**HANOI UNIVERSITY OF SCIENCE AND TECHNOLOGY**

# GRADUATION THESIS

## Thesis title: Improving Vietnamese Question Generation Using Reinforcement Learning

**Nguyen Hoang Dang**

dang.nh204873@sis.hust.edu.vn

**Major: Computer Science**
**Specialization: Data Science**

| | |
|---|---|
| **Supervisor:** | PhD. Tran The Hung |

_____

Signature

**Department:** Computer Science

**School:** School of Information and Communications Technology

**HANOI, 06/2024**

# ACKNOWLEDGMENT

# ABSTRACT

In the field of natural language processing, automatic question generation is reckoned to be a challenge. This is due to the ambiguity where many questions can lead to the same answer. We want our model, after being fine-tuned by on reinforcement learning task, to be able to generate questions that not only resemble the style in the data set but also appear natural and contextual. We introduce a selection method for the preference dataset which is then used for a preference-based reinforcement learning method to align the policy language model. Additionally, we enlarge the training dataset using synthetic data generated from a large language model (LLM) fine-tuned with instruction. Our training method enhances the model stability on different sampling methods, generating better-rated questions in terms of AI evaluation. Quantitatively, our method improves the model performance, showing an increase in Rouge-L (+3.44), BLEU-4 (+3.18), and Win rate (+12.73). The advancement made by our method on the model is also explained in detail.

Student

*(Signature and full name)*

# TABLE OF CONTENTS

# LIST OF FIGURES

# LIST OF TABLES

# LIST OF ABBREVIATIONS

| Abbreviation | Definition |
| --- | --- |
| API | Application Programming Interface |
| AQG | Automatic Question Generation |
| DPO | Direct Preference Optimization |
| LLM | Large Language Model |
| PPO | Proximal Policy Optimization |
| RL | Reinforcement Learning |
| RLAIF | Reinforcement Learning from Artificial Intelligence Feedback |
| RLHF | Reinforcement Learning from Human Feedback |
| SDPO | Self-play Fine-tuning with Direct Preference Optimization |
| SFT | Supervised Fine-tuning |