### Chapter 15: Information Extraction and Knowledge Harvesting

The Semantic Web is not a separate Web but an extension of the current one, in which information is given well-defined meaning.

-- Sir Tim Berners-Lee

-- Albert Einstein

The only source of knowledge is experience.

To attain knowledge, add things everyday. To attain wisdom, remove things every day

Information is not knowledge. Knowledge is not wisdom. Wisdom is not truth. Truth is not beauty. Beauty is not love. Love is not music. Music is the best.

**IRDM WS 2015** 

-- Frank Zappa









-- Lao Tse

### Outline

- 15.1 Motivation and Overview
- 15.2 Information Extraction Methods
- 15.3 Knowledge Harvesting at Large Scale



# **8.1 Motivation and Overview**

### What?

- extract entities and attributes from (Deep) Web sites
- mark-up entities and attributes in text & Web pages
- harvest relational facts from the Web to populate knowledge base Overall: lift Web and text to level of "crisp" structured data

### <u>Why?</u>

- compare values (e.g. prices) across sites
- extract essential info fields (e.g. job skills & experience from CV)
- more precise queries:
  - semantic search with/for "things, not strings"
  - question answering and fact checking
- constructing comprehensive knowledge bases
- sentiment mining (e.g. about products or political debates)
- context-aware recommendations
- business analytics

### **Use-Case Example: News Search**



### **Use-Case Example: News Search**

stics

📤 David Bowie 🗴 👘 jones



# Image: Top trending entitiesImage: John LennonImage: Kanye WestImage: Kanye West<

### Beenes Aires Benald.com The business of being Bowie

BuenosAiresHerald.com - Tue Jan 12 01:00:00 CET 2016

The business of being Bowie . Visionary bohemian was his own greatest product Tuesday, January 12, 2016 The business of being Bowie By leonid bershidsky, mark gilbert Bloomberg News (\*) By Leonid Bershidsky, Mark Gilbert David Bowie was that rare kind of rock star: you didn't have to like his music to admire him. Bowie was a business visionary like the ones who shaped Silicon Valley who just didn't see the point of building companies: he was his own greatest product. Bowie , born David Jones , died Sunday at 69, a bohemian who had amassed a vast fortune thanks in part to his many instinctive firsts. Cataloguing them would probably be a futile exercise, but some are worth recalling in the same spirit as his fans now play his songs to remember him. Bowie 's decision to get into rock music was the result of a conscious search for a way to blend business and creativity. Here's how he described it in a BBC interview: "I wanted to be thought of as someone who was very

much a trendy person, rathtrend, I wanted to be the in people on to new ideas an govern everything around use the easiest medium to and to add bits and pieces it, I was my own medium." 1



#### David Bowie

David Robert Jones (born 8 January 1947), known by his stage name David Bowie, is an ...

ጵ

of new stage personae, but not no ability arway a stay a nucleahead of the curve. In 1973, he killed off Ziggy Stardust, perhaps the most elaborate of his stage images; 29 years later, just as the

### http:/stics.mpi-inf.mpg.de

Ð

15-5

### **Use-Case Example: Biomedical Search**

	MEDIE — See what causes cancer?
	subject     verb     object     search     clear     stop       cause     colon cancer     advanced search
Resi » shi	ilts 1-10 for cause colon cancer
Sort	by ORank ODate Sort
ser	tence article table Show 10 results  Show 10 results  Subject verb object gene disease
sh	ow next »
1.	DIXDC1 targets p21 and cyclin D1 via PI3K pathway activation to promote colon cancer cell proliferation, »xmL Lei Wang, XI-Xi Cao, Gi Chen, Teng-Fang Zhu, Hong-Guang Zhu, Li Zheng, pp. 1801-8, Volume 100, Issue 10, Cancer science, 2009 [PMID:19572978] Both siRNA knockdown of DIXDC1 and blocking the PI3K pathway using a specific inhibitor cancer cells, «smul
2.	Induction and down-regulation of Sox17 and its possible roles during the course of gastrointestinal tumorigenesis, »XML Yu-Chen Du, Hiroko Oshima, Keisuke Oguma, Takanori Kitamura, Hiraku Itadani, Takashi Fujimura, Ying-Shi Piao, Tanihiro Yoshimato, Toshinari Minamoto, Hidehito Kotani, Makoto M Taketo, Masanobu Oshima, pp. 1346-57, Volume 137, Issue 4, Gastroenterology, 2009 [PMID:19549530]
	BACKGROUND & AIMS: The activation of Wnt/beta-catenin signaling causes the development of gastric and colon concern. »53ML
3.	A colorectal cancer expression profile that includes transforming growth factor beta inhibitor BAMBI predicts metastatic potential, <a href="https://www.wow.wow.com/static-automatic-back-com/static-back-com</td>
	In mice, overexpression of <b>BAMBI</b> caused colon concer cells to form tumors that metastasized more frequently to liver and lymph nodes than control cancer cells. »SYML
4.	Activation of aminoimidazole carcinogens by nitrosation: mutagenicity and nucleotide adducts, »xmL Terry V Zenser, Vijaya M Lakshmi, Herman A J Schut, Hui-jia Zhou, P David Josephy, pp. 109-15, Volume 673, Issue 2, Mutation research, 2009 [PMID:19449459]
	2-Amino-3-methylimidazo[4,5-f]quinoline (IQ) and 2-amino-3,8-dimethylimidazo[4,5-f]quinoxaline(MeIQx) are heterocyclic amines (HCAs) derived from high temperature cooking of meat and thought to cause color cause in humans.
5.	Inhibition of human sirtuins by in situ generation of an acetylated lysine-ADP-ribose conjugate. »xmL Tomomi Asaba, Takayoshi Suzuki, Rie Ueda, Hiroki Tsumoto, Hidehiko Nakagawa, Naoki Miyata, pp. 6989-96, Volume 131, Issue 20, Journal of the American Chemical Society, 2009 [PMID:19413317]
	Compound 2k also caused a dose-dependent increase of p53 acetylation in human colour cancer HCT116 cells, indicating inhibition of SIRT1 in the cells. »3000L
6.	Mismatch repair polymorphisms and risk of colon cancer, tumour microsatellite instability and interactions with lifestyle factors, »xmL P T Campbell, K Curtin, C M Ulrich, W S Samowitz, J Bigler, C M Velicer, B Caan, J D Potter, M L Slattery, pp. 661-7, Volume 58, Issue 5, Gut, 2009 [PMID:18523027]
	BACKGROUND: Germline mutations in DNA mismatch repair (MMR) genes cause Lynch synchrome colon cancers. asXML
7.	Preclinical development of the nicotinamide phosphoribosyl transferase inhibitor prodrug GMX1777, »XML Pierre Beauparlant, Dominique Bédard, Cynthia Bernier, Helen Chan, Karine Gilbert, Daniel Goulet, Michel-Olivier Gratton, Manon Lavoie, Anne Roulston, Emilie Turcotte, Mark Watson, pp. 346-54, Volume 20, Issue 5, Anti-cancer drugs, 2009 [PMID:19369827]
	Consistent with the requirement for a prolonged exposure for cytotoxicity in vitro, a dose of 75 mg/kg of GMX1777 administered as two bolus intravenous injections in 1 day were not effective at reducing the growth of multiple myeloma (IM-9) tumors over a 24 h intravenous infusion carcinome (HCT-116) model. as the compared of the growth of multiple myeloma (IM-9) tumors over a 24 h intravenous infusion carcinome (HCT-116) model.
	IRDM WS 2015 http://www.nactem.ac.uk/medie/search.cgi 15-6

