

ĐẠI HỌC QUỐC GIA HÀ NỘI  
TRƯỜNG ĐẠI HỌC CÔNG NGHỆ



CÔNG TRÌNH DỰ THI  
HỘI NGHỊ SINH VIÊN NGHIÊN CỨU KHOA HỌC

---

**MỘT PHƯƠNG PHÁP LAI TRÍCH XUẤT SỰ KIỆN  
VÀ ÁP DỤNG VÀO HỆ THỐNG THEO DÕI TIN TỨC  
TRỰC TUYẾN NewSOMoni**

---

Sinh viên thực hiện  
Nguyễn Minh Hoàng  
Nguyễn Sỹ Quân  
Ngô Quang Hiếu

Cán bộ hướng dẫn  
TS. Phan Xuân Hiếu  
ThS. Trần Mai Vũ

Hà Nội, Ngày 12 tháng 3 năm 2012

## **Tóm tắt nội dung**

Trích chọn thông tin luôn là vấn đề có vai trò cốt yếu khi xây dựng một hệ thống khai phá dữ liệu, đặc biệt trong các hệ thống theo dõi/giám sát thông tin, hệ thống tư vấn tin tức, hệ hỗ trợ ra quyết định. Một trong những bài toán cơ bản (và vô cùng quan trọng) của trích chọn thông tin là trích xuất sự kiện trên dữ liệu lớn. Sự kiện được lấy ra đúng đắn từ kho dữ liệu lớn sẽ giúp các hệ thống khai phá dữ liệu dễ dàng hơn trong việc thực thi nhiệm vụ của mình. Nghiên cứu này sẽ tập trung xem xét một phương pháp trích xuất sự kiện hiệu quả dành cho tiếng Việt với lượng dữ liệu lớn và cách thức áp dụng vào hệ thống theo dõi tin tức trực tuyến cùng những đánh giá để cho thấy phương pháp đưa ra có khả quan. Nhóm tác giả hy vọng kết quả của nghiên cứu sẽ góp phần vào sự phát triển của các hệ thống xử lý tin tức dành cho tiếng Việt.

# Mục lục

Tóm tắt nội dung	ii
Mục lục	iii
Danh sách hình vẽ	iv
Danh sách bảng	v
Danh sách ký hiệu và từ viết tắt	vi
Lời nói đầu	1
<b>1 Giới thiệu</b>	<b>3</b>
1.1 Động lực nghiên cứu . . . . .	3
1.2 Vấn đề nghiên cứu . . . . .	5
1.2.1 Bài toán . . . . .	5
1.2.2 Câu hỏi nghiên cứu . . . . .	5
1.3 Ý nghĩa . . . . .	6
1.3.1 Ý nghĩa khoa học . . . . .	6
1.3.2 Ý nghĩa thực tiễn . . . . .	6
1.4 Thách thức . . . . .	7
1.5 Nghiên cứu liên quan . . . . .	7
1.5.1 Một số nghiên cứu liên quan ở nước ngoài . . . . .	7
1.5.2 Một số nghiên cứu liên quan ở trong nước . . . . .	8
<b>2 Phương pháp trích xuất sự kiện</b>	<b>9</b>
2.1 Sự kiện là gì? . . . . .	9
2.2 Các phương pháp hướng nội dung . . . . .	9

2.3 Các phương pháp hướng dữ liệu . . . . .	9
<b>3 Phương pháp giải quyết</b>	<b>10</b>
3.1 Hệ thống theo dõi tin tức . . . . .	10
<b>A Danh sách các trang tin tức điện tử</b>	<b>11</b>
<b>B XYZ</b>	<b>15</b>
<b>Tài liệu tham khảo</b>	<b>16</b>

## Danh sách hình vẽ

# Danh sách bảng

# Bảng ký hiệu và từ viết tắt

Ký hiệu	Ý nghĩa
ACE	Automatic Content Extraction
DARPA	Defense Advanced Research Project Agency
MUC	Message Understanding Conferences
SIGIR	Special Interest Group on Information Retrieval
SIGKDD	International Conference on Knowledge Discovery and Data Mining
TDT	Topic Detection and Tracking

# Lời nói đầu

Được cộng đồng nghiên cứu khoa học trên toàn thế giới quan tâm rất sớm, trích xuất sự kiện được xem là một bài toán quan trọng trong lĩnh vực trích chọn thông tin (Information Extraction). Từ năm 1987, trích xuất sự kiện đã trở thành đề tài chủ chốt tại hội nghị *Message Understanding Conferences* ngay lần tổ chức đầu tiên (MUC-1) (RB96). Từ đó đến nay, nhiều phương pháp trích xuất sự kiện đã được đưa ra và áp dụng trong các hệ thống thực tế như BioCaster (<http://born.nii.ac.jp/>), HealthMap (<http://healthmap.org>), EpiSpider (<http://www.epispider.org/>), Metro Monitor (<http://www.metromonitor.com/>), ...

