# Get AIF-C01 Certified - Roadmap To Success

by Vladimir Raykov

```
┌─────────────────┐    ┌─────────────────┐    ┌─────────────────┐    ┌─────────────────┐
│ Intro + AI, ML, │ →  │ Mastering the   │ →  │ AWS AI Managed  │ →  │ Amazon          │
│ Deep Learning   │    │ Basics: AI & ML │    │ Services - Deep │    │ SageMaker AI    │
│ Overview        │    │ Concepts        │    │ Dive            │    │ Essentials And The │
│ (Section 1)     │    │ (Section 2)     │    │ (Section 3)     │    │ ML Dev. Lifecycle │
│                 │    │                 │    │                 │    │ (Section 4)     │
└─────────────────┘    └─────────────────┘    └─────────────────┘    └─────────────────┘

┌─────────────────┐    ┌─────────────────┐    ┌─────────────────┐    ┌─────────────────┐
│ Navigating      │ →  │ Business        │ →  │ Your GenAI Toolkit: │ → │ Applications of │
│ Generative AI: Core │ │ Applications of │    │ Essential AWS   │    │ FMs: Key Design │
│ Components,      │    │ Generative AI   │    │ Services and    │    │ Considerations  │
│ Model Types and │    │ (Section 6)     │    │ Features        │    │ (Section 8)     │
│ Lifecycle       │    │                 │    │ (Section 7)     │    │                 │
│ (Section 5)     │    │                 │    │                 │    │                 │
└─────────────────┘    └─────────────────┘    └─────────────────┘    └─────────────────┘
```

Intro + AI, ML, Deep Learning Overview
**(Section 1)**

Mastering the Basics: AI & ML Concepts
**(Section 2)**

AWS AI Managed Services - Deep Dive
**(Section 3)**

Amazon SageMaker Essentials & The ML Dev. Lifecycle
**(Section 4)**

Navigating Generative AI: Core Components, Model Types and Lifecycle
**(Section 5)**

Business Applications of Generative AI
**(Section 6)**

Your GenAI Toolkit: Essential AWS Services and Features
**(Section 7)**

Applications of FMs: Key Design Considerations
**(Section 8)**

Prompt Engineering Essentials And AI Vulnerabilities
**(Section 9)**

Fine-Tuning FMs - Deep Dive
**(Section 10)**

Evaluating FMs: Methods and Metrics
**(Section 11)**

Responsible AI Development: Key Concepts And Considerations
**(Section 12)**

Making AI Understandable
**(Section 13)**

Essential Security Practices and Tools on AWS
**(Section 14)**

Data Security and Governance for AI Systems
**(Section 15)**

AI Security, Compliance, and Governance
**(Section 16)**

Domain 1: Fundamentals of AI and ML
Domain 2: Fundamentals of Generative AI
Domain 3: Applications of Foundation Models
Domain 4: Guidelines for Responsible AI
Domain 5: Security, Compliance, and Governance for AI Solutions

# SECTION 1

# Key AWS Services And Concepts

1.  **Cloud Computing:** On-demand IT resources (compute, storage, databases) via the internet.

    a.   Benefits: Cost Efficiency, Scalability, Global Reach, Speed/Agility, Security.

2.  **AWS Global Infrastructure:**

    a.   <u>Regions</u>: Geographic data centers across the world.

    b.   <u>Availability Zones (AZs)</u>: Multiple data centers within each Region for redundancy.

    c.   <u>Edge Locations</u>: Local caches that speed up content delivery via Amazon CloudFront.

3.  **Core AWS Services:**

    a.   Computing: **Amazon EC2**: Virtual servers (instances) for various workloads (including AI/ML).

    b.   Storage: **Amazon S3:** Object storage (objects in buckets) for large datasets.

    c.   Database: **Amazon RDS**: Relational databases (structured data). **Amazon OpenSearch Service**: For semantic search.

    d.   Networking: **Amazon VPC**: Private network within AWS.

# AI vs ML vs Deep Learning Overview - Part 1

1. **AI is technology that <u>mimics</u> human intelligence.**
2. AI is the umbrella term that <u>includes</u> **Machine Learning** and **Deep Learning**.
3. AI learns and improves through <u>massive amounts of training data</u> (particularly true for Deep Learning Applications) - far more examples than humans need to learn similar tasks.
4. AI solves complex problems, with varying degrees of success, across many industries (e.g., healthcare, finance, retail, manufacturing, education, entertainment, transportation, energy, and more)
5. AI excels at <u>pattern recognition</u> and <u>processing large amounts of data</u>.
6. AI operates on **probability-based decision making**, providing confidence levels rather than absolute answers
7. Responsible AI development requires **transparency**, **fairness**, and **human oversight** to prevent misuse.

# AI vs ML vs Deep Learning Overview - Part 2

1. **Machine Learning (ML)** is the subset of AI focused specifically on learning from data.

    a. It transforms traditional programming by **learning patterns** rather than following explicit rules.

2. The **three main types** are:

    a. Supervised Learning

    b. Unsupervised Learning

    c. Reinforcement Learning

3. **Data quality and quantity** are even more crucial in ML than in general AI applications.

4. ML models improve continuously through exposure to more data.

5. The field faces unique challenges in **computation**, **algorithm selection**, and **explainability**.

# AI vs ML vs Deep Learning Overview - Part 3

1.  **Deep Learning (DL)** is a specialized subset of ML that uses **Neural Networks** <u>with multiple layers</u>.

    a.  It excels at **<u>processing unstructured data</u>** like images, text, and sound.

2.  Deep Learning can <u>automatically discover important features</u> in data **without** human guidance.

    a.  It requires <u>massive amounts of data</u> and <u>computational power</u> to train effectively.

3.  The technology powers many modern AI applications, including content generation and natural language processing.

4.  Deep Learning systems can be more complex and harder to interpret than traditional Machine Learning models.

    a.  Its architecture is inspired by the human brain's neural networks.

# Neural Networks Overview

1. **Neural Networks** consist of **input**, **hidden**, and **output layers**, each playing a crucial role in information processing.

2. Different <u>architectures</u> serve different purposes: **CNNs** excel at visual tasks, **RNNs** handle sequential data, and **Transformers** revolutionized language processing.

3. The Transformer architecture and GPT Models represent <u>significant advances</u> in AI, particularly in natural language understanding.

4. <u>Interpretability</u> remains a major challenge, especially as neural networks are deployed in critical applications.

5. Neural Networks can process various types of input data, from images to text to sensor readings.

6. The complexity of modern neural networks requires significant computational resources but <u>offers unprecedented capabilities</u>.

# Computer Vision & Natural Language Processing

1. **Computer Vision** enables machines to <u>interpret visual data</u>, using deep learning models (that use CNNs architecture) for tasks such as **image classification** and **object detection**.

2. **Natural Language Processing (NLP)** allows computers to understand and generate human language, leveraging architectures like RNNs and Transformers for tasks like **translation** and **sentiment analysis**.

3. AWS Services such as **<u>Amazon Rekognition</u>** and **<u>Amazon Comprehend</u>** provide powerful tools for implementing Computer Vision and NLP solutions.

4. Deep Learning is the driving force behind the advancements in both Computer Vision and NLP, enabling machines to learn from raw data and perform complex tasks.

# SECTION 2

# Intro: Mastering the Basics: AI & ML Concepts - 1

1. The **ML Process** consists of three main steps:

   a. **Training Data → ML Algorithm → Model**

2. **Amazon SageMaker is AWS's fully managed ML service that:**

   a. Helps _build_, _train_, and _deploy_ ML models.

   b. Provides _built-in algorithms_ and _pre-trained models_.

   c. Offers automated _model tuning_.

3. Remember SageMaker as an **ecosystem** with key services (not a complete list):

   a. **Studio**: Your one-stop shop for building, training, and deploying ML models.

      i. It provides a web-based interface for coding, debugging, and visualizing data.

   b. **Ground Truth**: Simplifies data labeling for ML training datasets.

      i. It offers automated and human-in-the-loop labeling options for various data types.

# Intro: Mastering the Basics: AI & ML Concepts - 2

a. **Data Wrangler**: Makes <u>data preparation</u> and <u>feature engineering</u> easier.

    i. It provides tools to import, transform, and analyze data for ML.

b. **Autopilot**: Automates ML model development by <u>exploring different algorithms</u> and <u>hyperparameters</u>.

    i. It helps find the best model without deep ML expertise.

c. **Clarify**: Helps understand how models make predictions and <u>detect potential biases</u>.

    i. It improves transparency and trust in ML models.

d. **Canvas**: Enables business users to build ML models <u>without writing code</u>.

    i. It provides a visual, point-and-click interface for simplified ML.

# Data Types in Machine Learning

1. The <u>ML process</u> starts with **collecting** and **processing** **training data**.
2. Data **quality** and **preparation** <u>are crucial for model success</u>.
3. **<u>Data Categories by Labels:</u>**
   a. **Labeled data**: Comes with predefined tags/labels (used in Supervised Learning, where the labels guide the learning process)
   b. **Unlabeled data**: No predefined labels (used in Unsupervised Learning to discover hidden structures or relationships within the data). The absence of labels requires different learning algorithms to find patterns.
4. **<u>Data Categories by Structure:</u>**
   a. **Structured data**: Organized in tabular formats (like SQL databases, spreadsheets). This organization makes it easy to query and analyze the data.
      i. **Time-series data**: Data points collected at successive points in time, used for analyzing trends and patterns over time.
5. **Unstructured data**: No predefined format (like social media posts, images, text, videos). This type of data often requires specialized techniques for processing and analysis.

# Learning Types - Supervised Learning

1.  **Supervised Learning** involves training algorithms on **labeled data** to predict outcomes for new data.
    a.  **Classification** assigns input data to predefined categories (labels).
    b.  **Regression** predicts continuous values (numbers).
2.  Labeled data is crucial for Supervised Learning, providing the necessary information for models to learn.
3.  Popular classification algorithms include Classification Decision Trees, Support Vector Machines, K-Nearest Neighbors, and Random Forest.
4.  Popular regression algorithms include Linear Regression, Decision Tree Regression, and Random Forest.
5.  Supervised Learning is one of the **three main categories** of ML, alongside **Unsupervised Learning** and **Reinforcement Learning**.

