

Topics in Applied Data Science for Social Scientists

Week 5

The Data Engineering Perspective

February 20th, 2019
Nana Yaw Essuman

Agenda

1. Introduction
2. Overview on Data Engineering
3. Accessibility of Data
4. Data Quality/Monitoring
5. Transforming Data
6. Model Outputs
7. Data - Model Lifecycle
8. Collaboration w/ Data Engineers
9. How to ask the right questions
10. Q&A

Over 7 years of Data Engineering Experience

Over 7 years of Data Engineering Experience



Over 7 years of Data Engineering Experience



Over 7 years of Data Engineering Experience



COMCAST

Over 7 years of Data Engineering Experience



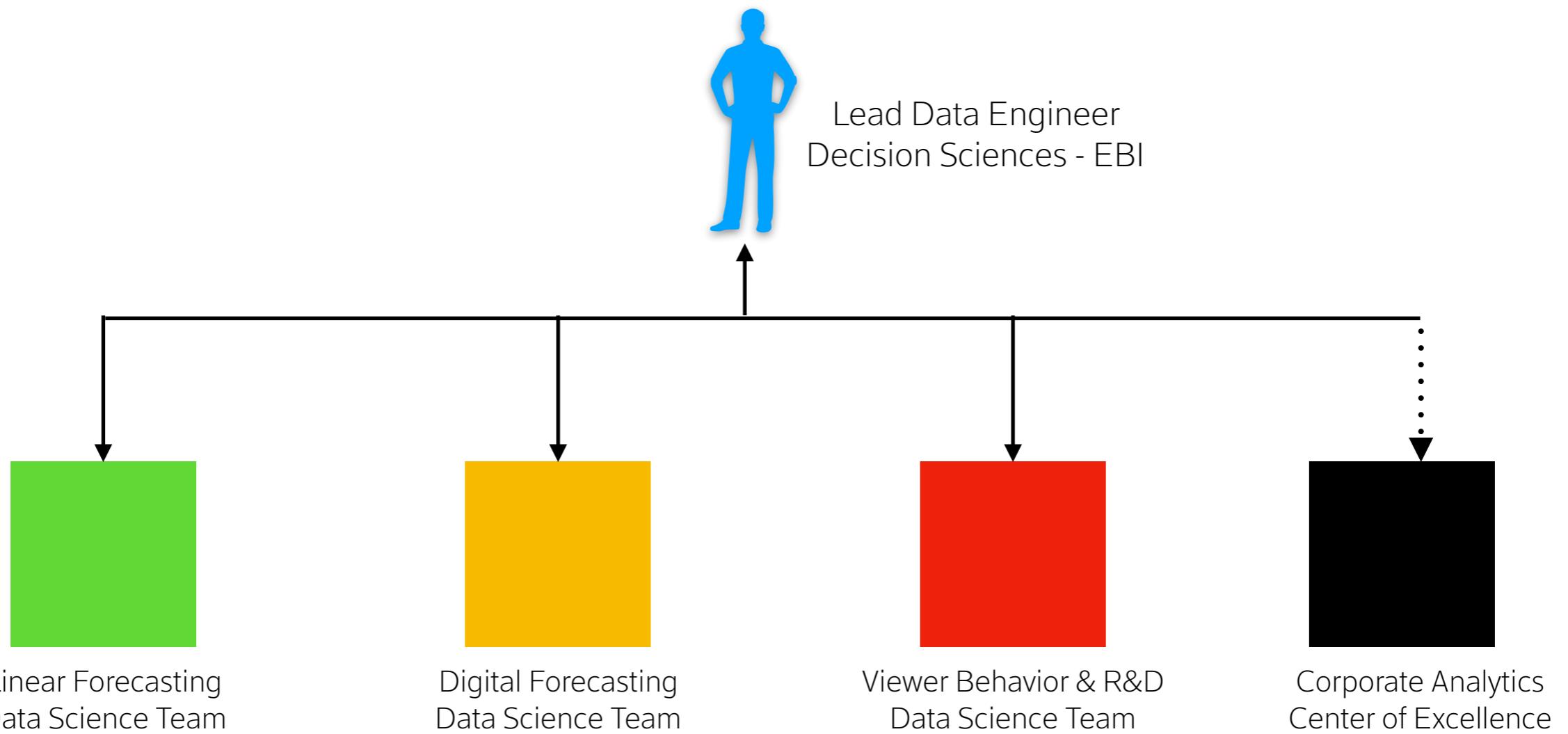
Over 7 years of Data Engineering Experience



Over 7 years of Data Engineering Experience



9am-5pm



What is Data Engineering?

What is Data Engineering?

Data engineering focuses on data, its curation, movement, storage, security, and processing. Data engineering is performed by data engineers. They work on building pipelines, applications, APIs, and systems that produce, process, and consume data to meet business needs. To compare; data engineering converts raw data into knowledge data. Data Scientists deals with using this data produced by data engineers to generate insights & to predict the future using data from the past.

Accessibility of Data

Accessibility of Data

Question: A Data Scientist has been tasked to predict the likelihood Show A will perform during a particular time slot on a network compared to the current time slot proposed by TV Network Executives. What should be the Data Scientists first worry?

Accessibility of Data

Question: A Data Scientist has been tasked to predict the likelihood Show A will perform during a particular time slot on a network compared to the current time slot proposed by TV Network Executives. What should be the Data Scientists first worry?

1. Where do I get my data from?
2. Is the data structured or unstructured?
3. How do I pull the data? One time or Continuous
4. Where do I store the data?
5. What is the right way to store the data?
6. How do I create various subsets of the data?
7. How do I access the data for my models?

