

B GIÁO D C VÀ ÀO T O
TR NG I H C CÔNG NGH TP.HCM



LÊ MINH HI U

PHÂN O N T TI NG VI T

LU N V N TH C S

Chuyên ngành: Công ngh thông tin

Mã s ngành:60480201

TP. H CHÍ MINH, tháng 01 n m 2015

BỘ GIÁO DỤC VÀ ĐÀO TẠO
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ TP. HCM



LÊ MINH HIU

PHÂN ÔN TẬP TIN HỌC VI T

LUẬN VĂN THẠC S

Chuyên ngành: Công nghệ thông tin

Mã số ngành: 60480201

CÁN BỘ HƯỚNG DẪN KHOA HỌC: TS. NGUYỄN CHÍ HIU

CÔNG TRÌNH C HOÀN THÀNH T I
TR NG I H C CÔNG NGH TP. HCM

Cán b h ng d n khoa h c: **TS. NGUY N CHÍ HI U**
(*Ghi rõ h , tên, h c hàm, h c v và ch ký*)

TS. Nguy n Chí Hi u

Lu n v n Th c s c b o v t i Tr ng i h c Công ngh TP. HCM
ngày 06 tháng 02 n m 2015

Thành ph n H i ng á nh giá Lu n v n Th c s g m:
(*Ghi rõ h , tên, h c hàm, h c v c a H i ng ch m b o v Lu n v n Th c s*)

TT	H và tên	Ch c danh H i ng
1	PGS.TS. Lê Hoài B c	Ch t ch
2	PGS.TS. Qu n Thành Th	Ph n bi n 1
3	TS. Võ ình B y	Ph n bi n 2
4	TS. L Nh t Vinh	y viên
5	TS. Cao Tùng Anh	y viên, Th ký

Xác nh n c a Ch t ch H i ng á nh giá Lu n v n sau khi Lu n v n ã c
s a ch a (n u có).

Chủ tịch Hội đồng đánh giá LV

PGS.TS. Lê Hoài B c

TP. HCM, ngày... tháng ... năm 20...

NHIỆM VỤ LUẬN VĂN THẠC SĨ

Họ tên học viên: **LÊ MINH HI** Giới tính: **NAM**

Ngày, tháng, năm sinh: **20/10/1985** Nơi sinh: **GIA LAI**

Chuyên ngành: **CÔNG NGHỆ THÔNG TIN** MSV: **1241860004**

I- Tên tài: Phân loại tiếng Việt

II- Nhiệm vụ và nội dung:

- Nghiên cứu các lý thuyết và xử lý ngôn ngữ tự nhiên.
- Khảo sát các nghiên cứu liên quan.
- Xây dựng mô hình phân loại tiếng Việt.
- Chạy thực nghiệm và đánh giá kết quả.

III- Ngày giao nhiệm vụ : (Ngày bắt đầu thực hiện LV ghi trong Q giao tài)

.....

IV- Ngày hoàn thành nhiệm vụ : (Ngày báo LV)

.....

V- Cán bộ hướng dẫn: TS. NGUYỄN CHÍ HIU

CÁN BỘ HƯỚNG DẪN KHOA QUẢN LÝ CHUYÊN NGÀNH

(Họ tên và chữ ký) (Họ tên và chữ ký)

TS. Nguyễn Chí Hi u

L I CAM OAN

Tôi xin cam oan đây là công trình nghiên c u c a riêng tôi. Các s li u, k t qu nêu trong Lu n v n là trung th c và ch a t ng c ai công b trong b t k công trình nào khác.

Tôi xin cam oan r ng m i s giúp cho vi c th c hi n Lu n v n này ã c c m n và các thông tin trích d n trong Lu n v n ã c ch rõ ngu n g c.

H c viên th c hi n Lu n v n

(Ký và ghi rõ h tên)

Lê Minh Hi u

L I C M N

V i t t c t m lòng, tôi xin g i l i c m n sâu s c nh t n th y giáo TS.Nguy n Chí Hi u – ng i th y ã t n tình h ng d n, ch b o và t o nh ng i u ki n t t nh t giúp tôi hoàn thành lu n v n này.

ng th i tôi xin g i l i c m n chân thành n toàn th quý th y cô tr ng i h c Công ngh Thành ph H Chí Minh ã trang b cho tôi nh ng ki n th c trong h c t p và nghiên c u khoa h c.

Tôi c ng xin chân thành c m n các thành viên trong tài “Nghiên c u phát tri n m t s s n ph m thi t y u v x lý ti ng nói và v n b n ti ng Vi t”, mã s KC01.01/06-10 ã cho phép tôi s d ng m t s d li u c a VietTreebank và Vietnamese Lexicon trong quá trình th c nghi m.

Cu i cùng, tôi xin g i l i c m n n gia ình, b n bè và các ng nghi p ã luôn ng viên và cho tôi nh ng l i khuyên b ích trong su t quá trình th c hi n lu n v n này.

Tp.H Chí Minh, tháng 01 n m 2015

Lê Minh Hi u

TÓM TẮT

Từ khóa

Phân o n t , phân gi i nh p nh ng, nh n d ng danh t riêng, thông tin t ng h .

Tóm t t

Không gi ng nh ti ng Anh, phân o n t trong ngôn ng ti ng Vi t, c ng nh h u h t các ngôn ng châu Á, là m t công vi c h t s c ph c t p. Vì b n thân ngôn ng không có nh ng d u hi u rõ ràng phân cách các t v i nhau, ch ng h n nh kho ng tr ng. ã có r t nhi u nghiên c u v i nhi u h ng ti p c n khác nhau v công vi c này. Tuy nhiên theo kh o sát, i a s các nghiên c u u xu t phát t ba h ng ti p c n chính: h ng ti p c n d a trên t i n, h ng ti p c n d a trên mô hình th ng kê và h ng ti p c n lai.

R t nhi u nghiên c u ã ch n h ng ti p c n d a trên t i n vì tính n gi n c a nó. H ng ti p c n này th ng s d ng t i n k t h p v i m t s thu t toán so kh p nh : Maximum matching (MM), Longest matching (LM), v.v... phân o n t . Tuy nhiên h ng ti p c n này th ng gây ra nhi u nh p nh ng khi phân o n và không th phân o n úng cho các t không có trong t i n.

H ng ti p c n d a trên th ng kê c n m t kho ng li u l n, ã tr i qua quá trình hu n luy n, k t h p v i các thu t toán th ng kê phân o n t . Có th k n m t s mô hình theo h ng ti p c n này nh : mô hình th ng kê N-gram, mô hình c c i hóa Entropy (ME), mô hình Conditional Random Fields (CRFs), mô hình cây quy t nh. u i m c a h ng ti p c n này là có th phát hi n c các t không có trong t i n và h n ch c nh p nh ng.

H ng ti p c n d a trên mô hình lai k t h p nhi u ph ng pháp khác nhau phân o n t . M t s mô hình phân o n t theo mô hình lai có th k n nh : mô hình so kh p Maximum matching k t h p v i SVMs, mô hình phân o n t s d ng WFST và m ng Neural, mô hình s d ng thu t toán Maximum matching và N-gram, mô hình k t h p CRFs và SVMs. H ng ti p c n này th ng ph c t p nh ng mang l i hi u qu cao.

Trong lu n v n này chúng tôi xu t m t mô hình phân o n t d a trên mô hình lai. Mô hình c a chúng tôi s d ng b n lu t phân gi i nh p nh ng c a h th ng MMSeg k t h p v i t i n, thông tin hu n luy n N-gram, thông tin h t ng và các bi u th c chính quy.

Th c nghi m trên v n b n g m 10,000 câu trích t VietTreebank cho k t qu F-measure t 91.74%.

ABSTRACT

Keywords

Vietnamese word segment, disambiguity, proper nouns identification, mutual information.

Abstract

Unlike in English, word segmentation in Vietnamese, as well as in many other Asian languages, is more complex because the language does not have any explicit word boundary delimiters, such as a space, to separate between each word. Many researchers with many approaches for the word segmentation task. However, these approaches can be classified into 3 major categories: dictionary-based, statistics-based and hybrid-based.

Most studies use dictionary-based approaches because of their simplicity. This approach type use dictionaries with matching methods as Maximum matching (MM), Longest matching (LM), ect for the word segmentation. However, most of the dictionary based approaches often get many ambiguous cases and can not detects new words.

Statistical approaches need a very large annotated training corpus for word segmentation. Some of studies based on this approaches are N-gram Language Model, Maximum Entropy (ME), Conditional Random Fields (CRFs), Decision Tree. This approach is usefull for detects new words and disambiguous.

Hybrid approaches combine different approaches to make use of individual advantages and overcome disadvantages. Some models are combination of Maximum matching and SVMs, WFST and Neural network, Maximum matching and Ngram language model, CRFs and SVMs. This approache are often complex however it give a high accuracy.

In this thesis, we propose a hybrid method for Vietnamese word segmentation. Our approach is base on four MMSegdisambiguity rules, dictionaries, ngram trained corpus, mutual information and regular expression.

Experiment on 10,000 sentences of VietTreebank corpus gives a result with an F-measure of 91.74%.

M C L C

L I C A M O A N	i
L I C M N	ii
T Ó M T T	iii
ABSTRACT	iv
M C L C	v
D A N H M C C Á C T V I T T T	vii
D A N H M C C Á C B Ñ G	viii
D A N H M C C Á C H Ì N H V	ix
G I I T H I U	1
1. t v n	1
2. Lý do ch n tài	2
3. M c tiêu và ph m vi nghiên c u	2
4. B c c c a l u n v n	3
C H Ñ G 1. T Ñ G Q U A N	4
C H Ñ G 2. C S L Ý T H U Ý T	7
2.1 C s lý thuy t v ngôn ng	7
2.1.1 Phân lo i ngôn ng	7
2.1.2 n v ch y u c a ngôn ng	10
2.1.3 C u trúc c a n v t ti ng Vi t	15
2.1.4 T v ng ti ng Vi t	19
2.1.5 V n nh p nh ng ngh a c a t	21
2.2 C s lý thuy t v ngôn ng h c th ng kê	24
2.2.1 T ng quan v ngôn ng h c th ng kê	24
2.2.2 M t s lý thuy t xác su t th ng kê trong x lý ngôn ng	25
C H Ñ G 3. G I I T H I U M Ô H Ì N H M M S E G	33
3.1 T ng quan v M M S e g	33

3.2	Áp dụng MMSEG vào tiếng Việt.....	35
3.3	Đánh giá MMSEG trên ngôn ngữ tiếng Việt.....	37
CHƯƠNG 4.	MÔ HÌNH XUẤT.....	39
4.1	Mô hình phân loại.....	39
4.2	Thiết kế thuật toán.....	40
4.2.1	Thuật toán xử lý văn bản.....	40
4.2.2	Thuật toán phân loại.....	42
4.2	Thời gian và khối lượng.....	43
4.3	Thử nghiệm.....	47
CHƯƠNG 5.	KẾT LUẬN.....	51
5.1	Nhìn xét chung.....	51
5.2	Kết quả thực nghiệm.....	52
5.3	Hiện trạng của tài liệu.....	52
5.4	Hướng phát triển của tài liệu.....	53
TÀI LIỆU THAM KHẢO.....		54

DANH MỤC CÁC T VI T T T

STT	T vi t t t	Di n gi i ti ng Anh	Di n gi i ti ng Vi t
1	ATS	Automatic Text Summarization	Tóm l c v n b n
2	CRFs	Conditional Random Fields	H c máy CRFs
3	DCB	Dictionary Based	D a trên t i n
4	DT	Decision Tree	Cây quy t nh
5	HMM	Hidden Markov Model	Mô hình Markov n
6	IR-IE	Information Retrieval and Extraction	Truy v n và khai thác thông tin
7	LM	Longest Matching	So kh p dài nh t
8	ME	Maximum Entropy	C c i hóa Entroy
9	MLB	Machine Learning Based	D a trên h c máy
10	MM	Maximal Matching	So kh p c c i
11	MT	Machine Translation	D ch máy
12	Q&A	Question and Answer	H th ng h i áp
13	SR	Speech Recognition	Nh n d ng ti ng nói
14	SVMs	SVMs	Máy h c vect h tr
15	WS	Word Segmentation	Phân o n t

DANH MỤC CÁC BẢNG

Bảng 2.1 Bảng minh họa ngôn ngữ hòa kết.....	13
Bảng 2.2 Bảng minh họa mọt trong tiếng Tschinuk.....	13
Bảng 2.3 Bảng phoneme.....	16
Bảng 2.4 Bảng phoneme cu i và bán nguyên âm.....	16
Bảng 2.5 Bảng nguyên âm.....	17
Bảng 2.6 Bảng liệt kê các ký hiệu thống nhất trong HMM.....	34
Bảng 3.1 Bảng liệt kê kết quả thực nghiệm MMSEG trên ngữ liệu tiếng Việt.....	40
Bảng 4.1 Danh sách mọt stop word trong tiếng Việt.....	45
Bảng 4.2 Bảng liệt kê số lượng tên cá nhân danh từ riêng.....	48
Bảng 4.3 Bảng liệt kê số lượng bài báo phục vụ cho việc huấn luyện mô hình.....	49
Bảng 4.4 Bảng liệt kê kết quả thực nghiệm của VNS so với MMS.....	51

DANH MỤC CÁC HÌNH V

Hình 2.1 Hình minh họa các nhân vật yêu cầu ngôn ngữ	15
Hình 2.2 Sơ đồ 3 tiêu chí khu biệt cho sáu âm vị thanh điệu.....	17
Hình 2.3 Hình minh họa biểu đồ thanh điệu.....	17
Hình 2.4 Hình minh họa bảng trật tự thanh điệu.....	18
Hình 4.1 Hình minh họa mô hình phân loại tiếng Việt (VNS).....	42
Hình 4.2 Hình minh họa cấu trúc từ tiếng Việt.....	47
Hình 4.3 Hình minh họa từ vựng riêng.....	48
Hình 4.4 Hình minh họa kết quả huấn luyện Uni-Gram.....	49
Hình 4.5 Hình minh họa kết quả huấn luyện Bi-Gram.....	50
Hình 4.6 Hình minh họa kết quả huấn luyện Tri-Gram.....	50
Hình 4.7 So sánh tham số Precision của mô hình VNS và MMS.....	51
Hình 4.8 So sánh tham số Recall của mô hình VNS và MMS	52
Hình 4.9 So sánh tham số F-Measure của mô hình VNS và MMS	52

GIỚI THIỆU

1. Tổng quan

Xử lý ngôn ngữ tự nhiên (NLP: Natural Language Processing) là một nhánh của trí tuệ nhân tạo, tập trung vào các ứng dụng trên ngôn ngữ con người. Xử lý ngôn ngữ tự nhiên góp phần trong việc làm cho máy móc có thể hiểu được ngôn ngữ con người, từ đó tạo ra các hệ thống thông minh.

Nghiên cứu về xử lý ngôn ngữ tự nhiên bao gồm nhiều lĩnh vực quan trọng như: dịch máy (MT: Machine Translation), truy vấn và khai thác thông tin (IR-IE: Information Retrieval and Extraction), hệ thống hỏi đáp (Q&A: Question and Answer), tóm tắt văn bản (ATS: Automatic Text Summarization), nhận dạng tiếng nói (SR: Speech Recognition), v.v... Tất cả những công nghệ này giúp máy tính hiểu được con người.

Phân đoạn từ (WS: Word Segmentation) là một bước quan trọng trong xử lý ngôn ngữ tự nhiên tiếng Việt, đặc biệt là xử lý văn bản. Phân đoạn từ là việc xác định ranh giới giữa các từ trong câu.