Công trình nghiên cứu **Một phương pháp lai trích xuất sự kiện và áp dụng vào hệ thống theo dõi tin tức trực tuyến NewSOMoni** khảo sát một số phương pháp trích xuất sự kiện tiêu biểu có hiệu quả tốt, đang được sử dụng trong nhiều hệ thống theo dõi thông tin. Dựa trên cơ sở đó, chúng tôi nghiên cứu và đề xuất một phương pháp lai nhằm mục đích trích xuất sự kiện trên miền tin tức tiếng Việt và thử nghiệm trên hệ thống theo dõi tin tức trực tuyến NewSOMoni. Phương pháp được đề xuất là sự kết hợp của phương pháp học máy Maximum Entropy và phương pháp trích xuất dựa trên luật với những cải tiến khi áp dụng cho dữ liệu tiếng Việt. Qua tiến hành thực nghiệm, chúng tôi đã thu được kết quả tương đối tốt và ổn định. Điều này chứng tỏ tính đúng đắn của phương pháp đề xuất cũng như tính thực tiễn trong hệ thống theo dõi tin tức trực tuyến, góp phần đưa thông tin đến với người dùng chính xác, kịp thời.

Báo cáo bao gồm bốn chương được mô tả như dưới đây.

**Chương 1.** *Giới thiệu* khái quát chung về động lực thực hiện nghiên cứu, mô tả về bài toán trích xuất sự kiện và cũng nêu một số nghiên cứu liên quan ở trong và ngoài nước.

**Chương 2.** *Phương pháp trích xuất sự kiện* đưa ra 3 phương pháp trích xuất sự kiện phổ biến và có độ chính xác cao. Hơn nữa, chúng tôi cũng phân tích những thuận lợi của từng phương pháp và cách áp dụng chúng vào mô hình giải quyết của mình để đạt được hiệu quả tốt hơn.

**Chương 3.** *Trích xuất sự kiện dựa trên luật kết hợp học máy và hệ thống theo dõi tin tức* trình bày phương pháp trích xuất sự kiện dựa trên luật kết hợp với phương pháp học máy Maximum Entropy—phương pháp chính trong mô hình giải quyết của nghiên cứu này. Đồng thời, mô hình hệ thống theo dõi tin tức cũng sẽ được nêu rõ và phân tích chi tiết.



**Chương 4.** *Thực nghiệm phương pháp trên hệ thống theo dõi tin tức* trình bày quá trình xây dựng hệ thống giám sát tin tức trên cơ sở áp dụng phương pháp đã đề xuất ở Chương 3. Kết quả thực nghiệm và đánh giá hiệu quả sẽ được mô tả kỹ lưỡng trong chương này.

**Phần kết luận** tổng kết, tóm lược nội dung của nghiên cứu và hướng phát triển tiếp theo.

# 1

## Giới thiệu

### 1.1 Động lực nghiên cứu

Thế giới đang thay đổi rất nhanh với sự tham gia của các phương tiện truyền thông xã hội. Mọi thông tin đều có thể đến với người dùng theo nhiều nguồn khác nhau. Tuy nhiên, sử dụng phương tiện truyền thông xã hội riêng lẻ khó có thể cập nhật được kịp thời và chính xác thông tin. Để đáp ứng nhu cầu đó, những hệ thống tổng hợp tin tức lần lượt ra đời giúp cho con người có thể dễ dàng nắm bắt thông tin. Khởi đầu bởi <tên hệ thống đầu tiên trên thế giới>, tiếp sau đó là <tên hệ thống khác 1>, <tên hệ thống khác 2>, ... Vào năm 2005, hệ thống tổng hợp tin tức tự động đầu tiên của Việt Nam ra đời dựa trên thành tựu nghiên cứu *Hệ thống thu thập và tách thông tin ICPS* của hai tác giả Nguyễn Thành Long và Nguyễn Phú Bình đạt giải nhì cuộc thi Trí Tuệ Việt Nam 2002. *Hệ thống xử lý tiếng Việt tự động ePi* được người dùng biết đến với tên BÁO MỚI <sup>1</sup> và nhanh chóng trở thành trang tin tức tổng hợp được nhiều người sử dụng bởi tính tiện lợi và cập nhật. Mặc dù có những ưu điểm như vậy, một hệ thống tổng hợp tin tức vẫn có những yếu điểm chưa thể khắc phục. Thứ nhất, thông tin được thu thập từ những nguồn tin định trước dựa trên giao diện cập nhật của nguồn tin, chưa phân tích sâu về ý nghĩa và tính chất của sự kiện chứa đựng trong thông tin. Thứ hai, tin tức không được trực quan hóa theo xu hướng quan tâm của người dùng. Thông thường, độ ưu tiên quan tâm của người dùng là: thời gian (when) > địa điểm (where) > thông tin gì(what). Hơn nữa, hệ thống tổng hợp tin tức xem xét tất cả các tin từ nguồn tin, sau đó phân lớp vào một lớp đã định nghĩa trước. Bởi tính phong phú của

---

<sup>1</sup>[www.baomoi.com](http://www.baomoi.com)

dạng thông tin, tính chính xác của quá trình phân lớp là một câu hỏi lớn chưa có lời giải đáp thỏa đáng!

