

英伟达助理级认证：生成式人工智能大型语言模型

NVIDIA-Certified Associate: Generative AI LLMs (NCA-GENL)

英伟达 GenAI 生态系统考点回顾

Cheat Sheet for NVIDIA GenAI Ecosystem

作者：DolbyUUU (<https://github.com/DolbyUUU>)

GPU 架构与核心 (GPU Architecture & Cores)

NVIDIA Hopper (H100, H200) / Blackwell (B100, B200, GB200)

这分别是英伟达当前和下一代的旗舰 AI 计算架构，为训练和推理生成式 AI 模型提供核心算力。Hopper 架构是目前市场的主力，而 Blackwell 架构通过系统级创新，为未来万亿参数级别的巨型模型提供了更高的性能和效率。

NVIDIA Tensor Cores

这是 GPU 内部专门为深度学习设计的硬件处理单元，其核心功能是极大地加速矩阵乘法与累加运算。由于 Transformer 等主流 AI 模型严重依赖此类计算，Tensor Cores 是实现高性能 AI 训练和推理的物理基础。

Transformer Engine

这是一项结合了硬件与软件的创新技术，内置于 Hopper 及更新的 GPU 架构中。它能够智能地在 FP8 和 FP16 等不同数值精度之间动态切换，从而在不牺牲模型准确率的前提下，显著提升 Transformer 模型的运行速度。

计算与网络系统(Compute & Networking Systems)

NVIDIA DGX / DGX SuperPOD

NVIDIA DGX 是集成了顶级 GPU、高速网络和全套 AI 软件的“一站式”AI 超级计算机，开箱即用。而 DGX SuperPOD 则是由多个 DGX 系统组成的集群参考架构，是企业构建大规模 LLM 训练中心的首选蓝图。

NVIDIA NVLink & NVSwitch

这是一套用于 GPU 之间高速通信的互联技术，旨在解决多 GPU 协同工作时的通信瓶颈。NVLink 提供 GPU 间的直接高速通道，而 NVSwitch 则像一个智能交换机，允许多达数百个 GPU 像一个整体一样全速通信。

NVIDIA BlueField DPU (Data Processing Unit)

这是一种可编程的数据处理器，其主要任务是将网络、存储和安全等基础设施负载从 CPU 和 GPU 上卸载下来。在 AI 集群中，DPU 可以独立处理数据通信任务，从而让宝贵的 CPU 和 GPU 资源能 100% 专注于核心的 AI 计算。

NVIDIA Spectrum-X Ethernet Platform

这是一个专为 AI 工作负载优化的以太网网络解决方案，旨在提供可与传统 InfiniBand 网络相媲美的超高性能。它通过将 Spectrum 交换机与 BlueField DPU 相结合，在开放的以太网生态中实现了 AI 所需的高带宽和低延迟。

NVIDIA Jetson AGX Orin / Thor

这是英伟达面向边缘设备（如机器人、自动驾驶汽车）的嵌入式计算平台系列，旨在将数据中心级的 AI 能力带到终端。它们功耗更低、体积更小，能够直接在设备上运行经过优化的生成式 AI 模型，实现实时智能交互。

核心平台与库(Core Platforms & Libraries)

NVIDIA CUDA (Compute Unified Device Architecture)

NVIDIA CUDA 是整个英伟达生态系统的核芯灵魂，它是一个并行计算平台和编程模型。通过 CUDA，开发者可以使用 C++ 或 Python 等通用编程语言，直接调用 GPU 强大的并行处理能力来执行各种计算任务，而不仅仅是图形渲染。

CUDA-X AI Libraries

CUDA-X AI 是英伟达提供的一系列针对 AI 和数据科学领域高度优化的 GPU 加速库的统称。这个集合包含了 cuDNN、NCCL、cuBLAS 等众多专业库，为开发者提供了构建高性能 AI 应用所需的、经过预先优化的软件模块。

NVIDIA cuDNN (CUDA Deep Neural Network library)