The Data

Structured



TERADATA



{API}



Unstructured



The Data

Structured



{API}

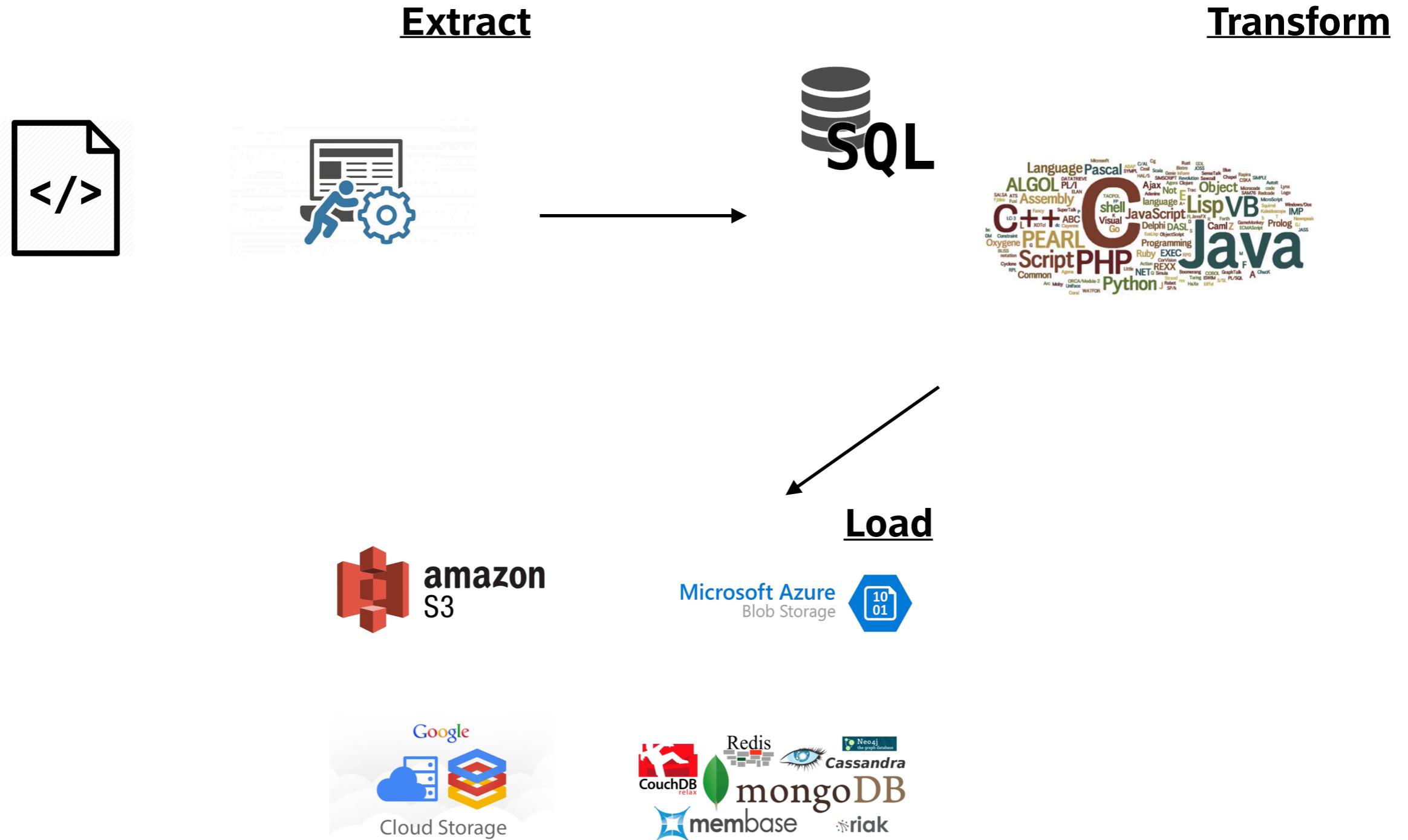


Unstructured

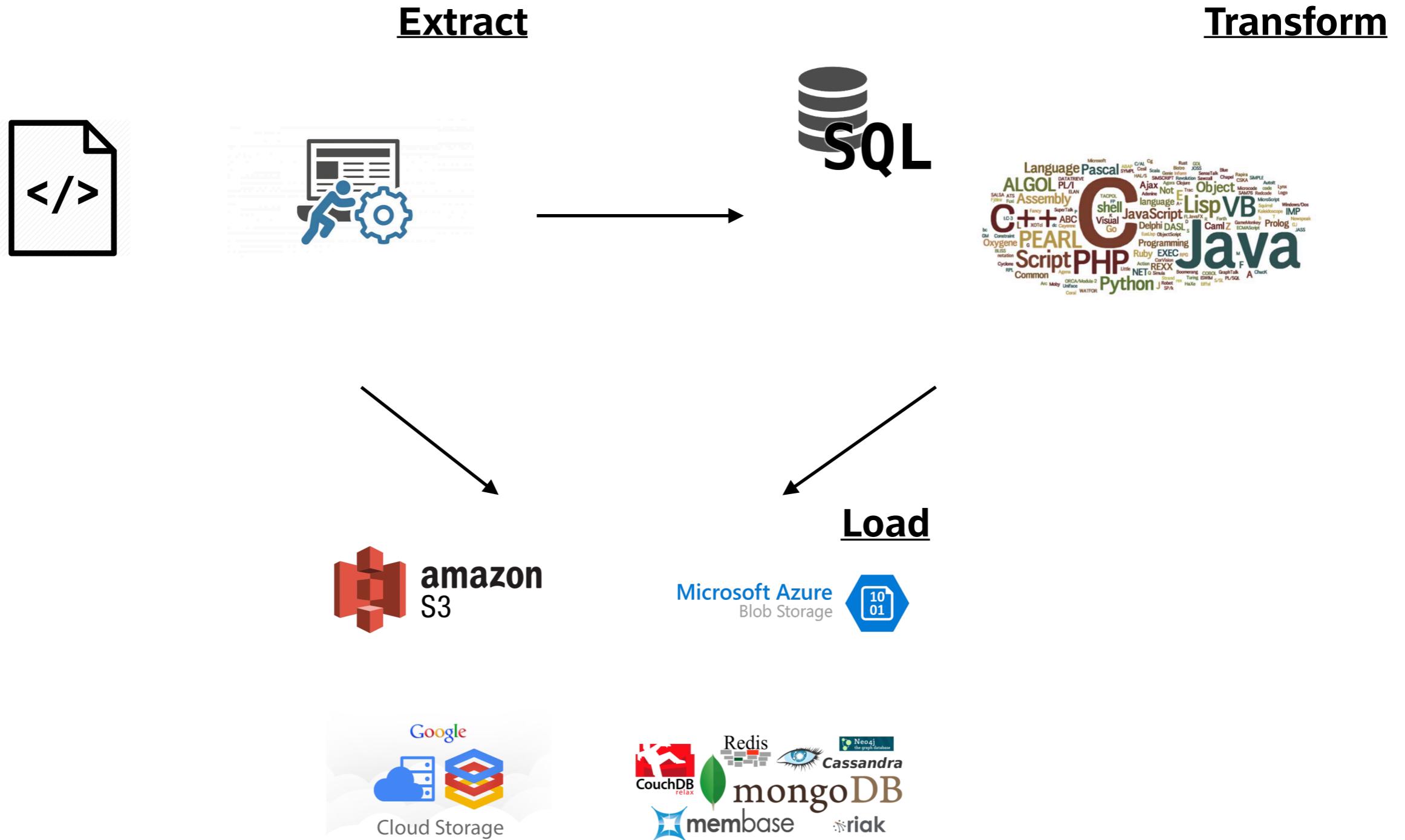


	Structured Data	Unstructured Data
Characteristics	<ul style="list-style-type: none">• Pre-defined data models• Usually text only• Easy to search	<ul style="list-style-type: none">• No pre-defined data model• May be text, images, sound, video or other formats• Difficult to search
Resides in	<ul style="list-style-type: none">• Relational databases• Data warehouses	<ul style="list-style-type: none">• Applications• NoSQL databases• Data warehouses• Data lakes
Generated by	Humans or machines	Humans or machines
Typical applications	<ul style="list-style-type: none">• Airline reservation systems• Inventory control• CRM systems• ERP systems	<ul style="list-style-type: none">• Word processing• Presentation software• Email clients• Tools for viewing or editing media
