

Rút trích thông tin trong văn bản (Information Extraction)

Nguyễn Trường Sơn
ntson@fit.hcmus.edu.vn

fit cdio™ KHOA CÔNG NGHỆ THÔNG TIN
TRƯỜNG ĐẠI HỌC KHOA HỌC TỰ NHIÊN



Nội dung

- Giới thiệu rút trích thông tin và ứng dụng
- Các cách tiếp cận
- Demo:
 - Nhận diện thực thể có tên (NER)



Rút trích thông tin

- **Information extraction (IE):**
 - Tự động rút trích các thông tin có cấu trúc từ những tài liệu không có cấu trúc / bán cấu trúc
 - Rút trích thông tin = rút trích từ các tài liệu đa phương tiện (multimedia): **văn bản**, hình ảnh, âm thanh, video

Yesterday, New York based Foo Inc. announced their acquisition of Bar Corp.



MergeBetween(Foo Inc, Bar Corp, Yesterday)

3



Các dạng bài toán đặc trưng

- Nhận diện thực thể đã định danh:
 - Named Entity Recognition (NER)
- Nhận diện quan hệ giữa các thực thể
 - Relation Extraction
 - Event Extraction
- Nhận diện đồng tham chiếu (Co-reference)
- ...





Named Entity Recognition and Coreference Resolution

Named Entity Recognition (NER):

Example:

The shiny red rocket was fired on Tuesday. It is the brainchild of Dr. Big Head.
Dr. Head is a staff scientist at We Build Rockets Inc.

→ <person> Dr. Big Head </person>

Output

<person> Dr. Head </person>
<organization> We Build Rockets Inc </organization>
<time> Tuesday </time>

Coreference resolution (anaphor resolution):

- Connect pronous etc. to subject/object of previous sentence

Examples:

- The shiny red rocket was fired on Tuesday. It is the brainchild of Dr. Big Head.
- ... or Tuesday. It <reference> The shiny red rocket </reference> is the ...
- Alas, poor Yorick, I knew him Horatio.



Tạo “ngữ nghĩa(Semantic)” cho dữ liệu

Hầu hết dữ liệu ở dạng HTML (hoặc PDF hoặc RSS hoặc ..
Hoặc từ các nguồn dữ liệu không có cấu trúc ..

Property Features

- Single Family
- Style: Ranch
- Interior features: Carpet, Clothes dryer, Clothes washer, Eat-in kitchen, Finished basement, Fireplace(s), Range and oven, Refrigerator, Utility Room, Wood fire
- Area: WAWARSING
- Year Built: 1965
- 3 total bedroom
- 2 car garage
- Heating features: Electric
- Exterior features: Sloped lot, Water supply from well(s), Wooded area(s)
- Lot size is between 5 and 10 acres
- School District: TRI ALLEY CENTRAL ELEMENTARY SCHOOL GRAHAMSVILLE
- Approximately 1640 sq ft.

Description

Approx. 235 Acres - WOW! Area: OutSide Area.
Community Name: Escalante.

Features: Lot Size: 235 Acre

Additional Information: Also features: * Single Family Property, * Area: OutSide Area, * Community Name: Escalante, * Year Built: 1973, * 6 total bedroom(s), * 4 total bath(s), * 3 total full bath(s), * 1 total half bath(s), *

Có thể truy xuất thông qua
Wrappers hoặc Web Service)

→ Luật, FSAs (biểu thức chính qui), ... → HMMs, MRFs, ...



Mục tiêu và tổng quan

Mục tiêu:

- gán nhãn(annotation) tài liệu văn bản hoặc trang Web (named entity recognition, html2xml, etc.)
- trích sự kiện (extract facts) từ tài liệu văn bản hay trang Web (relation learning)
- tìm sự kiện (find facts) trên Web (hoặc trên Wikipedia) để xây dựng các kho dữ liệu về quan hệ (thesaurus/ontology relations)
- làm giàu thông tin (information enrichment)
(vd. Cho phân tích kinh doanh)

Công nghệ:

- NLP (PoS tagging, chunk parsing, etc.)
- Lexicon lookups (name dictionaries, geo gazetteers, etc.)
- Pattern matching & rule learning (regular expressions, FSAs)
- Statistical learning (HMMs, CRFs, etc.)
- Text mining in general, deep learning *



Tạo “ngữ nghĩa(Semantic)” cho dữ liệu

Hầu hết dữ liệu ở dạng HTML (hoặc PDF hoặc RSS hoặc ..
Hoặc từ các nguồn dữ liệu không có cấu trúc ..

Country

The state borders France in the south and west, Luxembourg in the west and Rhineland-Palatinate in the north and the east.

State

It is named after the Saar River, which is an affluent of the Moselle River and runs through the state from the south to the northwest. Most inhabitants live in a city agglomeration on the French border, surrounding the capital of Saarbrücken.

River

The altitude above sea level of the city's area is between 100 m (on the westerly edge, toward the Rhine river) and 277.5 m (Turmburg in the east). Its geographical coordinates are: 49° 00' North 008° 04' East, which means that the 49th parallel (meridian) runs through the city center, its course being marked by a line of flag-stones in the Stadtgarten (city park).

