



PLURALSIGHT

Operations

Week 6



Proprietary and confidential

 PLURALSIGHT

Snowflake Topics

HIGH LEVEL TOPICS FOR THIS WEEK

- Automation
- Orchestration
- DevOps
- DataOps
- MLOps
- ModelOps
- LLMOps
- Data Governance
- Data Quality and Data Privacy/Security

DevOps

Definition:

- DevOps is a set of practices that combines software development (Dev) and IT operations (Ops) to shorten the systems development life cycle and provide continuous delivery with high software quality.

Key Concepts:

- Continuous Integration (CI)
- Continuous Delivery (CD)
- Infrastructure as Code (IaC)
- Automation
- Monitoring and Logging

Tools: Jenkins, Docker, Kubernetes, Ansible, Terraform, Prometheus, Grafana

DataOps

Definition:

- DataOps is an automated, process-oriented methodology used by analytic and data teams to improve the quality and reduce the cycle time of data analytics.

Key Concepts:

- Agile Methodologies
- Continuous Integration/Continuous Deployment (CI/CD) for Data
- Data Quality and Governance
- Collaboration and Communication

Tools: Apache NiFi, dbt, Airflow, DataKitchen, Prefect

MLOps

Definition: MLOps is a set of practices that aims to deploy and maintain machine learning models in production reliably and efficiently.

Key Concepts:

- Model Training and Serving
- Model Versioning
- Continuous Integration/Continuous Deployment (CI/CD) for ML
- Monitoring and Governance

Tools: MLflow, Kubeflow, TFX, Seldon, SageMaker, Databricks

MLOps

Definition: MLOps is a set of practices that aims to deploy and maintain machine learning models in production reliably and efficiently.

Key Concepts:

- Model Training and Serving
- Model Versioning
- Continuous Integration/Continuous Deployment (CI/CD) for ML
- Monitoring and Governance

Tools: MLflow, Kubeflow, TFX, Seldon, SageMaker, Databricks

LLMOps

Definition: LLMOps (Large Language Model Operations) refers to the practices and tools necessary to deploy, maintain, and scale large language models in production environments.

Key Concepts:

- Model Training and Fine-tuning
- Scalability and Efficiency
- Monitoring and Feedback Loops
- Ethical Considerations and Bias Mitigation

Tools: Hugging Face Transformers, OpenAI API, GPT-3, Anthropic's Claude, Microsoft Azure OpenAI

Their Evolution

- **DevOps:** Emerged from the need to improve the speed and efficiency of software development and delivery. It evolved from traditional IT operations and Agile software development practices.
- **DataOps:** Developed as a response to the growing complexity and scale of data operations. It evolved from Agile and DevOps practices but is specifically tailored to the data analytics lifecycle.
- **MLOps:** Grew out of the necessity to operationalize machine learning models. It evolved from DevOps practices but incorporates the unique requirements of machine learning workflows, such as data management and model versioning.
- **LLM Ops:** A recent development due to the rise of large language models. It builds on MLOps practices but focuses on the specific challenges associated with deploying and managing large-scale, complex language models.

Why we need them

- **DevOps:** To streamline and automate the software development lifecycle, improve collaboration between development and operations teams, and deliver software faster and more reliably.
- **DataOps:** To address the challenges of managing data workflows, ensure data quality, and reduce the time to insight by automating and streamlining data operations.
- **MLOps:** To handle the complexities of deploying machine learning models into production, ensure reproducibility, and maintain model performance over time.
- **LLM Ops:** To manage the deployment, scaling, and monitoring of large language models, ensuring they deliver accurate and unbiased results while maintaining efficiency and scalability.

Their Differences

- **DevOps:** Focuses on software development and IT operations, ensuring efficient and reliable software delivery.
- **DataOps:** Centers on the data lifecycle, aiming to improve data quality and analytics processes.
- **MLOps:** Targets the deployment and maintenance of machine learning models, addressing the unique needs of ML workflows.
- **LLM Ops:** Specifically designed for the challenges of managing large language models, encompassing aspects of both MLOps and additional considerations for model scale and ethical use.

