

AWS AI Practitioner Study Guide

Domain 2: Fundamentals of Generative AI

Task Statement 2.1: Basic Concepts of Generative AI

1. Foundational Generative AI Concepts

1.1 Tokens

Definition: The basic units of text that AI models use to understand and generate language.

Key Concepts:

- **Tokenization:** Process of breaking text into smaller units (words, subwords, characters)
- **Token Types:** Word-level, subword-level (BPE, WordPiece), character-level
- **Token Limits:** Most models have maximum token limits (e.g., GPT models have context windows)
- **Token Efficiency:** Affects cost and performance in API-based models

Examples:

- Word tokens: "Hello world" → ["Hello", "world"]
- Subword tokens: "unhappiness" → ["un", "happy", "ness"]
- Special tokens: <start>, <end>, <pad>, <unk>

AWS Context:

- **Amazon Bedrock:** Different models have varying token limits and pricing
- **Token counting:** Important for cost estimation and prompt design
- **Context window management:** Critical for long-form content generation

1.2 Chunking

Definition: Breaking large documents or datasets into smaller, manageable pieces for processing.

Purpose:

- Overcome model context window limitations
- Improve processing efficiency
- Enable better retrieval in RAG systems
- Maintain semantic coherence

Chunking Strategies:

- **Fixed-size chunks:** Equal token/character counts
- **Semantic chunks:** Based on paragraphs, sections, or topics
- **Sliding window:** Overlapping chunks to maintain context
- **Hierarchical chunking:** Multiple granularity levels

Best Practices:

- Maintain context boundaries (don't split mid-sentence)
- Include overlap between chunks (10-20%)
- Consider document structure (headers, paragraphs)
- Balance chunk size vs. semantic completeness

AWS Implementation:

- **Amazon Kendra:** Intelligent document chunking for search
- **Amazon OpenSearch:** Vector storage for chunked content
- **AWS Lambda:** Custom chunking logic for preprocessing

1.3 Embeddings

Definition: Numerical representations of text, images, or other data in high-dimensional vector space.

Key Properties:

- **Semantic similarity:** Similar content has similar embeddings
- **Dense vectors:** Typically 768, 1024, or 1536 dimensions
- **Contextual:** Modern embeddings capture context and meaning
- **Measurable distance:** Cosine similarity, Euclidean distance

Types of Embeddings:

- **Text embeddings:** Words, sentences, documents
- **Image embeddings:** Visual feature representations
- **Multimodal embeddings:** Combined text and image representations
- **Domain-specific embeddings:** Fine-tuned for specific industries

Applications:

- Semantic search and retrieval
- Recommendation systems
- Content similarity matching
- Retrieval-Augmented Generation (RAG)

AWS Services:

- **Amazon Bedrock:** Text and multimodal embeddings via Titan Embeddings
- **Amazon SageMaker:** Custom embedding model training and hosting
- **Amazon OpenSearch:** Vector database for embedding storage and search

1.4 Vectors

Definition: Mathematical representations of data points in multi-dimensional space.

Vector Operations:

- **Similarity calculation:** Cosine similarity, dot product
- **Distance metrics:** Euclidean, Manhattan, Hamming distance

- **Vector arithmetic:** Addition, subtraction for semantic operations
- **Dimensionality:** Higher dimensions capture more nuanced relationships

Vector Databases:

- **Purpose:** Efficient storage and retrieval of high-dimensional vectors
- **Features:** Approximate nearest neighbor search, indexing algorithms
- **Scalability:** Handle millions to billions of vectors

AWS Vector Solutions:

- **Amazon OpenSearch Service:** Vector search capabilities
- **Amazon RDS for PostgreSQL:** pgvector extension
- **Amazon MemoryDB for Redis:** Vector similarity search
- **Third-party options:** Pinecone, Weaviate on AWS

1.5 Prompt Engineering

Definition: The practice of designing effective inputs to guide AI model outputs.