Giải quyết nhược điểm của hệ thống tổng hợp tin tức tự động cần có một phương pháp trích xuất sự kiện phù hợp với tiếng Việt và hoạt động ổn định. Từ rất sớm, trích xuất sự kiện đã được cộng đồng khoa học máy tính đầu tư công sức nghiên cứu. Tiêu biểu có thể kể đến hội nghị **Message Understanding Conferences (MUC)** <sup>1</sup> tổ chức lần đầu tiên năm 1987 dưới sự hỗ trợ của DARPA (Quỹ nghiên cứu bộ quốc phòng Hoa Kỳ). Một trong những đóng góp quan trọng của hội nghị MUC là đưa ra phương pháp trích xuất sự kiện theo khung mẫu (*scenario template*) với mục đích chính là lấy ra được sự kiện cùng các thông tin liên quan: tổ chức, đối tượng tham gia (người, sự vật, sự việc). Độ chính xác và độ hồi tưởng của các nghiên cứu tham dự MUC nằm trong khoảng 50% tới 60 %. Để cải tiến hiệu suất trích xuất, nhiều phương pháp trích xuất sự kiện dựa trên luật ra đời như trong nghiên cứu của XYZ , .... Ngoài ra, chương trình nâng cao hiệu quả trích xuất sự kiện **Automatic Content Extraction (ACE)** <sup>2</sup> của Đại học Pennsylvania (Hoa Kỳ) cũng là một chương trình nổi tiếng, thu hút được nhiều nhóm nghiên cứu về trích xuất sự kiện tham gia và có những kết quả rất tích cực. Tuy nhiên, trích xuất sự kiện là một vấn đề mang đặc trưng ngôn ngữ học. Ngôn ngữ ảnh hưởng rất lớn tới hiệu quả của một phương pháp trích xuất. Theo tìm hiểu của chúng tôi, trích xuất sự kiện trên dữ liệu tiếng Việt chưa có nhiều nghiên cứu. Bởi vậy, phương pháp trích xuất sự kiện dành cho tiếng Việt vẫn còn hạn chế cả về chất lượng lẫn số lượng.

Một yếu tố khác đưa chúng tôi đến với đề tài nghiên cứu này là sự thú vị trong xử lý dữ liệu lớn. Theo xu hướng phát triển Công Nghệ Thông Tin hiện đại, thi hành hệ thống với dữ liệu lớn là tất yếu. Các công ty hàng đầu thế giới về Công Nghệ như Microsoft <sup>3</sup>, Google <sup>4</sup>, Oracle <sup>5</sup>, Facebook <sup>6</sup> đều có những chiến lược phát triển lâu dài về xử lý dữ liệu lớn. Cùng với đó, những trường đại học hàng đầu thế giới về khoa học máy tính đều đưa vào trường trình đào tạo của mình khoa học về xử lý dữ liệu lớn như Đại học

---

<sup>1</sup>[http://www-nlpir.nist.gov/related\\_projects/muc](http://www-nlpir.nist.gov/related_projects/muc)

<sup>2</sup><http://projects.ldc.upenn.edu/ace>

<sup>3</sup>[www.microsoft.com](http://www.microsoft.com)

<sup>4</sup>[www.google.com](http://www.google.com)

<sup>5</sup>[www.oracle.com](http://www.oracle.com)

<sup>6</sup>[www.facebook.com](http://www.facebook.com)

Princeton <sup>1</sup> (Hoa Kỳ) , Đại học Stanford <sup>2</sup> (Hoa Kỳ) , Đại học Carnegie Mellon <sup>3</sup> (CMU, Hoa Kỳ) hay Đại học tổng hợp Zurich <sup>4</sup> (Thụy Sĩ). Sự hỗ trợ tuyệt vời về dữ liệu và kỹ thuật từ phía ThS. Trần Mai Vũ đã giúp chúng tôi có thêm động lực và quyết tâm hoàn thành đề tài.

## 1.2 Vấn đề nghiên cứu

### 1.2.1 Bài toán

<Cần phân tích thêm> Những vấn đề phân tích ở phần 1.1 đã đưa nhóm nghiên cứu hướng tới ý tưởng đưa ra phương pháp trích xuất sự kiện phù hợp khi xử lý với dữ liệu tiếng Việt và xây dựng nên một hệ thống theo dõi tin tức trực tuyến mà trong đó trích xuất sự kiện là yếu tố trung tâm. Nghiên cứu đóng góp ở cả hai nội dung: khoa học và ứng dụng. Ý nghĩa của việc giải quyết vấn đề này được trình bày chi tiết ở mục 1.3.