Examples	<ul style="list-style-type: none">• Dates• Phone numbers• Social security numbers• Credit card numbers• Customer names• Addresses• Product names and numbers• Transaction information	<ul style="list-style-type: none">• Text files• Reports• Email messages• Audio files• Video files• Images• Surveillance imagery

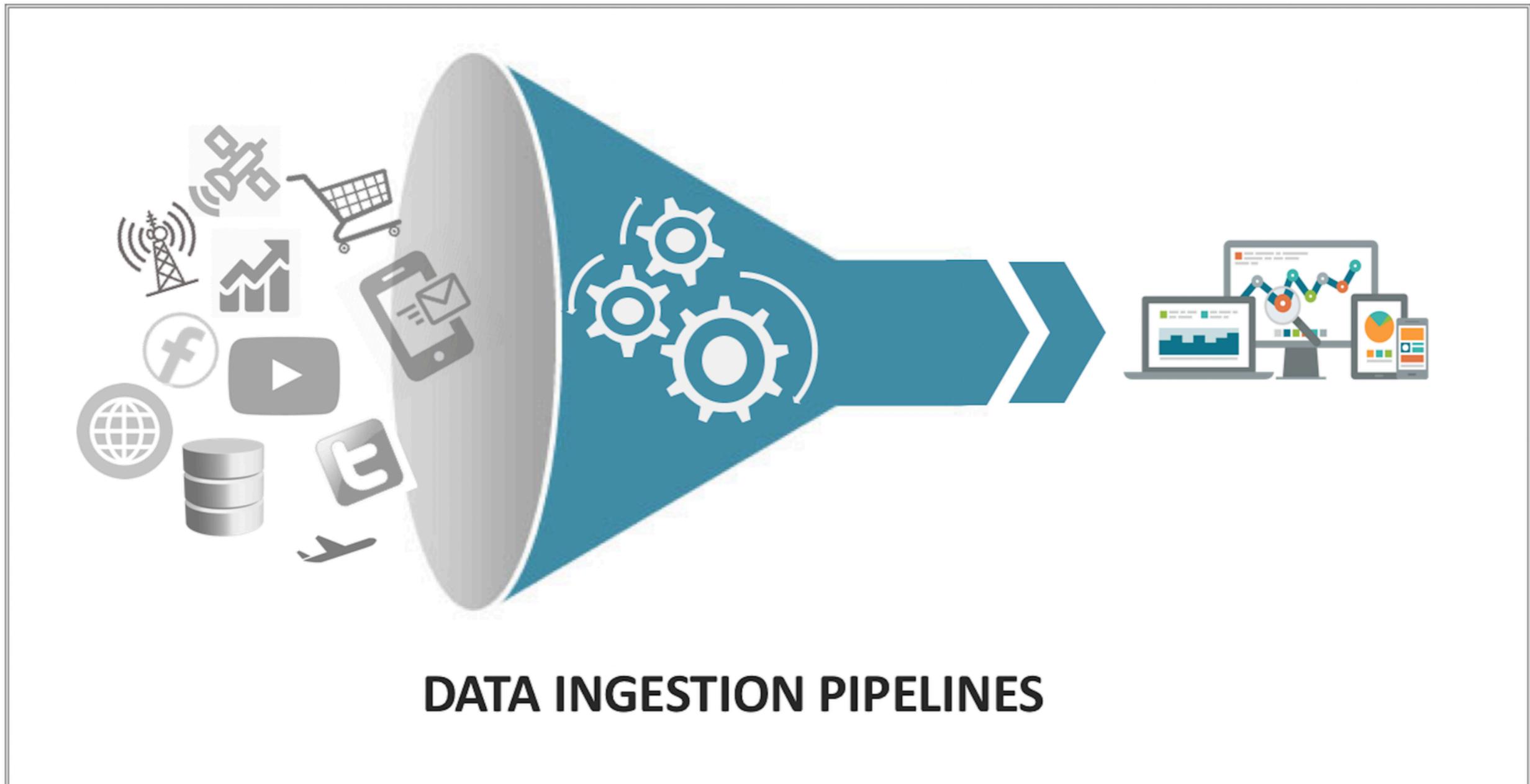
Extract-Transform-Load (ETL)



Extract-Transform-Load (ETL)



Automate & Scale



Data Quality & Monitoring

Effective data quality maintenance requires periodic data monitoring and cleaning. In general, data quality maintenance involves updating/standardizing data and deduplicating records to create a single data view.

Data Quality & Monitoring

Effective data quality maintenance requires periodic data monitoring and cleaning. In general, data quality maintenance involves updating/standardizing data and deduplicating records to create a single data view.