NVIDIA cuDNN 是一个专门用于加速深度学习基础运算的 GPU 库。它为卷积、池化、归一化和激活函数等神经网络常用操作提供了高度优化的实现，是 PyTorch 和 TensorFlow 等主流深度学习框架实现高性能的基石。

NVIDIA NCCL (NVIDIA Collective Communications Library)

NVIDIA NCCL 是专为多 GPU 和多节点分布式训练设计的通信库，被誉为分布式训练的“命脉”。它实现了如 All-Reduce 等高效的集体通信操作，确保在进行大规模模型训练时，各个 GPU 之间能够以极高的效率同步数据和梯度。

NVIDIA DALI (Data Loading Library)

NVIDIA DALI 是一个用于加速 AI 训练中数据加载和预处理流程的库。它通过将图像解码、数据增强等计算密集型任务从 CPU 转移到 GPU 上执行，有效解决了数据准备阶段的 CPU 瓶颈，从而提升了端到端的训练效率。

NVIDIA RAPIDS

NVIDIA RAPIDS 是一套开源的软件库，旨在将整个数据科学和分析工作流（从数据准备到机器学习）完全迁移到 GPU 上执行。通过利用 GPU 的并行计算能力，RAPIDS 能够将传统上由 CPU 处理的数据分析任务加速数倍甚至数十倍。

cuDF

cuDF 是 RAPIDS 套件中的核心组件之一，可以看作是著名数据分析库 Pandas 的 GPU 加速版本。它提供了与 Pandas 几乎相同的 API 接口，让数据科学家能够无缝地将其数据处理代码迁移到 GPU 上，以获得巨大的性能提升。

cuML

cuML 是 RAPIDS 套件中的机器学习库，其定位是流行机器学习库 Scikit-learn 的 GPU 加速版。它实现了多种经典的机器学习算法，使开发者能够在 GPU 上高速训练传统模型，极大地缩短了模型迭代和实验的时间。

NVIDIA cuBLAS / cuSPARSE

这两个库是 NVIDIA CUDA 中用于加速线性代数运算的基础组件，是所有上层 AI 框架的数学计算核心。cuBLAS 专注于加速密集的矩阵和向量运算，而 cuSPARSE 则专门优化稀疏矩阵的计算，两者共同构成了 GPU 高性能计算的基石。

LLM/GenAI 专用软件栈(LLM/GenAI-Specific Software Stack)

NVIDIA NeMo Framework

NVIDIA NeMo 是一个端到端的云原生框架，专为构建、定制和部署生成式 AI 模型而设计。它提供了一整套覆盖模型全生命周期的工具链，从数据处理、大规模分布式训练，到模型微调和高效推理，极大地简化了企业级 GenAI 应用的开发流程。

NVIDIA Megatron-LM

NVIDIA Megatron-LM 是一个开创性的框架，专门用于训练参数量巨大的 Transformer 模型。它通过实现张量并行、流水线并行等关键分布式技术，成功解决了单个 GPU 无法容纳超大模型的难题，是训练千亿乃至万亿参数级别模型的核心引擎。

NVIDIA NeMo Megatron

这是 NVIDIA NeMo 框架与 Megatron-LM 训练引擎的官方结合体，是英伟达当前主推的大语言模型训练解决方案。它将 Megatron-LM 强大的分布式训练能力与 NeMo 框架的工程化易用性相结合，为开发者提供了一个既性能卓越又便于管理的 LLM 开发平台。

NVIDIA FasterTransformer

NVIDIA FasterTransformer 是一个用于 Transformer 模型推理的高性能 C++/CUDA 库，通过深度优化 CUDA 内核来极致压缩推理延迟。尽管它非常成功，但其许多核心优化技术现已被整合进功能更全面的 TensorRT-LLM 库中，后者是英伟达目前的主力推理方案。

NVIDIA TensorRT-LLM

NVIDIA TensorRT-LLM 是一个专为加速和优化大语言模型推理而设计的开源库，是当前 LLM 推理性能优化的首选工具。它集成了 in-flight batching、paged attention、INT4/INT8 量化等最前沿技术，能显著提升 LLM 服务的吞吐量、降低延迟并减少显存占用。