# Learning Types - Unsupervised Learning

1. **Unsupervised learning** works with **unlabeled data** to discover <u>**hidden patterns**</u> and <u>**structures**</u>.

    a. **Clustering** groups similar data points together based on their characteristics (e.g., customer segmentation).

    b. **Dimensionality Reduction** simplifies complex data by reducing the number of features while retaining essential information (e.g., simplifying user preferences).

    c. **Anomaly Detection** identifies unusual data points or patterns that deviate from the norm (e.g., fraud detection).

    d. **Density Estimation** analyzes the distribution of data points to identify areas of high and low concentration (e.g., location planning).

2. A key difference from <u>Supervised Learning</u> is that <u>Unsupervised Learning</u> works <u>**without**</u> predefined labels or "correct answers."

# Learning Types - Reinforcement Learning & RLHF

1. **Reinforcement Learning [RL] is about learning through interaction and feedback**.
2. Key components include **agent**, **environment**, **state**, **action**, and **reward**.
3. RLHF incorporates <u>human preferences</u> into the learning process, specifically:
   a. Collects human feedback on model outputs.
   b. Trains a reward model based on human preferences.
   c. Fine-tunes the model to align with these preferences.
   d. Helps reduce undesirable behaviors.
   e. Particularly valuable for tasks where success is hard to define mathematically.
4. <u>[RL] Common applications include:</u>
   a. Self-driving cars
   b. Game AI
   c. Robotics
5. [RLHF] Important considerations:
   a. Human feedback can be valuable but also **expensive** and **subjective**.
   b. RLHF can be used for <u>fine-tuning</u> after <u>self-supervised learning</u>.
   c. The field combines elements of behavioral psychology and machine learning.

# Learning Types - Semi-Supervised Learning

1. **Semi-Supervised Learning** combines **labeled** and **unlabeled** data **simultaneously**.

2. It's cost-effective as it requires fewer labeled examples.

3. Perfect for real-world scenarios <u>with limited labeling resources.</u>

4. Works best when you have a small amount of labeled data and lots of unlabeled data

5. Examples of Semi-Supervised Learning applications:

    a. Fraud Identification

    b. Sentiment Analysis

    c. Document classification

# What Is Inference? - 1

1.  *Inference* is the process of using AI models to <u>make predictions</u> or <u>decisions</u> **based on new data.**
    a. **Batch processing:**
        i. Collects and processes data in bulk (high latency).
        ii. Suitable for large payloads (up to 100 MB per mini batch).
        iii. Uses **Amazon SageMaker Batch Transform** with S3.
        iv. More cost-effective as you only pay for processing time.
        v. Best for applications that can tolerate delays.
    b. **Real-time inferencing:**
        i. Processes data as it arrives (low latency).
        ii. Handles smaller payloads (up to 6 MB).
        iii. Uses **Amazon SageMaker Real-Time Endpoints**.
        iv. Costs more due to continuous running.
        v. Essential for applications requiring immediate responses.

# What Is Inference? - 2

2.  **The choice between batch and real-time depends on:**
    a.  Data volume.
    b.  Latency requirements.
    c.  Cost considerations.
    d.  Application needs.
3.  **Two Deployment Options**
    a.  **Serverless**
        i.  Handles smaller payloads (up to 6 MB).
        ii.  Low latency.
        iii.  Great for unpredictable traffic patterns but it has *"cold start"*.
    b.  **Asynchronous**
        i.  Handles larger payloads (up to 1 GB)
        ii.  Latency - up to 1h.

# When NOT To Use Artificial Intelligence (AI)

1. When you need 100% reproducible results, and AI's **probabilistic nature** could introduce variability.

    a. Some AI models are **deterministic**, but many, like deep learning, can produce slight variations.

2. For **simple problems** where traditional programming is more efficient.

3. When you need custom training but **lack adequate data**.

4. When your use case requires **complete transparency** in decision-making.

5. If the **cost-benefit analysis** doesn't justify an AI solution.

# SECTION 3

# Intro: AWS AI Managed Services Deep Dive

1. **AWS AI managed services** are categorized by their functions:
   a. *Text/Documents -* ***Comprehend, Translate, Textract***
   b. *Vision -* ***Rekognition***
   c. *Search -* ***Kendra***
   d. *Conversational AI -* ***Lex***
   e. *Speech -* ***Polly, Transcribe***
   f. *Business Solutions -* ***Fraud Detector, Personalize, Mechanical Turk***
   g. *Human Review -* ***Augmented AI (A2I)***
2. These are <u>managed services</u>, meaning **AWS handles the underlying infrastructure**.
3. **Pricing is consumption-based (pay-per-use):**
4. Each service has its own pricing metrics (tokens, minutes, images).
5. No deep ML expertise is required to use these services:
   a. They're **pre-trained** and **ready to use**.
   b. APIs make integration straightforward.

# Amazon Rekognition

1. **Amazon Rekognition** is a fully managed <u>computer vision service</u> that can analyze **images** and **videos** without ML expertise.
2. <u>Key Features to Remember</u>:
   a. *Face detection and analysis.*
   b. *Object and scene detection.*
   c. *Text extraction.*
   d. *Content moderation.*
   e. *Custom Labels for specialized use cases.*
3. **Content Moderation Capabilities:**
   a. *Detects inappropriate content.*
   b. *Helps maintain platform safety.*
   c. *Automated content filtering.*
4. **Custom Labels Important Points:**
   a. *Requires minimal training data.*
   b. *Builds specialized ML models.*
   c. *Perfect for unique business needs.*

# Amazon Transcribe

1. **Amazon Transcribe** is a <u>fully managed speech recognition service</u> that can **convert speech to text**.

2. **Key Features to Remember:**

    a. *Automatic speech recognition*

    b. *Speaker identification*

    c. *Custom Language Models*

    d. *Custom Vocabulary Capabilities*

3. It supports **batch** and **real-time (streaming) transcriptions**.

4. It supports <u>over 100 languages</u>.

# Amazon Translate

1. **Amazon Translate** is a <u>neural machine translation service</u> that utilizes **advanced deep learning models** to deliver accurate and natural-sounding translations.
2. It supports translation between **75 languages.**
3. **Customization Options:**
   a. *Custom Terminology*: *Define specific translations for unique terms, ensuring consistency for brand names and industry-specific vocabulary.*
   b. *Active Custom Translation*: *Influence translation output to match desired style and tone without building custom models.*
4. **Versatile Translation Modes:**
   a. **Real-Time Translation**: *Ideal for applications requiring immediate translations, such as live chats and customer support.*
   b. **Batch Translation**: *Efficiently handles large volumes of text or documents, streamlining extensive translation projects.*
5. **Automatic Language Detection:** *Identifies the source language of the input text, simplifying workflows when dealing with multilingual content.*
6. **Integration with AWS Services:** *Seamlessly integrates with other AWS services like S3 for storage, Lambda for serverless computing, and Polly for text-to-speech capabilities, enhancing functionality and ease of use.*

# Amazon Comprehend

1. **Amazon Comprehend** is an NLP service that **extracts insights** from text using machine learning.
2. <u>Core features include:</u>
   a. *Entity recognition*
   b. *Sentiment analysis*
   c. *Language detection*
   d. *Key phrase extraction*
3. Customization options allow businesses to create **custom classification models** and **entity recognition for industry-specific needs**.
4. Pay-as-you-go pricing applies, with a free tier available.
5. Seamless AWS integration enables scalable workflows with services like **S3, Lambda, DynamoDB, and QuickSight.**

# Amazon Lex

1. **Amazon Lex:** A service for building conversational AI interfaces, powered by the same technology as Alexa.

2. **Key Features:** Incorporates automatic speech recognition and natural language understanding; manages intents, utterances, and slots.

3. **Pricing**: Operates on a pay-as-you-go model with a free tier for the first year.

4. **Integration**: Offers seamless connectivity with other AWS services.

5. **Use Cases**: Ideal for developing interactive chatbots and virtual assistants across various domains.

# Amazon Polly

1. **Amazon Polly** is a **text-to-speech** service that turns written text into lifelike spoken audio.

2. **Amazon Polly's Four Engines:**
   a. *Standard: Fast and basic*
   b. *Neural: More natural*
   c. *Long-Form: Perfect for long readings, like audiobooks*
   d. *Generative: Most human-like, with emotions*

3. **Key Features:**
   a. *SSML: Controls how words are spoken*
   b. *Lexicons: Custom pronunciation dictionary*
   c. *Speech Marks: Timing information for each word*

# Amazon Fraud Detector

1. **Amazon Fraud Detector** is a fully managed service that leverages machine learning <u>and over 20 years of Amazon's fraud detection expertise</u> to help businesses **identify potentially fraudulent activities**.

2. **Automated Model Creation**: *Build custom fraud detection models without prior machine learning experience.*

3. **Real-Time Fraud Predictions**: *Evaluate events instantly through API calls, receiving <u>fraud risk scores</u> **between 0 and 1,000**.*

4. **Rule-Based Actions**: Define and implement rules to automate responses based on model scores and other variables.

5. **Seamless AWS Integration**: *Integrate effortlessly with services like **AWS CloudTrail** and A**mazon EventBridge** for enhanced monitoring and event-driven workflows.*

6. **Continuous Model Improvement**: *<u>Regularly retrain models</u> with new data to adapt to evolving fraud patterns.*

# Amazon Personalize

1. **Amazon Personalize** is a fully managed service that brings Amazon.com's recommendation technology to developers, enabling the **creation of personalized user experiences** without requiring machine learning expertise.

2. **Data Integration:** *Utilizes user interactions, user metadata, and item metadata to inform recommendations.*

3. **Pre-Built Recipes:** *Offers algorithms for user personalization, similar items, and personalized ranking.*

4. **Advanced Capabilities:** *Addresses cold start problems, provides real-time updates, and supports A/B testing.*

5. **AWS Integration:** *Seamlessly works with services like* *Amazon SageMaker Data Wrangler* *and* *AWS Amplify*.

6. **Flexible Pricing:** *Pay-as-you-go model with a generous free tier for initial exploration.*

# Amazon Kendra

1. **Amazon Kendra** is an **intelligent enterprise search service** that uses natural language processing to deliver accurate and context-aware answers.

