

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



KIỀU CÔNG CHÍNH

TỐI ƯU BẢNG CỤM TỪ ĐỂ CẢI TIẾN DỊCH MÁY THỐNG KÊ

LUẬN VĂN THẠC SĨ KHOA HỌC MÁY TÍNH

Thái Nguyên - 2015

MỤC LỤC

MỞ ĐẦU	1
CHƯƠNG I: DỊCH MÁY THỐNG KÊ TRÊN CƠ SỞ CỤM TỪ.....	3
1.1 Ngôn ngữ tự nhiên.....	3
1.2 Dịch máy.....	3
1.3 Dịch máy thống kê dựa vào cụm từ.....	4
<i>1.3.1 Cơ sở của phương pháp dịch máy thống kê</i>	<i>5</i>
<i>1.3.2 Giống hàng từ, giống hàng thống kê</i>	<i>6</i>
<i>1.3.3 Dịch máy thống kê dựa trên cơ sở cụm từ.....</i>	<i>8</i>
<i>1.3.4 Mục đích của việc dịch máy thống kê trên cơ sở cụm từ.....</i>	<i>11</i>
<i>1.3.5 Đảo cụm từ trong dịch máy thống kê.....</i>	<i>13</i>
<i>1.3.6 Bảng cụm từ trong dịch máy thống kê.....</i>	<i>13</i>
1.4 Mô hình ngôn ngữ	14
CHƯƠNG II: PHƯƠNG PHÁP TỐI ƯU BẢNG CỤM TỪ.....	16
2.1 Quy trình sinh bảng cụm từ	16
2.2 Phương pháp tối ưu bảng cụm từ.....	19
<i>2.2.1 Chỉ số cụm từ nguồn.....</i>	<i>19</i>
<i>2.2.2 Lưu trữ cụm từ mục tiêu.....</i>	<i>20</i>
<i>2.2.3 Nén ngữ liệu song ngữ.....</i>	<i>22</i>
<i>2.2.4 Nén bảng cụm từ.....</i>	<i>27</i>
<i>2.2.5 Mã hóa cụm từ</i>	<i>31</i>
<i>2.2.6 Giải mã cụm từ.....</i>	<i>33</i>
CHƯƠNG III: ĐÁNH GIÁ THỰC NGHIỆM BẢNG HỆ DỊCH MÁY THỐNG KÊ MOSES	36
3.1 Môi trường triển khai	36
3.2 Xây dựng chương trình dịch và thực hiện nén bảng cụm từ.	36
<i>3.2.1 Chuẩn hóa dữ liệu.....</i>	<i>36</i>
<i>3.2.2 Xây dựng mô hình ngôn ngữ, mô hình dịch</i>	<i>37</i>

3.2.3 <i>Nén bảng cụm từ</i>	37
3.2.4 <i>Đánh giá kết quả dịch</i>	38
3.3 Thực nghiệm và đánh giá kết quả dịch tiếng Anh sang tiếng Việt	39
3.3.1 <i>Thực nghiệm dịch với câu đơn giản</i>	43
3.3.2 <i>Thực nghiệm dịch 1 đoạn văn bản từ tiếng Anh-Tiếng Việt</i>	44
3.3.3 <i>Đánh giá kết quả dữ liệu huấn luyện bảng cụm từ</i>	44
3.3.4 <i>Đánh giá kết quả theo cỡ dữ liệu huấn luyện</i>	46
3.3.5 <i>Đánh giá kết quả theo thời gian tải bảng cụm từ</i>	47
PHỤ LỤC	50
1. <i>Kết quả dịch máy đối với câu đơn giản</i>	50
2. <i>Kết quả dịch máy đối với bộ dữ liệu</i>	51
3. <i>Một số công cụ tiền xử lý thường được hay sử dụng trong hệ dịch</i>	52
Tài liệu tham khảo	54