# **Use-Case Text Analytics: Disease Networks**



### But not so easy with:

diabetes mellitus, diabetis type 1, diabetes type 2, diabetes insipidus, insulin-dependent diabetes mellitus with ophthalmic complications, ICD-10 E23.2, OMIM 304800, MeSH *C18.452.394.750, MeSH* D003924, ...

K.Goh, M.Katsiek, D.Valle, B.Childs, M.Vidal, A.Barabasi: The Human Disease Network, PNAS, May-2007

Neurological
 Nutritional
 Ophthamological

Psychiatric
 Renal
 Respiratory

Skeletal
 multiple
 Unclassified

# **Methodologies for IE**

- Rules & patterns, especially regular expressions
- Pattern matching & pattern learning
- **Distant supervision** by dictionaries, taxonomies, ontologies etc.
- Statistical machine learning: classifiers, HMMs, CRFs etc.
- Natural Language Processing (NLP): POS tagging, parsing, etc.
- Text mining algorithms in general

# **IE Example: Web Pages to Entity Attributes**

realtor.com® ⊙ _ B	BUY RENT MORTGAGE Find REA	ALTORS <sup>®</sup> LOCAL NEWS & ADVICE		<b>Ytrulia</b> Buy Sell Rent	Mortgage Find an Agent M	Aore For Professionals		
San Diego, CA	Q Any Price - Any Beds -	Any Baths - More Filters -			San Diego, CA	Q Search Min Price 🗸	to Max Price 🗸	All Beds 👻 All
San Diego, CA Newest Hom California > San Diego County > Sar	n <b>e Listings</b> n Diego				3,032 San Diego, CA Hom Sort by: Featured ~	es For Sale & Real Estate Save search		List Map
5/T Homes Sort by Newest Listing:	S ↓ 14170 Rancho Vista Bnd, San Diego, North City \$2,749,000 4 bd • 4+ ba • 6,595 sq ft • 0.5 Single Family Home Brokered by Klein Real Estate	CA 92130 4 acres lot	List Map		NEW HOMES	1225 Pacific Beach Dr       Pacfic Beach, San Diego, CA 92109       2 bdi 2 bai 981 sqft       Condo       price       Call or email Helen Grebenc       B	Map \$2,; Sking for a great investment pro mary residence at the beach? Lc Fre	\$499,000 268/mo Get Pre-Approved perty, second hame or zok no further as this g View Details
realtor.com <sup>*</sup> ⊙	247 Orange Ave, Coronado, CA 9211 Coronado BUY RENT MORTGAGE Find REA	8 Altors <sup>ø</sup> local news&advice	~		Coogle Map data @2016 Goo	3+ Bd 1955+ sqft 2.5+ Ba New Community ogle	New Homes in Enc heart of Encini N	initas Located in the <b>Aore</b>
<ul> <li>K Back to search   California &gt; S</li> <li>Home For Sale – ACTIVE</li> <li>14170 Rancho V</li> <li>San Diego, CA 92130</li> <li>4 beds • 4 full , 1 half bath</li> </ul>	San Diego County -> San Diego -> 14170 Rancho Vi <b>ista Bnd</b> hs • <b>6,595</b> sq ft • <b>0.54</b> acres lot	sta Bnd	Free Moving Planner \$2,749,000 Estimate Payment   View Rates			FEATURED 3877 Riviera Dr Pactic Beach, San Diego, CA 92109 2 bd 2 ba 1,176 soft Single-Family Home be Call or cmail Jowille Team Keller Williams - Carlsbad	Map \$2, Crown Point/PB, across from Sai ach, this single story living 2 Bec	\$599,000 722/mo Get Pre-Approved II Bay & steps from the froom/ 2 Baths unit View Details
New Property Details Scho	ools & Neighborhood Payment Options Proper by the original owners.	ty History	ten 🖈		Sell Your Home Home Evaluation	Preferred Lender Guaranteed Sale Pr	ogram MORE - 🦨 R	Register 💄 Sign In
	Style Mediterranean Spanish     Single family home     Year built, 2004     Price/Sq FE S417     2 days on realtor com <sup>®</sup> Status. ACTIVE	Earl Warren Middle School     Torrey Pines High School	s		<b>1,645 Homes Fo</b> Sort By Highest Price • Per V <prev 1="" 138="" 2="" 3="" next2<="" td=""><td>Dr Sale in San Diego</td><td>6 7</td><td>Q ∷≣ SS Map List Grid</td></prev>	Dr Sale in San Diego	6 7	Q ∷≣ SS Map List Grid
Features	Bedrooms - Bedrooms: 4 - Master Bedroom Dimensions: 23x20 - Bedroom 2 Dimensions: 16x13	Bedroom 3 Dimensions: 16x13     Bedroom 4 Dimensions: 16x13	>	1,645 San Diego Homes for Sale           Search All         Type keywords           Separatemultiple terms with commun ()           Beds         Baths           1+         Any	Active	95,000 side covenant home is elegantly finished and offers a net kitchen w/a large island and nook which open to San Gorgonio Street, San Diego, CA 92106	Beds Bath MER in its own wing w/ separate his/ the g <u>View Full Details</u>	s SqFt /hers amenities, light sector 5,600
