# Module 1 - 6

An Architects Recap

```sql
BEGIN --get ready
    SELECT
        [Summary]
    FROM
        [Training]
    WHERE
        [Module]
        BETWEEN 1 AND 6;
```

# Module 1 to 6 Recap

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

# Agenda

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

# Question:
What is the answer to life, the universe and everything?

**Answer:**
42



**Answer:**
It depends!

**Question:**
What is big data?

**Answer:**
It depends! ✓

**Answer:**
Any data that you cannot process
        in the time that you have/want ✓
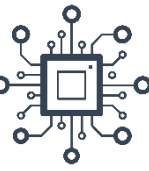                using the technology you have.

*- Buck Woody*

*@BuckWoodyMSFT*

# 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
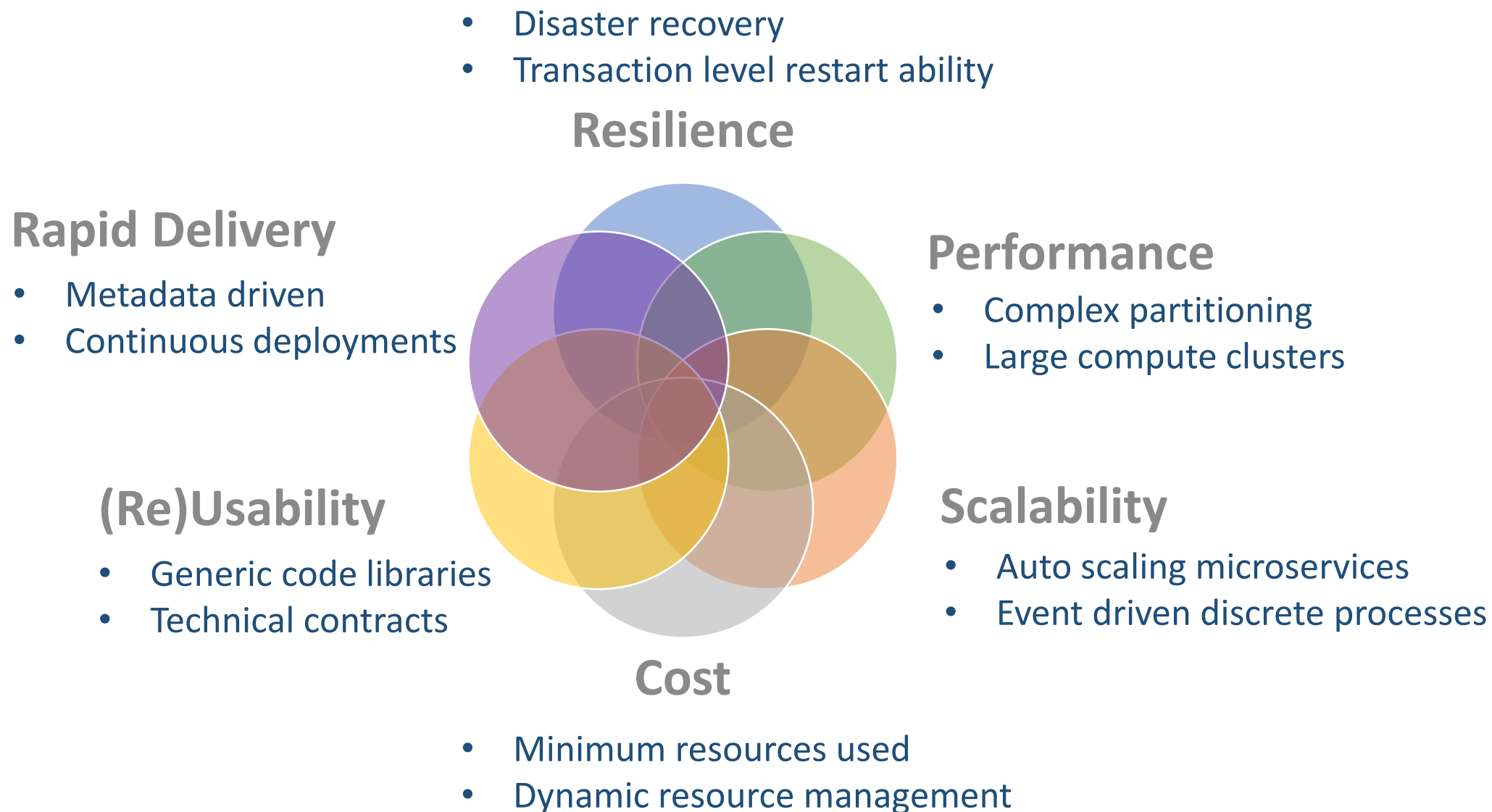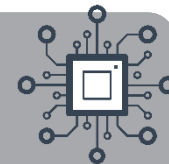Conform
Translate
Transform
Curate
Analyse
Model
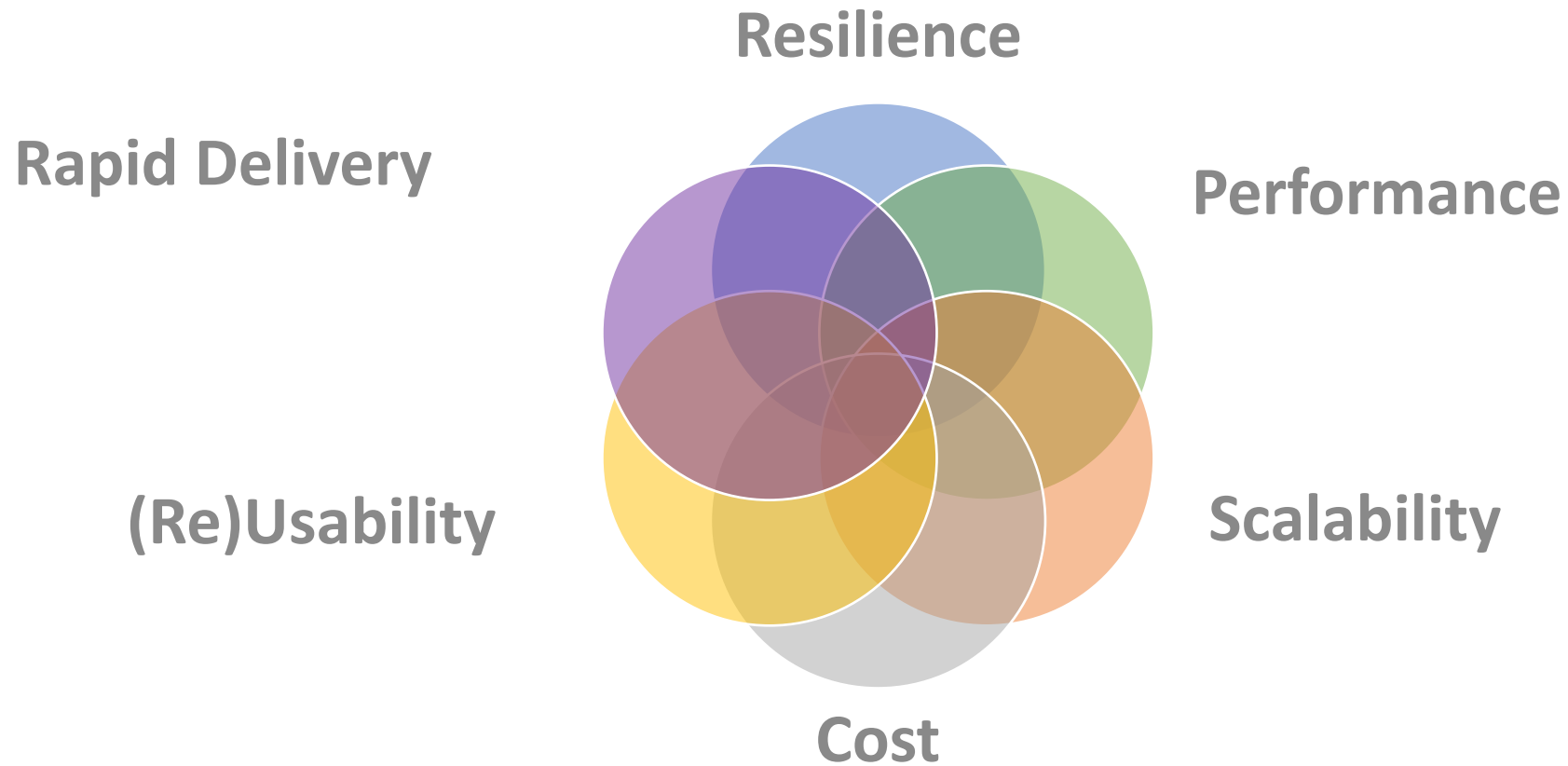Predict
Master

Data
Sources

Data
Warehouse

Data
Insights
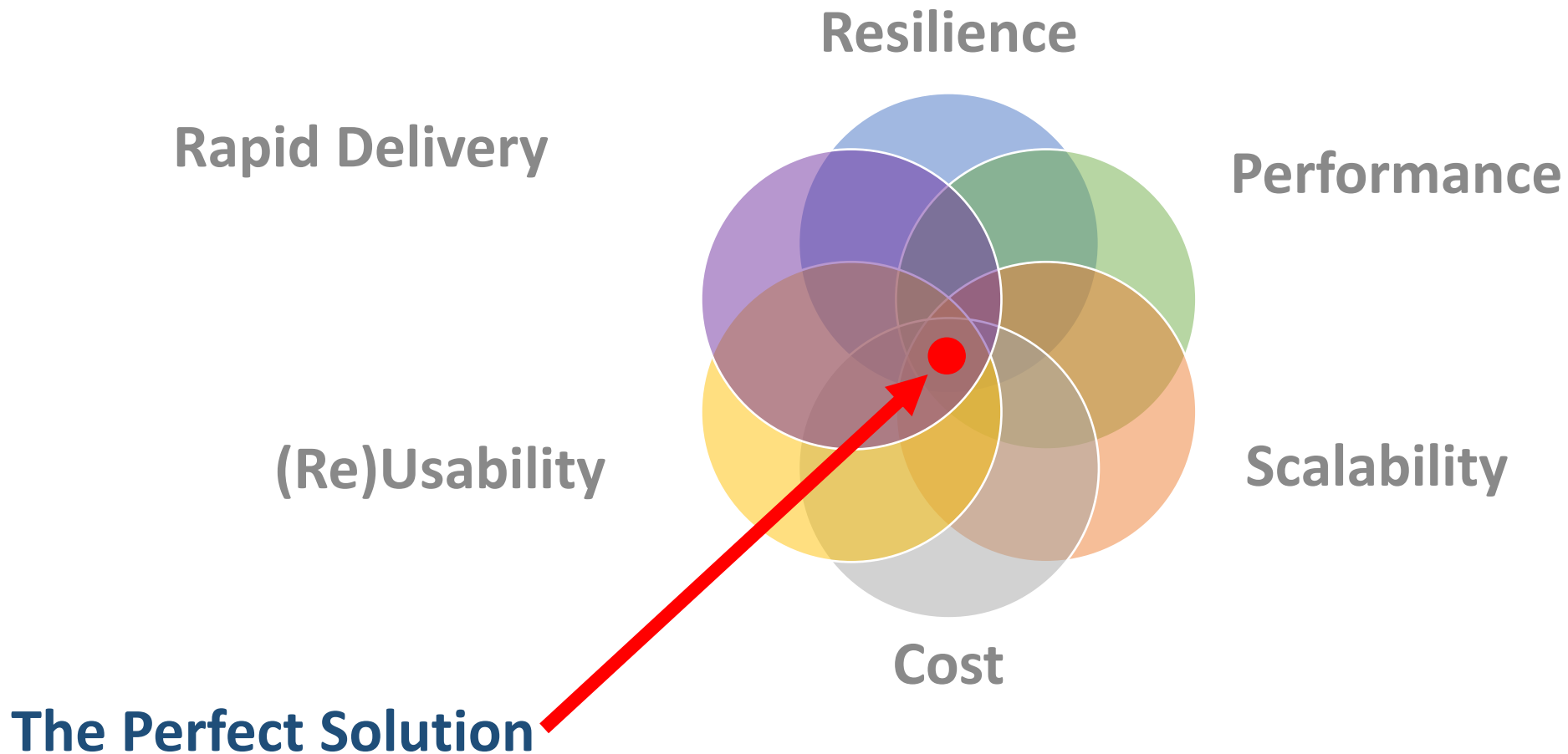
# What is your primary design focus?

