

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

**BÙI THANH THUY**

**NGHIÊN CỨU VỀ DỊCH MÁY THỐNG KÊ DỰA VÀO CỤM  
TỪ VÀ ỨNG DỤNG DỊCH TỪ TIẾNG VIỆT SANG TIẾNG  
ANH**

**LUẬN VĂN THẠC SĨ: KHOA HỌC MÁY TÍNH**

**Thái Nguyên - 2015**

*Số hoá bởi Trung tâm Học liệu – ĐHTN <http://www.lrc.tnu.edu.vn>*

## **LỜI CAM ĐOAN**

Tôi xin cam đoan toàn bộ nội dung trong luận văn này do tôi tự nghiên cứu, đọc, dịch tài liệu, tổng hợp và thực hiện. Trong luận văn tôi có sử dụng một số tài liệu tham khảo như đã trình bày trong phần tài liệu tham khảo.

Người viết luận văn

*Bùi Thanh Thủy*

## LỜI CẢM ƠN

Đầu tiên tôi xin gửi lời cảm ơn chân thành đến TS. **Nguyễn Văn Vinh** đã tận tình hướng dẫn, chỉ bảo cho tôi trong suốt quá trình làm luận văn. Em cũng xin cảm ơn anh **Trần Hồng Việt**, nghiên cứu sinh Trường đại học công nghệ, giảng viên Trường Đại học Kinh tế kỹ thuật công nghiệp đã giúp đỡ em trong quá trình làm luận văn

Tôi cũng xin gửi lời cảm ơn đến các thầy cô trường Đại học Công nghệ thông tin và Truyền thông – Đại học Thái Nguyên, các thầy cô Viện Công nghệ thông tin đã truyền đạt những kiến thức và giúp đỡ tôi trong suốt quá trình học của mình.

Tôi cũng xin gửi lời cảm ơn tới Ban giám hiệu, Phòng Đào tạo, các đồng nghiệp trường Cao đẳng nghề Phú Thọ, gia đình và bạn bè những người đã động viên tạo mọi điều kiện giúp đỡ tôi để hoàn thành luận văn.

# MỤC LỤC

<i>LỜI CAM ĐOAN</i> .....	1
<i>LỜI CẢM ƠN</i> .....	3
<i>MỤC LỤC</i> .....	4
<i>MỞ ĐẦU</i> .....	1
1. Lý do chọn đề tài .....	1
3. Hướng nghiên cứu của đề tài .....	2
4. Phương pháp nghiên cứu .....	2
5. Ý nghĩa khoa học của đề tài.....	3
6. Cấu trúc luận văn .....	3
<i>CHƯƠNG 1 – TỔNG QUAN VỀ DỊCH MÁY</i> .....	4
1.1. Khái niệm về hệ dịch máy .....	4
1.1.1. Định nghĩa .....	4
1.1.2. Vai trò của dịch máy .....	4
1.1.3. Sơ đồ tổng quan của một hệ dịch máy .....	5
1.2. Dịch máy thống kê là gì? .....	6
1.2.1. Tổng quan về dịch thống kê .....	6
1.2.1.1. Mô hình kênh nguồn .....	6
1.2.1.2. Cách tiếp cận Maximum và mô hình giống hàng .....	7
1.2.1.3. Nhiệm vụ trong dịch thống kê .....	7
1.2.1.4. Ưu điểm của phương pháp dịch thống kê .....	8
1.3. Phân loại dịch máy thống kê.....	12
1.3.1. Dịch máy thống kê dựa vào từ (word-based).....	12
1.3.2. Dịch máy thống kê dựa trên cụm từ (phrase-based).....	12
1.3.3. Dịch máy thống kê dựa trên cú pháp .....	13
1.3.4. Một số công cụ và các nhóm nghiên cứu trên Internet về SMT.....	13
<i>CHƯƠNG 2 – MÔ HÌNH DỊCH MÁY DỰA TRÊN CỤM TỪ VÀ ÁP DỤNG CHO NGÔN NGỮ VIỆT _ ANH</i> .....	15
2.1. Giới thiệu mô hình dịch máy dựa trên cụm từ.....	15
2.2. Kiến trúc của mô hình dịch dựa trên cụm từ .....	15
2.2.1. Mô hình log-linenear .....	16
2.2.2. Mô hình dịch .....	20
2.2.3. Mô hình ngôn ngữ.....	24

