



NVIDIA AI TECHNOLOGY CENTER (NVAITC)

Yang-Hsien Lin (yanghsienl@nvidia.com)

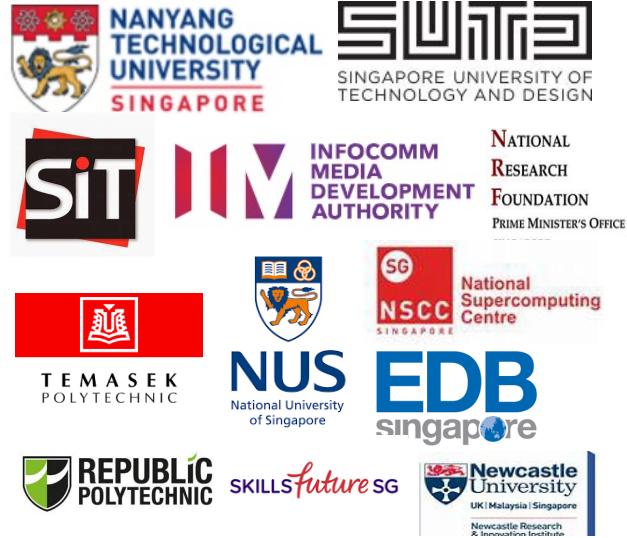
NVIDIA AI TECHNOLOGY CENTER (NVAITC)

Outcome-Focused Collaborations for Research & Talent Development



SELECTED NVAITC COLLABORATORS

Singapore



Europe



Taiwan



Hong Kong



Thailand



Brazil



India



Japan



NALA



Indonesia



NV AI & NV DOMAIN TECH CENTER ENGAGEMENT DRIVERS

Address critical needs across the whole national AI journey

Similar focus across countries

Science, Research and Development

Technology Adoption

Developing an AI workforce

}

Unique Local Implementations

Staffing Model

Government Agencies

Talent Development

Funding Sources

Supercomputing Centers

Level of Research Focus

Interns & Graduate Students

Universities

Collaborate vs Support Balance



Collaboration Model

NVAITC GENERAL RESEARCH GUIDELINES

Project Selection	<ul style="list-style-type: none">• Jointly agreed by Host and NVAITC• NV criteria include technology stack, computing scale, novelty• an agreed-upon “statement of work”
PI Contribution	<ul style="list-style-type: none">• Access to existing compute platform resources• Funding• Researchers
NV Contribution	<ul style="list-style-type: none">• GPU/ML/DL technology / adoption enablement• Staff effort is bounded in both total time and duration• No joint IP• Beyond NVAITC, other NVResearch or teams may or may not collaborate as per any other academic collaboration
Expected Outcome	<ul style="list-style-type: none">• scientific publication, with NVAITC acknowledgement or co-authorship• open-source code release

Collaborations Examples

NCHC X NVIDIA WORKSHOPS/BOOTCAMPS/HACKATHONS

Since 2020

Events

35+

Developers

1600+



RAG-BASED LLM COURSE

Start from Q1 2024



RAG-based LLM Course



Yang-Hsien Lin

We are thrilled to share our upcoming plans to partner with various academic institutions in delivering instruction on RAG-based LLM, while simultaneously offering students the chance to gain hands-on experience with RAG. We've harnessed NVIDIA's resources to create instructional [slides](#), incorporating [Generative AI Examples](#) provided by NVIDIA as our case studies (kindly note that these examples are subject to ongoing updates).

Recognizing that many academic institutions lack the necessary GPU capabilities to run LLM, we've opted for a Cloud-based LLM approach to fulfill our educational and experiential objectives. We warmly invite any innovative solutions or creative course content suggestions you may wish to contribute.

Outlined below are the tasks that students are expected to complete in preparation for the course.

This test can be roughly divided into four steps ([students must achieve the first three points before the class](#)):

1. Apply for an NVIDIA NGC account: The purpose of this is to obtain the NGC API Key and get the required docker image from NGC.
 - o Step 1: Go to the [NVIDIA NGC](#) page, click on "Explore NGC Catalog" in the upper right corner, and then proceed to Step 2.