**Resilience**

- Disaster recovery
- Transaction level restart ability

**Performance**

- Complex partitioning
- Large compute clusters

**Rapid Delivery**

- Metadata driven
- Continuous deployments

**Scalability**

- Auto scaling microservices
- Event driven discrete processes

**(Re)Usability**

- Generic code libraries
- Technical contracts

**Cost**

- Minimum resources used
- Dynamic resource management

# What is your primary design focus?

# Agenda

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

Resilience

Rapid Delivery

Performance

(Re)Usability

Scalability

Cost

# Agenda

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



Resilience

Rapid Delivery

Performance

(Re)Usability

Scalability

Cost

# Data Extraction & Ingestion

## Data Structure

CSV  TXT  XLS
HTML  TAR  JSON
XML  DAT  PNG  ZIP

## Data Source

SAP  salesforce

## Push or Pull

## Batch or Speed

## Public or Private Transfer

## Data Sensitivity

## Data Volume

!= Big

== Big

=> Big

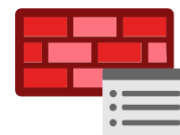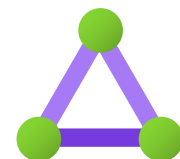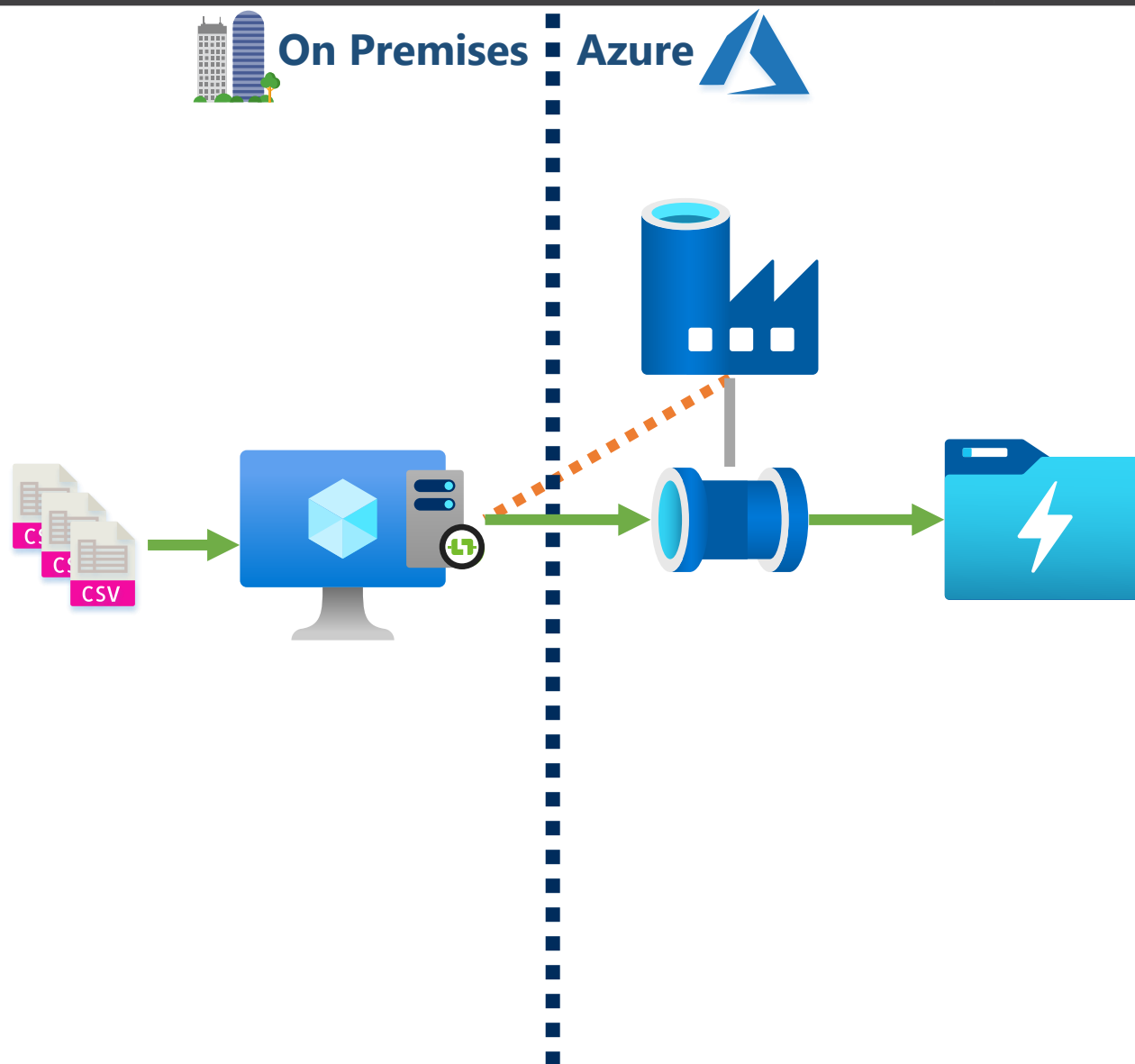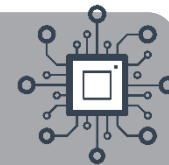# Data Extraction & Ingestion – Spec v1

## Data Structure

- CSV
- TXT
- XLS
- HTML
- TAR
- JSON
- XML
- ZIP
- DAT
- PNG

## Data Source

- SAP
- salesforce

## Push or Pull

## Batch or Speed

## Public or Private Transfer

## Data Sensitivity

## Data Volume
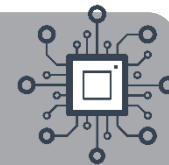
- != Big
- == Big
- => Big

**On Premises** | **Azure**

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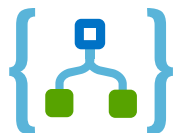- No PII data
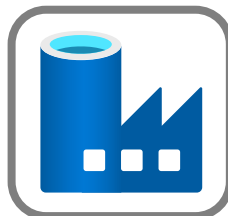- Small data volumes

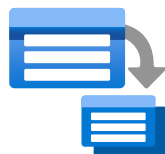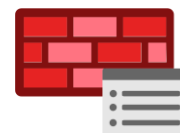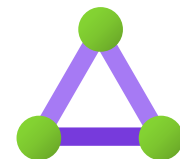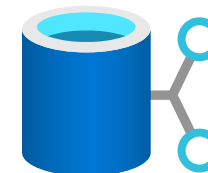# Data Extraction & Ingestion – Spec v2

## Data Structure

CSV  TXT  XLS
HTML  JSON
XML  TAR
DAT  PNG  ZIP

## Data Source

SAP
salesforce

## Push or Pull

## Batch or Speed
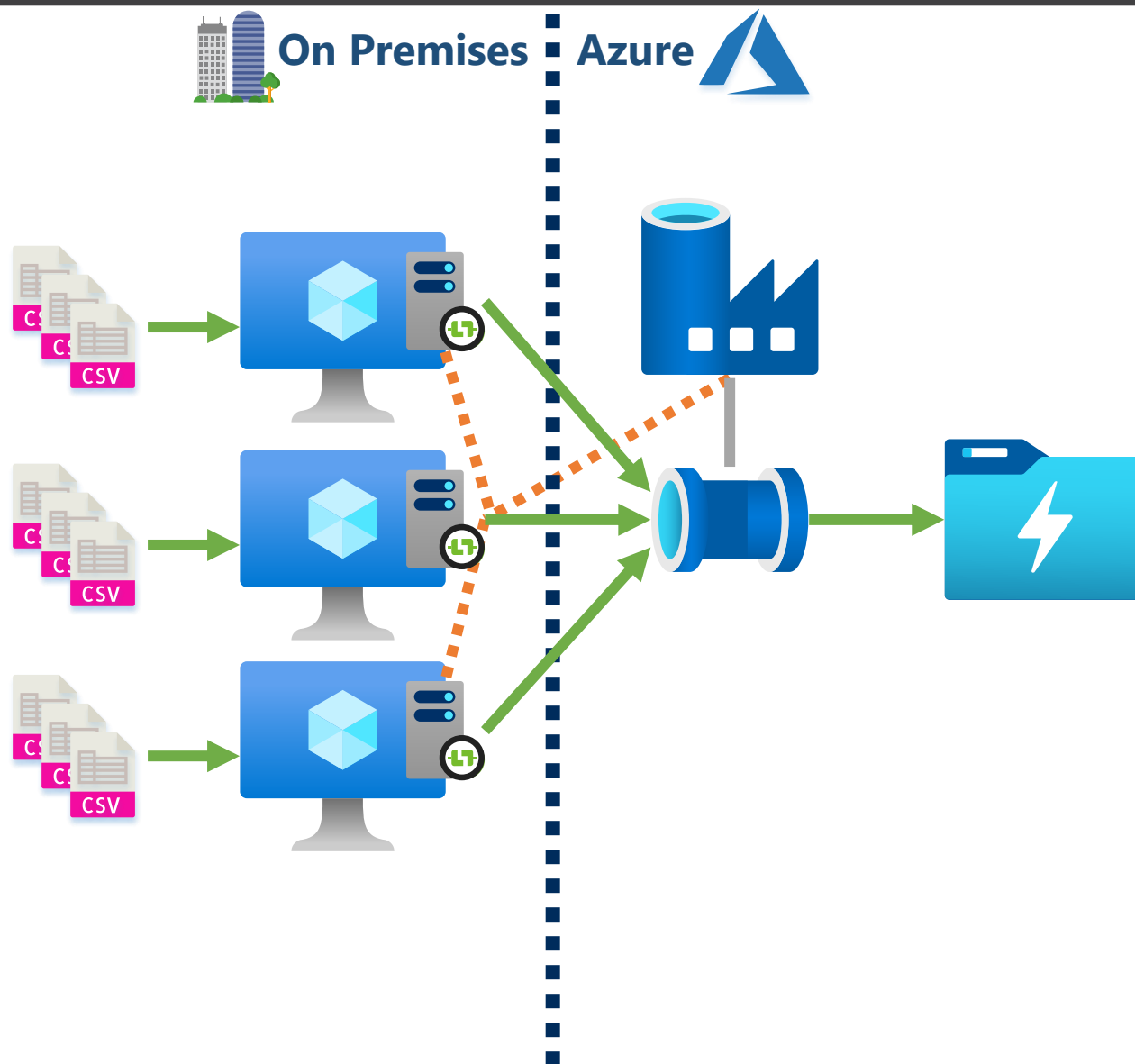
## Public or Private Transfer

## Data Sensitivity

## Data Volume

!= Big

== Big

=> Big

# Data Extraction & Ingestion – Solution 2

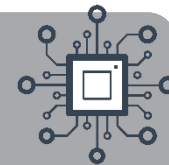**On Premises** | **Azure**

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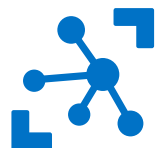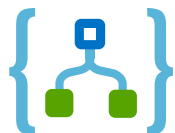- <u>Large</u> data volumes

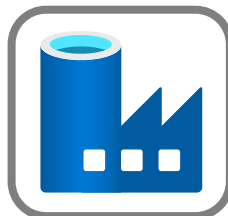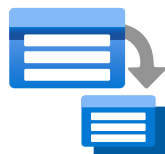# Data Extraction & Ingestion – Spec v3
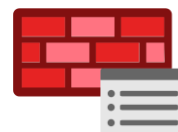
## Data Structure

CSV  TXT  XLS
HTML  JSON
XML  TAR
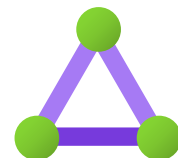DAT  PNG  ZIP