City

Transport

Elevation

GeoCoord

Globe

FIDEL

Coat of Arms of Karlsruhe



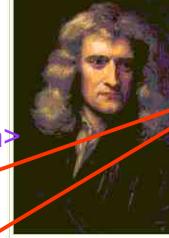
Tạo “ngữ nghĩa(Semantic)” cho dữ liệu

 Hầu hết dữ liệu ở dạng HTML (hoặc PDF hoặc RSS hoặc ..)
Hoặc từ các nguồn dữ liệu không có cấu trúc

Isaac Newton

From Wikipedia, the free encyclopedia.

Sir Isaac Newton (25 December 1642 – 20 March 1727 by the Julian calendar in use in England at the time; or 4 January 1643 – 31 March 1727 by the Gregorian calendar) was an English physicist, mathematician, astronomer, philosopher, and alchemist; who wrote the *Philosophiae Naturalis Principia Mathematica* (published 5 July 1687)¹, where he described universal gravitation and, via his laws of motion, laid the groundwork for classical mechanics. Newton also shares credit with Gottfried Wilhelm Leibniz for the development of differential calculus. However, their work was not a collaboration; they both discovered calculus separately but nearly contemporaneously.



Sir Isaac Newton in Kneller's portrait of 1689.

<TimePeriod>

<Scientist>

<Scientist>

<Publication>

Sir Isaac Newton in Kneller's portrait of 1689.

IRDW WS 2005

<Person>

<Person>

NLP-based IE từ Web



ANNIE Output for http://en.wikipedia.org/wiki/Che_Guevara

Annotation Key:
Person **Location** **Organization** **Date** **Address** **Money** **Percent**

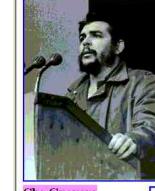
>> /**/ > /**/

Che Guevara

From Wikipedia, the free encyclopedia.

(Redirected from [Che Guevara](#))

Jump to: [navigation](#), [search](#)



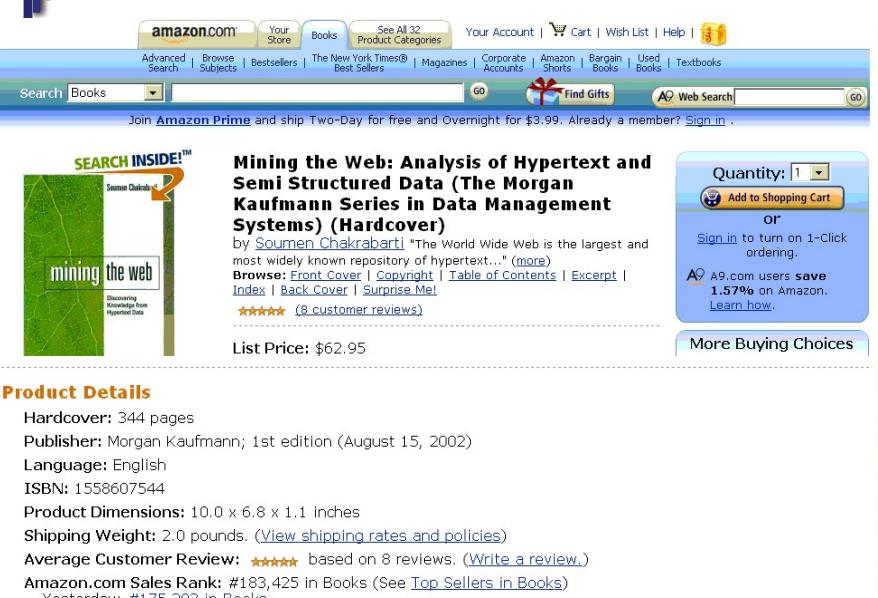
Che Guevara

Ernesto Rafael Guevara de la Serna (June 14, 1928¹¹ – October 9, 1967), commonly known as Che Guevara or el Che, was an Argentine-born Marxist revolutionary and Cuban guerrilla leader. Guevara was a member of Fidel Castro's "26th of July Movement" that seized power in Cuba in 1959. After serving in various important posts in the new government, Guevara left Cuba in 1965 with the hope of fomenting revolutions in other countries, first in the Congo-Kinshasa (currently the Democratic Republic of the Congo) and later in Bolivia, where he was captured in a CIA-organized military operation. It is believed by some that the CIA wished to keep Guevara alive for interrogation but, after his capture in the Yuro ravine, he died at the hands of the Bolivian Army in La Higuera near Vallegrande on October 9, 1967. Testimony by various individuals who were participants in, or

Leading open-source tool: GATE/ANNIE

<http://www.gate.ac.uk/annie/>

Rút trích các mẫu tin có cấu trúc từ Web (1)



The screenshot shows the product page for 'Mining the Web: Analysis of Hypertext and Semi Structured Data (The Morgan Kaufmann Series in Data Management Systems) (Hardcover)' by Soumen Chakrabarti. The page includes the book cover, a brief description, customer reviews, and purchase options like 'Add to Shopping Cart' and 'Sign in to turn on 1-Click ordering'. The price is listed as \$62.95.