Không giống như tiếng Anh và các ngôn ngữ Âu khác, tiếng Việt không sử dụng khoảng cách làm dấu hiệu xác định ranh giới từ. Ranh giới giữa các từ không có dấu hiệu rõ ràng mà cần phải dựa vào các yếu tố như: ngữ nghĩa, ngữ cảnh, văn phong, các từ lặp lại, v.v...

Ngoài ra, văn bản tiếng Anh, tiếng Việt ghép các từ gây nhiều khó khăn trong việc phân đoạn từ tiếng Việt. Phân đoạn từ có độ chính xác cao sẽ góp phần quan trọng vào các bài toán tiếp theo như: gán nhãn từ loại, kiểm tra cú pháp, dịch máy, v.v...

2. Lý do chọn tài

Vì các ngôn ngữ bị n hình nh tiếng Anh, Pháp, ... vì c nh n bi t ranh gi i gi a các t n gi n h n ti ng Vi t, ch y u đ a vào kho ng cách và các đ u câu. B n thân các t h u nh ã ph n ánh y hình thái, ngh a, th m chí ng pháp bên trong nó.

Tuy nhiên, ti ng Vi t là ngôn ngữ thu c h n l p, không bị n hình. V m t hình th c m t t có th c c u t o b i m t h o c nhi u âm ti t ghép l i. Kho ng tr ng ch dùng phân cách các âm ti t v i nhau. Có th ti n t i các x lý xa h n v x lý ngôn ngữ t nhiên tr c h t ta ph i làm t t bài toán phân o n t . T là n v c b n nh t phân tích cú pháp, ng ngh a c a ngôn ngữ .

Cho n nay, ã có r t nhi u công trình nghiên c u v phân o n t ti ng Vi t v i nh ng k t qu kh quan. Tuy nhiên các v n nh : hi n t ng phát sinh t m i, s nh p nh ng ng ngh a, v.v... ã nh h ng không ít n ch t l ng phân o n t . Vì v y phân o n t ti ng Vi t v n là ch c nhi u nhà nghiên c u quan tâm và là ng l c c a lu n v n này.

3. M c tiêu và ph m vi nghiên c u

Chúng tôi t ra m c tiêu nghiên c u chính c a lu n v n là xây d ng m t mô hình phân o n t ti ng Vi t đ a trên mô hình lai k t h p nhi u ph ng pháp nh m t ng c ng chính xác khi phân o n t .

Ph m vi c a tài t p trung nghiên c u phân o n t trên v n b n ti ng Vi t. V i u vào là m t v n b n ti ng Vi t, u ra là m t v n b n ti ng Vi t ã c phân o n thành các t .

V i m c tiêu nêu trên, lu n v n t p trung nghiên c u các v n sau ây:

- ❖ Nghiên c u t ng quan v x lý ngôn ngữ t nhiên. Kh o sát các công trình nghiên c u có liên quan n tài trong n c và qu c t .

- ❖ Nghiên cứu các lý thuyết về ngôn ngữ bao gồm: các loại hình ngôn ngữ, nguyên nhân của ngôn ngữ tiếng Việt, cấu trúc của ngôn ngữ tiếng Việt, nghiên cứu về từ vựng và hình thức ngữ pháp của tiếng Việt.
- ❖ Nghiên cứu các lý thuyết về ngôn ngữ học thống kê bao gồm: lý thuyết xác suất thống kê trong xử lý ngôn ngữ tự nhiên, mô hình Markov, mô hình thống kê N-Gram.
- ❖ Xây dựng kho ngữ liệu phục vụ các mô hình thống kê.
- ❖ Thu thập và xây dựng từ điển tiếng Việt, từ điển danh từ riêng.
- ❖ Nghiên cứu các phương pháp phân loại dựa trên từ điển.
- ❖ Nghiên cứu các phương pháp phân loại dựa trên mô hình thống kê.
- ❖ Nghiên cứu các phương pháp phát hiện từ misspelling trong kho ngữ liệu và thông tin tiếng Anh.
- ❖ Xây dựng mô hình phân loại tiếng Việt bằng cách kết hợp các phương pháp: phương pháp phân loại có tham khảo từ điển tiếng Việt, phương pháp nhận dạng danh từ riêng sử dụng từ điển danh từ riêng, phương pháp so sánh các mẫu dùng biểu thức chính quy và phương pháp phát hiện từ misspelling thông tin tiếng Anh.

4. B c c c a l u n v n

Luận văn kết cấu gồm có 5 chương. Chương 1: trình bày tổng quan về các hướng tiếp cận và các công trình nghiên cứu có liên quan liên quan tài. Chương 2: trình bày về các lý thuyết cơ bản, bao gồm các lý thuyết về ngôn ngữ và ngôn ngữ học thống kê. Chương 3: giới thiệu mô hình MMSEG – mô hình tham khảo chính của tài. Chương 4: giới thiệu mô hình phân loại dựa trên luận văn xuất. Chương 5: kết luận, đánh giá và nhận xét về những kết quả đạt được, những mặt còn hạn chế và hướng phát triển của tài.

CHƯƠNG 1. TÍNH QUAN

Không giống như tiếng Anh và các ngôn ngữ n-Âu s, độ khó của cách làm dù hi u phân cách t, h u h t các ngôn ngữ châu Á (nh tiếng Vi t, tiếng Thái, tiếng Nh t, v.v...) ph i d a vào nhi u y u t (nh ng ngh a, ng c nh, các t lân c n, v.v...) m i có th xác nh c ranh gi i gi a các t . Cho n nay ã có r tnh u công trình nghiên c u v phân o n t v i nhi u ph ng pháp khác nhau. Theo kh o sát c a chúng tôi các nghiên c u h u h t xu t phát t 3 h ng ti p c n chính sau ây: h ng ti p c n d a trên t i n (dictionary-based), h ng ti p c n d a trên th ng kê (statistics-based) và h ng ti p c n lai (hybrid-based).

H ng ti p c n d a trên t i n: ây là h ng ti p c n c b n nh t. c i m chung c a h ng ti p c n này là s d ng t i n t v ng k t h p v i cách thu t toán so kh p phân o n t . chính xác c a phân o n ph thu c vào tính y c a t i n. H ng ti p c n này có u i m: t c x lý nhanh, n gi n. Tuy nhiên có h n ch là không th xác nh c các t không có trong t i n, nh p nh ng phân o n có th x y ra l n.

H ng ti p c n d a trên th ng kê ho c th ng kê k t h p v i h c máy: h ng ti p c n này có c i m c n ph i xây d ng kho ng li u b ng cách thu th p d li u v ngôn ngữ , sau ó t i n hành th ng kê, h c máy trên kho ng li u thu th p c (g i là hu n luy n d li u), d a trên d li u hu n luy n và các thu t toán phân o n t . chính xác c a ph ng pháp ph thu c nhi u vào l n và bao quát c a kho ng li u. u i m c a h ng ti p c n này là có th phân o n c các t m i, h n ch c nh p nh ng phân o n nh ng có h n ch là t n nhi u th i gian, công s c xây d ng và x lý kho ng li u.

H ng ti p c n lai: s d ng k t h p cùng lúc nhi u ph ng pháp t ng c ng chính xác c a phân o n. u i m: chính xác c t ng c ng. Nh c i m: ph c t p l n.

Trong ph n ti p theo, chúng tôi nêu k t qu kh o sát và mô t m t s công trình nghiên c u có liên quan n tài. Các nghiên c u này c th c hi n trên ngôn ngữ tiếng Vi t ho c trên nh ng ngôn ngữ có c i m t ng ng v i tiếng Vi t.

Trên ngôn ngữ tiếng Myanmar, Hla Hla Htay và Kavi Narayana Murthy trong [14] sử dụng thuật toán số khớp dài nhất (LM: Longest Matching) phân ngôn tiếng Myanmar. Tác giả xây dựng bảng cách tiếp hợp kho ngữ 4550 âm tiết có trong ngôn ngữ, sau đó tiến hành gộp âm tiết tạo nên kho ngữ 800,000 tiếp hợp và các biến thể của tiếp hợp. Thử nghiệm các tiến hành trên 5000 câu (chứa 35049 tiếp hợp). Kết quả thu được 34,943 tiếp hợp và 34,633 tiếp hợp đúng. Chính xác F-measure đạt 98.95%.

Trên ngôn ngữ tiếng Hoa, Jin Kiat Low và cộng sự trong [18] sử dụng mô hình cực đại Entropy (ME: Maximum Entropy) có tham khảo tiến phân ngôn tiếng Trung Quốc. Tác giả sử dụng chứa kho ngữ 108.000 tiếp hợp. Thử nghiệm các tiến hành ngôn ngữ trên bốn corpus khác nhau: Academia Sinica (AS), City University of Hong Kong (CITYU), Microsoft Research (MSR) và Peking University (PKU). Kết quả F-measure đạt 95,6% - 96,9%.

Trên ngôn ngữ tiếng Nhật, Masaaki Nagata trong [20] xuất hiện mô hình phân ngôn tiếng Nhật dựa trên thống kê. Bảng chữ cái, mô hình sử dụng tiếp hợp các tiếp hợp là *word base*. Sau đó, tiến hành huấn luyện kho ngữ dựa trên vị trí tính toán tần suất hiện của các chữ cái trong tiếp hợp liên tiếp. Tiếp theo, *word base* được tạo nên bởi các tiếp hợp xác định trong quá trình huấn luyện. Cuối cùng, phương pháp thể hiện ánh xạ từ từ vựng không phù hợp trong *word base*. Khi kho ngữ liên tiếp lớn 3.9Mb với kho ngữ 1791 tiếp hợp, chính xác accuracy của phương pháp đạt 82,5%. Phương pháp này sử dụng *word base* như làm kinh nghiệm phân ngôn và không cần *word base* có kích thước lớn giai đoạn ban đầu.

Trên ngôn ngữ tiếng Thái, Thanaruk Theeramunkong và Sasiporn Usanavasin trong [24] xây dựng mô hình phân ngôn tiếng Thái dựa trên cây quy tắc không dùng từ vựng. Sử dụng cấu trúc từ vựng tiếng Thái làm dữ liệu để phân loại giai đoạn huấn luyện, tác giả tạo ra một corpus nhằm xây dựng cây quy tắc. Sau đó với ngôn ngữ tiếng Thái sử dụng phân ngôn dựa trên luật của cây quy tắc. Luật của cây quy tắc được xây dựng dựa vào những kết nối khác nhau không thể tách

ri, gi là “Thai character clusters - TCCs”. Th c nghi m trên kho ng li u ti ng Thái, k t qu chính xác accuracy t 87.41%.

Trên ngôn ng ti ng Vi t, nhi u mô hình phân o n t ã c nghiê n c u và xu t v i nh ng k t qu kh quan. Lê Trung Hi u và c ng s trong [13] xây d ng mô hình xác su t nh n d ng và phân tách t ti ng Vi t, ng th i áp d ng quá trình máy t h c xây d ng mô hình xác su t t i u. chính xác c a thu t toán phân tách t t trên 90%.

Tr n Ng c Anh và c ng s trong [3] xu t m t ph ng pháp phân o n t và x lý nh p nh ng phân o n d a trên mô hình lai. S d ng k thu t so kh p c c i (MM: Maximum Matching) phân o n t . Trong quá trình phân o n, tác gi s d ng ng th i ph ng pháp (FMM: Foward Maximum Matching) và (BMM: Backward Maximum Matching) nh m phát hi n nh p nh ng. Sau ó x lý nh p nh ng b ng cách k t h p nhi u ph ng pháp, bao g m: ph ng pháp th ng kê d a trên mô hình Bi-Gram trên t , mô hình N-Gram d a trên âm ti t, và ph ng pháp tham kh o t i n. Th c nghi m trên corpus ã c hu n luy n v i 2639 t p tin v n b n, v i 1,541,188 t . K t qu chính xác F-measure t 98.71% - 98.94%.

L u Tu n Anh và Yamamoto Kazuhide trong [2] xây d ng mô hình phân o n t v i h ng ti p c n Pointwise d a trên máy h c SVM. K t qu c a nghiê n c u c ng d ng xây d ng công c tách t có tên là ông Du v i chính xác 98,2 %.

Lê H ng Ph ng và c ng s trong [15] s d ng mô hình lai d a trên k thu t so kh p c c i k t h p automat h u h n tr ng thái và regular expression. Ngoài ra, x lý nh p nh ng, h th ng k t h p v i các th ng kê Uni-Gram và Bi-Gram hu n luy n trên t p v n b n tách t m u. K t qu nghiê n c u c ng d ng t o nên công c vnTokenizer v i chính xác F-measure t c g n 94%.

CHƯƠNG 2. C S LÝ THUYẾT

2.1 C s lý thuyết về ngôn ngữ

2.1.1 Phân loại ngôn ngữ

Xét theo loại hình ngôn ngữ, theo Nguyễn Thị Ngọc Giáp trong [9, tr 298–305] ngôn ngữ có thể chia làm 2 loại chính: ngôn ngữ n l p và ngôn ngữ không n l p. n l p có thể hiểu theo hai cách: n l p về ngữ âm và n l p về ngữ pháp. n l p về ngữ âm ghi nhận tính n t t c a t hay hình v. n l p về ngữ pháp nói nhận tính c l p c a t hoặc trong câu. Sự khác biệt cơ bản giữa 2 loại hình này là c i m c u t o c a t.

2.1.1.1 Ngôn ngữ không n l p

Ngôn ngữ không n l p có thể chia làm 3 loại chính: ngôn ngữ ch p dính, ngôn ngữ hòa kết và ngôn ngữ h n nh p.

❖ Ngôn ngữ ch p dính

C i m c a lo i ngôn ngữ này là sự dính r nhau các ph t c u t o t và biểu thị những mối quan hệ khác nhau. Mỗi ph t ch biểu thị cho một ý nghĩa ngữ pháp và ngữ c l i. Hình v trong các ngôn ngữ ch p dính có tính c l p l n và mối liên hệ giữa các hình v không chặt chẽ. Chính t có thể hoặc ngữ c l p.

Ví dụ, trong tiếng Thổ Nhĩ Kỳ:

- adam: người đàn ông
- adamlar: những người đàn ông
- kadin: người đàn bà
- kadinlar: những người đàn bà

Có thể liệt kê một số ngôn ngữ thuộc loại này như: tiếng Thổ Nhĩ Kỳ, tiếng Ugo-Ph n Lan, tiếng Bantu, v.v...

❖ Ngôn ngữ hòa kết

Còn có gì là ngôn ngữ chuyển đổi. Các hình ảnh ngôn ngữ này là có sự biến đổi nguyên âm và phụ âm trong hình vẽ mang ý nghĩa ngữ pháp. Ý nghĩa từ vựng và ý nghĩa ngữ pháp được dung hợp trong từng ngữ không thể tách bạch phần nào biểu thị ý nghĩa từ vựng, phần nào biểu thị ý nghĩa ngữ pháp. Một hình vẽ có thể mang nghĩa từ vựng và ngữ pháp. Các hình vẽ liên kết chặt chẽ với nhau.

Ví dụ :

Bảng 2.1 Bảng minh họa ngôn ngữ hòa kết

Tiếng Anh	foot: bàn chân– feet: nhiều bàn chân
Tiếng Ả Rập	balad: làng– bilād: nhiều làng

Các ngôn ngữ chuyển đổi gồm các tiếng Ấn-Âu như tiếng Pháp, tiếng Ý, tiếng Anh, tiếng Bungari, v.v...

❖ Ngôn ngữ hỗn hợp

Các hình ảnh các ngôn ngữ hỗn hợp là một từ có thể đứng ngữ vị một câu trong các ngôn ngữ khác. Nghĩa là từ ngữ hành động, trạng thái hành động không thể hiện bằng các thành phần câu chủ ngữ, tân ngữ, trợ ngữ, nhưng v.v... mà thể hiện bằng các phần khác nhau trong hình thái ngữ pháp.