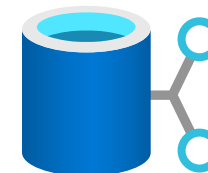## Data Source

SAP
salesforce

## Push or Pull

## Batch or Speed

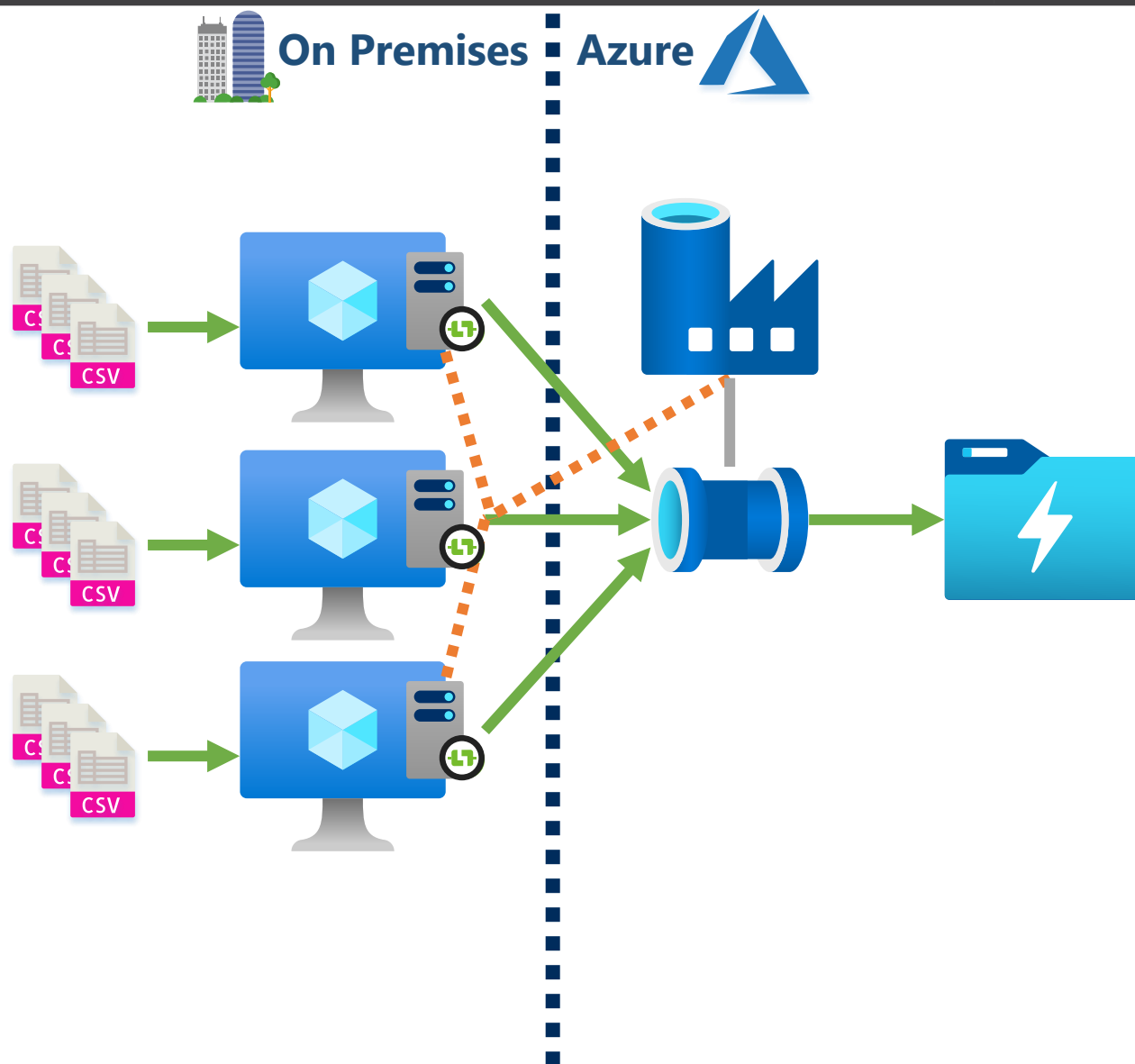## Public or Private Transfer

## Data Sensitivity

## Data Volume

!= Big

== Big

=> Big

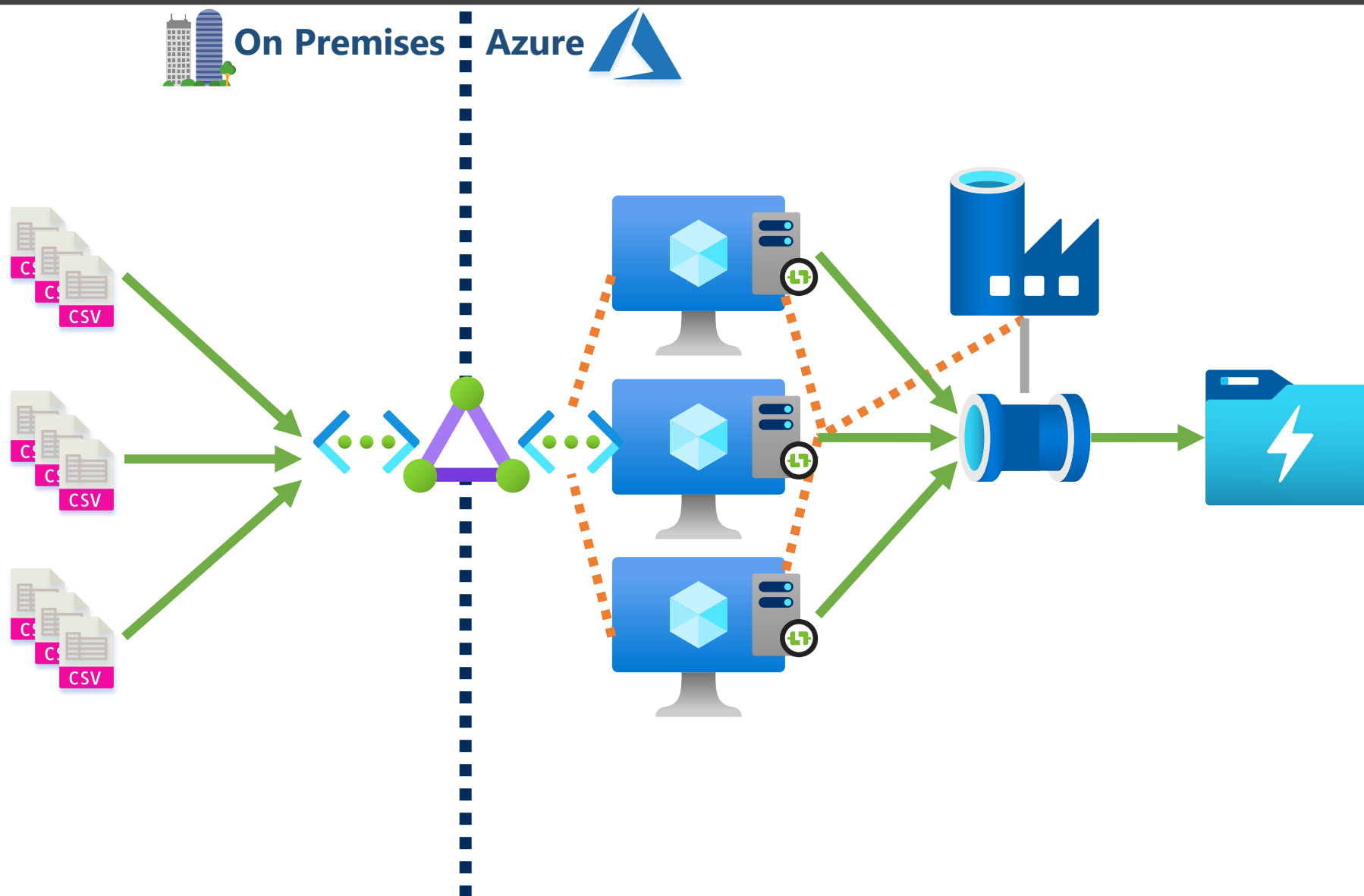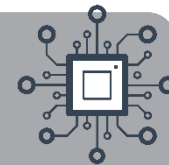# Data Extraction & Ingestion – Solution 3



**On Premises**     **Azure**

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

# Data Extraction & Ingestion – Solution 3
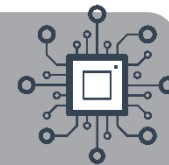
**On Premises** | **Azure**

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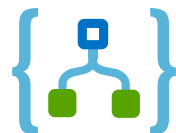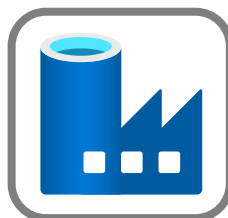- Large data volumes

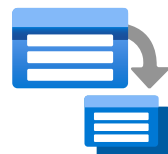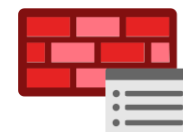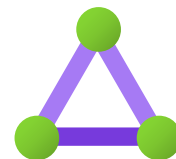# Data Extraction & Ingestion – Spec v4
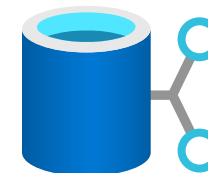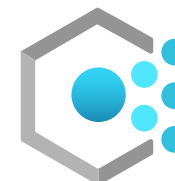
| Data Structure | Push or Pull | Batch or Speed | Public or Private Transfer | Data Sensitivity | Data Volume |
|---|---|---|---|---|---|

**Data Structure**

CSV, TXT, XLS, HTML, JSON, XML, TAR, DAT, PNG, ZIP

**Data Source**

SAP, salesforce, Twitter

!= Big

== Big

=> Big

# Data Extraction & Ingestion – Solution 4

**On Premises** | **Azure**

**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

CSV, TXT, XLS, HTML, JSON, XML, TAR, ZIP, DAT, PNG

## Data Source

SAP, salesforce

## Push or Pull

## Batch or Speed
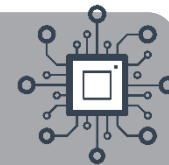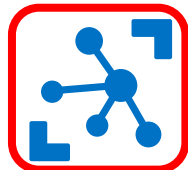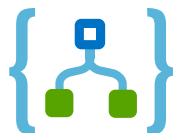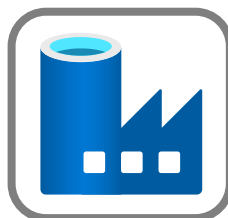
## Public or Private Transfer

## Data Sensitivity

## Data Volume

!= Big

== Big

=> Big

# Data Extraction & Ingestion – Solution 5

**On Premises** | **Azure**



**Requirements:**
- Flat files & JSON
- From local storage & database tables
- Pulled from source & pushed
- Batch load & streamed
- Private connections
- No PII data
- Large data volumes
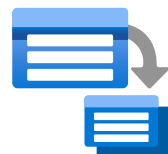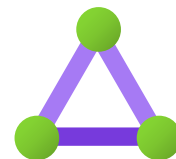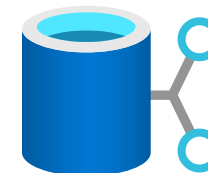
# Data Extraction & Ingestion – Spec v6

