

Study Guide

Quiz

Instructions: Answer each question in 2-3 sentences.

1. What is the primary difference between Machine Learning (ML) and Generative AI?
2. Explain the concept of "unstructured data" and provide one real-world example.
3. Describe the main purpose of the "Model Garden" within Vertex AI.
4. How does "Retrieval-Augmented Generation (RAG)" help improve the factual accuracy of an LLM's output?
5. What is a "reasoning loop" in the context of a generative AI agent?
6. List two considerations related to "needs" that a company should evaluate before starting a Gen AI project.
7. What is the function of "Gems" within the Gemini application?
8. Briefly explain the "Zero-shot" prompting technique.
9. Why is "human-in-the-loop (HITL)" important for responsible AI, particularly in generative AI solutions?
10. What specific Google model is designed for generating video content from text descriptions or still images?

Answer Key

1. **What is the primary difference between Machine Learning (ML) and Generative AI?** Machine Learning (ML) is a subfield of AI where machines learn from data to perform specific tasks, which can include prediction or classification. Generative AI is a specific application of ML that focuses on creating *new* content, such as text, images, or code, rather than just analyzing existing data.
2. **Explain the concept of "unstructured data" and provide one real-world example.** Unstructured data is information that lacks a predefined format or organization, making it challenging to search and analyze using traditional database methods. Examples include text documents, emails, images, audio files, or social media posts, which require sophisticated analysis techniques.
3. **Describe the main purpose of the "Model Garden" within Vertex AI.** The Model Garden in Vertex AI serves as a catalog where users can find and select from existing Google, third-party, or open-source machine learning models. Its purpose is to provide a starting point for building and deploying AI solutions, allowing for customization and integration.
4. **How does "Retrieval-Augmented Generation (RAG)" help improve the factual accuracy of an LLM's output?** RAG improves factual accuracy by allowing an LLM to retrieve relevant information from external, verifiable sources before generating a response. This retrieved information is then incorporated into the prompt, grounding the model's output in facts and reducing the likelihood of hallucinations or inaccuracies.
5. **What is a "reasoning loop" in the context of a generative AI agent?** A reasoning loop is an iterative process within a generative AI agent where it observes its environment, interprets information, reasons about the next steps, and then acts upon it. This loop often involves prompt engineering to guide the agent's decision-making and achieve its goals.
6. **List two considerations related to "needs" that a company should evaluate before starting a Gen AI project.** Before starting a Gen AI project, a company should consider the *scale* of the solution, determining how many users will interact with it. They should also assess the required *customization*, understanding how specialized the AI needs to be for their specific use case.

7. **What is the function of "Gems" within the Gemini application?** Gems are personalized AI assistants within the Gemini application that provide responses tailored to specific instructions. They streamline workflows by acting as templates, prompts, and guided interactions, offering personalized assistance for various tasks.
8. **Briefly explain the "Zero-shot" prompting technique.** Zero-shot prompting is a technique where the generative AI model is asked to complete a task without any prior examples or demonstrations. The model relies solely on its pre-trained knowledge to understand the instruction and generate a relevant output.
9. **Why is "human-in-the-loop (HITL)" important for responsible AI, particularly in generative AI solutions?** Human-in-the-loop (HITL) is crucial for responsible AI because it ensures human oversight and intervention in AI systems. It helps to maintain accuracy, mitigate bias, and ensure the safety and ethical use of generative AI by involving humans in data annotation, prompt review, and output validation.
10. **What specific Google model is designed for generating video content from text descriptions or still images?** The Google model specifically designed for generating video content from text descriptions or still images is Veo. It brings ideas to life in various cinematic and visual styles.

Essay Format Questions

1. Discuss how Google's "AI-first approach" and commitment to innovation provide competitive advantages in the generative AI space. Include specific examples of Google Cloud's offerings that demonstrate these advantages.
2. Analyze the critical factors to consider when choosing a generative AI model for a business solution. Explain how each factor (modality, context window, security, availability, reliability, cost, performance, fine-tuning) impacts strategic decision-making.
3. Elaborate on the importance of Responsible AI and transparency in the context of generative AI solutions. How do Google Cloud's principles and tools, including privacy considerations like data anonymization, contribute to building trust and ensuring ethical AI adoption?
4. Compare and contrast prompt engineering and fine-tuning as techniques to improve generative AI model output. Discuss the scenarios where each technique would be most appropriate, considering factors like speed, cost, customization level, and data privacy.
5. Outline a comprehensive strategy for developing and implementing a successful generative AI solution within an organization. Include the key pre-project considerations, elements of the strategic framework, and methods for measuring the value and ROI of the AI initiative.

Glossary of Key Terms

AI (Artificial Intelligence): Building machines that can perform tasks that typically require human intelligence, such as learning, problem-solving, and decision-making.

AI-first approach: Google's strategic philosophy emphasizing the integration of AI capabilities across all its products and services.

Agent (Generative AI): An application that tries to achieve a goal by observing the world and acting upon it using the tools it has at its disposal. Can be deterministic, generative, or hybrid.

