

**BỘ GIÁO DỤC VÀ ĐÀO TẠO**

**VIỆN HÀN LÂM KHOA HỌC**

**VÀ CÔNG NGHỆ VIỆT NAM**

**HỌC VIỆN KHOA HỌC VÀ CÔNG NGHỆ**



**NGUYỄN THỊ THU HIỀN**

**NGHIÊN CỨU PHƯƠNG PHÁP CHUẨN HOÁ VĂN BẢN  
VÀ NHẬN DẠNG THỰC THỂ ĐỊNH DANH  
TRONG NHẬN DẠNG TIẾNG NÓI TIẾNG VIỆT**

**LUẬN ÁN TIẾN SĨ NGÀNH MÁY TÍNH**

**HÀ NỘI - 2023**

BỘ GIÁO DỤC VÀ ĐÀO TẠO

VIỆN HÀN LÂM KHOA HỌC

VÀ CÔNG NGHỆ VIỆT NAM

HỌC VIỆN KHOA HỌC VÀ CÔNG NGHỆ

NGUYỄN THỊ THU HIỀN

NGHIÊN CỨU PHƯƠNG PHÁP CHUẨN HOÁ VĂN BẢN  
VÀ NHẬN DẠNG THỰC THỂ ĐỊNH DANH  
TRONG NHẬN DẠNG TIẾNG NÓI TIẾNG VIỆT

LUẬN ÁN TIẾN SĨ NGÀNH MÁY TÍNH

Chuyên ngành: Hệ thống thông tin

Mã số: 9 48 01 04

Xác nhận của Học viện

Người hướng dẫn 1

Người hướng dẫn 2

Khoa học và Công nghệ

(Ký, ghi rõ họ tên)

(Ký, ghi rõ họ tên)

## **LỜI CAM ĐOAN**

Tôi xin cam đoan đây là công trình nghiên cứu của riêng tôi. Các kết quả được viết chung với các tác giả khác đều được sự đồng ý của các đồng tác giả trước khi đưa vào luận án. Các kết quả nêu trong luận án là trung thực và chưa từng được công bố trong các công trình nào khác.

Tác giả

Nguyễn Thị Thu Hiền

## LỜI CẢM ƠN

Luận án của tác giả được thực hiện tại Học viện Khoa học và Công nghệ - Viện Hàn lâm Khoa học và Công nghệ Việt Nam, dưới sự hướng dẫn tận tình của PGS.TS. Lương Chi Mai và TS. Nguyễn Thị Minh Huyền. Tôi xin được bày tỏ lòng biết ơn sâu sắc đến hai Cô về những định hướng nghiên cứu, sự động viên và hướng dẫn tận tình giúp tôi vượt qua những khó khăn để hoàn thành luận án này.

Tôi cũng xin gửi lời cảm ơn chân thành đến các nhà khoa học, các đồng tác giả của các công trình nghiên cứu đã được trích dẫn trong luận án. Đây là những tư liệu quý báu có liên quan giúp tôi hoàn thành luận án.

Tôi xin chân thành cảm ơn đến Ban lãnh đạo Học viện Khoa học và Công nghệ, Viện Công nghệ Thông tin đã tạo điều kiện thuận lợi cho tôi trong quá trình học tập, nghiên cứu.

Tôi xin chân thành cảm ơn Ban giám hiệu trường Đại học Sư phạm - ĐH Thái Nguyên, Khoa Toán, Bộ môn Khoa học máy tính - Hệ thống thông tin và các đồng nghiệp đã giúp đỡ và tạo điều kiện thuận lợi để tôi có thể thực hiện kế hoạch nghiên cứu, hoàn thành luận án.

Tôi xin được bày tỏ tình cảm và lòng biết ơn vô hạn tới những người thân trong Gia đình, những người luôn dành cho tôi sự động viên, khích lệ, sẻ chia, giúp đỡ trong những lúc khó khăn.