	Bathrooms  Baths Full: 4  Kitchen and Dining	• Baths Half: 1		Min Price Max Price 150000 5000000 Property Type Any •	Active State	75,000 Ish style la playa mansion built in 2003 w/Unobstruct Ioma w/Ocean view guest quarter! includes infinity g Provided By Carolyn Yarbrough of Pacific Sotheby's int	Beds Bath ad panoramic views of downtown, b: . <u>View Full Details</u> ." <i>Realty (CA DRE: 01386981)</i>	s SqFt ay, coronado & the point of
17. Presen	Ixtrinen pumensions: 1xx24     Diting Room Dimensions: 13x19     Other rooms     Living Room Dimensions: 23x16     Family Room Dimensions: 31x24	Breakfrast Area Dimensions: 16x10     Great Room     Guest Maid	Street. val Estate	SEARCH More Options	6 25 ST02 S4,95 Overt Medit Medit	Meadows Del Mar, San Diego, CA 92130 50,000 coking luth coking luth coastal canyon and the waterway at the eramean manse offers a level of luxury rarely seen in g Provided By Greg Noonan of Berkshire Hathaway Hon	6 7 Beds Bath 17th hole of The Grand Golf Club, thi any home toda <u>View Full Details</u> neservice (CA DRE: 00655720)	7,995 is SqFt is magnificent
IRDM WS 2	• Extra Room 1 Dimensions: 18x16 • Bedroom Entry Level • Breakfast Area • Den Morev 2015	Library     Master Bdrm 2     Master Retreat     MBR Entry Level		San Diego Home Report MORE Overview Arenage Price 2945,820 Total Listings 1719	© 25 945 F Starting O Active	Font St 2201, San Diego, CA 92101 50,000 Le Penthouse opportunity in the highly sought after F a District. This North Tower delight will not cease to a g Provided By Tami Foller of Ascent Real Estate (CA DRE	5 6 Beds Bath Renalssance community located in the ma <u>View Full Details</u> 00000767)	s SqFt re heart of Downtown's 15-9

# **IE Example: Web Pages to Entity Attributes**



Acced 20 44 44270

# **IE Example: Text to Opinions on Entities**



Find Movies, TV shows, Celebrities and more... Movies, TV Celebs, Events News

& Photos





Search movies, TV, actors, more...

Q

#### Reviews & Ratings for Star Wars: Episode VII - The Force Awakens More at INDePro »

& Showtimes

Filter: Loved It • Hide Spoilers:

Interleaved...

Reviews from users who voted this title more than 8.4. Reviews from users who voted this title more than 8.4.

#### Page 4 of 98: 4 [1] [2] [3] [4] [5] [6] [7] [8] [9] [10] [11] [12] [13] [14] 🕨

Index 973 matching reviews (3317 reviews in total)

#### 7 out of 26 people found the following review useful:



I feel like this is brand new, with new actors and ideas being thrown in the mix. Yet it still feels like Star Wars. This burst of new life into Star Wars is defiantly worth a watch despite complaints of similarities to episode IV, whether a Star Wars fan or not.

I love how they have brought back the older characters (e.g Han Solo, Leia, ...) To pass on the mantle of Star Wars to some new characters. As much as I love this film it would be nice to have some more light sabre fight scenes but I feel as though this can be overlooked due to the fact that this is the first instalment of the series and brings you up to pace with the current affairs of a galaxy far far away a long time ago.

Was the above review useful to you? Yes No

6 out of 25 people found the following review useful:



Author: Alex Rivett from Peterborough, England 22 December 2015

#### \*\*\* This review may contain spoilers \*\*\*

Being a massive Star Wars fan, I went along to the midnight launch in 3D (UK) and I was completely blown away from start to finish I didn't get to sleep until 4am and I was due in work later that morning... I've been back to see it again, it was that good!

But I have to say this movie had all of the ingredients you need for a major blockbuster and I salute you Mr Abrams for putting together a stunning masterpiece for all of us fans. Not only has it re-energized the Star Wars universelfranchise but it has introduced us to some exciting new characters in Rey, Finn, Kylo Ren and BB-8 - to name a few as well as reuniting us with the faces from the past! Can I just mention that Daisy Ridley is truly fantastic in her role as Rey, she has made me fall in love with her character and I now have a new Star Wars character fave!