Ví dụ : trong tiếng Tschinuk Bắc Mỹ, từ **“inialudam”** đứng ngữ vị câu **“Tôi ăn cho cô cái này”**.

Bảng 2.2 Bảng minh họa một từ trong tiếng Tschinuk

Ký tự	i	n	i	a	l	u	d	a	m
V trí	[1]	[2]	[3]	[4]	[5]	[6]	[7]	[8]	[9]

Trong đó:

- Ph âm d[7] là ng t chính trong câu, có nghĩa là **cho**
- Ti n t i[1] bi u hi n thì quá kh , có nghĩa là **ã**
- Ph t n[2] bi u hi n ngôi th nh t s ít, có nghĩa là **tôi**
- Ph t i[3] bi u hi n tân ng gi i t , có nghĩa là **cái này**
- Ph t a[4] bi u hi n tân ng c a gi i t , có nghĩa là **cô**
- Ph t l[5] cho bi t tân ng c a gi i t **cô** là gián ti p
- Ph t u[6] ch ra r ng hành ng **th ch ng**
- Ph t am[8,9] nh n m nh **tính có m c ích** c a hành ng.

M t s ngôn ng n Nam M và ông Nam Xibêri v.v... c ng thu c lo i ngôn ng h n nh p.

2.1.1.2 Ngôn ng n l p

V m t ng pháp, t trong ngôn ng n l p không bi n i hình thái. C u t o t do c n t ho c s k t h p gi a các c n t t o thành. Quan h ng pháp và ý nghĩa ng pháp c th hi n b ng các ph ng ti n ngoài t nh : tr t t t , h t , ng i u, v.v...

Ví d :

- | | |
|---------------|---|
| Dùng h t | <ul style="list-style-type: none"> • Cu n sách – nh ng cu n sách • i – s i, ang i, ã i, m i i, ... |
| Dùng tr t t t | <ul style="list-style-type: none"> • C a tr c – tr c c a • Nhà n c – n c nhà • Xanh m t – m t xanh |

Ranh gi i gi a t ghép và c m t ôi khi khó phân bi t rõ ràng. Ví d : xe p, nhà ph , v.v...

V m t ng âm, ngôn ng n l p th hi n rõ m i quan h gi a hình v và âm ti t. Ranh gi i gi a hình v trùng v i âm ti t t o nên hình ti t. Hình ti t là m t n

v có v ng âm là âm ti t, có khi c dùng v i t cách m t t , có khi c dùng v i t cách là y u t c u t o t .

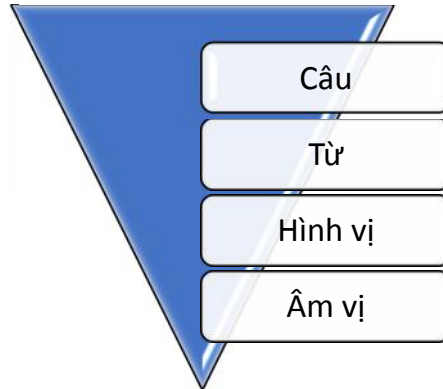
Âm ti t trong ngôn ngữ n l p có c u trúc ch t ch . M i âm v n m v trí nh t nh, có ch c n ng nh t nh.

Tiêu bi u cho ngôn ngữ n l p có th k n m t s ngôn ngữ nh : ti ng Hán, ti ng Thái, ti ng Dao, ti ng Mi n i n, ti ng Khmer, v.v...

Ti ng Vi t c ng thu c lo i hình ngôn ngữ n l p, không bi n i hình thái. S l ng v âm thanh mà ng i Vi t s d ng làm v ng âm cho hình v t i a kho ng 4 v n ti ng khác nhau. [5, tr.46]

2.1.2 n v ch y u c a ngôn ngữ

Theo Nguy n Thi n Giáp trong [11, tr 52-55] n v ch y u c a ngôn ngữ g m có:



Hình 2.1 Hình minh h a các n v ch y u c a ngôn ngữ

2.1.2.1 Âm v

Âm v còn c g i là âm ti t là n v t i thi u c a h th ng ng âm c a m t ngôn ngữ dùng c u t o và phân bi t v âm thanh c a các n v có ngh a c a ngôn ngữ .

Tiếng Việt thuộc loại hình ngôn ngữ có 6 thanh điệu. Vì vậy, khác với âm tiết của các ngôn ngữ châu Âu, âm tiết nào của tiếng Việt cũng mang một thanh điệu nhất định. Trong dòng lời nói, âm tiết tiếng Việt bao giờ cũng thể hiện khá rõ ràng, tách biệt và ngời ra thành từng khúc rõ ràng riêng biệt.

Theo Cao Xuân Hạo trong [12], hệ thống âm vị tiếng Việt bao gồm 22 phụ âm, 6 phụ âm cuội, 16 nguyên âm và 2 bán nguyên âm. Chi tiết hệ thống âm vị như sau:

- 22 phụ âm: /b, m, f, v, t, t', d, n, z, ẓ, s, ʃ, c, ʈ, ɲ, l, k, ʎ, ɣ, h, ʔ/
- 6 phụ âm cuội: /m, n, ɲ, p, t, k/
- 2 bán nguyên âm: /-w, -j/
- 16 nguyên âm: /i, e, ɛ, ɤ, ɤ̃, a, ɯ, ɤ, u, o, ɔ, ɔ̃, ɔ̃, ie, ɯɤ, uo/

Bảng 2.3 Bảng phụ âm u

Phương thức		Vị trí		Môi	Đầu lưỡi		Mặt lưỡi	Gốc lưỡi	Thanh hầu
					Bọt	Lưỡi			
Tắc	Ồn	Bật hơi			t'				
		Không bật hơi	Vô thanh		t	ʈ	c	k	ʔ
			Hữu thanh		b	d			
		Vang		m	n		ɲ		
Xát	Ồn	Vô thanh		f	s	ʃ		ç	h
		Hữu thanh		v	z	ẓ		ʝ	
		Vang			l				

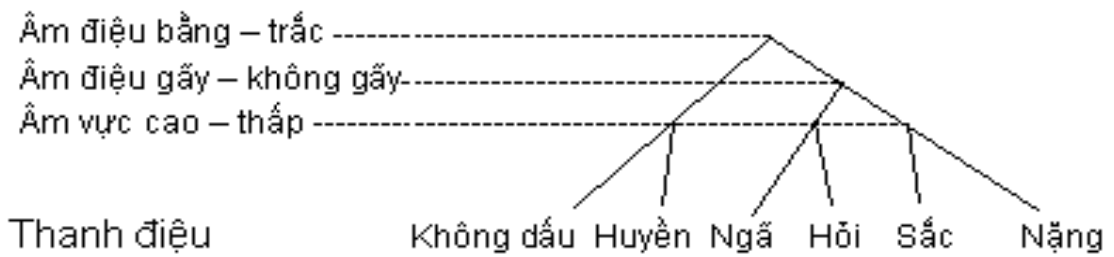
Bảng 2.4 Bảng phụ âm cuội và bán nguyên âm

Phương thức		Vị trí		Môi	Lưỡi	
					Đầu lưỡi	Gốc lưỡi
Ồn				p	t	k
Vang	Mũi			m	n	ɲ
	Không mũi			-w		-j

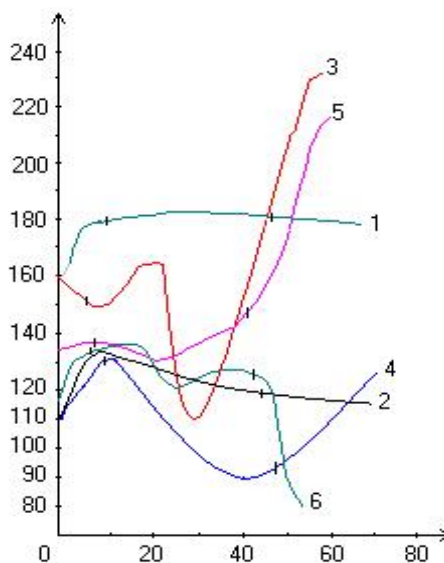
Bảng 2.5 Bảng nguyên âm

Âm sắc	Vị trí lưỡi, hình dáng môi	Trước, không tròn môi	Sau	
			Không tròn môi	Tròn môi
Có định	Nhỏ	i	u	u
	Lớn vừa	e	ɤ/ɛ	o
	Lớn	ɛ/ɛ̃	a/ã	ɔ/ɔ̃
Không có định		ɨɛ	ɯʉ	ɯo

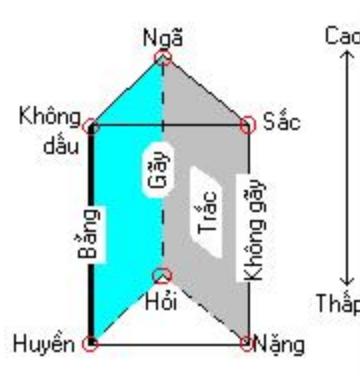
Theo Mai Ngọc Ch [7, tr 91-105] tiếng Việt có 6 thanh điệu: không dấu, huyền, ngã, hỏi, sắc, nặng.



Hình 2.2 Sơ đồ vị trí ba tiêu chí khu biệt cho sáu âm vị thanh điệu



Hình 2.3 Hình minh họa biểu đồ thanh điệu



Hình 2.4 Hình minh họa ngôn ngữ thanh điệu

2.1.2.2. Hình v

Hình v là một họ các chuỗi ký tự mà một vài âm vị, biểu thị một khái niệm. Hình v là những chuỗi có ý nghĩa. Chức năng của hình v là chức năng ngữ nghĩa.

Ví dụ, chuỗi "quốc gia" trong tiếng Việt gồm hai hình v: "quốc" là nước, "gia" là nhà; "Quốc gia" là hình v "Quốc gia" là hình v "Quốc gia", "Quốc gia" là chuyên chế, còn "- Quốc gia" là hình v "Quốc gia".

2.1.2.3. T

T là một khái niệm triết học rất thú vị của nhà ngôn ngữ học hiện đại. Theo H. L. P. C. i, triết gia ngôn ngữ Alexandri đã nói: "T là những chuỗi nói".

Theo E. Sapir thì: "T là một đơn vị có ý nghĩa, hoàn toàn có khả năng có lập và bản thân có thể làm thành câu nói".

Theo L. Bloomfield thì t là "một hình thái từ đơn".

Theo B. Golovin thì t là "những chuỗi có ý nghĩa của ngôn ngữ, có vẻ như được lập, tái hiện từ đó trong lời nói xây dựng nên câu". Đây chính là những đơn vị mà trong ngôn ngữ học hiện đại hay sử dụng.

Theo Solnev thì "T là những đơn vị ngôn ngữ có tính hai mặt: âm và nghĩa. T có khả năng có lập và cú pháp khi sử dụng trong lời".

Theo L c Chí V thì “T là n v nh nh t có th v n d ng t do trong câu”.

Theo V.G.Admoni thì “T là n v ng pháp, do hình v c u t o nên, dùng bi u th i t ng, quá trình, tính ch t và nh ng m i quan h trong hi n th c, có tính c thù rõ r t và có kh n ng ki n l p nhi u m i quan h a d ng nhau”.

Theo R.A.Bundagop thì “T là n v nh nh t và c l p, có hình th c v t ch t và có ngh a, có tính ch t bi n ch ng và l ch s ”.

Trong ngôn ng ti ng Vi t, c ng có nhi u nh ngh a c a ra. Theo quan i m c a Tr ng V n Trình và Nguy n Hi n Lê thì “T là âm có ngh a, dùng trong ngôn ng di n t m t ý n gi n nh t, ngh a là ý không th phân tích ra c”.

Theo Phan Khôi thì “T là m t l i t ra m t khái ni m trong khi nói”.. Theo Nguy n Lân thì “T là nh ng ti ng có ngh a, t c là m i khi nghe th y, trong óc chúng ta u có m t khái ni m”.

Theo Nguy n Kim Th n thì “T là n v c b n c a ngôn ng , có th tách kh i các n v khác c a l i nói v n d ng m t cách c l p và là m t kh i hoàn ch nh v m t ý ngh a (t v ng hay ng pháp) và c u t o”. Quan ni m c b n c a ông v “ n v c b n” là nh ng n v có s l ng h u h n và có ngh a. n v ó không th là câu (vì s l ng câu là vô h n) và c ng không th là âm ti t (vì nhi u âm ti t không có ngh a). V y n v c b n là cái gì ó nh h n câu và l n h n âm ti t.

Theo H Lê thì “T là n v ngôn ng có ch c n ng nh danh phi liên k t hi n th c, ho c có ch c n ng mô ph ng ti ng ng, có kh n ng k t h p t do, có tính v ng ch c v c u t o và tính nh t th v ý ngh a”. Theo ông t khác v i âm ti t ch y u v m t ý ngh a. T có kh n ng k t h p t do khi s d ng. T khác c m t b i tính v ng ch c v c u t o, tính nh t th v ý ngh a.

Nguy n Tài C n tuy không tr c ti p nh ngh a t trong ti ng Vi t, nh ng ông ã ch ng minh nh ng tính ch t c bi t c a “ti ng”, m t n v c ông xem nh là hình v và có tính n ng r t g n v i t , nó c ng có th là “t n” và là thành t tr c ti p t o nên “t ghép”. Theo ông, m i c thù v t c a ti ng Vi t b t ngu n t

tính n l p c a ti ng Vi t mà th hi n rõ nét nh t là qua m t n v c bi t, ó là “ti ng”. Quan i m này c ng c Cao Xuân H o ng tình.

K th a quan i m coi ti ng g n là t , Nguy n Thi n Giáp ã phát tri n t t ng này, ông coi “ti ng” trong ti ng Vi t chính là t trong các ngôn ng n-Âu.

Theo Mai Ng c Ch “T là n v nh nh t có ngh a, có k t c u v ng âm b n v ng, hoàn ch nh, có ch c n ng g i tên, c v n d ng c l p, tái hi n t do trong l i nói t o câu”.[10]

Có th th y, có r t nhi u quan i m v khái ni m “t ”, nh ng quan i m này tuy có nh ng khác bi t, nh ng h u nh không i l p mà b sung cho nhau. V y chính xác t là gì? Vi n s L.V.Sherba ã phát bi u: “Trong th c t , t là gì? Thi t ngh r ng trong các ngôn ng khác nhau, t s khác nhau. Do ó t t y u s không có khái ni m t nói chung”. [23]Cho n nay v n ch a có m t nh ngh a nào tr n v n v t , inh i n và H B o Qu c trong [10] ã ch ra nh ng nét c tr ng chính c a t nh sau:

- V hình th c: t ph i là m t kh i v c u t o (m t chính t , m t ng âm,...)
- V n i dung: t ph i có ý ngh a hoàn ch nh.
- V kh n ng: t có kh n ng ho t ng t do và c l p v cú pháp.

Chúng tôi th a nh n và s d ng nh ng nét c tr ng trên làm tiêu chí nh n di n t trong lu n v n này.

2.1.2.4. Câu

Câu là chu i k t h p c a m t hay nhi u t , ch c n ng c a nó là ch c n ng thông báo.

2.1.3 C u trúc c a n v t ti ng Vi t

2.1.3.1 Ti ng

“Ti ng” (còn g i là âm ti t) là n v c b n trong ti ng Vi t dùng c u t o các n v ngôn ng khác cao h n. S l ng ti ng trong ti ng Vi t không l n (kho ng 10.000), và chi u dài m i ti ng ng n (không quá 7 ch cái). Trong x lý

tiếng Việt tiếng bản máy tính, thì “tiếng” là ngôn ngữ tự nhiên nhưng mà máy tính đã dùng lưu trữ, nhận diện và xử lý. Tiếng chính là “tự chính tả”. [10]

Tiếng Việt có cấu tạo bằng các dùng một tiếng hoặc tập hợp các tiếng. Khi dùng một tiếng tạo nên từ các các từ. Khi tập hợp các tiếng tạo nên từ các các từ phức. Ngoài ra còn có một số như cấu tạo từ ghép từ các từ khác.