DataOps

Proprietary and confidential



DataOps

DataOps is a methodology that focuses on improving the communication, integration, and automation of data flows between data managers and data consumers across an organization. It aims to deliver data quickly, efficiently, and reliably.

Problems in Data and Analytics Industry

- **Data Silos:** Isolated data systems prevent easy access and integration.
- **Data Quality Issues:** Inconsistent, inaccurate, or incomplete data.
- **Slow Data Processing:** Long cycles for data preparation and processing.
- **Lack of Collaboration:** Poor communication between teams.
- **Manual Processes:** Labor-intensive and error-prone data handling.

DataOps

Root Cause: Organizational Complexities

- **Distributed Teams:** Teams working in isolation with different tools and methodologies.
- **Lack of Standardization:** Inconsistent processes and standards.
- **Complex Data Environments:** Diverse data sources and technologies.

What is DataOps?

DataOps is an automated, process-oriented methodology aimed at improving the quality and reducing the cycle time of data analytics. It involves the orchestration of people, processes, and technology to deliver high-quality data and insights.

The DataOps Manifesto

The DataOps Manifesto is a set of principles to guide the practice of DataOps:

- **Continuous Integration and Continuous Deployment (CI/CD):** Automated testing and deployment of data pipelines.
- **Collaboration and Communication:** Enhanced cooperation between teams.
- **Data Quality Management:** Ensuring data accuracy, completeness, and reliability.
- **Agility and Iteration:** Rapid, iterative development and deployment.
- **Orchestration and Automation:** Automated workflows and data processes.

The DataOps Principles

- **Agility:** Embrace agile methodologies for data analytics.
- **Collaboration:** Foster collaboration between data teams.
- **Automation:** Automate repetitive tasks to reduce errors and speed up processes.
- **Continuous Improvement:** Regularly review and improve data processes.
- **Data Governance:** Maintain data quality, security, and compliance.

The DataOps Challenges

The Core Challenges of Data for Operations

- **Data Integration:** Combining data from different sources.
- **Data Quality:** Ensuring the accuracy and reliability of data.
- **Data Processing:** Efficiently transforming and processing data.
- **Data Security:** Protecting data from unauthorized access and breaches.

Why DataOps

Purpose: DataOps aims to improve the speed, quality, and reliability of data analytics processes, enabling organizations to make better and faster data-driven decisions.

What is the Purpose of Data-Driven Operations?

- **Efficiency:** Streamline data operations to reduce time and effort.
- **Quality:** Ensure high-quality data for accurate insights.
- **Agility:** Quickly adapt to changing business needs and data requirements.

How Does It Work?

DataOps involves automating data pipelines, implementing CI/CD for data workflows, enhancing collaboration between teams, and continuously monitoring and improving data processes.

Why DataOps

What Problem Does DataOps Solve?

- **Data Silos:** Integrates disparate data sources.
- **Data Quality:** Ensures accurate and reliable data.
- **Slow Processing:** Accelerates data preparation and analysis.
- **Collaboration:** Improves communication between teams.

Best Practices

Know Your Data

- **Data Discovery:** Identify and understand data sources.
- **Data Profiling:** Assess the quality and structure of data.

Trust Your Data

- **Data Quality Management:** Implement processes to ensure data accuracy and reliability.
- **Data Lineage:** Track the origin and transformations of data.

Use Your Data

- **Data Integration:** Combine data from various sources.
- **Data Analytics:** Extract insights and make data-driven decisions.

Improve DataOps

- **Committing Code:** Regularly update and refine data workflows.
- **Code Migration:** Manage changes across development, UAT, and production environments.
- **Collaboration and Version Control:** Use tools like Git to manage code and track changes.