Product Details

- Hardcover: 344 pages
- Publisher: Morgan Kaufmann; 1st edition (August 15, 2002)
- Language: English
- ISBN: 1558607544
- Product Dimensions: 10.0 x 6.8 x 1.1 inches
- Shipping Weight: 2.0 pounds. ([View shipping rates and policies](#))
- Average Customer Review: ★★★★ based on 8 reviews. ([Write a review](#))
- Amazon.com Sales Rank: #183,425 in Books (See [Top Sellers in Books](#))
- Yesterdays: #175,203 in Books

Rút trích các mẫu tin có cấu trúc từ Web (1I)

```

<div class="buying"><b class="sans">Mining the Web: Analysis of Hypertext and  
Semi Structured Data (The Morgan  
Kaufmann Series in Data Management  
Systems) (Hardcover)</b><br />
<a href="/exec/obidos/search-handle-url/index=books&field-author-exact=Soumen%20Ch  
ake">Mining the Web: Analysis of Hypertext and  
Semi Structured Data (The Morgan  
Kaufmann Series in Data Management  
Systems) (Hardcover)</a>
<div class="buying" id="priceBlock">
<style type="text/css">
td.productLabel { font-weight: bold; text-align: right; white-space: nowrap; vertical-align: top; }
table.product { border: 0px; padding: 0px; border-collapse: collapse; }
</style>
<table class="product">
<tr>
<td class="productLabel">List Price:</td>
<td>$62.95</td>
</tr>
<tr>
<td class="productLabel">Price:</td>
<td><b class="price">$62.95</b></td>
</tr>
<tr>
<td colspan="2" style="text-align: center;">& this item ships for <b>FREE</b> with Super Saver Shipping</td>
</tr>

```

IRDW WS 2005

Rút trích các mẫu tin có cấu trúc từ Web (1II)

```
<a name="productDetails" id="productDetails"></a>
<hr noshade="noshade" size="1" class="bucketDivide
<table cellpadding="0" cellspacing="0" border="0">
<tr>
  <td class="bucket">
    <b class="h1">Product Details</b><br />
    <div class="content">
      <ul>
        <li><b>Hardcover:</b> 344 pages</li>
        <li><b>Publisher:</b> Morgan Kaufmann; 1st edition (2003)</li>
        <li><b>Language:</b> English</li>
        <li><b>ISBN:</b> 1558607544</li>
        <li><b>Product Dimensions:</b> 10.0 x 6.8 x 1.1 inches</li>
        <li><b>Shipping Weight:</b> 2.0 pounds. (<a href="http://www.amazon.com/gp/product/0131459054?%5Fencoding=UTF8&seller=ATVPDKIKX0D1Z">View on Amazon.com</a>)</li>
        <li><b>Average Customer Review:</b>  based on 8 reviews.
          (<a href="http://www.amazon.com/gp/customer-reviews/write-a-review.html/102-8395894-112-1000000-0000000?ie=UTF8&asin=B000000000&refRID=ATVPDKIKX0D1Z">Write a review</a>)
        <li><b>Amazon.com Sales Rank:</b> #183,425 in Books (See <a href="/exec/obidos/tg/new-books/b000000000">Books</a>)
      </ul>
    </div>
  </td>
</tr>
</table>
```

Sample dataset: VLSP NER

 NER for COVID-19

NER for COVID-19

Label	Definition
PATIENT_ID	Unique identifier of a COVID-19 patient in Vietnam. An PATIENT_ID annotation over "X" refers to as the X th patient having COVID-19 in Vietnam.
PERSON_NAME	Name of a patient or person who comes into contact with a patient.
AGE	Age of a patient or person who comes into contact with a patient.
GENDER	Gender of a patient or person who comes into contact with a patient.
OCCUPATION	Job of a patient or person who comes into contact with a patient.
LOCATION	Locations/places that a patient was presented at.
ORGANIZATION	Organizations related to a patient, e.g. company, government organization, and the like, with structures and their own functions.
SYMPTOM&DISEASE	Symptoms that a patient experiences, and diseases that a patient had prior to COVID-19 or complications that usually appear in death reports.
TRANSPORTATION	Means of transportation that a patient used. Here, we only tag the specific identifier of vehicles, e.g. flight numbers and bus/car plates.
DATE	Any date that appears in the sentence.

Table 1: Definitions of entity types in our annotation guidelines. We do not annotate nested entities.



Các ứng dụng IE

- So sánh giá và tư vấn mua sắm portals
ví dụ: consumer electronics, used cars, real estate, pharmacy, etc.
- Phân tích kinh doanh dựa trên hồ sơ khách hàng, báo cáo tài chính
e.g.: How was company X (the market Y) performing in the last 5 years?
- Thị trường/khách hàng, ứng dụng PR, phân tích truyền thông
e.g.: How are our products perceived by teenagers (girls)?
How good (and positive?) is the press coverage of X vs. Y?
Who are the stakeholders in a public dispute on a planned airport?
- Giới thiệu việc làm (applications/resumes, job offers)
e.g.: Who well does the candidate match the desired profile?
- Quản trị tri thức trong các công ty tư vấn
e.g.: Do we have experience and competence on X, Y, and Z in Brazil?
- Khai thác E-mail đã được lưu trữ
e.g.: Who knew about the scandal on X before it became public?
- Trích tri thức từ văn bản khoa học
e.g.: Which anti-HIV drugs have been found ineffective in recent papers?
- Phân tích tri thức tổng quan
Can we learn encyclopedic knowledge from text & Web corpora?



