

VU VIET ANH

12A/99 Trung Kinh Street, Yen Hoa Ward, Hanoi, Viet Nam
0989894650 | mrvietanh2@gmail.com



OBJECTIVES

Motivated AI Engineer with a solid foundation in Machine Learning, Deep Learning, and Large Language Models (LLMs). Passionate about developing and deploying innovative AI solutions while contributing to cutting-edge research. Eager to bridge the gap between academic research and practical applications in a dynamic environment, with a long-term vision of advancing AI technology and pursuing higher-level studies.

EDUCATION

Master of Computer Science

Hanoi University of Science and Technology, Hanoi, Vietnam

Jun 2025 – Present

Bachelor of Computer Science (Honours)

Hanoi University of Science and Technology, Hanoi, Vietnam

Oct 2020 – Jul 2024

CPA: 3.48/4.0

WORK EXPERIENCE

Prep Education, AI Engineer | Hanoi

October 2025 – Present

Project: Automated Essay Scoring System (AES) for Writing and Speaking

- Developed an AI-powered Automated Essay Scoring (AES) system for evaluating IELTS Writing Task 1 & 2 and Speaking performance with high accuracy.
- AES FOR SPEAKING: Designed and implemented the Speaking scoring system using Machine Learning models integrated with three specialized subsystems: APA (Automated Pronunciation Assessment) for pronunciation accuracy and fluency evaluation, GEC (Grammatical Error Correction) for detecting and correcting grammatical errors, and Dify for analyzing vocabulary richness and lexical diversity. Achieved performance metrics: Accuracy $\pm 0.5 \sim 82\%$, QWK and PCC $\sim 90\%$.
- AES FOR WRITING: Designed and trained deep learning models for cross-prompt essay scoring: PANN (Cross-Prompt Model base Deberta-V3) with 189M parameters, and quantized T5 model optimized to 120M parameters. Achieved strong performance: Accuracy $\pm 0.5 \sim 85\%$ (Wt1) and $\sim 87\%$ (Wt2); Accuracy $\pm 1 \sim 98\%$.
- Implemented end-to-end MLOps pipeline: model training → Docker containerization → GitLab CI/CD automation → ArgoCD deployment to production Kubernetes cluster.
- Optimized model serving with dynamic batching techniques and request queuing to handle high-throughput inference for thousands of concurrent student assessments.
- Deployed the system into production, enabling real-time automated scoring with low latency and high availability.

Project: English Lexical Grammar (ELG) Correction System

- Developed a hybrid English grammar correction system combining rule-based approaches with Large Language Models (LLMs) for enhanced accuracy.
- Implemented comprehensive error detection and correction for multiple linguistic aspects: grammatical errors, vocabulary misuse, spelling mistakes, collocation errors, idiomatic expressions abuse or misuse, lexical resource appropriateness, linking devices, and phrase usage errors, helping students identify and correct writing mistakes effectively.
- Designed rule-based modules for deterministic pattern matching alongside LLM-powered inference for complex grammatical constructs.
- Achieved F0.5 score of approximately 60% on private test sets, balancing precision and recall with emphasis on precision to minimize false corrections.
- Built production-grade serving infrastructure following MLOps best practices: Docker containerization for model packaging, GitLab CI/CD for automated testing and deployment, ArgoCD for Kubernetes orchestration.
- Implemented efficient request batching and caching strategies to optimize LLM inference latency and reduce computational costs while maintaining high throughput.
- Deployed the system to support thousands of students daily in improving their writing skills through real-time grammar and vocabulary correction.

Project: AI Healthcare Assistant for Hospital Management

- Developed an AI-integrated solution for healthcare, deployed within the National Hospital Management Software, particularly adopted by Bach Mai Hospital.
- Utilized advanced Machine Learning and Natural Language Processing (NLP) techniques, including Large Language Models (LLMs), for real-time medical support.
- Designed and implemented modules for patient data analysis, case management assistance, and initial diagnosis support.
- Built interactive AI agents to assist doctors and medical staff with medical history inquiry, prescription, and diagnostics.
- Ensured seamless integration with existing hospital information systems (HIS) for improved operational efficiency and patient care.

Project: Prescription Recommendation and Conflict Detection System

- Developed a predictive system to support doctors in prescribing medications based on patient information.
- Utilized data analysis techniques to preprocess and extract features from electronic health records.
- Applied traditional machine learning classification models (Decision Trees, Random Forests, SVM) and experimented with advanced classifiers for improved accuracy.
- Integrated Large Language Models (LLMs) to evaluate prescriptions, detect duplicate active ingredients, and identify potential drug-drug interactions.
- Delivered a decision-support tool that helps healthcare professionals reduce prescription errors and improve patient safety.

FPT Information System Company, AI Engineer Intern | Hanoi

August 2023 – May 2024

Project: Image Captioning

- Developed an image captioning system using advanced deep learning models.
- Processed image data to generate descriptive captions that accurately reflect the content of the images.
- Specifically utilized state-of-the-art models such as Transformers and Convolutional Neural Networks (CNNs) integrated with attention mechanisms to enhance caption generation.
- This project aims to create an effective solution for automated image description, leveraging cutting-edge techniques in the field of computer vision and natural language processing.

Project: OCR for Financial Report Recognition

- Developed an OCR system to recognize and extract key information from financial reports of banks.
- Applied state-of-the-art OCR models such as PaddleOCR, Tesseract, and TrOCR (Transformer-based OCR) integrated with NLP for data structuring.
- Utilized LayoutLM / LayoutLMv3 for document layout understanding to improve extraction of tables, figures, and key-value pairs.
- Implemented pre-processing for document images (deskewing, binarization, noise reduction) to improve text detection and recognition accuracy.
- Built post-processing pipeline with Regex + Named Entity Recognition (NER) to normalize extracted information.
- Deployed the model as part of an AI pipeline to assist financial data analysts.