GenAI and MLOps

What is Everyone Talking about in 2024

ChatGPT

**Artificial
Intelligence**

Chat Bots

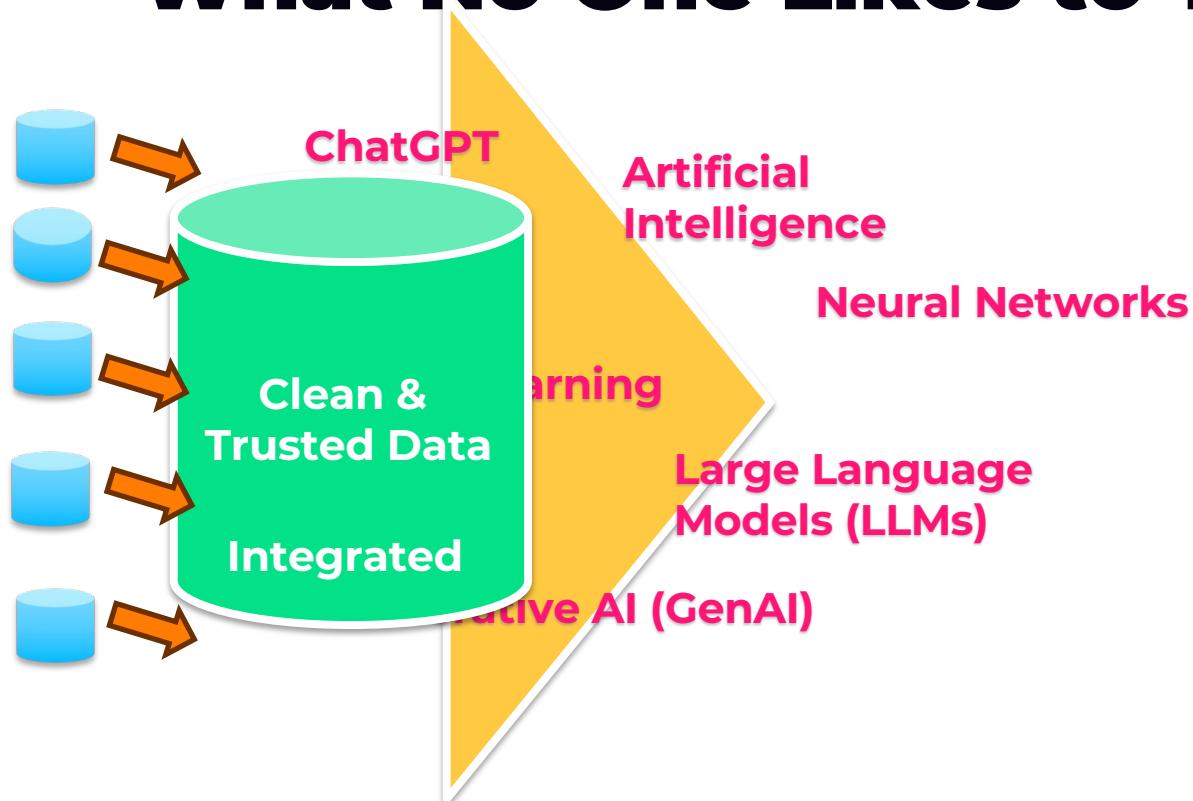
Neural Networks

Deep Learning

**Large Language
Models (LLMs)**

**Generative AI
(GenAI)**

What No One Likes to Talk About?



Role of Data Engineer in GenAI

- **Data Preparation:** Data engineers are responsible for collecting, cleaning, and organizing large datasets that are used to train generative AI models.
- **Data Pipelines:** They build and maintain data pipelines to ensure a continuous and efficient flow of data from various sources to the model training environments.
- **Data Storage:** They manage data storage solutions, ensuring that data is stored securely and efficiently, and is easily accessible for AI model training and inference.
- **Scalability:** They ensure that data infrastructure can scale to handle the large volumes of data required for training generative AI models.
- **Collaboration:** Data engineers work closely with data scientists and machine learning engineers to understand data requirements and optimize data workflows.

Role of Data Engineer in GenAI

Importance of the Data Engineering Role:

- **Data Quality:** Ensures high-quality data, which is critical for the performance of generative AI models.
- **Efficiency:** Optimizes data workflows to reduce time and resource consumption, enabling faster iterations and deployments.
- **Scalability:** Supports the scaling of AI models and data processing, crucial for handling large-scale AI applications.
- **Reliability:** Maintains the reliability of data pipelines and storage, ensuring consistent performance and availability of data for AI tasks.

