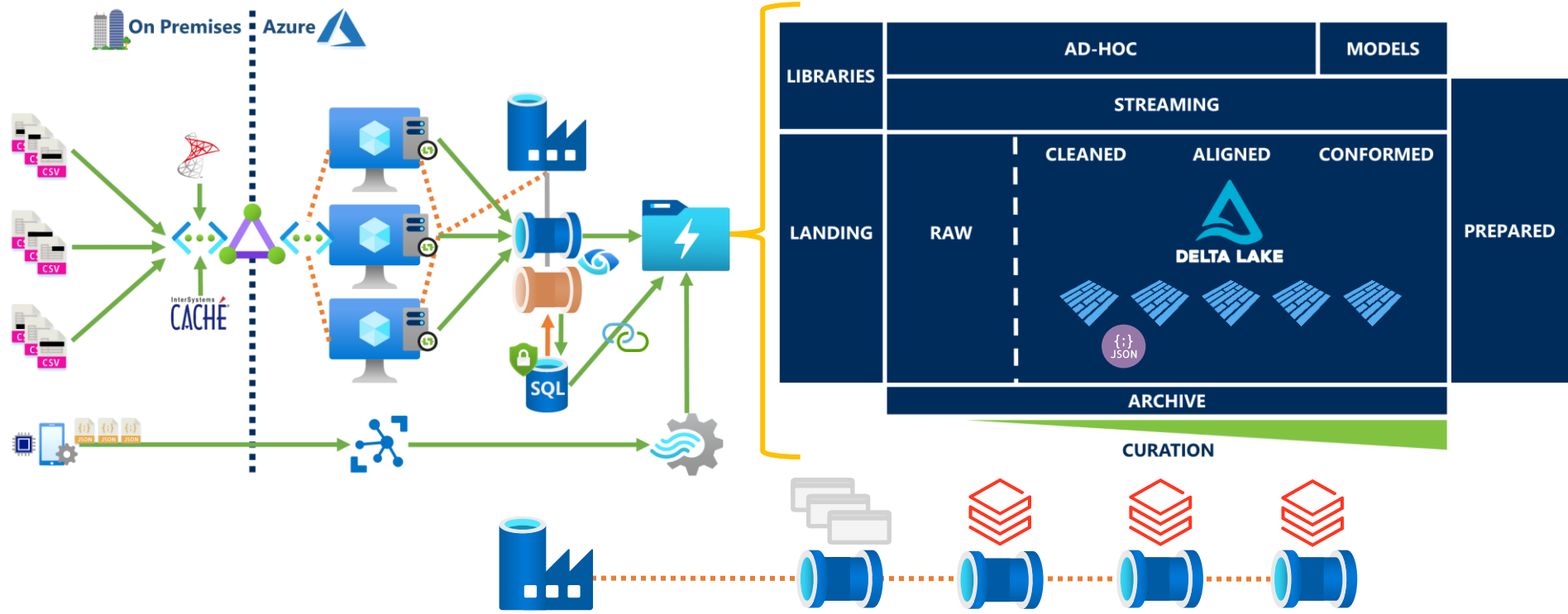
# Overall Architecture



## Extract

## Transform

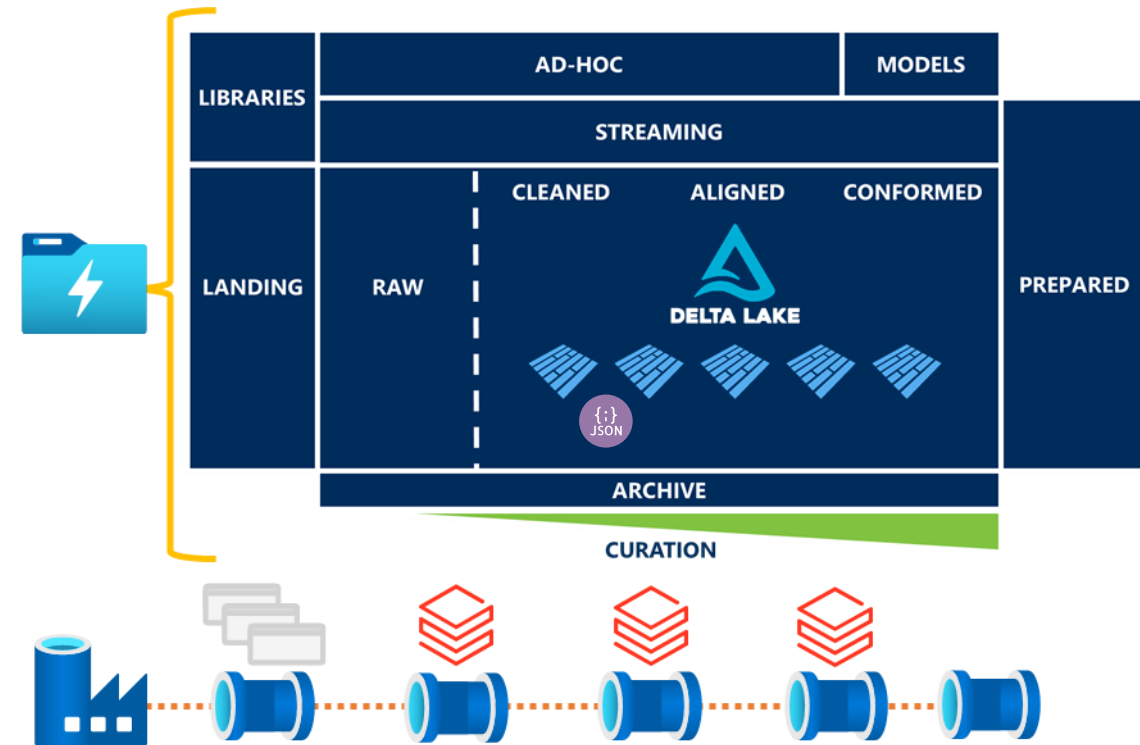