#### AUDIENCE REVIEWS FOR STAR WARS: EPISODE VII - THE FORCE AWAKENS

#### \*\*\*\*

Thought the scriptwriting was brilliant, but I didn't enjoy the actual sets and scenes so much.



Letitia Lew ★ Super Reviewer

#### \*\*\*

Rey, a young smuggler, is thrust into a battle between the First Order and the resistance when she teams up with a storm trooper who suffered a crisis of conscience.

The new entry into the Star Wars universe is profoundly derivative, essentially an updated retelling of A New Hope, and while ignoring the backstory about the First Order largely mutes

#### Read More



Jim Hunter ★ Super Reviewer

#### \*\*\*\*

Extraordinarily faithful to the tone and style of the originals, The Force Awakens brings back the Old Trilogy's heart, humor, mystery, and fun. Since it is only the first piece in a new three-part journey it can't help but feel incomplete. But everything that's already there, from the stunning visuals, to the thrilling action sequences, to the charismatic new characters,

#### Read More



Matthew Samuel Mirliani ★ Super Reviewer

#### \*\*\*

JJ Abrams is very good and knowing what his audience wants and giving just that to them. He is not great, however, because he rarely shows us something we didn't know we wanted. This film derives a lot from the first Star Wars, and just goes along as you might expect, yet it is still very enjoyable because it's Star Wars. The old faces were cool to see, and the new ones do

#### Read More



View All Audience Reviews (13449) 🔪

### **IE Example: Web Pages to Facts & Opinions**

LibraryThing	Home Groups Talk Zeitgeist		Member: g33kgrrl			
	1Q84 by Haruki Murakami		Collections	Your library (1,609), Read but unowned (19), Chicago (13), Cookbooks (3), Crafty (7), Textbooks (33), Currently reading (9), To read (67), New To Read (189), Favorites (176), On Loan (3), Nursing (40), Sold/Given Away (135), Ebook (57), Library Books (6), Finish unread books (207), Will's books (2), Children's Books (15), Zines (1), All collections (1,761)		
-	_,		Reviews	156 reviews		
HARUKI	Other authors: See the other authors section. Series: 1Q84 (1-3)		Tags	favorites (184), comic (173), non-fiction (114), read2008 (97), gift (97), read2007 (95), read2006 (88), read2013 (86), read2009 (83), read2010 (71) — see all tags		
			Media	Not set (12), Book (1,744), Paper Book (1,642), Audiobook (2), Ebook (45), Other (5), Software (1)		
Common Knowledge		Member revie	Clouds	tag cloud, author cloud, tag mirror		
You must log in to edit (	Common Knowledge data.	English (195) Spa	About me	Once when I was very young, I built a chair out of my Baby Sitter's Club books. I then sat in it to read more books.		
For more help see the C	ommon Knowledge help page.	Showing 1-5 of 195		I've since given up on making furniture out of books, but my dream is to someday have a library with in-wall shelves.		
Series (with order)	1Q84 (3.5 1-3)	I love Murakar	Groups	All Things Discworldian - The Guild of Pratchett Fans, Chicagoans		
		Jen.ODriscoll.Le	Venues	Favorites   Visited		
Canonical title	1Q84		Favorite bookstores	Book Loft, Brookline Booksmith, Chicago Comics, City Lit Books, Myopic Books, Uncharted Books		
Original title	1084 (ichi-kyû-hachi-yon)	I love Murakar	Favorite libraries	Billy Ireland Cartoon Library & Museum (The Ohio State University), Chicago Public Library - Logan Square Branch		
		Jen.ODriscoll.Le	Also on	Flickr, Last.fm, Ravelry, Twitter		
Alternative titles		I liked the stor	Membership	C LibraryThing Early Reviewers/Member Giveaway		
Original publication	2009 (vol 1-2)	a lot of unnece audio so I had	Location	Chicago IL		
date	2010 (vol 3)	eadieburke   Jai	Favorite authors	Not set		
	2011-10-25 (1-3 · English)		Account type	public, lifetime		
People/Characters	Masami Aomame	5 stars and a as dealing wit about Aomam	URLs	/profile/g33kgrrl (profile) /catalog/g33kgrrl (library)		
	Tengo Kawana	keep me think	Member since	Nov 18, 2005		
	Toshibaru Ushikawa	sashinka   Jan 1	Currently reading	House of Leaves: The Remastered Full-Color Edition by Mark Z. Danielewski The Complete Calvin and Hobbes (Calvin & Hobbes) by Bill Watterson		
	Tamaru	I actually thou		Just an Ordinary Day: The Uncollected Stories Of Shirley Jackson by Shirley Jackson American Lion: Andrew Jackson in the White House by Jon Meacham		
	Fairland	character for v descriptions of		The Nibelungenlied: Prose Translation (Penguin Classics) by Anonymous show all (9)		
		an editor in ge				
	(show all 9 items)	g33kgrrl   Jan 6,	, 2016   📾			
Important places	Tokyo, Japan	Showing 1-5 of 195	i (next   show all)			
Important events		Published revi	iew <i>s</i>			
Related movies		13 reviews		add a review		
Awards and honors	Man Asian Literary Prize Longlist (2011)	Murakami nam as long, it's als	e-drops George Orwell's so more fun to read.	laugh-riot 1984 several times. Both books deal with the concept of manipulated realities. And while Murakami's book		
IRDM WS 201:	5New York Times bestseller (Fiction, 2011)	Weekly Alibi, Jo	hn Bear (Jan 26, 2012)	15-12		

