

ĐẠI HỌC THÁI NGUYÊN
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ THÔNG TIN VÀ TRUYỀN THÔNG

MA NGỌC KHÁNH

**NGHIÊN CỨU PHƯƠNG PHÁP RÚT GỌN VĂN BẢN VÀ
CHUYỂN ĐỔI CÚ PHÁP NGÔN NGỮ KÝ HIỆU VIỆT NAM**

Chuyên ngành: Khoa học máy tính
Mã số: 8480101

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

Người hướng dẫn khoa học: PGS.TS. PHÙNG TRUNG NGHĨA

THÁI NGUYÊN, 2018

MỤC LỤC

MỞ ĐẦU	1
CHƯƠNG 1. TỔNG QUAN VỀ NGÔN NGỮ KÝ HIỆU VIỆT NAM	3
1.1. Tổng quan về ngôn ngữ ký hiệu.....	3
1.1.1. Khái niệm ngôn ngữ ký hiệu.....	3
1.1.2. Đặc điểm ngôn ngữ học của ngôn ngữ ký hiệu	6
1.1.3. Vai trò của ngôn ngữ kí hiệu với cộng đồng người khiếm thính....	8
1.2. Tổng quan về ngôn ngữ ký hiệu Việt Nam.....	9
1.3. Tính rút gọn trong ngôn ngữ ký hiệu Việt Nam	11
1.4. Trật tự cú pháp trong ngôn ngữ ký hiệu Việt Nam.....	12
1.5. Ứng dụng công nghệ thông tin trong dịch ngôn ngữ ký hiệu	13
1.6. Mục tiêu của luận văn.	15
CHƯƠNG 2. LUẬT RÚT GỌN VĂN BẢN VÀ CHUYỂN ĐỔI CÚ PHÁP ĐỐI VỚI NGÔN NGỮ KÝ HIỆU VIỆT NAM.....	16
2.1. Thu thập luật rút gọn trong ngôn ngữ ký hiệu Việt Nam	16
2.1.1. Thu thập luật rút gọn giới từ và liên từ	16
2.1.2. Thu thập luật rút gọn các từ tính thái.....	22
2.2. Thu thập luật chuyển đổi cú pháp trong ngôn ngữ ký hiệu Việt Nam. 25	
2.2.1. Vấn đề về xây dựng ngân hàng câu được chú giải cú pháp.....	25
2.2.2. Tổng kết những đặc điểm về trật tự cú pháp ngôn ngữ kí hiệu Việt Nam	31
2.3. Cơ sở dữ liệu văn bản tiếng Việt.....	33
2.3.1. Đặc trưng của văn bản tiếng Việt.....	33
2.3.2. Phân tích dữ liệu văn bản tiếng Việt.....	36
2.3.3. Các vấn đề về phân tích cú pháp trong Tiếng Việt.....	38

CHƯƠNG 3. XÂY DỰNG HỆ THỐNG RÚT GỌN VĂN BẢN VÀ CHUYỂN ĐỔI CÚ PHÁP NGÔN NGỮ KÝ HIỆU VIỆT NAM.....	47
3.1. Môi trường thực nghiệm hệ thống rút gọn văn bản và chuyển đổi cú pháp	47
3.2. Các công cụ hỗ trợ thực nghiệm	47
3.2.1. Công cụ TreeBank Editor	47
3.2.2. Bộ phân tích cú pháp Bikel	47
3.3. Cài đặt thuật toán rút gọn văn bản	48
3.3.1. Thuật toán rút gọn văn bản trong ngôn ngữ ký hiệu Việt Nam....	48
3.3.2. Đánh giá thực nghiệm	50
3.4. Cài đặt thuật toán chuyển đổi cú pháp	53
3.4.1. Xây dựng cây chuyển đổi cú pháp tương ứng trong ngôn ngữ kí hiệu.....	53
3.4.2. Cài đặt thuật toán	56
3.4.3. Đánh giá, kết quả thực nghiệm	57
3.5. Xây dựng phần mềm thực nghiệm rút gọn văn bản và chuyển đổi cú pháp	59
KẾT LUẬN	61
TÀI LIỆU THAM KHẢO	62