## Load



# Agenda



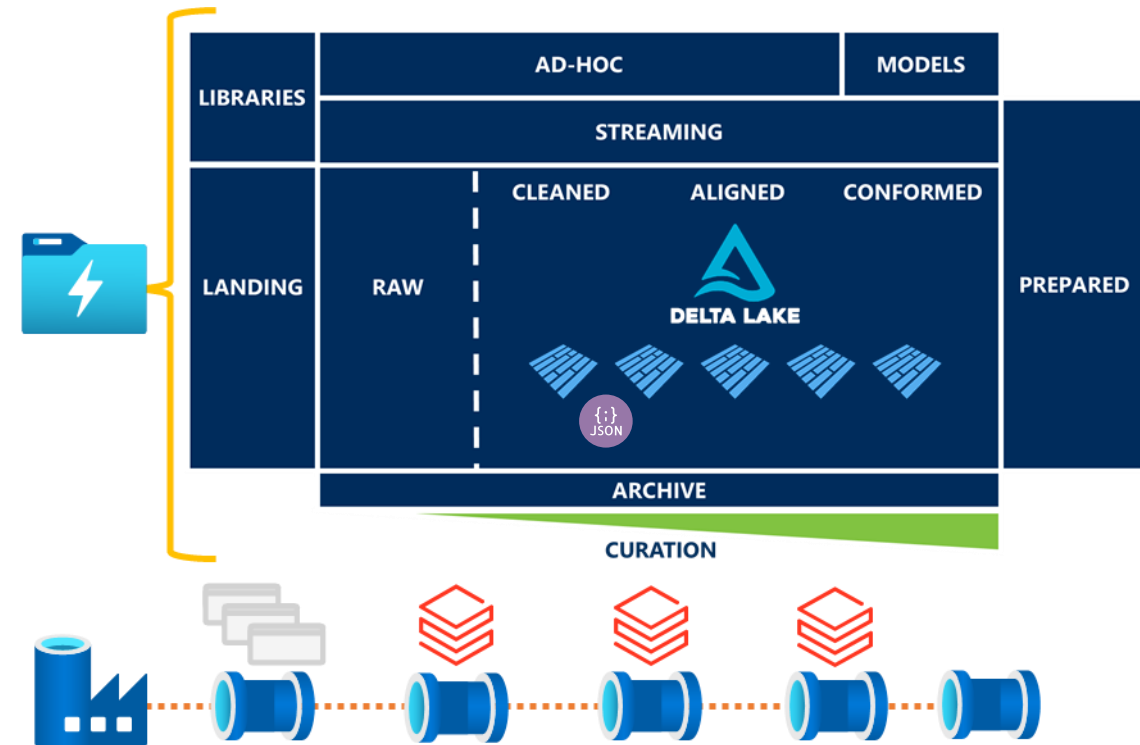
1. Design ✓
2. Extract ✓
3. Transform ✓
4. Load