Core Principles:

- **Clarity:** Clear, specific instructions
- **Context:** Relevant background information
- **Examples:** Few-shot learning with demonstrations
- **Structure:** Organized format for complex tasks

Prompt Techniques:

Zero-shot Prompting

- Direct instruction without examples
- Example: "Translate the following text to French: [text]"

Few-shot Prompting

- Provide examples to guide model behavior
- Example: Include 2-3 translation examples before the target text

Chain-of-Thought (CoT)

- Guide models through step-by-step reasoning
- Example: "Let's solve this step by step..."

Role-based Prompting

- Assign specific roles or personas
- Example: "Act as a financial advisor and provide investment guidance..."

Template-based Prompting

- Structured formats for consistent outputs
- Example: JSON schemas, specific formatting requirements

Advanced Techniques:

- **Prompt chaining:** Breaking complex tasks into steps
- **Self-consistency:** Multiple reasoning paths
- **Tree of thoughts:** Exploring multiple solution branches
- **Retrieval-augmented prompting:** Incorporating external knowledge

AWS Context:

- **Amazon Bedrock:** Optimize prompts for different foundation models
- **Cost optimization:** Efficient prompts reduce token usage
- **Model selection:** Different models respond better to different prompt styles

1.6 Transformer-based LLMs

Definition: Large Language Models built on the transformer architecture for natural language understanding and generation.

Transformer Architecture:

- **Self-attention mechanism:** Models relationships between all tokens
- **Encoder-decoder structure:** Input processing and output generation
- **Positional encoding:** Understanding token sequence order
- **Multi-head attention:** Parallel attention computations

Key Components:

- **Attention layers:** Focus on relevant parts of input
- **Feed-forward networks:** Process attention outputs
- **Layer normalization:** Stabilize training
- **Residual connections:** Enable deep network training

Popular LLM Architectures:

- **GPT family:** Decoder-only, autoregressive generation
- **BERT family:** Encoder-only, bidirectional understanding
- **T5:** Encoder-decoder, text-to-text transfer
- **PaLM, LaMDA:** Large-scale conversational models

Scaling Laws:

- **Model size:** More parameters generally improve performance
- **Training data:** More diverse, high-quality data improves capabilities
- **Compute resources:** Larger models require more computational power
- **Emergent abilities:** New capabilities arise at certain scales

AWS Foundation Models:

- **Amazon Titan:** AWS-developed foundation models

- **Anthropic Claude:** Available through Bedrock
- **AI21 Labs Jurassic:** Multilingual capabilities
- **Cohere Command:** Specialized for business applications

1.7 Foundation Models

Definition: Large-scale models trained on broad datasets that can be adapted for various downstream tasks.

Characteristics:

- **General-purpose:** Trained on diverse data sources
- **Transfer learning:** Can be fine-tuned for specific tasks
- **Emergent abilities:** Capabilities not explicitly programmed
- **Scale:** Billions to trillions of parameters

Types of Foundation Models:

Language Models

- **Text generation:** GPT, PaLM, Claude
- **Text understanding:** BERT, RoBERTa, DeBERTa
- **Code generation:** Codex, CodeT5, StarCoder

Vision Models

- **Image classification:** Vision Transformer (ViT), ResNet
- **Image generation:** DALL-E, Midjourney, Stable Diffusion
- **Image understanding:** CLIP, BLIP

Multimodal Models

- **Vision-language:** CLIP, BLIP, Flamingo
- **Text-to-image:** DALL-E 2, Stable Diffusion
- **Image-to-text:** BLIP, LLaVA

Foundation Model Benefits:

- **Reduced training time:** Pre-trained on massive datasets
- **Lower resource requirements:** Fine-tuning vs. training from scratch
- **Consistent performance:** Proven capabilities across tasks
- **Rapid deployment:** Faster time-to-market

AWS Foundation Model Access:

- **Amazon Bedrock:** Managed access to multiple foundation models
- **Model variety:** Text, image, and multimodal models
- **API access:** Pay-per-use pricing model
- **Custom models:** Fine-tuning and customization options

1.8 Multi-modal Models

Definition: AI models that can process and generate content across multiple modalities (text, images, audio, video).

