

NCA Generative AI LLM (NCA-GENL) Practice Exam (Sample)

Study Guide



Everything you need from our exam experts!

This is a sample study guide. To access the full version with hundreds of questions,

Copyright © 2026 by Examzify - A Kaluba Technologies Inc. product.

ALL RIGHTS RESERVED.

No part of this book may be reproduced or transferred in any form or by any means, graphic, electronic, or mechanical, including photocopying, recording, web distribution, taping, or by any information storage retrieval system, without the written permission of the author.

Notice: Examzify makes every reasonable effort to obtain from reliable sources accurate, complete, and timely information about this product.

SAMPLE

Table of Contents

Copyright	1
Table of Contents	2
Introduction	3
How to Use This Guide	4
Questions	6
Answers	9
Explanations	11
Next Steps	17

SAMPLE

Introduction

Preparing for a certification exam can feel overwhelming, but with the right tools, it becomes an opportunity to build confidence, sharpen your skills, and move one step closer to your goals. At Examzify, we believe that effective exam preparation isn't just about memorization, it's about understanding the material, identifying knowledge gaps, and building the test-taking strategies that lead to success.

This guide was designed to help you do exactly that.

Whether you're preparing for a licensing exam, professional certification, or entry-level qualification, this book offers structured practice to reinforce key concepts. You'll find a wide range of multiple-choice questions, each followed by clear explanations to help you understand not just the right answer, but why it's correct.

The content in this guide is based on real-world exam objectives and aligned with the types of questions and topics commonly found on official tests. It's ideal for learners who want to:

- Practice answering questions under realistic conditions,
- Improve accuracy and speed,
- Review explanations to strengthen weak areas, and
- Approach the exam with greater confidence.

We recommend using this book not as a stand-alone study tool, but alongside other resources like flashcards, textbooks, or hands-on training. For best results, we recommend working through each question, reflecting on the explanation provided, and revisiting the topics that challenge you most.

Remember: successful test preparation isn't about getting every question right the first time, it's about learning from your mistakes and improving over time. Stay focused, trust the process, and know that every page you turn brings you closer to success.

Let's begin.

How to Use This Guide

This guide is designed to help you study more effectively and approach your exam with confidence. Whether you're reviewing for the first time or doing a final refresh, here's how to get the most out of your Examzify study guide:

1. Start with a Diagnostic Review

Skim through the questions to get a sense of what you know and what you need to focus on. Don't worry about getting everything right, your goal is to identify knowledge gaps early.

2. Study in Short, Focused Sessions

Break your study time into manageable blocks (e.g. 30 - 45 minutes). Review a handful of questions, reflect on the explanations, and take breaks to retain information better.

3. Learn from the Explanations

After answering a question, always read the explanation, even if you got it right. It reinforces key points, corrects misunderstandings, and teaches subtle distinctions between similar answers.

4. Track Your Progress

Use bookmarks or notes (if reading digitally) to mark difficult questions. Revisit these regularly and track improvements over time.

5. Simulate the Real Exam

Once you're comfortable, try taking a full set of questions without pausing. Set a timer and simulate test-day conditions to build confidence and time management skills.

6. Repeat and Review

Don't just study once, repetition builds retention. Re-attempt questions after a few days and revisit explanations to reinforce learning.

7. Use Other Tools

Pair this guide with other Examzify tools like flashcards, and digital practice tests to strengthen your preparation across formats.

There's no single right way to study, but consistent, thoughtful effort always wins. Use this guide flexibly — adapt the tips above to fit your pace and learning style. You've got this!

SAMPLE

Questions

SAMPLE

- 1. Which of the following is NOT a feature of the NVIDIA AI Workbench?**
 - A. Data visualization tools**
 - B. Interactive model training environments**
 - C. Real-time performance monitoring**
 - D. Automated chat response generation**
- 2. Which training approach is less efficient due to the computational overhead it introduces?**
 - A. Synchronous Updates**
 - B. Asynchronous Updates**
 - C. Gradient Checkpointing**
 - D. Objective Function**
- 3. Which framework supports the creation and development of LLM microservices, enabling task division through APIs?**
 - A. Nvidia TensorFlow**
 - B. NeMo Microservices**
 - C. CUDA Toolkit**
 - D. Nvidia Triton Inference Server**
- 4. What does CUDA Graph with Fusion optimize in GPU tasks?**
 - A. Kernel launches and memory operations**
 - B. Data storage management**
 - C. Error detection and correction**
 - D. Performance monitoring**
- 5. Which mechanisms increase the interpretability of LLMs by indicating influential input data?**
 - A. Recurrent Neural Networks**
 - B. Attention Mechanisms**
 - C. Convolutional Layers**
 - D. Activation Functions**