# Agenda

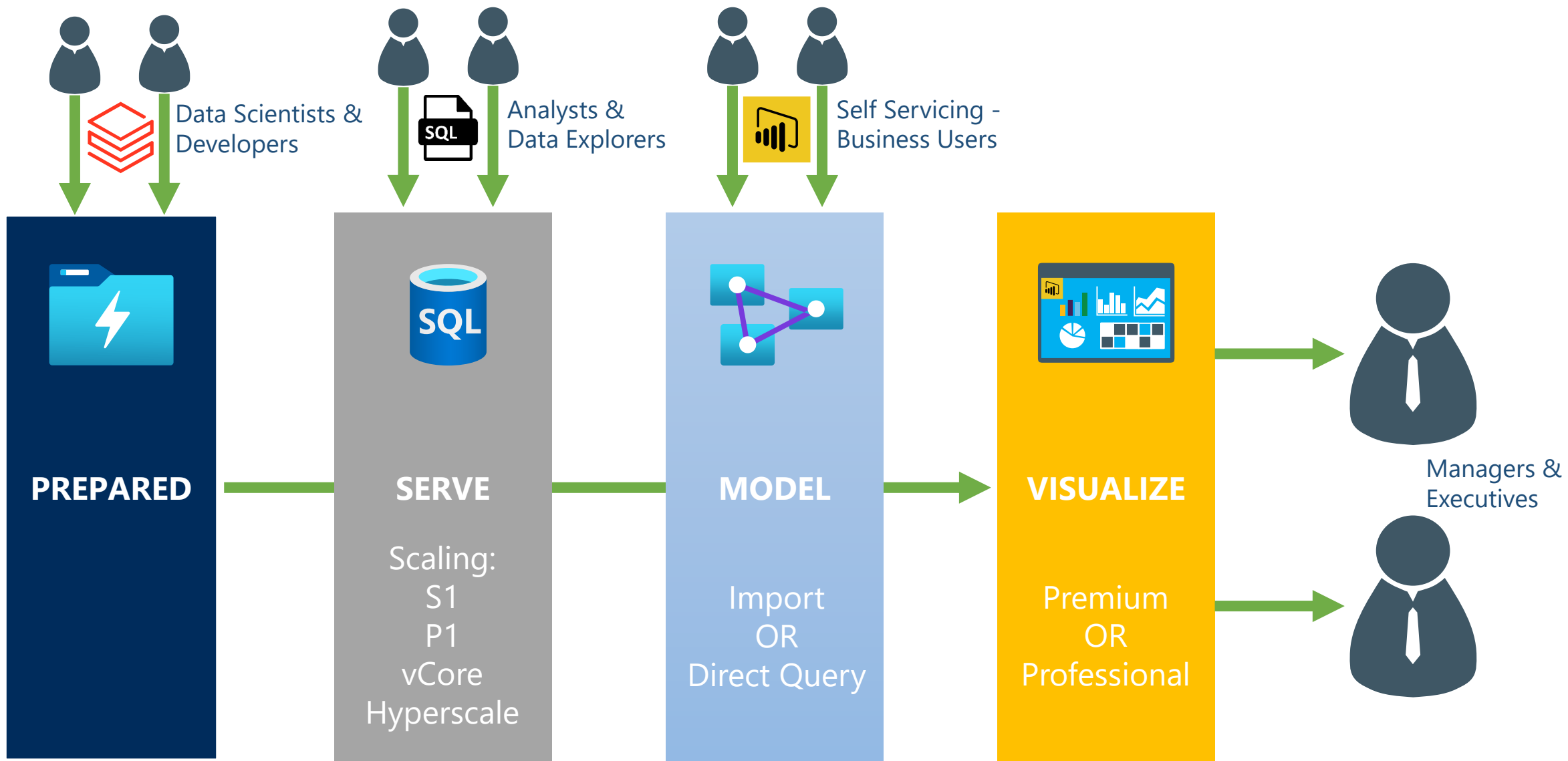
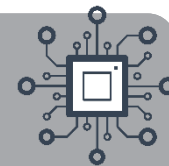


