

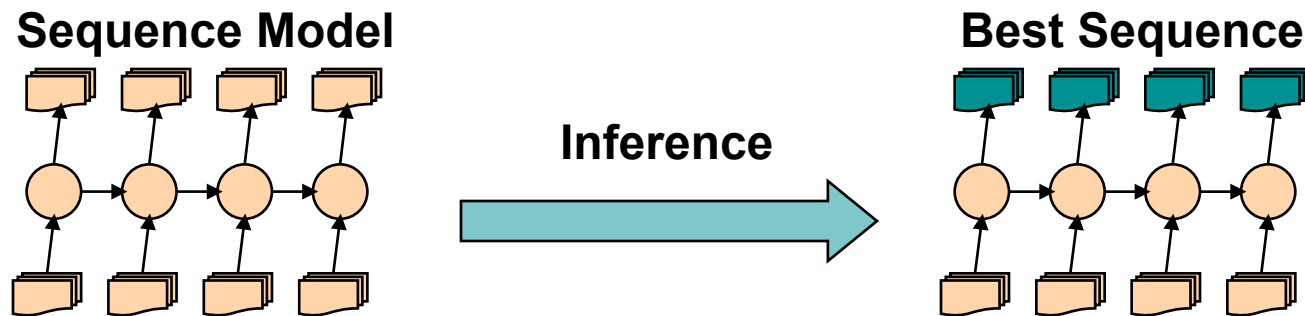
# MEMM inference in systems

- For a Conditional Markov Model (CMM) a.k.a. a Maximum Entropy Markov Model (MEMM), the classifier makes a single decision at a time, conditioned on evidence from observations and previous decisions.
- A larger space of sequences is explored via search

Local Context					Decision Point	Features	
-3	-2	-1	0	+1		$W_0$	22.6
DT	NNP	VBD	???	???		$W_{+1}$	%
The	Dow	fell	22.6	%		$W_{-1}$	fell
						$T_{-1}$	VBD
						$T_{-1}-T_{-2}$	NNP-VBD
						hasDigit?	true
						...	...

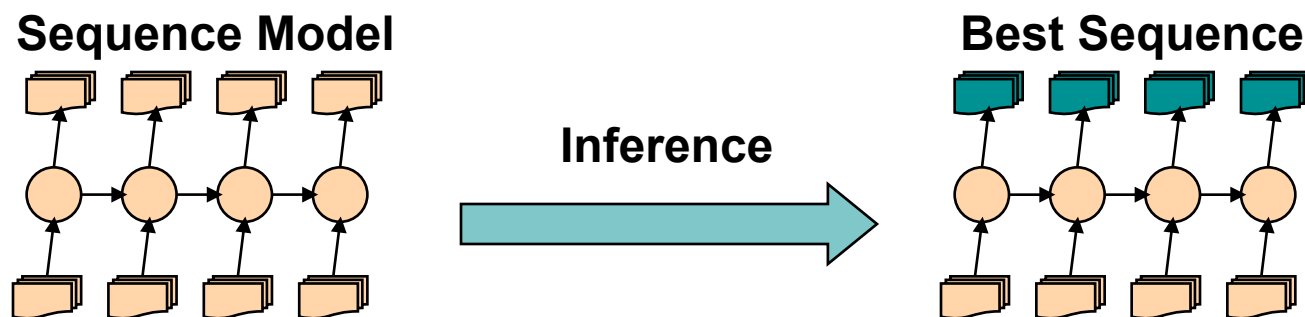
(Ratnaparkhi 1996; Toutanova et al. 2003, etc.)

# Beam Inference



- Beam inference:
  - At each position keep the top  $k$  complete sequences.
  - Extend each sequence in each local way.
  - The extensions compete for the  $k$  slots at the next position.
- Advantages:
  - Fast; and beam sizes of 3–5 are as good or almost as good as exact inference in many cases.
  - Easy to implement (no dynamic programming required).
- Disadvantage:
  - Inexact: the globally best sequence can fall off the beam.

# Viterbi Inference

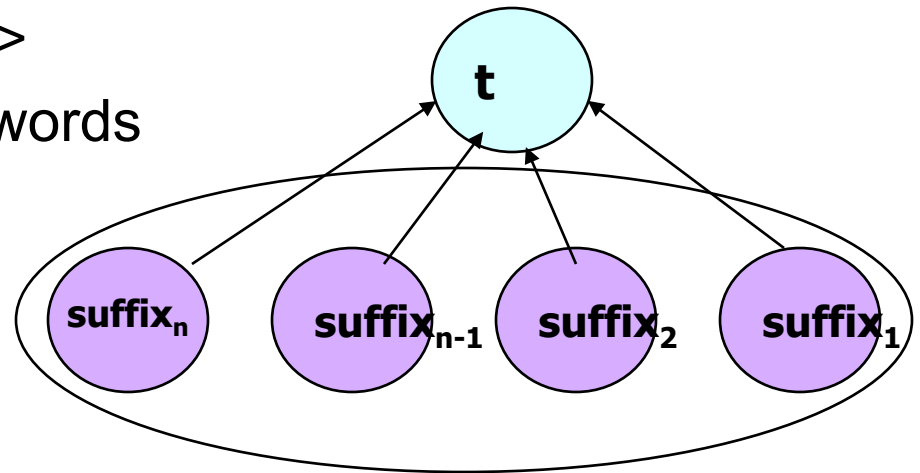


- Viterbi inference:
  - Dynamic programming or memoization.
  - Requires small window of state influence (e.g., past two states are relevant).
- Advantage:
  - Exact: the global best sequence is returned.
- Disadvantage:
  - Harder to implement long-distance state-state interactions (but beam inference tends not to allow long-distance resurrection of sequences anyway).