6. Which activation function is often used in output layers for binary classification problems?

- A. ReLU**
- B. Sigmoid**
- C. ELU**
- D. Tanh**

7. What learning strategy is essential for deep learning models to adapt without losing previous knowledge?

- A. Forgetting Theory**
- B. Replay Buffer**
- C. Fine-tuning**
- D. Transfer Learning**

8. What is the main characteristic of the RMSProp optimization algorithm?

- A. Maintains a fixed learning rate**
- B. Adapts learning rate based on gradient squares**
- C. Provides very slow convergence**
- D. Involves no gradient descent**

9. What distributed technique for LLM development involves splitting data across multiple GPUs?

- A. Model Parallelism**
- B. Data Parallelism**
- C. Hybrid Parallelism**
- D. Layer Parallelism**

10. What does the acronym LLM stand for in the context of AI?

- A. Large Language Model**
- B. Long Learning Model**
- C. Lightweight Language Model**
- D. Linear Learning Model**

Answers

SAMPLE

1. D
2. A
3. B
4. A
5. B
6. B
7. B
8. B
9. B
10. A

SAMPLE

Explanations

SAMPLE

1. Which of the following is NOT a feature of the NVIDIA AI Workbench?

- A. Data visualization tools**
- B. Interactive model training environments**
- C. Real-time performance monitoring**
- D. Automated chat response generation**

The correct choice indicates that automated chat response generation is not a feature of the NVIDIA AI Workbench. The NVIDIA AI Workbench is designed primarily to facilitate the processes involved in AI development, specifically for training and tuning models. It provides robust data visualization tools that allow users to better understand their datasets and model performance. Additionally, it supports interactive model training environments, enabling users to iteratively refine their models based on real-time feedback and results. Moreover, real-time performance monitoring is a crucial aspect of the NVIDIA AI Workbench, as it allows developers to track how their models are performing during training and make necessary adjustments to optimize results. In contrast, automated chat response generation is more closely associated with specific applications of generative AI models, rather than the infrastructure or tools provided by the AI Workbench for developing and managing AI systems. Thus, this makes automated chat response generation a feature that does not belong to the NVIDIA AI Workbench suite.

2. Which training approach is less efficient due to the computational overhead it introduces?

- A. Synchronous Updates**
- B. Asynchronous Updates**
- C. Gradient Checkpointing**
- D. Objective Function**

Synchronous updates are considered less efficient primarily due to the computational overhead involved in coordinating updates from all participating nodes before proceeding with the next round of model training. In this approach, each worker must wait for all other workers to complete their computations and share their gradients before any updates to the model can occur. This waiting time can lead to increased latency, particularly in large, distributed training scenarios where the communication overhead becomes significant. In contrast, other training approaches, like asynchronous updates, allow individual workers to update the model independently, which can lead to faster overall training times since there is no need for synchronization at every step. Gradient checkpointing and the objective function primarily relate to memory management and optimization, respectively. They do not inherently introduce the same level of coordination overhead that synchronous updates incur. Thus, the synchronous update mechanism's requirement for complete collaboration at each update step makes it less efficient in terms of computational resources.

3. Which framework supports the creation and development of LLM microservices, enabling task division through APIs?

- A. Nvidia TensorFlow
- B. NeMo Microservices**
- C. CUDA Toolkit
- D. Nvidia Triton Inference Server

The correct answer highlights NeMo Microservices as the framework specifically designed for the creation and development of large language model (LLM) microservices. This framework enables developers to divide tasks through APIs, which is essential for efficiently managing and scaling AI applications. NeMo Microservices allows users to build modular components that can independently handle specific tasks, making it easier to deploy and integrate various machine learning functionalities within larger systems. NeMo Microservices is part of the broader NVIDIA NeMo ecosystem, which is tailored for building and fine-tuning conversational AI models. This modular architecture facilitates seamless communication between different services, facilitating more effective collaboration and resource allocation in AI deployments. In contrast, other options, while related to AI and model training, do not primarily focus on the microservice architecture aimed at enabling task division through APIs. For instance, TensorFlow is a comprehensive machine learning framework, but it does not specifically emphasize microservices for LLMs. The CUDA Toolkit is oriented toward parallel computing and GPU acceleration, while the Nvidia Triton Inference Server is focused on model deployment and real-time inference rather than the creation and development of modular microservices for LLMs.