DANH MỤC HÌNH

<i>Hình 1.1: Sơ đồ của hệ dịch bằng phương pháp thống kê.....</i>	5
<i>Hình 1.2: Gióng hàng với những từ tiếng anh độc lập.....</i>	7
<i>Hình 1.3: Gióng hàng với những từ tiếng việt độc lập.....</i>	7
<i>Hình 1.4: Gióng hàng tổng quát</i>	7
<i>Hình 1.5: Mô hình dịch từ Tiếng Anh- Tiếng Việt.</i>	9
<i>Hình 1.6: Mô tả việc giải mã</i>	12
<i>Hình 2.1: Sơ đồ đơn giản hóa bảng cụm từ.....</i>	19
<i>Hình 2.2: Mô tả quá trình tạo cây Huffman</i>	26
<i>Hình 3.1: Dịch câu đơn giản với bảng cụm từ gốc.....</i>	51
<i>Hình 3.2: Dịch câu đơn giản với bảng cụm tối ưu</i>	51
<i>Hình 3.3: Điểm Bleu bộ dữ liệu bảng cụm từ gốc</i>	52
<i>Hình 3.4: Điểm Bleu bộ dữ liệu bảng cụm từ tối ưu.....</i>	52

DANH MỤC BẢNG

<i>Bảng 2.1 : Một số phần tử trong bảng cụm từ.....</i>	<i>18</i>
<i>Bảng 2.2: Bảng mã hóa huffman</i>	<i>27</i>
<i>Bảng 2.3: Bảng tùy chọn mã Simple 9.....</i>	<i>28</i>
<i>Bảng 3.1: Ngữ liệu tiếng việt.</i>	<i>40</i>
<i>Bảng 3.2: Ngữ liệu tiếng anh.</i>	<i>40</i>
<i>Bảng 3.5: Dữ liệu đầu vào</i>	<i>42</i>
<i>Bảng 3.3: So sánh kết quả dịch máy với một câu đơn.</i>	<i>43</i>
<i>Bảng 3.4: So sánh hai phương pháp dịch với đầu vào là một văn bản</i>	<i>44</i>
<i>Bảng 3.5: So sánh dữ liệu bảng cụm từ gốc và bảng cụm sau khi nén</i>	<i>45</i>

DANH MỤC BIỂU ĐỒ

<i>Biểu đồ 3.1: Biểu đồ so sánh 1</i>	45
<i>Biểu đồ 3.2: Biểu đồ so sánh 2.</i>	46
<i>Biểu đồ 3.3: Biểu đồ so sánh 3</i>	48

DANH SÁCH CÁC TỪ VIẾT TẮT

Viết tắt	Đầy đủ
PB-SMT	Cụm từ base Statistical Machine Translation
SMT	Statistical Machine Translation
PR-Enc	Cụm từ Rank Encoding

MỞ ĐẦU

Hiện nay trên thế giới có khoảng 5650 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. Mặt khác, với việc bùng nổ Internet như hiện nay, có một khối lượng văn bản khổng lồ trên Internet mà phần lớn là bằng tiếng Anh. Do tính đa dạng của nó mà việc hiểu các văn bản này hoàn toàn không dễ chút nào. Do đó việc có một hệ dịch tự động Anh-Việt là hết sức cần thiết. Với những khó khăn như vậy người ta đã phải dùng đến một đội ngũ phiên dịch khổng lồ, để dịch các văn bản, tài liệu, lời nói từ tiếng nước này sang tiếng nước khác. Những công việc đó mang tính chất thủ công, nặng nhọc trong khi khối lượng văn bản cần dịch ngày càng nhiều. Để khắc phục những nhược điểm trên hiện nay có rất nhiều những hệ thống tự động dịch miễn phí trên mạng như: systran, google translate, vietgle, vdict... 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) [1]. Điều ấy cho thấy sự phát triển của dịch máy càng ngày càng tiến gần hơn đến ngôn ngữ tự nhiên của con người.