# HMM Part-of-speech Tagging Models - Brants 2000

- Highly competitive with other state-of-the-art models
- Trigram HMM with smoothed transition probabilities
- Capitalization feature becomes part of the state – each tag state is split into two e.g.  
    NN → <NN,cap>, <NN,not cap>
- Suffix features for unknown words

$$P(w | tag) = P(suffix | tag)(w | suffix) \\ \approx \hat{P}(suffix) \tilde{P}(tag | suffix) / \hat{P}(tag)$$



$$\tilde{P}(tag | suffix_n) = \lambda_1 \hat{P}(tag | suffix_n) + \lambda_2 \hat{P}(tag | suffix_{n-1}) + \dots + \lambda_n \hat{P}(tag)$$

# MEMM Tagging Models -II

- Ratnaparkhi (1996): local distributions are estimated using maximum entropy models
  - Previous two tags, current word, previous two words, next two words, suffix, prefix, hyphenation, and capitalization features for unknown words
- Toutanova et al. (2003)
  - Richer features, bidirectional inference, better smoothing, better unknown word handling

Model	Overall Accuracy	Unknown Words
HMM (Brants 2000)	96.7	85.5
MEMM (Ratn. 1996)	96.63	85.56
MEMM (T. et al 2003)	97.24	89.04

## CRFs [Lafferty, Pereira, and McCallum 2001]

- Another sequence model: Conditional Random Fields (CRFs)
- A whole-sequence conditional model rather than a chaining of local models.

$$P(c | d, \lambda) = \frac{\exp \sum_i \lambda_i f_i(c, d)}{\sum_{c'} \exp \sum_i \lambda_i f_i(c', d)}$$

- The space of  $c$ 's is now the space of sequences
  - But if the features  $f_i$  remain local, the conditional sequence likelihood can be calculated exactly using dynamic programming
- Training is slow, but CRFs avoid causal-competition biases
- These (or a variant using a max margin criterion) are seen as the state-of-the-art these days

# Summary of Tagging

For tagging, the change from generative to discriminative model **does not by itself** result in great improvement

One profits from discriminative models for specifying dependence on **overlapping features of the observation** such as spelling, suffix analysis, etc

A CMM allows integration of rich features of the observations, but can suffer strongly from assuming independence from following observations; this effect can be relieved by adding dependence on following words

This additional power (of the CMM ,CRF, Perceptron models) has been shown to result in improvements in accuracy

The **higher accuracy** of discriminative models comes at the price of **much slower training**

# Biomedical NER Motivation

- The biomedical field contains a large body of information, which is growing rapidly.
  - MEDLINE, the primary research database serving the biomedical community, currently contains over 12 million abstracts, with 60,000 new abstracts appearing each month.
  - There is also an impressive number of biological databases containing information on genes, proteins, nucleotide and amino acid sequences, including *GenBank*, *Swiss-Prot*, and *Fly-Base*; each contains entries numbering from the thousands to the millions and are multiplying rapidly.



# Motivation

- Currently, all of these resources are curated by hand by expert annotators at enormous expense.
- The information overload from the massive growth in the scientific literature has shown the necessity to automatically locate, organize and manage facts relating to experimental results
- Natural Language Processing can aid researchers and curators of biomedical databases by automating these tasks.

# Named Entity Recognition

- General NER vs. Biomedical NER

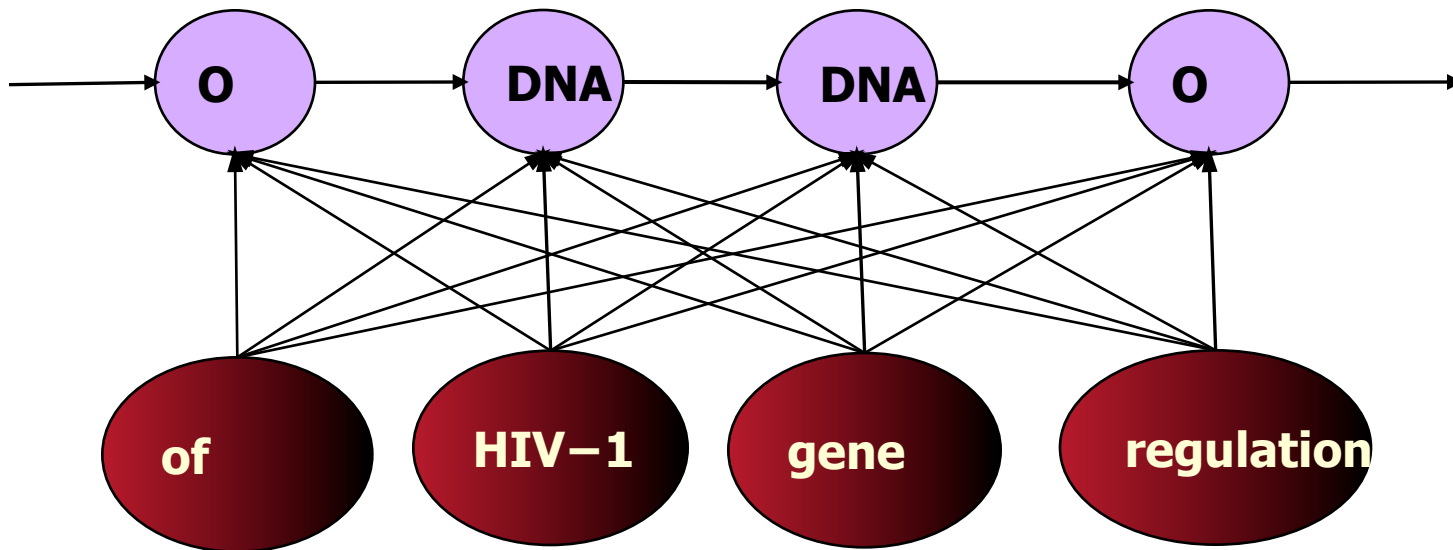
<PER> Christopher Manning </PER> is a professor at <ORG> Stanford University </ORG>, in <LOC> Palo Alto </LOC>.

<RNA> TAR </RNA> independent transactivation by <PROTEIN> Tat </PROTEIN> in cells derived from the <CELL> CNS </CELL> - a novel mechanism of <DNA> HIV-1 gene </DNA> regulation.