Modality Combinations:

- **Vision-Language:** Text and images
- **Audio-Language:** Speech and text
- **Video-Language:** Video content with text descriptions
- **Tri-modal:** Text, image, and audio together

Key Capabilities:

- **Cross-modal understanding:** Relate concepts across modalities
- **Multi-modal generation:** Create content in multiple formats
- **Modal translation:** Convert between different modalities
- **Unified representations:** Single model handling multiple inputs

Popular Multi-modal Models:

- **CLIP:** Connects text and images
- **DALL-E:** Text-to-image generation
- **Flamingo:** Few-shot learning across modalities
- **GPT-4V:** Vision-enabled language model
- **Whisper:** Speech recognition and translation

Applications:

- **Content creation:** Generate images from text descriptions
- **Accessibility:** Alt-text generation, audio descriptions
- **Search:** Find images using natural language queries
- **Education:** Interactive learning with multiple media types

AWS Multi-modal Solutions:

- **Amazon Bedrock:** Multi-modal foundation models
- **Amazon Rekognition:** Image and video analysis with text output
- **Amazon Transcribe:** Speech-to-text conversion
- **Amazon Polly:** Text-to-speech synthesis

1.9 Diffusion Models

Definition: Generative models that create data by learning to reverse a noise-adding process.

How Diffusion Models Work:

Forward Process (Noise Addition)

1. Start with real data (e.g., images)
2. Gradually add Gaussian noise over multiple steps
3. Eventually reach pure noise

Reverse Process (Denoising)

1. Start with random noise
2. Predict and remove noise step by step
3. Generate new data samples

Key Advantages:

- **High-quality outputs:** Superior image generation quality
- **Training stability:** More stable than GANs
- **Controllable generation:** Can guide generation process
- **Diverse outputs:** Generate varied samples from same input

Popular Diffusion Models:

- **DDPM:** Denoising Diffusion Probabilistic Models
- **Stable Diffusion:** Open-source text-to-image model
- **DALL-E 2:** OpenAI's image generation model
- **Imagen:** Google's text-to-image model

Applications:

- **Image generation:** Create images from text prompts
- **Image editing:** Inpainting, outpainting, style transfer
- **Super-resolution:** Enhance image quality and resolution
- **3D generation:** Create 3D models and scenes

AWS Integration:

- **Amazon Bedrock:** Access to Stable Diffusion models
- **Amazon SageMaker:** Host custom diffusion models
- **Amazon EC2:** GPU instances for diffusion model training

2. Potential Use Cases for Generative AI

2.1 Content Generation

Image Generation

Applications:

- **Marketing materials:** Advertisements, social media content, product images
- **Art and design:** Creative artwork, logos, illustrations
- **Product visualization:** Mockups, prototypes, concept designs
- **Synthetic data:** Training data for computer vision models

AWS Implementation:

- **Amazon Bedrock:** Stable Diffusion, Titan Image Generator
- **Custom models:** Fine-tuned diffusion models on SageMaker

- **Integration:** Lambda functions for automated image generation

Video Generation

Applications:

- **Content creation:** Short videos, animations, promotional content
- **Education:** Explainer videos, training materials
- **Entertainment:** Synthetic actors, scene generation
- **Personalization:** Customized video content for users

Considerations:

- **Computational requirements:** High GPU and memory needs
- **Quality vs. speed:** Trade-offs in generation time
- **Copyright and ethics:** Ensuring appropriate content generation

Audio Generation

Applications:

- **Music composition:** Background music, jingles, soundtracks
- **Voice synthesis:** Audiobooks, podcasts, voice assistants
- **Sound effects:** Game audio, media production
- **Accessibility:** Text-to-speech for visually impaired users

AWS Services:

