

Free PDF Quiz 2026 NVIDIA NCA-GENL: NVIDIA Generative AI LLMs First-grade Latest Dump



What's more, part of that PassExamDumps NCA-GENL dumps now are free: <https://drive.google.com/open?id=1K3X7PeRASaaMaWbGxbvaN6DSgYyw2SWv>

The NCA-GENL study materials of our company is the study tool which best suits these people who long to pass the exam and get the related certification. So we want to tell you that it is high time for you to buy and use our NCA-GENL Study Materials carefully. Now we are glad to introduce the study materials from our company to you in detail in order to let you understanding our study products.

We are committed to helping you pass the exam and get the certificate as soon as possible. NCA-GENL exam bootcamp of us have the questions and answers, and it not only have quality but also contain certain quantity, it will be enough for you to deal with your exam. With the pass rate more than 98.65%, we can ensure you pass your exam. NCA-GENL Exam Dumps also have most of knowledge points of the exam, and they may help you a lot. We offer you free update for 365 days after you purchase the NCA-GENL exam bootcamp.

>> NCA-GENL Latest Dump <<

Pass Guaranteed 2026 NVIDIA NCA-GENL Latest Latest Dump

The test software used in our products is a perfect match for Windows' NCA-GENL learning material, which enables you to enjoy the best learning style on your computer. Our NCA-GENL certification guide also use the latest science and technology to meet the new requirements of authoritative research material network learning. Unlike the traditional way of learning, the great benefit of our NCA-GENL learning material is that users can flexibly adjust their learning plans. We hope that our new design of NCA-GENL test questions will make the user's learning more interesting and colorful.

NVIDIA Generative AI LLMs Sample Questions (Q27-Q32):

NEW QUESTION # 27

Your company has upgraded from a legacy LLM model to a new model that allows for larger sequences and higher token limits. What is the most likely result of upgrading to the new model?

- A. The number of tokens is fixed for all existing language models, so there is no benefit to upgrading to higher token limits.
- **B. The newer model allows larger context, so outputs will improve, but you will likely incur longer inference times.**
- C. The newer model allows for larger context, so the outputs will improve without increasing inference time overhead.
- D. The newer model allows the same context lengths, but the larger token limit will result in more comprehensive and longer outputs with more detail.

Answer: B

Explanation:

Upgrading to a new LLM with larger sequence lengths and higher token limits, as discussed in NVIDIA's Generative AI and LLMs course, typically allows the model to process larger contexts, leading to improved output quality due to better understanding of

extended dependencies in text. However, handling larger sequences increases computational requirements, often resulting in longer inference times, especially on the same hardware. This trade-off is a key consideration in LLM deployment. Option A is incorrect, as token limits vary across models, and higher limits offer benefits. Option B is wrong, as larger context processing typically increases inference time. Option C is inaccurate, as higher token limits primarily enable larger context, not just longer outputs. The course notes: "Larger sequence lengths in LLMs allow for improved output quality by capturing more context, but this often comes at the cost of increased inference times due to higher computational demands." References: NVIDIA Building Transformer-Based Natural Language Processing Applications course; NVIDIA Introduction to Transformer-Based Natural Language Processing.

NEW QUESTION # 28

In the context of fine-tuning LLMs, which of the following metrics is most commonly used to assess the performance of a fine-tuned model?

- A. Number of layers
- **B. Accuracy on a validation set**
- C. Model size
- D. Training duration

Answer: B

Explanation:

When fine-tuning large language models (LLMs), the primary goal is to improve the model's performance on a specific task. The most common metric for assessing this performance is accuracy on a validation set, as it directly measures how well the model generalizes to unseen data. NVIDIA's NeMo framework documentation for fine-tuning LLMs emphasizes the use of validation metrics such as accuracy, F1 score, or task-specific metrics (e.g., BLEU for translation) to evaluate model performance during and after fine-tuning.

These metrics provide a quantitative measure of the model's effectiveness on the target task. Options A, C, and D (model size, training duration, and number of layers) are not performance metrics; they are either architectural characteristics or training parameters that do not directly reflect the model's effectiveness.

References:

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

NEW QUESTION # 29

When designing prompts for a large language model to perform a complex reasoning task, such as solving a multi-step mathematical problem, which advanced prompt engineering technique is most effective in ensuring robust performance across diverse inputs?

- A. Retrieval-augmented generation with external mathematical databases.
- **B. Chain-of-thought prompting with step-by-step reasoning examples.**
- C. Zero-shot prompting with a generic task description.
- D. Few-shot prompting with randomly selected examples.

Answer: B

Explanation:

Chain-of-thought (CoT) prompting is an advanced prompt engineering technique that significantly enhances a large language model's (LLM) performance on complex reasoning tasks, such as multi-step mathematical problems. By including examples that explicitly demonstrate step-by-step reasoning in the prompt, CoT guides the model to break down the problem into intermediate steps, improving accuracy and robustness.

NVIDIA's NeMo documentation on prompt engineering highlights CoT as a powerful method for tasks requiring logical or sequential reasoning, as it leverages the model's ability to mimic structured problem-solving. Research by Wei et al. (2022) demonstrates that CoT outperforms other methods for mathematical reasoning. Option A (zero-shot) is less effective for complex tasks due to lack of guidance. Option B (few-shot with random examples) is suboptimal without structured reasoning. Option D (RAG) is useful for factual queries but less relevant for pure reasoning tasks.

References:

NVIDIA NeMo Documentation: <https://docs.nvidia.com/deeplearning/nemo/user-guide/docs/en/stable/nlp/intro.html> Wei, J., et al. (2022). "Chain-of-Thought Prompting Elicits Reasoning in Large Language Models."

NEW QUESTION # 30

When comparing and contrasting the ReLU and sigmoid activation functions, which statement is true?

- A. ReLU is more computationally efficient, but sigmoid is better for predicting probabilities.
- B. ReLU and sigmoid both have a range of 0 to 1.
- C. ReLU is a linear function while sigmoid is non-linear.
- D. ReLU is less computationally efficient than sigmoid, but it is more accurate than sigmoid.

Answer: A

Explanation:

ReLU (Rectified Linear Unit) and sigmoid are activation functions used in neural networks. According to NVIDIA's deep learning documentation (e.g., cuDNN and TensorRT), ReLU, defined as $f(x) = \max(0, x)$, is computationally efficient because it involves simple thresholding, avoiding expensive exponential calculations required by sigmoid, $f(x) = 1/(1 + e^{-x})$.

P.S. Free 2026 NVIDIA NCA-GENL dumps are available on Google Drive shared by PassExamDumps:
<https://drive.google.com/open?id=1K3X7PeRASaaMaWbGxbvaN6DSgYyw2SWv>