# Why is this difficult?

- The list of biomedical entities is growing.
  - New genes and proteins are constantly being discovered, so explicitly enumerating and searching against a list of known entities is not scalable.
  - Part of the difficulty lies in identifying previously unseen entities based on contextual, orthographic, and other clues.
- Biomedical entities don't have strict naming conventions.
  - Common English words such as *period*, *curved*, and *for* are used for gene names.
  - Entity names can be ambiguous. For example, in FlyBase, “clk” is the gene symbol for the “Clock” gene but it also is used as a synonym of the “period” gene.
- Biomedical entity names are ambiguous
  - Experts only agree on whether a word is even a gene or protein 69% of the time. (Krauthammer *et al.*, 2000)
  - Often systematic polysemies between gene, RNA, DNA, etc.

# Maximum Entropy Markov Model



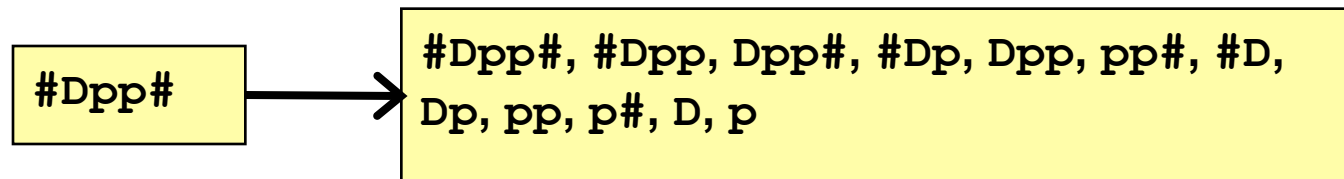
$$P(t | h) = \frac{\exp\left(\sum_{j=1}^m f_j(h, t) \lambda_j\right)}{\sum_{k=1}^K \exp\left(\sum_{j=1}^m f_j(h, t_k) \lambda_j\right)}$$

# Interesting Features

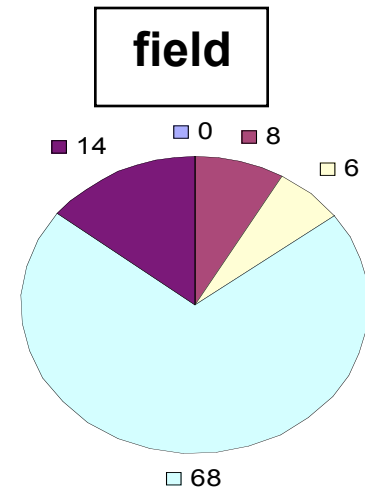
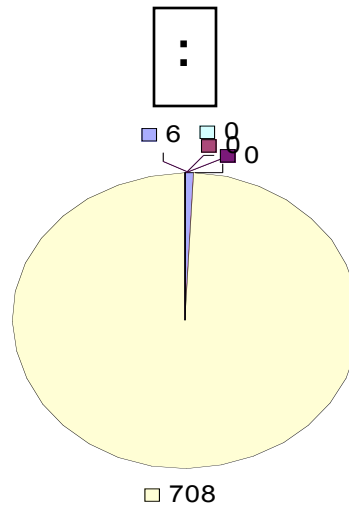
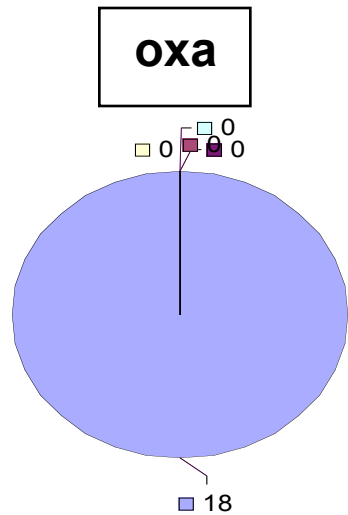
- Word, and surrounding context
- Word Shapes
  - Map words to simplified representation that encodes attributes such as length, capitalization, numerals, Greek letters, internal punctuation, etc.

Varicella-zoster	Xx-xxx
mRNA	xXXX
CPA1	XXXd

- Character substrings



# Features: What's in a Name?



**Cotrimoxazole**

**Wethersfield**

**Alien Fury: Countdown to Invasion**

# Interesting Features

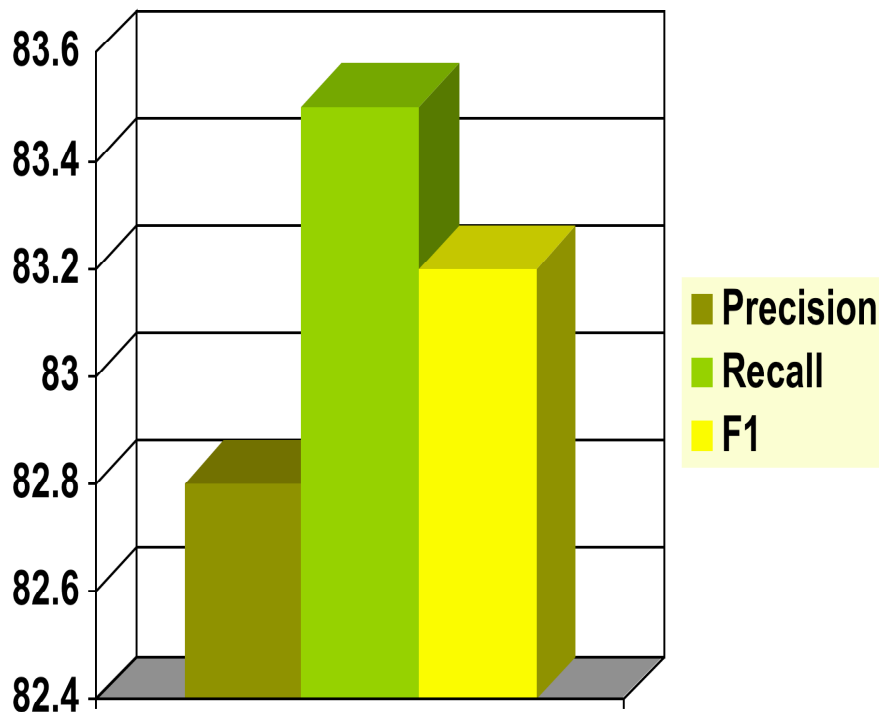
- Part-of-Speech tags
- Parsing information
- Searching the web for the word in a given context
  - *X gene, X mutation, X antagonist*
- Gazetteer
  - list words whose classification is known
- Abbreviation extraction (Schwartz and Hearst, 2003)
  - Identify short and long forms when occurring together in text

... Zn finger homeodomain 2 (Zfh 2)

...