The screenshot shows the NVIDIA NGC homepage. At the top, there is a navigation bar with links for Shop, Drivers, Support, and a search icon. Below the navigation bar, there is a main menu with links for NGC, Industries, Solutions, Software, and Resource Center. The main content area features a large image of server racks with the text "NVIDIA NGC" overlaid. Below the image, the text reads: "A portal of enterprise services, software, and support for AI, digital twins, and high-performance computing." There is a "Watch Video" button at the bottom of the main content area. A green arrow points to the "Explore NGC Catalog" button in the top right corner of the main menu bar.

Document

- Increase by 1000+ students

- NGC accounts
- NVIDIA AI Endpoint users

- Promote RAG to campus

- National Taiwan University
- Chung Yuan Christian University
- National Tsing Hua University
- China Medical University

- Experience RAG

- Based on our public [Generative AI examples](#)
- NVIDIA Cloud based LLM



Llama 2 13B

Text Generation

Llama 2 is a large language AI model capable of generating text and code in response to prompts.

[View Labels](#) [Learn More](#)

NVAITC

Outcome Highlighted

Journals & Magazines > IEEE Access > Volume: 12 ⓘ

Solar Irradiance Forecasting Using a Hybrid Quantum Neural Network: A Comparison on GPU-Based Workflow Development Platforms

Publisher: IEEE [Cite This](#) [PDF](#)

Ying-Yi Hong ⓘ; Dylan Josh Domingo Lopez ⓘ; Yun-Yuan Wang ⓘ [All Authors](#)

IEEE Access

PAPER • OPEN ACCESS

Validating large-scale quantum machine learning: efficient simulation of quantum support vector machines using tensor networks

IOP Science

Kuan-Cheng Chen*, Tai-Yue Li, Yun-Yuan Wang, Simon See, Chun-Chieh Wang, Robert Wille, Nan-Yow Chen, An-Cheng Yang and Chun-Yu Lin

Published 20 February 2025 • © 2025 The Author(s). Published by IOP Publishing Ltd

Machine Learning: Science and Technology, Volume 6, Number 1

Home > International Conference on Biomedical and Health Informatics 2024 > Conference paper

Pilot Study of Retrieval-Augmented Generation Model in Recommending Traditional Chinese Medicine Formulations

Conference paper | First Online: 28 March 2025
pp 331–341 | Cite this conference paper

Ya-Chuan Chan, Po-Yu Huang, Zhi-Liang Chen, Chih-Nung Wang, Wen-Chen Lin, Jung-Peng Chiu, Yi-Chun Chiu, Yang-Hsien Lin, Eddie T. C. Huang, Simon See & Kang-Ping Lin ⓘ

IFMBE

Maximum Entropy Reinforcement Learning via Energy-Based Normalizing Flow

NeurIPSChen-Hao Chao^{*1,2}Chien Feng^{*1}Wei-Fang Sun²Cheng-Kuang Lee²Simon See²Chun-Yi Lee^{*1}¹ Elsa Lab, National Tsing Hua University, Hsinchu City, Taiwan² NVIDIA AI Technology Center, NVIDIA Corporation, Santa Clara, CA, USA

Harnessing Stable Diffusion: A Leap Forward in Rib Suppression for Chest X-Rays (C-10746)

ECR 2025

was presented within the framework of the Educational and Scientific Programme at

ECR 2025
FEBRUARY 26 – MARCH 2, 2025



GTC 2025 Posters

Driving Smarter Financial Intelligence: The Power of Deeply Trained Domain-Specific Llama Models [P71458]

We present FIN, an answering, stock trading, and fine-tuning system that provides a new way to accelerate genetic diagnosis.

Topic: Generative AI
Industry Segments: I-Chan Chiu, Po Hsuan Tseng, Yu-Rou Chiu, Ting-Chun Lin, Wei-Ting Wang

Pathogenic or Benign? AI-Powered Predictions in Genomic Variant Classification [P73461]

Discover how we're pushing the boundaries of precision. We've accelerated the diagnostic process by providing a new way to answer questions, stock trading, and fine-tuning.