**Đầu vào** của bài toán là một bản ghi tin tức về một trong ba lĩnh vực: tai nạn giao thông, hình sự, cháy nổ. Mỗi bản ghi bao gồm các thông tin: tiêu đề, tóm tắt nội dung, toàn văn tin tức. Tổng số bản ghi là XYZ (số tin tức thu thập được) thu thập từ trang tổng hợp tin tức BÁO MỚI <sup>5</sup>.

**Kết quả mong muốn** của bài toán là có hay không có sự kiện trong bản ghi tin tức. Nếu có thì phải đưa ra được các thông tin liên quan tới sự kiện gồm có: tên sự kiện, thời gian, địa điểm, người, sự vật, sự việc. Sự kiện thu được cũng phải được trực quan hóa trên hệ thống theo dõi tin tức trực tuyến.

### 1.2.2 Câu hỏi nghiên cứu

Nghiên cứu sẽ trả lời ba câu hỏi.

**Thứ nhất** thế nào là trích xuất sự kiện tin tức và những phương pháp thường được sử dụng để làm điều đó?

---

<sup>1</sup><http://www.cs.princeton.edu/courses/archive/spr02/cs493>

<sup>2</sup><http://www.stanford.edu/class/cs246>

<sup>3</sup><http://www.cs.cmu.edu/neill/courses/90866.html>

<sup>4</sup><http://las.ethz.ch/courses/datamining-s12>

<sup>5</sup>[www.baomoi.com](http://www.baomoi.com)

**Thứ hai** tồn tại những khó khăn nào khi áp dụng những phương pháp từ câu hỏi trên vào dữ liệu tiếng Việt và cách giải quyết những khó khăn này?

**Và cuối cùng** một hệ thống theo dõi tin tức có khả thi không?

## 1.3 Ý nghĩa

### 1.3.1 Ý nghĩa khoa học

Về mặt khoa học, chúng tôi đề xuất phương pháp trích xuất sự kiện dựa trên luật kết hợp học máy để thu được sự kiện xảy ra hằng ngày thông qua dữ liệu tin tức tiếng Việt thu thập từ một số nguồn thông tin tin cậy dưới sự cho phép của Bộ Thông Tin và Truyền Thông <sup>1</sup>(xem danh sách trang báo điện tử cung cấp tin tức tại phụ lục A, trang 11). Trong bối cảnh vấn đề trích xuất sự kiện ở trong nước chưa có nhiều nghiên cứu, công trình của chúng tôi sẽ góp phần thôi thúc đề tài thú vị này được quan tâm nhiều hơn. <nói về trích xuất sự kiện ở Việt Nam>.

### 1.3.2 Ý nghĩa thực tiễn

Xét tới phương diện ứng dụng, chúng tôi tiến hành xây dựng một hệ thống theo dõi thông tin trực tuyến. Như đã nói ở mục 1.1, một hệ thống tổng hợp tin tức tự động chưa đủ thông minh để đáp ứng nhu cầu ngày càng cao của người dùng. Bởi thế, trong nghiên cứu này chúng tôi muốn xây dựng một hệ thống theo dõi, giám sát thông tin sự kiện. Bởi quy mô của một công trình sinh viên nghiên cứu khoa học, nhóm chúng tôi tập trung vào ba loại sự kiện thường xảy ra hằng ngày: tai nạn giao thông, hình sự và cháy nổ. Một cách rõ ràng nhất, sự kiện thuộc ba dạng trên sẽ được trích xuất theo các thông tin: tên sự kiện, thời gian/địa điểm diễn ra sự kiện, các nhân tố tham gia sự kiện. Sau đó, sự kiện được trực quan hóa trên bản đồ giúp cho người sử dụng dễ dàng theo dõi. Theo khảo sát của nhóm nghiên cứu, một hệ thống như đã mô tả chưa xuất hiện ở Việt Nam. Đề tài nghiên cứu đóng góp vào việc phổ biến hình thức nắm bắt tin tức mới dễ dùng và trực quan hơn so với các hệ thống cung cấp tin tức truyền thống.

<sup>1</sup><http://mic.gov.vn/vbqppl/Lists/Luat-cong-nghe-thong-tin>

## 1.4 Thách thức

Mặc dù được các nhà khoa học quan tâm nghiên cứu từ rất sớm, trích xuất sự kiện vẫn còn những khó khăn cần phải vượt qua.

Trích xuất sự kiện liên quan mật thiết tới các nghiên cứu về ngôn ngữ học. Lĩnh vực xử lý ngôn ngữ tự nhiên nói chung và xử lý tiếng Việt nói riêng tương đối rộng, tồn tại nhiều bài toán chưa được giải quyết triệt để mà trong đó có xử lý nhập nhằng ngữ nghĩa (Word Sense Disambiguation), bài toán đồng tham chiếu (Co-references) hay việc nhận dạng tính đa hình cấu trúc ngữ pháp trong tiêu đề tin tức (Syntactically Ambiguous Headlines). Ba bài toán trên là những khó khăn cơ bản nhất mà chúng tôi phải giải quyết để đưa ra được phương pháp trích xuất sự kiện phù hợp.