Các quan điểm và cách tiếp cận IE

IE là học các biểu thức chính qui (**regular expressions**)
(wrapping pages with common structure from Deep-Web source)

IE là học các quan hệ (**relations**)

(rules for identifying instances of n-ary relations)

IE là học các sự kiện (**fact boundaries**)

IE là học phân đoạn văm bản (**segmentation**) (HMMs etc.)

IE là học các mẫu ngữ cảnh (**contextual patterns**) (graph models etc.)

IE là phân tích ngôn ngữ tự nhiên (**natural-language analysis**) (NLP)

IE là khai phá kho văn bản lớn để khám phá tri thức
(combination of tools incl. Web queries)



Đánh giá chất lượng của IE

Sử dụng các độ đo chuẩn của IR:

- precision
- Recall
- F1 measure

Các bộ dữ liệu chuẩn từ các cuộc thi, hội nghị:

Tiếng Anh và các ngôn ngữ khác:

- MUC (Message Understanding Conference)
- ACE (Automatic Content Extraction), <http://www.nist.gov/speech/tests/ace/>
- TREC Enterprise Track, <http://trec.nist.gov/tracks.html>
- Enron e-mail mining, <http://www.cs.cmu.edu/~enron>

Tiếng Việt:

- VLSP Shared Tasks: 2018, 2019, 2020



MUC Information Extraction Example

Bridgestone Sports Co. said Friday it had set up a joint venture in Taiwan with a local concern and a Japanese trading house to produce golf clubs to be supplied to Japan. The joint venture, Bridgestone Sports Taiwan Co., capitalized at 20 million new Taiwan dollars, will start production in January 1990 with production of 20,000 iron and "metal wood" clubs a month.

JOINT-VENTURE-1

- Relationship: TIE-UP
- Entities: "Bridgestone Sport Co.", "a local concern", "a Japanese trading house"
- Joint Ent: "Bridgestone Sports Taiwan Co."
- Activity: ACTIVITY-1
- Amount: NT\$20 000 000

ACTIVITY-1

- Activity: PRODUCTION
- Company: "Bridgestone Sports Taiwan Co."
- Product: "iron and 'metal wood' clubs"
- Start date: DURING: January 1990



Các thách thức trong IE

- Phụ thuộc lĩnh vực
 - Landscape of IE Tasks**
 - Nhiều cấp độ định dạng khác nhau
 - Sự đa dạng
 - Sự phức tạp
 - Cấp độ rút trích khác nhau



Landscape of IE Tasks (1/4): Degree of Formatting

Text paragraphs
without formatting

Astro Teller is the CEO and co-founder of BodyMedia. Astro holds a Ph.D. in Artificial Intelligence from Carnegie Mellon University, where he was inducted as a national Hertz fellow. His M.S. in symbolic and heuristic computation and B.S. in computer science are from Stanford University. His work in science, literature and business has appeared in international media from the New York Times to CNN to NPR.

Non-grammatical snippets, rich formatting & links

Barto, Andrew G.	(413) 545-2109	barto@cs.umass.edu	CS276
Professor.	Computational neuroscience, reinforcement learning, adaptive motor control, artificial neural networks, adaptive learning control, motor development.	 	
Berger, Emery D.	(413) 577-4211	emery@cs.umass.edu	CS344
Assistant Professor.		 	
Brock, Oliver	(413) 577-0334	oliver@cs.umass.edu	CS246
Assistant Professor.		 	
Clarke, Lori A.	(413) 545-1328	clarke@cs.umass.edu	CS304
Professor.	Software verification, testing, and analysis; software architecture and design.	 	
Cohen, Michael R.	(413) 545-3638	cohen@cs.umass.edu	CS278
Professor.	Planning, simulation, natural language, agent-based systems, intelligent decision analysis, intelligent user interfaces.	 	

Grammatical sentences

Dr. Steven Minton - Founder/CTO
Dr. Minton is a fellow of the American Association of Artificial Intelligence and was the founder of the Journal of Artificial Intelligence Research. Prior to founding Fetch, Minton was a faculty member at USC and a project leader at USC's Information Sciences Institute. A graduate of Yale University and Carnegie Mellon University, Minton has been a Principal Investigator at NASA Ames and taught at Stanford, UC Berkeley and USC.