Tác giả

Nguyễn Thị Thu Hiền

**MỤC LỤC**

	<b>Trang</b>
<b>LỜI CAM ĐOAN .....</b>	<b>i</b>
<b>LỜI CẢM ƠN .....</b>	<b>ii</b>
<b>MỤC LỤC .....</b>	<b>iii</b>
<b>DANH MỤC TỪ VIẾT TẮT.....</b>	<b>v</b>
<b>DANH MỤC BẢNG BIỂU .....</b>	<b>vii</b>
<b>DANH MỤC HÌNH VẼ .....</b>	<b>viii</b>
<b>MỞ ĐẦU .....</b>	<b>1</b>
<b>CHƯƠNG 1: TỔNG QUAN VẤN ĐỀ NGHIÊN CỨU.....</b>	<b>7</b>
1.1. Xử lý ngôn ngữ tự nhiên.....	7
1.2. Nhận dạng tiếng nói.....	11
1.3. Chuẩn hóa văn bản .....	16
1.4. Nhận dạng thực thể định danh.....	24
1.5. Tổng quan về dữ liệu .....	34
1.6. Kết luận Chương 1.....	36
<b>CHƯƠNG 2: KIẾN THỨC CƠ SỞ.....</b>	<b>37</b>
2.1. Mô hình xử lý chuỗi .....	37
2.2. Mô hình biểu diễn từ .....	44
2.3. Mô hình gán nhãn chuỗi .....	50
2.4. Học đa tác vụ .....	53
2.5. Kết luận chương 2 .....	56
<b>CHƯƠNG 3: CHUẨN HÓA VĂN BẢN ĐẦU RA CỦA HỆ THỐNG NHẬN DẠNG TIẾNG NÓI TIẾNG VIỆT.....</b>	<b>57</b>
3.1. Bài toán.....	57
3.2. Xây dựng dữ liệu .....	58
3.3. Kiến trúc mô hình.....	60
3.4. Kết quả thực nghiệm.....	68
3.5. Kết luận Chương 3.....	73

<b>CHƯƠNG 4: NHẬN DẠNG THỰC THỂ ĐỊNH DANH CHO VĂN BẢN ĐÀU RA CỦA HỆ THỐNG NHẬN DẠNG TIẾNG NÓI TIẾNG VIỆT....</b>	<b>75</b>
4.1. Bài toán.....	75
4.2. Tổng quan dữ liệu.....	76
4.3. Nhận dạng thực thể định danh theo hướng tiếp cận Đường ống.....	77
4.4. Nhận dạng thực thể định danh theo hướng tiếp cận E2E.....	87
4.5. Kết luận Chương 4.....	98
<b>KẾT LUẬN .....</b>	<b>99</b>
<b>DANH MỤC CÔNG TRÌNH CỦA TÁC GIẢ .....</b>	<b>101</b>
<b>TÀI LIỆU THAM KHẢO .....</b>	<b>103</b>

**DANH MỤC TỪ VIẾT TẮT**

<b>STT</b>	<b>Từ viết tắt</b>	<b>Từ tiếng Anh</b>	<b>Ý nghĩa tiếng Việt</b>
1	ASR	Automatic Speech Recognition	Nhận dạng tiếng nói tự động
2	BERT	Bidirectional Encoder Representations from Transformers	Mã hóa biểu diễn hai chiều dựa trên Transformers
3	BiLSTM	Bidirectional Long Short Term Memory	Mô hình bộ nhớ ngắn-dài hạn hai chiều
4	BPE	Byte-Pair-Encoding	Mã hoá cặp byte
5	CaPu	Recovering Capitalization and Punctuation model	Mô hình khôi phục dấu câu và chữ hoa
6	CBOW	Continuous Bag of Words	Mô hình nhúng từ “Túi từ liên tục”
7	CNN	Convolutional Neural Network	Mạng nơ-ron tích chập
8	CRF	Conditional Random Fields	Trường ngẫu nhiên có điều kiện
9	DL	Deep Learning	Học sâu
10	DNN	Deep Neural Networks	Mạng nơ-ron sâu
11	ELMO	Embeddings from Language Model	Nhúng từ từ mô hình ngôn ngữ
12	E2E	End-to-End	Mô hình đầu - cuối
13	GloVe	Global Vectors for Word Representation	Mô hình nhúng từ dựa trên biểu diễn từ
14	GRU	Gated Recurrent Unit	Mạng hồi tiếp có cổng

15	GPT	Generative pre-trained transformer	Mô hình biến đổi được huấn luyện trước
16	HMM	Hidden Markov Model	Mô hình Markov ẩn
17	LM	Language Model	Mô hình ngôn ngữ
18	LSTM	Long Short Term Memory	Mô hình bộ nhớ ngắn-dài hạn
19	ME	Maximum Entropy	Mô hình Entropy cực đại
20	MEMM	Maximum Entropy Markov Model	Mô hình Markov Entropy cực đại
21	MTL	Multi-Task Learning	Học đa tác vụ
22	NER	Named Entity Recognition	Nhận dạng thực thể định danh
23	OOV	Out-of-Vocabulary	Từ nằm ngoài từ điển
24	RNN	Recurrent Neural Network	Mạng nơ-ron hồi quy
25	Seq2seq	Sequence-to-Sequence	Mô hình ánh xạ từ chuỗi sang chuỗi
26	SLU	Spoken Language Understanding	Hiểu ngôn ngữ nói
27	SVM	Support Véc-tơ Machine	Máy véc-tơ hỗ trợ
28	VLSP	Vietnamese Language and Speech Processing	Hội thảo xử lý ngôn ngữ và tiếng nói tiếng Việt
29	XLNNTN		Xử lý ngôn ngữ tự nhiên
30	TTS	Text To Speech	Hệ thống chuyển văn bản sang tiếng nói
31	WER	Word Error Rate	Tỉ lệ lỗi từ