- 1. Design ✓
- 2. Extract ✓
- 3. Transform ✓
- 4. Load



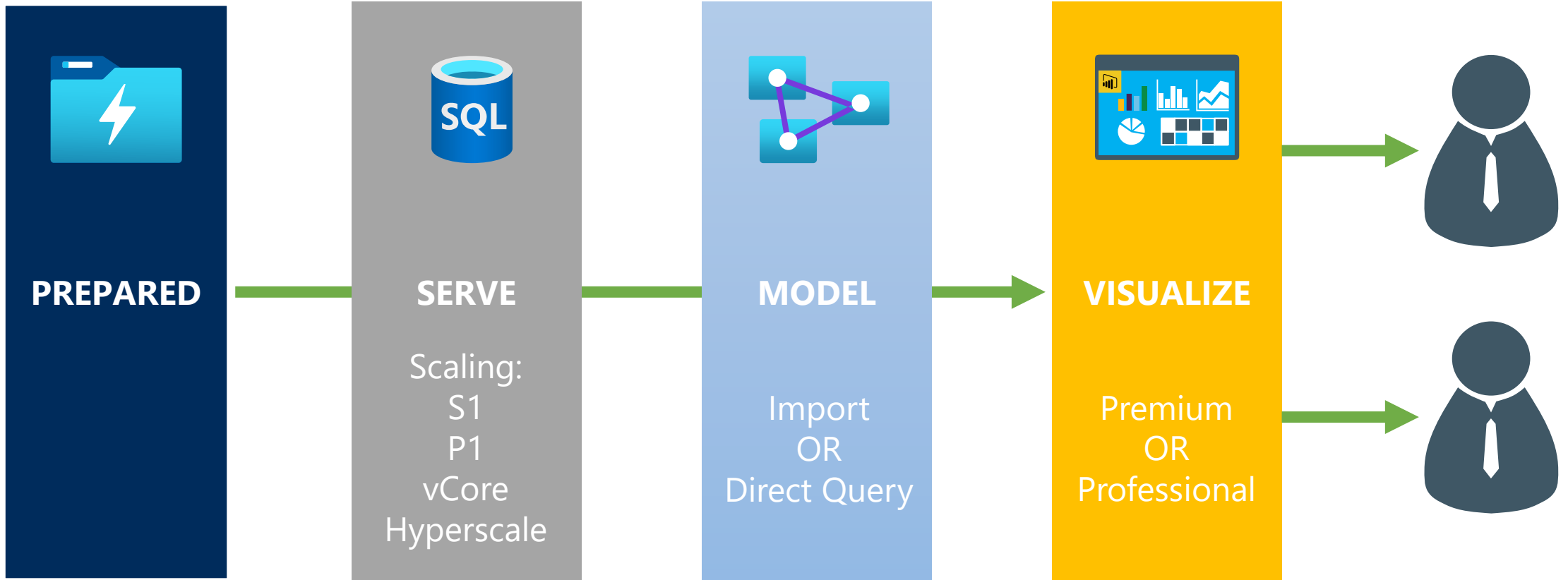


# Loading & Consuming Data



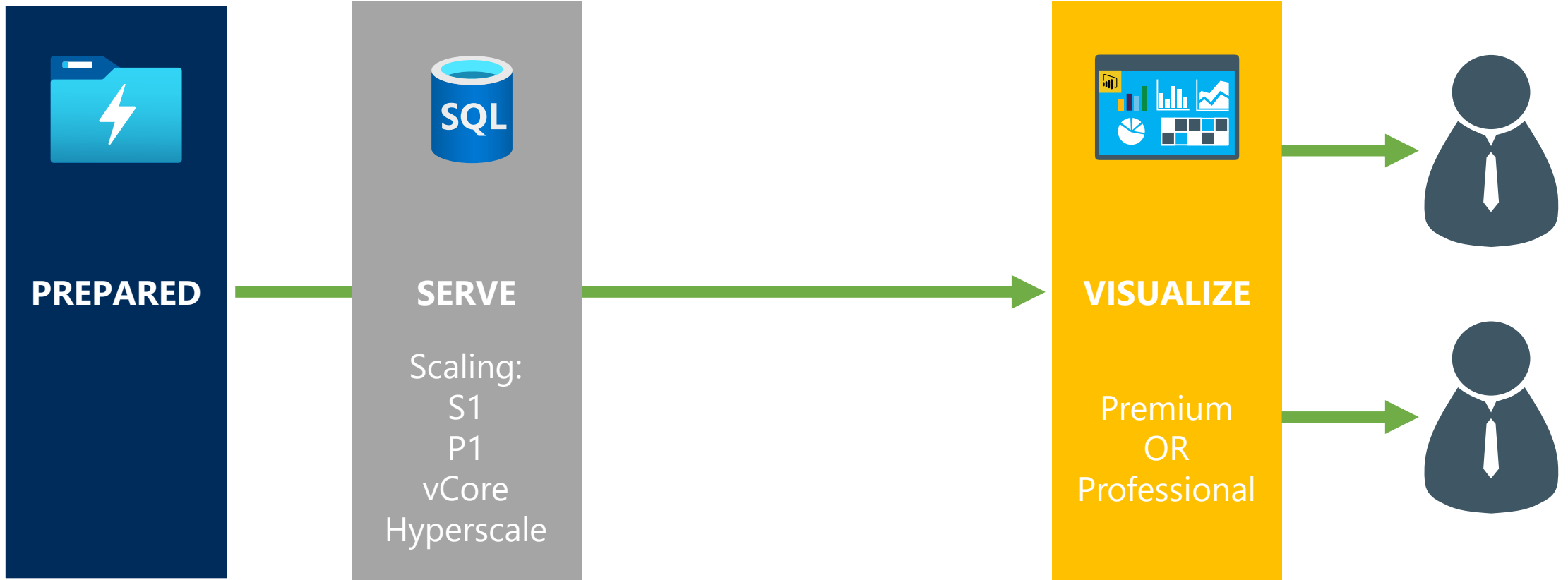