4. What does CUDA Graph with Fusion optimize in GPU tasks?

- A. Kernel launches and memory operations**
- B. Data storage management
- C. Error detection and correction
- D. Performance monitoring

The correct answer focuses on the optimization of kernel launches and memory operations in GPU tasks through CUDA Graph with Fusion. This approach is particularly significant because it allows developers to define a series of operations as a single task graph. By doing this, multiple operations can be executed without the overhead typically associated with launching kernels and managing memory transfers. CUDA Graphs can minimize the number of interactions with the GPU by allowing the execution of multiple operations in one go, thereby reducing launch overhead and optimizing memory transfers. The fusion aspect refers to merging several operations into a single kernel execution, which further enhances performance by leveraging the GPU's resources more effectively. This optimization is crucial for improving the performance of GPU-accelerated applications, as it directly affects the efficiency with which data is processed and the overall speed of computations. Other options like data storage management, error detection and correction, and performance monitoring, while important in the context of GPU tasks, do not specifically align with the primary focus of CUDA Graph and Fusion, which is centered on optimizing kernel and memory operation efficiencies.

5. Which mechanisms increase the interpretability of LLMs by indicating influential input data?

A. Recurrent Neural Networks

B. Attention Mechanisms

C. Convolutional Layers

D. Activation Functions

Attention mechanisms play a crucial role in increasing the interpretability of large language models (LLMs) because they provide a way to visualize and understand the contributions of different input tokens to the model's predictions. In essence, attention allows the model to focus on specific words or phrases in the input when generating an output, highlighting which parts of the input are most influential in determining the response. By examining the attention weights that are assigned to each part of the input during processing, researchers and practitioners can gain insights into how the model makes decisions and which elements are driving those decisions. This capability is particularly valuable in natural language processing tasks, where understanding context and relationships between words is vital. Because attention mechanisms explicitly quantify the importance of different inputs, they serve as a powerful tool for interpreting model behavior. In contrast, the other options do not inherently provide this level of interpretability. Recurrent neural networks, while effective for sequence data, don't provide direct insights into which inputs influenced a specific output. Convolutional layers are primarily used for spatial data and lack mechanisms for direct interpretability in sequential tasks such as language. Activation functions, on the other hand, serve to introduce non-linearity in the model but do not offer any direct explanation of input influence or decision-making

6. Which activation function is often used in output layers for binary classification problems?

A. ReLU

B. Sigmoid

C. ELU

D. Tanh

In the context of binary classification problems, the sigmoid activation function is particularly well-suited for the output layer. This is due to its characteristic of producing an output in the range of 0 to 1, which aligns perfectly with the interpretation of binary outcomes, such as representing probabilities. When a model needs to determine membership in one of two classes, the sigmoid function can effectively map any real-valued input into this range, allowing the model to output a probability score. A threshold is then typically applied (commonly set at 0.5) to decide which class the input should belong to based on the predicted probability. In contrast, while other activation functions like ReLU, ELU, and Tanh have their respective strengths in different contexts, they do not serve the same purpose for binary classification output layers. For instance, ReLU and ELU can produce outputs that extend beyond the $[0, 1]$ interval, making them unsuitable for probability estimation. Tanh outputs values in the range of -1 to 1, which does not fit the requirement for a probability representation in binary classification tasks. Thus, the sigmoid activation function is the ideal choice for producing a clear and interpretable output in binary classification scenarios.