### **IE Example: Web Pages to Facts on Entities**



Max Planck Medal (1929) Known for Hawking radiation Instruments Vocals · guitar · bass guitar · Occupation Actress Nails • Mick Jagger • Tina Turner Goethe Prize (1945) Penrose-Hawking theorems keyboards · drums · percussion <sup>1987-present</sup> 15-13 Ooi Hoe Soeng (1996-2010) Years active Mott the Hoople A Brief History of Time (1988) Spouse RDM<sup>rc</sup>WS<sup>7</sup>2015 Lahels Verve · Bizarre · Straight · Spouse(s) Website davidbowie.com 🚱 DiscReet · Zappa · Barking

### **IE Example: Text to Relations**



Max Planck The Nobel Prize in Physics 1918

Max Karl Ernst Ludwig Planck was born in Kiel, Germany, on April 23, 1858, the son of Julius Wilhelm and Emma (*née* Patzig) Planck.

Planck studied at the Universities of Munich and Berlin where his teachers included Kirchhoff and Helmholtz, and received his doctorate of philosophy at Munich in spouse (Max Planck, Marie Merck) He was Privatdozent in Munich from 1880 to 1885, the Associate Professor of Theoretical Physics at Kiel until 1889,

in which year he succeeded Kirchhoff as Professor at-Berlin University, where he remained until his retiren Afterwards he became President of the Kaiser Wilhel for the Promotion of Science, a post he held until 193

He was also a gifted pianist and is said to have at one considered music as a career.

Planck was twice married. Upon his appointment, in to Associate Professor in his native town Kiel he married a friend of his childhood, Marie Merck, w in 1909. He remarried her cousin Marga von Hösslin. Three of his children died young, leaving him with ty bornOn (Max Planck, 23 April 1858) bornIn (Max Planck, Kiel) type (Max Planck, physicist)

advisor (Max Planck, Kirchhoff) advisor (Max Planck, Helmholtz) AlmaMater (Max Planck, TU Munich) plays (Max Planck, piano) spouse (Max Planck, Marga Hösslin)

Person	<b>BirthDate</b>	<b>BirthPlace</b>	•••
Max Planck	4/23, 1858	Kiel	
Albert Einstein	3/14, 1879	Ulm	
Mahatma Gandhi	10/2, 1869	Porbandar	

Person	Award
Max Planck	Nobel Prize in Physics
Marie Curie	Nobel Prize in Physics
Marie Curie	Nobel Prize in Chemistry

### IE Example: Text to Annotations Named Entity Recognition with ANNIE

<u>GATE</u> is an open source infrastructure for developing and deploying software components that process human language. GATE excels at text analysis of all shapes and sizes. From large corporations to small startups, from multi-million research consortia to undergraduate projects. More than €5 million has been invested in GATE development and our objective is to make sure that this continues to be money well spent for all GATE's users.

GATE is distributed with an example Information Extraction system, known as <u>ANNIE</u>, which has formed the basis of many commercial and research systems. While ANNIE is capable of recognising a number of different entity types this simple demo focuses on the annotion of **example**, **O** locations, and **e** organizations.

To try the demo please enter either a URL:

or some free text to process:

Planck was twice married. Upon his appointment, in 1885, to Associate Professor in his native town Kiel he married a friend of his childhood, Marie Merck, who died in 1909. He remarried her cousin Marga von Hösslin. Three of his children died young, leaving him with two sons.

↓ Process Text ↓

▲ Max Karl Ernst ▲ Ludwig Planck was born in ♥ Kiel, ♥ Germany, on April 23, 1858, the son of ▲ Julius Wilhelm and ▲ Emma (ne Patzig) ▲ Planck. ▲ Planck studied at the Universities of ♥ Munich and ♥ Berlin, where his teachers included Kirchhoff and Helmholtz, and received his doctorate of philosophy at ♥ Munich in 1879. He was Privatdozent in ♥ Munich from 1880 to 1885, then Associate Professor of Theoretical Physics at ▲ Kiel until 1889, in which year he succeeded Kirchhoff as Professor at ■ Berlin University, where he remained until his retirement in 1926. Afterwards he became President of the ■ Kaiser Wilhelm Society for the Promotion of Science, a post he held until 1937. He was also a gifted pianist and is said to have at one time considered music as a career. ▲ Planck was twice married. Upon his appointment, in 1885, to Associate Professor in his native town ▲ Kiel he married a friend of his childbood. ● Maria Marck, who died in 1909. He remarried her cousin Marca. ● yon Hösselin. Three of his childbood ★

### http://services.gate.ac.uk/annie/



# **IE Example: Text to Annotations**

### **Open Calais Demo**

Open Calais demo is best viewed in Google Chrome

	30
Language: 🕕	
English	
Topics: ()	
Education (M:F)	100%
Germany (G:3D)	95%
Key Personnel Changes (E:4B)	77%
Retirement / Old Age (M:18)	75%
Film (M:H1)	70%
Music (M:H0)	56%
Living / Lifestyle (M:1T)	56%
Human Interest	50%
Children / Youth Issues (M:NW)	50%
Science (M:V)	47%
Womens Issues (M:J2)	33%
Entertainment Production (TRBC) (B:95)	) 31%
Health / Medicine (M:P)	4%
Entities: 1)	
► Facility	
▼	
Berlin University     Kaiser Wilhelm Society for     Person	
Julius Wilhelm     Karl Ernst Ludwig Planck     Marga von Hösslin     Marie Merck	

View RDF C Back Max Karl Ernst Ludwig Planck was born in Kiel, Germany, on April 23, 1858, the son of Julius Wilhelm and Emma (née Patzig) Planck. Planck studied at the Universities of Munich and Berlin, where his teachers included Kirchhoff and Helmholtz, and received **his** doctorate of philosophy at **Munich** in 1879. He was Privatdozent in Munich from 1880 to 1885, then Associate Professor of Theoretical Physics at Kiel until 1889, in which year he succeeded Kirchhoff as Professor at Berlin University, where he remained until his retirement in 1926.