2. **Key features include:**

   a. ***Natural Language Understanding:*** *Interprets and responds to conversational questions.*

   b. ***Contextual Answers:*** *Provides precise information extracted from relevant documents.*

   c. ***Adaptive Learning:*** *Improves search relevance based on user interactions.*

   d. ***Built-In Domain Knowledge:*** *Offers expertise across multiple industries.*

   e. ***Diverse Data Source Support:*** *Integrates with various platforms and repositories.*

   f. ***Robust Security Controls:*** *Ensures users access only authorized information.*

   g. ***Direct Answers and Document Links:*** *Delivers concise responses with source references.*

   h. ***Continuous Improvement:*** *Learns and adapts from user behavior to enhance accuracy.*

# Amazon Textract

1.  **Amazon Textract** is an __intelligent document processing service__ that surpasses basic OCR by understanding document structures.

2.  __Key capabilities include:__
    a.  *Extracting text and handwriting.*
    b.  *Processing forms and tables.*
    c.  *Identifying key-value pairs.*
    d.  *Comprehending document layouts.*
    e.  *It supports various document types such as PDFs, images, and scanned files.*
    f.  *Textract provides confidence scores, maintains spatial information.*
    g.  *It offers both synchronous and asynchronous processing options.*

# Amazon Forecast

1. **Amazon Forecast** is a fully managed **forecasting** service using ML.

2. **Key capabilities:**

    a. *Automatic algorithm selection.*

    b. *Weather integration.*

    c. *Probabilistic forecasts.*

    d. *Forecast explainability.*

3. **Common use cases:**

    a. *Retail demand planning.*

    b. *Supply chain optimization.*

    c. *Energy consumption prediction.*

    d. *Financial planning.*

# Amazon Mechanical Turk

1. **Amazon Mechanical Turk (MTurk)** is a <u>crowdsourcing marketplace</u> for human intelligence tasks.
2. **<u>Key capabilities:</u>**
   a. *Global workforce access*
   b. *Quality control tools*
   c. *Worker qualification systems*
   d. *Automated task distribution*
3. **<u>Common use cases:</u>**
   a. *Data labeling*
   b. *Content moderation*
   c. *Survey completion*
   d. *Data verification*
   e. *Image/video processing*
   f. *Data collection*

# Amazon Augmented AI (A2I)

1. **Amazon Augmented AI (A2I)** improves the accuracy of ML predictions by <u>**incorporating human review into your workflows**</u>, particularly beneficial for **low-confidence** predictions or complex data.

2. A2I automates the routing of these uncertain predictions to human reviewers, streamlining workflows and minimizing manual effort.

3. A2I integrates seamlessly with **Amazon Textract** for document processing and **Amazon Rekognition** for image and video analysis, providing pre-built human review workflows for common use cases.

4. Its API also allows integration with custom ML models, offering flexibility for unique applications.

5. **A2I's pay-as-you-go pricing model** provides a **cost-effective solution** for maintaining high-quality results while efficiently managing expenses.

SECTION 4

# Amazon SageMaker AI - Overview

1. **Amazon SageMaker AI** is a fully managed service for <u>building machine learning (ML) models and foundation models (FMs)</u>.
2. **Unified Studio** integrates data preparation, analytics, training, and deployment.
3. **Pre-built tools:** Tools like **JumpStart** and **Partner AI Apps** accelerate workflows.
4. **HyperPod:** Enables efficient large-scale training by leveraging thousands of GPUs or accelerators.
5. **Automated Model Tuning:** SageMaker simplifies hyperparameter tuning to optimize model performance.
   a. *<u>Hyperparameters</u> are settings defined <u>**before training**</u> that control the learning process, such as the learning rate or the number of layers in a neural network.*
   b. *<u>Model parameters</u> are internal variables like weights in a neural network. These parameters are learned automatically <u>**during training**</u>.*
6. **Built-in Security:** Features like **SageMaker Catalog** ensure compliance with organizational policies.
7. **Scalable Deployment:** Deploy models at scale with cost-effective, pay-as-you-go pricing.
8. **Generative AI Assistance:** Tools like **Q Developer** simplify workflows with natural language guidance.

# Machine ML Development Lifecycle - Recap 1

1. The Machine Learning Development Lifecycle is **an iterative process** that guides the development of ML models. It includes several phases that build upon one another.
   a. *Each phase—such as defining the business goal, processing data, developing models, and deployment—feeds into the next, and model **monitoring** and **retraining** <u>ensure long-term success</u>.*
   b. *The process is NOT linear, meaning that <u>you might revisit earlier steps</u> based on the insights and results you get at later stages.*
2. **Business Goal Identification:**
   a. *Start by <u>clearly defining the business problem</u> you're solving.*
   b. *Ensure <u>alignment with stakeholders</u> and determine if ML  is the right approach.*
3. **ML Problem Framing:**
   a. *Translate the business problem into an ML problem.*
   b. *Choose the appropriate approach (e.g., classification or regression) and define success metrics.*

# Machine ML Development Lifecycle - Recap 2

4.  **Data Processing:**
    a.  ***Data Collection:*** *Identify and gather the relevant data, ensuring its quality and compliance with privacy standards (**AWS S3**).*
    b.  ***Data Preparation***: *Clean and transform data to make it suitable for modeling (**AWS SageMaker Data Wrangler**).*
    c.  ***Feature Engineering***: *Create and select features that will help the model learn patterns (**AWS SageMaker Feature Store**).*
5.  **Model Development:**
    a.  **Training:**
        i.   Split your data into training (80%), validation (10%), and test (10%) sets.
        ii.  Train the model on the training data while using the **validation set** to *fine-tune hyperparameters*.
        iii. Key concepts (hyperparameters):
            1.  **Batch size:** Determines **how many samples** the model processes at once. *Smaller batch sizes can improve accuracy but take longer; larger batches speed up training but require more memory.*
            2.  **Learning rate:** Controls **how much the model adjusts** after each batch. *A low learning rate results in slow learning, while a high learning rate results in overshooting the optimal solution.*

# Machine ML Development Lifecycle - Recap 3

        3.     **Epochs:** <u>**Refers to one full pass through the entire dataset.**</u> *Too few epochs can result in underfitting; too many can cause overfitting.*

    b.  <u>**Tuning:**</u> *Adjust hyperparameters (e.g., batch size, learning rate) to optimize model performance.*

    c.  <u>**Evaluation:**</u> *Assess the model against the <u>**test set**</u> to ensure it meets business metrics.*

    d.  **SageMaker Model Registry** *stores and tracks model versions, making it easy to manage and transition your models from training to deployment.*

6.  **Model Deployment:**

    a.  Choose between **self-hosted API** or **managed API** **deployment** based on the need for flexibility and technical resources.

    b.  Deployment Options (AWS SageMaker):

        i.  **Real-Time Inference:** *Instant processing with payload size limits (<u>**up to 6MB**</u>). Processing time: up to 60 sec. Latency: <u>**millisecond**</u>.*

        ii.  **Serverless Inference:** *Scales automatically based on demand and only charges for compute time used (<u>**up to 4MB payload**</u>). Processing time: up to 60 sec. Latency is low but may experience <u>**cold starts**</u>.*

# Machine ML Development Lifecycle - Recap 4

       iii.    **Batch Transform:** Designed for offline processing of large datasets available upfront. Payload size: **<u>Gigabytes (GBs)</u>**. Processing time: **could be days**.

       iv.    **Asynchronous Inference:** Suitable for large payloads that do not require immediate responses (**<u>up to 1GB</u>**). Processing time: up to 1 hour. Latency: **<u>near real-time</u>**.

7. **Model Monitoring:**

    a.    *Monitor for potential issues like **<u>data drift</u>** (changes in input data), **<u>concept drift</u>** (changes in the relationship between input and output data), and performance degradation.*

    b.    Key types of monitoring:

       iii.    **<u>Data Quality Monitoring:</u>** Tracks data changes (e.g., invalid inputs).

       iv.    **<u>Model Performance Monitoring</u>**: Tracks performance changes over time.

       v.    **<u>Bias Monitoring</u>**: Ensures fairness across different groups.

       vi.    **<u>Feature Attribution Drift</u>**: Monitors how features impact model decisions.

    c.    **Amazon SageMaker Model Monitor** provides automated monitoring for these issues.

8. **Model Retraining:**

    a.    *Retrain models based on performance degradation or changes in data.*

    b.    *Incorporate new data and automate the retraining process with CI/CD pipelines for continuous improvement.*

# ML Development Lifecycle vs ML Pipeline

1. **ML Development Lifecycle**
   a. Covers the entire ML project journey, from strategy to deployment.
   b. Includes business strategy and stakeholder involvement.
   c. <u>Flexible and iterative process</u>.
   d. Focuses on *"what"* and *"why."*
2. **ML Pipeline**
   a. **<u>Automates technical steps</u>** like data preprocessing, model training, and deployment.
   b. Standardized, reproducible workflows focusing on *"how".*
3. <u>**Key Differences**</u>
   a. **Scope**: Lifecycle is strategic; Pipeline is tactical.
   b. **Automation**: Pipeline automates technical tasks; Lifecycle involves manual processes.
   c. **Purpose**: Lifecycle guides the project; Pipeline executes the technical steps.
4. <u>**SageMaker Pipelines Highlights**</u>
   a. *CI/CD service for ML workflow automation.*
   b. *Operates on serverless infrastructure.*
   c. *Automates and scales workflows intelligently.*
   d. *Tracks model lineage for compliance and transparency.*

# MLOps Fundamentals

1. MLOps is the practice of systematically managing the ML lifecycle, **making ML projects reliable and scalable**.
2. **Key principles include:**
   a. *Version control for tracking all changes*
   b. *Automation of repetitive tasks*
   c. *CI/CD for seamless testing and deployment*
   d. *Model governance for control and compliance*
3. **Major benefits are:**
   a. *Faster time to market*
   b. *Increased team productivity*
   c. *Enhanced reliability and repeatability*
   d. *Complete auditability*
   e. *Improved data and model quality*
   f. *Scalable systems that grow with your needs*
   g. *Reduced technical debt*
   h. *Production-ready deployments*