7. What learning strategy is essential for deep learning models to adapt without losing previous knowledge?

- A. Forgetting Theory**
- B. Replay Buffer**
- C. Fine-tuning**
- D. Transfer Learning**

The concept of a replay buffer is pivotal in deep learning models, particularly in the context of continual learning. A replay buffer enables the model to store and reuse previous experiences or data samples, which helps maintain knowledge that may otherwise be lost when adapting to new information. This mechanism is especially useful in scenarios where the model needs to learn new tasks over time while still retaining skills from earlier tasks. By periodically revisiting and reinforcing past knowledge through samples held in the replay buffer, the model can effectively reduce the risk of catastrophic forgetting, where learning new information disrupts or erases earlier knowledge. In contrast, other strategies such as fine-tuning, transfer learning, and forgetting theory focus on different aspects of model adaptability and learning. Fine-tuning usually involves making adjustments to an already trained model on a new dataset, but it doesn't inherently include mechanisms for preserving old knowledge. Transfer learning refers to applying knowledge gained in one context to a different but related context, which may not directly support continuous adaptation without forgetting previous knowledge. Forgetting theory offers insights into the nature and mechanisms of learning and memory decay but does not provide a direct strategy for retaining knowledge while learning new information. The replay buffer, therefore, stands out as the essential strategy for enabling deep learning models to

8. What is the main characteristic of the RMSProp optimization algorithm?

- A. Maintains a fixed learning rate**
- B. Adapts learning rate based on gradient squares**
- C. Provides very slow convergence**
- D. Involves no gradient descent**

The main characteristic of the RMSProp optimization algorithm is its ability to adapt the learning rate based on the magnitude of the recent gradients. This method involves maintaining a moving average of the squares of past gradients, which allows the algorithm to adjust the learning rate dynamically. When the gradients are large, the learning rate is decreased, and when the gradients are small, the learning rate is increased. This adaptive approach helps stabilize the optimization process, especially in the presence of noisy or varying gradients, ultimately leading to improved convergence properties in training neural networks. The RMSProp algorithm is particularly beneficial in dealing with the challenges posed by non-stationary objectives, as it helps mitigate issues related to the selection of an appropriate static learning rate. This makes it especially effective in scenarios involving complex loss surfaces and diverse training datasets.

9. What distributed technique for LLM development involves splitting data across multiple GPUs?

- A. Model Parallelism**
- B. Data Parallelism**
- C. Hybrid Parallelism**
- D. Layer Parallelism**

The technique for distributed LLM (Large Language Model) development that involves splitting data across multiple GPUs is known as data parallelism. In this approach, each GPU holds a complete copy of the model while processing different subsets of the data simultaneously. This allows for efficient training because multiple batches of data can be handled at once, significantly speeding up the overall training process. When using data parallelism, gradients from each GPU are aggregated, which means that after each training step, the weights of the model are updated based on the collective information gained from all the processed data. This method maximizes the utilization of multiple GPUs by focusing on handling data input efficiently rather than changing the structure of the model itself. Consequently, it is particularly effective for large datasets, where the entire dataset may not fit onto a single GPU. Although there are other techniques like model parallelism, which divides the model architecture across GPUs, or hybrid and layer parallelism that explore different dimensions of parallelization, they do not specifically focus on distributing data batches across multiple GPUs in the same way that data parallelism does.

10. What does the acronym LLM stand for in the context of AI?

- A. Large Language Model**
- B. Long Learning Model**
- C. Lightweight Language Model**
- D. Linear Learning Model**

In the context of AI, LLM stands for Large Language Model. This term refers to a type of artificial intelligence that has been trained on vast amounts of text data to understand and generate human-like language. Large Language Models leverage deep learning techniques, particularly using architectures such as transformers, to process and generate text that is coherent and contextually relevant. The significance of "large" in Large Language Models indicates the size of the training dataset and the complexity of the model itself—typically characterized by millions or even billions of parameters. This scale allows these models to capture nuanced patterns in language, making them effective for a variety of applications such as text generation, translation, summarization, and conversational agents. Understanding this term is crucial for grasping the current advancements in AI language processing, as LLMs represent a significant leap in the capabilities of natural language understanding and generation compared to previous models.

Next Steps

Congratulations on reaching the final section of this guide. You've taken a meaningful step toward passing your certification exam and advancing your career.

As you continue preparing, remember that consistent practice, review, and self-reflection are key to success. Make time to revisit difficult topics, simulate exam conditions, and track your progress along the way.

If you need help, have suggestions, or want to share feedback, we'd love to hear from you. Reach out to our team at hello@examzify.com.

Or visit your dedicated course page for more study tools and resources:

<https://ncagenl.examzify.com>

We wish you the very best on your exam journey. You've got this!

SAMPLE