VIETNAM NATIONAL UNIVERSITY, HANOI
**VNU UNIVERSITY OF ENGINEERING AND TECHNOLOGY**
_____

**SOCIALIST REPUBLIC OF VIETNAM**
**Independence – Freedom – Happiness**
_____

## INFORMATION ON DOCTORAL THESIS

1. Full name : Ngoc-Khuong Nguyen            2. Sex: Male

3. Date of birth: 22/10/1984            4. Place of birth: Haiphong

5. Admission decision number: 1006/QĐ-CTSV .. Dated: 07/12/2015

6. Changes in academic process:

- Adjust the Supervisors for PhD student Nguyen Ngoc Khuong according to Decision No. 467/QD-DT dated May 18, 2017, by the Rector of the VNU University of Engineering and Technology.

- Change the title of the doctoral thesis for PhD student Nguyen Ngoc Khuong according to Decision No. 1023/QD-DT dated October 24, 2017, by the Rector of the VNU University of Engineering and Technology.

7. Official thesis title: Studying sequence-to-sequence models using deep learning and their applications in natural language processing.

8. Major: Computer Science            9. Code: 9480101.01

10. Supervisors:

- Nguyen Viet Ha, *Assoc. Prof., PhD,* VNU University of Engineering and Technology.

- Le Anh Cuong, *Assoc. Prof., PhD,* Ton Duc Thang University.

11. Summary of the **new findings** of the thesis:

In this thesis, we focus on research and development of Seq2Seq deep learning models for natural language processing, specifically addressing the tasks of paraphrase generation and text summarization. The Seq2seq model, in its general form, consists of two components: an encoder to encode information from the input sequence and a decoder to generate responses. The thesis concentrates on enhancing the existing machine learning model structures, with a primary emphasis on improving the encoding stage of the input sequence to adequately represent hierarchical semantic relationships within the input text. The second issue the dissertation addresses is the improvement of the decoding component to generate output sequences that align with the objectives of each specific task, meeting content requirements and imposed constraints. Thus, the thesis is dedicated to the enhancement of Seq2seq machine learning models by effectively encoding information

from the input text and generating output text that satisfies the content and other constraints. In addressing these challenges, the thesis has obtained results and contributions including:

- Proposed a hierarchical text representation method in the sequence-to-sequence model for the abstractive summarization task. The research focuses on integrating hidden representation layers with hierarchical representation levels of structural components in the text, from low to high, to enhance the understanding and modeling of relationships between structural components in the input text. The proposed model with experimental results have been published in the proceedings of the International Conference on Knowledge and Systems Engineering (KSE), 2021.

- Proposed improvement in attention mechanisms within the sequence-to-sequence model. The first proposal is an attention penalty coefficioent mechanism that enhances the attention scoring model by introducing penalties to better differentiate the roles of words in the input text for text paraphrasing. The proposed model with experimental results have been published in the proceedings of the International Symposium on Integrated Uncertainty in Knowledge Modelling and Decision Making (IUKM), 2018. The second proposal is a hierarchical attention mechanism based on the hierarchical structure of the input text in the sequence-to-sequence model for text paraphrasing. Sentence-level attention and word-level attention mechanisms make the representation layers more comprehensive in capturing relationships between structural components in the text. The proposed model with experimental results have been published in the proceedings of the 12th Multi-disciplinary International Conference on Artificial Intelligence (MIWAI), 2018. The third proposal is a model that combines local attention and global attention mechanisms by defining and leveraging attention gates, aiming to reduce computational complexity while enhancing the ability to understand the input text and generate output text flexibly and accurately. The proposed model with experimental results have been published in the proceedings of the 5th Asia Pacific Information Technology Conference (APIT), 2023.

- Proposed an End-to-End Seq2Seq model for the task of generating length-constrained abstractive summarization. In the encoding phase, length information is encoded as a vector, combining the word embedding, positional embedding, and length information embedding. During the decoding phase, the desired length is utilized to supplement information for the output words, similar to the words in the input, and a head attention mechanism in the Transformer architecture is employed to effectively control the length of the output sequence. The proposed model with experimental results have been published in International Journal of Intelligent Automation & Soft Computing (IASC), 2023.

12. Further research directions, if any:

At present, research in natural language processing has undergone significant transformations with the emergence of Large Language Models (LLMs). Our subsequent research endeavors will be oriented towards advancing the proposed techniques in the thesis within the context of LLMs' development. We aim to continue refining the results outlined in the thesis by exploring methods that leverage hierarchical text structures in conjunction with Reinforcement Learning with Human Feedback (RLHF) approaches, particularly in resource-constrained environments. Additionally, we intend to extend the proposed model to incorporate constraints on text generation, including length constraints as well as other constraints such as morphological, lexical, and content-related constraints. This research direction is geared towards pushing the boundaries of text generation models and addressing various constraints to enhance the practical applicability of the proposed methods.

13. Thesis-related publications

- [1] Khuong Nguyen-Ngoc, Anh-Cuong Le and Viet-Ha Nguyen. "A Hierarchical Conditional Attention-based Neural Networks for Paraphrase Generation", the 12th Multi-disciplinary International Conference on Artificial Intelligence (MIWAI), 2018, pp 161 - 174, DOI:10.1007978-3-030-03014-8_14.

- [2] Khuong Nguyen-Ngoc, Anh-Cuong Le and Viet-Ha Nguyen. "An Attention-based Long-Short-Term-Memory Model for Paraphrase Generation", the 6th International Symposium on Integrated Uncertainty in Knowledge Modelling and Decision Making (IUKM), 2018, pp.166-178, DOI:10.1007978-3-319-75429-1_14.

- [3] Khuong Nguyen-Ngoc, Anh-Cuong Le and Viet-Ha Nguyen. "A Hierarchical Encoder-Decoder Long Short-Term Memory Model for Abstractive Summarization", 13th International Conference on Knowledge and Systems Engineering (KSE), 2021, pp 281-286.

- [4] Ngoc-Khuong Nguyen, Anh-Cuong Le and Viet-Ha Nguyen. "A Local Attention-based Neural Networks for Abstractive Text Summarization", 5th Asia Pacific Information Technology Conference (APIT), 2023, pp 152-159.

- [5] Ngoc-Khuong Nguyen, Viet-Ha Nguyen, Dac-Nhuong Le and Anh-Cuong Le. "A Method of Integrating Length Constraints into Encoder-Decoder Transformer for Abstractive Text Summarization", Jounal of Intelligent Automation & Soft Computing (IASC), 2023, doi:10.32604/iasc.2023.037083.

Date: …………………..                        Date: ……………………..
Signature: ………………                   Signature: …………………
Full name: Nguyen Viet Ha                   Full name: Nguyen Ngoc Khuong