# Mapping The AWS SageMaker AI Features - 1

1. **General Tools For All Stages:**
   a. **Amazon SageMaker Studio:** *A web-based IDE for end-to-end ML workflows, supporting tasks like data preparation, training, deployment, and monitoring within a single interface.*
   b. **Amazon SageMaker Unified Studio:** *An expanded platform that integrates ML with broader AWS data and analytics services like Glue, Redshift, and Bedrock, ideal for enterprise-scale workflows.*
   c. **Amazon SageMaker Canvas:** *A no-code visual interface that enables users to prepare data, build, and deploy ML models without writing code. It streamlines the end-to-end ML lifecycle, making it accessible to business analysts and domain experts.*
2. **Data Preparation:**
   a. **SageMaker Ground Truth** *for building highly accurate training datasets for ML through data labeling.*
   b. **Amazon SageMaker Data Wrangler** *for visual data preparation and transformation.*
   c. **Amazon SageMaker Processing Jobs** *for large-scale data processing tasks.*
   d. **Amazon SageMaker Feature Store** *for centralized feature storage and sharing across projects.*

# Mapping The AWS SageMaker AI Features - 2

1. **Training and Tuning:**
   a. **Amazon SageMaker Training Jobs** *for model training and scaling infrastructure.*
   b. **Amazon SageMaker Autopilot** *for automated model building and hyperparameter tuning.*
   c. **Amazon SageMaker's integration with MLflow** *for tracking and comparing different experiments.*
   d. **Amazon SageMaker Processing Jobs** *for model evaluation and performance analysis.*
   e. **Amazon SageMaker Model Registry** *for versioning and managing models.*

2. **Deploy and Manage:**
   a. **Amazon SageMaker Endpoints** *for real-time inference with automatic scaling.*
   b. **Amazon SageMaker Batch Transform** *for batch processing scenarios.*
   c. **Amazon SageMaker Edge Manager** *for deploying models to edge devices (smartphones, cameras, sensors).*
   d. **Amazon SageMaker Model Monitor** *for performance monitoring, including data quality, model quality, bias drift, and feature attribution drift.*
   e. **Amazon SageMaker Clarify:** *Detects bias in predictions and provides insights into model behavior.*

# Amazon SageMaker AI -Model Sources And Selection

1. **Amazon SageMaker AI** offers **four approaches**, ranging from <u>least to most complex</u>:
   a. Pre-trained Models (via JumpStart)
   b. Built-in Algorithms
   c. Pre-configured Framework Containers (e.g. PyTorch, TensorFlow, MXNet)
   d. Custom Containers

2. **When analyzing model selection questions, look for these keywords:**
   a. *"Least effort" or "quick deployment"* → **Pre-trained models**
   b. *"Standard ML tasks" or "balanced approach"* → **Built-in algorithms**
   c. *"Custom training code" but "managed infrastructure"* → **Pre-configured frameworks**
   d. *"Complete control" or "specific requirements"* → **Custom containers**

3. <u>**When considering operational overhead:**</u>
   a. *Least overhead: Pre-trained models (but remember: fine-tuning adds some overhead)*
   b. *Moderate overhead: Built-in algorithms*
   c. *High overhead: Pre-configured frameworks*
   d. *Highest overhead: Custom containers*

# Tech. Perform. Metrics - For Classification Problems

1. **Accuracy:** *The percentage of total correct predictions (both positive and negative).*

2. **Precision:** *How many of our positive predictions were actually positive.*

3. **Recall:** *How many actual positive cases we caught.*

4. **F1-Score:** *The balance between precision and recall.*

5. **AUC-ROC:** *How well our model can distinguish between classes (0.5 is random, 1.0 is perfect).*

# Tech. Perform. Metrics - For Regression Problems

1.  **MAE:** *Average absolute difference between predictions and actual values.*

2.  **MSE:** *Average squared difference (penalizes big mistakes more).*

3.  **RMSE:** *Square root of MSE (brings us back to original scale).*

4.  **R-squared:** *How much variation our model explains (0 to 1).*

5.  **MAPE:** *Average percentage difference between predictions and actual values.*

# Business Performance Metrics

1. **Cost per User**
   a. *Measures efficiency and scalability*
   b. *Calculated by dividing total infrastructure costs by active users*
   c. *Should decrease over time with optimization*
2. **Development Costs**
   a. *Include data collection, computing resources, training, and maintenance*
   b. *Often higher than initially estimated*
   c. *Critical for accurate ROI calculations*
3. **Return on Investment (ROI)**
   a. *Key metric for stakeholders*
   b. *Considers both tangible and intangible benefits*
   c. *Calculated as (Benefits - Costs) / Costs*
4. **Customer Feedback**
   a. *Provides qualitative insights*
   b. *Measures real-world impact*
   c. *Collected through surveys, NPS, and usage analytics*

SECTION 5

# Foundation Models (FMs) - Overview

1. **Definition:** *FMs are large-scale, <u>pre-trained models</u> adaptable to a wide range of tasks through fine-tuning. They learn general patterns from massive datasets.*
2. **Key Architectures:** *Common architectures include Transformers (especially for language), CNNs (for images), RNNs (for sequential data), and GNNs (for graph data).*
   a. ***Transformers** are crucial for many modern FMs.*
3. **Training Process:** FMs are pre-trained using <u>self-supervised learning</u> on <u>large amounts of unlabeled data</u>. This allows them to learn general representations, which are then adapted through <u>fine-tuning with smaller labeled datasets for specific tasks.</u>
4. **Multimodality:** *Modern FMs are increasingly multimodal, processing and generating information across modalities like text, images, and code (e.g., Google's Gemini family, Amazon Nova).*
5. **Key Capabilities:** *FMs excel at language processing, visual comprehension, code generation, human-centered engagement (chatbots), and speech-to-text.*
6. **Prompt Engineering:** *Effective prompt engineering is crucial for eliciting desired outputs from FMs.*
7. ***AWS and Bedrock:*** *<u>Amazon Bedrock</u> is your <u>key AWS service</u> for accessing various FMs, including Amazon Titan and models from other providers.*
8. **Key Considerations:** *FMs require significant resources for training, can have reliability and accuracy issues, and raise ethical concerns regarding bias.*

# Large Language Models (LLMs) - Overview

1. **What are LLMs?** Large Language Models (LLMs) are specialized Foundation Models designed to **understand and generate human language**.
2. **How do LLMs work?**
   a. **Tokenization:** *LLMs break text into smaller units called **tokens**, like words or parts of words, to process language more effectively.*
   b. **Transformer Architecture:** *Powered by **self-attention**, transformers excel at understanding word relationships, keeping context over long text, and processing information quickly.*
   c. **Pre-training:** *LLMs learn language patterns and context by processing massive datasets using **self-supervised learning**.*
3. **Why are LLMs important?** *They have transformed technology by enabling advanced applications like:*
   a. **Content Creation:** *Writing, summarizing, and analyzing text.*
   b. **Code Generation:** *Tools like GitHub Copilot and Amazon CodeWhisperer assist developers with coding tasks.*
   c. **Language Translation:** *Accurate and context-aware translations.*
   d. **Customer Service:** *Chatbots and virtual assistants for natural, human-like interactions.*
4. **AWS Connection:**
   a. **Amazon Bedrock** makes it easy to access LLMs, including **Amazon Titan Text** and third-party models, for various applications.

# Tokens, Embeddings and Vectors

1. **Tokens**: *The smallest text units used by LLMs for processing, often broken into subword pieces for efficiency.*
2. **Embeddings**: *Numerical vector representations of tokens that capture meaning and relationships in a high-dimensional space.*
3. **Vector Space**: *The mathematical environment where embeddings exist, allowing models to compare word relationships and perform similarity searches.*
4. **Semantic Search**: *Uses embeddings to find relevant content based on meaning, not just keywords—key for AI-powered search applications.*
5. <u>**AWS Services:**</u>
6. **Amazon OpenSearch** *enables vector-based search for semantic matching.*
7. **Amazon Bedrock** provides access to <u>**embedding models**</u>, such as **Amazon Titan Embeddings G1** and models from Cohere.
8. **Chunking**: *A method for handling long texts that exceed token limits by <u>splitting</u> them into overlapping segments while preserving context.*

# Multimodal Models

1. **Multimodal Models** can **process** and **generate <u>multiple types of data simultaneously</u>**, such as text, images, audio, and video.

2. They learn the relationships between **<u>different modalities</u>**, allowing them to combine and understand data in a more holistic way.

3. They differ from **LLMs (unimodal models)**, which are limited to **<u>text-based input and output</u>**.

4. They represent a major step forward in AI, enabling more human-like understanding and interaction with the world.

# Diffusion Models

1. **Diffusion Models** start with random noise and gradually refine it into meaningful outputs, like images, text, or audio.

2. They operate in two main steps:

    a. **Forward Diffusion:** Gradually adds noise to structured data until it becomes pure noise.

    b. **Reverse Diffusion:** Gradually removes noise from random data to generate a coherent output.

3. Popular services using diffusion models include **Stable Diffusion, DALL·E**, and **MidJourney**.

4. Their ability to learn patterns through noise makes them **incredibly powerful and versatile**.

5. **Important:** Diffusion models **cannot** interpret image content.