Topic: Data Science
Industry Segments: I-Chan Chiu, Po Hsuan Tseng, Yu-Rou Chiu, Ting-Chun Lin, Wei-Ting Wang

RTeachBot: Transforming Health Education in Radiation Oncology [P73437]

RTeachBot addresses patient education while providing a new way to answer questions, stock trading, and fine-tuning.

Topic: Generative AI
Industry Segments: I-Chan Chiu, Po Hsuan Tseng, Yu-Rou Chiu, Ting-Chun Lin, Wei-Ting Wang

Adaptive Distribution Generator: Powering Split-Step Quantum Walks with CUDA-Q [P72789]

The poster explores how the Adaptive Distribution Generator leverages quantum computing's capabilities to power split-step quantum walks with CUDA-Q.

Topic: Generative AI
Industry Segments: I-Chan Chiu, Po Hsuan Tseng, Yu-Rou Chiu, Ting-Chun Lin, Wei-Ting Wang

CARDIA3D: Comprehensive Anatomical Recognition and Deep-Learning Imaging Algorithm for 3D Cardiac CT [P72960]

Join us as we explore new solutions for multiclass semantic segmentation in 3D cardiac CT images to address the critical lack of models capable of annotating both large structures — like atria and ventricles — and fine details such as coronary arteries, fat, and plaques. By enhancing recognition of multiple classes in high-resolution CT data, we aim to improve feature extraction for heart-focused AI applications and pave the way for future developments, such as training a cardiac-specific language model to interpret 3D CT images. This work is timely and important, filling a significant gap in cardiac imaging tools essential for advancing precision medicine.

Topic: Computer Vision / Video Analytics - Medical Imaging
Industry Segment: Healthcare & Life Sciences

Po Hsuan Tseng, Professor, National Taipei University of Technology
Shu-Yu Hsu, Graduate Researcher, National Taipei University of Technology

Conferences > ICASSP 2025 - 2025 IEEE Inter... ⓘ

3D Gaussian Splatting with Grouped Uncertainty for Unconstrained Images

Publisher: IEEE [Cite This](#) [PDF](#)

Hao-Yu Hou; Chia-Chi Hsu; Yu-Chen Huang; Mu-Yi Shen; Wei-Fang Sun; Cheng Sun [All Authors](#)

ICASSP

Conferences > 2024 International Conference... ⓘ

Quantum-Classical-Quantum Workflow in Quantum-HPC Middleware with GPU Acceleration

Publisher: IEEE [Cite This](#) [PDF](#)

Kuan-Cheng Chen; Xiaoren Li; Xiaotian Xu; Yun-Yuan Wang; Chen-Yu Liu [All Authors](#)

IEEE QCNC

PCIe Bandwidth-Aware Scheduling for Multi-Instance GPUs

Yan-Mei Tang
National Tsing Hua University
Hsinchu, Taiwan
ymtang@lsalab.cs.nthu.edu.tw

Wei-Fang Sun
NVIDIA AI Technology Center
Santa Clara, USA
johnsons@nvidia.com

Ming-Hung Chen
IBM Research
Yorktown Heights, NY, USA
minghungchen@ibm.com

I-Hsin Chung
IBM Research
Yorktown Heights, NY, USA
ihchung@us.ibm.com

Jerry Chou
National Tsing Hua University
Hsinchu, Taiwan
jchou@lsalab.cs.nthu.edu.tw

HPC Asia

Yan-Mei Tang
National Tsing Hua University
Hsinchu, Taiwan
hsutzu.ting@lsalab.cs.nthu.edu.tw



Best Paper Award

PCIe Bandwidth-Aware Scheduling for Multi-Instance GPUs

Yan-Mei Tang, Wei-Fang Sun, Hsu-Tzu Ting, Ming-Hung Chen, I-Hsin Chung, Jerry Chou

HPC ASIA 2025
Chip-based Exploration for HPC



NVIDIA Tech