- Press
 - Contact**
 - General information
 - Directions maps

Tables						
Invited Talk: Plausibility Measures: A General Approach for Representing Uncertainty Joseph Y. Halpern, Cornell University						
Coffee Break						
Technical Paper Sessions:						
Cognitive Robotics	Logic Programming	Natural Language Generation	Complexity Analysis	Neural Networks	Games	
739: A Logical Approach of Causal Temporal Maps <i>Emre and Benjamin Kupers</i>	116: A-System: Problem Solving through Abduction <i>Rong Jin and Alexander G. Hauptmann</i>	758: Title Generation for Machine-Translated Documents <i>Denecker, Antonis Kazas, and Koenraad Van Nuffelen</i>	417: Let's go Complexity of Nested Circumscriptive Abnormality Theories <i>Thomas Colai, Thomas Eiter, and Georg Gottlob</i>	179: Knowledge Extraction and Summarization from Local Function Descriptions <i>Kenneth McCrary, Stefano Ferri, and John MacIntyre</i>	71: Iterative Widening Prisoner's Dilemma <i>Yannick Viossat</i>	
S40:	Online-Execution of Programs <i>Hilke Grossenbacher</i>	131: A Comparative Study of Programmatic Languages <i>Yannick Viossat</i>	246: Dealing with Dependencies in Planning and Scheduling Programs with Constraints <i>Yannick Viossat</i>	470: A Perspective on Violation-Guided Compilation <i>Yannick Viossat</i>	258: Temporal Difference Learning Applied to a Constrained Planning Problem <i>Yannick Viossat</i>	



 IE is different in different domains!

Example: on web there is less grammar, but more formatting & linking

Newswire

Apple to Open Its First Retail Store in New York City

MACWORLD EXPO, NEW YORK—July 17, 2002—Apple's first retail store in New York City will open in Manhattan's SoHo district on Thursday, July 18 at 8:00 a.m. EDT. The SoHo store will be Apple's largest retail store to date and is a stunning example of Apple's commitment to offering customers the world's best computer shopping experience.

"Fourteen months after opening our first retail store, our 31 stores are attracting over 100,000 visitors each week," said Steve Jobs, Apple's CEO. "We hope our SoHo store will surprise and delight both Mac and PC users who want to see everything the Mac can do to enhance their digital lifestyles."

Landscape of IE Tasks (2/4): Intended Breadth of Coverage



Web site specific

Formatting

The screenshot shows the Amazon.com homepage. At the top, there's a search bar with the placeholder "Search Amazon.com Books & DVDs". Below it, a green banner features the text "WELCOME YOUR FRIENDS" and "SEARCH BROWSE SUBJECTS". A large "amazon.com." logo is centered above a navigation bar with categories: HOME & KITCHEN, BOOKS, ELECTRONICS, DVD & BLU-RAY, and CORPORATE ACCOUNT. A "VIEW CART" button is also present. The main content area has a large "amazon.com." logo. Below it, a promotional banner for "Machine Learning" by Michael L. Jordan offers a \$5 off discount. To the right, another "VIEW CART" button is shown. The page also features a "NEW Super Saver Shipping FREE" offer. On the left, there's a sidebar for "Machine Learning" with a "Look Inside" button and a "Buy Now Pay Later" button. The main content area contains text about learning graphical models and a "Look Inside" button for "Graphical Models". A "Used & New" section shows a used copy for \$20.00. The "Edition" section indicates it's a Paperback. A "See more product details" link is also provided.

Genre specific

Layout Resumes

<p>Jason D. M. Rennie</p> <p>Massachusetts Institute of Technology MIT AI Lab NE43-102 200 Technology Sq. Cambridge, MA 02139</p> <p>Research Interests My main interests lie in the automated analysis of data for the purposes of classification and the acquisition of knowledge. I have both interests in applying such methods to real-world problems.</p>	<p>jrennie@mit.edu http://www.ai.mit.edu/people/jrennie/ (617) 253-5319</p>
<p>L. Douglas Baker</p> <p>Address mobile phone: 703-272-1002 Office Address Stanford Research Institute Computer Mathematics Laboratory 5000 Forbes Avenue Pittsburgh, PA 15213 (412) 485-4016 Home Page http://www.homesite.com/~lbaker</p>	<p>Education Carnegie-Mellon University Ph.D. Computer Science, in progress University of California at Berkeley Technical University of Berlin Exchange Fellow, 1990-1993 University of Michigan M.S.E. Computer Science and Engineering, 1984 B.S.E. Computer Engineering, Summa Cum Laude, 1982</p>
<p>Research Experience</p> <p>Carnegie-Mellon University</p>	<p>Pittsburgh, PA Berlin, Germany</p>
<p>I am currently pursuing my dissertation research: a hierarchical system for automatically detecting objects in 3D. This work is being done as part of the Fraunhofer Institute for Computer Graphics and Vision's project on 3D object recognition.</p>	<p>1994</p>