# Finkel et al. (2004) Results

- BioCreative task – Identify genes and proteins



Precision	Recall	F1
81.3%	86.1%	83.6%

$$\text{precision} = \text{tp} / (\text{tp} + \text{fp})$$

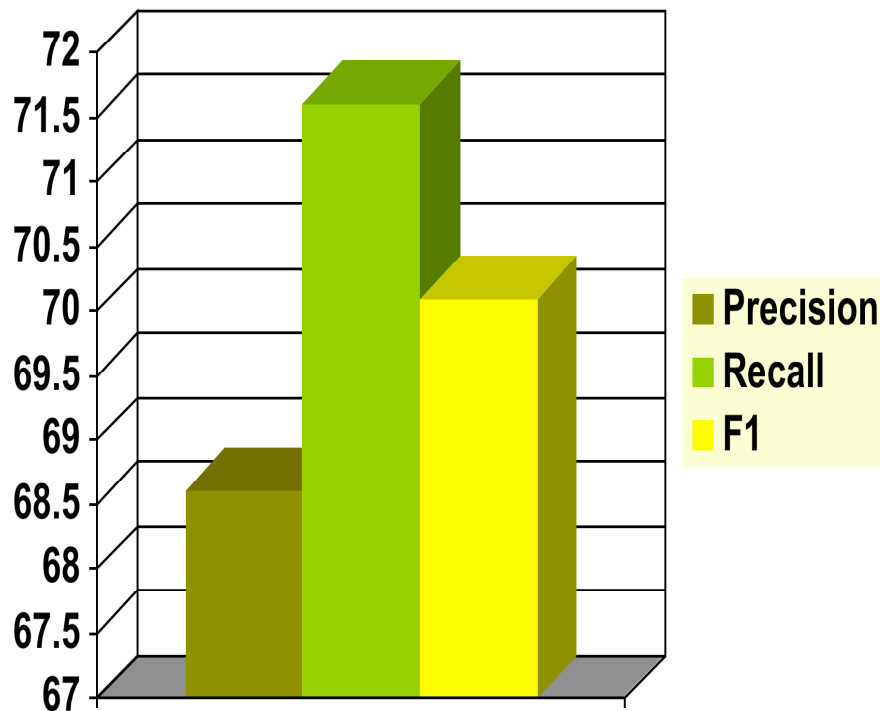
$$\text{recall} = \text{tp} / (\text{tp} + \text{fn})$$

$$\text{F1} = \frac{2(\text{precision})(\text{recall})}{(\text{precision} + \text{recall})}$$



# Finkel et al. (2004) Results

- BioNLP task – Identify genes, proteins, DNA, RNA, and cell types



Precision	Recall	F1
68.6%	71.6%	70.1%

$$\text{precision} = \text{tp} / (\text{tp} + \text{fp})$$

$$\text{recall} = \text{tp} / (\text{tp} + \text{fn})$$

$$\text{F1} = \frac{2(\text{precision})(\text{recall})}{(\text{precision} + \text{recall})}$$

# Information Extraction and Integration

Following slides from:

William Cohen

Andrew McCallum

Eugene Agichtein

Sunita Sarawagi

# The Value of Text Data

- “Unstructured” text data is the primary source of human-generated information
  - Citeseer, comparison shopping, PIM systems, web search, data warehousing
- Managing and utilizing text: information extraction and integration
- Scalability: a bottleneck for deployment
- Relevance to data mining community

# Example: A Solution

job search find employment careers @ FlipDog.com free! - Microsoft Internet Explorer

Address <http://www.flipdog.com/home.html> Go File Edit View Favorites Tools Help Links

**FlipDog.com**

Home Find Jobs Your Account Resource Center Support Employers

Job Search at FlipDog.com: Employment & Career Management

 **647,514**  
Job Opportunities  
from **53,641** Employers

[Find a Job!](#)

[Post Your Resume](#)

**Employers**  
click here for  
Products & Services 

**Pigskin Places**

- Health Care in NY [2,770](#)
- Health Care in MD [1,262](#)
- Sales in NY [3,751](#)
- Sales in MD [958](#)
- Computing in NY [8,050](#)
- Computing in MD [4,114](#)

**Jobs for Sports Fans**

- [Head Football Coach](#)
- [Football Coach](#)
- [Asst. Football Coach](#)
- [High School Football Coach](#)
- [Univ. Asst. Football Coach](#)

**Job Seeker Newsletter**

Enter your e-mail address:

[Sign Me Up!](#)

**Job Seekers: Find your dream job!**

- Check our 'Best Places to Find a Job' [January report](#).
- Open your [FREE account](#) and put your [resume online](#).
- Search 24x7 with our FREE automatic [JobHunters™](#).
- Research our database of over [50,000 employers](#).
- Get [expert advice](#) at our new [Resource Center](#).
- Access [salary surveys/calculators](#), [relocation tools](#), [networking opportunities](#), & [training/testing](#) tools.
- Use FlipDog.com to search jobs right from your desktop! Download [Snippets](#) today!

**Showcase Jobs**

  
Management Recruiters  
of Charlotte North

We provide total staffing solutions in the areas of Human Resources, Compensation, Web-based HR self-service, and Customer Management Systems.


[Learn More](#)





Looking for a Vice President of Academic Affairs to oversee planning, operation and evaluation of the college's academic programs.

[Learn More](#)

powered by **WhizBang!**

 "Top 100 Web Sites"  
PC Magazine, Nov. 2000

 "Top 10 Career Web Site"  
Media Metrix, Sept. 2000

 "Top 10 Job Site"

Start | Microsoft PowerPoint - [sta... | job search find employmen... | 12:12 AM

# Extracting Job Openings from the Web

OPUS International, Inc., an executive search firm focusing on the Food Science industry. - Microsoft Internet Explorer

File Edit View Favorites

Back Forward Stop

Address [http://www.foodscience.com/jobs\\_midwest.html#top](http://www.foodscience.com/jobs_midwest.html#top)

Links AMEX Rewards Time DogHouse My Yahoo!