# The Foundation Model Lifecycle - 1

1. The **Foundation Model Lifecycle** consists of the following stages:

    a. **Data Collection & Preprocessing:** *Gathering and preparing high-quality, diverse data. Key considerations: Relevance, Diversity, Quality.*

    b. **Model Development (Architecture Selection & Configuration):** *Choosing the right architecture for the task. E.g. Transformers, CNNs.*

    c. **Pre-training:** *Training the model on large datasets to learn general patterns. E.g. a language model learns grammar, sentence structure, and word relationships.*

    d. **Fine-tuning:** *Adapting the model to specific tasks with smaller datasets.*

        i. **Key point:** *Fine-tuning is typically a **supervised learning process**, where the model learns from **labeled** examples.*

        ii. Additional technique: ***Reinforcement Learning from Human Feedback (RLHF)** can be used to align the model's behavior with human values and expectations.*

# The Foundation Model Lifecycle - 2

    e.   **Evaluation & Validation:** *Testing the model for **accuracy**, **robustness**, and **fairness**.*

        i.   Human evaluation and Benchmark datasets (e.g., GLUE, SuperGLUE) are often used for evaluation.

    f.   **Deployment:** *Integrating the model into real-world applications. E.g. embedding the model in a web application.*

    g.   **Monitoring, Feedback & Iteration:** *Continuously improving the model based on user feedback and new data. Retraining or fine-tuning the model to address issues.*

2.   AI models are not static - **they need to evolve** to remain **effective** and **reliable**.

# SECTION 6

# Disadvantages of Generative AI

1. **Regulatory Violations:** GenAI may break laws. → **Mitigation:** *Follow governance policies, monitor compliance, and involve legal teams.*

2. **Social Risks:** GenAI can spread misinformation or bias. → **Mitigation:** *Select training data carefully, audit for bias, and set ethical guidelines.*

3. **Data Security & Privacy:** Sensitive data might be exposed. → **Mitigation:** *Use anonymized data, encrypt information, and conduct security audits.*

4. **Toxicity:** AI can generate offensive content. → **Mitigation:** *Filter harmful data, use moderation tools, and apply human review.*

5. **Hallucinations:** AI may generate false information. → **Mitigation:** *Verify outputs, use RAG, browsing, and fine-tuning.*

6. **Interpretability:** AI decisions can be unclear. → **Mitigation:** *Use simpler models and explainability techniques.*

7. **Nondeterminism:** AI can give inconsistent answers. → **Mitigation:** *Use deterministic algorithms and perform extensive testing.*

# Advantages of Generative AI

1. **Key Advantages of GenAI:**

   a. **Adaptability**: *Learns and adapts to new contexts.*

   b. **Responsiveness**: *Provides real-time, immediate feedback.*

   c. **Simplicity**: *Simplifies complex tasks for users of all skill levels.*

   d. **Creativity**: *Enables creativity and innovative problem-solving.*

   e. **Data Efficiency**: *Operates efficiently with minimal input data.*

   f. **Personalization**: *Delivers personalized, context-aware outputs.*

   g. **Scalability**: *Scales to handle global-level enterprise needs.*

# Model Selection Factors for GenAI

1.   When selecting a GenAI Model, consider the following:

    a.   **Define your task:** *Know exactly what the model needs to accomplish.*

    b.   **Choose the right model type:** *Text, image, or multimodal capabilities.*

    c.   **Set performance expectations:** *Match the model's reliability and scalability to your needs.*

    d.   **Account for constraints:** *Assess your resources and deployment environment.*

    e.   **Evaluate capabilities:** *Look for features that meet your business requirements.*

    f.   **Ensure compliance:** *Prioritize ethical and regulatory standards.*

    g.   **Consider cost:** *Balance performance with budget and scalability.*

# Business Value Assessment

1. **User Satisfaction:** *Measure user happiness with AI outputs*

2. **Cross-Domain Performance:** *Assess AI performance across business areas.*

3. **Efficiency:** *Track time and effort saved by the AI.*

4. **Conversion Rate:** *Evaluate the AI's impact on customer conversion.*

5. **Average Revenue Per User (ARPU):** *Monitor revenue generated per user.*

6. **Customer Lifetime Value (CLV):** *Measure total customer lifetime revenue.*

# SECTION 7

# Amazon Bedrock - 1

1. Amazon Bedrock is a **fully managed service** that provides access to **foundation models** through a unified API, requiring **zero infrastructure management**.

2. The service offers **foundation models (FMs)** from leading AI companies including AI21 Labs, Anthropic, Cohere, Meta, Mistral AI, Stability AI and more.

3. **Three core features that make Bedrock powerful:**

   a. **[Important] Private model customization with your data.**

   b. Enterprise-grade **security** and **privacy**.

   c. Seamless integration with AWS services.

4. **Two flexible pricing options:**

   a. **On-demand:** Pay per use (tokens/images).

   b. **Provisioned throughput:** Reserved capacity for consistent workloads (mandatory for fine-tuning pre-trained models).

# Amazon Bedrock - 2

5.   **Key capabilities include:**

    a.   Model experimentation through **interactive playgrounds**.

    b.   **Knowledge Base** creation for enhanced responses.

    c.   **Agent** building for automated tasks.

    d.   **Fine-tuning** options for model customization.

    e.   Built-in **safety features** and **guardrails**.

6.   Amazon Bedrock is **continuously evolving**, with <u>new models</u> and <u>features</u> being added regularly.

> **Disclaimer:** This overview highlights key features of Amazon Bedrock. For a comprehensive understanding, please refer to the official AWS documentation and watch the rest of the videos related to Amazon Bedrock.

https://aws.amazon.com/documentation-overview/bedrock/

https://docs.aws.amazon.com/bedrock/latest/userguide/what-is-bedrock.html

# Why Should You Use the AWS infrastructure for GenAI?

1. AWS offers significant a**dvantages and benefits <u>for developing generative AI applications</u>**:

2. **Key Advantages:**

    a. **Accessibility:** *Broad access for all, from startups to enterprises, via pre-built services.*

    b. **Ease of Use:** *Lower barrier to entry with simplified tools and resources.*

    c. **Efficiency:** *Scalable infrastructure for faster model training and data processing.*

    d. **Cost-Effectiveness:** *Pay-as-you-go pricing for optimized cost management.*

    e. **Speed and Agility:** *Faster time to market and flexible business alignment.*

3. **Key Benefits:**

    a. **Security:** *Comprehensive security features, infrastructure protection, data safeguards, and AI-specific protections.*

    b. **Compliance:** *Adherence to global compliance standards for regulated industries.*

    c. **Responsibility:** E*mphasis on ethical and transparent AI practices.*

    d. **Safety:** *Tools for risk mitigation, bias detection, and content moderation.*

# PartyRock

1. PartyRock is a **no-code platform** for building AI applications.

2. You can get started without an AWS account.

3. **It's powered by Amazon Bedrock's foundation models.**

4. PartyRock is perfect for experimenting with AI applications in a user-friendly environment.

5. It includes a free trial period to explore its features.

6. And it offers a **smooth transition** to Amazon Bedrock for production-level projects.

# Amazon SageMaker JumpStart

1.  JumpStart simplifies ML workflows by offering access to over **350 built-in algorithms, pre-trained models, and pre-built solutions.**
2.  It includes **popular foundation models** like Stable Diffusion, BLOOM and more.
3.  **Integration with Amazon Bedrock enhances its capabilities.**
4.  Users are **responsible for managing the underlying infrastructure** and operational tasks.
5.  **Private hubs** allow curated model sharing within organizations but **don't support uploading custom models.**
6.  **Security and access control** are managed through **AWS Identity and Access Management (IAM) roles and policies**
7.  **SageMaker Instances:** *Virtual machines used for developing and training your machine learning models. You only pay when they are running.*
8.  **SageMaker Endpoints:** *Deployments of your trained models, accessible via a URL for making predictions. You pay for continuous availability.*

# Amazon Q

1. **Amazon Q Offers Two Versions:**
   a. **Amazon Q Business:**
      i. *For everyday users to access company documents, create content, and answer questions.*
      ii. *Integrates with tools like SharePoint, Jira, and Amazon WorkDocs.*
      iii. *Available as browser extensions (Chrome, Firefox, Edge).*
      iv. *Works with QuickSight for Business Intelligence (BI), focusing on data visualization.*
   b. **Amazon Q Developer:**
      i. *Supports developers (e.g. code conversion, troubleshooting, and AWS best practices).*
2. <u>**Core Features Across Both:**</u>
   a. **[Important]** *Built on Amazon Bedrock with multiple foundation models.*
   b. *Enterprise-grade security with permission-based access.*
   c. *Natural language interaction.*
   d. *Reads, creates, and moves data.*
   e. *Visual content analysis (diagrams, charts, PDFs).*
   f. *No-code application creation via Amazon Q Apps.*

# AWS AI/ML Service Stack - 1

1. **ML Frameworks Layer:**
   a. **Amazon SageMaker AI:** *For custom ML model development and end-to-end workflows.*
2. **AI/ML Services Layer (categorized by functionality):**
   a. **Text and Documents:**
      i. **Amazon Comprehend:** *Content understanding.*
      ii. **Amazon Translate:** *Language translation.*
      iii. **Amazon Textract:** *Document scanning and data extraction.*
   b. **Conversational AI:**
      i. **Amazon Lex:** *Chatbot development.*
   c. **Speech Services:**
      i. **Amazon Polly:** *Text-to-speech.*
      ii. **Amazon Transcribe:** *Speech-to-text.*
   d. **Vision Services:**
      i. **Amazon Rekognition:** *Image and video analysis.*
   e. **Enterprise Search:**
      i. **Amazon Kendra:** *Intelligent document search.*
   f. **Recommendations:**
      i. **Amazon Personalize:** *Personalized recommendation systems.*

# AWS AI/ML Service Stack - 2

3.   **Generative AI Layer:**
   a.   **Amazon SageMaker JumpStart:** *Quick access to pre-built models and solutions.*
   b.   **Amazon Bedrock:** *API access to foundation models.*
   c.   **PartyRock:** *Generative AI platform for specialized tasks.*
   d.   **Amazon Q Business:** *Workplace AI assistant for business users.*
   e.   **Amazon Q Developer:** *AI assistant tailored for developers and technical teams.*

# Cost Considerations for AWS AI/LM Services

1. AWS offers both **on-demand (pay-as-you-go)** and **provisioned** pricing models.

2. Higher **responsiveness and availability** come with increased costs, especially for <u>multi-region deployments</u>.

3. **Redundancy** across multiple locations provides reliability but requires additional resource costs.

4. **Performance options** vary from CPU to GPU, with costs increasing for higher-performance hardware.

5. **Token-based pricing** means you pay for the amount of content processed.

6. **Provisioned throughput** offers **guaranteed capacity** at a premium price.

7. **Custom model development and deployment** incur additional costs compared to using pre-trained models.