## Data Structure

CSV · TXT · XLS · HTML · JSON · XML · TAR · ZIP · PNG · DAT

## Data Source

SAP · salesforce

## Push or Pull

## Batch or Speed

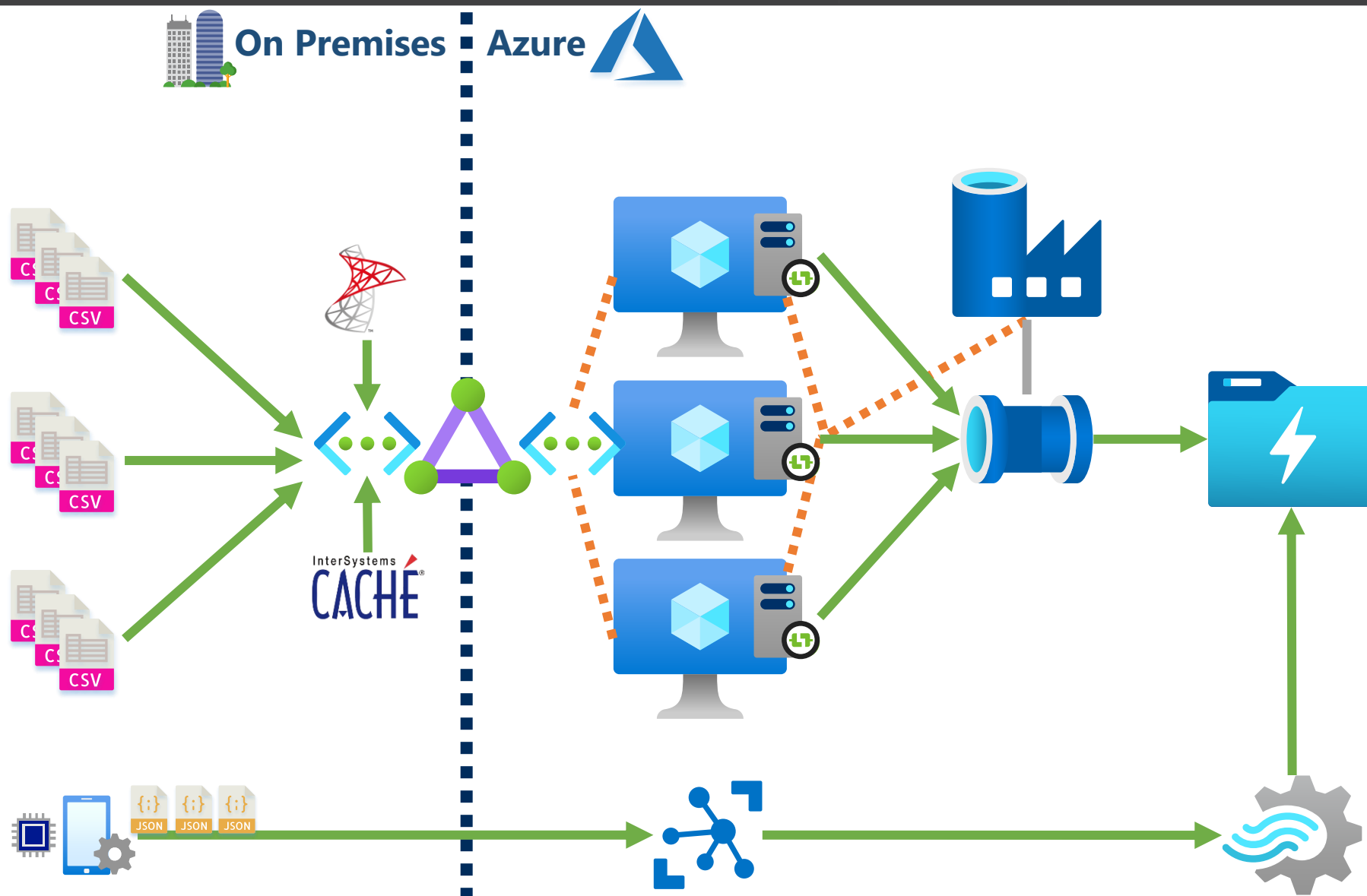## Public or Private Transfer

## Data Sensitivity

## Data Volume
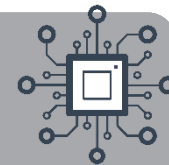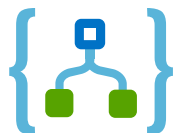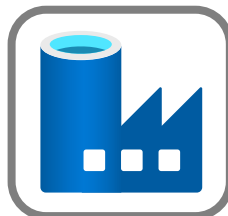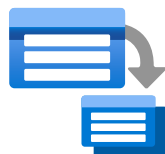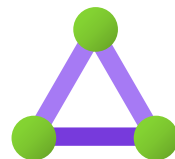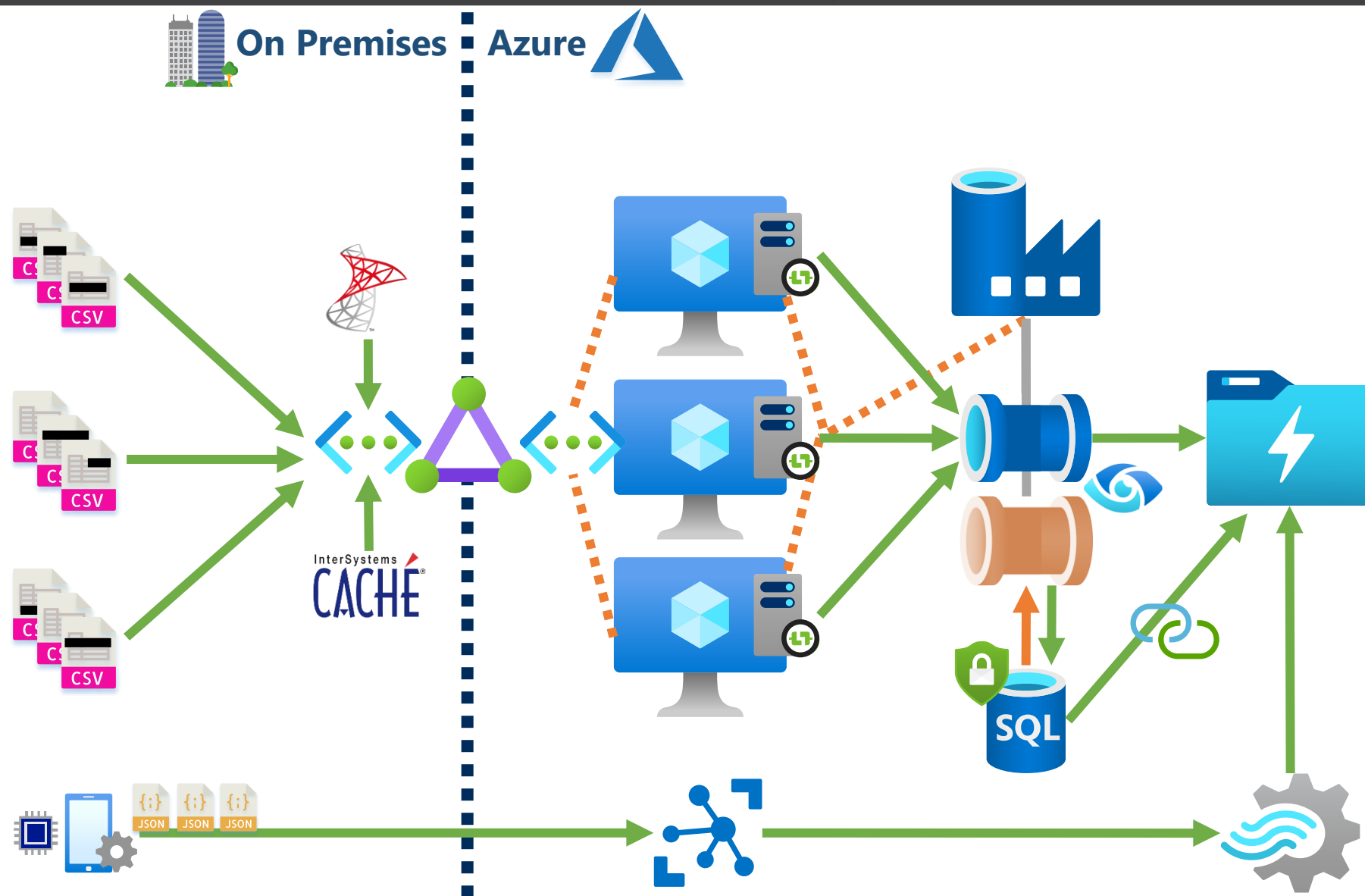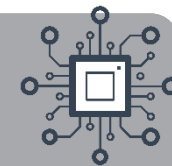
!= Big

== Big

=> Big

# Data Extraction & Ingestion – Solution 6

**On Premises** **Azure**

CSV
CSV
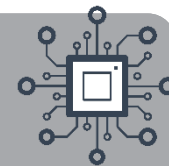CSV

InterSystems CACHÉ®

JSON JSON JSON

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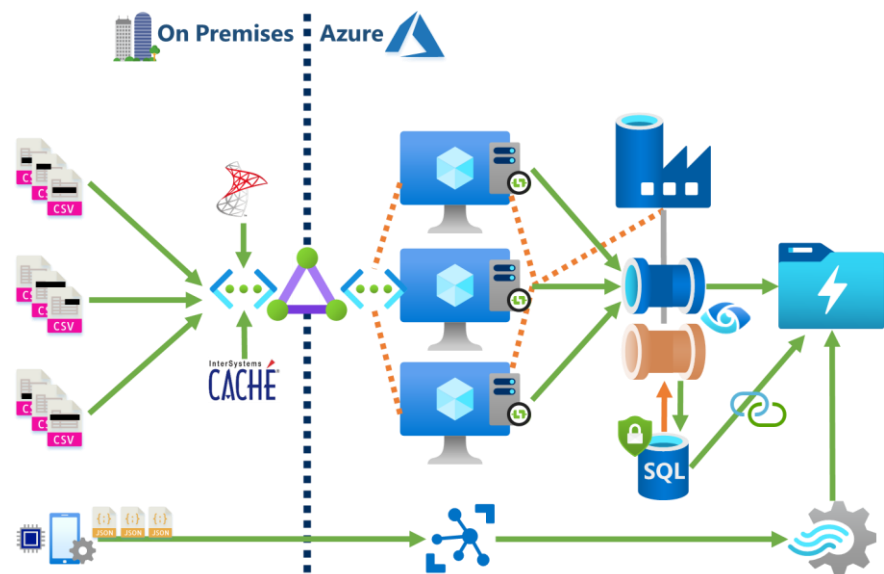- Large data volumes
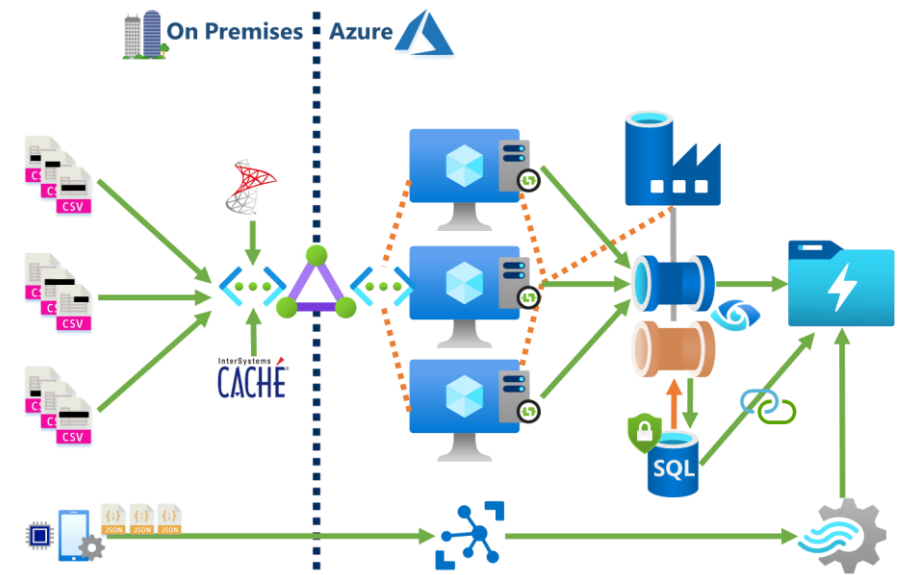
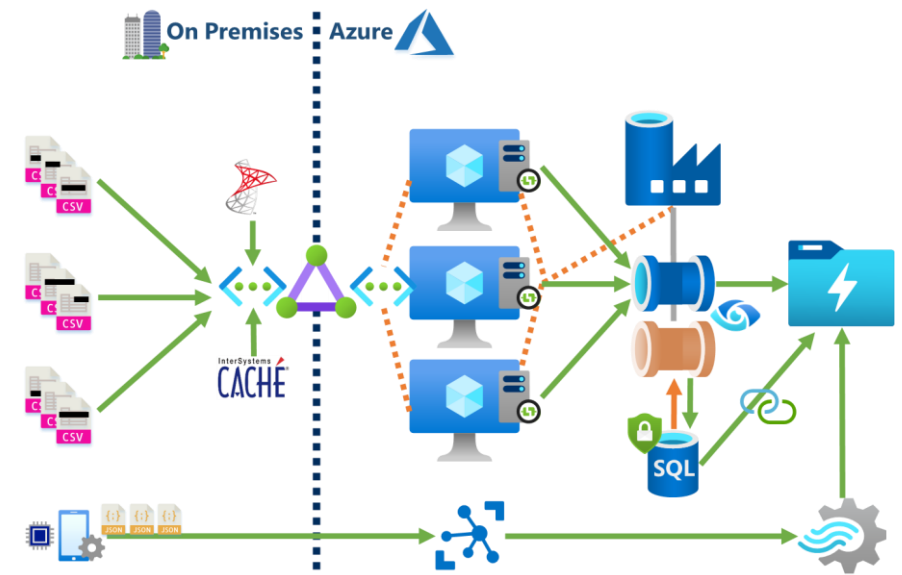SQL

# Extract | Transform | Load

# Agenda



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

# Agenda

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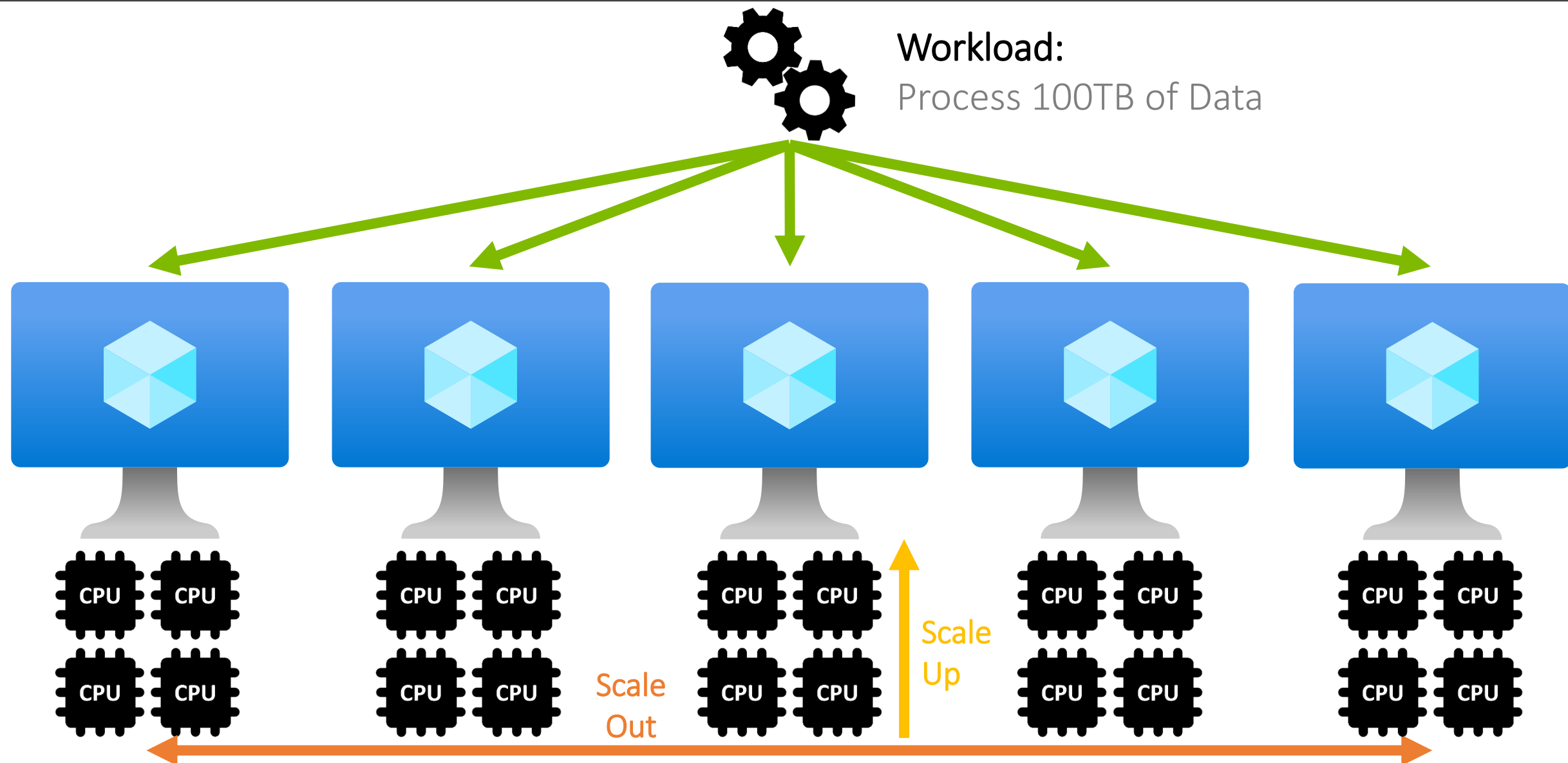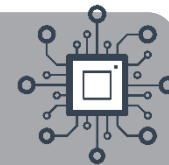2. Extract ✓
3. Transform
4. Load