- **Amazon Polly:** Neural text-to-speech
- **Amazon Transcribe:** Speech recognition and processing
- **Custom models:** Audio generation models on SageMaker

2.2 Text Processing and Generation

Summarization

Types:

- **Extractive:** Select key sentences from original text
- **Abstractive:** Generate new summary text
- **Multi-document:** Summarize across multiple sources
- **Query-focused:** Summarize based on specific questions

Applications:

- **Document processing:** Legal documents, research papers, reports
- **News aggregation:** Daily briefings, topic summaries
- **Customer feedback:** Review summaries, survey insights
- **Meeting notes:** Automated meeting summaries

AWS Implementation:

- **Amazon Bedrock:** Claude, Jurassic models for summarization
- **Amazon Comprehend:** Extract key phrases and sentiment
- **Amazon Textract:** Extract text from documents for summarization

Chatbots and Virtual Assistants

Capabilities:

- **Natural conversation:** Human-like dialogue
- **Context awareness:** Maintain conversation history
- **Task completion:** Booking, ordering, information retrieval
- **Multilingual support:** Multiple language conversations

Implementation Patterns:

- **Retrieval-Augmented Generation (RAG):** Combine knowledge bases with LLMs
- **Function calling:** Integrate with external APIs and services
- **Conversation management:** Handle complex, multi-turn dialogues
- **Personalization:** Adapt responses to user preferences

AWS Solutions:

- **Amazon Bedrock:** Foundation models for conversational AI
- **Amazon Lex:** Purpose-built chatbot service
- **Amazon Connect:** Contact center integration
- **Amazon Kendra:** Enterprise search integration

Translation

Types:

- **Machine translation:** Automated language translation
- **Neural translation:** Context-aware, fluent translations
- **Document translation:** Preserve formatting and structure
- **Real-time translation:** Live conversation translation

Applications:

- **Global business:** International communication, content localization
- **Customer support:** Multilingual customer service
- **Content publishing:** Website and document translation
- **Travel and tourism:** Real-time travel assistance

AWS Services:

- **Amazon Translate:** Neural machine translation service
- **Amazon Bedrock:** Multilingual foundation models
- **Integration:** Combine with text extraction and processing services

2.3 Code and Development

Code Generation

Capabilities:

- **Code completion:** Autocomplete code snippets
- **Function generation:** Create functions from descriptions
- **Bug fixing:** Identify and suggest fixes for code issues
- **Code explanation:** Document and explain existing code

Use Cases:

- **Developer productivity:** Faster coding and debugging
- **Learning assistance:** Help developers learn new languages
- **Code review:** Automated code quality assessment
- **Documentation:** Generate code comments and documentation

AWS Integration:

- **Amazon CodeGuru:** AI-powered code reviews
- **Amazon SageMaker:** Train custom code generation models
- **Amazon Bedrock:** Code-capable foundation models

2.4 Business Applications

Customer Service Agents

Features:

- **24/7 availability:** Round-the-clock customer support
- **Instant responses:** Immediate answers to common questions
- **Escalation handling:** Route complex issues to human agents
- **Personalization:** Customized responses based on customer history

Implementation:

- **Knowledge base integration:** Access to company information
- **CRM integration:** Customer history and preferences
- **Multi-channel support:** Web, mobile, voice, chat
- **Analytics:** Performance monitoring and improvement

AWS Solutions:

- **Amazon Connect:** Cloud contact center platform
- **Amazon Lex:** Conversational interfaces
- **Amazon Bedrock:** Advanced language understanding
- **Amazon Kendra:** Enterprise search for knowledge retrieval

Search Enhancement

Capabilities:

- **Semantic search:** Understand search intent, not just keywords
- **Conversational search:** Natural language queries
- **Personalized results:** Tailored to user preferences and history
- **Multi-modal search:** Search across text, images, and other media

Implementation Patterns:

- **Vector search:** Embedding-based similarity matching
- **Hybrid search:** Combine keyword and semantic search
- **Query expansion:** Enhance queries with related terms
- **Result ranking:** AI-powered relevance scoring