Tính tới thời điểm thực hiện công trình, Việt Nam chưa có nhiều nghiên cứu về trích xuất sự kiện. <cần viết tiếp, đang thiếu dữ kiện về tình hình trong nước>

Ngoài ra, khó khăn trong xử lý dữ liệu lớn cũng là một thách thức mà nhóm nghiên cứu phải đối mặt. Để có thể trích chọn được sự kiện từ tập dữ liệu lớn cần phải tối ưu thuật toán đảm bảo rằng hệ thống có thể hoạt động tốt trong điều kiện tài nguyên cho phép.

## 1.5 Nghiên cứu liên quan

### 1.5.1 Một số nghiên cứu liên quan ở nước ngoài

Kể từ hội nghị MUC lần đầu tiên (1987) cho tới nay, hàng ngàn nghiên cứu về trích xuất sự kiện đã được công bố trong những hội nghị, chương trình có uy tín cao như MUC, SIGKDD<sup>1</sup>, ACM SIGIR<sup>2</sup>, TDT<sup>3</sup>, ACE. Theo Hogenboom. F và các cộng sự, tựu chung lại các công bố này có thể phân loại theo ba hướng tiếp cận chính: phân tích ngữ nghĩa (còn gọi là hướng theo nội dung), học máy-thống kê (hướng theo dữ liệu) và cuối cùng là kết hợp hai cách tiếp cận trên (FFU11).

Giai đoạn cuối thập niên tám mươi, đầu thập niên chín mươi, sự kiện được trích xuất chủ yếu dựa trên các mẫu được tạo sẵn (*scenario template*) (BS92). Mẫu là các bản ghi còn thiếu thông tin sự kiện. Thông tin về sự kiện còn thiếu này sẽ được bổ sung từ dữ liệu căn cứ vào những thông tin đã định nghĩa trên mẫu. Một cách thuần túy

---

<sup>1</sup>International Conference on **K**nowledge **D**iscovery and **D**ata Mining

<sup>2</sup>**S**pecial **I**nterest **G**roup on **I**nformation **R**etrieval

<sup>3</sup>**T**opic **D**etection and **T**racking

thì đây là bài toán tìm kiếm các từ được định nghĩa trước rồi lấy thông tin đi kèm với chúng để điền vào mẫu. Độ chính xác của phương pháp này ở mức trung bình nằm trong khoảng 50%–60% (MW11). Cách giải quyết bài toán hết sức đơn giản mà về sau, trong các chương trình nghiên cứu TDT hay ACE vẫn còn sử dụng nhưng với những định nghĩa mẫu tổng quát và trên nhiều miền lĩnh vực khác nhau. Hơn nữa, đây cũng là sự khởi đầu của các phương pháp đi theo hướng tiếp cận đầu tiên kể ở trên: sử dụng luật phân tích ngữ nghĩa.

Trong nghiên cứu của Nishihara và cộng sự, ba thông tin: địa điểm, đối tượng, hành vi của sự kiện được lấy ra từ trang cá nhân <sup>1</sup> (YKW09) sử dụng các *luật lexico-syntactic* <sup>2</sup> để tìm kiếm các câu chứa sự kiện trong từng bài viết <sup>3</sup>. Cùng với cách tiếp cận này, Aone.C và Ramos.M đã trích chọn các sự kiện về tài chính và chính trị. Hai tác giả tập trung đưa ra các luật biểu diễn quan hệ giữa sự kiện với các thông tin xung quanh nhằm mục đích khai thác tối đa thuộc tính của sự kiện, và giữa các sự kiện để lấy được tập các sự kiện liên quan tới nhau (CM00). Nghiên cứu của Xu và cộng sự cũng sử dụng các *luật lexico-syntactic* trên dữ liệu bản tin về sự kiện giải thưởng Nobel. Nhưng thay vì các luật được áp dụng ngay trên dữ liệu, một tập luật được tạo ra sau đó sử dụng học máy không giám sát để huấn luyện tập luật này trên tập các bản tin đã được gán nhãn. Sau đó mô hình học sẽ được áp dụng với các bản tin còn lại (FHH06).

Một điểm yếu của *luật lexico-syntactic* là không thể phủ hết được trạng thái quan hệ giữa các sự kiện, có nghĩa là không thể nhận biết hai sự kiện có trùng nhau hay không. Do đó, giám sát quá trình tiến triển của một sự kiện là tương đối khó khi sử dụng cách tiếp cận này. Để khắc phục điều này, *luật lexico-semantic* <sup>4</sup> được đề xuất.