Afterwards he became President of the Kaiser Wilhelm Society for the Promotion of Science, a post he held until 1937.

He was also a gifted pianist and is said Karl Ernst Ludwig Planck (Person) Relevance: 80% considered music as a career. Count: 20 forenduserdisplay: true Planck was twice married. Upon his app persontype: N/A nationality: N/A to Associate Professor in his native to confidencelevel: 0.99 firstname: Karl he married a friend of his childhood, Ma middlename: Ernst Ludwig lastname: Planck in 1909. He remarried her cousin Marga commonname: Karl Ernst Ludwig Planck

Three of **his** children died young, leaving **him** with two sons

http://www.opencalais.com/opencalais-demo/

# Info Extraction vs. Knowledge Harvesting

Surajit obtained his PhD in CS from Stanford University under the supervision of Prof. Jeff Ullman. He later joined HP and worked closely with Umesh Dayal ...

one source



instanceOf (Surajit, scientist) inField (Surajit, computer science) hasAdvisor (Surajit, Jeff Ullman) almaMater (Surajit, Stanford U) workedFor (Surajit, HP) friendOf (Surajit, Umesh Dayal)

• targeted: hasAdvisor, almaMater



• open: worked for, affiliation, employed by, romance with, affair with, ...

yield-centric harvesting

precision !
 recall

### many sources

IRDM WS 2015

	***	
hasA	dviso	r

Student	Advisor
Surajit Chaudhuri	Jeffrey Ullman
Alon Halevy	Jeffrey Ullman
Jim Gray	Mike Harrison
2	

### almaMater

Student	University
Surajit Chaudhuri	Stanford U
Alon Halevy	Stanford U
Jim Gray	UC Berkeley

# 15.2.1 IE with Rules on Patterns (aka. Web Page Wrappers)

Goal: Identify and extract entities and attributes in regularly structured HTML page, to generate database records

Rule-driven regular expression matching

- regex over alphabet Σ of tokens:
   ε, σ∈Σ, (expr1|expr2), (expr)\*
- Interpret pages from same source (e.g. Web site to be wrapped) as regular language (FSA, Chomsky-3 grammar)
- Specify rules by regex's for detecting and extracting attribute values and relational tuples



Title	Year
The Shawshank Redemption	1994
The Godfather	1972
The Godfather - Part II	1974
Pulp Fiction	1994
The Good, the Bad, and the Ugly	1966

# LR Rules: Left and Right Tokens

L token (left neighbor) pre-filler pattern

fact token *filler pattern*  **R token** (right neighbor) *post-filler pattern* 

<u>Example</u>:  $L = \langle B \rangle, R = \langle B \rangle$   $\rightarrow MovieTitle$   $L = \langle I \rangle, R = \langle I \rangle$   $\rightarrow Year$ produces relation with <HTML> <TITLE>Top-250 Movies</TITLE> <BODY> <B>Godfather 1</B><I>1972</I><BR> <B>Interstellar</B><I>2014</I><BR> <B>Titanic</B><I>1997</I><BR> </BODY> </HTML>

tuples: <Godfather 1, 1972>, <Interstellar, 2014>, <Titanic, 1997>

Rules can be combined and generalized  $\rightarrow$  RAPIER [Califf and Mooney '03]

# Advanced Rules: HLRT, OCLR, NHLRT, etc

<u>Idea</u>: Limit application of LR rules to proper context (e.g., to skip over HTML table header)

<TABLE>

<TR><TH><B>Country</B></TH><TH><I>Code</I></TH></TR> <TR><TD><B>Godfather 1</B></TD><TD><I>1972</I></TD></TR> <TR><TD><B>Interstellar</B></TD><TD><I>2014</I></TD></TR> <TR><TD><B>Titanic</B></TD><TD><I>1997</I></TD></TR> </TABLE>

- **HLRT rules** (head left token right tail) apply LR rule only if inside HT (e.g., H = <TD> T = </TD>)
- OCLR rules (open (left token right)\* close):
   O and C identify tuple, LR repeated for individual elements
- **NHLRT** (nested HLRT): apply rule at current nesting level, open additional levels, or return to higher level

# **Rules for HTML DOM Trees**

- Use HTML tag paths from root to target element
- Use more powerful operators for matching, splitting, extracting

INTL BUS MACHINE (NYSE:IBM) - More Info: <u>News</u> , <u>SEC</u> , <u>Msgs</u> , <u>Profile</u> , <u>Research</u> , <u>Insider</u>								
Last Trade	Change		Prev Cls	Volume	Div Date	300 IBM 26-May-1999 (C) Yahoo!		
2:54PM · 114 <sup>7</sup> / <sub>16</sub>	-3 <sup>11</sup> / <sub>N</sub> (-3	8.12%)	236 <sup>1</sup> / <sub>4</sub>	8,390,700	May 26	299		
Dev's Range	Bid	Ask	Open	Avg Vol	Ex-Div			
112 <sup>5</sup> / <sub>8</sub> 116 <sup>7</sup> / <sub>8</sub>	N/A	N/A	116 <sup>11</sup> / <sub>16</sub>	5,444,363	May 27	Jul Sep Nov Jan Mar May		
52-week Range	Earn/Shr	P/E	Mkt Cap	Div/Shr	Yield	Small: [ <u>1d   5d</u>   1y   <u>none</u> ]		
53 - 123	3.53	33.46	103.8B	0.48	0.41	Big: [1d   5d   3m   1y   2y   5y   max ]		