Welcome

About OPUS

Executive Staff

**Job Listings**

Résumé Form

Job Hunt Hints

Academic Links

Science Fair Help

Industry Assocs.

FAQs

Contact Us

Site Map

e-mail

OPUS INTERNATIONAL INC.

About | Staff | Job

File Edit View Favorites Tools Help

Back Forward Stop Refresh Home Search Favorite

Address [http://www.foodscience.com/jobs\\_midwest.html#top](http://www.foodscience.com/jobs_midwest.html#top)

Links AMEX Rewards Time DogHouse My Yahoo!

Welcome

About OPUS

Executive Staff

**Job Listings**

Résumé Form

Job Hunt Hints

Academic Links

Science Fair Help

Industry Assocs.

FAQs

Contact Us

Site Map

e-mail

**Ice Cream Guru**

If you dream of cold creamy chocolate or coochy boochy cookie, there's a great opportunity for you to maintain and expand this major corporation's high-end ice cream brand. Will be based in the Upper Midwest for about a year. After that, California here I come! Requires a BS in Food Science or dairy, plus ice cream formulation experience. Will consider entry level with an MS and an internship. Contact Susana e-mail 1-800-488-2611

foodscience.com-Job2

JobTitle: Ice Cream Guru

Employer: foodscience.com

JobCategory: Travel/Hospitality

JobFunction: Food Services

JobLocation: Upper Midwest

Contact Phone: 800-488-2611

DateExtracted: January 8, 2001

Source: [www.foodscience.com/jobs\\_midwest.html](http://www.foodscience.com/jobs_midwest.html)

OtherCompanyJobs: foodscience.com-Job1



**Job Openings:**  
**Category = Food Services**  
**Keyword = Baker**  
**Location = Continental U.S.**

The screenshot shows the FlipDog.com website interface. At the top, there are navigation links: Home, Find Jobs, Your Account, and Resource Center. Below these are search options: Return to Results, Modify Search, and New Search. A banner for 'The University Alliance' offers degrees online. Another banner promotes a resume service, ResumeZapper.com. A third banner advertises a 'Breakthrough ebook' about job applications. The main content area shows search results for 'Baker' in 'Food Services' in the 'Continental U.S.' location. The results are displayed in a table with columns for job title, employer, date, and location. The first 15 results are listed below.

Job Title	Employer	Date	Location
<a href="#">Food Pantry Workers</a>	<a href="#">Lutheran Social Services</a>	October 11, 2002	<a href="#">Archbold, OH</a>
<a href="#">Cooks</a>	<a href="#">Lutheran Social Services</a>	October 11, 2002	<a href="#">Archbold, OH</a>
<a href="#">Bakers Assistants</a>	<a href="#">Fine Catering by Russell Morin</a>	October 11, 2002	<a href="#">Attleboro, MA</a>
<a href="#">Baker's Helper</a>	<a href="#">Bird-in-Hand</a>	October 11, 2002	United States
<a href="#">Assistant Baker</a>	<a href="#">Gourmet To Go</a>	October 11, 2002	<a href="#">Maryland Heights, MO</a>
<a href="#">Host/Hostess</a>	<a href="#">Sharis Restaurants</a>	October 10, 2002	<a href="#">Beaverton, OR</a>
<a href="#">Cooks</a>	<a href="#">Alta's Rustler Lodge</a>	October 10, 2002	<a href="#">Alta, UT</a>
<a href="#">Line Attendant</a>	<a href="#">Sun Valley Coporation</a>	October 10, 2002	<a href="#">Huntsville, UT</a>
<a href="#">Food Service Worker II</a>	<a href="#">Garden Grove Unified School District</a>	October 10, 2002	<a href="#">Garden Grove, CA</a>
<a href="#">Night Cook / Baker</a>	<a href="#">SONOCO</a>	October 10, 2002	<a href="#">Houma, LA</a>
<a href="#">Cooks/Prep Cooks</a>	<a href="#">GrandView Lodge</a>	October 10, 2002	<a href="#">Nisswa, MN</a>
<a href="#">Line Cook</a>	<a href="#">Lone Mountain Ranch</a>	October 10, 2002	<a href="#">Big Sky, MT</a>
<a href="#">Production Baker</a>	<a href="#">Whole Foods Market</a>	October 08, 2002	<a href="#">Willowbrook, IL</a>
<a href="#">Cake Decorator/Baker</a>	<a href="#">Mandalay Bay Hotel and Casino</a>	October 08, 2002	<a href="#">Las Vegas, NV</a>
<a href="#">Shift Supervisors</a>	<a href="#">Brueggers Bagels</a>	October 08, 2002	<a href="#">Minneapolis, MN</a>

# What is “Information Extraction”

**As a task:** **Filling slots in a database from sub-segments of text.**

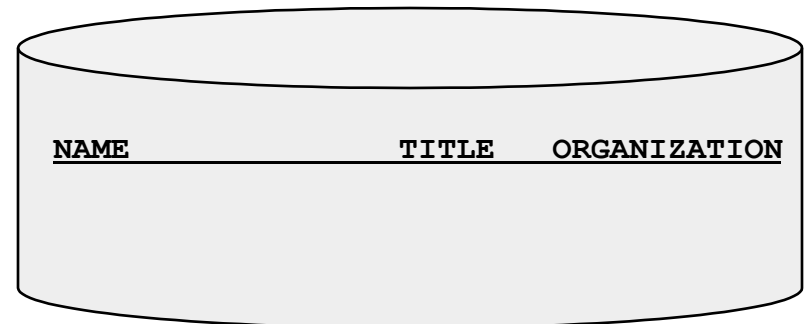
October 14, 2002, 4:00 a.m. PT

For years, Microsoft Corporation CEO Bill Gates railed against the economic philosophy of open-source software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.

Today, Microsoft claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. Gates himself says Microsoft will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"We can be open source. We love the concept of shared source," said Bill Veghte, a Microsoft VP. "That's a super-important shift for us in terms of code access."