Hanoi University of Science and Technology, Computer Science Student | Hanoi

September 2022 – Present

Project: Adapting to climate change by improving extreme weather forecasts (WIDS Datathon 2023 - Kaggle)

- This project involves developing a climate change prediction model based on machine learning models.
- Processed data in a time-series format and applied ensemble techniques to enhance the performance of the machine learning models.
- Specifically utilized the latest prediction models such as Gradient Boosting, Light Gradient Boosting, and CatBoost, and combined them.
- This project was part of a competition held at the end of 2022 on the Kaggle platform aimed at predicting climate change.

Project: Vietnamese Speaker Verification (VLSP 2021 competition)

- This project involves developing a Vietnamese speaker verification system using deep learning models.
- Processed voice samples to create speaker-specific acoustic features for enrollment and verification.
- Specifically utilized models such as ECAPA_TDNN, ECAPA_CNN_TDNN, and RawNet3, and experimented with loss functions like triplet, softmax, and proto.
- This project is part of the VLSP competition focusing on speaker verification, aiming to accurately determine if two voice samples are from the same person.

Project: Aspect-Based Sentiment Analysis (ABSA - NLP LAB)

- Built a Vietnamese Aspect-Based Sentiment Analysis system using SA-Transformer.
- Combined NLP techniques such as POS tagging and Biaffine Parsing to extract aspect-sentiment pairs.
- Collected and annotated domain-specific data to train and evaluate the model.
- Achieved improved accuracy in aspect-level sentiment classification compared to baseline methods.

RESEARCH EXPERIENCE

Student Research with Dr. Kate Han and Dr. Thanh Nguyen – March 2023 to Present Machine Learning & Ensemble Learning

VISTA: Variable-Length Genetic Algorithm and LSTM-Based Surrogate Assisted Ensemble Selection algorithm in Multiple Layers Ensemble System

- This project focused on developing the VISTA method, which integrates a Variable-Length Genetic Algorithm (VLGA) with an LSTM-based surrogate model for optimizing ensemble selection in Multiple Layers Ensemble Systems (MLES).
- Researched existing ensemble learning methods and surrogate-assisted evolutionary algorithms (SAEA).
- Proposed a new method combining VLGA with LSTM to transform variable-length encoding into fixed-size representations for fitness prediction.
- Implemented a multi-layer ensemble system where each layer's ensemble of classifiers (EoC) is trained on both original data and predictions from the previous layer.
- Conducted experiments on 15 popular datasets, demonstrating that VISTA outperforms benchmark algorithms.
- Successfully developed and validated the VISTA method, highlighting its efficacy in improving the performance of MLES.

Imputation Surveys

- Researched Imputation Surveys in UCI Machine Learning Repository (15 datasets).
- Deployed some Imputation methods (SimpleImputer, KNN, MICE, GINN, etc.) for missing data and evaluation in the Classification domain.

Student Research with Assoc.Prof. Pham Van Hai – September 2023 to Present Computer Vision & Deep Learning

Application of deep learning network in underwater image recognition and segmentation (Manuscript under review)

- This project involves applying deep learning networks for underwater image recognition and segmentation.
- Processed underwater images to improve clarity and segment aquatic organisms accurately.
- Specifically utilized ResNet for image denoising, followed by Unet_v2 and V_UNet for segmentation, combining advanced architectures like Unet with Spatial Dropout and Vgg16 for enhanced performance.
- The project aims to address the challenges of underwater image analysis by leveraging modern deep learning techniques to achieve accurate recognition and segmentation of marine life.

Enhancing Facial Expression Recognition with Lightweight Attention-Based CNN (Manuscript under review)

- Designed a lightweight deep learning model for facial expression recognition optimized for deployment on edge devices with limited computational resources.
- Utilized a truncated MobileNetV2 backbone for efficient feature extraction from facial images.
- Proposed a Patch Extraction Block to divide the feature maps into non-overlapping regions, enabling the model to focus on local facial features even under occlusion or head pose variations.
- Integrated an Attention Classifier with self-attention mechanism to enhance feature discrimination and classification accuracy.
- Achieved competitive results on standard datasets (RAF-DB, FER2013, FERPlus) and a custom real-world challenge dataset with occlusion and pose variations.

TECHNICAL SKILLS

Programming Languages	Python, Bash, C++, SQL
Machine Learning Frameworks	Keras, PyTorch
Computer Vision	Object detection, Segmentation, Image processing algorithms
Natural Language Processing	Large Language Models (LLMs), NLP techniques, Question Answering systems
OCR & Document Processing	PaddleOCR, Tesseract, TrOCR, LayoutLM, LayoutLMv3
MLOps & Model Deployment	Docker, Kubernetes, GitLab CI/CD, ArgoCD, Model serving optimization, Dynamic batching, Caching strategies, Production monitoring
Areas of Expertise	Deep Learning, Computer Vision, Large Language Models, Healthcare AI applications, Education

ACHIEVEMENTS

- Second Prize Winner – WIDS DATATHON 2023 (Kaggle, 2023)

Adapting to climate change by improving extreme weather forecasts

LANGUAGES

English	VSTEP B2 (March 25, 2025). Strong ability to read and comprehend English documents. Basic communication skills.
---------	---

PUBLICATIONS

- Kate Han, Truong Thanh Nguyen, **Viet Anh Vu**, Alan Wee-Chung Liew & Tien Thanh Nguyen. 2024. "VISTA: A Variable Length Genetic Algorithm and LSTM-Based Surrogate Assisted Ensemble Selection algorithm in Multiple Layers Ensemble System" – IEEE_SCCI 2024.

CERTIFICATIONS

IBM Python for Data Science Certification completed