

# Latest NCA-GENL Material | NCA-GENL Test Dates



BONUS!!! Download part of Prep4sures NCA-GENL dumps for free: <https://drive.google.com/open?id=1Agf-UA5ZLqm4nupPn1vaY3UwZB9JMmO>

There are some prominent features that are making the NVIDIA Generative AI LLMs (NCA-GENL) exam dumps the first choice of NCA-GENL certification exam candidates. The prominent features are real and verified NVIDIA Generative AI LLMs (NCA-GENL) exam questions, availability of NVIDIA NVIDIA exam dumps in three different formats, affordable price, 1 year free updated NVIDIA NCA-GENL Exam Questions download facility, and 100 percent NVIDIA NCA-GENL exam passing money back guarantee.

All time and energy you devoted to the NCA-GENL preparation quiz is worthwhile. With passing rate up to 98 percent and above, our NCA-GENL practice materials are highly recommended among exam candidates. So their validity and authority are unquestionable. Our NCA-GENL Learning Materials are just starting points for exam candidates, and you may meet several challenging tasks or exams in the future about computer knowledge, we can still offer help. Need any help, please contact with us again!

>> **Latest NCA-GENL Material** <<

## Reliable NCA-GENL Practice Materials - NCA-GENL Real Study Guide - Prep4sures

Prep4sures, as a provider, specializing in providing all candidates with NCA-GENL exam-related materials, focus on offering the most excellent dumps for the candidates. In contrast with other websites, Prep4sures is more trustworthy. Why? Because Prep4sures has many years of experience and our NVIDIA experts have been devoted themselves to the study of NVIDIA certification exam and summarize NCA-GENL Exam rules. Thus, Prep4sures exam dumps have a high hit rate. Meanwhile, it guarantees the qualification rate in the exam. Therefore, Prep4sures got everyone's trust.

### NVIDIA NCA-GENL Exam Syllabus Topics:

Topic	Details
Topic 1	<ul style="list-style-type: none"><li>• <b>Alignment:</b> This section of the exam measures the skills of AI Policy Engineers and covers techniques to align LLM outputs with human intentions and values. It includes safety mechanisms, ethical safeguards, and tuning strategies to reduce harmful, biased, or inaccurate results from models.</li></ul>
Topic 2	<ul style="list-style-type: none"><li>• <b>Python Libraries for LLMs:</b> This section of the exam measures skills of LLM Developers and covers using Python tools and frameworks like Hugging Face Transformers, LangChain, and PyTorch to build, fine-tune, and deploy large language models. It focuses on practical implementation and ecosystem familiarity.</li></ul>
Topic 3	<ul style="list-style-type: none"><li>• <b>Experiment Design</b></li></ul>

Topic 4	<ul style="list-style-type: none"> <li>• <b>Prompt Engineering:</b> This section of the exam measures the skills of Prompt Designers and covers how to craft effective prompts that guide LLMs to produce desired outputs. It focuses on prompt strategies, formatting, and iterative refinement techniques used in both development and real-world applications of LLMs.</li> </ul>
Topic 5	<ul style="list-style-type: none"> <li>• <b>Data Analysis and Visualization:</b> This section of the exam measures the skills of Data Scientists and covers interpreting, cleaning, and presenting data through visual storytelling. It emphasizes how to use visualization to extract insights and evaluate model behavior, performance, or training data patterns.</li> </ul>
Topic 6	<ul style="list-style-type: none"> <li>• <b>Experimentation:</b> This section of the exam measures the skills of ML Engineers and covers how to conduct structured experiments with LLMs. It involves setting up test cases, tracking performance metrics, and making informed decisions based on experimental outcomes.:</li> </ul>

## NVIDIA Generative AI LLMs Sample Questions (Q59-Q64):

### NEW QUESTION # 59

When implementing data parallel training, which of the following considerations needs to be taken into account?

- A. The model weights are kept independent for as long as possible increasing the model exploration.
- B. The model weights are synced across all processes/devices only at the end of every epoch.
- **C. A ring all-reduce is an efficient algorithm for syncing the weights across different processes/devices.**
- D. A master-worker method for syncing the weights across different processes is desirable due to its scalability.