DANH MỤC HÌNH

Hình 1.1. Ngôn ngữ kí hiệu trong hệ thống Arthrological.....	5
Hình 2.1. Quá trình gán nhãn.....	31
Hình 2.2. Cây cú pháp của câu "tôi nhìn cô gái với chiếc ống nhòm"	39
Hình 2.3. Dẫn xuất phân tích top - down.....	42
Hình 2.4. Dẫn xuất phân tích bottom - up.....	45
Hình 3.1. Sơ đồ thuật toán rút gọn văn bản	50
Hình 3.2. Cấu trúc cây cú pháp chuyển đổi tương ứng sang dạng NNKH câu đơn	53
Hình 3.3. Cấu trúc cây cú pháp chuyển đổi tương ứng sang dạng NNKH câu phủ định dạng 1.....	53
Hình 3.4. Cấu trúc cây cú pháp chuyển đổi tương ứng sang dạng NNKH câu phủ định dạng 2.....	54
Hình 3.5. Cấu trúc cây cú pháp chuyển đổi tương ứng sang dạng NNKH câu nghi vấn dạng 1	54
Hình 3.6. Cấu trúc cây cú pháp chuyển đổi tương ứng sang dạng NNKH câu nghi vấn dạng 2	54
Hình 3.7. Cấu trúc cây cú pháp chuyển đổi tương ứng sang dạng NNKH câu đơn có bao gồm số từ.	55
Hình 3.8. Sơ đồ thuật toán chuyển đổi cú pháp.....	56
Hình 3.9. Kết quả dịch tự động câu tiếng Việt sang dạng câu đúng ngữ pháp trong ngôn ngữ kí hiệu Việt Nam.....	58
Hình 3.10. Giao diện phần mềm thực nghiệm rút gọn và chuyển đổi cú pháp.....	59
Hình 3.11. Giao diện phần mềm thực nghiệm rút gọn và chuyển đổi cú pháp.....	60
Hình 3.12. Giao diện phần mềm thực nghiệm rút gọn và chuyển đổi cú pháp.....	60

DANH MỤC BẢNG BIỂU

Bảng 1.1. So sánh câu tiếng việt và câu ngôn ngữ ký hiệu	12
Bảng 2.1. Một số mẫu câu rút gọn giới từ và liên từ	22
Bảng 2.2. Tập nhãn từ loại	26
Bảng 2.3. Tập nhãn cụm từ	28
Bảng 2.4. Nhãn mệnh đề	28
Bảng 2.5. Nhãn chức năng cú pháp.....	29
Bảng 2.6. Nhãn chức năng trạng ngữ.....	30
Bảng 2.7. Bảng các thành phần âm tiết.....	34
Bảng 3.1. Điểm số của BLEU.....	52
Bảng 3.2. Điểm BLEU đánh giá tập dữ liệu của thuật toán chuyển đổi cú pháp NNKH	58

LỜI CAM ĐOAN

Tôi là: **Ma Ngọc Khánh**

Lớp: CK15

Khoá học: 2016 - 2018

Chuyên ngành: Khoa học máy tính

Mã số chuyên ngành: 8480101

Cơ sở đào tạo: Trường Đại học Công nghệ thông tin và Truyền thông Thái Nguyên.

Giáo viên hướng dẫn: **PGS.TS. Phùng Trung Nghĩa**

Tôi xin cam đoan luận văn “*Nghiên cứu phương pháp rút gọn văn bản và chuyển đổi cú pháp ngôn ngữ ký hiệu Việt Nam*” này là công trình nghiên cứu của riêng tôi dưới sự hướng dẫn của PGS.TS. Phùng Trung Nghĩa. Các số liệu sử dụng trong luận văn là trung thực. Các kết quả nghiên cứu được trình bày trong luận văn chưa từng được công bố tại bất kỳ công trình nào khác.