AWS Services:

- **Amazon Kendra:** Intelligent enterprise search
- **Amazon OpenSearch:** Full-text and vector search
- **Amazon Bedrock:** Query understanding and enhancement

Recommendation Engines

Types:

- **Content-based:** Recommend based on item features
- **Collaborative filtering:** Use user behavior patterns
- **Hybrid approaches:** Combine multiple recommendation methods
- **Real-time recommendations:** Dynamic, context-aware suggestions

Applications:

- **E-commerce:** Product recommendations
- **Content platforms:** Article, video, music recommendations
- **Learning platforms:** Course and content suggestions
- **News and media:** Personalized content feeds

AWS Implementation:

- **Amazon Personalize:** Managed recommendation service
- **Amazon Bedrock:** Content understanding and generation
- **Amazon SageMaker:** Custom recommendation models

3. Foundation Model Lifecycle

3.1 Data Selection

Purpose: Choose appropriate training data that will enable the model to learn desired capabilities.

Data Types:

- **Text corpora:** Books, articles, web pages, documentation
- **Code repositories:** GitHub, programming tutorials, technical documentation

- **Conversational data:** Dialogues, Q&A pairs, chat logs
- **Multimodal data:** Image-text pairs, video transcripts, audio-text alignments

Data Quality Considerations:

- **Relevance:** Data should match intended use cases
- **Diversity:** Broad representation of topics, languages, styles
- **Quality:** High-quality, well-written, factually accurate content
- **Freshness:** Recent data for current knowledge and trends
- **Balance:** Avoid overrepresentation of specific domains or viewpoints

Data Challenges:

- **Copyright and licensing:** Ensure legal use of training data
- **Privacy:** Remove or anonymize personal information
- **Bias:** Address potential biases in training data
- **Scale:** Managing petabytes of training data
- **Preprocessing:** Cleaning, filtering, and formatting data

AWS Data Services:

- **Amazon S3:** Massive data storage and management
- **AWS Glue:** Data cataloging and ETL for preprocessing
- **Amazon EMR:** Big data processing for data preparation
- **AWS Data Exchange:** Access to third-party datasets

3.2 Model Selection

Purpose: Choose the appropriate model architecture and size for the intended application.

Architecture Considerations:

- **Model type:** Encoder-only, decoder-only, encoder-decoder
- **Size trade-offs:** Parameter count vs. computational requirements
- **Capabilities needed:** Text generation, understanding, reasoning, multimodal
- **Performance requirements:** Latency, throughput, accuracy

Popular Foundation Model Families:

Text Models

- **GPT family:** Excellent for text generation and completion
- **BERT family:** Strong text understanding and classification
- **T5:** Versatile text-to-text transfer learning
- **PaLM:** Large-scale reasoning and mathematical capabilities

Multimodal Models

- **CLIP:** Vision-language understanding
- **DALL-E:** Text-to-image generation
- **Flamingo:** Few-shot multimodal learning

Selection Criteria:

- **Task alignment:** Model capabilities match use case requirements
- **Resource constraints:** Available compute, memory, and storage
- **Latency requirements:** Real-time vs. batch processing needs
- **Cost considerations:** Training, inference, and operational costs
- **Customization needs:** Ability to fine-tune or adapt the model

AWS Model Options:

- **Amazon Bedrock:** Multiple foundation models to choose from
- **Model comparison:** Evaluate different models for specific tasks
- **Cost analysis:** Compare pricing across different model options

3.3 Pre-training

Purpose: Train foundation models on large-scale datasets to learn general language and reasoning capabilities.