Richard Stallman, founder of the Free Software Foundation, countered saying...



# What is “Information Extraction”

**As a task:** Filling slots in a database from sub-segments of text.

October 14, 2002, 4:00 a.m. PT

For years, [Microsoft Corporation](#) [CEO](#) [Bill Gates](#) railed against the economic philosophy of open-source software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.

Today, Microsoft claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. Gates himself says Microsoft will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"We can be open source. We love the concept of shared source," said [Bill Veghte](#), a [Microsoft](#) [VP](#). "That's a super-important shift for us in terms of code access."

[Richard Stallman](#), [founder](#) of the [Free Software Foundation](#), countered saying...



<u>NAME</u>	<u>TITLE</u>	<u>ORGANIZATION</u>
Bill Gates	CEO	Microsoft
Bill Veghte	VP	Microsoft
Richard Stallman	founder	Free Soft..



# What is “Information Extraction”

As a family  
of techniques:

Information Extraction =  
segmentation + classification + clustering + association

October 14, 2002, 4:00 a.m. PT

For years, Microsoft Corporation CEO Bill Gates railed against the economic philosophy of open-source software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.

Today, Microsoft claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. Gates himself says Microsoft will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"We can be open source. We love the concept of shared source," said Bill Veghte, a Microsoft VP. "That's a super-important shift for us in terms of code access."

Richard Stallman, founder of the Free Software Foundation, countered saying...

Microsoft Corporation

CEO

Bill Gates

Microsoft

Gates

Microsoft

Bill Veghte

Microsoft

VP

Richard Stallman

founder

Free Software Foundation

aka “named entity  
extraction”

# What is “Information Extraction”

As a family  
of techniques:

Information Extraction =  
segmentation + classification + association + clustering

October 14, 2002, 4:00 a.m. PT

For years, [Microsoft Corporation](#) [CEO](#) [Bill Gates](#) railed against the economic philosophy of open-source software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.

Today, [Microsoft](#) claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. [Gates](#) himself says [Microsoft](#) will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"We can be open source. We love the concept of shared source," said [Bill Veghte](#), a [Microsoft](#) [VP](#). "That's a super-important shift for us in terms of code access."

[Richard Stallman](#), [founder](#) of the [Free Software Foundation](#), countered saying...

[Microsoft Corporation](#)

[CEO](#)

[Bill Gates](#)

[Microsoft](#)

[Gates](#)

[Microsoft](#)

[Bill Veghte](#)

[Microsoft](#)

[VP](#)

[Richard Stallman](#)

[founder](#)

[Free Software Foundation](#)

# What is “Information Extraction”

As a family  
of techniques:

Information Extraction =  
segmentation + classification + association + clustering

October 14, 2002, 4:00 a.m. PT

For years, [Microsoft Corporation](#) [CEO](#) [Bill Gates](#) railed against the economic philosophy of open-source software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.

Today, [Microsoft](#) claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. [Gates](#) himself says [Microsoft](#) will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"We can be open source. We love the concept of shared source," said [Bill Veghte](#), a [Microsoft](#) [VP](#). "That's a super-important shift for us in terms of code access."

[Richard Stallman](#), [founder](#) of the [Free Software Foundation](#), countered saying...

[Microsoft Corporation](#)  
[CEO](#)  
[Bill Gates](#)

[Microsoft](#)  
[Gates](#)

[Microsoft](#)  
[Bill Veghte](#)  
[Microsoft](#)  
[VP](#)

[Richard Stallman](#)  
[founder](#)  
[Free Software Foundation](#)

# What is “Information Extraction”

As a family  
of techniques:

Information Extraction =  
segmentation + classification + association + clustering

October 14, 2002, 4:00 a.m. PT

For years, [Microsoft Corporation](#) [CEO](#) [Bill Gates](#) railed against the economic philosophy of open-source software with Orwellian fervor, denouncing its communal licensing as a "cancer" that stifled technological innovation.

Today, [Microsoft](#) claims to "love" the open-source concept, by which software code is made public to encourage improvement and development by outside programmers. [Gates](#) himself says [Microsoft](#) will gladly disclose its crown jewels--the coveted code behind the Windows operating system--to select customers.

"We can be open source. We love the concept of shared source," said [Bill Veghte](#), a [Microsoft](#) [VP](#). "That's a super-important shift for us in terms of code access."

[Richard Stallman](#), [founder](#) of the [Free Software Foundation](#), countered saying...

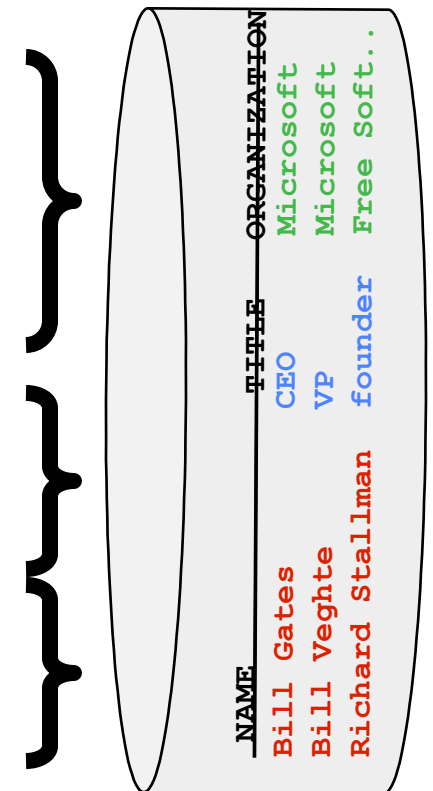
\* [Microsoft Corporation](#)  
[CEO](#)  
[Bill Gates](#)

\* [Microsoft](#)  
[Gates](#)

\* [Microsoft](#)  
[Bill Veghte](#)

\* [Microsoft](#)  
[VP](#)

[Richard Stallman](#)  
[founder](#)  
[Free Software Foundation](#)



# IE is different in different domains!

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

## Newsire

### 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."

The directory structure, link structure, formatting & layout of the Web is its own new grammar.