# SECTION 8

# Factors To Consider When Choosing A FM

1. Foundation models (FMs) are **pre-trained AI systems** that can be adapted for various tasks.
2. Here are the **top ten considerations** when choosing an FM:

    2.1. *Cost: Balance initial and operational expenses.*

    2.2. *Modality: Match the model to your data type (text, images, etc.).*

    2.3. *Latency: Prioritize fast response for real-time applications.*

    2.4. *Multi-lingual Support: Ensure the model handles required languages.*

    2.5. *Model Size: Fit the model size to your computational resources.*

    2.6. *Model Complexity: Choose sophistication based on task needs.*

    2.7. *Customization: Determine if fine-tuning is required.*

    2.8. *Input/Output Limits: Check data processing capacity.*

    2.9. *Responsibility: Mitigate bias, ensure ethics, and comply with regulations.*

    2.10. *Deployment & Integration: Ensure compatibility with existing systems.*

# Temperature, Top K, Top P, Input/Output Length - 1

1. AI models generate text by assigning probabilities to possible word choices. **<u>Inference parameters</u>** control how the model makes these selections:

2. **Temperature (Creativity Control):**
   a. *Low (0.2): Predictable outputs (e.g., reports, summaries).*
   b. *Medium (0.7): Balanced creativity and reliability (e.g., emails, ideas).*
   c. *High (1.0): Creative, diverse outputs (e.g., creative writing).*

3. **Top P (Nucleus Sampling): Controls cumulative probability.**
   a. *Low (0.3): Uses only top options until 30% cumulative probability.*
   b. *High (0.9): Includes more diverse options.*

4. **Top K: Limits the number of options considered.**
   a. *Top K = 5: Uses only top 5 options.*
   b. *Top K = 50: More variety, controlled choices.*

# Temperature, Top K, Top P, Input/Output Length - 2

5. **Context Window & Length:**
    a.  *Total input and output tokens must fit within the model's context window.*
    b.  *Example: With a 4,000-token window, 3,000 tokens input leaves room for 1,000 tokens output.*
6. **Common Combinations:**
    a.  *Creative Writing: High Temperature (0.7–1.0), High Top P (0.9)*
    b.  *Factual/Technical: Low Temperature (0.2–0.4), Low Top P (0.3)*
    c.  *Balanced: Medium Temperature (0.5–0.7), Medium Top P (0.7)*

# What Is Retrieval Augmented Generation (RAG)? - 1

1. Core Retrieval Augmented Generation (RAG) Concepts:
   a. **<u>RAG enhances foundation models (FM) without modifying their weights.</u>**
   b. Provides **real-time context** during inference **without** altering the model's "brain."
   c. RAG is an **architectural pattern**, not a training method.
2. **Key RAG Process** (Important for exam scenarios):
   a. *Ingest: Documents → Chunks → Embeddings → Vector Store*
   b. *Retrieve: Query → Vector → Similarity Search*
   c. *Augment: Combine Query + Retrieved Context*
   d. *Generate: FM produces contextually accurate response*
3. **AWS Services for RAG:**
   a. **Amazon S3:** Document storage
   b. **AWS OpenSearch:** Vector database operations
   c. **Amazon Titan Embeddings:** Vector generation
   d. **Amazon Bedrock / SageMaker AI**: FM inference

# What Is Retrieval Augmented Generation (RAG)? - 2

4. **Key Technical Benefits:**
    a. *No retraining required.*
    b. *Real-time knowledge updates.*
    c. *Reduced hallucinations (more accurate, grounded responses).*
    d. ***Cost-effective compared to fine-tuning.***

5. **Business Value:**
    a. *Adapts general-purpose FMs to specific domains.*
    b. *Maintains up-to-date knowledge.*
    c. *Provides traceable sources for responses.*
    d. *Scales efficiently with growing knowledge bases.*

# Amazon Bedrock - RAG and Knowledge Bases

1.  **Amazon Bedrock Knowledge Bases** **fully manages the entire RAG workflow**, from data ingestion to response generation.

2.  **Pre-processing** *involves document splitting, embedding creation, and vector storage.*

3.  **Runtime** *handles query processing, information retrieval, and response generation, including multi-turn conversations.*

4.  Multiple **data sources** (Amazon S3, Confluence, Salesforce, and SharePoint) and **vector database options** (OpenSearch Serverless, Amazon Aurora PostgreSQL, Pinecone, Redis Cloud, and MongoDB Atlas) are supported.

# Vector Database Options For Storing Embeddings

1. Key Databases for Vector Embedding Storage:
   a. **Amazon OpenSearch Service:** *Best for storing millions of embeddings, high-performance similarity search, horizontal scaling, and low-latency operations.*
   b. **Amazon Neptune:** *Ideal if you need both graph database features and vector search, maintaining complex data relationships.*
   c. **Amazon Aurora & Amazon RDS for PostgreSQL:** *Suitable for moderate vector storage needs, especially if you're already using relational databases and cost optimization is important. The pgvector extension enables vector search.*
   d. **Amazon DocumentDB:** *Best when your primary data model is document-based, and you need MongoDB compatibility along with vector search.*

2. **Exam Tips:**

3. *Focus on **similarity search**, **scalability**, and **performance requirements** when choosing the right database for a use case.*

4. *OpenSearch is your best bet for large-scale vector search, Neptune for graph-based needs, and Aurora/RDS/PostgreSQL for moderate workloads with relational data.*

# Foundation Model Customization Methods - 1

1. **Most Expensive: Pre-training**
   a. *Requires massive computational resources*
   b. *Involves extensive data processing*
   c. *Demands significant technical expertise*
2. **Moderately Expensive: Fine-tuning**
   a. Needs specialized hardware for training
   b. Requires domain expertise
   c. Involves ongoing maintenance
3. **Moderately Expensive: Continued Pre-training**
   a. *Uses unlabeled data*
   b. *Requires secure environment setup*
   c. *Needs domain-specific data preparation*
4. **More Affordable: RAG**
   a. *Primary costs come from storage and retrieval*
   b. *Requires vector database setup*
   c. *Lower computational demands*

# Foundation Model Customization Methods - 2

5.  **Most Cost-effective: In-context Learning**
    a.   No additional training costs
    b.   Uses existing model capabilities
    c.   Minimal technical expertise needed

6.  *Key Points to Remember*
    a.   *Pre-training builds foundational knowledge*
    b.   *Fine-tuning modifies model weights for specific tasks*
    c.   *Continued pre-training enhances domain knowledge using unlabeled data*
    d.   *RAG enhances responses without changing model weights*
    e.   *In-context learning requires no parameter updates*
    f.   *Cost increases with customization complexity*

# Amazon Bedrock: Agents

1.  **Amazon Bedrock Agents** <u>**automate multi-step tasks**</u>, freeing up your teams for higher-value work.

2.  They use **foundation models**, **APIs**, and **data** to break down complex workflows.

3.  <u>**Key features include:**</u>

    a.  *Easy setup*

    b.  *Memory retention*

    c.  *Built-in security*

    d.  *Multi-agent collaboration*

4.  <u>**Agents play crucial roles as:**</u>

    a.  *Intermediaries*

    b.  *Action launchers*

    c.  *Providers of valuable feedback for model improvement*

5.  **Action Groups** **simplify action management** and <u>**promote reusability**</u>.

# SECTION 9

# What Is Prompt Engineering?

1. **Prompt Engineering** is **optimizing textual input** to AI Models to **get desired responses**.
   a. *Think of it as effective communication with AI models.*
2. **Latent space** is the model's internal mathematical workspace where it processes prompts and generates responses.
3. **A well-crafted prompt has four key components:**
   a. *Instructions: specific tasks/directions for the model.*
   b. *Context: background information and frameworks.*
   c. *Input Data: specific information to process.*
   d. *Output Indicator: desired format/type of response.*
4. **Negative Prompting technique:**
   a. *Explicitly tells the model what to avoid/exclude.*
   b. *Used to eliminate unwanted content.*
   c. *Helps avoid specific writing styles.*
   d. *Prevents model assumptions.*
   e. *Maintains focus on relevant information.*

# Prompt Engineering Techniques

1. **Zero-Shot Prompting:**

   a. *Relies on the LLM's existing knowledge; best for simple, common-sense tasks.*

2. **Few-Shot Prompting (including Single-Shot)**:

   a. *Provides a few (or just one in the case of Single-Shot) examples to guide the model; effective for tasks requiring specific formats or styles.*

3. **Chain-of-Thought Prompting**:

   a. *Encourages step-by-step reasoning; crucial for complex problem-solving.*

4. **Prompt Templates:**

   a. *Predefined prompt structures; ensure consistency and streamline repetitive tasks.*

# Optimizing Prompts for Quality and Specificity

1.  **Clarity and Conciseness:** *Use simple, direct language for unambiguous instructions.*

2.  **Context is Key:** *Provide relevant background information for more accurate and relevant responses.*

3.  **Considering the Desired Output:** *Specify the format, style, and tone to guide the LLM to deliver results that meet your needs.*

4.  **Breaking Down Complex Tasks:** *Divide complex requests into smaller prompts for more effective processing and higher-quality outputs. This also facilitates more effective chain-of-thought prompting.*

5.  **Experimentation and Iteration:** *Refine your prompts through different phrases and approaches to discover what works best.*

6.  **Using Prompt-Level Controls:** *Fine-tune the LLM's response by specifying output length, restricting* topics, or defining the desired tone.

    a.  *Remember, while helpful for refining outputs, responsible AI relies primarily on **platform-level guardrails** and **ethical AI practices.***

# Understanding AI Vulnerabilities

1. **Exposure:** *Sensitive data used to train or operate an AI model is unintentionally revealed. This could include personally identifiable information (PII) like names, addresses, or purchase history.*

2. **Poisoning:** *Malicious data is injected into the training dataset to manipulate the model's behavior. This can cause the model to produce incorrect or biased outputs.*

3. **Model Hijacking (Backdoor Attack):** *A trained AI model is secretly modified to perform unintended tasks while still functioning normally. For example, an image recognition model for pets could be hijacked to also identify faces.*

4. **Prompt Injection:** *Carefully crafted prompts trick an LLM into bypassing its intended purpose or safety guidelines. This can lead to the generation of harmful, biased, or inappropriate content.*

5. **Jailbreaking:** *Safety restrictions in an LLM are bypassed, allowing it to generate content that would normally be prohibited. This can be achieved through specific prompts, exploiting internal logic weaknesses, or other techniques.*

SECTION 10

# Methods For Fine-Tuning An FM

1. **Instruction Tuning:** *Best when you want the model to follow specific instructions and generalize to new tasks with clear commands. It focuses on providing explicit guidance through natural language.*

2. **RLHF:** *Ideal for tasks where defining a clear objective function is difficult, such as generating creative text or engaging in open-ended dialogues. It uses human feedback to align the model's outputs with human preferences.*

3. **Transfer Learning (within Fine-Tuning):** *Fine-tuning itself leverages pre-trained knowledge. The key is deciding which layers to adjust: fine-tuning all layers maximizes adaptability (but is costly and can overfit), while freezing earlier layers is more efficient and prevents overfitting, especially for similar tasks.*

4. **Continuous Pre-training (as a precursor):** *Best when you have abundant unlabeled data in a specific domain. It enhances domain knowledge before fine-tuning for better results. This differs from standalone continuous pre-training.*

# Data Preparation For Fine-Tuning an FM

1. **To successfully fine-tune a model**, remember that data **quality**, **relevance**, and **specificity** are fundamental.