2.1.3.2 Từ đơn

Từ đơn là từ mà bộ phận không thể chia nhỏ, nghĩa, có thể đứng độc lập trong câu. Nói cách khác, từ đơn là từ chỉ có một thành tố, mỗi thành tố là một hình vị.

Căn cứ vào số lượng âm tiết có thể chia từ đơn ra làm hai loại: từ đơn âm và từ đơn a âm.

- Từ đơn âm: gồm một hình vị, mỗi hình vị là một âm tiết. Ví dụ: trái, t, nhà, quán, vui, v.v...
- Từ đơn a âm: gồm một hình vị, hình vị này gồm 2 âm tiết trở lên, ví dụ: cà vạt, cà phê, tivi, ...

2.1.3.3 Từ phức

Từ phức là từ cấu tạo từ hai hay nhiều hình vị kết hợp. Ví dụ: hi quân, bình, nhà nóc, Thơ Cầm Viên, v.v... Từ phức có thể chia làm 2 loại: từ láy và từ ghép.

2.1.3.3.1 Từ láy

Từ láy là những từ có một hình vị gốc và một hình vị láy. Hình vị láy có dạng âm trùng lặp hoàn toàn hoặc bộ phận với hình vị gốc. Có thể phân chia từ láy dựa vào số lượng âm tiết thành các loại như sau:

- Láy đôi: oai, chầm chầm, lung linh, ...
- Láy ba: sầm sầm, tu tu tu tu, ...
- Láy từ: bích bích, khắp khắp khiêng, ...

Ngoài ra từ láy còn có thể chia theo láy toàn phần và láy bộ phận:

- Láy toàn phần: xanh xanh, vui vui, buồn buồn, móm, ...
- Láy bộ phận: gồm láy âm và láy vần
 - o Láy âm: phần âm của hình vị gốc lặp lại trong hình vị láy, còn vần thì thay đổi. Ví dụ: d dàng, gòn gàng, chim chóc, m nh m , ...
 - o Láy vần: phần vần của hình vị gốc lặp lại trong hình vị láy, phần phần âm thì thay đổi. Ví dụ: lúng túng, b i r i, bèo nhèo, ...

2.1.3.3.2 T ghép

❖ Khái niệm:

T ghép là từ có ít nhất hai thành tố kết hợp với nhau. Mỗi thành tố có thể là một hình vị nguyên âm, a âm hay thậm chí là một từ hợp hình vị. Trong thì mỗi thành tố đều phải có nghĩa.

Ví dụ: xe đạp: có 2 thành tố là [xe, đạp]. Cà phê chè: gồm 2 thành tố là [cà phê, chè]. Trong đó cà phê: hình vị nguyên âm, chè: hình vị nguyên âm.

Có những từ hợp hình thức có vị ngữ từ láy, nhưng thực ra là từ ghép vì các thành tố đều có nghĩa riêng. Ví dụ: bóng bay, song sọt, v.v...

❖ Tính chất bản chất ghép:

Số kết hợp các hình vị tạo thành từ ghép phải chính xác về mặt hình thức và về mặt ngữ nghĩa.

- Chính xác về hình thức: thể hiện chính trong tiếng Việt các hình vị trong từ ghép luôn đi với nhau, từng hình vị không thể tách rời nhau. Nếu không thỏa mãn yêu cầu thì ý nghĩa ban đầu của từ ghép bị phá vỡ.

Ví dụ: các từ nhà gỗ, bàn cây, ghế sắt chúng ta có thể diễn giải như sau: cái nhà làm bằng gỗ, cái bàn làm bằng cây, cái ghế làm bằng sắt. Có thể bỏ đi các từ “làm bằng” thì nghĩa vẫn không thay đổi. Nhưng từ nhà gỗ nếu bỏ đi 1 trong 2 từ [nhà, gỗ] thì ý nghĩa sẽ khác hẳn.

- Ch t ch v ng ngh a: có nhi u m c khác nhau, m c cao nh t là có tính thành ng . Tính thành ng c hi u khi ý ngh a c a m t t h p không th gi i thích b ng cách gi i ngh a c a t ng y u t t o nên nó.

Ví d : khi nói v m t th y thu c có tài ch a b nh, ta có th dùng t : mát tay. Nh ng khi tách các t ra gi i ngh a thì ta l i thu c ý ngh a khác v i ban u. M t s ví d khác nh : l m mi ng (nhi u chuy n), y u tim (nhát), ...

2.1.3.4C m t c nh

C m t c nh có th chia thành 3 lo i: thành ng , quán ng và ng c nh nh danh. [8],[17]

2.1.3.4.1 Thành ng

Thành ng là c m t c nh, hoàn ch nh v c u trúc và ý ngh a. Ngh a c a chúng có tính hình t ng ho c/và g i c m. Ví d : ba c c ba ng, chó c n áo rách, nhà ngói cây mít, bán bò t u nh ng, méo mi ng òi n xôi vò, ông m t c a kia bà chìa c a n , ng nh nh ch nh trôi sông, v.v...

2.1.3.4.2 Quán ng

Quán ng là nh ng c m t c dùng l p i l p l i trong các lo i di n t thu c phong cách khác nhau. Ch c n ng c a chúng là a y, rào ón, nh n m nh ho c liên k t trong di n t . Ví d : c a áng t i, b ngoài tai, nói tóm l i, k t c c là, nói cách khác, v.v...

2.1.3.4.3. Ng c nh nh danh

Ng c nh nh danh là nh ng c m t c nh, nh danh, g i tên s v t, là nh ng n v n nh v c u trúc và ý ngh a h n các quán ng , nh ng ý ngh a mang tính hình t ng ch a c nh thành ng . Trong m i c m t nh v y th ng có m t thành t chính và m t vài thành t ph miêu t s v t c nêu thành t chính.

Ví d : lông mày lá li u, m t b câu, tr m ng, tóc r tre, con gái r u, bàn m u tính k , v.v...

2.1.4 T v ng ti ng Vi t

2.1.4.1 Gi i thi u

T v ng là t p h p t t c các t và nh ng n v t ng ng v i t c a m t ngôn ng . T v ng là ch t li u c n thi t, c b n nh t ki n t o nên m t ngôn ng mà n u thi u nó chúng ta không th hình dung c ngôn ng này [4].

2.1.4.2 S phân l p t v ng

H th ng t v ng c a m t ngôn ng th ng r t l n, thu n ti n cho vi c h c t p và nghiên c u ng i ta có th phân chia thành các l p t v ng riêng mang ý ngh a khái quát chung c a m t t p h p, m t nhóm. Theo nghiên c u c a Mai Ng c Ch và c ng s trong [8] trong ngôn ng ti ng Vi t các l p t v ng c phân l p d a vào các tiêu chí sau:

- Theo ngu n g c
- Theo ph m vi s d ng
- Theo t n s s d ng
- Theo phong cách s d ng

2.1.4.2.1 Theo tiêu chí ngu n g c

Theo tiêu chí ngu n g c, t v ng th ng c chia làm hai l p: l p t thu n và l p t ngo i lai. L p t thu n là l p t v n có c a ngôn ng ó; còn l p t ngo i lai là l p t vay m n c a ngôn ng khác trong quá trình giao thao v n hóa.

Trong ti ng Vi t, có l p t thu n Vi t và l p t có ngu n g c t ti ng Hán (g m Hán Vi t và Hán c), g c n - Âu (ti ng Anh, ti ng Pháp, ti ng Nga v.v...). Có th li t kê m t s t thông d ng nh : kh n mùi xoa, xà phòng, sô cô la, ti vi, mít tin, c n tin, cà v t, bi ông, ...

2.1.4.2.2 Theo tiêu chí ph m vi s d ng

Theo tiêu chí ph m vi s d ng, t v ng ti ng Vi t c chia thành các l p: t ph thông, t a ph ng, t ngh nghi p, thu t ng , ti ng lóng.

Từ phổ thông: là từ có tính phổ biến trong cộng đồng ngôn ngữ. Mỗi ngôn ngữ đều có từ phổ này, đóng vai trò cơ bản trong hệ thống từ vựng của một ngôn ngữ. Từ phổ này chính là từ vựng chủ yếu trong văn viết của ngôn ngữ đó.

Từ địa phương: là từ thuộc một phương ngữ, thường dùng trong giao tiếp hàng ngày, chủ yếu bị hạn chế trong lãnh thổ, phạm vi địa phương đó. Ví dụ: thây, u, m n, tía, má, ...

Thuật ngữ: là những từ dùng làm tên gọi cho các khái niệm, các đối tượng xác định trong một ngành, một lĩnh vực khoa học. Có tính chất: chính xác, chủ yếu, hệ thống và cụ thể hóa. Ví dụ: thuật ngữ trong lĩnh vực hóa học: nguyên tử, hợp chất, phân tử, nguyên tố, v.v...

Terminology: là từ bao gồm những từ ngữ chuyên ngành phân bố trong phạm vi ngành nghề nào đó. Ví dụ: lĩnh vực nghề làm mứt có những thuật ngữ: bào cóc, bào xoa, mứt, v.v... Lĩnh vực hát tuồng có: ào, kép, v.v...

Tên gọi: là từ do những nhóm người trong xã hội dùng để gọi tên những sự vật, hiện tượng, hành động văn hóa đã có tên gọi trong văn từ vựng chung. Ví dụ: lính phòng không (trai chày), hồi cá (lực lượng khác khi xảy ra sự cố), xe ô tô (xe ô tô) v.v...

2.1.4.2.3 Theo tiêu chí từ vựng

Theo tiêu chí từ vựng, từ vựng tiếng Việt được phân thành hai loại: từ tích cực và từ tiêu cực.

Từ tích cực: là những từ mang tính sử dụng mới mẻ, mới mẻ, có tính xu hướng cao, phân bố rộng. Đây là thành phần cơ bản của từ vựng.

Từ tiêu cực: là những từ có tính sử dụng thấp. Chia làm hai loại:

+ **Từ mới:** là những từ xu hướng bù đắp sự thiếu hụt của từ vựng. Khi mới xuất hiện, từ mới thường không được sử dụng rộng rãi nên thuộc từ vựng tiêu cực. Khi từ mới được chấp nhận và được sử dụng phổ biến thì trở thành từ

ng tích cực. Ví dụ: m t s t ng tr c ây thu c l p tiêu c c: t ch c (làm tí c), xây d ng (l p gia ình), ph n m m, ph n c ng, v.v...

+ T c : là nh ng t b lo i d n kh i h th ng t v ng hi n t i b i các nguyên nhân l ch s , xã h i, v n hóa, v.v... Ví dụ : i n trang (trang tr i l n), thái thú (m t ch c quan), dân cày (ng i làm ru ng), gác b u (cái ch n bìn), v.v...

2.1.4.2.4 Theo tiêu chí phong cách s d ng

Theo tiêu chí phong cách s d ng, t v ng ti ng Vi t c phân thành ba l p: l p t kh u ng , l p t thu c phong cách vi t và l p t trung hòa.

L p t kh u ng : là nh ng t ng dùng trong giao ti p, th ng có nh ng c i m sau ây: t do, phóng túng, c ng i u, th ng dùng kèm thành ng , quán ng , các t th a g i, v.v... Ví dụ : lo th t ru t, ch m t, ánh s c ti t, ch y b h i tai, v.v...

L p t thu c phong cách vi t: là nh ng t c ch n l c, trau d i, g n bó v i chu n t c nghiêm ng t. Có c i m chung là g n li n v i n i dung c a m t s phong cách ch c n ng c th nh : phong cách khoa h c, hành chính s v , chính lu n báo chí, v n h c. Không mang tính thông t c. Mang tính khái quát, tr u t ng, v.v... tu theo ph m vi riêng c a m i phong cách ch c n ng. Th ng dùng nhi u các t có g c Hán, n-Âu c du nh p. Ví dụ : phong cách khoa h c: âm v , hình v , ng pháp, v.v... ; phong cách hành chính s v : công v n, v n th , t t ng, v.v... ; phong cách v n h c: m u i, l ng l y, v.v...

L p t trung hòa: là nh ng t không mang d u hi u c tr ng c a l p t kh u ng ho c l p t thu c phong cách vi t. Ví dụ : au bu n, l ng l , i t n b , v.v...

2.1.5 V n nh p nh ng ngh a c a t

2.1.5.1 Gi i thi u

Nh p nh ng là hi n t ng m h v ng ngh a, không phân nh r ch rời ranh gi i gi a các t do hi n t ng a ngh a, a t lo i c a t , ho c do s k th p c a các âm ti t t c nh nhau t o thành nh ng t khác nhau, v.v... ây là hi n t ng

thông ngữ pháp khi xử lý ngôn ngữ tự nhiên. Xử lý ngữ pháp thông qua hình thức xử lý căn cứ vào ứng ngữ pháp trong câu. Ví dụ con ngữ thì đây không phải là v-n-l-n, vì con ngữ có thể ứng ngữ pháp trong câu như vào các yếu tố khác như: ngữ cảnh, ngữ nói, ngữ nghe, hoàn cảnh lịch sử, v.v... Nhưng vì máy tính thì đây lại là một vấn đề nghiêm trọng, vì máy tính không phải là con ngữ.

Ví dụ cho câu: Con bò c p con bò c p.

Câu này có thể hiểu theo nhiều cách:

- Con bò c p | con bò c p
- Con bò | c p | con bò c p
- Con bò | c p | con bò | c p

2.1.5.2 M t s h i n t ngữ pháp ngữ

2.1.5.2.1 Nh p nh ngữ ranh giới

V ranh giới, sự phân ngữ cấu trúc Việt có thể chia thành 2 kiểu sau:

- Ngữ pháp chéo: chủ ngữ “abc” c ngữ là ngữ pháp chéo như ngữ “ab”, “bc” u xuất hiện trong từ ngữ Việt.

Ví dụ như trong câu “ông già nhanh quá” thì chủ ngữ “ông già ” bị phân ngữ chéo vì các từ “ông già” và “già ” u có trong từ ngữ.

- Ngữ pháp kết hợp: chủ ngữ “abc” c ngữ là ngữ pháp kết hợp như từ “a”, “b”, “ab” u xuất hiện trong từ ngữ Việt.

Ví dụ như trong câu: “Bàn là này còn r t m i” thì chủ ngữ “bàn là” bị phân ngữ kết hợp, do các từ “bàn”, “là”, “bàn là” u có trong từ ngữ.

2.1.5.2.2 Nh p nh ngữ t a ngh a

B t c ngôn ngữ nào cũng có t a ngh a, nguyên nhân là vì r t nhiều khái niệm có các sắc thái ý nghĩa tuy không hoàn toàn trùng kh p nhau nhưng lại có nhiều nét t ng ng.

Ví dụ : Cho 2 câu

- Câu 1: Nó ăn bánh này

- Câu 2: Nó ăn cơm này.

Có thể thấy chữ “n” câu 1 và chữ “n” câu 2 có ý nghĩa hoàn toàn khác nhau. Ngoài ra cách dùng cũng khác nhau. Trong câu 1, chữ “n” là một từ, chỉ hành động ăn uống. Trong khi đó câu 2, chữ “n” lại là một âm tiết trong từ “cơm”, chỉ hành động ăn.

Hiện tượng này gây cản trở khá lớn trong xử lý ngôn ngữ tự nhiên như phân loại, dịch thuật, ...