Pre-training Objectives:

- **Language modeling:** Predict next tokens in sequences
- **Masked language modeling:** Predict masked tokens in sentences
- **Contrastive learning:** Learn representations by comparing examples
- **Multi-task learning:** Train on multiple objectives simultaneously

Training Process:

1. **Data preprocessing:** Tokenization, chunking, formatting
2. **Model initialization:** Set up architecture and parameters
3. **Distributed training:** Use multiple GPUs/machines for scale
4. **Loss optimization:** Minimize prediction errors across the dataset
5. **Checkpointing:** Save model state at regular intervals
6. **Monitoring:** Track training metrics and model behavior

Infrastructure Requirements:

- **Compute resources:** Thousands of high-end GPUs
- **Storage:** Petabytes of training data and model checkpoints
- **Networking:** High-bandwidth connections between compute nodes
- **Time:** Weeks to months for large model training

AWS Pre-training Infrastructure:

- **Amazon SageMaker:** Managed training with distributed computing
- **Amazon EC2 P4/P5 instances:** High-performance GPU instances
- **AWS ParallelCluster:** HPC clusters for large-scale training
- **Amazon FSx:** High-performance file systems for training data

3.4 Fine-tuning

Purpose: Adapt pre-trained foundation models for specific tasks, domains, or use cases.

Types of Fine-tuning:

Supervised Fine-tuning (SFT)

- **Task-specific training:** Train on labeled examples for specific tasks
- **Domain adaptation:** Adapt to specific industries or domains
- **Format training:** Teach models to follow specific output formats
- **Instruction following:** Train models to follow human instructions

Reinforcement Learning from Human Feedback (RLHF)

- **Human preference training:** Align model outputs with human preferences
- **Safety improvements:** Reduce harmful or inappropriate outputs
- **Helpfulness training:** Improve the quality and usefulness of responses
- **Reward modeling:** Learn to predict human preferences

Parameter-Efficient Fine-tuning

- **LoRA (Low-Rank Adaptation):** Fine-tune small parameter subsets
- **Adapters:** Add small modules without changing base model
- **Prompt tuning:** Learn optimal prompts for specific tasks
- **In-context learning:** Use examples within the prompt

Fine-tuning Process:

1. **Data preparation:** Create high-quality task-specific datasets
2. **Model setup:** Load pre-trained model and configure for fine-tuning
3. **Training configuration:** Set learning rates, batch sizes, epochs
4. **Training execution:** Run fine-tuning on prepared dataset
5. **Validation:** Monitor performance on held-out validation set
6. **Hyperparameter tuning:** Optimize training parameters

AWS Fine-tuning Services:

- **Amazon Bedrock:** Custom model fine-tuning capabilities
- **Amazon SageMaker:** Full fine-tuning infrastructure and tools
- **SageMaker JumpStart:** Pre-configured fine-tuning workflows

3.5 Evaluation

Purpose: Assess model performance, capabilities, and limitations before deployment.

Evaluation Dimensions:

Task Performance

- **Accuracy metrics:** Task-specific performance measures
- **Benchmark datasets:** Standardized evaluation sets

- **Comparative analysis:** Performance vs. other models
- **Error analysis:** Understanding failure modes and limitations

Safety and Alignment

- **Bias detection:** Identify unfair or discriminatory outputs
- **Toxicity testing:** Detect harmful or inappropriate content
- **Factual accuracy:** Verify correctness of generated information
- **Alignment assessment:** Measure adherence to human values

Robustness

- **Adversarial testing:** Resistance to malicious inputs
- **Out-of-distribution:** Performance on unseen data types
- **Consistency:** Reliability across similar inputs
- **Edge case handling:** Behavior in unusual situations

Evaluation Methods:

- **Automated metrics:** BLEU, ROUGE, perplexity, F1 scores
- **Human evaluation:** Expert assessments and user studies
- **A/B testing:** Comparative evaluation with baseline models
- **Red team exercises:** Systematic attempts to find failure modes

AWS Evaluation Tools:

- **Amazon SageMaker Clarify:** Model bias and explainability analysis
- **Amazon Bedrock:** Built-in model evaluation capabilities
- **Custom evaluation:** SageMaker notebooks for comprehensive testing

3.6 Deployment

Purpose: Make trained models available for production use in applications and services.