Key data quality components are as follows:

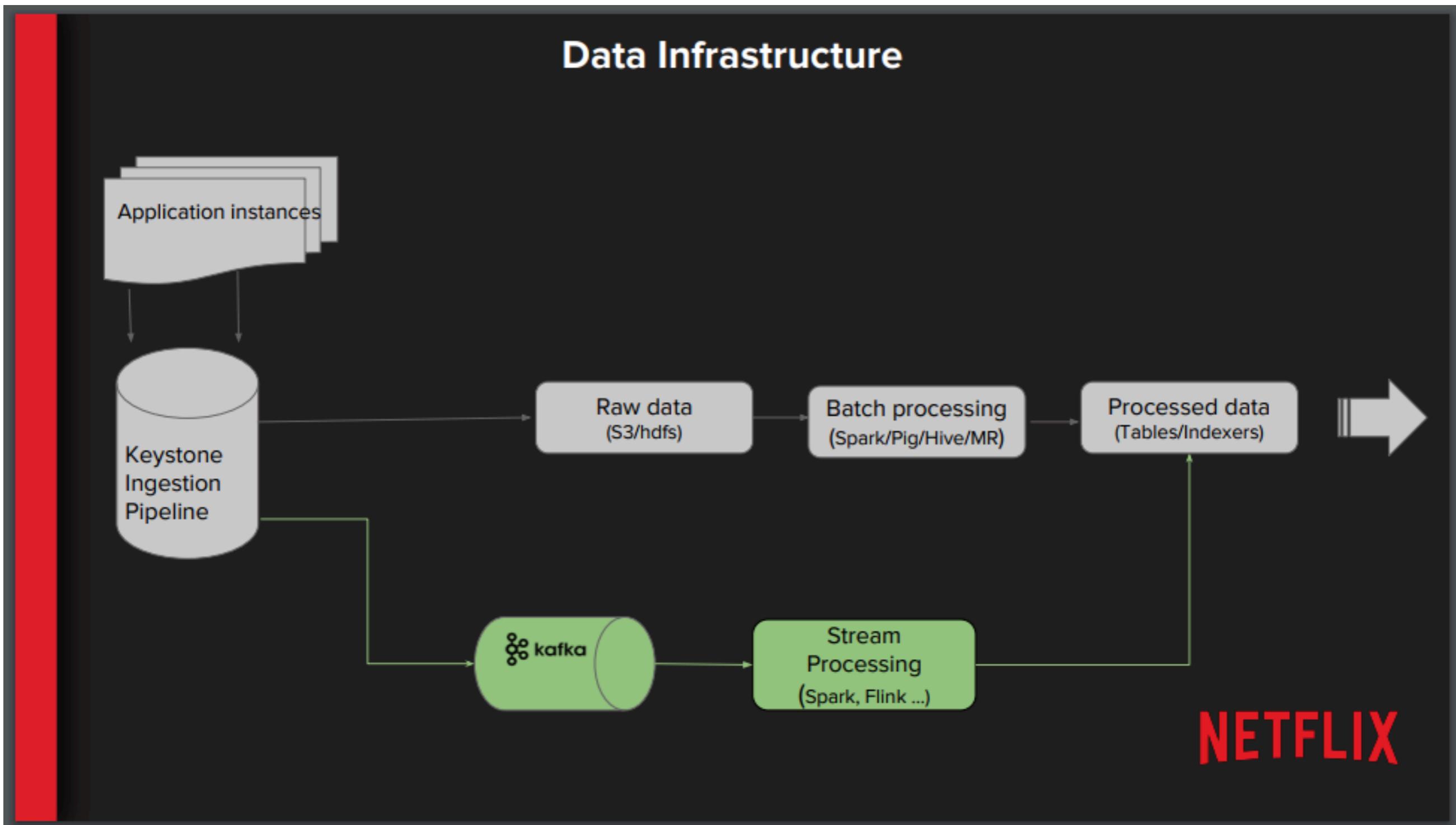
Data Quality & Monitoring

Effective data quality maintenance requires periodic data monitoring and cleaning. In general, data quality maintenance involves updating/standardizing data and deduplicating records to create a single data view.

Key data quality components are as follows:

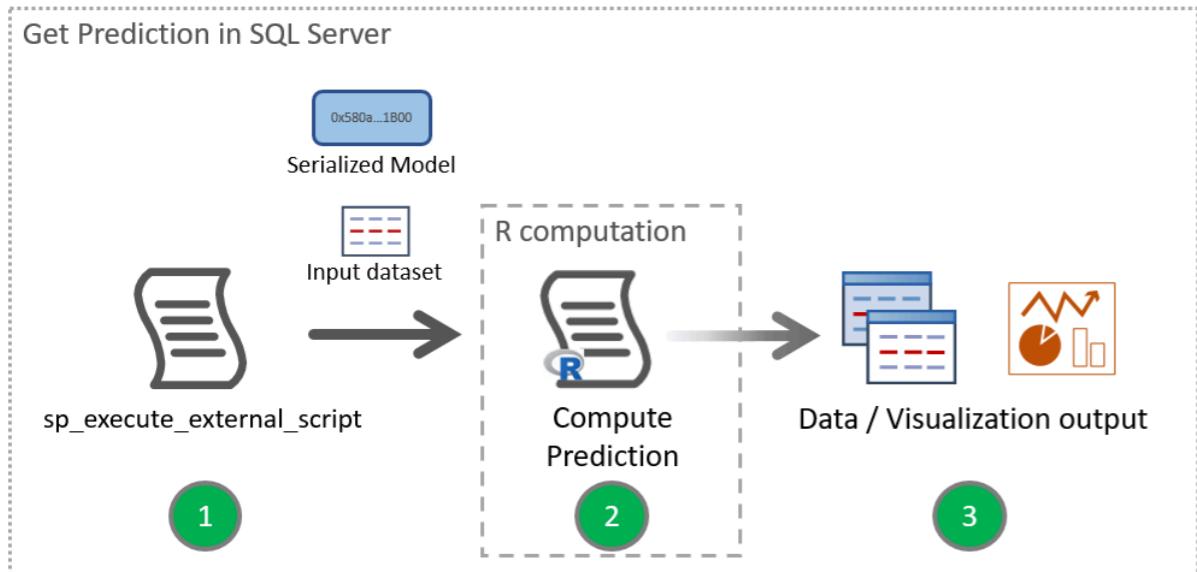
- Completeness: Level at which desired data attributes are supplied. Data does not need to be 100 percent complete
- Accuracy: Represents data's real world status. May be calculated by using an automated method with the help of various lists and mapping
- Credibility: Extent to which data is considered credible and true. May differ from by source
- Timeliness (age of data): Extent to which data is adequately updated for a current venture
- Consistency: Assesses whether various dataset facts match
- Integrity: Assesses reference validity and accurate joining of various datasets