Source: A. Sahuguet, F. Azavant: Looking at the Web through <XML> glasses, http://db.cis.upenn.edu/research/w4f.html

### Example: extract the volume table.tr[1].td[\*].txt, **match /Volume/** extract the % change table.tr[1].td[1].txt, **match /[(](.\*?)[)]/** extract the day's range for the stock: table.tr[2].td[0].txt, **match/Day's Range (.\*)/, split /-/**

match /.../, split /.../ return lists of strings

# Learning Regular Expressions (aka. Wrapper Induction)

<u>Input</u>: Hand-tagged examples of a regular language

<u>Output</u>: (Restricted) regular expression for the language of a finitestate transducer that reads sentences of the language and outputs token of interest

Example:

This apartment has 3 bedrooms. <BR> The monthly rent is \$ 995.
This apartment has 4 bedrooms. <BR> The monthly rent is \$ 980.
The number of bedrooms is 2. <BR> The rent is \$ 650 per month.
yields <sup>\*</sup> <digit> \* "<BR>" \* "\$" <digit>+ \*
as learned pattern
Problem: Grammar inference for general regular languages is hard.
→ restricted class of regular languages
(e.g. WHISK [Soderland 1999], LIXTO [Baumgartner 2001])



Source: R. Baumgartner, Datalog-related Aspects in Lixto Visual Developer, 2010, IRDM WS 201http://datalog20.org/slides/baumgartner.pdf



Source: R. Baumgartner, Datalog-related Aspects in Lixto Visual Developer, 2010, IRDM WS 20nttp://datalog20.org/slides/baumgartner.pdf

### **Limitations and Extensions of Rule-Based IE**

- Powerful for wrapping **regularly structured web pages** (e.g., template-based from same Deep Web site / CMS)
- Many **complications** with real-life HTML (e.g., misuse of tables for layout)
- Extend flat view of input to **trees**:
  - hierarchical document structure (DOM tree, XHTML)
  - extraction patterns for restricted regular languages on trees (e.g. fragements and variations of XPath)
- **Regularities with exceptions** are difficult to capture
  - Identify positive and negative cases and use statistical models

# **15.2.2 IE with Statistical Learning**

For heterogeneous Web sources and for natural-language text

- NLP techniques (PoS tagging, parsing) for tokenization
- Identify patterns (regular expressions) as features
- **Train statistical learners** for segmentation and labeling (e.g., HMM, CRF, SVM, etc.) augmented with lexicons
- Use learned model to **automatically tag** new input sequences
- Example for labeled training data: *The WWW conference in 2007 takes place in Banff in Canada. Today 's keynote speaker is Dr. Berners-Lee from W3C.*  with tags of the following kinds: *event, person, location, organization, date*

# **IE as Boundary Classification**

<u>Idea</u>: Learn **classifiers** to **recognize start token** and **end token** for the facts under consideration. Combine multiple classifiers (ensemble learning) for more robust output.

Example: There will be a talk by Alan Turing at the University at 4 PM. Prof. Dr. James Watson will speak on DNA at MPI at 6 PM. The lecture by Francis Crick will be in the IIF at 3:15 today.

Trained classifiers test each token (with PoS tag, LR neighbor tokens, etc. as features) for two classes: begin-fact, end-fact erson

**!place** 

# **IE as Text Segmentation and Labeling**

<u>Idea</u>: Observed text is concatenation of structured record with limited reordering and some missing fields

Examples: Addresses and bibliographic records



Source: S. Sarawagi: Information Extraction, 2008

### → Hidden Markov Model (HMM) !

### **HMM Example: Postal Address**

<u>Goal</u>: Label the tokens in sequences *Max-Planck-Institute, Stuhlsatzenhausweg 85* with the labels **Name, Street, Number** 

 $\Sigma = {$  "MPI", "St.", "85" $}$ S = {Name, Street, Number} // output alphabet
// (hidden) states



### **HMM Example: Postal Addresses**



Source: Eugene Agichtein and Sunita Sarawagi, Tutorial at KDD 2006

### **Basics from NLP for IE** (in a Nutshell)

Surajit Chaudhuri obtained his PhD from Stanford University under the supervision of Prof. Jeff Ullman



### NLP: Part-of-Speech (POS) Tagging

Tag each word with its grammatical role (noun, verb, etc.) Use HMM or CRF trained over large corpora

### POS Tags (Penn Treebank):

CC coordinating conjunction CD cardinal number DT determiner EX existential *there* FW foreign word IN preposition or subordinating conjunction JJ adjective JJR adjective, comparative JJS adjective, superlative LS list item marker MD modal NN noun NNS noun, plural NNP proper noun NNPS proper noun, plural PDT predeterminer POS possessive ending PRP personal pronoun

PRP\$ possessive pronoun **RB** adverb RBR adverb, comparative RBS adverb, superlative **RP** particle SYM symbol TO to UH interjection VB verb, base form VBD verb, past tense VBG verb, gerund or present participle VBN verb, past participle VBP verb, non-3rd person singular present VBZ verb, 3rd person singular present WDT wh-determiner (which ...) WP wh-pronoun (what, who, whom, ...) WP\$ possessive wh-pronoun WRB wh-adverb

http://www.lsi.upc.edu/~nlp/SVMTool/PennTreebank.html

# **HMM for Part-of-Speech Tagging**



How to find the best sequence of POS tags for sentence "We can buy a can"?  $(PRP) \rightarrow (MD) \rightarrow (VB)$ 