Agent Assist: A component of Google's Customer Engagement Suite that supports live human contact center agents.

Agentspace: A Google Cloud offering that allows integration of customized search and conversation agents to access and understand data from various internal sources for internal websites or dashboards.

AutoML: A feature within Vertex AI Model Builder that allows users to create and train models with minimal technical knowledge and effort.

Availability and reliability: Factors affecting whether a generative AI solution can be consistently used in a production environment.

BigQuery: Google Cloud's fully managed, serverless data warehouse optimized for scalable data analysis.

Bias (in AI): The tendency of an AI model to produce outputs that reflect imbalances or prejudices present in its training data.

Chain-of-Thought (CoT) Prompting: A prompt engineering technique that guides a Large Language Model (LLM) through a problem-solving process by providing examples with intermediate reasoning steps, improving accuracy and transparency.

Cloud Functions: A serverless execution environment within Google Cloud for running event-driven code.

Cloud Storage: A Google Cloud service for storing data.

Completeness (Data Quality): Ensuring that all necessary data is present and accounted for.

Consistency (Data Quality): Ensuring data is uniform and reliable across all systems.

Contact Center as a Service (CCaaS): An enterprise-grade contact center solution native to the cloud, forming part of Google's Customer Engagement Suite.

Context window: How much information a generative AI model can process at once, crucial for tasks like document analysis.

Conversational Agents: Effective chatbots designed to interact with customers as part of Google's Customer Engagement Suite.

Conversational Insights: A component of Google's Customer Engagement Suite that provides insights into customer communications.

Cost: A factor impacting budget considerations for generative AI solutions.

Customization options: The degree to which a generative AI model can be adapted to specific needs, including fine-tuning.

Customer Engagement Suite: A set of Google tools to support companies in engaging with customers effectively, built on top of CCaaS.

Data accessibility: Ensuring data for model training is readily available, usable, and in the proper format.

Data ingestion and preparation: The process of collecting, cleaning, and transforming raw data into a usable format for analysis or model training.

Data Quality: Ensuring data is accurate, complete, consistent, and relevant.

Data stores: Tooling components that provide access to information for AI agents.

Data (information): Information that can come in many forms: numbers, dates, text descriptions, images, or sounds.

Data dependency (Foundation Model Limitation): The reliance of a model's performance on the quality and completeness of its training data.

Deep learning: A subset of Machine Learning (ML) that uses artificial neural networks with many layers to extract complex patterns from data.

Deterministic Agents: AI agents that are built with predefined paths and actions.

Diffusion models: A type of generative AI model capable of generating high-quality images.

Drift monitoring (Vertex AI Model Monitoring): The process of watching for changes in a model's accuracy over time to ensure consistent performance.

Edge AI: Running AI solutions on infrastructure (devices or servers) closer to where the action is happening.

Ethical AI: Ensuring AI applications don't cause harm and are used in an ethical manner.

Extensions: Tooling components that connect to external services (via APIs) for an AI agent.

Fairness (Responsible AI): Ensuring AI systems treat all individuals and groups equitably and do not perpetuate or amplify biases.

Few-shot prompting: Providing a generative AI model with multiple examples to learn from before generating an output.

Fine-tuning: The process of retraining a pre-trained foundation model on specific, proprietary data for deep customization and higher accuracy on specialized tasks.

Foundation models: Powerful Machine Learning (ML) models trained on massive amounts of unlabeled data, allowing them to develop a broad understanding of the world.

Functions: Tooling components that define specific actions or tasks for an AI agent.

Gemma: A family of lightweight, open AI models from Google.

Gemini: Google's multimodal generative AI chatbot, integrated across Workspace apps and Google Cloud, providing assistance with writing, planning, learning, and more.

Gemini for Google Cloud: Google's AI assistant for Google Cloud, helping with code, application management, data analysis, and security.

Gemini for Google Workspace: Integrates generative AI into familiar Workspace apps (e.g., Gmail, Slides, Meet).

Gems (in Gemini): Personalized AI assistants within Gemini that provide tailored responses and streamline workflows.

Generative AI: An application of Machine Learning (ML) that focuses on creating new content, such as text, images, or code.

Generative Agents: AI agents that are defined with natural language using LLMs to give a real conversational feel to a chatbot.

Google AI Studio: A free platform intended for quick AI prototyping and experimentation with the Gemini API.

Google Cloud: Google's unified cloud computing platform providing infrastructure and services, including AI offerings.

Grounding: Connecting the AI's output to verifiable sources of information to reduce hallucinations and improve factual accuracy.

GPUs (Graphics Processing Units): Specialized hardware used to accelerate AI/ML model training and inference.

Hallucinations: Instances where a generative AI model produces inaccurate or fabricated information, despite sounding plausible.

Human-in-the-Loop (HITL): Google-recommended practice for continuous human oversight and intervention in AI systems to ensure accuracy, reduce bias, and maintain ethical standards.

Hypercomputer: Google Cloud's AI-optimized infrastructure, combining TPUs, GPUs, and other resources.

Identity and Access Management (IAM): A Google Cloud security tool used to manage permissions and access to resources.