Role of Data Engineer in GenAI

Importance of the Data Engineering Role:

- **Data Quality:** Ensures high-quality data, which is critical for the performance of generative AI models.
- **Efficiency:** Optimizes data workflows to reduce time and resource consumption, enabling faster iterations and deployments.
- **Scalability:** Supports the scaling of AI models and data processing, crucial for handling large-scale AI applications.
- **Reliability:** Maintains the reliability of data pipelines and storage, ensuring consistent performance and availability of data for AI tasks.

What is Generative AI?

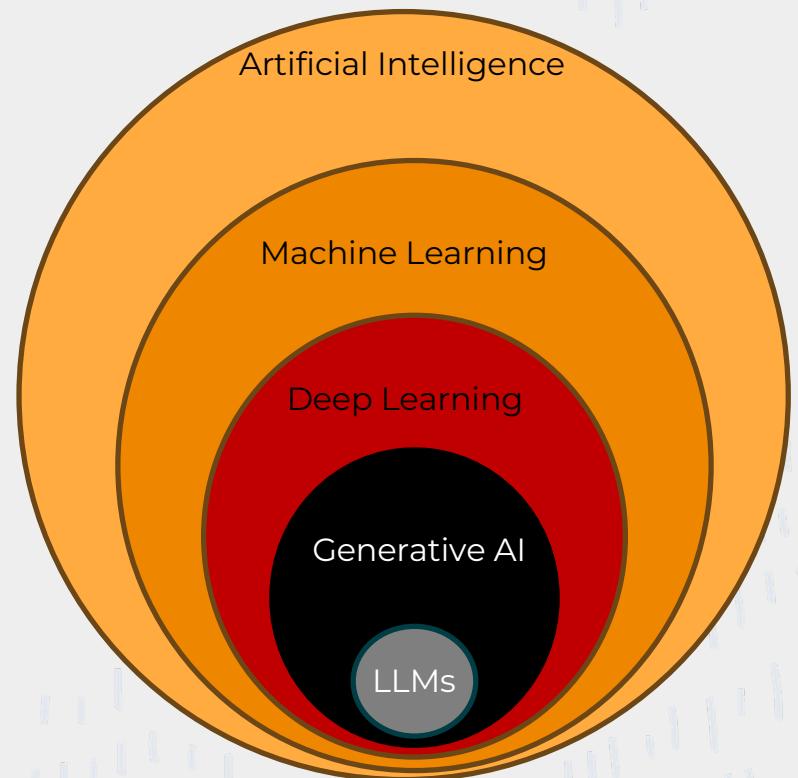
AI refers to the broad concept of machines or computers performing tasks that typically require human intelligence. This includes reasoning, learning, problem-solving, perception, language understanding, etc.

ML is a subset of AI focused on the idea that machines can learn from data, identify patterns, and make decisions with minimal human intervention

DL is a subset of ML that uses neural networks with many layers (deep networks) to model complex patterns in data.

Generative AI refers to a class of AI, often realized through DL, that focuses on generating new content or data that is similar to but distinct from the training data.

LLMs are a type of deep learning model designed to understand, generate, and interact with human language at a large scale. They are trained on vast amounts of text data.



Difference Between ML, DL, and GenAI

Machine Learning (ML):

- **Definition:** ML involves algorithms that learn patterns from data to make predictions or decisions without being explicitly programmed.
- **Examples:** Linear regression, decision trees, clustering.
- **Applications:** Fraud detection, recommendation systems, predictive analytics.

Deep Learning (DL):

- **Definition:** DL is a subset of ML that uses neural networks with many layers (deep neural networks) to model complex patterns in large datasets.
- **Examples:** Convolutional Neural Networks (CNNs) for image recognition, Recurrent Neural Networks (RNNs) for time series analysis.
- **Applications:** Image and speech recognition, natural language processing (NLP), autonomous driving.