Ngay từ khi xuất hiện chiếc máy tính điện tử đầu tiên người ta đã tiến hành nghiên cứu về dịch máy. Công 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 ích lợi to lớn về mặt chiến lược và kinh tế, đồng thời các vấn đề liên quan đến dịch máy cũng là những chủ đề quan trọng của ngành khoa học máy tính, bởi chúng liên quan đến vấn đề xử lý ngôn ngữ tự nhiên, một trong những vấn đề có ý nghĩa nhất mà trí tuệ nhân tạo có khả năng giải quyết. Người ta tin rằng việc xử lý ngôn ngữ tự nhiên trong đó có dịch máy sẽ là giải pháp cho việc mở rộng cánh cửa đối thoại người-máy, lúc đó con người không phải tiếp xúc với

máy qua những dòng lệnh cứng nhắc nữa mà có thể giao tiếp một cách trực tiếp với máy.

Với sự phát triển mạnh mẽ của dịch máy tự động thì dịch máy thống kê (Statistical Machine Translation) đã chứng tỏ là một hướng tiếp cận đầy tiềm năng bởi ưu điểm vượt trội so với các phương pháp dịch máy dựa trên cú pháp truyền thống. Kết quả thực tế của hệ thống dịch máy thống kê tốt hơn, ngôn ngữ dịch càng ngày càng gần với ngôn ngữ của người, 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. Để dịch hiệu quả thì bảng cụm từ phải lớn chính vì vậy việc lưu trữ và tìm kiếm trong bảng cụm từ là rất quan trọng. Chính vì thế, luận văn này tôi lựa chọn và thực hiện đề tài “**Tối ưu bảng cụm từ để cải tiến dịch máy thống kê**”.

CHƯƠNG I: DỊCH MÁY THỐNG KÊ TRÊN CƠ SỞ CỤM TỪ

Hiện nay dịch máy thông kê dựa trên cơ sở cụm từ là một trong những hướng phát triển đang được rất nhiều người quan tâm. Dịch máy thông kê dựa trên cụm từ nhằm mục đích dịch một văn bản từ ngôn ngữ nguồn sang ngôn ngữ đích dựa vào bảng ngữ cụm từ sau khi thực hiện việc giống hàng từ, giống hàng thống kê, đảo cụm từ... kết hợp với mô hình ngôn ngữ.

1.1 Ngôn ngữ tự nhiên

Ngôn ngữ tự nhiên là những ngôn ngữ được con người sử dụng trong các giao tiếp hàng ngày nghe, nói, đọc, viết. Mặc dù con người có thể dễ dàng hiểu và học các ngôn ngữ tự nhiên, việc làm cho máy hiểu được ngôn ngữ tự nhiên không phải là chuyện dễ dàng. Sở dĩ có khó khăn là do ngôn ngữ tự nhiên có các bộ luật, cấu trúc ngữ pháp phong phú hơn nhiều các ngôn ngữ máy tính, hơn nữa để hiểu đúng nội dung các giao tiếp, văn bản trong ngôn ngữ tự nhiên cần phải nắm được ngữ cảnh của nội dung đó.

Do vậy, để có thể xây dựng được một bộ ngữ pháp, từ vựng hoàn chỉnh, chính xác để máy có thể hiểu ngôn ngữ tự nhiên là một việc rất tốn công sức và đòi hỏi người thực hiện phải có hiểu biết sâu về ngôn ngữ học. Do đó cần phải tìm ra một phương pháp dịch tự động tối ưu để làm giảm công sức trong vấn đề về dịch ngôn ngữ nói chung.

1.2 Dịch máy

Dịch tự động hay còn gọi là dịch máy là một trong những ứng dụng quan trọng của xử lý ngôn ngữ tự nhiên, là sự kết hợp của ngôn ngữ, dịch thuật và khoa học máy tính. Như tên gọi dịch tự động là việc thực hiện dịch một ngôn ngữ đầu vào (ngôn ngữ này gọi là ngôn ngữ nguồn) sang một hoặc nhiều ngôn ngữ khác (gọi là ngôn ngữ đích) bằng các công cụ, phần mềm trên máy tính đã được lập trình sẵn mà không cần có sự can thiệp của con người.