# Loading & Consuming Data



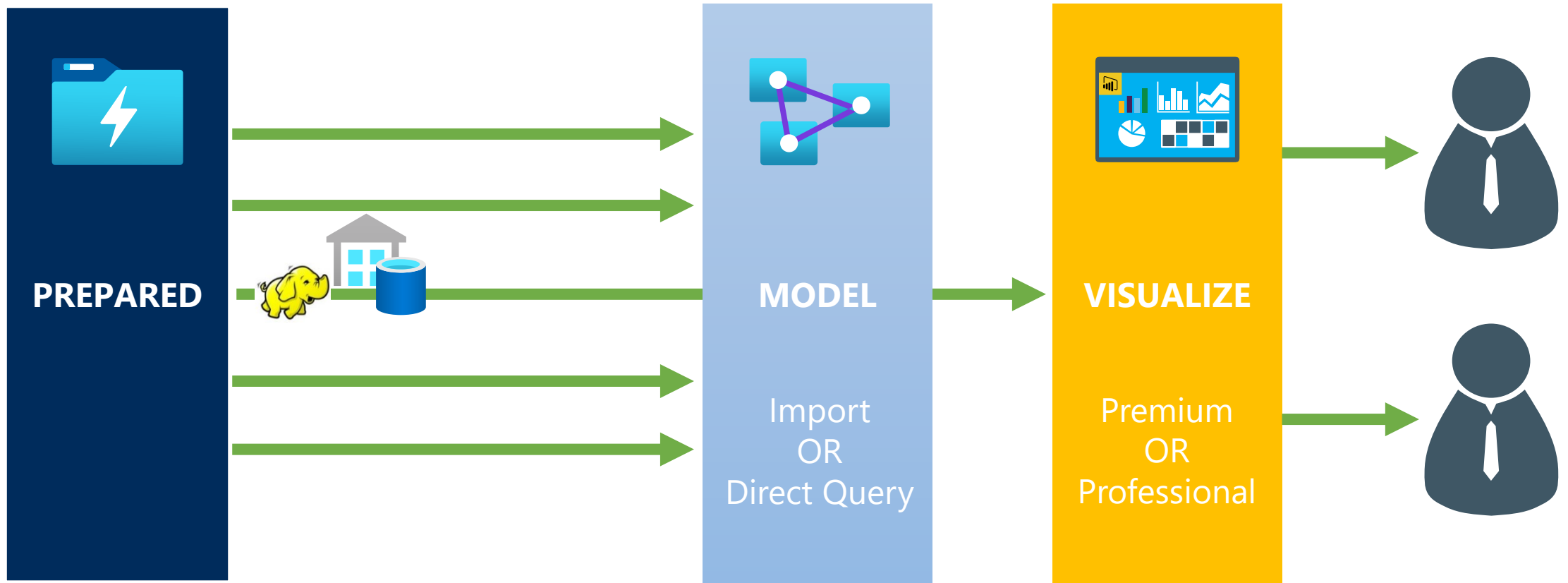


# Loading & Consuming Data





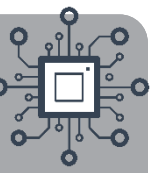
# Loading & Consuming Data








# Loading & Consuming Data



**PREPARED**



DELTA LAKE