Thái Nguyên, ngày 30 tháng 5 năm 2018

HỌC VIÊN

Ma Ngọc Khánh

LỜI CẢM ƠN

Học viên xin gửi lời cảm ơn chân thành tới Thầy hướng dẫn PGS.TS. Phùng Trung Nghĩa, Trường Đại học Công nghệ Thông tin và Truyền thông - Đại học Thái Nguyên, người đã tận tình hướng dẫn giúp học viên hoàn thành luận văn tốt nghiệp.

Học viên cũng xin gửi lời cảm ơn sâu sắc đến các thầy cô giáo của Trường Đại học Công nghệ Thông tin và Truyền thông - Đại học Thái Nguyên, cùng các thầy cô giáo của Viện Công nghệ thông tin - Viện khoa học Việt Nam đã nhiệt tình giảng dạy, truyền đạt kiến thức cho học viên trong suốt 2 năm học để học viên có thể hoàn thành được luận văn của mình.

Ma Ngọc Khánh

MỞ ĐẦU

Hiện nay theo thống kê, Việt Nam có khoảng trên 2.5 triệu người khiếm thính [1]. Do khả năng nghe bị suy giảm nên khả năng giao tiếp bằng lời ở cộng đồng người khiếm thính bị hạn chế rất nhiều. Để thay thế cho khả năng giao tiếp bằng tiếng nói, ngôn ngữ ký hiệu, là ngôn ngữ tượng hình sử dụng biểu diễn, chuyển động của bàn tay, cơ thể, và sắc thái biểu cảm của khuôn mặt đã ra đời. Tuy nhiên, việc sử dụng ngôn ngữ ký hiệu chưa phát huy được hiệu quả giúp người khiếm thính hòa nhập được với xã hội do việc giao tiếp giữa người khiếm thính và người nghe tốt còn gặp nhiều khó khăn.

Trên thế giới hiện nay đã và đang nghiên cứu phát triển và đưa ra nhiều dịch vụ thông dịch và sản phẩm công nghệ nhằm hỗ trợ người khiếm thính trong giao tiếp xã hội như máy trợ thính dành cho người nghe kém, gắng tay chuyển đổi ngôn ngữ ký hiệu thành giọng nói [9], các phần mềm dịch từ văn bản/giọng nói sang ngôn ngữ ký hiệu hay các từ điển tra cứu ngôn ngữ ký hiệu online [12], v.v... Tuy nhiên mỗi một nghiên cứu hay sản phẩm đều có những hạn chế và chưa đáp ứng được việc hỗ trợ trong giao tiếp hai chiều giữa người khiếm thính và người nghe tốt trong thực tế.

Việc nghiên cứu xử lý ngôn ngữ ký hiệu trên máy tính ở nước ta còn rất mới mẻ. Chúng ta chưa thực sự có một hệ thống ngôn ngữ đồng nhất cho ngôn ngữ ký hiệu tiếng Việt [6]. Bên cạnh vấn đề ngôn ngữ học, việc phát triển sản phẩm ứng dụng công nghệ để phát huy ngôn ngữ ký hiệu nhằm nâng cao trình độ, tiếp nhận thông tin, khả năng giao tiếp cho người khiếm thính lại càng ít và kém hiệu quả.

Với sự quan tâm đặc biệt của Đảng và Nhà nước, đã có nhiều trường học, trung tâm hỗ trợ dạy học và việc làm riêng cho người khiếm thính. Vì vậy việc nghiên cứu về các thuật toán và xây dựng phần mềm rút gọn văn

bản, chuyển đổi cú pháp đối với ngôn ngữ ký hiệu Việt Nam là cần thiết [2]. Do đó tôi chọn đề tài ***“Nghiên cứu phương pháp rút gọn văn bản và chuyển đổi cú pháp ngôn ngữ ký hiệu Việt Nam”***

Mục tiêu của luận văn là Nghiên cứu các lý thuyết đã có để phân tích, đánh giá về các tính chất, các luật rút gọn, chuyển đổi cú pháp đối với ngôn ngữ ký hiệu Việt Nam. Dựa trên các cơ sở lý thuyết và các phân tích, đánh giá sẽ nghiên cứu cài đặt thực nghiệm các thuật toán rút gọn, chuyển đổi cú pháp này và xây dựng phần mềm hỗ trợ rút gọn và chuyển đổi cú pháp ngôn ngữ ký hiệu Việt Nam.