## Web

www.apple.com/retail

Coming Soon

[Millenia](#)  
Orlando, FL  
Grand Opening, October 19

In the News

[Jaguar Launch Event](#)  
All across the country, thousands of people came to Apple Stores for the nighttime Jaguar launch, lining up in anticipation of the release of Mac OS X v10.2. See what they wore and what they did on this special evening.

Now Open

Arizona [Chandler Fashion Center](#)  
Chandler

Florida [The Falls](#)  
Miami

New York [Crossgates](#)  
Albany

[Biltmore](#)  
Phoenix

[Wellington Green](#)  
Wellington

[Palisades](#)  
West Nyack

[Grand Opening at the Grove](#)  
See pictures from the grand opening weekend of The Grove, the new Apple store in Los Angeles.

[Roosevelt Field](#)  
Garden City

www.apple.com/retail/soho

you to digital cameras, music, email and the Internet. Join us Saturday mornings for a free Getting Started Workshop for new Mac owners.

### [Theater Events](#)

#### Address:

SoHo  
103 Prince Street  
New York, NY 10012  
212-226-3126

#### Store Hours:

Monday - Saturday  
10 a.m. to 8 p.m.  
Sunday  
11 a.m. to 6 p.m.

www.apple.com/retail/soho/theatre.html

Made on a Mac

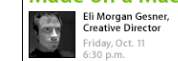
Presentation	Presented By	Date	Time
Andy Milburn Filmmaker	Apple	Wed Oct 16	6:30 p.m.
Jean Miele Landscape Photographer	Apple	Thu Oct 17	6:30 p.m.
William Levin Cartoon Animator	Apple	Mon Oct 21	6:30 p.m.
David Chalk Photographer, Illustrator and Animator	Apple	Thu Oct 24	6:30 p.m.
Day in the Life of Africa David Cohen-Publisher David Turnley-Photographer Douglas Kirkland-Photographer	Apple	Thu Oct 29	6:30 p.m.

Theater

Presentation	Presented By	Date	Time
Getting Started on a Mac -Introduction and Basics -Advanced	Apple	Every Sat	9 a.m. 10 a.m.
Mac OS X v10.2 Jaguar Workshop -Introduction and Basics	Apple	Every Sun	11:00 a.m.
Digital Photography Workshop	Apple	Every Sun	3:00 p.m.

In the News

### Made on a Mac



Eli Morgan Gesner,  
Creative Director  
Friday, Oct. 11  
6:30 p.m.

**Andy Milburn**  
Andy Milburn of the filmmaking partnership tomandandy discusses their groundbreaking audio technology called Q MIX. October 16, 6:30 p.m.

**Jean Miele**  
New York photographer Jean Miele discusses how he creates his large-scale black-and-white landscape photographs using his Power Mac G4, iBook, and three other Mac computers as replacements for the traditional darkroom. October 17, 6:30 p.m.

**William Levin**  
William "Macboy" Levin presents his animated Flash cartoons and discusses the process of their creation. October 21, 6:45 p.m.

# 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.

## Grammatical sentences and some formatting & links

**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

**Frank Huybrechts** - COO  
Mr. Huybrechts has over 20 years of

## Non-grammatical snippets, rich formatting & links

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

## Tables

8:30 - 9:30 AM	<b>Invited Talk: Plausibility Measures: A General Approach for Representing Uncertainty</b> <i>Joseph Y. Halpern, Cornell University</i>					
9:30 - 10:00 AM	Coffee Break					
10:00 - 11:30 AM	Technical Paper Sessions:					
<b>Cognitive Robotics</b>	<b>Logic Programming</b>	<b>Natural Language Generation</b>	<b>Complexity Analysis</b>	<b>Neural Networks</b>	<b>Games</b>	
739: A Logical Account of Causal and Topological Maps <i>Emilio Remolina and Benjamin Kuipers</i>	116: A-System: Problem Solving through Abduction <i>Marc Denecker, Antonis Kakas, and Bert Van Nuffelen</i>	758: Title Generation for Machine-Translated Documents <i>Rong Jin and Alexander G. Hauptmann</i>	417: Let's go Nats: Complexity of Nested Circumscription and Abnormality Theories <i>Marco Cadoli, Thomas Eiter, and Georg Gottlob</i>	179: Knowledge Extraction and Comparison from Local Function Networks <i>Kenneth McGarry, Stefan Wermter, and John MacIntyre</i>	71: Iterative Widening <i>Tristan Cazenave</i>	
549: Online-Execution of ccGolog Plans <i>Henrik Grosskreutz and Gerhard Lakemeyer</i>	131: A Comparative Study of Logic Programs with Preference <i>Torsten Schaub and Kewen</i>	246: Dealing with Dependencies between Content Planning and Surface Realisation in a Pipeline Generation	470: A Perspective on Knowledge Compilation <i>Adnan Darwiche and Pierre Marquis</i>	258: Violation-Guided Learning for Constrained Formulations in Neural-Network Time-Series	353: Temporal Difference Learning Applied to a High Performance Game-Playing	

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

## Web site specific

Formatting

Amazon.com Book Pages

The screenshot shows the Amazon.com interface for the book 'Learning in Graphical Models' by Michael Irwin Jordan (Editor). The page features a navigation bar with categories like 'WELCOME', 'YOUR STORE', 'BOOKS', 'ELECTRONICS', 'DVD', and 'TOYS & GAMES'. A search bar is visible. The book cover is displayed with a 'LOOK INSIDE!' feature. The price is listed as \$60.00, with a 'NEW Super Saver Shipping FREE' offer. A 'Great Buy' banner is present at the bottom, suggesting a bundle with 'Probabilistic Reasoning in Intelligent Systems' for a total price of \$128.95.

## Genre specific

Layout

Resumes

The screenshot shows two resumes. The top resume is for Jason D. M. Rennie, listing his affiliation with MIT AI Lab, his contact information, and his research interests in automated data analysis. The bottom resume is for L. Douglas Baker, detailing his education at Carnegie Mellon University and the Technical University of Berlin, his professional experience at CMU, and his current dissertation research on probabilistic models for text detection.