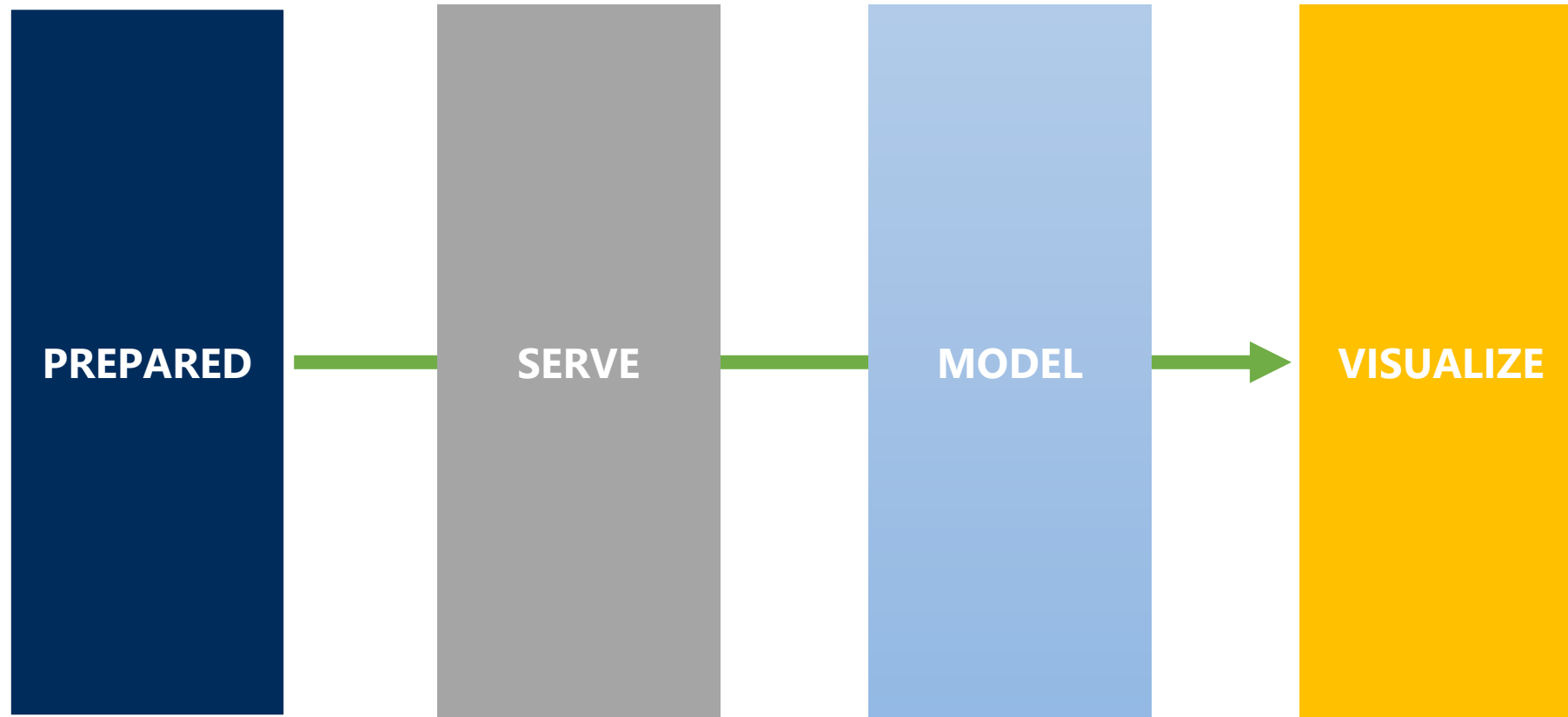
**VISUALIZE**

Premium  
OR  
Professional



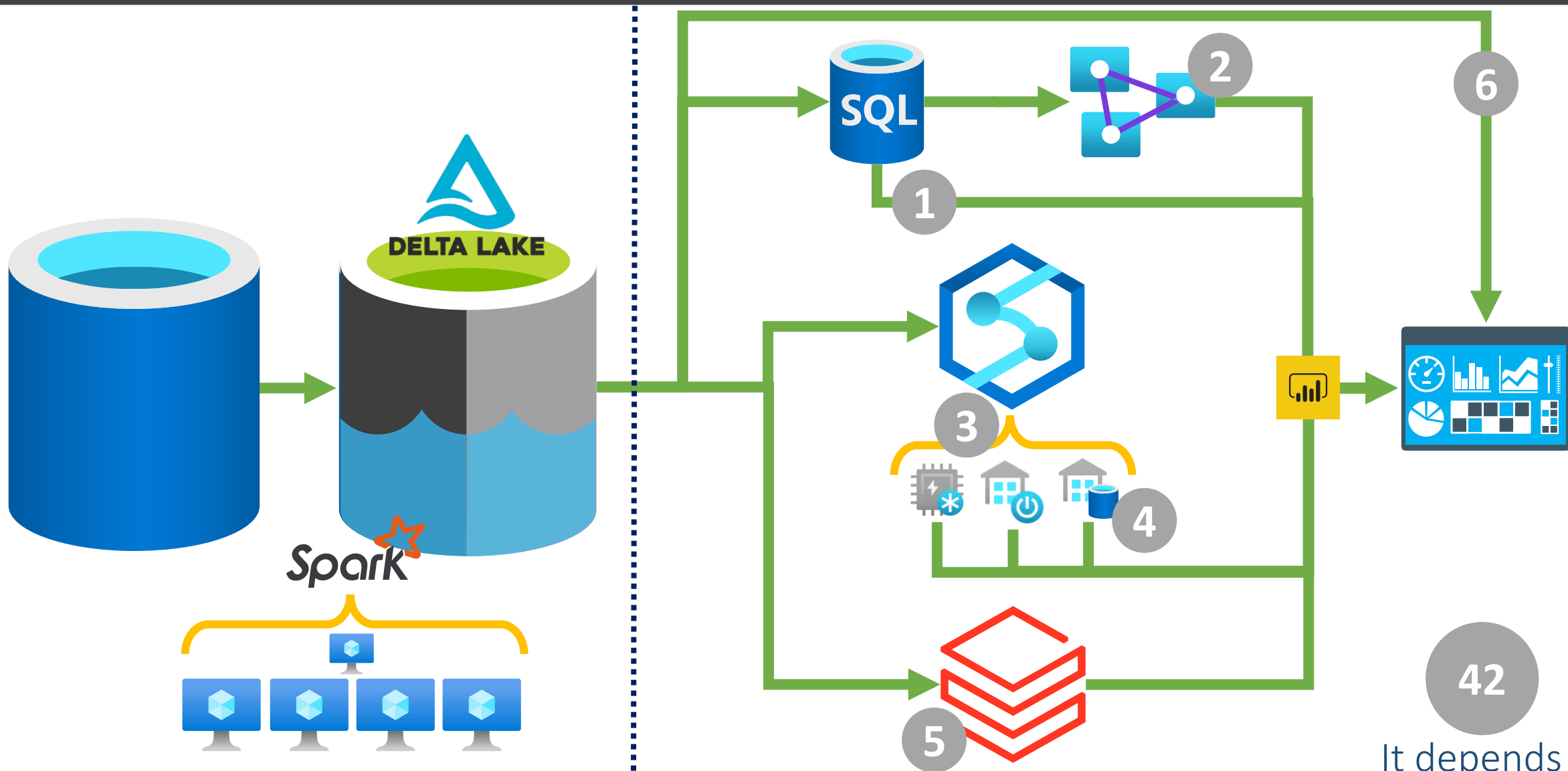
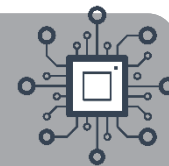