# Agenda

1. Design ✓
2. Extract ✓
3. Transform
4. Load
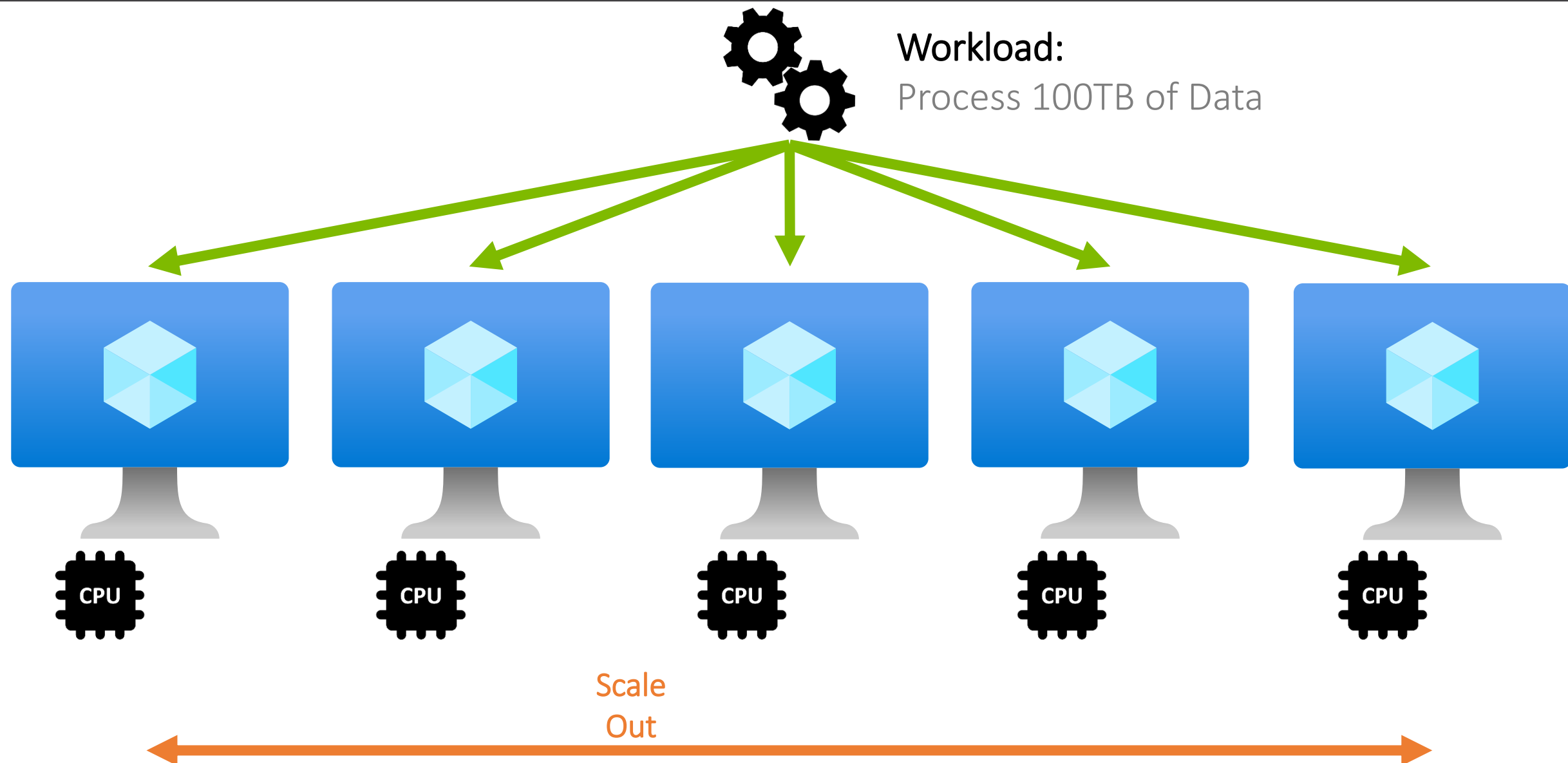
Compute
Storage, Structure
& Data Format
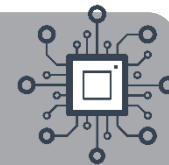
# Scaling Up and/or Scaling Out

**Workload:**
Process 100TB of Data
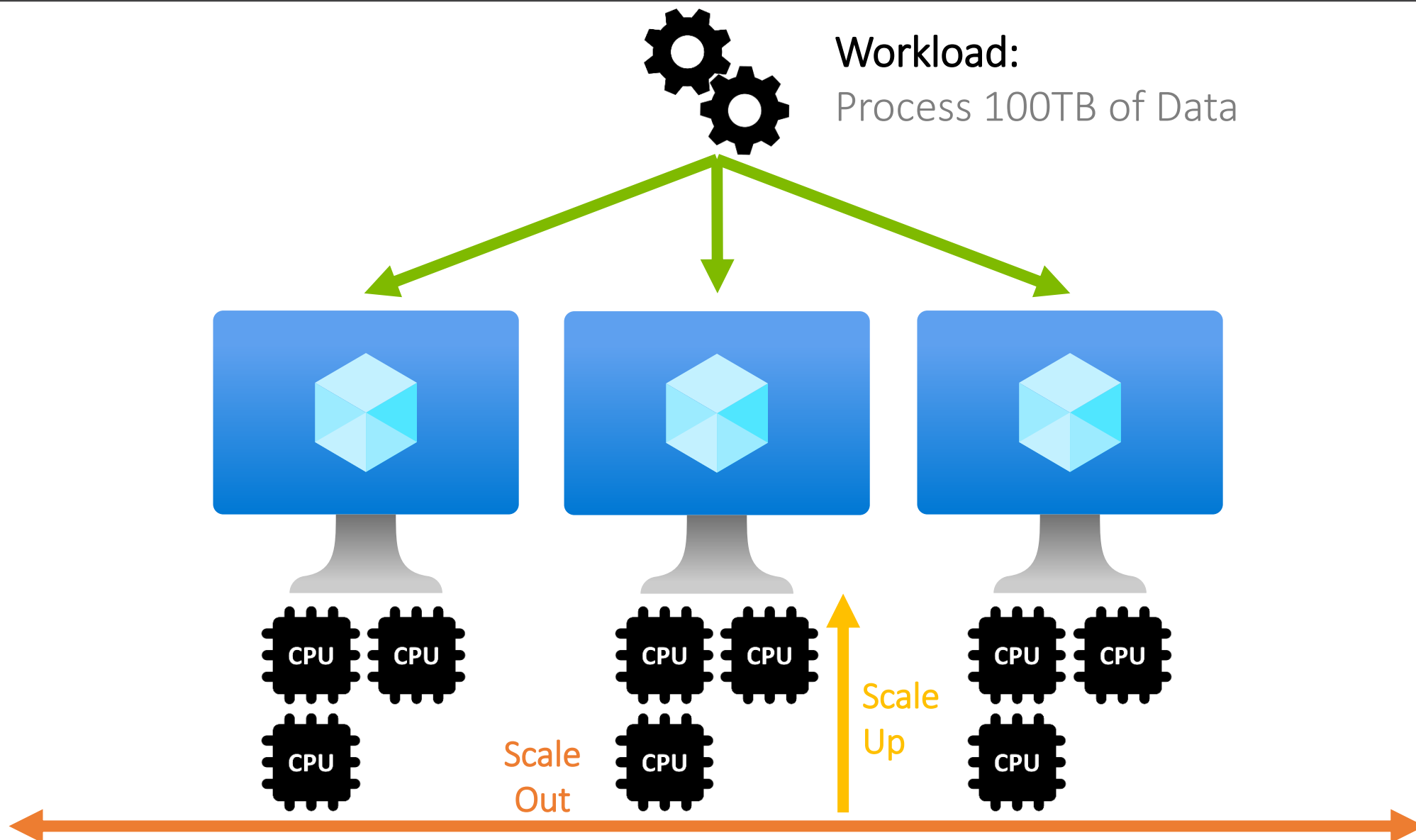
Scale Up

Scale Out

# Scaling Up and/or Scaling Out

Workload:

Process 100TB of Data

CPU    CPU    CPU    CPU    CPU

Scale Out

**Workload:**
Process 100TB of Data

CPU CPU CPU CPU CPU CPU

CPU CPU CPU

Scale Out

Scale Up

**Workload:**

Process 100TB of Data

**P**latform

**I**nfrastructure

As

A

Service

**Workload:**

Process 100TB of Data

**P**latform

As
A
Service

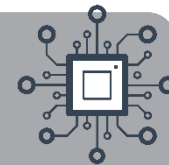| IaaS | PaaS |
|------|------|
| Applications | Applications |
| Data | Data |
| Runtime | Runtime |
| Middleware | Middleware |
| Operating System | Operating System |
| Virtualization | Virtualization |
| Servers | Servers |
| Storage | Storage |
| Networking | Networking |

# Data Transformation – Compute

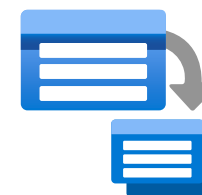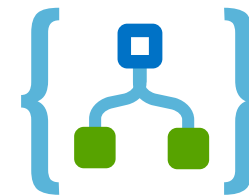| Data Lake Analytics | HDInsight | Relational Database | Synapse – SQL Pools or Spark Pools | Databricks | Batch Service | Data Explorer |
|---|---|---|---|---|---|---|

| Automation | Cosmos | Functions | Power BI Data Flows | Logic Apps | Data Factory Data Flows | Analysis Services |
|---|---|---|---|---|---|---|

# Data Transformation – Compute

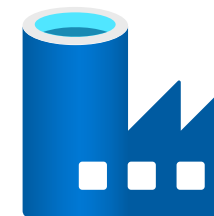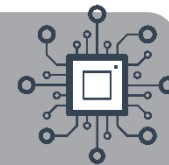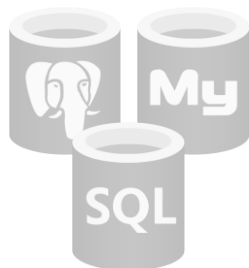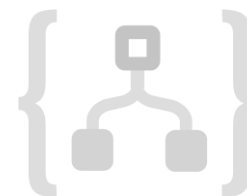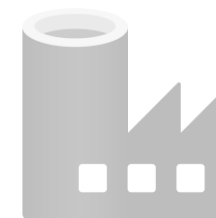| Data Lake Analytics | HDInsight | Relational Database | Synapse – SQL Pools or Spark Pools | Databricks | Batch Service | Data Explorer |
|---|---|---|---|---|---|---|