**Answer: C**

Explanation:

In data parallel training, where a model is replicated across multiple devices with each processing a portion of the data, synchronizing model weights is critical. As covered in NVIDIA's Generative AI and LLMs course, the ring all-reduce algorithm is an efficient method for syncing weights across processes or devices. It minimizes communication overhead by organizing devices in a ring topology, allowing gradients to be aggregated and shared efficiently. Option A is incorrect, as weights are typically synced after each batch, not just at epoch ends, to ensure consistency. Option B is wrong, as master-worker methods can create bottlenecks and are less scalable than all-reduce. Option D is inaccurate, as keeping weights independent defeats the purpose of data parallelism, which requires synchronized updates. The course notes: "In data parallel training, the ring all-reduce algorithm efficiently synchronizes model weights across devices, reducing communication overhead and ensuring consistent updates." References: NVIDIA Building Transformer-Based Natural Language Processing Applications course; NVIDIA Introduction to Transformer-Based Natural Language Processing.

### NEW QUESTION # 60

Which of the following prompt engineering techniques is most effective for improving an LLM's performance on multi-step reasoning tasks?

- A. Retrieval-augmented generation without context
- B. Few-shot prompting with unrelated examples.
- **C. Chain-of-thought prompting with explicit intermediate steps.**
- D. Zero-shot prompting with detailed task descriptions.

**Answer: C**

Explanation:

Chain-of-thought (CoT) prompting is a highly effective technique for improving large language model (LLM) performance on multi-step reasoning tasks. By including explicit intermediate steps in the prompt, CoT guides the model to break down complex problems into manageable parts, improving reasoning accuracy. NVIDIA's NeMo documentation on prompt engineering highlights CoT as a powerful method for tasks like mathematical reasoning or logical problem-solving, as it leverages the model's ability to follow structured reasoning paths. Option A is incorrect, as retrieval-augmented generation (RAG) without context is less effective for reasoning tasks. Option B is wrong, as unrelated examples in few-shot prompting do not aid reasoning. Option C (zero-shot prompting) is less effective than CoT for complex reasoning.

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 # 61

What are the main advantages of instructed large language models over traditional, small language models (< 300M parameters)? (Pick the 2 correct responses)

- A. Smaller latency, higher throughput.
- **B. Single generic model can do more than one task.**
- C. Trained without the need for labeled data.
- **D. Cheaper computational costs during inference.**
- E. It is easier to explain the predictions.

**Answer: B,D**

Explanation:

Instructed large language models (LLMs), such as those supported by NVIDIA's NeMo framework, have significant advantages over smaller, traditional models:

\* Option D: LLMs often have cheaper computational costs during inference for certain tasks because they can generalize across multiple tasks without requiring task-specific retraining, unlike smaller models that may need separate models per task.

References:

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

/intro.html

Brown, T., et al. (2020). "Language Models are Few-Shot Learners."

### NEW QUESTION # 62

In the context of evaluating a fine-tuned LLM for a text classification task, which experimental design technique ensures robust performance estimation when dealing with imbalanced datasets?

- **A. Stratified k-fold cross-validation.**
- B. Bootstrapping with random sampling.
- C. Single hold-out validation with a fixed test set.
- D. Grid search for hyperparameter tuning.

**Answer: A**

Explanation:

Stratified k-fold cross-validation is a robust experimental design technique for evaluating machine learning models, especially on imbalanced datasets. It divides the dataset into k folds while preserving the class distribution in each fold, ensuring that the model is evaluated on representative samples of all classes.

NVIDIA's NeMo documentation on model evaluation recommends stratified cross-validation for tasks like text classification to obtain reliable performance estimates, particularly when classes are unevenly distributed (e.g., in sentiment analysis with few negative samples). Option A (single hold-out) is less robust, as it may not capture class imbalance. Option C (bootstrapping) introduces variability and is less suitable for imbalanced data. Option D (grid search) is for hyperparameter tuning, not performance estimation.

References:

NVIDIA NeMo Documentation: [https://docs.nvidia.com/deeplearning/nemo/user-guide/docs/en/stable/nlp/model\\_fineting.html](https://docs.nvidia.com/deeplearning/nemo/user-guide/docs/en/stable/nlp/model_fineting.html)

### NEW QUESTION # 63

Why do we need positional encoding in transformer-based models?

- **A. To represent the order of elements in a sequence.**
- B. To increase the throughput of the model.
- C. To prevent overfitting of the model.
- D. To reduce the dimensionality of the input data.

**Answer: A**

Explanation:



UA5ZLqm4nupPn1vaY3UwZB9JMmO