# (Linear-Chain) Conditional Random Fields (CRFs)

- Extend HMMs in several ways:
- exploit **complete input sequence** for predicting state transition, not just last token
- use **features** of input tokens
  - (e.g. hasCap, isAllCap, hasDigit, isDDDD, firstDigit,
  - isGeoname, hasType, afterDDDD, directlyPrecedesGeoname, etc.)
- For token sequence  $x=x_1...x_k$  and state sequence  $y=y_1..y_k$ HMM models joint distr.  $P[x,y] = \prod_{i=1..k} P[y_i|y_{i-1}] * P[x_i|y_i]$ CRF models conditional distr. P[y|x]with conditional independence of non-adjacent  $y_i$ 's given x





# **CRF** Training and Inference



.2

graph structure of conditional-independence assumptions leads to:

$$P[y|x] = \frac{1}{Z(x)} exp\left(\sum_{j=1}^{m} \lambda_j \sum_{t=1}^{T} f_j(y_{t-1}, y_t, x)\right)$$

where j ranges over feature functions and Z(x) is a normalization constant

parameter estimation with n training sequences: MLE with regularization

$$\log L(\theta) = \sum_{i=1}^{n} \sum_{t=1}^{T} \sum_{j=1}^{m} \lambda_j f_j(y_{t-1}^{(i)}, y_t^{(i)}, x_t^{(i)}) - \sum_{i=1}^{n} \log Z(x^{(i)}) - \sum_{j=1}^{m} \frac{\lambda_j^2}{2\sigma^2}$$

inference of most likely (x,y) for given x: dynamic programming (similar to Viterbi)

CRFs can be further generalized to undirected graphs of coupled random variables (aka. MRF: Markov random field) IRDM WS 2015

### **NLP: Deep Parsing for Constituent Trees**



- Construct syntax-based tree of sentence constituents
- Use **non-deterministic context-free grammars** natural ambiguity
- Use probabilistic grammar (PCFG): likely vs. unlikely parse trees (trained on corpora)



The bright student who works hard will pass all exams.

### **Extensions and variations:**

- Lexical parser: enhanced with lexical dependencies (e.g., only specific verbs can be followed by two noun phrases)
- Chunk parser: simplified to detect only phrase boundaries

# NLP: Link-Grammar-Based Dependency Parsing 🗞

Dependency parser based on grammatical rules for left & right connector



- Parser finds all matchings that connect all words into planar graph (using dynamic programming for search-space traversal)
- Extended to probabilistic parsing and error-tolerant parsing



**O(n<sup>3</sup>)** algorithm with many implementation tricks, and grammar size n is huge

IRDM WS 2015

### **Dependency Parsing Examples (1)**

http://www.link.cs.cmu.edu/link/



LEFT-WALL Alice and Bob met.v in New York on New Year 's.p Eve .

### Selected tags (CMU Link Parser), out of ca. 100 tags (plus variants):

MV connects verbs to modifying phrases like adverbs, time expressions, etc.

O connects transitive verbs to direct or indirect objects

J connects prepositions to objects

B connects nouns with relative clauses

### **Dependency Parsing Examples (2)**

http://nlp.stanford.edu/software/lex-parser.shtml

### Your sentence

The bright student who works hard will pass all exams.

### Tagging

The/DT bright/JJ student/NN who/WP works/VBZ hard/JJ will/MD pass/VB all/DT exams/NNS ./.

### Parse

### Typed dependencies

. . .

```
(ROOT
 (S
    (NP
      (NP (DT The) (JJ bright) (NN student))
      (SBAR
      (WHNP (WP who))
      (S
        (VP (VBZ works)
            (ADJP (JJ hard))))))
 (VP (MD will)
      (VP (VB pass)
        (NP (DT all) (NNS exams))))
 (...)))
```

```
det(student-3, The-1)
amod(student-3, bright-2)
nsubj(works-5, student-3)
nsubj(pass-8, student-3)
rel(works-5, who-4)
rcmod(student-3, works-5)
acomp(works-5, hard-6)
aux(pass-8, will-7)
det(exams-10, all-9)
dobj(pass-8, exams-10)
```

### Selected tags (Stanford Parser), out of ca. 50 tags:

nsubj: nominal subject rel: relative dobj: direct object det: determiner IRDM WS 2015 amod; adjectival modifier rcmod: relative clause modifier acomp: adjectival complement poss: possession modifier

### **Additional Literature for 15.2**

- S. Sarawagi: Information Extraction, Foundations & Trends in Databases 1(3), 2008
- H. Cunningham: Information Extraction, Automatic. in: Encyclopedia of Language and Linguistics, 2005, http://www.gate.ac.uk/ie/
- M.E. Califf, R.J. Mooney: Bottom-Up Relational Learning of Pattern Matching Rules for Information Extraction, JMLR 2003
- S. Soderland: Learning Information Extraction Rules for Semi-Structured and Free Text, Machine Learning Journal 1999
- N. Kushmerick: Wrapper induction: Efficiency and expressiveness, Art. Intelligence 2000
- A. Sahuguet, F. Azavant: Building light-weight wrappers for legacy web data-sources using W4F, VLDB 1999
- R Baumgartner et al.: Visual Web Information Extraction with Lixto, VLDB 2001
- G. Gottlob et al.: The Lixto data extraction project, PODS 2004
- B. Liu: Web Data Mining, Chapter 9, Springer 2007
- C. Manning, H. Schütze: Foundations of Statistical Natural Language Processing, MIT Press 1999