Imagen: A text-to-image diffusion model from Google that generates high-quality images from textual descriptions.

Infrastructure: The foundation upon which AI rests, providing core computing resources (physical hardware like servers, GPUs, TPUs, and essential software).

Knowledge cutoff: The point in time after which a foundation model has not been trained on new information, leading to an inability to answer questions about more recent events.

Labeled data: Data that is organized with predefined tags or categories (labels) associated with each input, used for supervised learning.

Large Language Models (LLMs): A type of foundation model that is designed to understand and generate human language.

Latency: The response time of an AI system, a key consideration for user interaction and real-time applications.

Lite Runtime (LiteRT): A Google tool to help developers deploy AI models on edge devices.

Machine learning (ML): A subfield of AI where machines learn from data to perform specific tasks.

Metaprompting: Using prompting to guide the AI model to generate, modify, or interpret other prompts.

Modality: The type of data a generative AI model can handle (e.g., numbers, dates, text, images, sounds, code).

Model (AI system component): The "brains" of the AI system, consisting of various algorithms that learn patterns from data and can make predictions or generate new content.

Model Builder (Vertex AI): A tool within Vertex AI that allows users to train and use their own custom models.

Model deployment: The process of making a trained ML model available for use.

Model management: The process of managing and maintaining ML models over time.

Model training: The process of creating an ML model using data.

Multimodal: The ability of generative AI applications to process and generate different types of data simultaneously (e.g., text, images, audio, code).

Natural Language Processing (NLP): A field of AI that focuses on enabling computers to understand, interpret, and generate human language.

NotebookLM: A Google tool that acts as a research assistant, summarizing key points, answering questions, and generating ideas from uploaded source material.

One-shot prompting: Providing the generative AI model with one example to learn from before completing a task.

Performance: Determines the speed and quality of a generative AI model's output.

Platform (AI initiatives): The foundation for building and scaling AI initiatives, such as Vertex AI.

Plugins: Tooling components that add new skills and integrations to an AI agent.

Privacy risks: Potential threats to sensitive data in AI systems.

Prompt chaining: Continuing conversations within the same chatbot to maintain context and build upon previous interactions.

Prompt engineering: The process of designing and refining inputs (prompts) to generative AI models to elicit desired outputs.

Reason and Act (ReAct): A prompt engineering technique that allows an LLM to interleave reasoning steps with actions for more reliable answers.

Reasoning loop (Agent component): An iterative process where an AI agent observes, interprets, reasons, and acts, often using prompt engineering.

Reinforcement learning: An ML approach where the model learns through interaction and feedback to maximize rewards and minimize penalties.

Relevance (Data Quality): Ensuring that the data is pertinent and applicable to the task at hand.

Responsible AI: Principles and practices focused on developing and deploying AI systems in a fair, transparent, and accountable manner, considering ethical implications.

Retrieval-Augmented Generation (RAG): A technique that involves retrieving relevant information from external sources and augmenting the LLM's prompt with this context before generating a response, improving factual accuracy.

Role prompting: Assigning a persona to the model to influence its style, tone, and focus.

Sampling parameters: Settings that allow users to influence an AI model's behavior and customize its output (e.g., token count, temperature, Top-p).

Secure AI: Protecting AI applications from harm, including malicious attacks and misuse.

Secure AI Framework (SAIF): Google's framework to help organizations manage AI/ML model risks and ensure security throughout the ML lifecycle.

Secure-by-design infrastructure: A principle of building security directly into the foundation of AI systems and infrastructure.

Security Command Center: A Google Cloud security tool used for centralized security management and risk detection.

Structured data: Data that is organized and easy to search, often stored in relational databases.

Supervised learning: An ML approach that trains models on labeled data to predict outputs for new inputs.

Temperature (Sampling Parameter): Controls the creativity versus focus of an AI model's responses. Higher temperature leads to more creative/random outputs.

Token count: Represents meaningful chunks of text in AI models, used to control output length.

Tooling (Generative AI): Functionalities that allow an AI agent to interact with its environment, such as accessing and processing data or interacting with software or hardware.

Transparency (Responsible AI): Providing clear information about data use and AI system behavior to foster trust.

TPUs (Tensor Processing Units): Google's custom-designed chips optimized for machine learning workloads, excelling at parallel processing.

Unlabeled data: Data that lacks predefined tags or categories, used for unsupervised learning to find natural groupings and patterns.

Unstructured data: Data that lacks a predefined structure and requires sophisticated analysis techniques.

Veo: A Google model that generates video content based on text descriptions or still images.

Vertex AI: Google Cloud's unified machine learning (ML) platform designed to streamline the entire ML workflow, providing infrastructure, tools, and pre-trained models.

Vertex AI Search: Google Cloud service for search and recommendation solutions for businesses, often used in conjunction with RAG.

Vertex AI Studio: An enterprise-grade platform within Vertex AI for building and deploying production-ready AI applications.

Versioning (Models): Managing different iterations of AI models over time.

Zero-shot prompting: Asking the model to complete a task with no prior examples.