| 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**

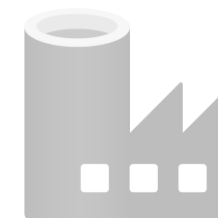**Automation**

**Cosmos**

**Functions**

**Power BI Data Flows**

**Logic Apps**

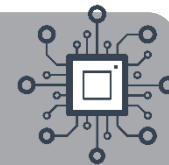**Data Factory Data Flows**

**Analysis Services**

# Agenda

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

Compute ✓
Storage, Structure & Data Format

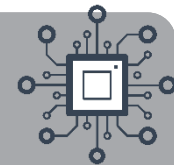# Data Transformation – Storage & Format

Azure Storage Account

Azure Data Lake Gen2
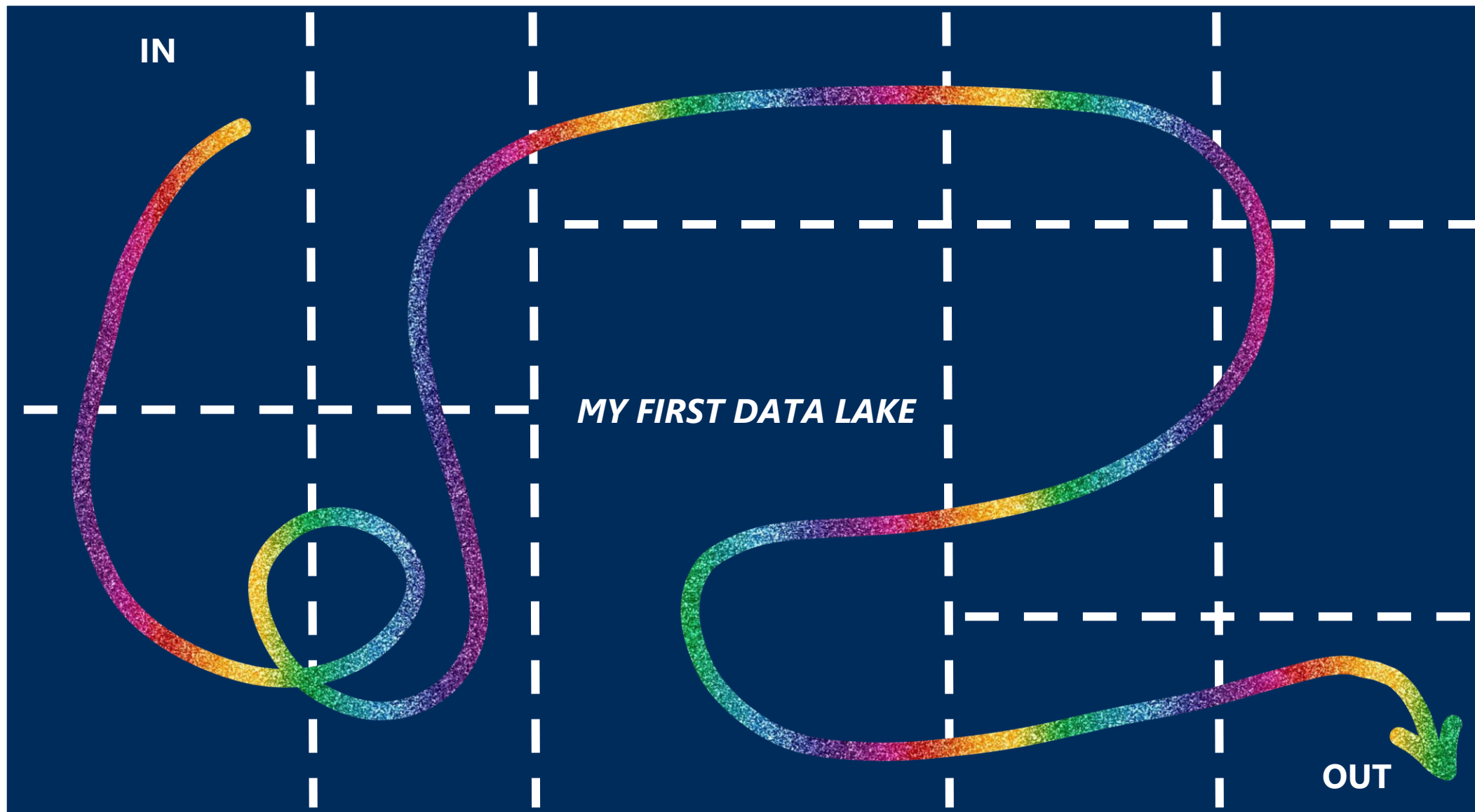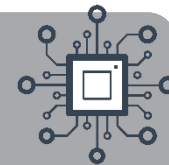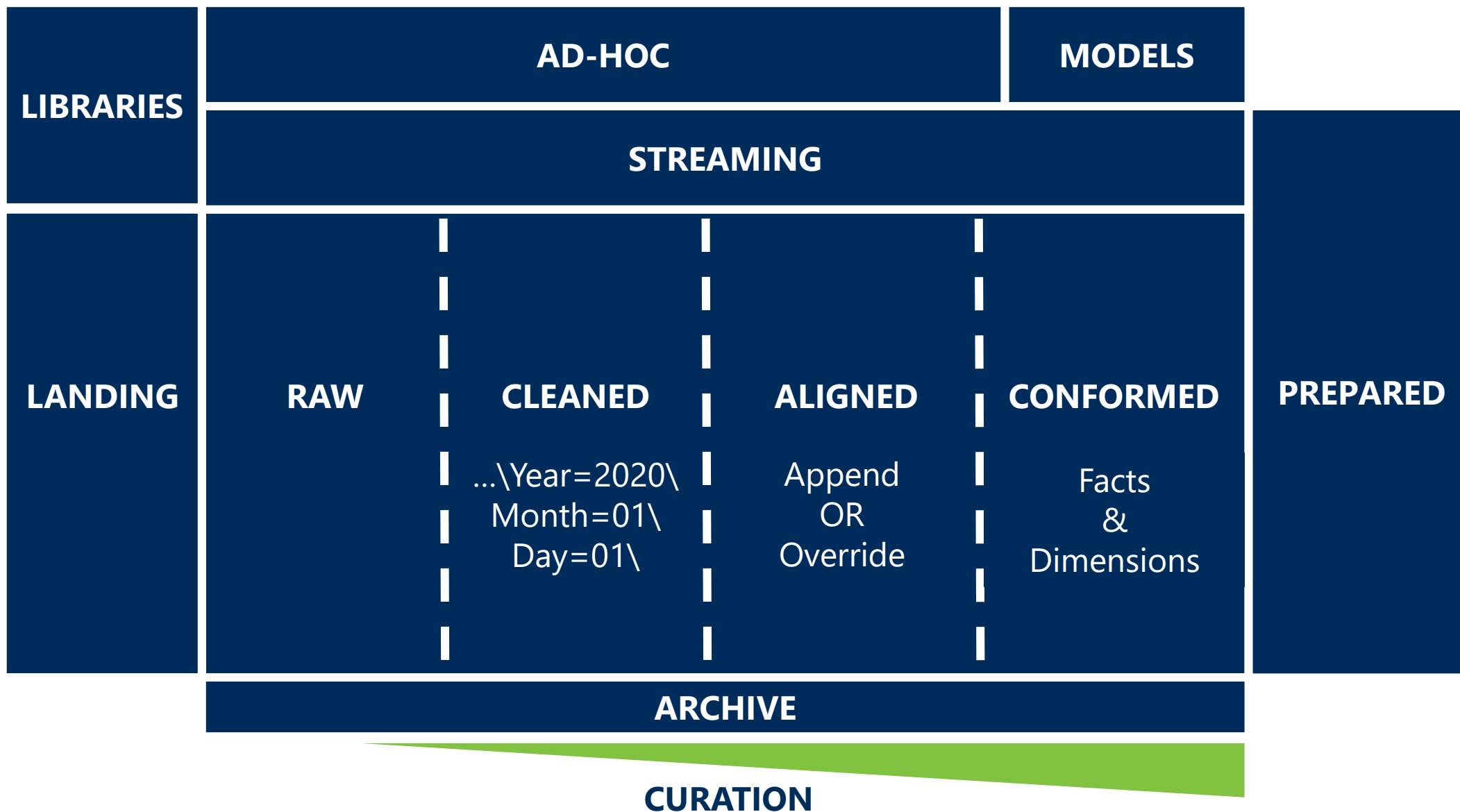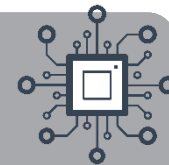
Hadoop Distributed File System ( HDFS )

# Data Transformation – Storage & Format

# Data Transformation – Storage & Format

| LIBRARIES | AD-HOC | | | | MODELS |
|---|---|---|---|---|---|
| | STREAMING | | | | |
| LANDING | RAW | CLEANED<br><br>...\Year=2020\<br>Month=01\<br>Day=01\ | ALIGNED<br><br>Append<br>OR<br>Override | CONFORMED<br><br>Facts<br>&<br>Dimensions | PREPARED |
| | ARCHIVE | | | | |

**CURATION**

# Data Transformation – Storage & Format
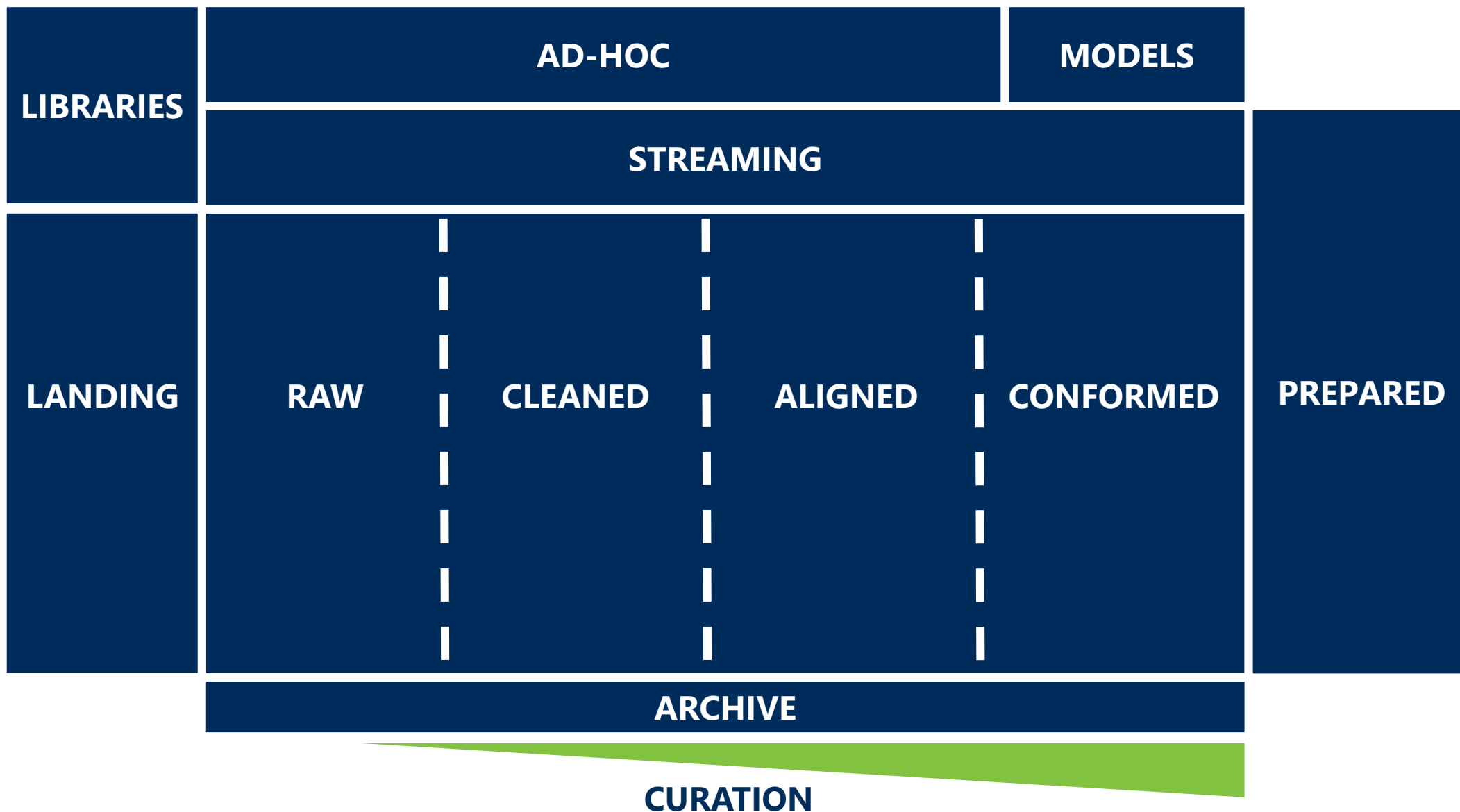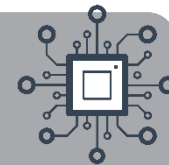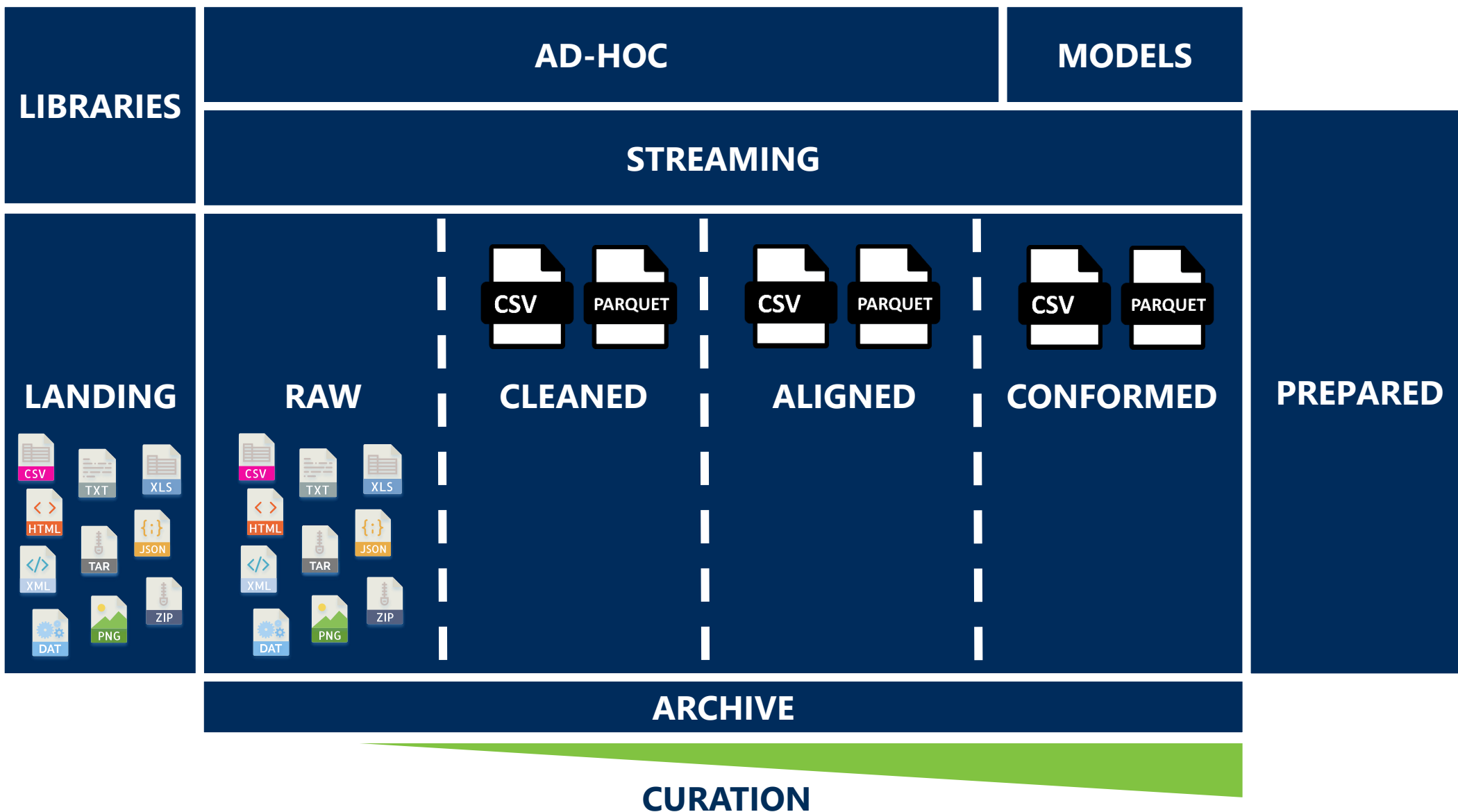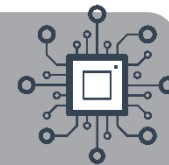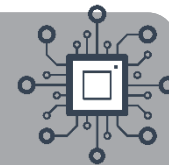
# Data Transformation – Storage & Format