CUDA Graph

CUDA Graph 是 CUDA 平台的一项核心性能优化功能，它允许将一系列 GPU 操作“录制”下来，然后作为一个整体“重放”。对于计算流程固定的 LLM 推理任务，该机制可以消除 CPU 逐个分派任务的开销，从而显著降低延迟，提升最终性能。

NeMo Guardrails

NeMo Guardrails 是一个开源工具包，用于为基于大语言模型的应用程序建立可编程的“安全护栏”。它帮助开发者精确控制模型的对话行为，确保输出内容的主题相关性、安全性和事实准确性，有效防止模型产生不当或偏离轨道的回答。

推理与部署 (Inference & Deployment)

NVIDIA TensorRT

NVIDIA TensorRT 是一个用于高性能深度学习推理的软件开发工具包（SDK），其核心作用是深度优化已经训练好的模型。它通过层融合、精度量化和内核自动调整等技术，为模型创建一个在特定 GPU 上运行速度最快、效率最高的运行时引擎，是实现极致推理性能的关键工具。

NVIDIA Triton Inference Server

NVIDIA Triton Inference Server 是一款开源的推理服务软件，专为在生产环境中大规模、可靠地部署 AI 模型而设计。它支持来自 TensorRT、PyTorch 等多种框架的模型，并提供动态批处理、模型版本管理等企业级功能，是构建标准化、高吞吐量 AI 推理服务的核心组件。

开发、分析与管理工具 (Development, Profiling & Management Tools)

NVIDIA Nsight Suite

NVIDIA Nsight 是一套功能强大的开发者工具集，旨在帮助开发者调试、分析并优化在 NVIDIA GPU 上运行的应用程序。它包含了多个专业工具，如 Nsight Systems 和 Nsight Compute，分别用于从系统全局和单个内核层面进行性能剖析。

Nsight Systems

Nsight Systems 是 Nsight 套件中的系统级性能分析工具，专注于捕捉和可视化应用程序在 CPU、GPU、操作系统和网络之间的交互。它的主要目标是帮助开发者识别系统层面的瓶颈，例如 CPU 等待 GPU 或数据传输延迟，从而进行宏观层面的优化。

Nsight Compute

Nsight Compute 是 Nsight 套件中用于 CUDA 内核深度剖析的工具，它为开发者提供了关于 GPU 内核执行情况的极其详细的反馈和分析。通过这个工具，开发者可以深入检查单个内核的内存访问模式、计算单元利用率等微观指标，从而进行代码级别的极致优化。

NVIDIA System Management Interface (nvidia-smi)

nvidia-smi 是一个基础且广泛使用的命令行工具，用于实时监控和管理系统中的 NVIDIA GPU 设备。它能够快速查询 GPU 的利用率、显存占用、功耗、温度和正在运行的进程等关键状态信息，是每个 GPU 用户和管理员的必备工具。

NVIDIA DCGM (Data Center GPU Manager)

NVIDIA DCGM 可以被看作是 nvidia-smi 的企业级和集群级版本，它是一个专为数据中心环境设计的 GPU 管理和监控工具套件。相比 nvidia-smi，DCGM 提供了更丰富的性能指标、主动的健康诊断以及与集群管理系统的集成能力，适用于大规模 GPU 集群的运维。

NVIDIA Fabric Manager

NVIDIA Fabric Manager 是一款专门用于配置、监控和管理多 GPU 系统中 NVLink 高速互联结构的软件。它能确保 NVLink 网络（Fabric）在系统启动时被正确初始化并保持健康运行，这对于 DGX 等多 GPU 节点实现最佳通信性能至关重要。

NVIDIA Base Command Manager

NVIDIA Base Command Manager 是一个用于简化 AI 基础设施集群生命周期管理的软件平台。它能够自动化地完成从裸金属服务器的操作系统部署、软件栈安装与更新，到整个集群的健康监控和管理等一系列复杂运维工作。

平台、应用与服务 (*Platforms, Applications & Services*)