2. **Data Curation:** *Use **Amazon SageMaker Data Wrangler** for cleaning, transforming, and preparing your data, and Amazon S3 for storage.*

3. **Data Governance:** *Leverage **AWS Lake Formation** and **IAM** for data security and access control.*

4. **Data Size:** *Use **Amazon Athena** to query and analyze large datasets.*

5. **Data Labeling:** *Utilize **Amazon SageMaker Ground Truth** for efficient and accurate labeling (necessary for fine-tuning).*

6. **Data Representativeness:** *Use **Amazon QuickSight** for data visualization.*

7. **Feedback Integration:** *Monitor performance with **Amazon CloudWatch**.*

# SECTION 11

# Different Approaches to Evaluate Foundation Models

1. Evaluating Foundation Models is crucial for understanding their performance and comparing different models. We explored **two primary methods**:

2. **Human evaluation:** *Involves human assessors judging model outputs based on criteria like fluency, creativity, and accuracy. It captures nuanced aspects of quality <u>but can be subjective and resource-intensive.</u>*

3. **Benchmark datasets:** *Standardized datasets with labeled data used for <u>objective and quantifiable measurements.</u> They allow for direct comparison between models but may not perfectly represent real-world scenarios.*

4. <u>**Key Benchmark Examples:**</u>
   a. ***GLUE:*** *This dataset evaluates tasks like text classification, question answering, and natural language inference.*
   b. ***SuperGLUE:*** *A more challenging version of GLUE, focusing on tasks requiring deeper reasoning.*
   c. ***SQuAD:*** *Specifically designed for question answering tasks, where models must identify the answer within a given text passage.*
   d. ***ImageNet:*** *Used for image classification and object detection.*
   e. ***CodeXGLUE:*** *Focuses on code-related tasks like translation, completion, and summarization.*

5. <u>**Combined approach:**</u> Often, both **human evaluation** and **benchmark datasets** **are used together** to provide a comprehensive assessment of a Foundation Model's performance.

# Foundation Model Evaluation Metrics

1.  **Perplexity:** *Measures how "surprised" a language model is by a sequence of words; lower is better (typical range: 5-20 for good models).*
    a.  <u>**FMEval**</u> (an open-source library by AWS) in Amazon Bedrock and SageMaker AI can track perplexity.
2.  **BLEU:** *Used for <u>**machine translation**</u>, measures n-gram overlap with reference translations (an n-gram is a sequence of n words). Scores range from 0 to 1 (1 is perfect), with scores above 0.5 generally considered good.*
    a.  <u>**SageMaker Clarify**</u> can calculate BLEU.
3.  **ROUGE:** *Used for text generation tasks like <u>**summarization**</u> and <u>**translation**</u>, focuses on recall. Scores range from 0 to 1 (1 is perfect), and scores above 0.5 are generally considered good.*
    a.  <u>**SageMaker Clarify**</u> can calculate ROUGE.
4.  **BERTscore:** *Uses contextual embeddings to <u>**measure semantic similarity**</u> between generated and reference text. With scores above 0.8 often considered indicative of high semantic similarity.*
    a.  <u>**SageMaker Model Monitor**</u> can track BERTscore.
5.  **Accuracy & F1-score:** *Relevant when <u>**FMs are fine-tuned for classification tasks**</u>. Scores range from 0 to 1 (1 meaning perfect performance).*
    a.  <u>**SageMaker Model Monitor**</u> offers real-time monitoring of these metrics.

# SECTION 12

# Responsible AI - Key Concepts

1. Key Concepts For **Responsible AI**
    a. **Fairness**: *AI systems should be unbiased.* **Amazon SageMaker Clarify** *helps detect bias in data and models.*
    b. **Model Transparency**: Transparency is <u>key</u> for regulatory needs.
        i. **Interpretability**: *Understanding <u>how a model works</u> (internal mechanics).*
        ii. **Explainability**: *Understanding <u>what a model does</u>.*
        iii. <u>**AWS AI Service Cards**</u> provide transparency documentation for AWS AI services.
    c. **Performance-Transparency Tradeoff**: *Simpler models offer more transparency but potentially less performance.*
    d. **Privacy and Security**: *Protecting sensitive data is crucial.* **Amazon Comprehend** *helps with PII detection and removal.*
    e. **Veracity and Robustness:**
        i. **Veracity**: *Accurate and reliable data.*
        ii. **Robustness**: *Withstanding unexpected inputs.*
        iii. **Amazon Bedrock Guardrails** *help ensure robustness by filtering harmful content.*
    f. **Governance**: *Frameworks for responsible AI development.*
    g. **Safety**: *Minimizing harm to humans and the environment.*
    h. **Controllability**: *Maintaining human oversight.*
        i. **Amazon SageMaker Model Monitor** *detects model drift, ensuring continued control.*

# Responsible Model Selection: Environmental and Sustainability Considerations

1. **Energy Consumption:** *Large models can be energy-intensive. Choose efficient models and optimize training processes. Consider purpose-built hardware like **AWS Inferentia** and **Trainium**, and the fact that AWS has already achieved its goal of powering its operations with 100% renewable energy.*

2. **Resource Utilization:** *Minimize data storage and processing needs through efficient model selection. Leverage tools like **Amazon SageMaker Feature Store** and efficient processor architectures like **Graviton**.*

3. **Environmental Impact Assessment:** *Evaluate the potential environmental consequences of model deployment, considering the full hardware lifecycle and the efficiency of cloud infrastructure compared to on-premises.*

4. **Economic Considerations:** *Balance the economic benefits of AI with potential **social** and **environmental** costs like **job displacement.***

5. **Sustainability:** *Strive for AI systems that are **socially, environmentally, and economically sustainable** in the long term.*

# Generative AI: Legal and Ethical Concerns

1. **Toxicity:** *Generative AI can produce harmful outputs, including hate speech and offensive language, if not properly mitigated.*
   a. ***Example:*** *Abusive chatbot responses.*
2. **Hallucinations:** *Models can confidently generate incorrect or fabricated information.*
   a. ***Example:*** *Inventing sources in a research paper.*
3. **Intellectual Property:** *Questions arise about copyright ownership of AI-generated content.*
   a. ***Example:*** *Generating a melody similar to a copyrighted song.*
4. **Plagiarism/Cheating:** *Using AI to generate academic work raises ethical concerns.*
   a. ***Example:*** *Using AI to write an entire thesis.*
5. **Loss of Customer Trust:** *Ethical lapses can damage a company's reputation.*
   a. ***Example:*** *A news website publishing AI-generated articles with factual errors.*
6. **Biased Model Outputs:** *Models can perpetuate and amplify biases present in the training data.*
   a. ***Example:*** *An image generator struggling to depict women in leadership roles.*

# Balance - An Essential Characteristic of High-Quality Datasets

1. **Representativeness:** *Balanced datasets accurately reflect real-world scenarios for better generalization – meaning the model performs well on new, unseen data, not just the data it was trained on.*

2. **Bias Mitigation:** ***Balanced datasets*** *are crucial for preventing biased models and promoting fairness.*

3. **Data Strategies:** Achieve balance through **inclusive data collection**, **careful curation**, **data augmentation**, and **re-sampling**.

4. **Critical Applications:** *Balanced datasets are essential in sensitive areas like hiring, lending, and criminal justice.*

5. **AWS Tools:** *Use **Amazon SageMaker Clarify** and **Data Wrangler** for bias detection and data preparation.*

# Model Fit: Bias and Variance

1. **Statistical Bias:** *When a model makes overly simple assumptions.*

    a. *It misses important patterns in the data.*

    b. **High bias leads to underfitting.**

2. **Variance:** *When a model is too sensitive to small fluctuations in the training data.*

    a. *It performs well on training data but poorly on new data.*

    b. **High variance leads to overfitting.**

3. **Underfitting** *means the model is too simple and performs poorly on both training and new data.*

4. **Overfitting** *means the model is too complex and performs well on training data but poorly on new data.*

5. The goal is to **find a balance between bias and variance**, achieving a **good fit** that generalizes well to new data.

    a. *This balance requires both* **low bias** *and* **low variance**.