Nội dung chính của luận văn bao gồm 3 chương:

Chương 1. Tổng quan về ngôn ngữ ký hiệu Việt Nam

Chương 2. Luật rút gọn văn bản và chuyển đổi cú pháp đối với ngôn ngữ ký hiệu Việt Nam.

Chương 3. Xây dựng hệ thống rút gọn văn bản và chuyển đổi cú pháp ngôn ngữ ký hiệu Việt Nam.

Khi viết báo cáo này học viên đã cố gắng để đạt được những mục tiêu và định hướng nghiên cứu đề ra ban đầu, song điều kiện thời gian và năng lực còn hạn chế nên không tránh khỏi thiếu sót. Học viên mong nhận được sự góp ý của thầy giáo hướng dẫn, thầy cô giáo để học viên có được những kinh nghiệm thực tế và bổ ích để sau này có thể xây dựng được một chương trình hoàn thiện hơn.

CHƯƠNG 1

TỔNG QUAN VỀ NGÔN NGỮ KÍ HIỆU VIỆT NAM

1.1. Tổng quan về ngôn ngữ ký hiệu

1.1.1. Khái niệm ngôn ngữ ký hiệu

Ngôn ngữ ký hiệu (hay ngôn ngữ dấu hiệu, thủ ngữ) là ngôn ngữ dùng những biểu hiện của bàn tay thay cho âm thanh của tiếng nói. Ngôn ngữ ký hiệu do người điếc tạo ra nhằm giúp họ có thể giao tiếp với nhau trong cộng đồng của mình và tiếp thu tri thức của xã hội. Việc thay thế âm thanh của tiếng nói có thể liên quan đến đồng thời sự kết hợp các hình dạng tay, hướng và chuyển động của bàn tay, cánh tay hoặc cơ thể, và nét mặt để thể hiện trôi chảy những suy nghĩ của người nói. Ngôn ngữ kí hiệu có nhiều điểm tương đồng với ngôn ngữ nói (đôi khi được gọi là "ngôn ngữ bằng miệng" - mà phụ thuộc chủ yếu vào âm thanh), đó là lý do tại sao ngôn ngữ học xem xét cả hai dạng ngôn ngữ là ngôn ngữ tự nhiên. Tuy nhiên cũng có một số khác biệt đáng kể giữa các ngôn ngữ ký hiệu và ngôn ngữ nói. Đặc biệt không nên nhầm lẫn ngôn ngữ kí hiệu với ngôn ngữ cơ thể, là một loại giao tiếp phi ngôn ngữ.

Bất cứ đâu trong cộng đồng người khiếm thính trên thế giới, ngôn ngữ ký hiệu đều được phát triển. Ngôn ngữ kí hiệu không chỉ được sử dụng bởi người điếc mà nó cũng được sử dụng bởi những người có thể nghe thấy, nhưng thể chất bị hạn chế để có thể nói chuyện bình thường. Ngôn ngữ kí hiệu có những thuộc tính ngôn ngữ riêng biệt. Hiện nay, hàng trăm ngôn ngữ ký hiệu được sử dụng trên thế giới và phát triển trong cộng đồng người khiếm thính ở tất cả các quốc gia. Một số ngôn ngữ ký hiệu có được công nhận pháp lý, trong khi một số khác thì chỉ mang tính cục bộ, địa phương.

Một quan niệm sai lầm phổ biến là tất cả các ngôn ngữ ký hiệu là trên toàn thế giới là hoàn toàn giống nhau hoặc ngôn ngữ ký hiệu là một ngôn ngữ