Deployment Strategies:

Managed API Deployment

- **Serverless inference:** On-demand model serving
- **Auto-scaling:** Dynamic resource allocation based on demand
- **Multiple regions:** Deploy across geographical locations
- **Version management:** Support multiple model versions

Self-managed Deployment

- **Container deployment:** Docker containers on Kubernetes or ECS
- **Edge deployment:** Deploy models on edge devices
- **Batch inference:** Process large datasets offline
- **Hybrid deployment:** Combine cloud and on-premises infrastructure

Deployment Considerations:

- **Latency requirements:** Real-time vs. batch processing needs
- **Throughput needs:** Expected request volume and concurrency
- **Cost optimization:** Balance performance and operational costs
- **Security:** Protect model IP and ensure data privacy
- **Compliance:** Meet regulatory and industry requirements

AWS Deployment Services:

- **Amazon Bedrock:** Fully managed foundation model deployment
- **Amazon SageMaker Endpoints:** Custom model hosting
- **AWS Lambda:** Lightweight model serving
- **Amazon ECS/EKS:** Container-based deployments

3.7 Feedback and Continuous Improvement

Purpose: Collect user feedback and performance data to continuously improve model performance and alignment.

Feedback Collection:

- **User ratings:** Thumbs up/down, star ratings, quality scores
- **Usage analytics:** Track how models are used and where they fail
- **A/B testing:** Compare different model versions or configurations
- **Expert review:** Domain expert evaluation of model outputs

Feedback Integration:

- **Reinforcement learning:** Use feedback to improve model behavior
- **Fine-tuning updates:** Periodic retraining with new feedback data
- **Prompt optimization:** Improve prompts based on performance data
- **Configuration updates:** Adjust model parameters and settings

Monitoring and Observability:

- **Performance metrics:** Track accuracy, latency, and throughput
- **Data drift detection:** Monitor changes in input data distribution
- **Model degradation:** Identify declining performance over time
- **Cost monitoring:** Track inference costs and resource usage

Continuous Improvement Process:

1. **Data collection:** Gather feedback and performance metrics
2. **Analysis:** Identify improvement opportunities and issues
3. **Experimentation:** Test potential improvements
4. **Validation:** Ensure improvements don't introduce new issues
5. **Deployment:** Roll out improvements to production
6. **Monitoring:** Track impact of changes

AWS Feedback and Monitoring:

- **Amazon CloudWatch:** Monitor model performance and costs
 - **Amazon SageMaker Model Monitor:** Detect data drift and model degradation
 - **Custom feedback systems:** Build feedback collection into applications
-

Key Study Points for AWS AIF-C01

Conceptual Understanding

1. **Token economics:** Understand how tokenization affects costs and performance
2. **Embedding applications:** Know when and how to use embeddings for various tasks
3. **Prompt engineering strategies:** Master different prompting techniques
4. **Model selection criteria:** Choose appropriate models for different use cases
5. **Deployment trade-offs:** Understand managed vs. self-hosted options

AWS Service Mapping

- **Amazon Bedrock:** Primary service for foundation model access
- **Amazon SageMaker:** Custom model development and deployment
- **Supporting services:** S3, Lambda, OpenSearch, Kendra, etc.
- **Cost considerations:** Understand pricing models for different approaches

Practical Applications

- **Use case identification:** Match generative AI capabilities to business needs
- **Implementation patterns:** RAG, fine-tuning, prompt engineering
- **Performance optimization:** Balance quality, speed, and cost
- **Ethical considerations:** Bias, safety, and responsible AI practices

Foundation Model Lifecycle

- **End-to-end understanding:** From data selection to continuous improvement
- **AWS integration:** How AWS services support each lifecycle stage
- **Best practices:** Industry standards and AWS recommendations
- **Evaluation methods:** Comprehensive model assessment approaches

This study guide provides the foundational knowledge needed for the generative AI section of the AWS AI Practitioner certification, focusing on both conceptual understanding and practical AWS implementation.