**LIBRARIES**

**AD-HOC**

**MODELS**

**STREAMING**

**LANDING**

CSV
TXT
XLS
HTML
JSON
XML
TAR
DAT
PNG
ZIP

**RAW**

CSV
TXT
XLS
HTML
JSON
XML
TAR
DAT
PNG
ZIP

**CLEANED**

CSV | PARQUET

**ALIGNED**

CSV | PARQUET

**CONFORMED**

CSV | PARQUET

**PREPARED**

**ARCHIVE**

**CURATION**

# Data Transformation – Storage & Format

# Agenda

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

Compute ✓
Storage, Structure
& Data Format ✓

# Agenda

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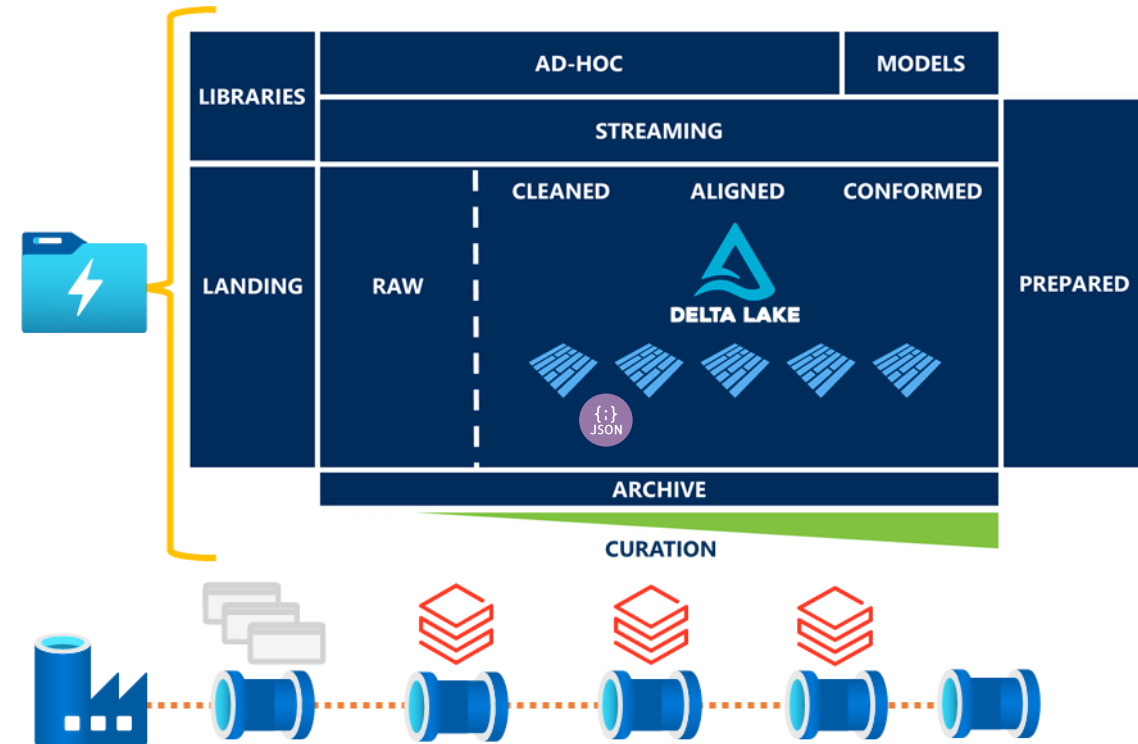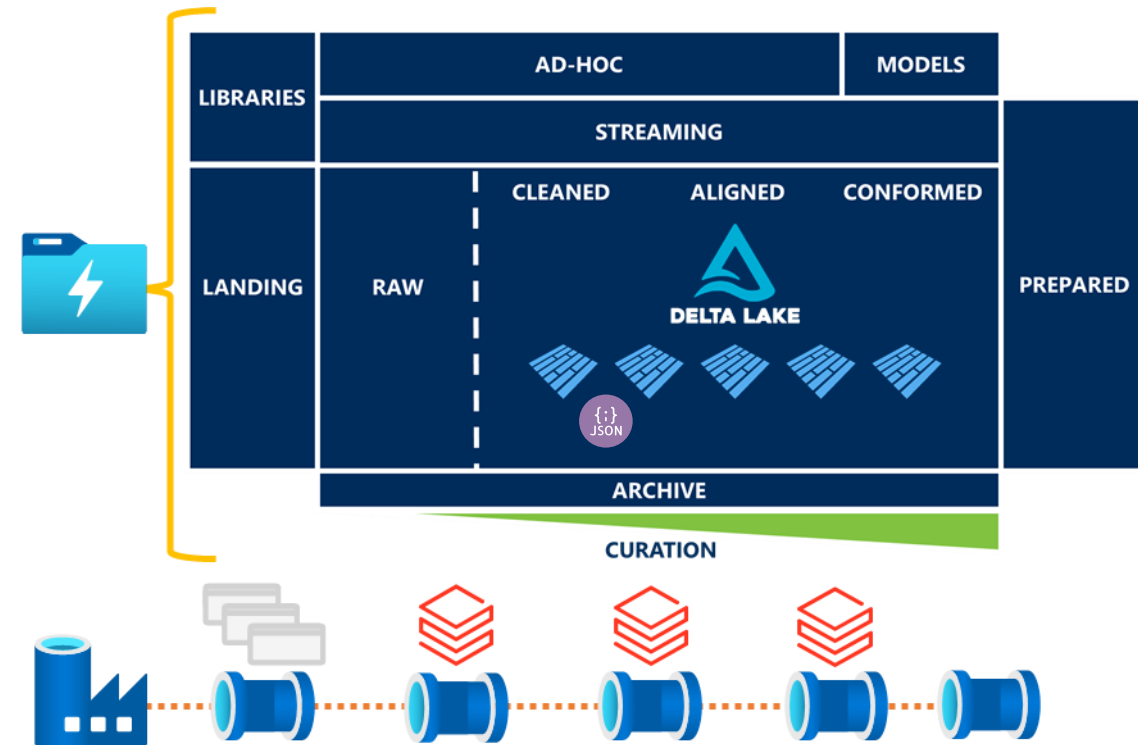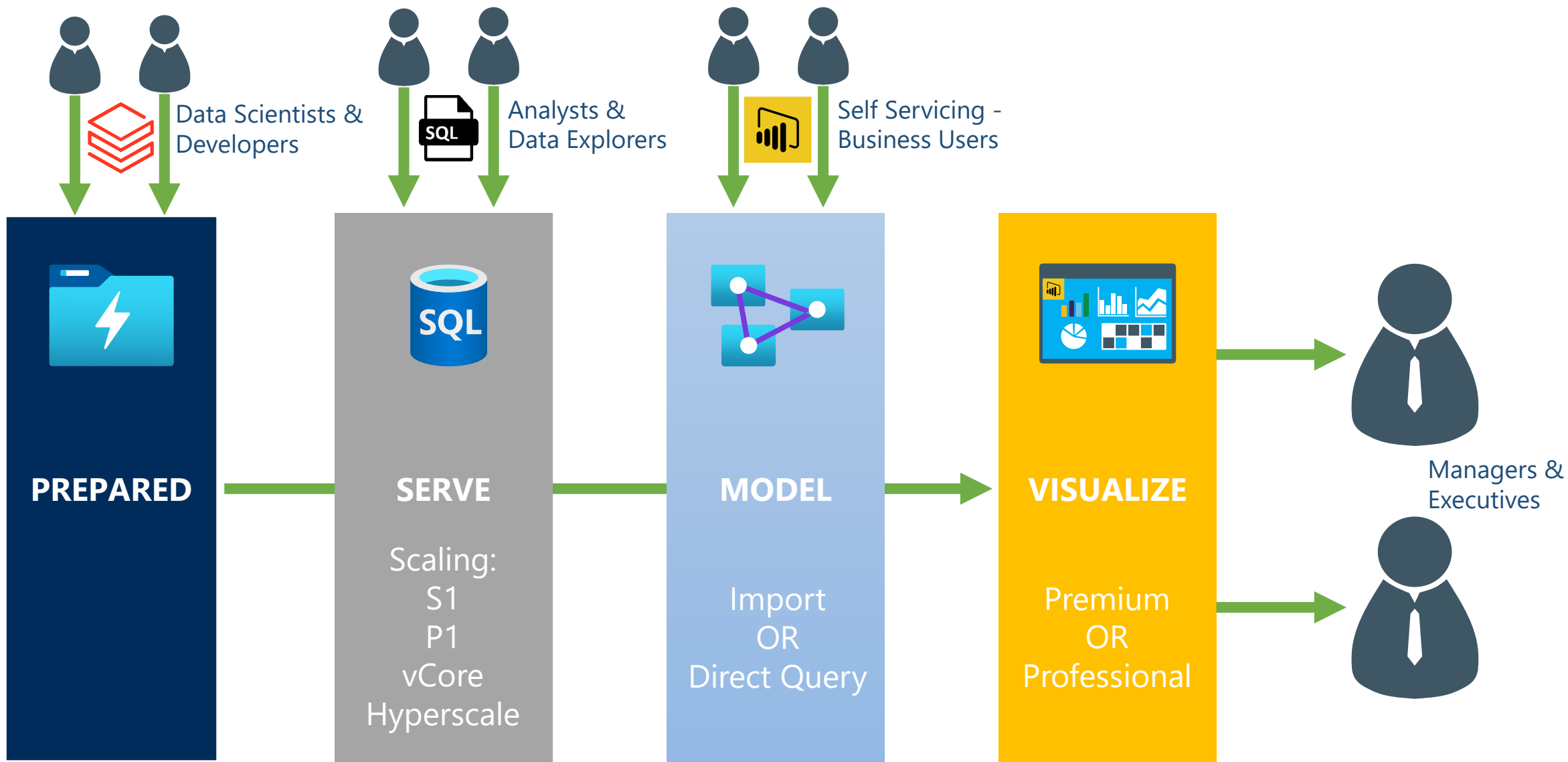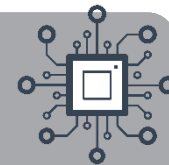3. Transform ✓
4. Load