2.1.5.2.3 Nhận diện ngữ âm

Hai từ ngữ âm vị tự nhau nghĩa là hai từ có âm vị ngữ tự nhau nhưng mang nghĩa khác nhau, còn ngữ pháp là hai từ về mặt ký hiệu là ngữ tự nhau nhưng nghĩa khác nhau. Do tính đa nghĩa của từ ngữ âm vị ngữ tự nhau là từ ngữ tự, các ngôn ngữ khác hai hiện tượng này không trùng khớp nhau. Công cụ phân biệt từ ngữ tự vị ngữ tự, trong từ ngữ tự các nghĩa đều có chung một ngữ gốc và do vậy luôn có nét ngữ tự ngữ tự trong khi đó trong từ ngữ tự chúng không có liên hệ ngữ gốc vị ngữ tự, nghĩa của chúng khác nhau rõ rệt. Ví dụ từ kim trong hai câu sau đây là hai từ ngữ tự :

- Anh ta sẽ đem kim rơm để luyên.
- Kim này bây giờ khó lắm.

Vì xác định nghĩa chính xác của từ ngữ tự dạng hình tự ngữ tự vị ngữ tự khác nhau nên việc nghĩa của chúng giúp phân biệt ra các tiêu chuẩn từ ngữ tự phân biệt.

2.1.5.2.4 Nhận diện ngữ loại

Từ loại là một yếu tố quan trọng trong việc xác định nghĩa chính xác của từ và sắp xếp các từ thành câu hoàn chỉnh trong dịch thuật. Nhưng vì có nghĩa là từ loại giúp phân biệt ngữ, nghĩa chính bản thân nó trong một chuỗi ngữ pháp cũng như ngữ.

Phân loại các ngôn ngữ biến hình từ loại xác nhận từ ngữ dễ dàng vì khi chuyển từ loại thì từ của ngôn ngữ biến hình của nó ví dụ trong tiếng Anh từ free là tính từ có nghĩa là tự do, chuyển từ loại thành danh từ có thêm hậu tố “dom” thành freedom nghĩa là sự tự do. Điều này trở nên dễ dàng vì các ngôn ngữ biến hình từ loại có cách thức biến hình các từ như nhau nên biến hình từ loại trở nên dễ dàng.

Các ngôn ngữ không biến hình như tiếng Việt vẫn xác nhận từ loại yêu cầu các thuật toán phân tích cú pháp, mặt khác ngay trong nội bộ ngành ngôn ngữ vẫn chưa có sự thống nhất về phân loại từ loại cho tiếng Việt.

2.2 Các lý thuyết về ngôn ngữ học thống kê

2.2.1 Tiếng quan về ngôn ngữ học thống kê

Ngôn ngữ học thống kê là một ngành khoa học có truyền thống lâu đời, ra đời trên cơ sở nghiên cứu về ngôn ngữ kết hợp với lý thuyết xác suất thống kê. Ngay từ thế kỷ 18, F.Kaeding đã áp dụng phương pháp thống kê trong ngôn ngữ xây dựng từ điển tiếng Đức. Năm 1913, nhà toán học Nga A.A.Markov đã dùng phương pháp xác suất thống kê nghiên cứu quy luật nối tiếp nhau của các phụ âm và nguyên âm trong tiếng Nga.

Từ thập niên 1950 trở lại đây, ngôn ngữ học thống kê đã liên tục phát triển và trở thành một thành tựu quan trọng của các lĩnh vực ngôn ngữ học như: ngữ âm học, từ vựng học, ngữ pháp học, ngữ nghĩa học, v.v... và đặc biệt là trong lĩnh vực máy học (ML: Machine Learning).

Ngôn ngữ học thống kê đã được áp dụng trong nhiều bài toán ngôn ngữ như:

- Nghiên cứu từ loại hình của ngôn ngữ.
- Xây dựng từ điển từ vựng, từ điển tiếng nước ngoài.
- Xác định phong cách tác giả thông qua các tác phẩm.
- Xử lý thông tin tự nhiên: tách câu, tách từ, dịch máy, soạn lịch chính tả, nhận diện tiếng nói.
- Xác định niên đại của ngôn ngữ, v.v...

Cơ sở toán học của ngôn ngữ học thống kê chính là lý thuyết xác suất thống kê. Trong phần tiếp theo chúng tôi sẽ trình bày một số lý thuyết thống kê thường dùng trong xử lý ngôn ngữ.

2.2.2 Một số lý thuyết xác suất thống kê trong xử lý ngôn ngữ

2.2.2.1 Hàm xác suất

Hàm xác suất của một biến ngẫu nhiên E là một ánh xạ từ miền xác định của E (không gian các giá trị E có thể nhận) vào không gian thực $[0,1]$.

Giả sử E có thể nhận các giá trị phân biệt e_1, e_2, \dots, e_n .

Hàm xác suất phải thỏa các tính chất sau:

1. $P(e_i) \geq 0, \forall i$
2. $P(e_i) \leq 1, \forall i$
3. $\sum_{i=1}^n P(e_i) = 1$

2.2.2.2 Xác suất điều kiện

Cho các biến ngẫu nhiên X và Y , xác suất điều kiện được định nghĩa:

$$P(X|Y) = \frac{P(XY)}{P(Y)}$$

Trong đó sử dụng ký hiệu:

1. $P(X)$ thay cho $P(X=x)$.
2. $P(XY)$ có nghĩa là đồng thời có $X=x$ và $Y=y$.

❖ Xác suất đồng thời

$$P(XY) = P(X) * P\left(\frac{Y}{X}\right)$$

❖ Định luật Bayes

$$P\left(\frac{X}{Y}\right) = \frac{P\left(\frac{Y}{X}\right) * P(X)}{P(Y)}$$

❖ Biến ngẫu nhiên rời rạc

Hai biến ngẫu nhiên X, Y độc lập khi và chỉ khi:

—

Tổ hợp suy ra nếu X, Y độc lập thì:

❖ **Kỳ vọng và phương sai**

Kỳ vọng là giá trị trung bình của biến ngẫu nhiên. Giả sử X là biến ngẫu nhiên, thì kỳ vọng là:

$$E(X) = \sum_x xP(x)$$

Phương sai của biến ngẫu nhiên A là một số không âm dùng để mô tả phân tán các giá trị của biến ngẫu nhiên xung quanh giá trị trung bình của nó.

$$Var(X) = E\left(\left(X - E(X)\right)^2\right) = E(X^2 - E^2(X))$$

2.2.2.3 Xác suất cổ điển

❖ **Xác suất chính xác**

Trong lý thuyết xác suất, nếu có một dãy dữ liệu ta có thể tính toán xác suất chính xác của một biến ngẫu nhiên. Chẳng hạn ví dụ số liệu thống kê từ năm 01-01-2014 đến ngày hôm nay (31-05-2014) về số kiện “trị mả trong ngày”, có 39 ngày mả trên tổng số 150 ngày. Như vậy ta tính được xác suất chính xác xảy ra số kiện “trị mả trong ngày” là 39/150 trong khoảng thời gian từ năm 01-01-2014 đến hết ngày hôm nay (31-05-2014).

Tuy nhiên, xác suất chính xác không phải là nội dung chính của lý thuyết xác suất thống kê. Vì ta chỉ cần tính được xác suất “trị mả trong ngày” cho những ngày tiếp theo, thì những thông tin chính xác.

❖ **Phương pháp ước lượng hợp lý cực đại (MLE: Maximum Likelihood Estimator)**

Là cách ước lượng dùng xác suất chính xác để ước lượng cho một biến ngẫu nhiên. Cách ước lượng MLE có độ chính xác tùy thuộc vào lượng dữ liệu: càng nhiều càng chính xác.

❖ **Phương pháp ước lượng kỳ vọng hợp lý (ELE: Expect Likelihood Estimator)**

Đây là phương pháp ước lượng thích hợp cho biến ngẫu nhiên có tính đối xứng. Xét biến ngẫu nhiên X , giá trị V_i là số lần xuất hiện $X = x_i$.

Khi thu thập được ELE tính xác suất theo công thức sau:

$$P(X = X_i) \cong \left(\frac{V_i}{\sum V_i} \right)$$

Để tránh vấn đề số bằng không, (xác suất bằng không), ta cần thêm một giá trị V_i , chẳng hạn:

$$V_i = |X_i| + 0.5, \text{ trong đó } |X_i| \text{ là số lần } X = x_i$$

Chúng ta xét tình huống dưới đây phân biệt MLE và ELE:

Giả sử trong thí nghiệm không xuất hiện trong khoảng li u và ta cần ước lượng xác suất xuất hiện các w_1, w_2, \dots, w_{40} .

Như vậy ta có một biến ngẫu nhiên X , với $X = x_i$ chỉ khi t xuất hiện trong $t = w_i$.

Với công thức ước lượng MLE, ta có $P(X = x_i)$ không xác định ($=0$), nghĩa là khoảng li u không cung cấp thông tin gì về số xuất hiện các t trong $t = w_i$.

Ngược lại, với công thức ELE, ta có $V_i = 0.5, \forall i = 1 \dots 40$, do đó:

$$P\left(\frac{W_i}{t}\right) \cong \frac{0.5}{0.5 * 40} = 0.025$$

Công thức này thể hiện thông tin ước lượng số xuất hiện các t trong w_i , mặc dù trong khoảng li u không hề có t .

2.2.2.4 Mô hình Markov ẩn và các mô hình N-Gram

❖ Mô hình Markov ẩn

Mô hình Markov ẩn (HMM: Hidden Markov Model)[1][26] là một tiến trình ngẫu nhiên kép. Tiến trình ngẫu nhiên ẩn đầu tiên là tiến trình Markov bất biến được biểu diễn bởi một lược đồ chuyển trạng thái. Một trạng thái là một quan sát có thể của tiến trình Markov và xác suất chuyển trạng thái A sang trạng thái B là $P(S_{t+1} = B | S_t = A)$ – xác suất chuyển trạng thái B thì ở $t+1$ với vị trí ẩn thì ở t trạng thái A. Tiến trình ngẫu nhiên thứ hai là tập các xác suất nhúng của trạng thái. Tiến trình xác suất thứ hai này tạo ra một bảng mã, tức là vị trí dãy ký hiệu quan sát được, dãy trạng thái sinh ra dãy ký hiệu âm thanh (không quan sát được). Nó cũng là vị trí quan sát.

HMM ẩn và ứng dụng phổ biến cho các mô hình thống kê, đặc biệt là trong lĩnh vực xử lý ngôn ngữ tự nhiên. HMM thực chất là một hàm xác suất của quá trình Markov. Các quá trình Markov được phát triển đầu tiên bởi Andrei A. Markov, là nhà tiên phong trong việc mô hình hóa các chuỗi ký tự trong tác phẩm văn học Nga (Markov 1913) – những sau đó mô hình này đã được phát triển thành một công cụ thống kê tổng quát.

Các mô hình Markov có thể sử dụng khi cần mô hình hóa xác suất của một dãy tùy tính các sự kiện. Chẳng hạn, chúng được sử dụng trong xử lý ngôn ngữ tự nhiên mô hình hóa dãy các ký tự (xử lý văn bản), âm thanh (xử lý tiếng nói), v.v...

Trong một HMM, ta không biết dãy trạng thái mà mô hình đi qua, nhưng biết một hàm xác suất của nó. Sau đây là định nghĩa tổng quát của một HMM:

HMM là một bộ 5 (S, K, π, A, B)

Trong đó:

1. S: tập các trạng thái
2. K: tập các ký tự output
3. π : xác suất khởi đầu trạng thái ban đầu
4. A: xác suất chuyển chuyển trạng thái

5. B: xác suất phát ra ký hiệu

Bảng 2.6 Bảng liệt kê các ký hiệu thông dụng trong HMM

Tập các trạng thái	$\overline{S} = \{\overline{s_1}, \dots, \overline{s_N}\}$
Mũ ký t output	$\overline{K} = \{\overline{k_1}, \dots, \overline{k_M}\} = \{1, \dots, \overline{M}\}$
Các xác suất trạng thái ban đầu	$\pi = \{\pi_i\}, i \in S$
Các xác suất chuyển đổi trạng thái	$\overline{A} = \{a_{ij}\}, i, j \in S$
Các xác suất phát ra bit ngữ	$\overline{B} = \{b_{ijk}\}, i, j \in S, k \in K$
Dãy trạng thái	$\overline{X} = (\overline{x_1}, \dots, \overline{x_T}) \in \overline{K}^T \rightarrow \{1, \dots, \overline{M}\}^T$
Dãy output	$\overline{O} = (\overline{o_1}, \dots, \overline{o_T}) \in \overline{K}^T$

Vì mô hình HMM cho trước, có thể dùng mô hình ngẫu nhiên có a m t quá trình Markov và t o ra m t dãy output.

❖ Các mô hình N-Gram

N-Gram [26], [19] là mô hình Markov n c dùng trong các gi i thu t x lý ngôn ngữ t nhiên s d ng nh ng ph ng pháp th ng kê.

Gi s ta c n tính xác suất xuất hiện c a chu i T_1, T_2, \dots, T_Q

Theo công thức tính xác suất đ ãng th i, ta có:

$$P(T_1 T_2) = P(T_1) P\left(\frac{T_2}{T_1}\right)$$

$$P(T_1 T_2 T_3) = P(T_1 T_2) P\left(\frac{T_3}{T_1 T_2}\right)$$

Tổng quát ta có:

$$P(T_1 T_2 \dots T_Q) = P(T_1) P\left(\frac{T_2}{T_1}\right) P\left(\frac{T_3}{T_1 T_2}\right) \dots P\left(\frac{T_Q}{T_1 T_2 \dots T_{Q-1}}\right)$$

Như vậy tính xác suất của chuỗi $T_1 T_2 \dots T_Q$, ta cần tính các xác suất liên tiếp trong vế phải của công thức trên. Khi áp dụng vào bài toán xử lý ngôn ngữ tự nhiên thì ngữ cảnh $T_1 T_2 \dots T_Q$ có thể xem là sự xuất hiện của một từ tiếng Việt có Q từ ngữ. Thuật ngữ sinh ra sự quá tải và không thể thể hiện được sự liên tiếp của từ trong từ là không bị tắc. Tuy nhiên theo nghiên cứu của nhóm tác giả Mai Ngọc Chấn, Vũ Ngọc Nghiêu, Hoàng Trọng Phiến trong [8], từ tiếng Việt có từ là 4 từ ngữ. Ví dụ: chènghĩa xã hội, ngành nghề, v.v... nên hoàn toàn có thể áp dụng mô hình này.

Giả sử từ trong từ, công thức tính xác suất dùng N-Gram như sau:

$$P\left(\frac{T_j}{T_1 T_2 \dots T_{j-1}}\right) \cong \left(\frac{T_j}{T_{j-n+1} T_{j-n+2} \dots T_{j-1}}\right)$$

Có nghĩa là sự xuất hiện của từ thứ j chỉ phụ thuộc vào $n-1$ từ ngữ trước. Thuật ngữ N-Gram là mô hình Markov, trong đó giá trị xác suất của một quan sát chỉ phụ thuộc vào $n-1$ quan sát trước nó.

Trong bài toán này, chúng ta sẽ dùng sự phân loại của các từ ngữ trước, sử dụng *lch s* tiên đoán cho từ ngữ xuất hiện tiếp theo. Tuy nhiên, chúng ta không thể xem xét từng *lch s* một cách riêng biệt, vì giữa các câu trong ngôn ngữ tự nhiên khi xuất hiện trong văn bản không hoàn toàn giống với các câu đã xuất hiện trước đó. Thuật ngữ trong từ ngữ hợp phần của câu đã có trong *lch s* thì phần cuối của nó vẫn hoàn toàn mới. Vì vậy không có một *lch s* nào có thể làm việc tiên đoán một cách chính xác. Do đó chúng ta cần một phương pháp thành lập các nhóm *lch s* từ ngữ làm việc cho từ ngữ kế tiếp. Mô hình phân loại sử dụng nguyên tắc: mỗi *lch s* có cùng $(n-1)$ từ ngữ xếp vào một lớp từ ngữ, đây

là mô hình Markov cấp $(n-1)$ và cũng chính là mô hình N-Gram. Tính cụ thể cùng trong N-Gram là tính cần tiên đoán sự xuất hiện của nó.