Wide, non-specific

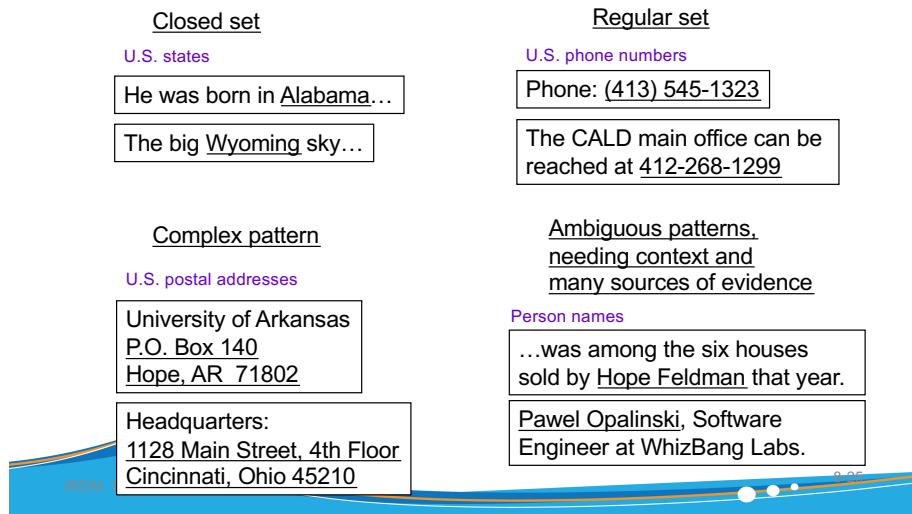
Language University Names

8:30 - 9:30 AM	Invited Talk: Plausibility Measures: A General Approach Joseph Y. Halpern, Cornell University
9:30 - 10:00 AM	Coffee Break
10:00 - 11:30 AM	Technical Paper Sessions:
	Cognitive Science
	Neural Language Generation
	Complexity Analysis
	No W
11:30 - 12:30 PM	Programmatic Generations
	Problem Generation
	Text Generation
	Nats:
	Ex
12:30 - 1:30 PM	Account of Causality and Topological Problem Solving
	Semantics Translated
	Documents
	Abstraction
	Rong Jin and Alexander G.
	Emilio Remedios and Benjamin Davies
	Hausmann
	Circumscription
	Fu
	Albinorwsky
	Theories
	M
	Antonis Kakas, and Bart Van
	Marco Cadoli, and Peter M.
	W
	Decker
	Steinbock
	W
	• Press
	Contact
	General information
	Directions maps



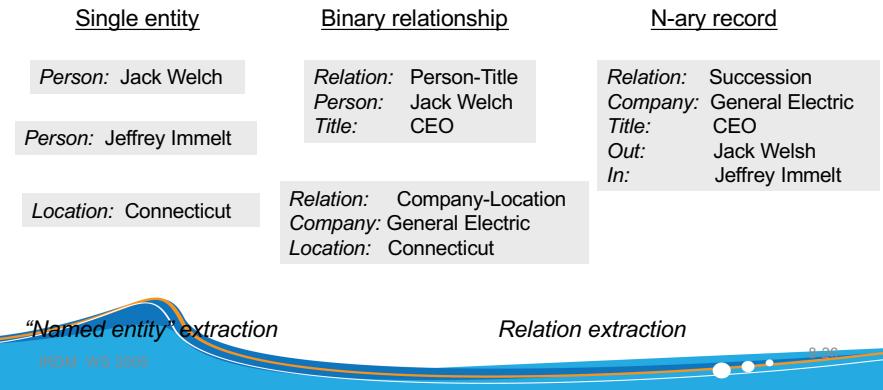
Landscape of IE Tasks (3/4): Complexity

E.g. word patterns:

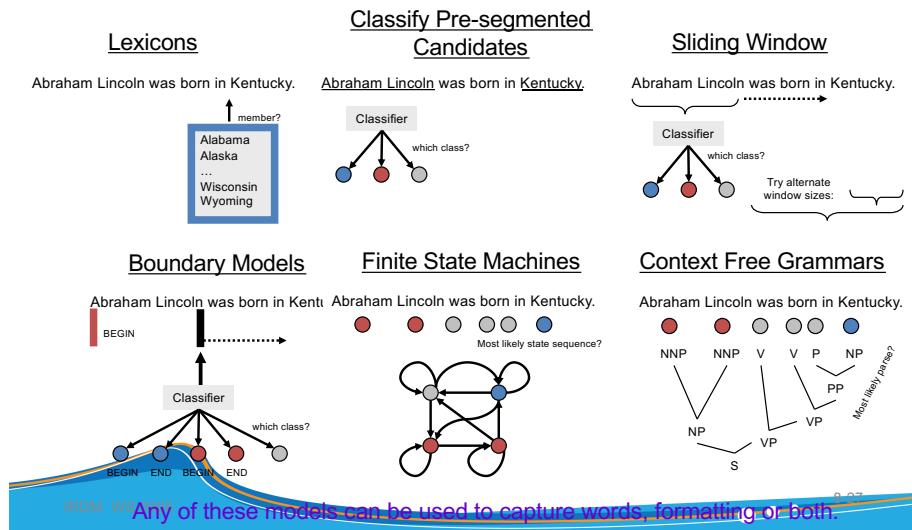


Landscape of IE Tasks (4/4): Single Field vs. Record

Jack Welch will retire as CEO of General Electric tomorrow. The top role at the Connecticut company will be filled by Jeffrey Immelt.