Transforming your Data

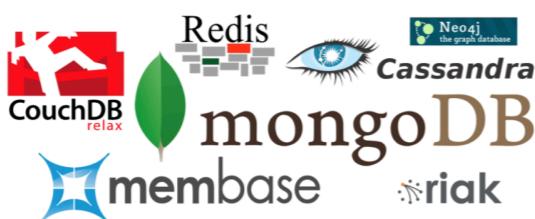
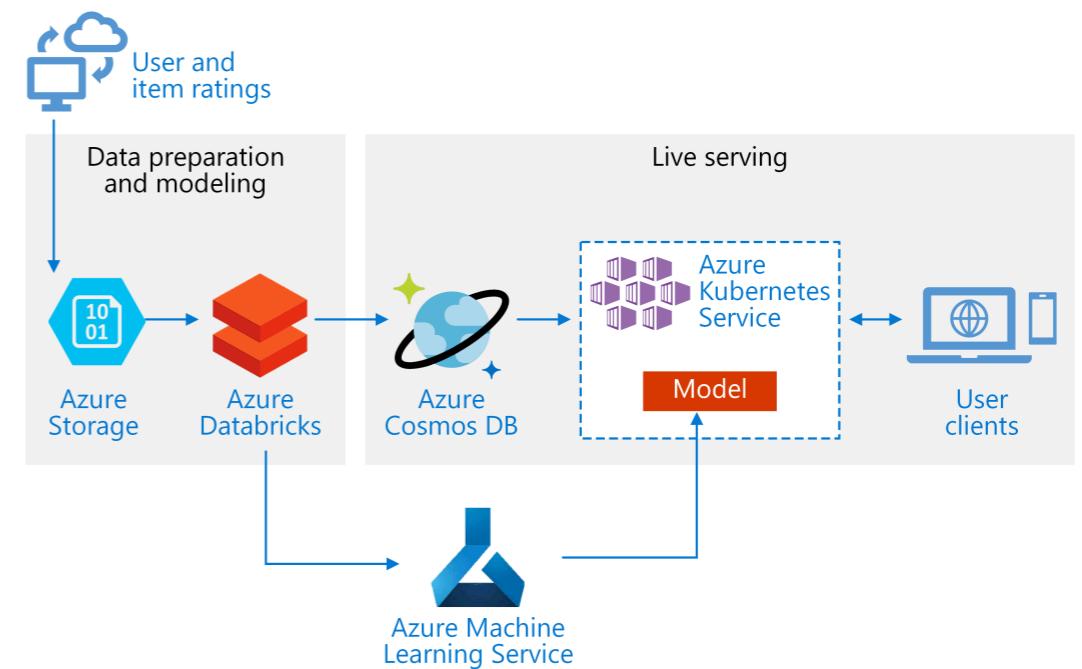