Generative AI (GenAI):

- **Definition:** GenAI involves models that can generate new content, such as text, images, or music, based on the patterns learned from the training data.
- **Examples:** Generative Adversarial Networks (GANs), Variational Autoencoders (VAEs), Transformer models (e.g., GPT-3).
- **Applications:** Content creation, style transfer, text generation, synthetic data generation.

Discriminative vs. Generative Models

Discriminative Models:

- **Purpose:** Classify input data into categories.
- **Examples:** Logistic regression, Support Vector Machines (SVM), standard neural networks.
- **Function:** Learn the boundary between classes (e.g., spam vs. non-spam emails).

Generative Models:

- **Purpose:** Generate new data samples that resemble the training data.
- **Examples:** GANs, VAEs, LLMs.
- **Function:** Learn the underlying distribution of the data to create new, similar instances (e.g., generating realistic images of faces).

LLMs (Large Language Models)

What are LLMs?

Large Language Models (LLMs) are a type of generative AI model specifically designed to understand and generate human-like text based on the patterns and structures learned from vast amounts of textual data. Examples include GPT-3, BERT, and T5.

Importance of Data Quality

- **Accuracy:** Ensures that the data used for training models is correct and representative of the real world.
- **Completeness:** Guarantees that all necessary data is included, preventing models from learning incomplete or biased patterns.
- **Consistency:** Maintains uniformity in data formats and definitions across different datasets.
- **Reliability:** Ensures that data is dependable and available when needed for model training and inference.

Importance of Data Engineering Role

- **Data Accessibility:** Facilitates easy access to data for data scientists and analysts.
- **Efficiency:** Streamlines data workflows to reduce latency and improve productivity.
- **Scalability:** Ensures that data infrastructure can handle growing data volumes and processing needs.
- **Security:** Protects data from unauthorized access and breaches, ensuring compliance with regulations.

DevOps vs. DataOps vs. MLOps vs. LLMOps

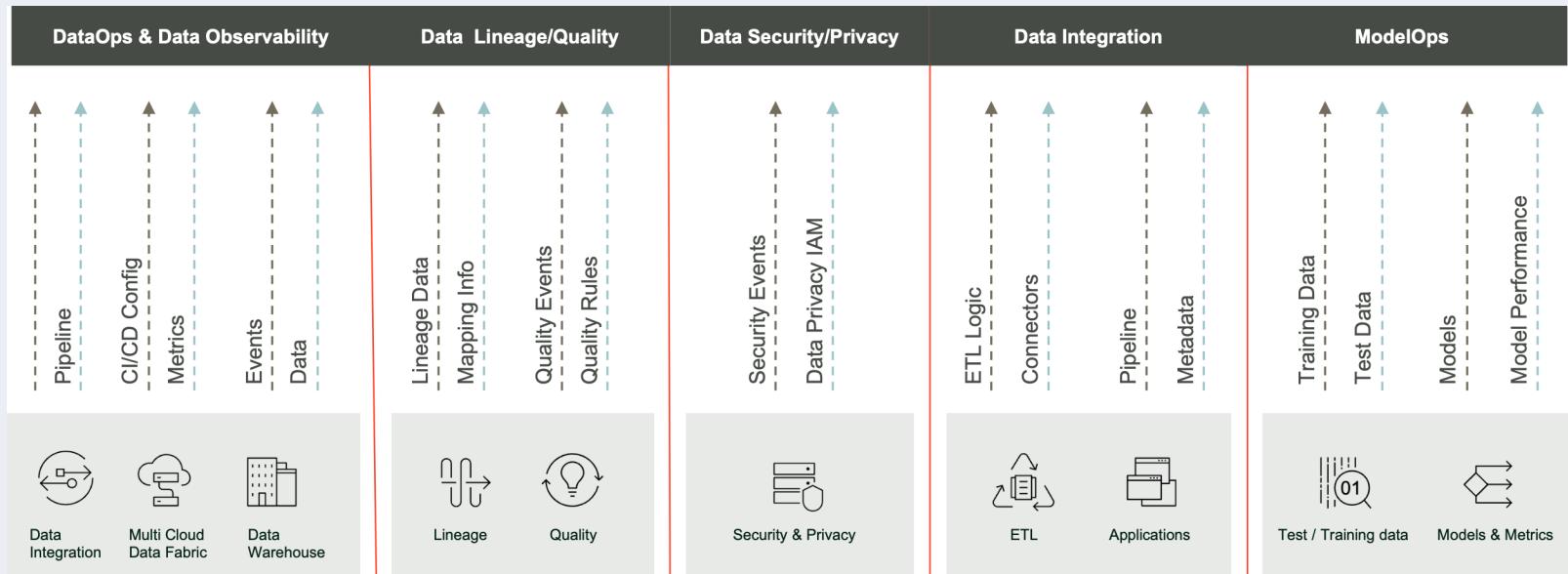
What is Common:

- **Automation:** All practices emphasize automating repetitive tasks to improve efficiency.
- **Collaboration:** Focus on enhancing collaboration between different teams (e.g., developers, operations, data engineers, data scientists).
- **Continuous Improvement:** Aim to continuously refine and improve processes.

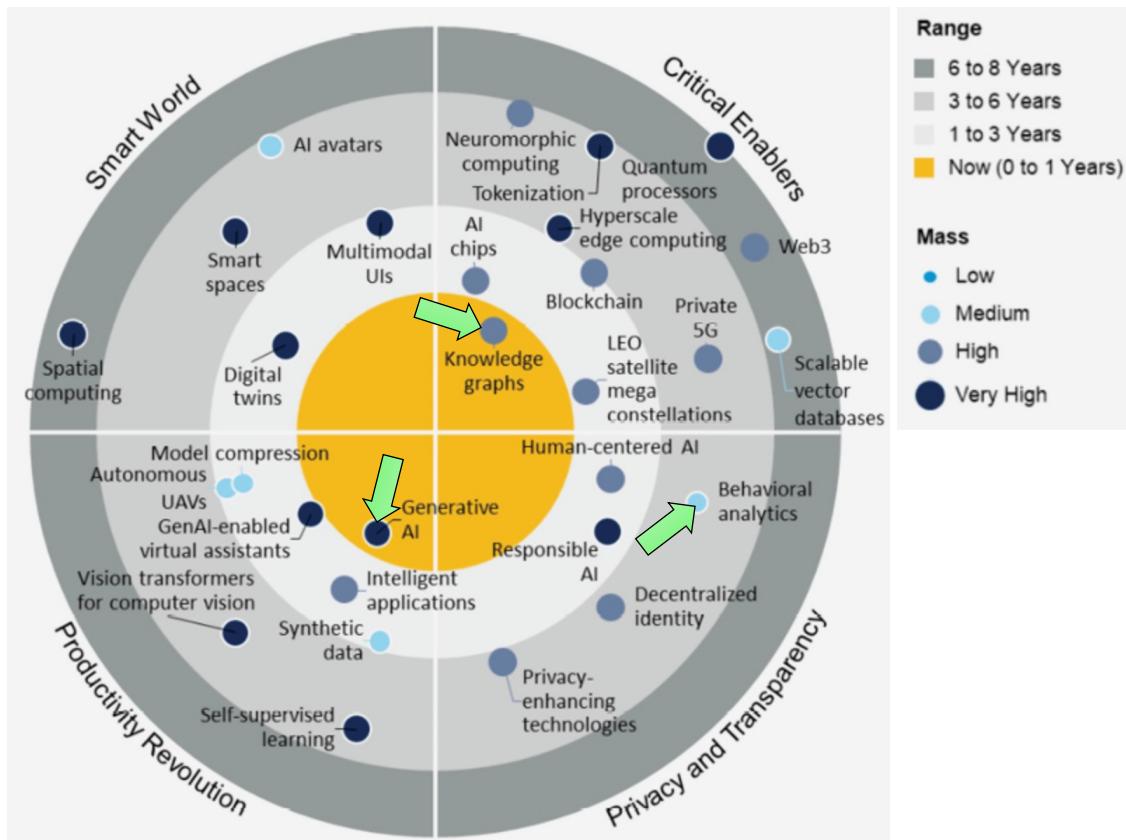
What is Different:

- **DevOps:** Concentrates on software development and IT operations.
- **DataOps:** Focuses on data management and analytics processes.
- **MLOps:** Addresses the deployment and maintenance of machine learning models.
- **LLM Ops:** Specializes in managing large language models, including ethical considerations and scalability challenges.