NVIDIA AI Enterprise

NVIDIA AI Enterprise 是一个经过企业级认证和支持的端到端 AI 软件平台，它将 NVIDIA 生态中的众多关键软件（如 NeMo、TensorRT、Triton 等）打包在一起。该平台极大地简化了企业在私有云、公有云或混合云环境中部署、管理和扩展 AI 工作负载的复杂性，并提供了稳定性和安全性保障。

NVIDIA NGC Catalog

NVIDIA NGC Catalog 是一个官方的软件和模型中心，可以被看作是 AI 和 HPC 领域的“应用商店”。开发者可以从中轻松获取经过 NVIDIA GPU 优化的软件容器、预训练模型、Helm 图表和行业 SDK，从而大幅缩短开发周期，快速构建高性能应用。

NVIDIA DGX Cloud

NVIDIA DGX Cloud 是一项 AI 超级计算云服务，让企业能够直接租用 NVIDIA DGX 集群的顶级算力，而无需承担自建和运维数据中心的巨大成本与复杂性。这项服务是 NVIDIA 与主流云服务商合作推出的，为大规模 AI 训练和推理任务提供了便捷、弹性的基础设施。

NVIDIA AI Foundation Models

NVIDIA AI Foundation Models 是 NVIDIA 官方提供的一系列高质量的预训练生成式 AI 模型，涵盖了语言、视觉等多个领域。企业可以直接在 NGC 上获取这些模型用于推理任务，或者更普遍地，将它们作为强大的基础，通过微调来快速构建满足自身业务需求的定制化模型。

NVIDIA Riva

NVIDIA Riva 是一个 GPU 加速的软件开发工具包（SDK），专门用于构建实时的多模态对话式 AI 应用。它提供了一整套高性能的 API，涵盖了自动语音识别（ASR）、自然语言处理（NLP）和文本转语音（TTS）等功能，使开发者能够轻松创建流畅、自然的语音 AI 助手和服务。

生态合作伙伴与关键开源技术

PyTorch

PyTorch 是由 Meta AI 主导的开源深度学习框架，凭借其灵活性和易用性，已成为 AI 研究和开发领域的事实标准。它与 NVIDIA 的 CUDA 平台深度集成，能够无缝利用 cuDNN 等加

速库，并且 NVIDIA 最新的硬件特性（如 FP8 混合精度）往往会最先在 PyTorch 中得到支持。

TensorFlow / JAX

这两款由 Google 开发的框架在 AI 生态中扮演着关键角色：TensorFlow 作为一个成熟的端到端平台，在企业生产环境中部署广泛；而 JAX 则以其高性能和函数式编程范式在尖端研究领域备受青睐。两者都将 NVIDIA GPU 作为其高性能计算的基石，深度依赖 CUDA 来实现加速。

FlashAttention

FlashAttention 是一种革命性的高效注意力算法，它通过重构计算过程来优化 GPU 内存的读写方式，避免了在显存中实例化巨大的注意力矩阵。这项技术显著减少了长序列 Transformer 模型在训练和推理过程中的显存占用，并大幅提升了计算速度，现已成为 LLM 领域的标配。

vLLM

vLLM 是一个广受欢迎的开源库，专为实现极致的 LLM 推理吞吐量而设计。其核心创新 PagedAttention 算法，借鉴了操作系统的虚拟内存思想来管理 KV 缓存，有效解决了显存碎片问题，从而实现了更高的 GPU 利用率和更灵活的批处理，是当前部署 LLM 服务的领先优化方案之一。

Kubernetes

Kubernetes 是业界公认的容器编排开源标准，在 AI 领域，它被广泛用于自动化部署、扩展和管理大规模 GPU 集群上的 AI 工作负载。NVIDIA 通过提供 GPU Operator 等关键工具，使得 Kubernetes 能够原生感知并高效调度 GPU 资源，从而将强大的容器管理能力无缝延伸至 AI 基础设施。

InfiniBand

InfiniBand 是一种提供极高带宽和超低延迟的计算机网络通信标准，是构建大规模 AI 超级计算机（如 NVIDIA DGX SuperPOD）的神经中枢。它负责在成百上千个 GPU 节点之间进行高速数据交换，是保证分布式训练任务能够高效并行、克服通信瓶颈的关键技术。