# Consuming Our Lake House in Azure



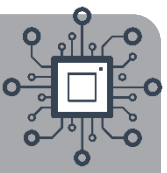


# Consuming Our Lake House in Azure





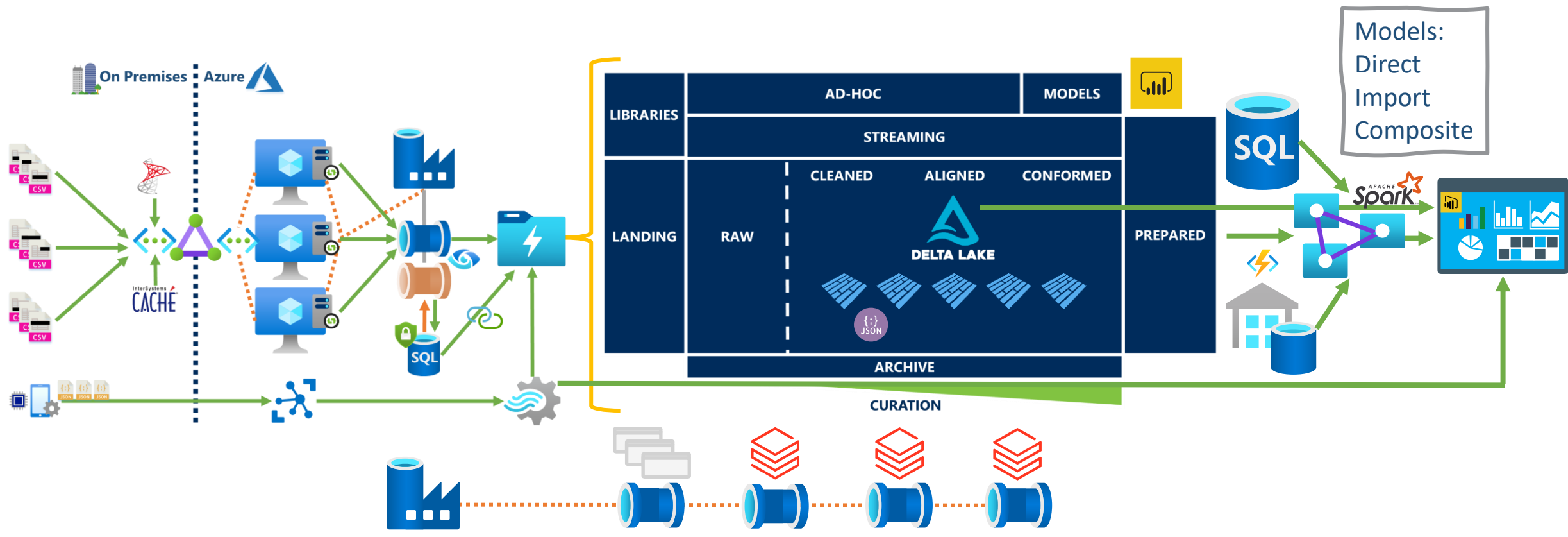
# Overall Architecture



## Extract

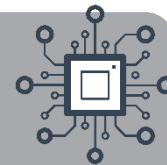
## Transform

## Load





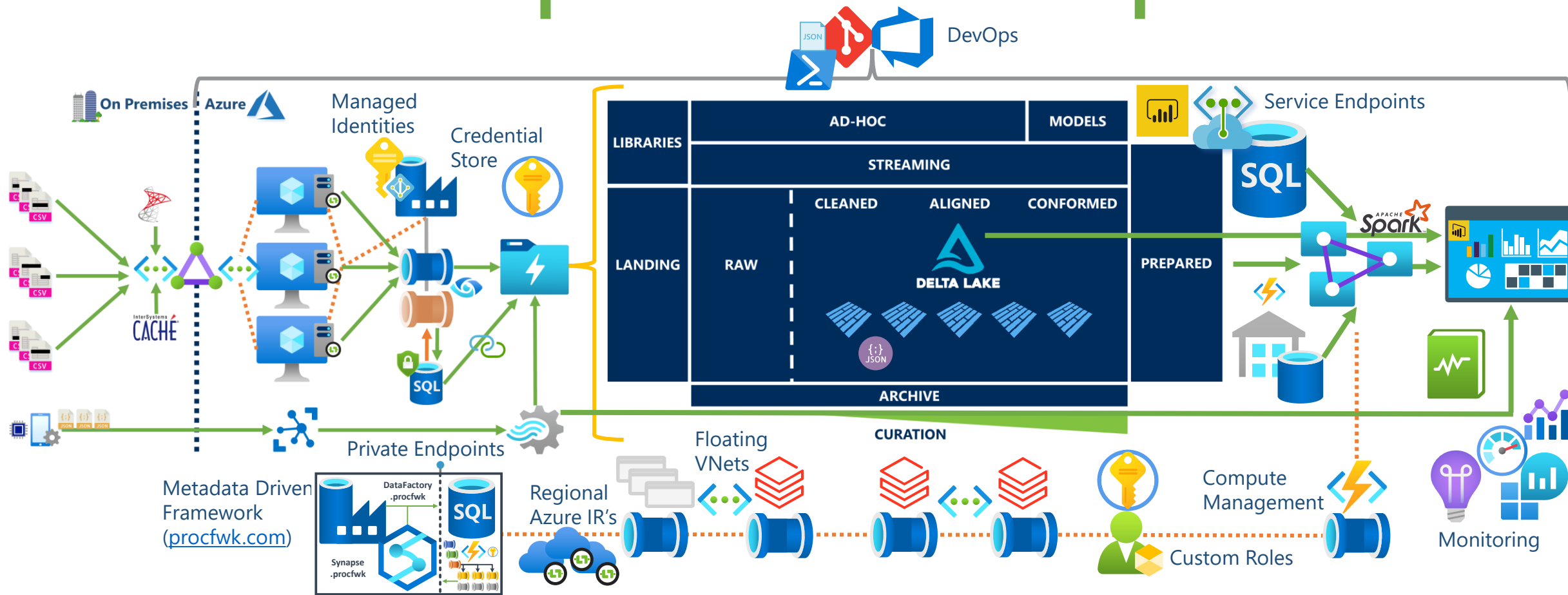
# Overall Architecture



## Extract

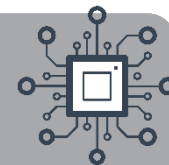
