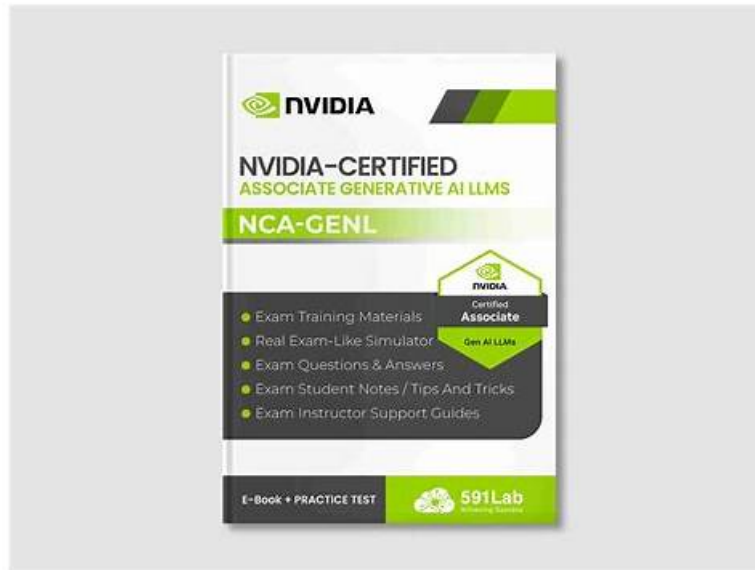


# New NCA-GENL Test Syllabus, NCA-GENL Valid Study Materials



BTW, DOWNLOAD part of TestSimulate NCA-GENL dumps from Cloud Storage: <https://drive.google.com/open?id=1z0xPVdfWpIm4J-DlHQWmAn8vYlhd9KZt>

TestSimulate will provide exam prep and NVIDIA NCA-GENL Exam Simulations you will need to take a certification examination. About NVIDIA NCA-GENL test, you can find related dumps from different websites or books, however, TestSimulate has the advantage of perfect contents, strong logicity and complete supporting facilities. TestSimulate original questions and test answers can not only help you to pass an exam, can also save you valuable time.

If you are worried about your NCA-GENL real exam and you are not prepared so, now you don't need to take any stress about it. Get most updated NVIDIA dumps torrent with 100% accurate answers. Our website is considered one of the best website where you can save extra money by getting one-year of free updates after buying the NCA-GENL Dumps PDF files.

>> **New NCA-GENL Test Syllabus** <<

## NCA-GENL Valid Study Materials - Best NCA-GENL Vce

Our NCA-GENL study materials have a high quality which is mainly reflected in the pass rate. Our product can promise a higher pass rate than other study materials. 99% people who have used our NCA-GENL study materials passed their exam and got their certificate successfully, it is no doubt that it means our NCA-GENL study materials have a 99% pass rate. So our product will be a very good choice for you. If you are anxious about whether you can pass your exam and get the certificate, we think you need to buy our NCA-GENL Study Materials as your study tool, our product will lend you a good helping hand. If you are willing to take our NCA-GENL study materials into more consideration, it must be very easy for you to pass your exam in a short time.

### NVIDIA NCA-GENL Exam Syllabus Topics:

Topic	Details
Topic 1	<ul style="list-style-type: none"><li>Software development: Covers the programming practices and coding skills required to build, maintain, and deploy generative AI applications.</li></ul>
Topic 2	<ul style="list-style-type: none"><li>LLM integration and deployment: Addresses connecting LLMs into real-world applications and deploying them reliably across production environments.</li></ul>
Topic 3	<ul style="list-style-type: none"><li>Experiment design: Focuses on structuring controlled tests and workflows to systematically evaluate LLM performance and outcomes.</li></ul>

Topic 4	<ul style="list-style-type: none"> <li>• <b>Alignment:</b> Addresses methods for ensuring LLM behavior is safe, accurate, and consistent with human intentions and values.</li> </ul>
Topic 5	<ul style="list-style-type: none"> <li>• <b>Data preprocessing and feature engineering:</b> Covers preparing raw data through cleaning, transformation, and feature selection to make it suitable for model training.</li> </ul>
Topic 6	<ul style="list-style-type: none"> <li>• <b>Fundamentals of machine learning and neural networks:</b> 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.</li> </ul>

## NVIDIA Generative AI LLMs Sample Questions (Q82-Q87):

### NEW QUESTION # 82

In the context of language models, what does an autoregressive model predict?

- A. The probability of the next token using a Monte Carlo sampling of past tokens.
- B. The probability of the next token by looking at the previous and future input tokens.
- **C. The probability of the next token in a text given the previous tokens.**
- D. The next token solely using recurrent network or LSTM cells.