N-Gram trong bài toán này được áp dụng với $n=1$ (Uni-Gram), $n=2$ (Bi-Gram) và $n=3$ (Tri-Gram).

❖ Mutual information

Mutual information hay còn gọi là thông tin tương hỗ là một ứng dụng của lý thuyết thống kê. Thông tin này dùng để đo lường thông tin thu được về một biến ngẫu nhiên thông qua giá trị của một biến ngẫu nhiên khác. Trong xử lý ngôn ngữ tự nhiên mutual information có thể dùng để xác định sự liên kết giữa các âm tiết, làm cơ sở xác định các từ, các cụm từ.

Trong luận văn này, chúng tôi sử dụng mutual information để phát hiện các từ mới (chưa có trong từ điển), chúng tôi sử dụng mutual information theo công thức của Arabi Zhang, J., et al in [13]:

Mutual information MI(x,y) của một Bi-Gram (x, y) được tính như sau:

$$MI(x, y) = \frac{f(x, y)}{f(x) + f(y) - f(x, y)}$$

Mutual information MI(x,y,z) của một Tri-Gram (x,y,z) được tính như sau:

$$MI(x, y, z) = \frac{f(x, y, z)}{f(x) + f(y) + f(z) - f(x, y, z)}$$

V i:

- $f(x)$, $f(y)$ và $f(z)$ là t n s xu t hi n c a Uni-Gram(x), Uni-Gram(y) và Uni-Gram(z) trong ng li u hu n luy n.
- $f(x,y)$ là t n s xu t hi n c a Bi-Gram(x,y) trong ng li u hu n luy n.
- $f(x,y,z)$ là t n s xu t hi n c a Tri-Gram (x,y,z) trong ng li u hu n luy n.

M tN-Gram c xác nh là m t t khi mutual information \geq *threshold*. V i *threshold* là giá tr ng ng, giá tr ng ng này có th khác nhau tùy thu c vào ng li u. Trong th c nghi m này, chúng tôi tìm *threshold* là 0.03 cho Bi-Gram và 0.02 cho Tri-Gram.

CHƯƠNG 3. GIỚI THIỆU MÔ HÌNH MMSEG

3.1 Tổng quan về MMSEG

MMSeg là một hệ thống phân loại tiếng Hoa được xuất bản bởi Chih-Hao Tsai[6]. Đây là hệ thống phân loại tiếng Hoa sử dụng hai dạng cấu trúc toán số khớp Maximum Matching kết hợp với tính toán và bản đồ giúp phân tích ngữ pháp. Theo kết quả tác giả đã công bố, khi thử nghiệm trên kho ngữ liệu chứa 1.013 từ, hệ thống cho kết quả rất khả quan (98.41%). Cấu trúc bản đồ phân tích ngữ pháp như sau:

- Luật 1: Maximum matching - số khớp tối đa:
 - Áp dụng thuật toán Maximum matching để tìm kiếm: lý thuyết có chi phí dài nhất.
 - Áp dụng thuật toán Maximum matching để phân tích: lý thuyết ưu tiên tối đa 3 từ có chi phí dài nhất. Nếu có nhiều hơn một từ 3 từ có chi phí dài nhất thì áp dụng luật tiếp theo.
- Luật 2: Chi phí trung bình của từ nhất:
 - Lý thuyết ưu tiên của từ 3 từ có chi phí trung bình nhất.
 - Nếu có hơn một từ 3 từ có chi phí trung bình nhất thì áp dụng các quy tắc tiếp theo.

Ví dụ:

- Trường hợp 1 (TH1): _C1_C2_C3_
- Trường hợp 2 (TH2): _C1C2C3_

Theo luật 2, ta lấy từ C1C2C3 là TH2.

- Luật 3: Phân tích ngữ pháp chi phí dài:
 - Chọn từ ưu tiên trong từ 3 từ có phân tích ngữ pháp chi phí dài nhất.
 - Nếu có nhiều hơn một từ 3 từ có phân tích ngữ pháp chi phí dài thì áp dụng các quy tắc tiếp theo.

Ví dụ :

- TH1: _C1C2_C3C4_C5C6_
- TH2: _C1C2C3_C4_C5C6_

Theo luật số 3, ta lấy C1C2 trong TH1.

- Luật số 4: tính logarit n của tích do hình vẽ của các t:
 - Công thức cơ bản để tính toán tính logarit n của tích do hình vẽ là tính logarit n của các t trong b.
 - Luật số 4 cho phép lấy tích của a b có tính logarit n của a và b.
 - Khi có tích của hai t có cùng tính do hình vẽ nên theo tác giả này nên xem như tích quy t.

Ví dụ :

- TH1: _C1_C2_C3C4_
- TH2: _C1_C2C3_C4_

3.2 Áp dụng MMSEG vào tiếng Việt

Trong phần này chúng tôi áp dụng MMSEG trên ngôn ngữ tiếng Việt để đánh giá mức hiệu quả của hệ thống xử lý tiếng Việt. Thử nghiệm được tiến hành trên kho ngữ liệu gồm 10.000 câu trích từ VietTreebank[22]. Ngữ liệu được sắp xếp ngẫu nhiên sau đó chia thành 5 phần cho 5 lần thử nghiệm.

Để đánh giá, chúng tôi sử dụng các tham số Precision(1), Recall(2) và F-measure (3) được tính theo các công thức sau:

ở đây:

$$Precision = \frac{CorrectWords}{FoundWords} \quad (1)$$

$$Recall = \frac{CorrectWords}{StandardWords} \quad (2)$$

$$F\text{-measure} = 2 * \frac{Recall * Precision}{Recall + Precision} \quad (3)$$

V i:

- StandardWords: s t chu n theo ng li u.
- CorrectWords: s t phân o n úng so v i s t chu n
- FoundWords: s t do h th ng tìm c.

Sau ây là k t qu thu c:

B ng 3.1 B ng li t kê k t qu th c nghi m MMSeg trên ng li u ti ng Vi t

L n th c nghi m	S t chu n	S t tìm c	S t úng	Precision(%)	Recall(%)	F-measure(%)
		MMS	MMS	MMS	MMS	MMS
L n 1	44,030	43,030	38,538	89.56	87.53	88.53
L n 2	38,724	38,004	33,440	87.99	86.35	87.16
L n 3	35,570	34,823	30,479	87.53	85.69	86.6
L n 4	32,668	32,035	29,210	91.18	89.41	90.29
L n 5	31,581	30,925	27,744	89.71	87.85	88.77
Trung bình	36,515	35,763	32,282	89.19	87.37	88.27

3.3 ánh giá MMSeg trên ngôn ng ti ng Vi t

C n c vào k t qu th c nghi m, chúng tôi nh n th y MMSeg thu c k t qu t t khi áp d ng trên ngôn ng ti ng Hoa (F-measure: 98.41%) nh ng th p h n nhi u khi áp d ng trên ti ng Vi t (F-measure: 88.27%). Sau ây là m t s nguyên nhân chúng tôi tìm c:

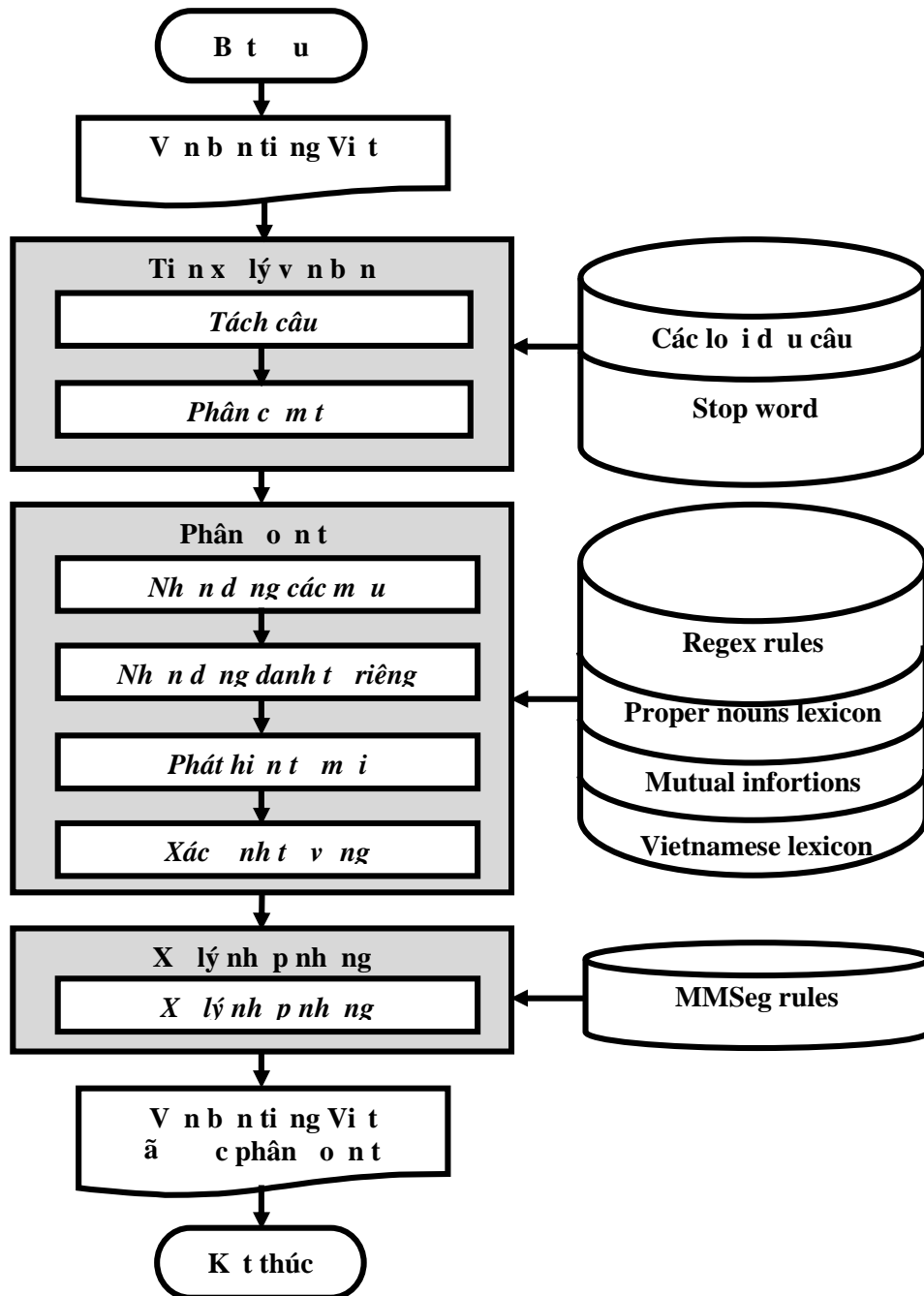
- T i n ti ng Vi t không y , t m i phát sinh, t m n ti ng n c ngoài t ng i nhi u nh ng ch a c c p nh t trong t i n.
- Các danh t riêng nh tên, a danh, ... th ng c s d ng r t ph bi n nh ng khó c li t kê t t c trong t i n.

- Nhi u tên riêng có ngu n g c t danh t (Hoa, Tùng, Lan, Ngân Hà, ...), tính t (Th ng M n, Tu n Tú, Lung Linh, ...), ho c c m t (Hai Bà Tr ng, Ph Hàng Bông, ...) v.v... gây ra nhi u nh p nh ng khi phân o n.
- Các m u c bi t nh s th p phân (m t ph n t , ba ph n tám, ...), ngày tháng (tháng ch p, tháng giêng, ...), v.v... th ng không c phân o n chính xác.
- Ngoài ra, vì c k t h p các tỉ ng t o nên t trong tỉ ng Vi t r t a d ng c ng gây nhi u khó kh n khi phân o n t .

CHƯƠNG 4. MÔ HÌNH XUẤT

4.1 Mô hình phân o n t

Trong phần này chúng tôi xuất m t mô hình phân o n t t i ng Vi t t tên là Vietnamese Segmentation(VNS) d a trên mô hình lai.



Hình 4.1 Hình minh h a mô hình phân o n t t i ng Vi t (VNS)

Mô hình được xây dựng bằng các kỹ thuật phân tích pháp âm để nghiên cứu và áp dụng trong các lĩnh vực: mô hình so sánh dựa trên từ vựng, mô hình thống kê dựa trên N-Gram, mô hình nhận dạng từ dựa trên thông tin thống kê, mô hình nhận dạng danh từ riêng dựa trên từ vựng và mô hình so sánh dựa trên các luật. Ngoài ra chúng tôi còn tiếp tục nghiên cứu thêm khả năng phân tích ngữ pháp bằng các kỹ thuật dựa trên từ vựng và bản đồ phân tích ngữ pháp của hệ thống MMSEG.

4.2 Thiết kế hệ thống

Đầu vào của mô hình là một văn bản tiếng Việt, đầu ra là văn bản tiếng Việt đã được phân loại. Mô hình gồm có ba bước xử lý chính. Bước 1: tiền xử lý văn bản. Bước 2: phân loại. Bước 3: xử lý ngữ pháp. Trong khuôn khổ luận văn này, chúng tôi chỉ xây dựng thuật toán cho bước 1 và bước 2. Thuật toán phân tích ngữ pháp bước 3 chúng tôi không xây dựng mà sử dụng lại thuật toán của hệ thống MMSEG trong [6]. Phân tích theo sơ đồ trình bày các hệ thống.

4.2.1 Hệ thống tiền xử lý văn bản

Input:

Văn bản tiếng Việt cần phân loại, danh sách dấu câu, danh sách stop word.

Output:

Danh sách các cụm từ.

Hệ thống:

```

1. Open file VanBan.txt as Text for "Read"
2. Open file DSDauCau.txt as SignList for "Read"
3. Open file DSStopWord.txt as StopWordList for "Read"
4. SentenceList = [] // Ch a danh sách các câu
5. ChunkingList = [] // Ch a danh sách các c m t
6. foreach(sign in SignList) do
7. Sentence = DoSentenceChunking(Text, sign)
8. add Sentence to SentenceList
9. end
10. for each (sentence in SentenceList) do
11. for each (stopWord in StopWordList) do
12. Chunking = DoPhraseChunking(sentence, sign)
13. add Chunking to ChunkingList
14. end
15. end
16. return ChunkingList

```

Trong giai đoạn này chúng tôi tiến hành phân chia câu và phân chia các từ dựa trên các dấu câu và các stop word. Mục đích nhằm rút ngắn thời gian xử lý cho các bước tiếp theo.

Dấu câu là các loại dấu câu thường dùng trong tiếng Việt như: dấu chấm, dấu phẩy, dấu chấm phẩy, dấu chấm than, v.v...