Example Framework

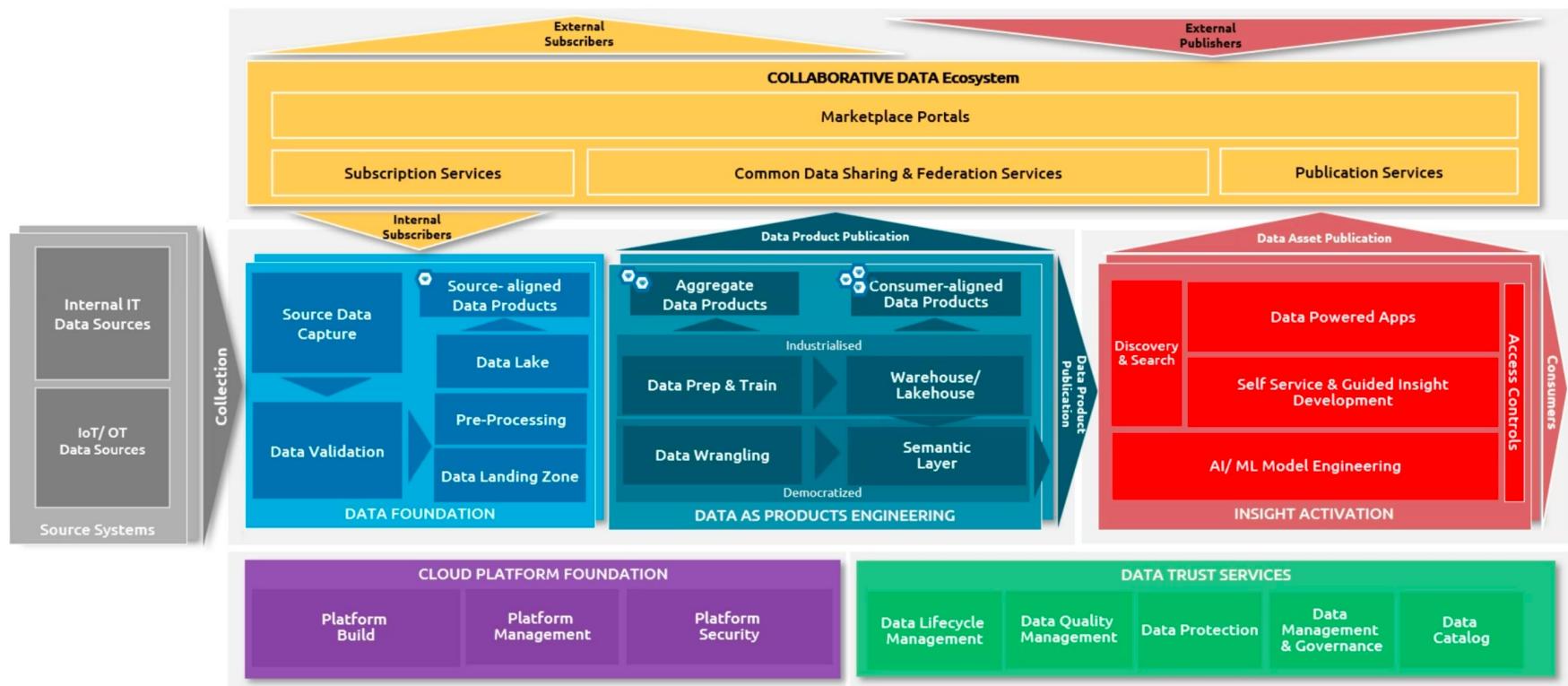


2024 Impact Radar



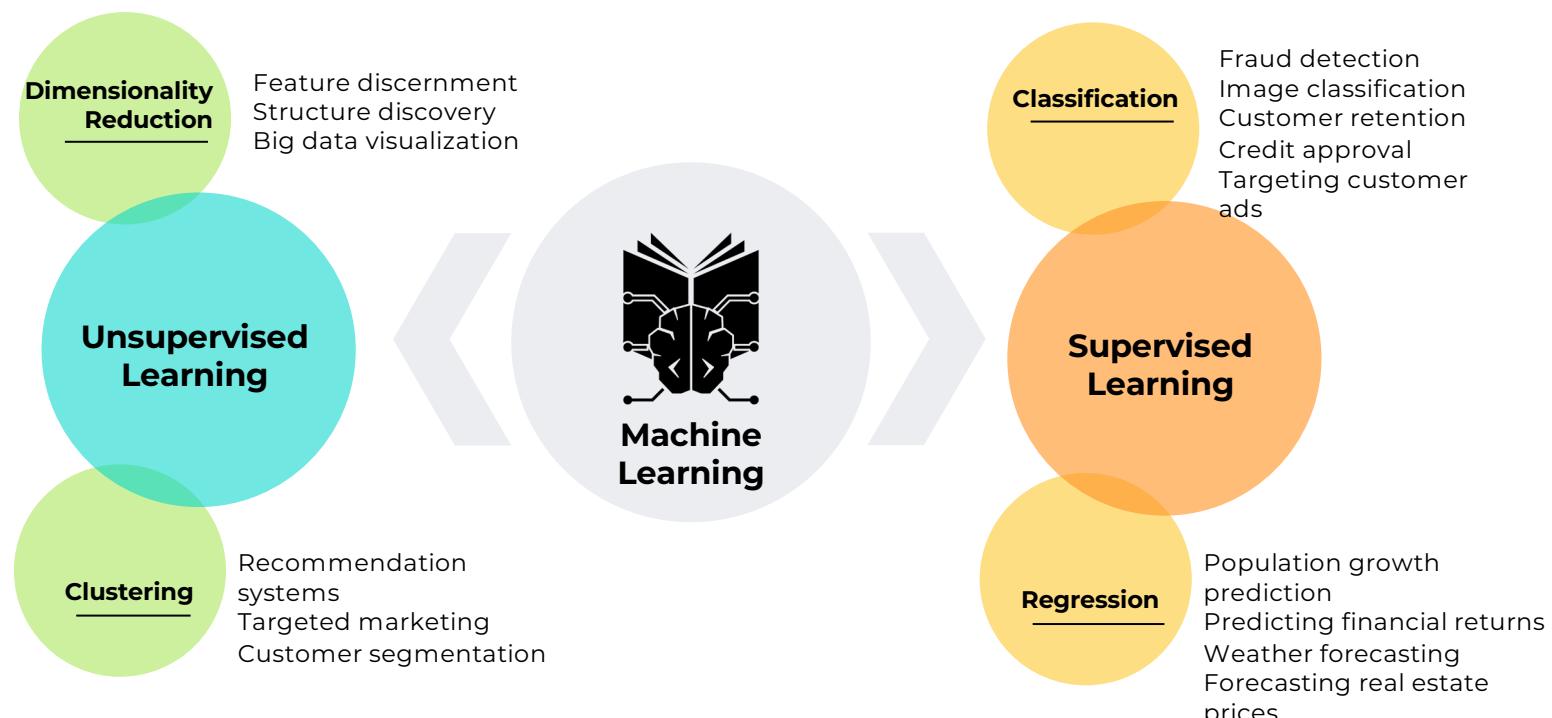
Modernize Infrastructure and Modern Data Architecture

Building a Data-Centric Architecture for the Future of



Machine Learning

Generally in Machine Learning there are **two main approaches:**



Data Governance

Proprietary and confidential



AWS for Data Governance

- The AWS Glue/Lake Formation technical data catalog
- AWS Glue DataBrew for profiling datasets
- AWS Glue Data Quality
- AWS Key Management Service (KMS) for data encryption
- Amazon Macie for detecting PII data in Amazon S3 objects
- The AWS Glue Studio Detect PII transform for detecting PII data in datasets
- Amazon GuardDuty for detecting threats in an AWS account
- AWS Identity and Access Management (IAM) service