# Agenda


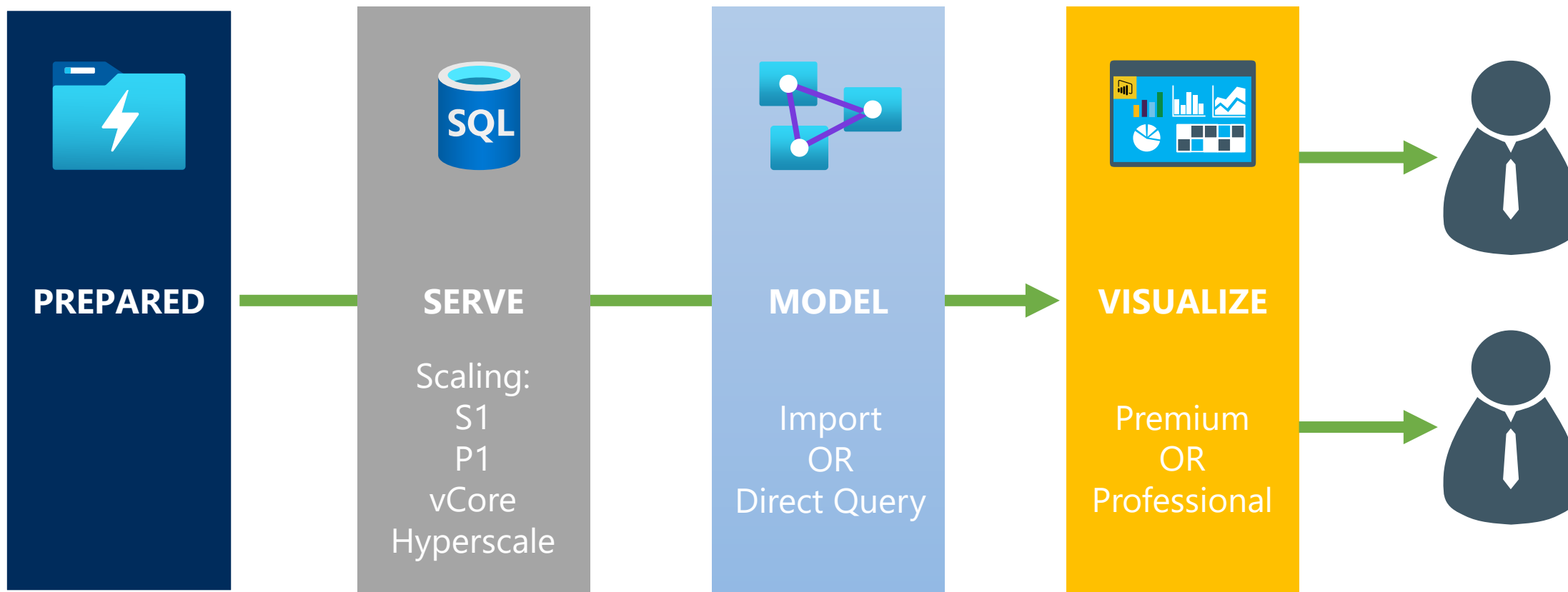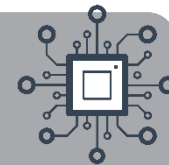
1. Design ✓
2. Extract ✓
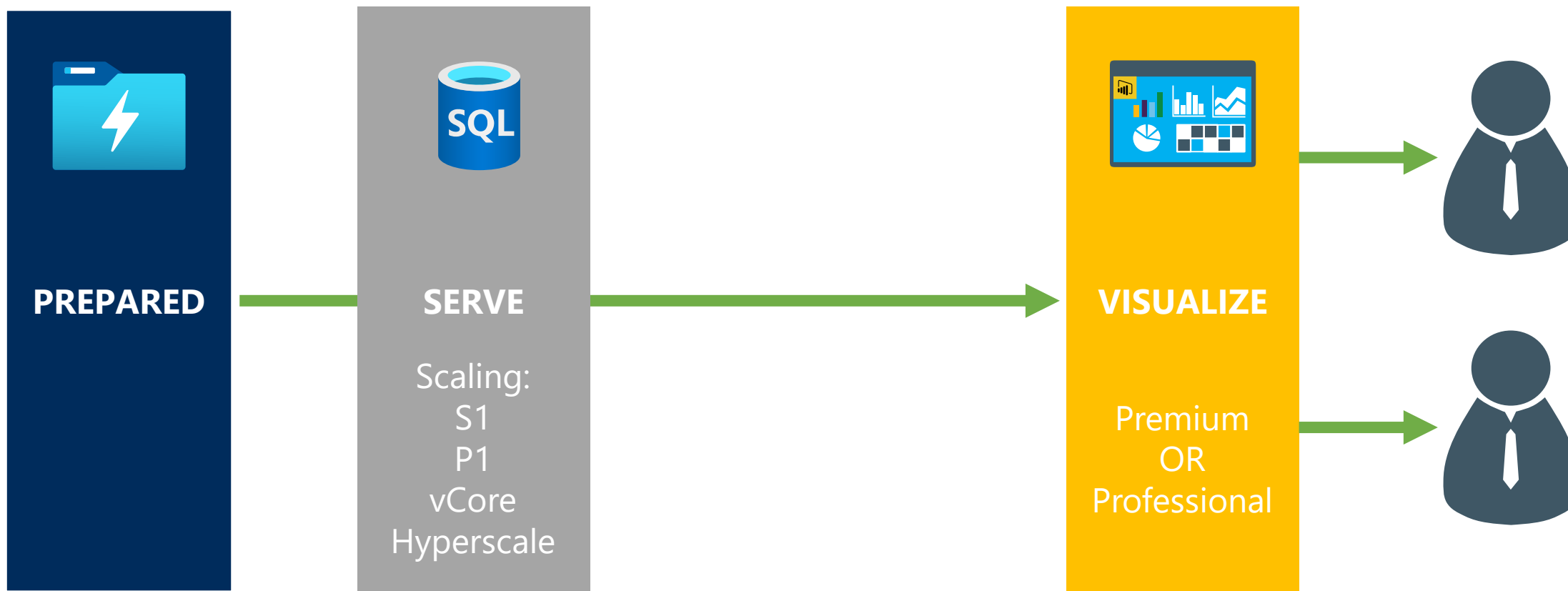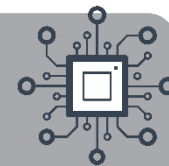3. Transform ✓
4. Load

# Loading & Consuming Data

Data Scientists & Developers

Analysts & Data Explorers

Self Servicing - Business Users

**PREPARED**

**SERVE**

Scaling:
S1
P1
vCore
Hyperscale

**MODEL**

Import
OR
Direct Query

**VISUALIZE**

Premium
OR
Professional

Managers & Executives

# Loading & Consuming Data

**PREPARED**

**SERVE**

Scaling:
S1
P1
vCore
Hyperscale

**MODEL**

Import
OR
Direct Query

**VISUALIZE**

Premium
OR
Professional

# Loading & Consuming Data

**PREPARED**

**SERVE**

Scaling:
S1
P1
vCore
Hyperscale

**VISUALIZE**

Premium
OR
Professional

# Loading & Consuming Data

PREPARED

MODEL

Import
OR
Direct Query

VISUALIZE

Premium
OR
Professional

# Loading & Consuming Data

**PREPARED**

**VISUALIZE**

Premium
OR
Professional

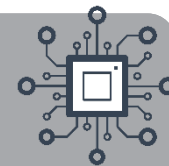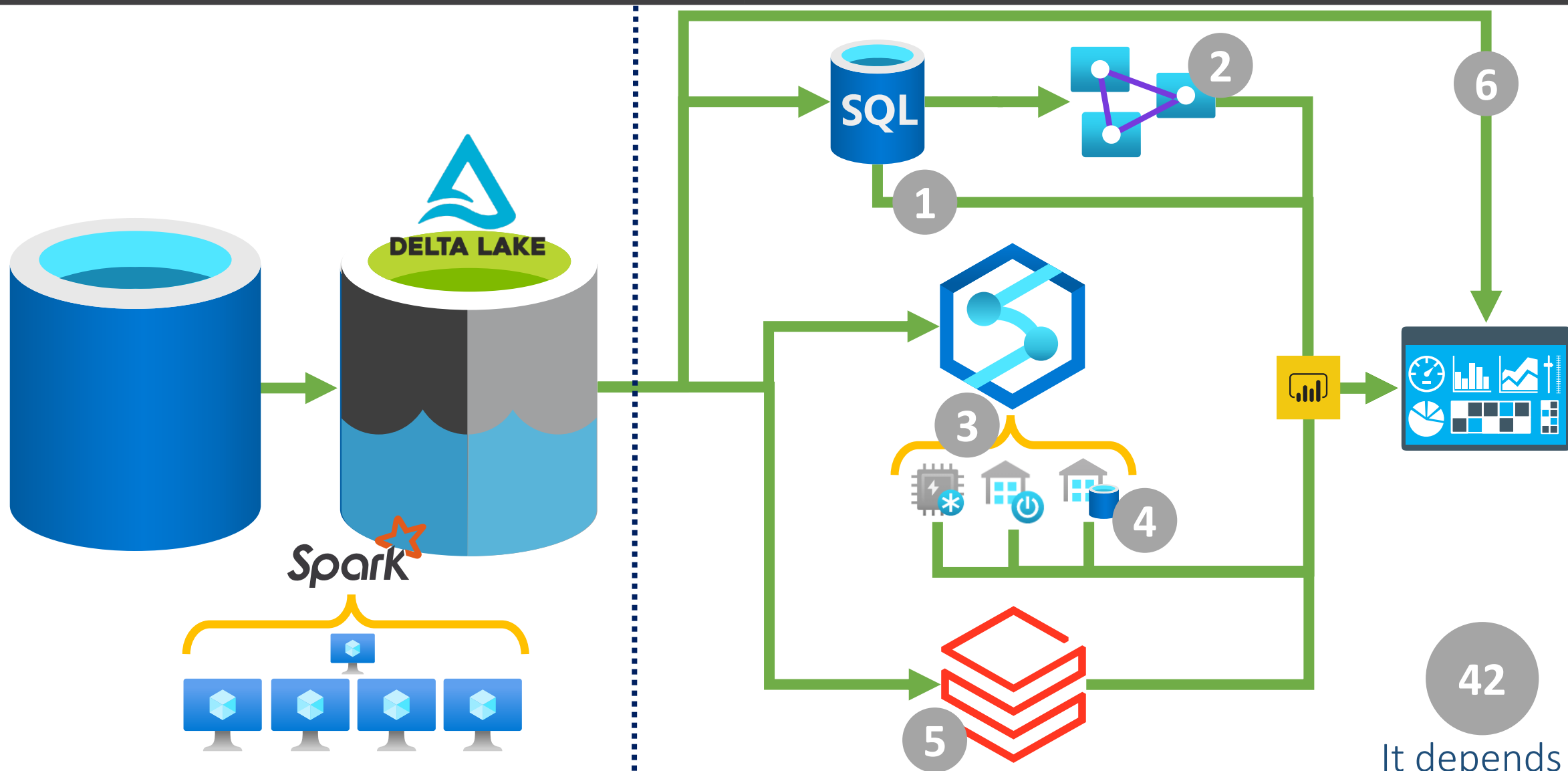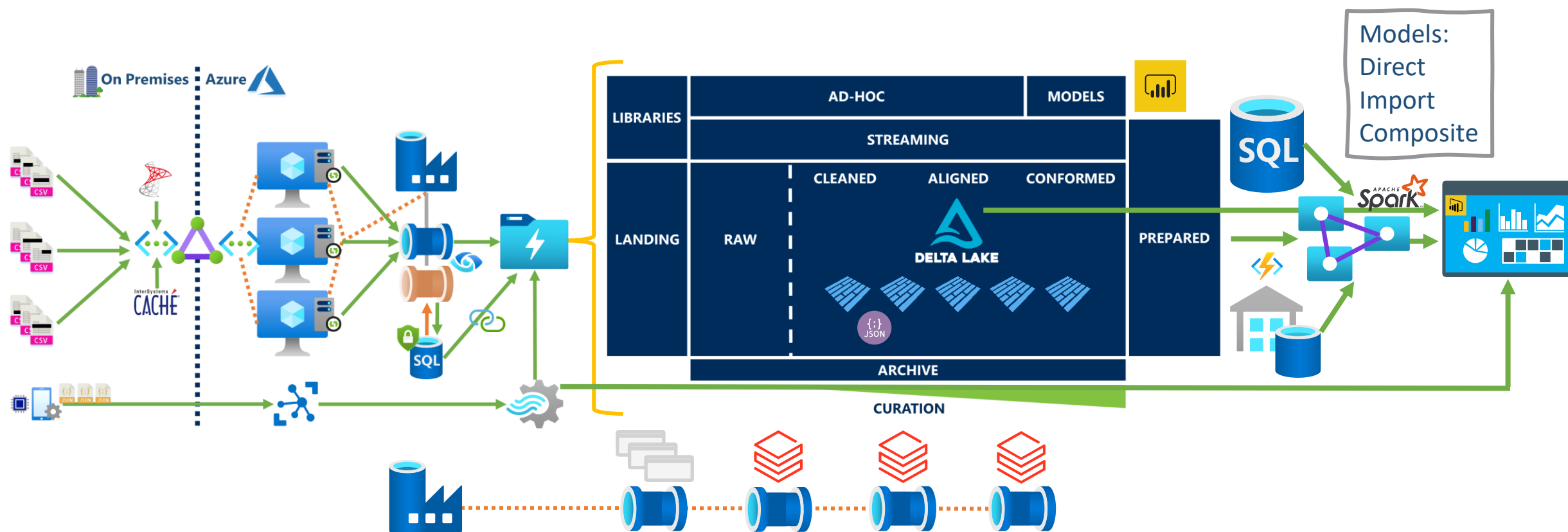# Consuming Our Lake House in Azure

Consuming Our Lake House in Azure

It depends!

# Module 1 - 6

An Architects Recap

```sql
SELECT
    [Summary]
FROM
    [Training]
WHERE
    [Module]
    BETWEEN 1 AND 6;
END; --module, fetch next
```