Stop word là những từ thường xuất hiện trong hầu hết các tài liệu, thường dùng liên kết các mệnh đề, các từ và các mệnh đề với nhau. Luận văn sử dụng danh sách stop word do tác giả kê b i tác giả Nguyễn Văn Dân trong [9] phân chia từ. Sau đây là danh sách một số stop word luận văn sử dụng:

Bảng 4.1 Danh sách m t s stop word trong từ ng Vi t

và	còn	hay	ho c	không
không nh ng	không ch	mà	còn	n u
thì	nên	h	tuy	nh ng
v l i	giá	vì	b i	t i
do	song	d u	m c d u	dù
d u	d u cho	ch ng l	làm nh	th mà
v y mà	có i u	h n n a	hu ng h	hu ng gì
hu ng n a	ngay	c ng	chính	c

4.2.2 Gi i thu t phân o n t

Input: danh sách các c m t (*ChunkingList*), danh sách bi u th c chính quy, t i n danh t riêng, thông tin t ng h , t i n từ ng Vi t

Output: danh sách các t

Gi i thu t:

```

1. Open file BieuThucChinhQuy.txt as RegexRuleList for "Read"
2. Open file TuDienDanhTu.txt as ProperNounDic for "Read"
3. Open file ThongTinTuongHo.txt as MutualInforDic for "Read"
4. Open file TuDienTiengViet.txt as VNLexicon for "Read"
5. WordList = [] // Ch a danh sách các t
6. for each (chunk in ChunkingList) do
7. wordsMatch1 = DoMatchingRegex(chunk, RegexRuleList, ref chunk2)
8. wordsMatch2 = DoMatchingProperNoun(chunk2, ProperNounDic, ref chunk3)
9. wordsMatch3 = DoDetectNewWord(chunk3, MutualInforDic, ref chunk4)
10. wordsMatch4 = DoDetectWord(chunk4, VNLexicon)
11. add wordsMatch1, wordsMatch2, wordsMatch3, wordsMatch4 to WordList
10. end
11. return WordList;

```

Mô t gi i thu t

u vào c a gi i thu t là danh sách các c m t và các thông tin tham kh o, ti n hành t o m t danh sách các t g i ý theo các b c:

- ❖ B c 1: t o m t danh sách các t g i ý s d ng ph ng pháp so trùng m u, v i các lu t c nh ngh a b ng các bi u th c chính quy.
- ❖ B c 2: t o m t danh sách các t g i ý s d ng t i n danh t riêng.
- ❖ B c 3: t o m t danh sách các t g i ý s d ng t i n t v ng.
- ❖ B c 4: t o m t danh sách các t g i ý s d ng thông tin t ng h đ a trên ng li u hu n luy n N-Gram, giúp xác nh các t m i.

u ra c a gi i thu t là m t danh sách các t . Danh sách này s c dùng nh t i n t v ng ph c v cho gi i thu t phân gi i nh p nh ng.

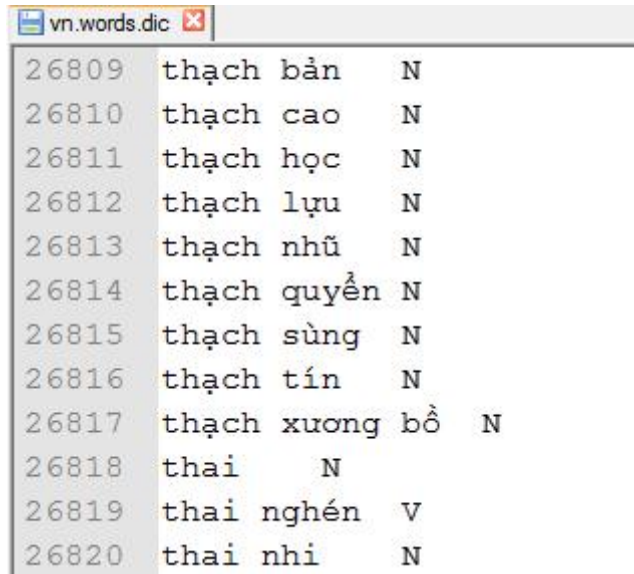
Ph n ti p theo chúng tôi trình bày v quá trình thu th p và xây d ng các t i n ph c v cho th c nghi m.

4.2 T i n và kho ng li u

th c nghi m, chúng tôi ti n hành thu th p và xây d ng t i n t v ng, t i n danh t riêng, kho ng li u hu n luy n N-Gram và t p lu t nh n d ng các m u. Sau ây là k t qu thu c:

❖ T i n t v ng

T i n t v ng lu n v n này s d ng là t i n ti ng Vi t - Vietnamese Lexicon [16], m t s n ph m c a đ án KC01.01/06-10, ã tùy bi n l i c u trúc cho phù h p. Phiên b n lu n v n s d ng có 3,1243 t .



Code	Word	Part of Speech
26809	thạch bản	N
26810	thạch cao	N
26811	thạch học	N
26812	thạch lựu	N
26813	thạch nhũ	N
26814	thạch quyển	N
26815	thạch sùng	N
26816	thạch tín	N
26817	thạch xương bồ	N
26818	thai	N
26819	thai ghen	V
26820	thai nhi	N

Hình 4.2 Hình minh họa cấu trúc từ vựng tiếng Việt

❖ Từ vựng riêng

Theo khảo sát của chúng tôi, trong khuôn khổ tài liệu nhà nước KC.01.21 – “Nghiên cứu các kỹ thuật xây dựng và khai thác thông tin Web có nghĩa”, hiện nhóm tác giả Cao Hoàng Trĩ và đồng nghiệp đang xây dựng một cơ sở dữ liệu về các thực thể có tên tiếng Việt. Dự án này có tên là VN-KIM [25]. VN-KIM Ontology bao gồm 347 lớp thực thể và 114 quan hệ và thuộc tính. Về kho ngữ 85,767 tên riêng.

VN-Kim thực sự phù hợp làm từ vựng riêng trong mô hình xuất. Tuy nhiên, vì chia sẻ liên hệ với tác giả xin phép sử dụng dữ liệu nên chúng tôi xây dựng từ vựng riêng phục vụ cho việc thực nghiệm. Từ vựng có hai loại: tên riêng và tên địa danh. Dữ liệu được thu thập từ danh bạ internet trực tuyến *nhungtrangvang.com.vn* và bách khoa toàn thư *vi.wikipedia.org*. Sau khi loại bỏ các tên trùng nhau, các tên không hợp lệ chúng tôi thu được tổng cộng 1,068,435 tên. Cấu trúc:

Bảng 4.2 Bảng liệt kê số lượng tên cá nhân danh từ riêng

Loại tên	Số lượng
Tên người	1,065,613
Tên họ danh	2,822
Tổng cộng	1,068,435



Hình 4.3 Hình minh họa tên cá nhân danh từ riêng

❖ Ng liệu hu n luy n N-Gram

Ng liệu hu n luy n cho mô hình N-Gram c chúng tôi thu th p t ng t các báo i n t tr c tuyền nh : vietnamnet.vn, vnexpress.net, dantri.com.vn, tuoitre.com.vn, vnthuquan.net, v.v... Tr c khi hu n luy n chúng tôi làm s ch v n b n b ng cách: bóc tách toàn b html, gi l i ph n n i dung chính, xóa các n i dung mang tính l p l i nh : menu, logo, qu ng cáo, các liên k t liên quan, T ng dung l ng thu c là 26MB v i 3356 bài báo. C th :

Bảng 4.3 Bảng liệt kê số lượng bài báo phân bố cho các chủ đề nghiên cứu

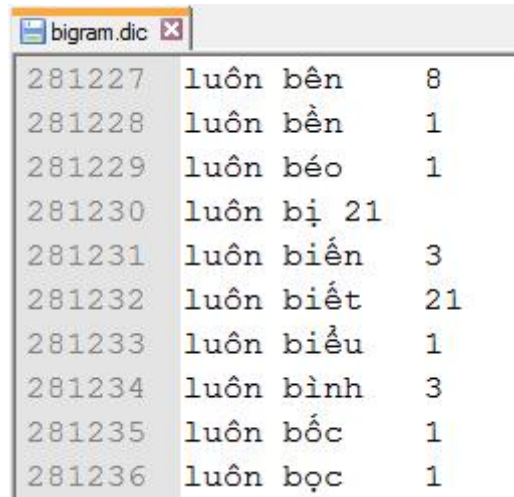
Loại bài báo	Số lượng
Thị trường	641
Giáo dục	367
Kinh tế	340
Văn hóa	352
Xã hội	525
Công nghệ	231
Nghiên cứu	450
Thể thao	316
Sports	134
Tổng cộng	3356

Sau khi huấn luyện chúng tôi thu được các chỉ số huấn luyện như sau:

- Uni-Gram: 16.340 m/c
- Bi-Gram: 571.821 m/c
- Tri-Gram: 1.769.290 m/c

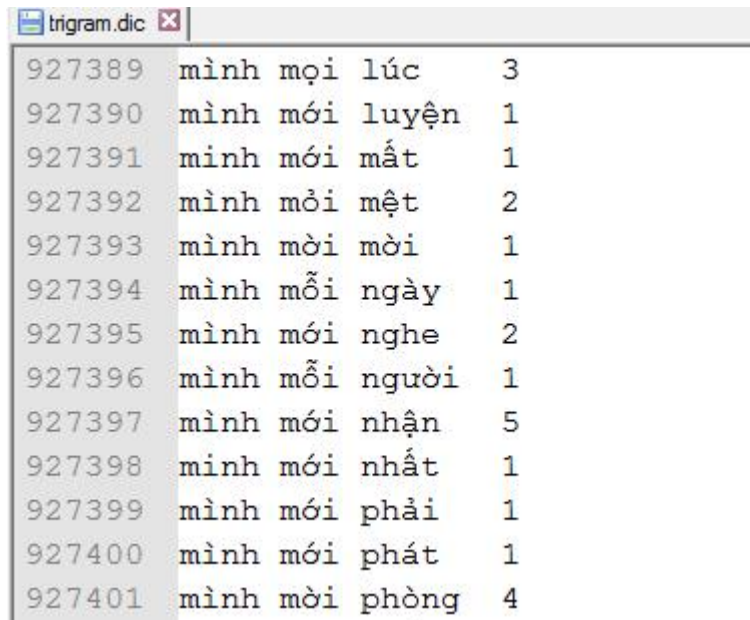
Index	Word	Count
10119	nhằm	36
10120	nhằm	75
10121	nhậm	787
10122	nhân	254
10123	nhân	1
10124	nhân	285
10125	nhân	15871
10126	nhân	233
10127	nhân	406
10128	nhân	247
10129	nhân	310
10130	nhân	38

Hình 4.4 Hình minh họa kết quả huấn luyện Uni-Gram



bigram.dic			
281227	luôn	bên	8
281228	luôn	bền	1
281229	luôn	béo	1
281230	luôn	bị	21
281231	luôn	biển	3
281232	luôn	biết	21
281233	luôn	biểu	1
281234	luôn	bình	3
281235	luôn	bốc	1
281236	luôn	bọc	1

Hình 4.5 Hình minh họa kết quả huấn luyện Bi-Gram



trigram.dic				
927389	mình	mọi	lúc	3
927390	mình	mới	luyện	1
927391	mình	mới	mất	1
927392	mình	mỏi	mệt	2
927393	mình	mời	mời	1
927394	mình	mỗi	ngày	1
927395	mình	mới	nghe	2
927396	mình	mỗi	người	1
927397	mình	mới	nhận	5
927398	mình	mới	nhất	1
927399	mình	mới	phải	1
927400	mình	mới	phát	1
927401	mình	mời	phòng	4

Hình 4.6 Hình minh họa kết quả huấn luyện Tri-Gram

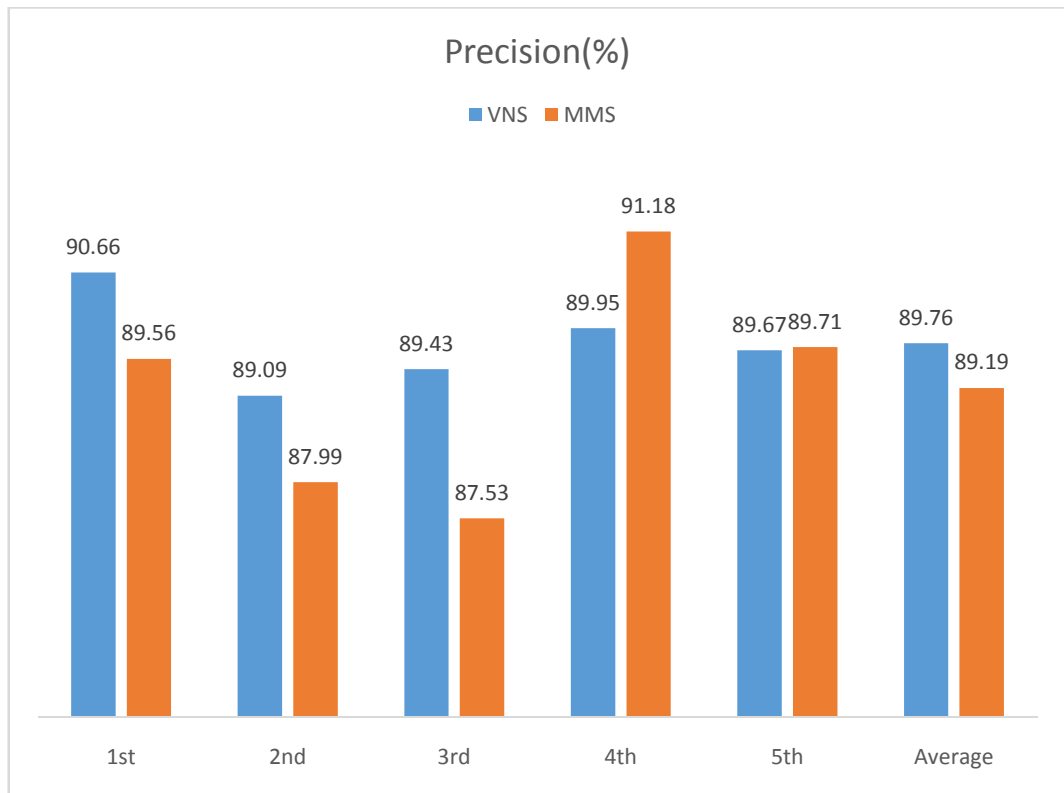
4.3 Thử nghiệm

Thử nghiệm được tiến hành trên các mẫu ngữ liệu đã dùng trong mục 5.2 và sử dụng một phương pháp đánh giá nhằm so sánh hai mô hình.

