

New NCA-GENL Test Braindumps & NCA-GENL Valid Exam Labs



P.S. Free & New NCA-GENL dumps are available on Google Drive shared by ValidBraindumps: <https://drive.google.com/open?id=1QsegzHWUJYK9M0kU9lQR8qtKEvS7k4QV>

Before we start develop a new NCA-GENL real exam, we will prepare a lot of materials. After all, we must ensure that all the questions and answers of the NCA-GENL exam materials are completely correct. First of all, we have collected all relevant reference books. Most of the NCA-GENL Practice Guide is written by the famous experts in the field. And we also add the latest knowledge points into the content of the NCA-GENL learning questions, so that they are always being up to date.

As far as we know, in the advanced development of electronic technology, lifelong learning has become more accessible, which means everyone has opportunities to achieve their own value and life dream though some ways such as the NCA-GENL certification. With over a decade's endeavor, our NCA-GENL practice materials successfully become the most reliable products in the industry. There is a great deal of advantages of our NCA-GENL exam questions you can spare some time to get to know.

>> New NCA-GENL Test Braindumps <<

Excel in Your NVIDIA NCA-GENL Exam with ValidBraindumps: The Quick Solution for Success

We aim to leave no misgivings to our customers so that they are able to devote themselves fully to their studies on NCA-GENL guide materials and they will find no distraction from us. I suggest that you strike while the iron is hot since time waits for no one. With our NCA-GENL Exam Questions, you will be bound to pass the exam with the least time and effort for its high quality. With our NCA-GENL study guide for 20 to 30 hours, you will be ready to take part in the exam and pass it with ease.

NVIDIA NCA-GENL Exam Syllabus Topics:

Topic	Details
Topic 1	<ul style="list-style-type: none">Experimentation: Explores running and evaluating trials to test model behavior, compare approaches, and validate generative AI solutions.
Topic 2	<ul style="list-style-type: none">Data preprocessing and feature engineering: Covers preparing raw data through cleaning, transformation, and feature selection to make it suitable for model training.
Topic 3	<ul style="list-style-type: none">Prompt engineering: Focuses on techniques for designing and refining input prompts to effectively guide LLM outputs toward desired results.
Topic 4	<ul style="list-style-type: none">Experiment design: Focuses on structuring controlled tests and workflows to systematically evaluate LLM performance and outcomes.

Topic 5	<ul style="list-style-type: none"> • Fundamentals of machine learning and neural networks: Covers the core concepts of how machine learning models learn from data, including the structure and function of neural networks that underpin large language models.
Topic 6	<ul style="list-style-type: none"> • Alignment: Addresses methods for ensuring LLM behavior is safe, accurate, and consistent with human intentions and values.
Topic 7	<ul style="list-style-type: none"> • Software development: Covers the programming practices and coding skills required to build, maintain, and deploy generative AI applications.
Topic 8	<ul style="list-style-type: none"> • Python libraries for LLMs: Covers key Python frameworks and tools — such as LangChain, Hugging Face, and similar libraries — used to build and interact with LLMs.

NVIDIA Generative AI LLMs Sample Questions (Q48-Q53):

NEW QUESTION # 48

"Hallucinations" is a term coined to describe when LLM models produce what?

- A. Outputs are only similar to the input data.
- B. Grammatically incorrect or broken outputs.
- C. Images from a prompt description.
- D. Correct sounding results that are wrong.

Answer: D

Explanation:

In the context of LLMs, "hallucinations" refer to outputs that sound plausible and correct but are factually incorrect or fabricated, as emphasized in NVIDIA's Generative AI and LLMs course. This occurs when models generate responses based on patterns in training data without grounding in factual knowledge, leading to misleading or invented information. Option A is incorrect, as hallucinations are not about similarity to input data but about factual inaccuracies. Option B is wrong, as hallucinations typically refer to text, not image generation. Option D is inaccurate, as hallucinations are grammatically coherent but factually wrong. The course states: "Hallucinations in LLMs occur when models produce correct-sounding but factually incorrect outputs, posing challenges for ensuring trustworthy AI." References: NVIDIA Building Transformer-Based Natural Language Processing Applications course; NVIDIA Introduction to Transformer-Based Natural Language Processing.

