

権威のある NVIDIA NCA-GENL 受験体験 & 合格スムーズ NCA-GENL 最速合格 | 信頼できる NCA-GENL トレーニング



さらに、PassTest NCA-GENL ダンプの一部が現在無料で提供されています: https://drive.google.com/open?id=1LPI6EV7zok2_qDbyMvdB5_U5M14wJQxK

試験準備のための学習資料を見つけている場合、当社の資料は検索を終了します。私たちの NCA-GENL 試験トレントは、あなたが期待できない高品質を持っています。NCA-GENL トレントは時間を大幅に節約するのに役立ち、あなたがやりたいことをする自由が増えると思います。私たちの NCA-GENL テスト問題集の使用について後悔がないことを保証できます。アクションの時間が来たら、思考を止めて、入って、私たちの NCA-GENL 試験トレントを試してください。NCA-GENL 試験に合格し、短時間で証明書を取得する必要があります。

NVIDIA NCA-GENL 認定試験の出題範囲:

トピック	出題範囲
トピック 1	<ul style="list-style-type: none">LLM の統合と展開: LLM を実際のアプリケーションに接続し、実稼働環境全体に確実に展開する方法について説明します。
トピック 2	<ul style="list-style-type: none">実験設計: LLM のパフォーマンスと結果を体系的に評価するために、制御されたテストとワークフローの構造化に重点を置きます。
トピック 3	<ul style="list-style-type: none">整合: LLM の動作が安全かつ正確であり、人間の意図や価値観と一致していることを保証する方法について説明します。
トピック 4	<ul style="list-style-type: none">データの事前処理と特徴エンジニアリング: クリーニング、変換、特徴選択を通じて生データを準備し、モデルのトレーニングに適したものにします。
トピック 5	<ul style="list-style-type: none">LLM 用の Python ライブラリ: LangChain、Hugging Face などの、LLM の構築と操作に使用される主要な Python フレームワークとツールについて説明します。
トピック 6	<ul style="list-style-type: none">ソフトウェア開発: 生成 AI アプリケーションの構築、保守、展開に必要なプログラミング手法とコーディングスキルをカバーします。
トピック 7	<ul style="list-style-type: none">実験: モデルの動作をテストし、アプローチを比較し、生成 AI ソリューションを検証するためのトライアルの実行と評価を検討します。
トピック 8	<ul style="list-style-type: none">プロンプトエンジニアリング: LLM 出力を効果的に望ましい結果に導くために、入力プロンプトを設計および改良する手法に焦点を当てます。

- データ分析と視覚化: データセットを解釈し、視覚的なツールを通じて洞察を提示して、情報に基づいたモデル開発の意思決定をサポートします。

>> NCA-GENL受験体験 <<

NVIDIA NCA-GENL最速合格 & NCA-GENLトレーニング

NCA-GENL試験に合格すると多くのメリットが得られることは誰もが知っていますが、NVIDIAすべての受験者がそれを達成するのは容易ではありません。NCA-GENLガイド急流は、すべての受験者が試験に合格するのを支援することを目的としたツールです。私たちの試験資料は、コンピュータと人の量に制限なしでインストールおよびダウンロードできます。弊社が提供するNCA-GENL学習資料が有用であり、テストに合格するのに役立つことを保証します。製品を購入すると、便利な方法を使用して、いつでもどこでもNCA-GENL試験トレントを学習できます。そのため、購入の前後に安心して、NCA-GENL学習教材にウイルスがないことを信頼してください。NVIDIA Generative AI LLMs当社の製品PassTestに慣れるために、NCA-GENL学習教材の機能と利点を次のようにリストします。

NVIDIA Generative AI LLMs 認定 NCA-GENL 試験問題 (Q87-Q92):

質問 #87

Which of the following is a feature of the NVIDIA Triton Inference Server?

- A. Dynamic batching
- B. Model pruning
- C. Model quantization
- D. Gradient clipping

正解: A

解説:

The NVIDIA Triton Inference Server is designed to optimize and deploy machine learning models for inference, and one of its key features is dynamic batching, as noted in NVIDIA's Generative AI and LLMs course. Dynamic batching automatically groups inference requests into batches to maximize GPU utilization, reducing latency and improving throughput for real-time applications. Option A, model quantization, is incorrect, as it is typically handled by frameworks like TensorRT, not Triton. Option C, gradient clipping, is a training technique, not an inference feature. Option D, model pruning, is a model optimization method, not a Triton feature. The course states: "NVIDIA Triton Inference Server supports dynamic batching, which optimizes inference by grouping requests to maximize GPU efficiency and throughput." References: NVIDIA Building Transformer-Based Natural Language Processing Applications course; NVIDIA Introduction to Transformer-Based Natural Language Processing.

質問 #88

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. Zero-shot prompting with detailed task descriptions.
- C. Few-shot prompting with unrelated examples.
- D. Chain-of-thought prompting with explicit intermediate steps.

正解: D

解説:

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."

質問 # 89

In neural networks, the vanishing gradient problem refers to what problem or issue?

- A. The problem of overfitting in neural networks, where the model performs well on the training data but poorly on new, unseen data.
- B. The issue of gradients becoming too large during backpropagation, leading to unstable training.
- C. The problem of underfitting in neural networks, where the model fails to capture the underlying patterns in the data.
- **D. The issue of gradients becoming too small during backpropagation, resulting in slow convergence or stagnation of the training process.**

正解: D

解説:

The vanishing gradient problem occurs in deep neural networks when gradients become too small during backpropagation, causing slow convergence or stagnation in training, particularly in deeper layers. NVIDIA's documentation on deep learning fundamentals, such as in CUDA and cuDNN guides, explains that this issue is common in architectures like RNNs or deep feedforward networks with certain activation functions (e.g., sigmoid). Techniques like ReLU activation, batch normalization, or residual connections (used in transformers) mitigate this problem. Option A (overfitting) is unrelated to gradients. Option B describes the exploding gradient problem, not vanishing gradients. Option C (underfitting) is a performance issue, not a gradient-related problem.

References:

NVIDIA CUDA Documentation: <https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html> Goodfellow, I., et al. (2016). "Deep Learning." MIT Press.

質問 # 90

What is the Open Neural Network Exchange (ONNX) format used for?

- A. Reducing training time of neural networks
- B. Compressing deep learning models
- **C. Representing deep learning models**
- D. Sharing neural network literature

正解: C

解説:

The Open Neural Network Exchange (ONNX) format is an open-standard representation for deep learning models, enabling interoperability across different frameworks, as highlighted in NVIDIA's Generative AI and LLMs course. ONNX allows models trained in frameworks like PyTorch or TensorFlow to be exported and used in other compatible tools for inference or further development, ensuring portability and flexibility.

Option B is incorrect, as ONNX is not designed to reduce training time but to standardize model representation. Option C is wrong, as model compression is handled by techniques like quantization, not ONNX. Option D is inaccurate, as ONNX is unrelated to sharing literature. The course states: "ONNX is an open format for representing deep learning models, enabling seamless model exchange and deployment across various frameworks and platforms." References: NVIDIA Building Transformer-Based Natural Language Processing Applications course; NVIDIA Introduction to Transformer-Based Natural Language Processing.

質問 # 91

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

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

正解: D