### 1.5.2 Một số nghiên cứu liên quan ở trong nước

---

<sup>1</sup>blog

<sup>2</sup>*Luật lexico-syntactic* là sự kết hợp giữa biểu thức chính quy với từ vựng thuộc miền lĩnh vực và các quy tắc ngữ pháp của ngôn ngữ để sinh luật

<sup>3</sup>entry

<sup>4</sup>*luật lexico-semantic* là sự kết hợp giữa biểu thức chính quy, tập từ vựng thuộc miền lĩnh vực và vai trò ngữ nghĩa của từ vựng trong ngôn ngữ để sinh luật

## 2

# Phương pháp trích xuất sự kiện

2.1 Sự kiện là gì?

2.2 Các phương pháp hướng nội dung

2.3 Các phương pháp hướng dữ liệu



## 3

# Phương pháp giải quyết

### 3.1 Hệ thống theo dõi tin tức

## Phụ lục A

# Danh sách các trang tin tức điện tử

24h.com.vn - <http://www21.24h.com.vn> VietnamNet (2Sao) - <http://2sao.vietnamnet.vn>  
aFamily - <http://afamily.channelvn.net> Alobacsi.vn - <http://alobacsi.vn> Báo An Ninh Thủ Đô (ANTĐ) - <http://www.anninhthudo.vn> Báo An Ninh Thế Giới (ANTG) - <http://antg.cand.com.vn> Báo Công An Nhân Dân (ANTGCT) - <http://antg.cand.com.vn> Archi.vn (Archi) - <http://archi.vn> ATPVietnam - <http://atpvietnam.com> Autonet - <http://www.autonet.com.vn> AutoPro - <http://autopro.channelvn.net> Báo Biên phòng - <http://www.bienphong.com.vn> Báo Bóng Đá - <http://http://www.baobongda.com.vn> Báo Công Lý - <http://congly.com.vn> Báo Công Thương - <http://baocongthuong.com.vn> Báo Đất Việt - <http://www.baodatviet.vn> Báo Đất Việt (Báo Đất Việt - Khoa học) - <http://khoaoc.baodatviet.vn> Báo Đất Việt (Báo Đất Việt - Quốc phòng) - <http://quocphong.baodatviet.vn> Báo Giáo dục Việt Nam - <http://giaoduc.net.vn> Báo Giao Thông Vận Tải (Báo GTVT) - <http://giaothongvantai.com.vn> Báo Khoa học Phát triển - <http://khoaocphattrien.com.vn> Báo Người cao tuổi - <http://nguocaotui.org.vn> Báo Nông nghiệp VN - <http://nongnghiep.vn> Báo Phụ Nữ (Báo Phụ Nữ Online) - <http://www.phunuonline.com.vn> Báo Thế giới Việt nam - <http://www.tgvn.com.vn> Báo Tia sáng - <http://www.tiasang.com.vn> Thông Tấn Xã VN (Báo Tin tức) - <http://baotintuc.vn> Báo Thể thao Văn Hóa (Báo TTVH) - <http://thethaovanhoa.vn> Báo Thể Thao VN (Báo TTVN) - <http://www.thethaovietnam.com.vn> Báo Văn hóa - <http://www.baovanhoa.vn> Báo Khoa học Đời sống (Bee.net.vn) - <http://bee.net.vn/> Bóng Đá 24H - <http://www.bongda24h.vn> Bóng đá số - <http://www.bongdaso.com> YTT (Bongda.com.vn) - <http://www.bongda.com.vn> CafeF - <http://cafef.vn> Báo Công