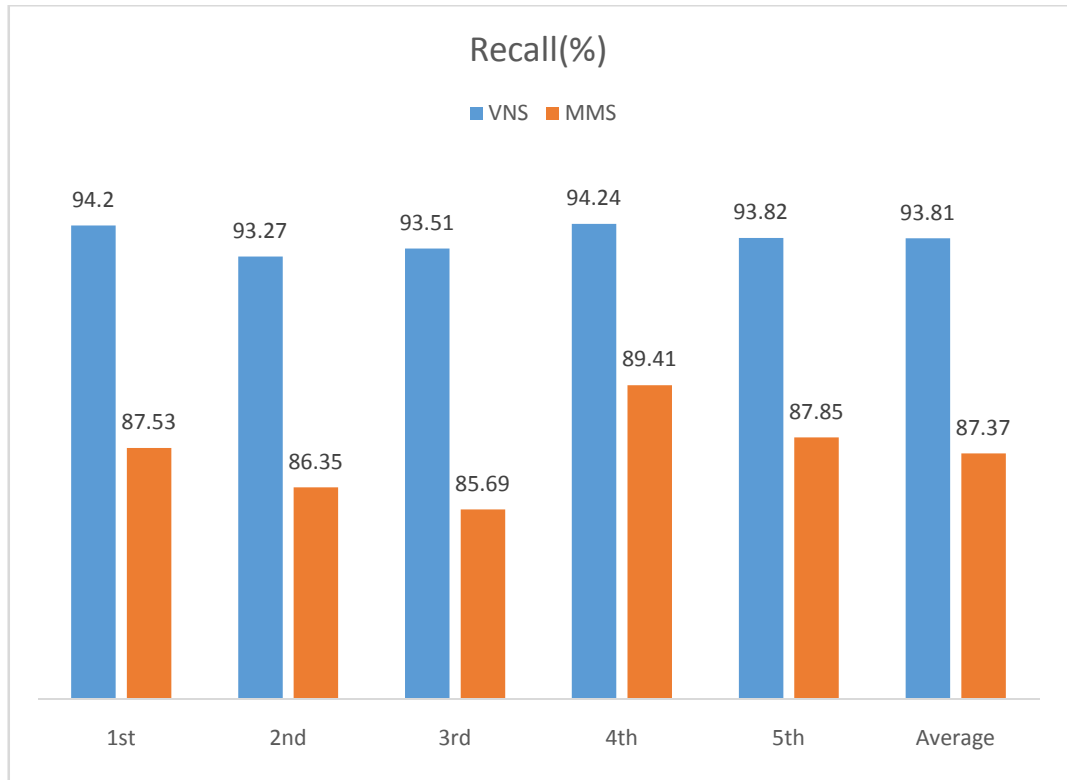
Bảng 4.4 trình bày kết quả của mô hình do chúng tôi đề xuất (VNS) so với mô hình của MMSeg (MMS) sau quá trình thử nghiệm:

Bảng 4.4 Bảng liệt kê kết quả thực nghiệm của VNS so với MMS

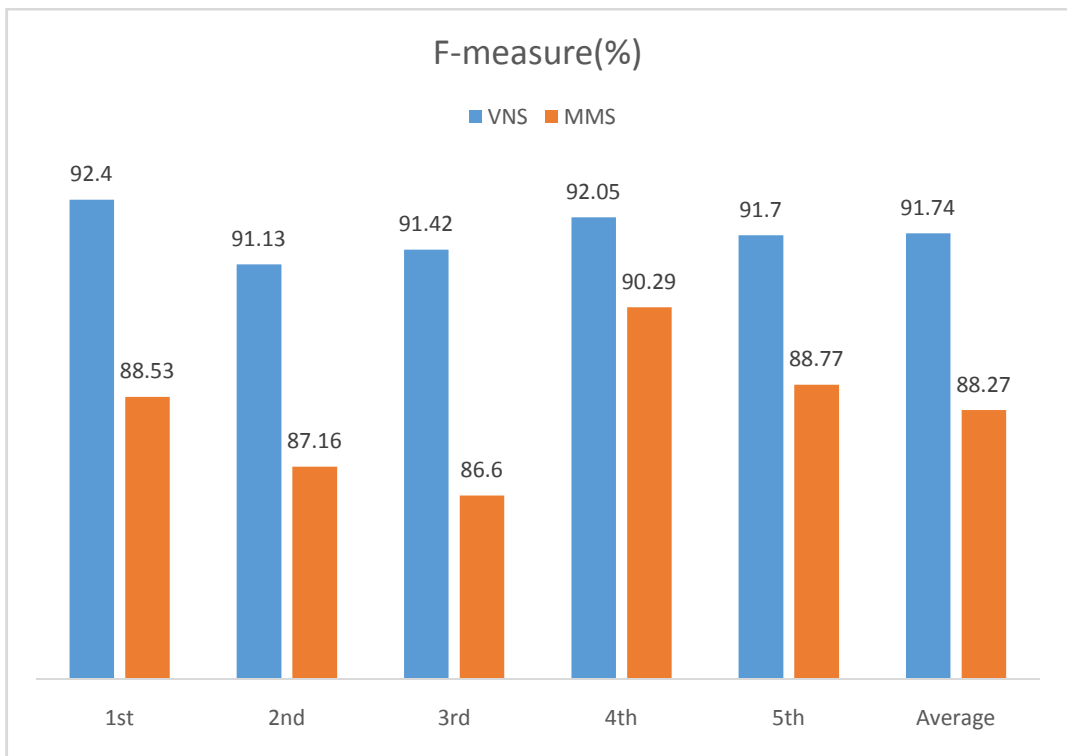
Lần thực nghiệm	Số chu n	Số tìm c		Số úng		Precision(%)		Recall(%)		F-measure(%)	
		VNS	MMS	VNS	MMS	VNS	MMS	VNS	MMS	VNS	MMS
1	44,030	45,751	43,030	41,477	38,538	90.66	89.56	94.2	87.53	92.4	88.53
2	38,724	40,541	38,004	36,117	33,440	89.09	87.99	93.27	86.35	91.13	87.16
3	35,570	37,195	34,823	33,263	30,479	89.43	87.53	93.51	85.69	91.42	86.6
4	32,668	34,225	32,035	30,787	29,210	89.95	91.18	94.24	89.41	92.05	90.29
5	31,581	33,041	30,925	29,629	27,744	89.67	89.71	93.82	87.85	91.7	88.77
Trung bình	36,515	38,151	35,763	34,255	32,282	89.76	89.19	93.81	87.37	91.74	88.27



Hình 4.7 So sánh tham số Precision của mô hình VNS và MMS



Hình 4.8 So sánh tham số Recall của mô hình VNS và MMS



Hình 4.9 So sánh tham số F-Measure của mô hình VNS và MMS

Kết quả sau năm lần thử nghiệm, về tổng quát các tham số Precision, Recall và F-measure trung bình của mô hình xuất VNS đều cao hơn so với mô hình của MMS. Tuy nhiên, nếu xét riêng tham số Precision, có thể nhận thấy rằng độ chính xác Precision của VNS thấp hơn MMS. Cụ thể là lần thử nghiệm thứ 4: VNS: 89,95 – MMS: 91,18; lần thử nghiệm thứ 5: VNS: 89,67 – MMS: 89,71.

Qua khảo sát dữ liệu, chúng tôi nhận thấy VNS tỏ ra hiệu quả hơn khi xử lý những văn bản có chứa nhiều danh từ riêng, văn bản có chứa những từ chệch có trong từ điển, những từ viết tắt xuyên trong tập huấn luyện, văn bản có chứa các mục lục. Tóm lại thì ý nghĩa của những kết quả thí nghiệm trong mô hình của VNS so với mô hình của MMS là có hiệu quả.

CHƯƠNG 5. KẾT LUẬN

5.1 Nhìn xét chung

Trong xử lý ngôn ngữ tự nhiên tiếng Việt, phân đoạn là một trong những công việc hết sức quan trọng. Do tính chất đặc biệt của ngôn ngữ tiếng Việt, công việc này có tính khó cao và phức tạp. Các quy tắc bài toán này sẽ làm tiền đề cho các bài toán liên quan như: dịch máy, tóm tắt văn bản, xử lý ngữ nghĩa, v.v...

Chúng tôi có nhiều mô hình phân đoạn của arabic nhiều nhà nghiên cứu trước. Mỗi phương pháp đều có ưu và nhược điểm riêng. Nhìn chung việc phân đoạn tự động của từ tự nhiên chính xác, ngay cả khi còn ngữ cảnh thì phần lớn [21]. Nguyên nhân một phần là do khả năng xử lý ngôn ngữ trên máy tính của con người còn hạn chế, một phần khác là do bản thân ngôn ngữ luôn biến đổi và phát triển không ngừng theo thời gian, đòi hỏi người làm xử lý ngôn ngữ phải luôn nghiên cứu và cải tiến phương pháp.

Trong các mô hình phân đoạn mà luận văn tham khảo, mô hình sử dụng từ điển kết hợp với các luật phân giải như phương pháp của họthanh MMSeg là khá phù hợp với ngôn ngữ tiếng Việt nên chúng tôi chọn làm hướng tham khảo chính. Tuy nhiên, khi áp dụng trên ngôn ngữ tiếng Việt, vì những đặc thù riêng, MMSeg sẽ không thể sử dụng được.

Qua quá trình nghiên cứu và thực nghiệm, luận văn sẽ xây dựng một mô hình phân đoạn mới, dựa trên mô hình lai với các luật phân giải như phương pháp của họthanh MMSeg kết hợp với các phương pháp khác như: thuật toán phát hiện từ mới, nhận dạng danh từ riêng, so trùng mục, phương pháp theo kê.

Kết quả thực nghiệm sẽ chứng minh mô hình mới có hiệu suất kết quả khả quan, là tiền đề cho những nghiên cứu tiếp theo.

5.2 Kết quả thực nghiệm

Trong khuôn khổ luận văn cao học ngành Công nghệ thông tin, luận văn đã nghiên cứu và đưa ra mô hình phân loại tiếng Việt. Bảng dưới đây là một số kết quả sau:

- Xây dựng cơ sở dữ liệu riêng với 1.068.435 tên. Trong đó có 1,065,613 tên người và 2,822 tên địa danh Việt Nam.
- Xây dựng cơ sở dữ liệu gồm có 3.356 bài báo chèn vào trong 26Mb dữ liệu.
- Thực nghiệm và đánh giá hiệu suất của MMSeg trên ngôn ngữ tiếng Việt với dữ liệu gồm 10,000 câu trích từ VietTreebank.
- Xây dựng mô hình phân loại cho tiếng Việt dựa trên mô hình lai sử dụng nhiều phương pháp tích hợp: phương pháp dựa trên từ vựng, phương pháp dựa trên thống kê, phương pháp so trùng mẫu và phương pháp phân tích ngữ pháp dựa trên các luật.

5.3 Hạn chế và tài liệu

Vì thời gian có hạn, nên tài liệu còn hạn chế trong một số vấn đề sau:

- Nguồn dữ liệu huấn luyện cho mô hình N-Gram còn nhỏ, chỉ giới hạn một số lĩnh vực nhất định, không mang tính phổ quát, phần nào ảnh hưởng đến chính xác của mô hình.
- Từ vựng riêng xây dựng cơ sở dữ liệu còn phải bổ sung thêm nhiều.
- Số lần thực nghiệm và dữ liệu thực nghiệm còn ít nên phần đánh giá còn mang tính chất quan sát.

Tuy nhiên, các thành quả nghiên cứu, những hạn chế và tài liệu còn nhiều thiếu sót, em rất mong tiếp tục nhận được sự chỉ dẫn của quý thầy cô.

5.4 Hướng phát triển của tài

Qua quá trình khảo sát và thực nghiệm, chúng tôi nhận thấy bản luật phân gii nh p nh ng c a h th ng MMSeg mà mô hình đang áp dụng là chưa cho ngôn ngữ tiếng Việt. Các phân o n sai do nh p nh ng còn r t nhi u và ph c t p. Các luật hi n t i ch y u ch x lý nh p nh ng đ a trên xác xu t c a t và c m t .

Trong tương lai, chúng tôi muốn tiếp tục thêm các luật phân gii nh p nh ng m i có xét n khía c nh ng ngh a và ng pháp. C th , trong nh ng n l c ti p theo, chúng tôi muốn áp dụng thêm cây phân tích cú pháp VietTreebank vào quá trình phân gii nh p nh ng, nh m nâng cao h n chính xác khi phân o n t .

TÀI LI U THAM KH O

- [1] Tr n Ng c Anh, Nguy n Nh t An. (2011). L a ch n t p gán nh n ranh gi i t cho mô hình Markov n trong bài toán tách t ti ng Vi t.
- [2] Luu Tuan Anh, Yamamoto Kazuhide. (2012). A pointwise approach for Vietnamese Diacritics Restoration. 2012 International Conference on Asian Language Processing, pp.189 – 192.
- [3] Ngoc Anh Tran, Thanh Tinh Dao, Phuong Thai Nguyen. (2012). An effective context-based method for Vietnamese-word segmentation. IEEE 9th, pp.34-40.
- [4] D ng H u Biên. (2010). Giáo trình c s ngôn ng h c, Hà L t.
- [5] Nguy n Tài C n. (1975). Ng pháp ti ng Vi t, Ti ng - T ghép - o n ng , Nxb Khoa h c xã h i, Hà N i.
- [6] Chih-Hao Tsai. (1996). MMSEG: A Word Identification System for Mandarin Chinese Text Based on Two Variants of the Maximum Matching Algorithm. www.casper.beckman.uiuc.edu/~ctsai4/chinese/wordseg/mmseg.html.
- [7] Mai Ng c Ch . (1997). C s ngôn ng h c và ti ng Vi t, Nxb Giáo d c, trang 91–105.
- [8] Mai Ng c Ch , V c Nghi u và Hoàng Tr ng Phi n. (1997). C s ngôn ng h c và ti ng Vi t. Nxb Giáo d c, trang 142–152.
- [9] Nguy n c Dân. (1987). Logic ng ngh a cú pháp. NXB H&TH chuyên nghi p, Hà N i.
- [10] inh i n, H B o Qu c. (2008). V n v ranh gi i t trong ng li u song ng Anh-Vi t
- [11] Nguy n Thi n Giáp. (1998). D n lu n Ngôn ng h c, Nxb Giáo d c, trang 298–305.
- [12] Cao Xuân H o. (2003). Ti ng Vi t - M y v n Ng âm, Ng pháp, Ng ngh a. Nxb Khoa h c xã h i.

- [13] Lê Trung Hi u, Lê Anh V , Lê Trung Kiên. (2013). Áp d ng xác su t th ng kê và quá trình máy t h c cho bài toán phân tách t v n b n ti ng Vi t. T p chí Khoa h c & Công ngh i h c Duy Tân s 6, trang 32-38.
- [14] Hla Hla Htay, Kavi Narayana Murthy. (2008). Myanmar Word Segmentation using Syllable level Longest Matching. Proceedings of the 6th Workshop on Asian Language Resources, pp.41-48.
- [15] H. P. Lê, T. M. H. Nguyen, A. Roussanaly and T. V. Ho. (2008). A hybrid approach to word segmentation of Vietnamese texts. In 2nd International Conference on Language and Automata Theory and Applications, Tarragona, Spain, pp.240-249.
- [16] Thi Minh Huyen Nguyen, Laurent Romary, Mathias Rossignol, Xuan Luong Vu. (2006). A lexicon for Vietnamese language processing. Language Reseourse Evaluation - Volume 40, pp.291-309.
- [17] Nguy n Th Minh Huy n, Hoàng Th Tuy n Linh, V Xuân L ng. (2009). H ng d n nh n di n n v t trong v n b n ti ng Vi t.
- [18] Jin Kiat Low, Hwee Tou Ng and Wenyan Guo. (2005). A Maximum Entropy Approach to Chinese Word Segmentation. In Proceedings of the Fourth SIGHAN Workshop on Chinese Language Processing, pp.161-164.
- [19] Jurafsky and Martin. (2009). Speech and Language Processing: An Introduction to Speech Recognition. Computational Linguistics and Natural Language Processing, SE, Prentice Hall, pp.934.
- [20] Masaaki Nagata. (1997). A self-organizing Japanese word segmenter using heuristic word identification and re-estimation. In Joe Zhou and Kenneth Church, editors, Proceedings of the Fifth Workshop on Very Large Corpora, pp.203-215.
- [21] Richard Sproat, Chilin Shih, William Gale, Nancy Chang. (1994). A stochastic finite-state word-segmentation algorithm for Chinese. ACL '94 Proceedings of the 32nd annual meeting on Association for Computational Linguistics, pp.66-73.

- [22] Phuong-Thai Nguyen, Xuan-Luong Vu, Thi-Minh-Huyen Nguyen, Van-Hiep Nguyen, Hong-Phuong Le. (2009). Building a large syntactically-annotated corpus of Vietnamese. Proceedings of the Third Linguistic Annotation Workshop, Suntec, Singapore, pp.182-185.
- [23] Nguyễn Kim Thìn. (1997). Nghiên cứu ngữ pháp tiếng Việt. NXB GD, trang 28.
- [24] Theeramunkong, T., Usanavasin, S. (2001). Non-dictionary-based Thai word segmentation using decision trees. The first international conference on Human language technology research. New Jersey, USA (2001), pp.1-5.
- [25] Truc-Vien T. Nguyen., Tru H. Cao. (2007). VN-KIM IE: Automatic extraction of vietnamese named-entities on the web. New Generation Computing May 2007, Volume 25, Issue 3, pp 277-292.
- [26] Trần Ngọc Tuấn. (2002). Phân loại tiếng Việt dùng Corpus và các mô hình thống kê, luận văn thạc sĩ, Viện Công nghệ Thông tin và Truyền thông Đại học Bách Khoa TP.H Chí Minh.