Operationalize Model Outputs

Batch

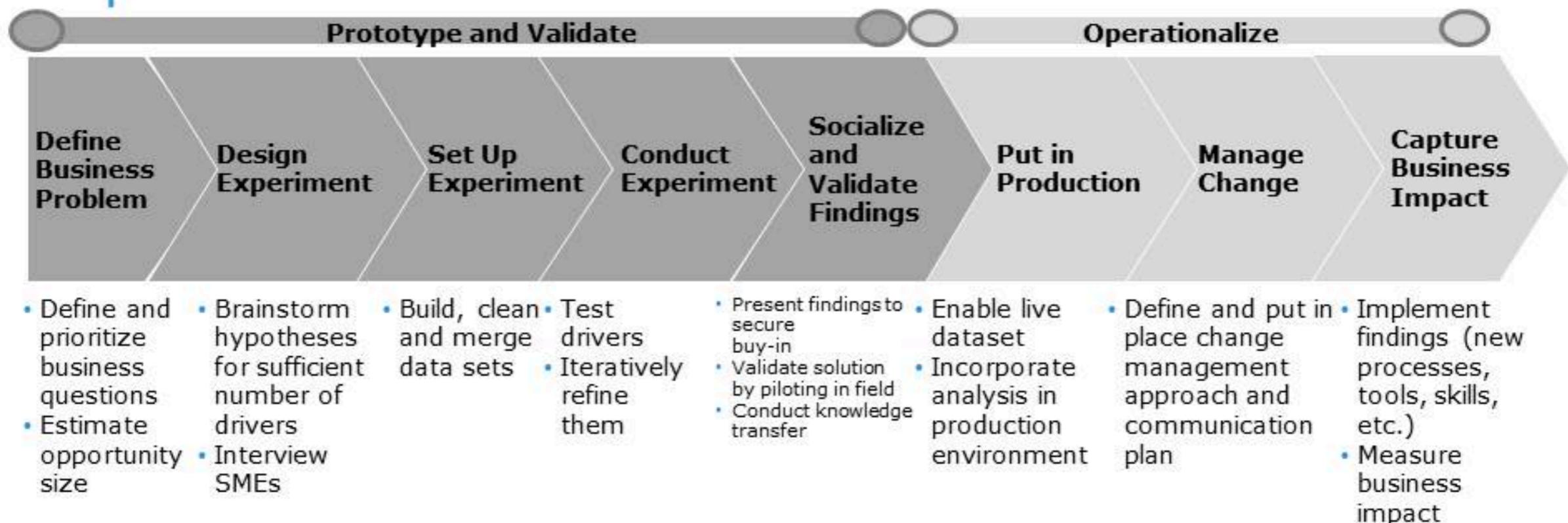


Real-Time

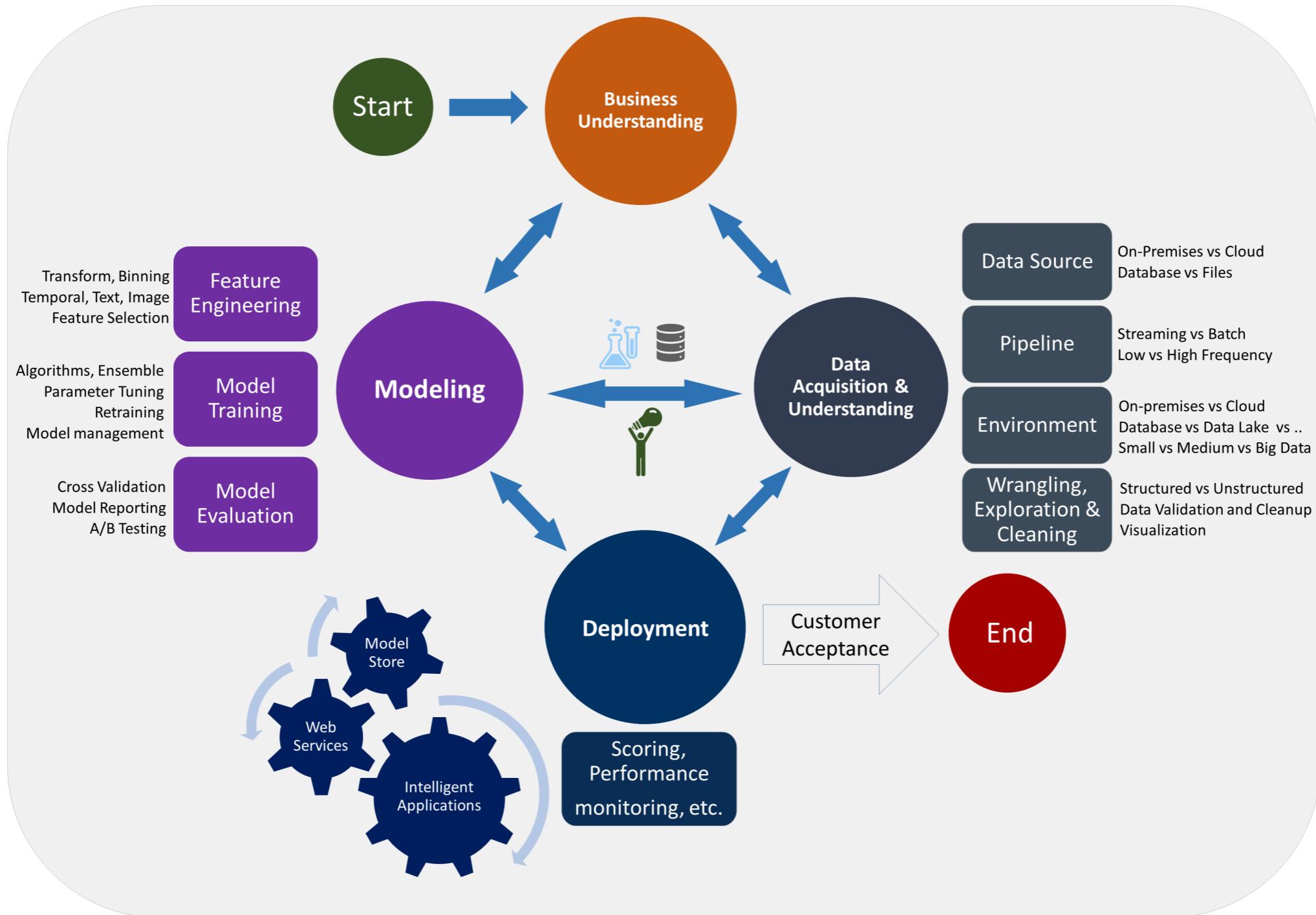


Data - Model Lifecycle

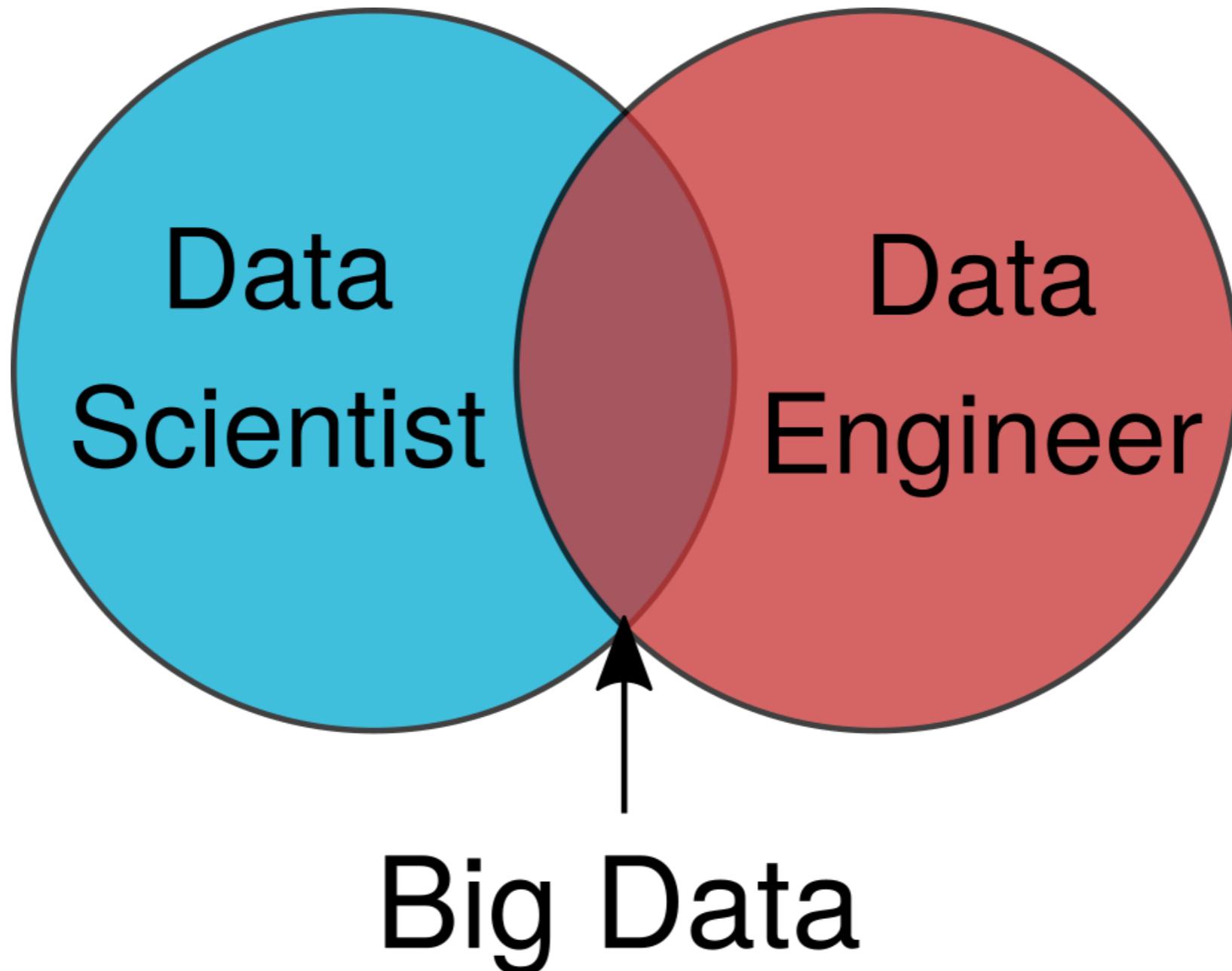
Establish a Process: Insight to Operational