## DANH MỤC BẢNG BIỂU

Bảng 1.1: Điểm khác biệt giữa văn bản đầu ra ASR và văn bản viết dạng chuẩn .....	13
Bảng 1.2: Tỷ lệ lỗi từ của một số hệ thống nhận dạng tiếng nói tiếng Việt....	15
Bảng 3.1: Thông tin bộ dữ liệu .....	59
Bảng 3.2: Số lượng tham số của các mô hình.....	69
Bảng 3.3: Các tham số huấn luyện mô hình .....	69
Bảng 3.4: So sánh kết quả mô hình Transformer Encoder - CRF khi áp dụng và không áp dụng hợp nhất chồng lán .....	71
Bảng 3.5: So sánh tốc độ xử lý (tokens/second) .....	73
Bảng 4.1: Tham số cấu trúc và huấn luyện mô hình ViBERT .....	81
Bảng 4.2: Thống kê bộ dữ liệu NER của VLSP 2018 .....	83
Bảng 4.3: Đánh giá các mô hình NER dựa trên bộ dữ liệu NER của VLSP 2018.....	85
Bảng 4.4: Đánh giá mô hình NER đề xuất theo cách tiếp cận đường ống với các kiểu văn bản đầu vào khác nhau .....	85
Bảng 4.5: Tỷ lệ lỗi của TTS-ASR và REC-ASR trên dữ liệu kiểu số, dữ liệu ngoại lai và các lỗi khác .....	95
Bảng 4.6: Đánh giá mô hình NER đề xuất theo cách tiếp cận E2E với các kiểu văn bản đầu vào khác nhau .....	97
Bảng 4.7: So sánh mô hình E2E với mô hình đường ống.....	97

## DANH MỤC HÌNH VẼ

Hình 1.1: Minh họa các vấn đề cần thực hiện để tăng chất lượng văn bản đầu ra của ASR .....	14
Hình 1.2: Mô hình NER dựa trên học sâu.....	30
Hình 2.1: Mô hình Transformer [34] .....	40
Hình 2.2: Minh họa hoạt động của CBOW và Ship-Gram.....	45
Hình 2.3: Tổng thể quy trình tiền huấn luyện và tinh chỉnh cho BERT [35].	48
Hình 2.4: Tinh chỉnh BERT cho nhiệm vụ NER [35] .....	49
Hình 2.5: Mô hình Conditional Random Fields.....	51
Hình 2.6: Mô hình phương pháp chia sẻ tham số cứng .....	54
Hình 2.7: Mô hình phương pháp chia sẻ tham số mềm .....	55
Hình 3.1: Minh họa đầu vào, đầu ra của khôi phục dấu câu, chữ hoa đối với văn bản đầu ra ASR.....	58
Hình 3.2: Kiến trúc mô hình .....	60
Hình 3.3: Mô hình xử lý chuỗi đầu vào, đầu ra thông thường.....	61
Hình 3.4: Đề xuất mô hình phân chia/hợp nhất đoạn chồng lấn.....	62
Hình 3.5: Mô tả phân chia đoạn chồng lấn .....	63
Hình 3.6: Ví dụ phân chia đoạn chồng lấn với $l = 10$ và $k = 5$ .....	63
Hình 3.7: Mô tả cách ghép nối .....	64
Hình 3.8: Hợp nhất các đoạn chồng chéo dựa trên tham số $c$ .....	65
Hình 3.9: Mô hình CaPu đề xuất cho văn bản đầu ra của ASR tiếng Việt.....	66
Hình 3.10: Mô tả đầu ra nhận dạng dạng văn bản và dạng nhãn.....	68
Hình 3.11: Kết quả của các mô hình sử dụng và không sử dụng hợp nhất đoạn chồng lấn .....	70
Hình 3.12: Kết quả của các mô hình với đầu ra là dạng văn bản hoặc dạng nhãn .....	71
Hình 3.13: Ma trận lỗi cho mô hình Transformer Encoder - CRF .....	72
Hình 4.1: Mô tả kiến trúc NER tổng quát theo cách tiếp cận đường ống.....	78