---

An Nhân Dân (CAND Portal) - <http://www.cand.com.vn> Công An TP.HCM (CAT-PHCM) - <http://www.congan.com.vn> CTTĐT Chính phủ (Chinhphu.vn) - <http://baodientu.chinhphu.vn> Báo Công An Nhân Dân (CSTC) - <http://cstc.cand.com.vn> Báo Đại đoàn kết (Đại Đoàn Kết) - <http://baodaidoanket.net> Báo Dân Trí (Dân Trí) - <http://www.dantri.com.vn> Báo Nông thôn ngày nay (Dân Việt) - <http://www.danviet.vn> Báo Đầu Tư Chứng Khoán (Đầu tư CK) - <http://www.tinnhanhchungkhoan.vn> Báo Điện tử Đảng cộng sản VN (ĐCSVN) - <http://www.cpv.org.vn> Địa ốc Online - <http://www.diaonline.vn> Báo Diễn Đàn Doanh Nghiệp (Diễn đàn Doanh nghiệp) - <http://www.dddn.com.vn> Điện Tử Tiêu Dùng - <http://dientutieudung.vn> Doanh nhân 360 - <http://doanhnhan360.com> Doanh nhân Sài Gòn - <http://doanhnhansaigon.vn> Báo Đời sống Pháp luật (Đời sống Pháp luật) - <http://www.doisongphapluat.com.vn> Dothi.net - <http://dothi.net> Báo Doanh nhân Việt Nam toàn cầu (DVT.vn) - <http://dvt.vn> VnExpress (eBank) - <http://ebank.vnexpress.net> eFinance - <http://www.taichinhdientu.vn> 24H.COM.VN (Eva.vn) - <http://www.eva.vn> VnExpress (eVăn) - <http://evan.vnexpress.net> Gafin.vn - <http://gafin.vn> Game4V - <http://news.game4v.vn> GameK - <http://gamek.channelvn.net> Gamethu.net - <http://gamethu.net> Báo Gia đình Xã hội (Giadinh.net) - <http://giadinh.net.vn> Báo GDĐTĐ (Giáo dục Thời đại) - <http://giaoducthoidai.vn> Báo Hà Nội Mới (Hà Nội Mới) - <http://hanoimoi.com.vn> Báo Hoa Học Trò (HHT) - <http://www.hoahoctro.vn> Báo Bưu Điện (ICTNews) - <http://ictnews.vn> ICTPress - <http://ictpress.vn> Infonet - <http://infonet.vn> InfoTV - <http://infotv.vn> VnExpress (iOne.net) - <http://ione.net> Kênh 14 - <http://kenh14.channelvn.net> KhoaHoc.com.vn - <http://khoaoc.com.vn> Báo Kinh tế Đô Thị (KTĐT) - <http://www.ktdt.com.vn> KTNT - <http://kinhtenongthon.com.vn> LandToday - <http://landtoday.net/> Báo Lao Động (Lao Động) - <http://laodong.com.vn> Báo Mực Tím (Mực tím) - <http://muctim.com.vn> MUST.vn - <http://www.must.vn> NDHMoney.vn - <http://ndhmoney.vn> Ngoisao.net - <http://ngoisao.net> Báo Người Lao Động (Người Lao Động) - <http://nld.com.vn> Báo Đời sống Pháp luật (Nguoiduatin.vn) - <http://nguoiduatin.vn> Nhà báo Công luận - <http://congluan.vn> Báo Nhân Dân (Nhân dân) - <http://www.nhandan.com.vn> Nhịp Cầu Đầu Tư - <http://nhipcaudautu.vn> Tạp chí PCWorld VN (PCWorld VN) - <http://pcworld.com.vn> Petrotimes.vn (Petrotimes) - <http://www.petrotimes.vn> Báo Pháp luật Xã hội (Pháp luật Xã hội) - <http://www.phapluatxahoi.vn> Báo Pháp luật TPHCM (Pháp luật TPHCM) - <http://www.phapluattp.vn> Pháp luật VN - <http://www.phapluatvn.vn> Báo Đời sống Pháp luật (Phunutoday.vn) - <http://phunutoday.vn> Báo Quân Đội Nhân Dân (QĐND) - <http://www.qdnd.vn> Saga.vn - <http://www.saga.vn>

---

SaigonNews - <http://www.saigonnews.vn> SaigonTimes Online (SaigonTimes) - <http://www.thesaigontimes.com>  
Sàn OTC - <http://news.sanotc.com> Báo Sài Gòn Giải Phóng (SGGP) - <http://sggp.org.vn>  
Báo Sài Gòn Tiếp Thị (SGTT) - <http://www.sggt.com.vn> Sohoa.net - <http://sohoa.net>  
StockBiz - <http://stockbiz.vn> Truyền thông Tài chính StoxPlus (StoxPlus) - <http://stox.vn/>  
Báo Sức khỏe Đời sống (Sức Khỏe Đời Sống) - <http://suckhoedoisong.vn> Sức Sống Mới -  
<http://www.sucsongmoi.net> Sinh viên Việt Nam (SVVN) - <http://www.svv.vn> Tamn-  
hin.net - <http://www.tamnhin.net> Tạp chí ĐẸP Online (Tạp chí ĐẸP) - <http://www.dep.com.vn>  
Tạp chí Hoạt Động Khoa Học (Tạp chí HDKH) - <http://www.tchdkh.org.vn> Tạp chí Tài  
chính - <http://tapchitaichinh.vn/> Báo Thanh Niên (Thanh Niên) - <http://www.thanhnien.com.vn>  
Báo Thanh Niên (Thanh Niên - Thể thao) - <http://www.thanhnien.com.vn/thethao>  
Báo Thanh Niên (Thanh Niên - Tuần san) - <http://www.thanhnien.com.vn/tnotuan>  
Báo Thanh Niên (Thanh niên - WC 2010) - <http://www.thanhnien.com.vn/worldcup2010>  
Báo điện tử Thế Giới Điện Ảnh (Thế Giới Điện Ảnh) - <http://thegioidienanh.vn> Thi-  
ennhien.net - <http://www.thiennhien.net> Thongtinduan.vn - <http://thongtinduan.vn>  
Báo Tiền Phong (Tiền Phong) - <http://www.tienphongonline.com.vn> Tiin.vn - <http://tiin.vn>  
Tin tức Du lịch - <http://www.dulichvn.org.vn> Tin Tức Online - <http://tintuonline.com.vn>  
Tinhte.com - <http://tinhte.com> YTT (TinTheThao) - <http://www.tinthethao.com.vn>  
Báo Tổ Quốc (Tổ quốc) - <http://www.toquoc.gov.vn> Thông Tin Công Nghệ (TTCN) -  
<http://www.thongtincongngh.com> Tuần Vietnamnet (Tuần Việt Nam) - <http://tuanvietnam.net>  
Báo Tuổi Trẻ (Tuổi Trẻ) - <http://tuoitre.vn> Báo Tuổi Trẻ (Tuổi trẻ - Cuối tuần) -  
<http://tuoitre.vn/Tuoi-tre-cuoi-tuan/index.html> Báo Tuổi Trẻ (Tuổi Trẻ - Địa Ốc) -  
<http://diaoc.tuoitre.vn> Báo Tuổi Trẻ (Tuổi trẻ - Du lịch) - <http://dulich.tuoitre.com.vn>  
Báo Tuổi Trẻ (Tuổi trẻ - Nhịp sống số) - <http://nhipsongso.tuoitre.com.vn> Báo Tuổi  
Trẻ (Tuổi Trẻ - Thể Thao) - <http://thethao.tuoitre.vn> Báo Tuổi Trẻ (Tuổi trẻ - Tuyển  
sinh) - <http://chuyentrang.tuoitre.vn/tuyensinh/> Báo Tuổi Trẻ (Tuổi trẻ - Việc làm) -  
<http://vieclam.tuoitre.vn/> Báo Tuổi Trẻ (Tuổi trẻ - WC 2010) - [http://chuyentrang.tuoitre.vn/WorldC](http://chuyentrang.tuoitre.vn/WorldCup2010)  
VietnamNet (VEF) - <http://vef.vn> Vietnam Economic News Online (VEN) - <http://ven.org.vn>  
Thông Tấn Xã VN (Vietnam Plus) - <http://www.vietnamplus.vn> VietnamNet - <http://vietnamnet.vn>  
VietnamNet (VietnamNet - Thể Thao) - <http://thethao.vietnamnet.vn> Vietstock -  
<http://www.vietstock.com.vn> Báo Đầu Tư (VIR) - <http://www.baodautu.vn> Vitinfo  
- <http://vitinfo.com.vn> Báo Công An Nhân Dân (VNCA) - <http://vnca.cand.com.vn>  
Thời báo Kinh Tế (VnEconomy) - <http://vneconomy.vn> VnExpress - <http://vnexpress.net>  
VnMedia - <http://vnmedia.vn> VnRock - <http://vnrock.com/> Đài Tiếng Nói TP.HCM