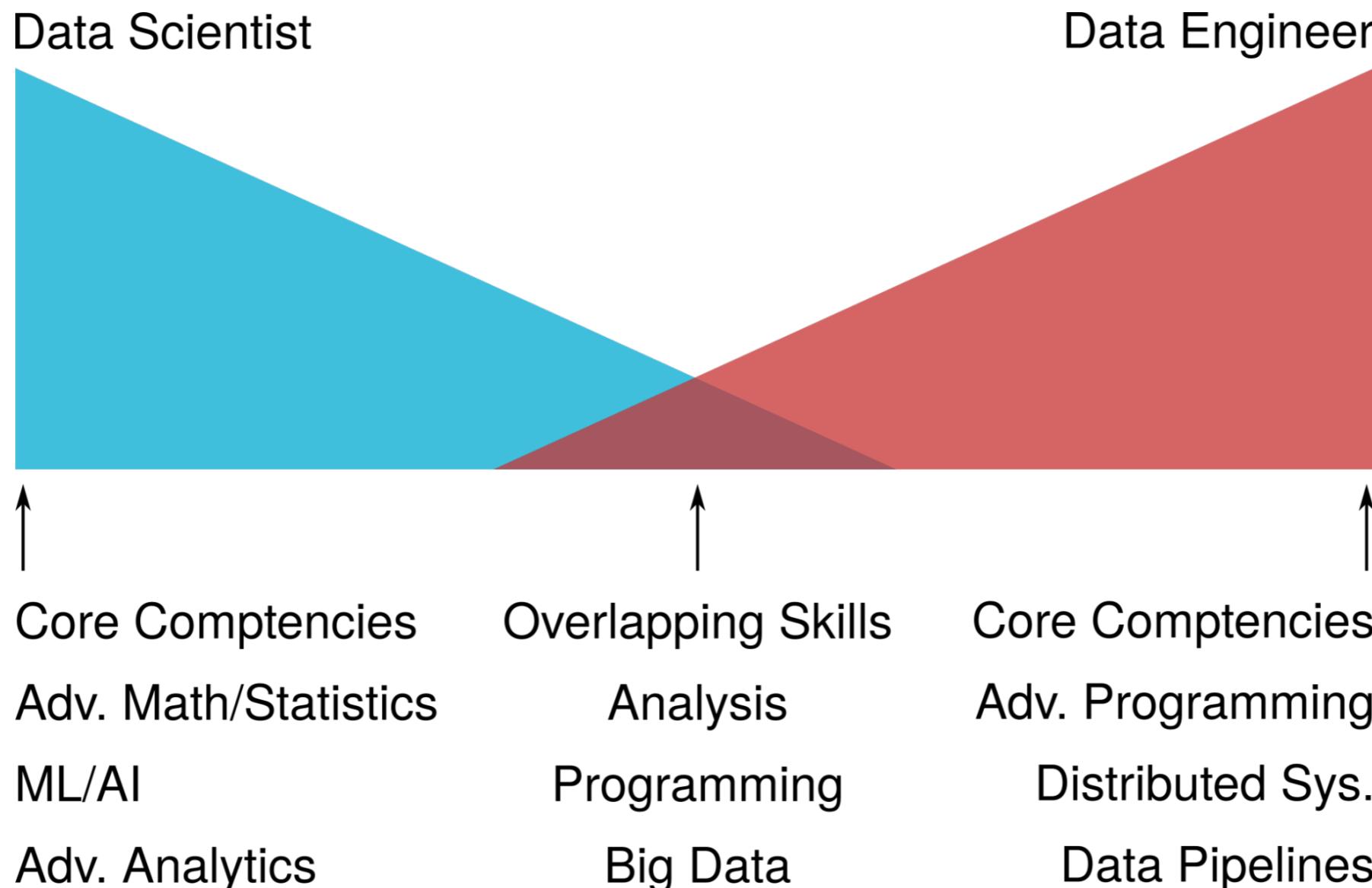
Data - Model Lifecycle



Collaboration w/ Data Engineers



Collaboration w/ Data Engineers



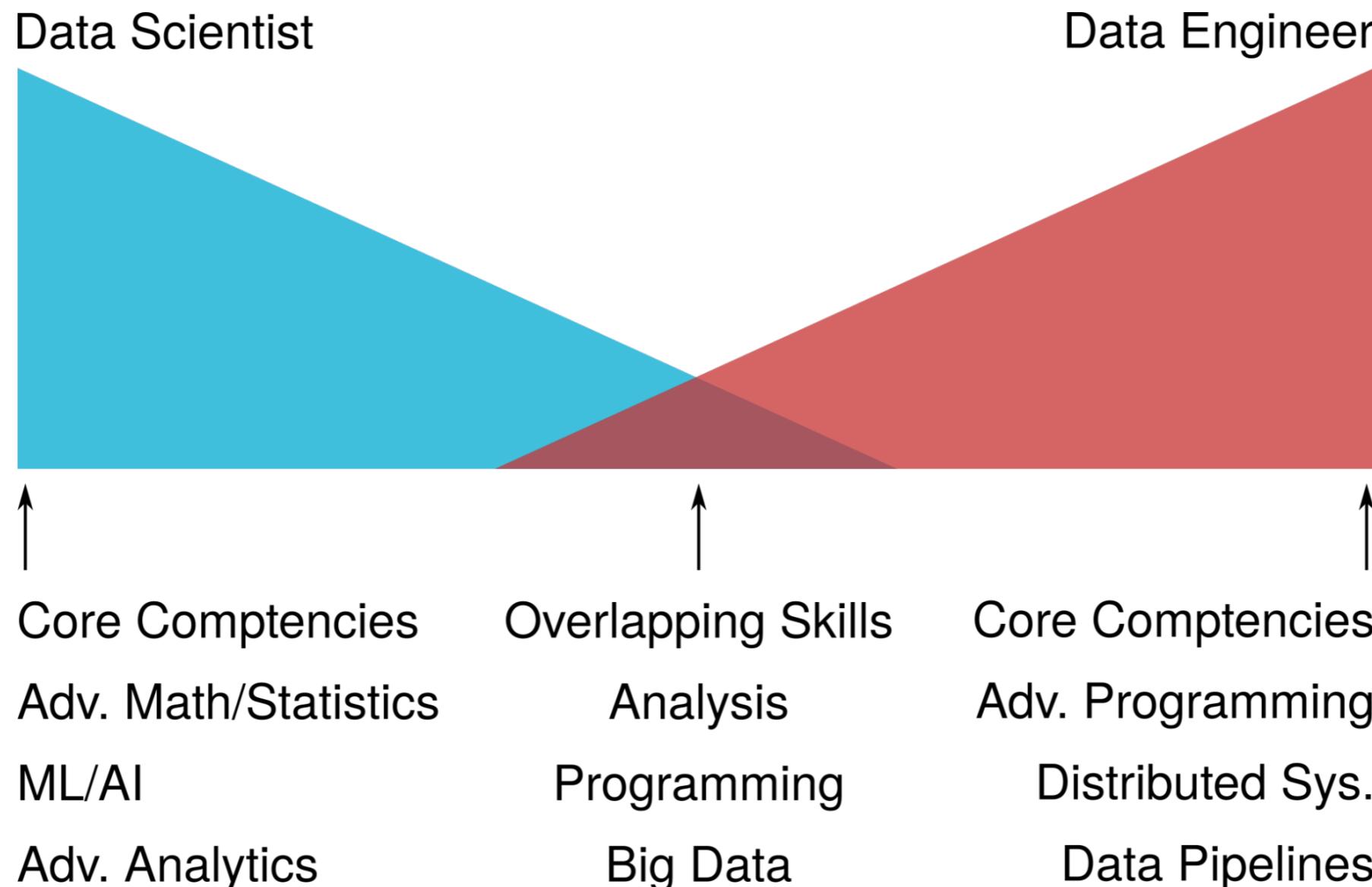
BDI

BIG DATA INSTITUTE

For more information go to <http://bigdatainstitute.io>

Attribution-NoDerivatives 4.0 International (CC BY-ND 4.0) Licensed

Collaboration w/ Data Engineers



Varies based on size of organization

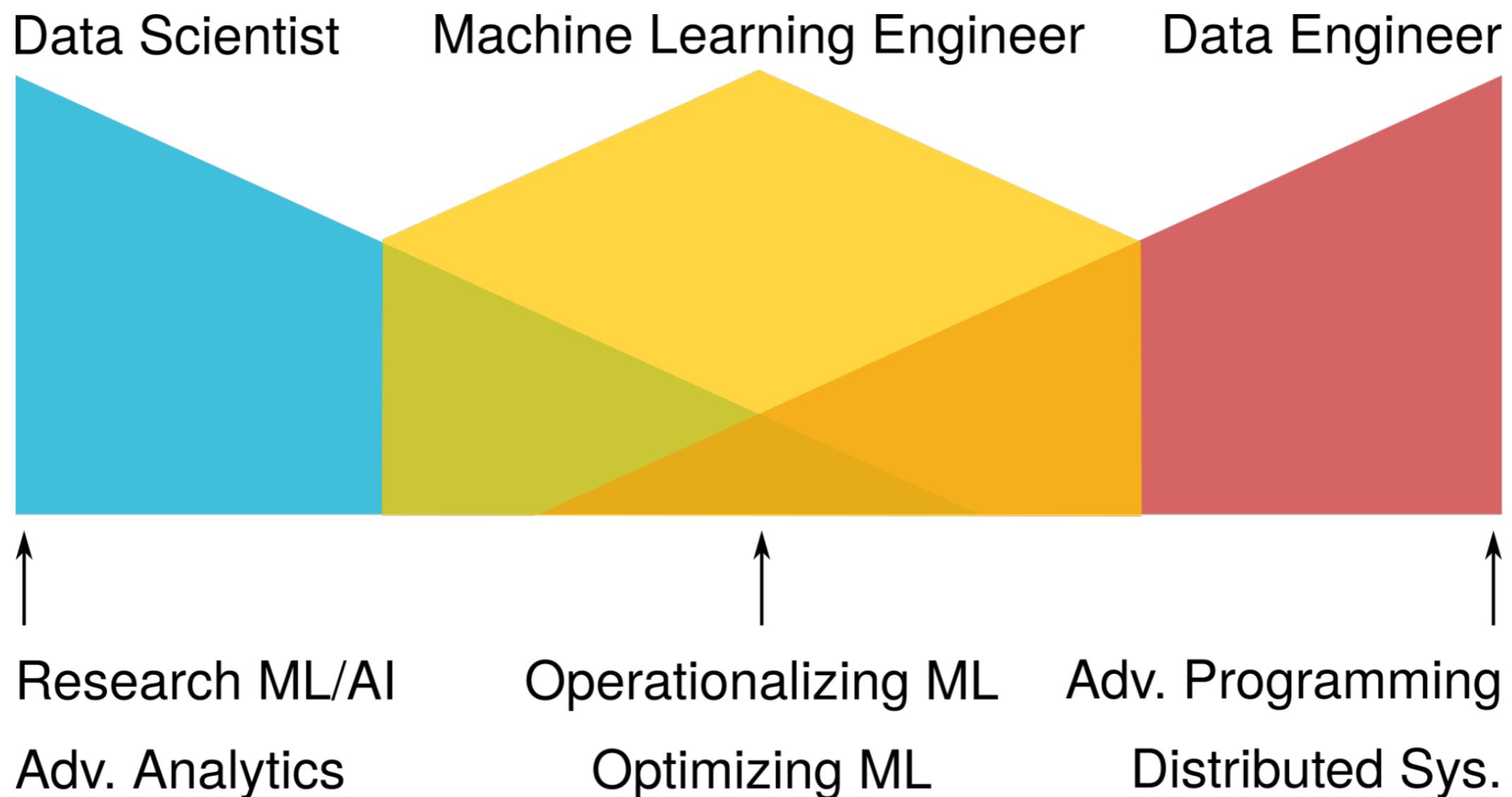
BDI

BIG DATA INSTITUTE

For more information go to <http://bigdatainstitute.io>

Attribution-NoDerivatives 4.0 International (CC BY-ND 4.0) Licensed

Collaboration w/ Data Engineers



BDI

BIG DATA INSTITUTE

For more information go to <http://bigdatainstitute.io>

Attribution-NoDerivatives 4.0 International (CC BY-ND 4.0) Licensed

What to ask your Data Engineer or yourself

What to ask your Data Engineer or yourself

1. Where do I get my data from?
2. Is the data structured or unstructured?
3. How do I pull the data? One time or Continuous
4. Where do I store the data?
5. What is the right way to store the data?
6. How do I create various subsets of the data?
7. How do I access the data for my models?
8. How do I deploy this model to production? Real-Time or Batch
9. How do I monitor my model metrics?

Resources

<https://www.oreilly.com/ideas/data-engineers-vs-data-scientists>

<https://www.datamation.com/big-data/structured-vs-unstructured-data.html>

<https://www.infoworld.com/article/3250852/data-science/handing-off-models-from-data-science-to-it.html>

<https://www.datascience.com/blog/predictive-data-models-from-data-cleaning-to-model-deployment>

Contact



nanayaw.essuman@nbcuni.com
nyessuman@gmail.com