2.3. Giải mã.....	29
2.3.1. Đặt vấn đề.....	29
2.3.2. Mô tả thuật toán.....	30
2.4. Đánh giá chất lượng dịch.....	33
2.5. Phần mềm mã nguồn mở Moses.....	34
2.6. Quá trình giải mã.....	37
2.6.1. Huấn luyện cực tiểu sai số (MERT).....	37
2.7. Áp dụng với cặp ngôn ngữ Việt – Anh.....	40
2.7.1. Xây dựng ngữ liệu (corpus).....	40
2.7.1.1. Tạo corpus thô.....	40
2.7.1.2. Tạo corpus song ngữ.....	42
2.7.2. Phân đoạn từ trong corpus tiếng Việt (Segmentation).....	42
2.7.2.1. Phương pháp Maximum Matching.....	43
2.7.2.2. Phương pháp Transformation-based Learning (TBL).....	43
2.7.2.3. Phương pháp dựa trên thống kê từ Internet và thuật giải di truyền.....	44
2.7.3. Đánh giá theo dữ liệu huấn luyện.....	44
2.7.4. Đánh giá theo mô hình giống hàng từ trong văn bản.....	44
<b>CHƯƠNG 3 – THỬ NGHIỆM VÀ ĐÁNH GIÁ.....</b>	<b>46</b>
3.1. Công cụ tiền xử lý cho hệ dịch.....	46
3.1.1. Môi trường triển khai.....	46
3.1.2. Chuẩn bị dữ liệu đầu vào cho hệ dịch.....	46
3.1.3. Huấn luyện mô hình dịch.....	46
3.2. Kết quả thực nghiệm.....	47
3.2.1. Dữ liệu đầu vào.....	47
3.2.2. Quá trình chuẩn bị dữ liệu và huấn luyện.....	48
3.2.2.1. Chuẩn bị dữ liệu.....	48
<b>KẾT LUẬN.....</b>	<b>53</b>
<b>TÀI LIỆU THAM KHẢO.....</b>	<b>54</b>
<b>Tài liệu tiếng Việt.....</b>	<b>54</b>
<b>Tài liệu tiếng Anh.....</b>	<b>54</b>

## DANH MỤC CÁC HÌNH

<i>Hình 1.1: Sơ đồ tổng quan của hệ dịch máy.....</i>	<i>6</i>
<i>Hình 1.2: Chu kì phát triển của hệ thống dịch thống kê.....</i>	<i>10</i>
<i>Hình 2.1. Kiến trúc mô hình dịch dựa trên cụm từ.....</i>	<i>15</i>
<i>Hình 2.2: Ví dụ về mô hình dóng hàng.....</i>	<i>20</i>
<i>Hình 2.3: Thuật toán giải mã <math>A^*</math> cho dịch máy.....</i>	<i>31</i>
<i>Hình 2.4: Giải thuật tìm kiếm beam sử dụng đa ngăn xếp trong Pharaoh.....</i>	<i>32</i>

## MỞ ĐẦU

### 1. Lý do chọn đề tài

Trong quá trình phát triển và hội nhập văn hóa, kinh tế thế giới. Quá trình giao lưu giữa người Việt Nam và người nước ngoài ngày càng nhiều dẫn đến khó khăn trong quá trình giao tiếp và sử dụng văn bản tài liệu tiếng Anh. Hiện nay có nhiều hệ thống tự động dịch miễn phí trên mạng như: google translate, vietgle, vdict, lạc việt,... Những hệ thống này cho phép dịch tự động các văn bản với một cặp ngôn ngữ chọn trước (ví dụ dịch từ tiếng Anh sang tiếng Việt). Điều ấy cho thấy sự phát triển của dịch máy càng ngày càng tiến gần hơn đến ngôn ngữ tự nhiên của con người.

Vào những năm gần đây, dịch máy nói chung, dịch máy thống kê nói riêng được phát triển mạnh và ứng dụng rộng rãi. Kết quả thực tế của hệ thống dịch này rất tốt. Ngôn ngữ của máy dịch ngày càng gần với ngôn ngữ của người. Ngoài ra cùng với hệ thống dịch máy thống kê, các sản phẩm ứng dụng ngày càng nhiều giúp con người trao đổi thông tin dễ dàng hơn, tốc độ nhanh hơn và cùng với nhiều ngôn ngữ hơn.

Hiện nay, phương pháp dịch thống kê dựa trên cụm từ là phương pháp cho kết quả dịch tốt nhất hiện nay. Điều này được thể hiện của qua các hệ dịch máy của Google, Vietgle. Hơn nữa việc dịch giữa tiếng Việt sang tiếng Anh là rất cần thiết khi khối lượng văn bản tiếng Anh ngày càng lớn trong thời kỳ Việt Nam hội nhập sâu rộng với quốc tế.