---

(VOH) - <http://voh.com.vn> Đài Tiếng Nói VN (VOV Online) - <http://vov.vn> VTC News (VTC) - <http://vtc.vn> VTC News (VTC - Bạn đọc) - <http://vtc.vn/trangbandoc/> VTC News (VTC - Bảo vệ NTD) - <http://vtc.vn/bvntd/> VTC News (VTC - Công nghệ) - <http://vtc.vn/congnghe/> VTC News (VTC Games) - <http://vtc.vn/thegioigame> Đài TH VN (VTV) - <http://www.vtv.vn> Vzone - <http://vzone.vn> Tạp chí Xã Hội Thông Tin (XHTT) - <http://xahoithongtin.com.vn> Xinhxinh.com.vn - <http://xinhxinh.com.vn> XZone - <http://xzone.vn> Zing - <http://www.zing.vn> Zing (Zing - Thể thao) - <http://thethao.zing.vn/news>

**Phụ lục B**

**XYZ**

# Tài liệu tham khảo

*Tiếng Việt*

*Tiếng Anh*

- [FFU11] Frederik Hogenboom, Flavius Frasincar, Uzay Kaymak, Franciska De Jong. *An Overview of Event Extraction from Text*. Workshop on Detection, Representation, and Exploitation of Events in the Semantic Web, DeRiVE, 2011. 7
- [MW11] Martin Wunderwald. *NewsX–Event Extraction from News Articles*. Master Thesis. Dresden University of Technology, Germany, 2011. 8
- [HLH10] Sheng-Hao Hung, Chia-Hung Lin Jen-Shin Hong. *Web mining for event-based commonsense knowledge using lexico-syntactic pattern matching and semantic role labeling*. Journal Expert Systems with Applications, vol. 37, 2010.
- [YKW09] Yoko Nishihara, Keita Sato, Wataru Sunayama. *Event Extraction and Visualization for Obtaining Personal Experiences from Blogs*. Human Interface and the Management of Information. Information and Interaction. LNCS, vol. 5839, Springer–Verlag, 2009. 8
- [FHH06] Feiyu Xu, Hans Uszkoreit, Hong Li. *Automatic Event and Relation Detection with Seeds of Varying Complexity*. AAAI Workshop on Event Extraction and Synthesis, 2006. 8
- [CM00] Chinatsu Aone, Mila Ramos-Santacruz. *REES: a large-scale relation and event extraction system*. Applied Natural Language Processing Conference, 6<sup>th</sup>, ANLP00, 2000. 8

- [RB96] Ralph Grishman, Beth Sundheim. *Message Understanding Conference - 6: A Brief History*. MUC-6, 1996. 1
- [BS92] Beth Sundheim. *Overview of the fourth message understanding evaluation and conference*. MUC-4, 1992. 7