NEW QUESTION # 49

Which calculation is most commonly used to measure the semantic closeness of two text passages?

- A. Euclidean distance
- B. Cosine similarity
- C. Hamming distance
- D. Jaccard similarity

Answer: B

Explanation:

Cosine similarity is the most commonly used metric to measure the semantic closeness of two text passages in NLP. It calculates the cosine of the angle between two vectors (e.g., word embeddings or sentence embeddings) in a high-dimensional space, focusing on the direction rather than magnitude, which makes it robust for comparing semantic similarity. NVIDIA's documentation on NLP tasks, particularly in NeMo and embedding models, highlights cosine similarity as the standard metric for tasks like semantic search or text similarity, often using embeddings from models like BERT or Sentence-BERT. Option A (Hamming distance) is for binary data, not text embeddings. Option B (Jaccard similarity) is for set-based comparisons, not semantic content. Option D (Euclidean distance) is less common for text due to its sensitivity to vector magnitude.

References:

NVIDIA NeMo Documentation: <https://docs.nvidia.com/deeplearning/nemo/user-guide/docs/en/stable/nlp/intro.html>

NEW QUESTION # 50

Which of the following is a parameter-efficient fine-tuning approach that one can use to fine-tune LLMs in a memory-efficient fashion?

- A. NeMo
- B. TensorRT
- C. Chinchilla
- **D. LoRA**

Answer: D

Explanation:

LoRA (Low-Rank Adaptation) is a parameter-efficient fine-tuning approach specifically designed for large language models (LLMs), as covered in NVIDIA's Generative AI and LLMs course. It fine-tunes LLMs by updating a small subset of parameters through low-rank matrix factorization, significantly reducing memory and computational requirements compared to full fine-tuning. This makes LoRA ideal for adapting large models to specific tasks while maintaining efficiency. Option A, TensorRT, is incorrect, as it is an inference optimization library, not a fine-tuning method. Option B, NeMo, is a framework for building AI models, not a specific fine-tuning technique. Option C, Chinchilla, is a model, not a fine-tuning approach. The course emphasizes: "Parameter-efficient fine-tuning methods like LoRA enable memory-efficient adaptation of LLMs by updating low-rank approximations of weight matrices, reducing resource demands while maintaining performance." References: NVIDIA Building Transformer-Based Natural Language Processing Applications course; NVIDIA Introduction to Transformer-Based Natural Language Processing.

NEW QUESTION # 51

You are in need of customizing your LLM via prompt engineering, prompt learning, or parameter-efficient fine-tuning. Which framework helps you with all of these?

- A. NVIDIA Triton
- B. NVIDIA DALI
- **C. NVIDIA NeMo**
- D. NVIDIA TensorRT

Answer: C

Explanation:

The NVIDIA NeMo framework is designed to support the development and customization of large language models (LLMs), including techniques like prompt engineering, prompt learning (e.g., prompt tuning), and parameter-efficient fine-tuning (e.g., LoRA), as emphasized in NVIDIA's Generative AI and LLMs course.

NeMo provides modular tools and pre-trained models that facilitate these customization methods, allowing users to adapt LLMs for specific tasks efficiently. Option A, TensorRT, is incorrect, as it focuses on inference optimization, not model customization. Option B, DALI, is a data loading library for computer vision, not LLMs. Option C, Triton, is an inference server, not a framework for LLM customization. The course notes:

"NVIDIA NeMo supports LLM customization through prompt engineering, prompt learning, and parameter-efficient fine-tuning, enabling flexible adaptation for NLP tasks." References: NVIDIA Building Transformer-Based Natural Language Processing Applications course; NVIDIA NeMo Framework User Guide.

NEW QUESTION # 52

Which of the following is a key characteristic of Rapid Application Development (RAD)?

- **A. Iterative prototyping with active user involvement.**
- B. Minimal user feedback during the development process.
- C. Extensive upfront planning before any development.
- D. Linear progression through predefined project phases.

Answer: A

Explanation:

Rapid Application Development (RAD) is a software development methodology that emphasizes iterative prototyping and active user involvement to accelerate development and ensure alignment with user needs.

NVIDIA's documentation on AI application development, particularly in the context of NGC (NVIDIA GPU Cloud) and software