**Answer: C**

Explanation:

Autoregressive models are a cornerstone of modern language modeling, particularly in large language models (LLMs) like those discussed in NVIDIA's Generative AI and LLMs course. These models predict the probability of the next token in a sequence based solely on the preceding tokens, making them inherently sequential and unidirectional. This process is often referred to as "next-token prediction," where the model learns to generate text by estimating the conditional probability distribution of the next token given the context of all previous tokens. For example, given the sequence "The cat is," the model predicts the likelihood of the next word being "on," "in," or another token. This approach is fundamental to models like GPT, which rely on autoregressive decoding to generate coherent text. Unlike bidirectional models (e.g., BERT), which consider both previous and future tokens, autoregressive models focus only on past tokens, making option D incorrect. Options B and C are also inaccurate, as Monte Carlo sampling is not a standard method for next-token prediction in autoregressive models, and the prediction is not limited to recurrent networks or LSTM cells, as modern LLMs often use Transformer architectures. The course emphasizes this concept in the context of Transformer-based NLP: "Learn the basic concepts behind autoregressive generative models, including next-token prediction and its implementation within Transformer-based models." References: NVIDIA Building Transformer-Based Natural Language Processing Applications course; NVIDIA Introduction to Transformer-Based Natural Language Processing.

### NEW QUESTION # 83

Which metric is commonly used to evaluate machine-translation models?

- A. Perplexity
- B. ROUGE score
- C. F1 Score
- **D. BLEU score**

**Answer: D**

Explanation:

The BLEU (Bilingual Evaluation Understudy) score is the most commonly used metric for evaluating machine-translation models. It measures the precision of n-gram overlaps between the generated translation and reference translations, providing a quantitative measure of translation quality. NVIDIA's NeMo documentation on NLP tasks, particularly machine translation, highlights BLEU as the standard metric for assessing translation performance due to its focus on precision and fluency. Option A (F1 Score) is used for classification tasks, not translation. Option C (ROUGE) is primarily for summarization, focusing on recall. Option D (Perplexity) measures language model quality but is less specific to translation evaluation.

References:

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

Papineni, K., et al. (2002). "BLEU: A Method for Automatic Evaluation of Machine Translation."

#### NEW QUESTION # 84

What is Retrieval Augmented Generation (RAG)?

- A. RAG is a method for manipulating and generating text-based data using Transformer-based LLMs.
- **B. RAG is a methodology that combines an information retrieval component with a response generator.**
- C. RAG is an architecture used to optimize the output of an LLM by retraining the model with domain- specific data.
- D. RAG is a technique used to fine-tune pre-trained LLMs for improved performance.

**Answer: B**

Explanation:

Retrieval-Augmented Generation (RAG) is a methodology that enhances the performance of large language models (LLMs) by integrating an information retrieval component with a generative model. As described in the seminal paper by Lewis et al. (2020), RAG retrieves relevant documents from an external knowledge base (e.g., using dense vector representations) and uses them to inform the generative process, enabling more accurate and contextually relevant responses. NVIDIA's documentation on generative AI workflows, particularly in the context of NeMo and Triton Inference Server, highlights RAG as a technique to improve LLM outputs by grounding them in external data, especially for tasks requiring factual accuracy or domain- specific knowledge. Option A is incorrect because RAG does not involve retraining the model but rather augments it with retrieved data. Option C is too vague and does not capture the retrieval aspect, while Option D refers to fine-tuning, which is a separate process.

References:

Lewis, P., et al. (2020). "Retrieval-Augmented Generation for Knowledge-Intensive NLP Tasks." NVIDIA NeMo Documentation: <https://docs.nvidia.com/deeplearning/nemo/user-guide/docs/en/stable/nlp/intro.html>

#### NEW QUESTION # 85

Which principle of Trustworthy AI primarily concerns the ethical implications of AI's impact on society and includes considerations for both potential misuse and unintended consequences?

- A. Certification
- **B. Accountability**
- C. Legal Responsibility
- D. Data Privacy

**Answer: B**

Explanation:

Accountability is a core principle of Trustworthy AI that addresses the ethical implications of AI's societal impact, including potential misuse and unintended consequences. NVIDIA's guidelines on Trustworthy AI, as outlined in their AI ethics framework, emphasize accountability as ensuring that AI systems are transparent, responsible, and answerable for their outcomes. This includes mitigating risks of bias, ensuring fairness, and addressing unintended societal impacts. Option A (Certification) refers to compliance processes, not ethical implications. Option B (Data Privacy) focuses on protecting user data, not broader societal impact. Option D (Legal Responsibility) is related but narrower, focusing on liability rather than ethical considerations.

References:

NVIDIA Trustworthy AI: <https://www.nvidia.com/en-us/ai-data-science/trustworthy-ai/>

#### NEW QUESTION # 86

In the development of trustworthy AI systems, what is the primary purpose of implementing red-teaming exercises during the alignment process of large language models?

- **A. To identify and mitigate potential biases, safety risks, and harmful outputs.**
- B. To increase the model's parameter count for better performance.
- C. To optimize the model's inference speed for production deployment.
- D. To automate the collection of training data for fine-tuning.

**Answer: A**

Explanation:

Red-teaming exercises involve systematically testing a large language model (LLM) by probing it with adversarial or challenging inputs to uncover vulnerabilities, such as biases, unsafe responses, or harmful outputs. NVIDIA's Trustworthy AI framework