## Transform

## Load





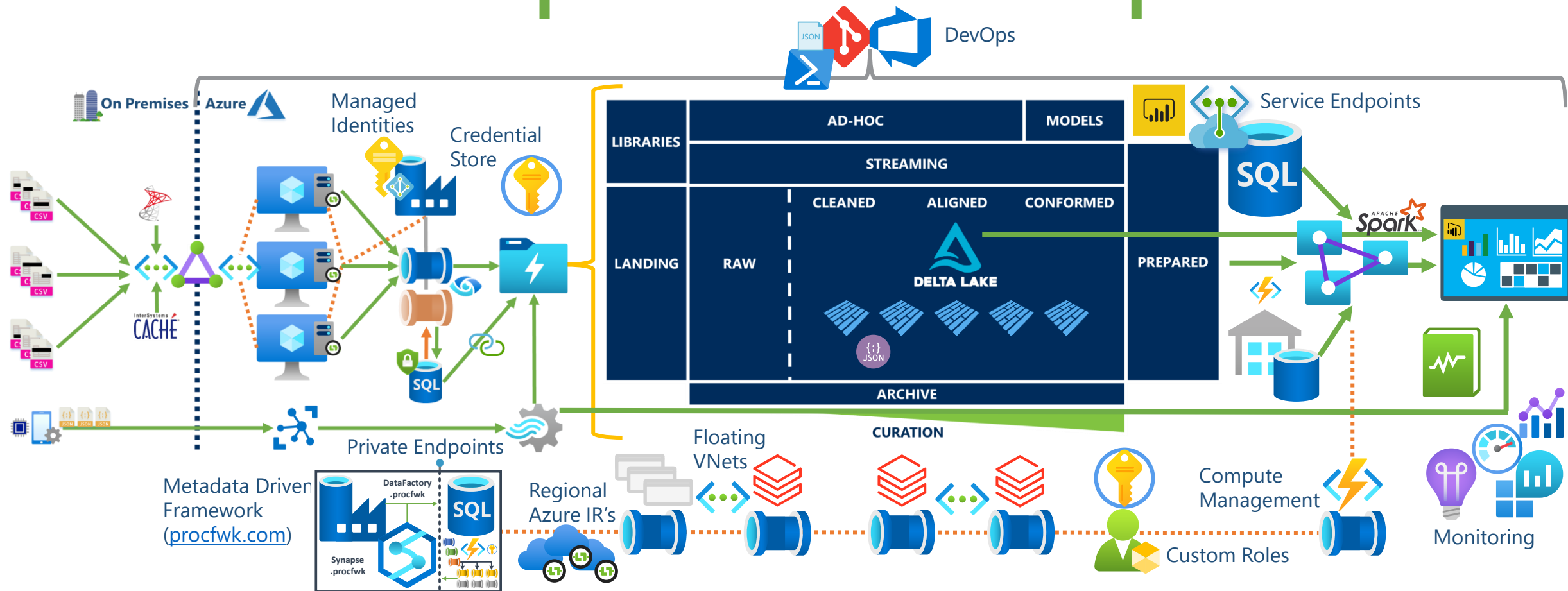
# Overall Architecture



## Extract

## Transform

## Load



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

# Module 13 - *Bonus*

An Architects Recap



```
SELECT
    *
FROM
    [Training]
WHERE
    [Module]
    BETWEEN 1 AND 12;
```

```
END;
```

```
RETURN; --complete
```