Các mô hình sử dụng trong IE



Học các biểu thức chính qui

Input: Các mẫu biểu thức chính qui được gán nhãn thủ công

learn: (restricted) biểu thức chính qui hoặc một mạng trạng thái hữu hạn dùng để đọc một câu và xuất ra các mẫu mong muốn

Ví dụ:

This apartment has 3 bedrooms.
 The monthly rent is \$ 995.

This appartment has 3 bedrooms.
 The monthly rent is \$ 995.

The number of bedrooms is 2.
 The rent is \$ 675 per month.

Learned pattern: * Digit „
“ * „\$“ Number *

Input sentence: There are 2 bedrooms.
 The price is \$ 500 for one month.

Output tokens: Bedrooms: 1, Price: 500

but: grammar inference for full-fledged regular languages is hard → focus on restricted fragments of the class of regular languages

implemented in WHISK (Soderland 1999) and a few other systems

IRDM WS 2005

8.20

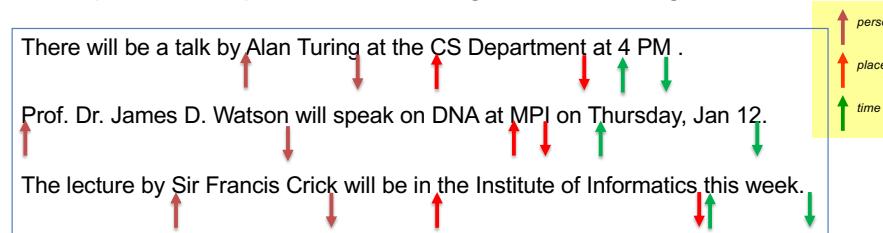
IE như bài toán phân loại các giới hạn



Ý tưởng chính:

Bộ phân loại (ví dụ: SVMs) nhận dạng ra token bắt đầu và token kết thúc của một đối tượng được xét đến

Tổ hợp nhiều bộ phân loại để nâng cao hiệu năng



Bộ phân loại kiểm tra mỗi token (với nhãn từ loại (PoS), các token xung quanh ... như các đặc trưng cho việc phân thành hai lớp: begin-fact, end-fact)



IE như là bài toán gán nhãn từng token trong chuỗi – Sequence Labeling

Ý tưởng chính:

Bộ phân loại xác định mỗi token thuộc về một **loại đặc biệt** thể hiện cho loại đối tượng (dựa vào các đặc trưng của token đang xét)

There will be a talk by Alan Turing at the CS Department at 4 PM .
O O O O O B-PER I-PER O O B-PLACE I-PLACE 4 PM O B-TIME I-TIME O

Prof. Dr. James D. Watson will speak on DNA at MPI on Thursday, Jan 12.

The lecture by Sir Francis Crick will be in the Institute of Informatics this week.



IOB Notation / CONLL format

ORG MISC
EU rejects German call to boycott British lamb . MISC



PER
Peter Blackburn

LOC
BRUSSELS 1996-08-22

IOB annotation: một số biến thể: IOBE, IOBES, IOE

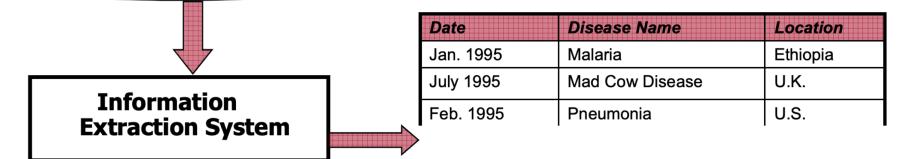
	token	postag	label	
0	EU	NNP	B-ORG	ORGANIZATION:EU
1	rejects	VBZ	O	
2	German	JJ	B-MISC	MISC:German
3	call	NN	O	
4	to	TO	O	
5	boycott	VB	O	
6	British	JJ	B-MISC	MISC: British
7	lamb	NN	O	
8	.	.	O	
9				
10	Peter	NNP	B-PER	PERSON: Peter Blackburn
11	Blackburn	NNP	I-PER	
12				
13	BRUSSELS	NNP	B-LOC	LOCATION: BRUSSELS
14	1996-08-22	CD	O	
15				



Relation Extraction: a task in IE

Tìm mối quan hệ giữa các thực thể:

May 19 1995, Atlanta -- The Centers for Disease Control and Prevention, which is in the front line of the world's response to the deadly Ebola epidemic in Zaire , is finding itself hard pressed to cope with the crisis...