# SECTION 13

# AWS AI Service Cards

1. **AWS AI Service Cards** are essential resources for anyone using Amazon's AI services. They provide:

    a. **Intended Use Cases:** *Clear insights into what the service is designed for and its limitations.*

    b. **Responsible Design Principles:** *Information on how the service incorporates fairness, explainability, and other responsible AI practices.*

    c. **Best Deployment Practices:** *Practical guidance on deploying and using the service responsibly to ensure ethical and effective AI development.*

# SageMaker Model Cards

1. **Core Purpose: <u>SageMaker Model Cards</u>** serve as the **single source of truth** for your models. Consider the following key aspects:

   a. **Responsible AI:** *Promotes transparency and accountability.*

   b. **Key Information:** *Captures intended use, risk ratings, metrics, evaluations, and recommendations.*

   c. **Lifecycle Management:** *Supports the entire model lifecycle from development to deployment.*

   d. **[Important] Audit and Communication:** *Aids in audits and facilitates stakeholder communication.*

   e. **Version Control:** *Ensures an immutable record of all changes.*

   f. **Complementary to AI Service Cards:** *While AI Service Cards focus on entire AWS services, Model Cards document individual models built within SageMaker.*

# AWS SageMaker Clarify

1.  **Core Purpose:** <u>SageMaker Clarify</u> is designed to ensure **fairness**, **transparency**, and **accountability** in your machine learning models.
2.  <u>**Three Core Functionalities:**</u>
    a.  *Bias Detection: Identifies and mitigates biases in both training data and model predictions.*
    b.  *Model Explainability: Provides insights into how models make decisions.*
    c.  *Model Monitoring: Tracks model performance over time to detect drifts and degradation.*
3.  <u>**Key Features and Tools:**</u>
    a.  *Data Bias Analysis: Analyzes training data for imbalances in feature distributions.*
    b.  *Model Bias Analysis: Quantifies biases in model predictions across different groups.*
    c.  *Feature Importance: Identifies the most influential features.*
    d.  *Partial Dependence Plots: Visualizes the relationship between features and model output.*
    e.  *Individual Prediction Explanations: Provides insights into individual predictions.*
    f.  *Concept Drift Detection: Detects changes in input data distribution.*
    g.  *Performance Degradation Detection: Monitors key performance metrics.*
4.  <u>**Foundation Model Evaluation:**</u> *Includes specialized metrics for evaluating foundation models (e.g., toxicity, bias in generated text, factual accuracy).*

# The 3 Key Principles Of Human-Centered Design For Responsible AI

1. **Design for Amplified Decision-Making:** *AI acts as a tool that enhances human judgment.*

2. **Design for Unbiased Decision-Making:** *Implementing bias mitigation in practice.*

3. **Design for Human and AI Learning:** *Creating systems that improve through mutual feedback.*

# SECTION 14

# AI System Security & The AWS Shared Responsibility Model

1. **AI systems present unique security challenges**, *including vulnerabilities specific to large language models, as highlighted by the **OWASP (Open Web Application Security Project) Top 10 for LLMs***.

2. **A layered security approach (defense in depth)** *is essential for AI workloads, providing multiple layers of protection against various threats.*

3. The **AWS Shared Responsibility Model** divides security responsibilities between AWS and the customer.
   a. **AWS is responsible for "Security of the Cloud",** *managing the underlying infrastructure, physical security, and host operating systems.*
   b. **You are responsible for "Security in the Cloud"**, *managing guest operating systems, applications, data, access controls, and configurations like Security Groups.*

# Identity and Access Management (IAM)

1. **Identity And Access Management (IAM)** **is the security system for your AWS account.**
   a. **Users** *represent individuals or applications and have long-term credentials.*
   b. **Groups** *are collections of users for easier permission management. Groups cannot be nested.*
   c. **Policies** *define what actions are allowed or denied using JSON. Permissions are additive.*
   d. **Permissions** *are the specific abilities granted by a policy.*
   e. **Roles** *provide temporary permissions for users or AWS services and do not have credentials.*
2. Use **MFA** for enhanced security.
3. Use **groups** to manage permissions efficiently.
4. Create **individual IAM users** for each person or application.
5. **[IMPORTANT]** **Never use the Root User for daily tasks.**
6. Always follow **The Principle of Least Privilege**.
7. Use **IAM Access Analyzer** to validate your policies.

# Network Security for AI Workloads - AWS PrivateLink

1. **Private Connectivity:** *PrivateLink* provides **private connectivity** *between your VPC and AWS services, services hosted by other AWS accounts, and your own services in another VPC, without traversing (accessing) the public internet.*
2. **Endpoints:** *It uses endpoints within your VPC, which are virtual network interfaces with private IP addresses, to connect to the services.*
3. **Enhanced Security:** *By bypassing the internet,* *PrivateLink* *significantly reduces the attack surface and protects sensitive data.*
4. **Simplified Network Architecture:** It simplifies network configurations by eliminating the need for internet gateways, NAT devices, and public IP addresses for service communication.
5. **AI Use Cases:** *It's particularly valuable for AI workloads involving services like* *Amazon SageMaker AI* *and* *Amazon S3*, *ensuring secure data access and model training within your VPC.*
6. Remember that **your VPC is your private network in the cloud**, and **PrivateLink extends that private connectivity to AWS services.**

# SECTION 15

# Secure Data Engineering Best Practices

1. **Data Quality:** Focus on these core metrics:
   a. *Completeness, Accuracy, Timeliness, and Consistency: Ensuring your data is a reliable foundation for your AI models.*
2. **Privacy:** Protect sensitive data with:
   a. *Masking, Obfuscation, Differential Privacy, Encryption, and Tokenization.*
3. **Data Access Control:** Manage access and monitor activity using:
   a. *Data Governance Framework: Establishing clear policies.*
   b. *Role-Based Access Control (**AWS IAM**): Managing permissions using users, groups, and policies.*
   c. *Multi-Factor Authentication (MFA): Adding an extra layer of security.*
   d. *Monitoring and Logging (**AWS CloudTrail & Amazon CloudWatch**): Capturing API calls with CloudTrail and analyzing logs with CloudWatch to detect suspicious activity.*
4. **Data Integrity:** Ensure data accuracy, reliability, and consistency using:
   a. *Data Validation: Checking for errors.*
   b. *Backup and Recovery (**AWS Backup**): Implementing robust backup strategies.*
   c. *Transaction Management: Guaranteeing data consistency during operations.*
   d. *Data Lineage Tracking (**AWS Glue Data Catalog, AWS Lake Formation**): Understanding data's origin and transformations.*

# Data Provenance and Lineage

1. **Data Provenance:** *This refers to the **origin or source of your data**, directly impacting data quality.*

2. **Data Lineage:** *This tracks the **data's journey** – how it was collected, transformed, and used. This is closely related to data integrity and the need for audit trails.*

   a. **Importance:** *Data provenance and lineage are essential for **transparency**, **accountability**, and building trust in AI systems. These principles support the broader goals of secure data engineering.*

3. **Key Practices:** ***Source Citation** and **Data Cataloging** are crucial for establishing clear data provenance and lineage.*

4. **AWS CloudTrail & CloudWatch:** *These services provide the **Audit Trails** necessary for tracking data access and changes, contributing to both secure data engineering and data lineage.*

# Data Governance Strategies

1. **Data Governance:** *Encompasses the processes and policies ensuring data is fit for purpose within AI initiatives. It defines roles, responsibilities, and standards for data usage. Data Governance Strategies include:*
   a. **Data Lifecycles:** *Managing data from creation to disposal, focusing on quality and compliance. This includes establishing data retention policies.*
   b. **Data Logging:** *Systematically recording data processing activities for auditing and tracking, leveraging services like AWS CloudTrail and Amazon CloudWatch Logs. This is essential for maintaining audit trails and supporting data lineage.*
   c. **Data Residency:** *Understanding and adhering to regulations regarding the physical location of data storage.*
   d. **Data Monitoring:** *Continuous observation and analysis of data used in AI to maintain quality and performance.*

# SECTION 16

# Security and Privacy Considerations for AI Systems

1. **Application Security:** *Implementing <u>access controls</u> and security measures within the AI application.*

2. **Threat Detection:** *Actively **monitoring** for malicious activity and leveraging **AWS Shield** and **AWS WAF** (Web Application Firewall).*

3. **Vulnerability Management:** *Regularly **assessing** and **updating systems** to address weaknesses.*

4. **Data Protection:** *Implementing <u>encryption</u> and following the <u>Shared Responsibility Model.</u>*

5. **Specialized AI Protections:** *Addressing unique AI vulnerabilities like **prompt injections** through filtering and validation.*

# Regulatory Compliance Standards for AI Systems

1. **ISO (specifically ISO/IEC 27002):** *A globally recognized standard providing best practices for information security management. Think of it as a broad set of security guidelines, like anonymizing data.*
2. **SOC Reports:** *Independent audits of AWS's controls, demonstrating their commitment to security. Leverage these reports for your own compliance efforts. For example, verifying data protection in S3.*
3. **NIST 800-53:** *A robust set of security controls often used by U.S. federal information systems and as a benchmark for others. For example, controlling access to AI models.*
4. **EU AI Act:** *An emerging regulation categorizing AI by risk level, with banned, high-risk, and unregulated tiers.*
5. **Algorithmic Accountability Laws:** *Emerging regulations focused on fairness and transparency in automated decision-making. Stay informed about the laws relevant to you, such as explaining AI decision-making.*
6. **AWS Audit Manager:** *Automates evidence collection, offers pre-built and custom frameworks, and provides continuous auditing and reporting to simplify compliance.*

# AWS Services for Governance and Compliance

1. **AWS Config:** *Tracks changes to your AWS resources, enabling you to audit configurations and detect non-compliance.*
2. **Amazon Inspector:** *Automates security assessments of your applications and infrastructure, identifying vulnerabilities and deviations from best practices.*
3. **AWS Audit Manager:** *Automates evidence collection for audits, simplifying compliance with pre-built and custom frameworks.*
4. **AWS Artifact:** *Provides on-demand access to AWS compliance reports and agreements.*
5. **AWS CloudTrail:** *Records API calls made within your AWS account, providing a detailed audit trail of all actions.*
6. **AWS Trusted Advisor:** *Provides recommendations for optimizing your AWS infrastructure across cost, performance, security, fault tolerance, and service limits.*

# Governance Protocols and Frameworks for GenAI

1. **Generative AI Security Scoping Matrix:** *A structured approach for assessing and mitigating risks by evaluating factors like data sensitivity, model sensitivity, deployment environment, access control, and output validation. This is a key tool for proactive risk management.*
2. **Policies:** *Overarching rules and guidelines covering data privacy, ethical use, security, and regulatory compliance.*
3. **Review Cadence/Strategies:** *Regular reviews (code reviews, security audits, ethical reviews) to ensure ongoing compliance and identify emerging risks.*
4. **Transparency Standards:** *Being open about how your AI models work and the data they use to build trust.*
5. **Team Training:** *Equipping your team with the knowledge and skills to develop and deploy AI responsibly.*