Chính vì lý do đó, tôi lựa chọn và thực hiện đề tài “Nghiên cứu về dịch thống kê dựa vào cụm từ và áp dụng cho dịch từ tiếng Việt sang tiếng Anh”.

## 2. Đối tượng và phạm vi nghiên cứu

*Đối tượng nghiên cứu:*

- Nghiên cứu về các phương pháp, mô hình dịch máy thống kê
- Thử nghiệm và đánh giá kết quả dịch từ tiếng Việt sang tiếng Anh

*Phạm vi nghiên cứu:*

Đề tài tập trung vào nghiên cứu phương pháp dịch thống kê dựa vào cụm từ và ứng dụng dịch tài liệu, văn bản tiếng Việt, tiếng Anh.

## 3. Hướng nghiên cứu của đề tài

- Nghiên cứu, tìm hiểu, phân tích về dịch máy thống kê trên cơ sở cụm từ.
- Cài đặt thử nghiệm tối ưu hóa cụm từ bằng hệ dịch máy thống kê Moses

## 4. Phương pháp nghiên cứu

- Tìm hiểu các hệ dịch tự động đã có để tìm ra các phương pháp dịch máy mà các hệ dịch đang sử dụng.
- Nghiên cứu và đánh giá các phương pháp dịch máy, những ưu điểm và hạn chế, sau đó tìm ra phương pháp có hiệu quả và đề xuất áp dụng cho bài toán đề tài đặt ra.
- Nghiên cứu các phương pháp đánh giá chất lượng dịch máy để đánh giá hiệu quả dịch cho hệ thống đề tài đã xây dựng.



## **5. Ý nghĩa khoa học của đề tài**

*Ý nghĩa khoa học:*

Dịch máy dựa vào cụm từ là một trong những phương pháp dịch máy hiệu quả nhất hiện nay. Hơn nữa dữ liệu văn bản ngày càng lớn và đa dạng. chính vì vậy nghiên cứu về hệ dịch dựa vào cụm từ và ứng dụng cho dịch Việt – Anh có ý nghĩa khoa học cũng như thực tiễn

## **6. Cấu trúc luận văn**

- + Chương 1: Tổng quan về dịch máy
- + Chương 2: Dịch máy thống kê dựa vào cụm từ và áp dụng cho ngôn ngữ Việt \_ Anh
- + Chương 3: Thực nghiệm, đánh giá
- + Kết luận

## CHƯƠNG 1 – TỔNG QUAN VỀ DỊCH MÁY

### 1.1. Khái niệm về hệ dịch máy

#### 1.1.1. Định nghĩa

Các hệ dịch máy (machine translation system-MT) là các hệ thống sử dụng máy tính để dịch từ một thứ tiếng (trong ngôn ngữ tự nhiên) sang một hoặc vài thứ tiếng khác.

Ngôn ngữ của văn bản cần dịch được gọi là ngôn ngữ nguồn, ngôn ngữ của văn bản đã dịch ra được gọi là ngôn ngữ đích.

#### 1.1.2. Vai trò của dịch máy

Hiện nay trên thế giới có khoảng hơn 5000 ngôn ngữ khác nhau, với một số lượng ngôn ngữ lớn như vậy đã gây ra rất nhiều khó khăn trong việc trao đổi thông tin, trong giao tiếp, đồng thời ngăn cản sự phát triển của thương mại và mậu dịch quốc tế.

Với những khó khăn như vậy con người đã phải dùng đến một đội ngũ phiên dịch khổng lồ, để dịch các văn bản, tài liệu, lời nói, ngôn ngữ từ tiếng nước này sang tiếng nước khác. Những công việc đó mang tính chất thủ công, tỉ mỉ đòi hỏi người dịch phải làm mất rất nhiều thời gian và công sức, trong khi khối lượng văn bản cần dịch ngày càng nhiều.

Để khắc phục được những nhược điểm trên con người đã nghĩ đến việc thiết kế một mô hình tự động trong công việc dịch ngôn ngữ, do đó ngay từ khi xuất hiện chiếc máy tính điện tử đầu tiên ( năm 1946) người ta đã tiến hành nghiên cứu về dịch máy. Việc đưa ra mô hình tự động cho việc dịch đã và đang được phát triển, mặc dù chưa giải quyết được triệt để lớp ngôn ngữ tự nhiên. Nhưng sự ra đời của chúng đã khẳng định được lợi ích to lớn về mặt chiến lược và phát triển kinh tế, đồng thời các vấn đề liên quan đến dịch máy