Relation Extraction: Method

- Manually engineered rules:
 - Rules defined over words/ Rules defined over words/entities:
 - “<company> located in located in <location>”
 - Rules defined over parsed text:
 - ((Obj) (Verb located) (*) (Subj))
- Machine Learning-based
 - Binary Relation Association as Binary Classification



Nhận diện thực thể có tên bằng CRF

Huấn luyện



Sử dụng và đánh giá



CRFs: J. Lafferty, A. McCallum, and F. Pereira. [Conditional random fields: Probabilistic models for segmenting and labeling sequence data](#). In Proc. of ICML, pp.282-289, 2001

Các cách tiếp cận

- Dựa trên luật (Rule-based)
 - Các mô hình cho sequence labeling sử dụng học máy:
 - HMM
 - CRF *
 - RNN / LSTM
 - LSTM-CRF *
 - BERT *
 - Other deep learning models
 - Kết hợp các trên tri thức ngôn ngữ học:
 - Đặc trưng từ loại (POS tagger), ngữ nghĩa (tùy theo bài toán)
-

Nhận diện thực thể có tên bằng CRF

- Các tập dữ liệu:
 - CONLL 2002 (Spain)
 - CONLL 2003 (English)
 - VLSP 2018 (Vietnamese)
 - Covid NER
 - https://github.com/VinAIResearch/PhoNER_COVID19
 - NER and slot filling:
 - <https://github.com/VinAIResearch/JointIDSF>
 - Nhãn ORG, PER, LOC, MISC, (DATE, NUM), ...
-



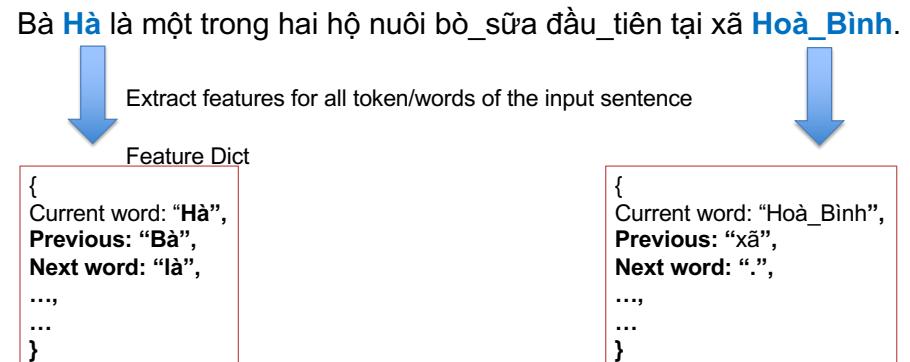


Nhận diện thực thể có tên bằng CRF

- Các công cụ thư viện cần thiết:
 - CRF: CRF++ (C++),
 - Công cụ rút trích đặc trưng:
 - nltk : Tiếng Anh và các ngôn ngữ khác
 - VnCoreNLP: Tiếng Việt
- Tip : Nên biểu diễn dữ liệu đầu vào dưới dạng CONLL format



Các đặc trưng thường dùng cho NER



Các mô hình NER khác

- Sử dụng các mô hình học sâu:
 - Cho kết quả tốt vượt trội
- Một số mô hình tiêu biểu:
 - RNN / LSTM
 - LSTM-CRF *
 - BiLSTM-CRF *
 - BERT *
- Google: Tài liệu Hướng dẫn khá đầy đủ
 - https://github.com/guillaumegenthial/tf_ner
 - <https://github.com/google-research/bert>

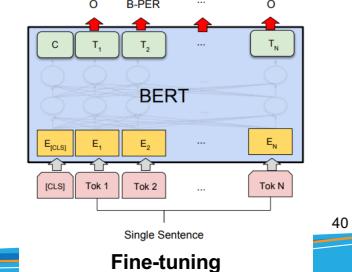
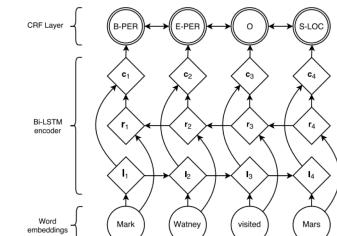
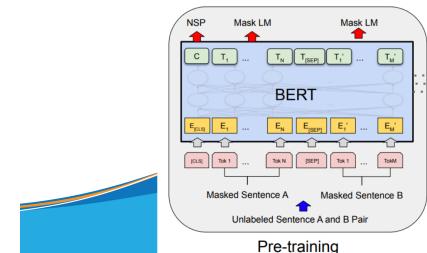


Các mô hình NER khác

Neural Architectures for Named Entity Recognition
 Guillaume Lample^{*} Miguel Ballesteros^{**}
 Sandeep Subramanian^{*} Kazuya Kawakami^{*} Chris Dyer^{*}
^{*}Carnegie Mellon University ^{**}NLP Group, Pompeu Fabra University
 {glample, sandeeps, kkawakan, cdyer}@cs.cmu.edu, miguel.ballesteros@upf.edu

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova
 Google AI Language
 {jacobjdevlin, mingweichang, kentonl, kristout}@google.com



Tham khảo

- https://web.stanford.edu/class/cs124/lec/Information_Extraction_and_Named_Entity_Recognition.pdf
- <https://www.slideshare.net/SHUBH177/text-categorization-72584299>
- <https://www.slideshare.net/kanimozhiu/text-datamining-txtcat>
- https://scikit-learn.org/stable/tutorial/text_analytics/working_with_text_data.html

- <https://www.coursera.org/learn/text-mining-analytics/lecture/QmxDT/2-1-description-of-possible-project-ideas>
- <https://www.coursera.org/learn/text-mining-analytics/home/week/3>