## Wide, non-specific

Language

University Names

The screenshot shows a technical paper session schedule for 8:30 AM to 11:30 AM, listing topics like 'Invited Talk: Plausibility Measures', 'Coffee Break', and 'Technical Paper Sessions' in various categories such as Robotics, Logic Programming, Natural Language Generation, Complexity Analysis, and Neural Networks. Below the schedule is a contact card for Dr. Steven Minton, founder/CTO of Fetch, detailing his background and providing contact information for press, general information, and directions.

8:30 - 9:30 AM	<b>Invited Talk: Plausibility Measures: A General Approach for Repl</b> <i>Joseph Y. Halpern, Cornell University</i>			
9:30 - 10:00 AM	Coffee Break			
10:00 - 11:30 AM	Technical Paper Sessions:			
<b>Cognitive Robotics</b>	<b>Logic Programming</b>	<b>Natural Language Generation</b>	<b>Complexity Analysis</b>	<b>Neural Networks</b>
739: A Logical Account of Causal and Topological Maps <i>Emilio Remolina and Benjamin Kuipers</i>	116: A-System: Problem Solving through Abduction <i>Marc Denecker, Antonis Kakas, and Bert Van</i>	758: Title Generation for Machine-Translated Documents <i>Rong Jin and Alexander G. Hauptmann</i>	417: Let's go Nats: Complexity of Nested Circumscription and Abnormality Theories <i>Marco Cadoli,</i>	179: Knowledge Extraction and Comparison from Local Function Networks <i>Kenneth McGarry, Stefan Wernter, and Maclnty</i>

**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.

**Frank Huybrechts - COO**  
 Mr. Huybrechts has over 20 years of

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



# Landscape of IE Tasks (3/4): Complexity

E.g. word patterns:

## Closed set

U.S. states

He was born in Alabama...

The big Wyoming sky...

## Complex pattern

U.S. postal addresses

University of Arkansas  
P.O. Box 140  
Hope, AR 71802

Headquarters:  
1128 Main Street, 4th Floor  
Cincinnati, Ohio 45210

## Regular set

U.S. phone numbers

Phone: (413) 545-1323

The CALD main office can be reached at 412-268-1299

## Ambiguous patterns, needing context and many sources of evidence

Person names

...was among the six houses sold by Hope Feldman that year.

Pawel Opalinski, Software Engineer at WhizBang Labs.



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

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

## Single entity

*Person:* Jack Welch

*Person:* Jeffrey Immelt

*Location:* Connecticut

## Binary relationship

*Relation:* Person-Title

*Person:* Jack Welch

*Title:* CEO

*Relation:* Company-Location

*Company:* General Electric

*Location:* Connecticut

## N-ary record

*Relation:* Succession

*Company:* General Electric

*Title:* CEO

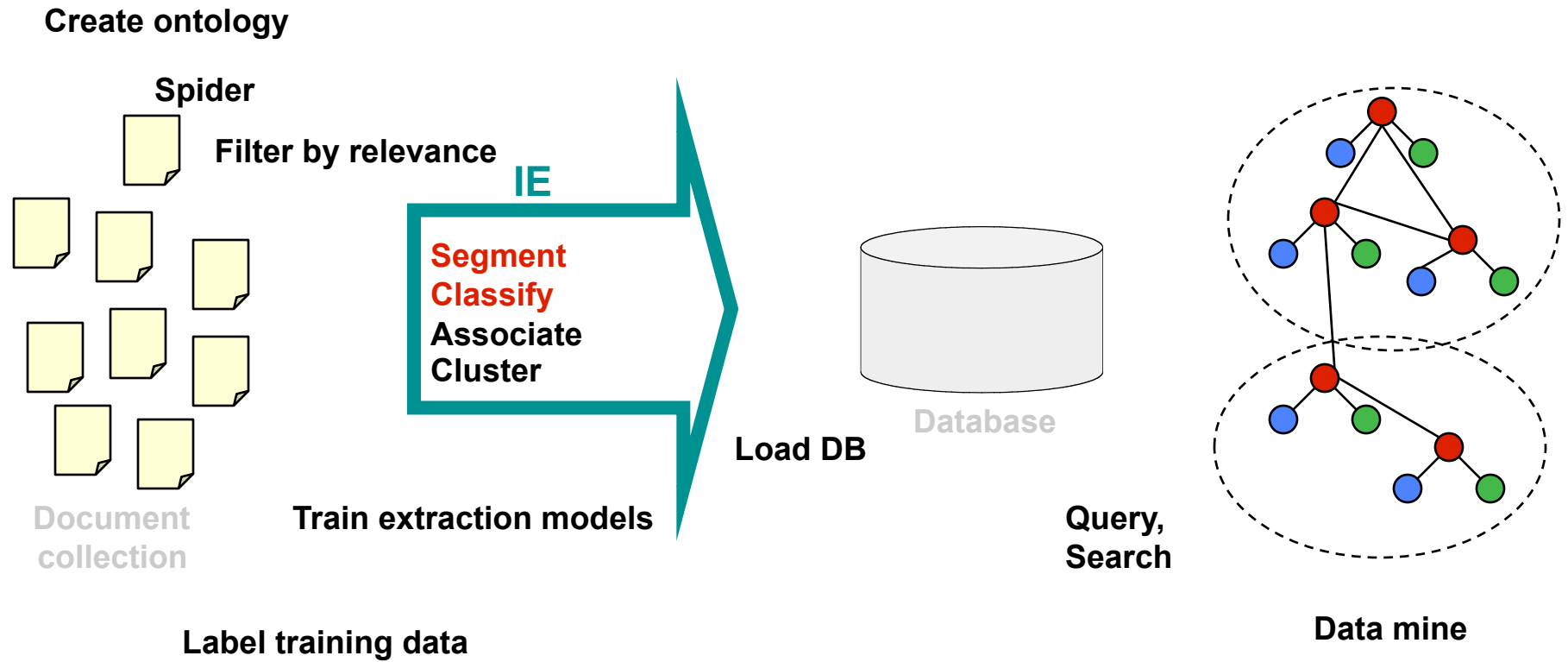
*Out:* Jack Welch

*In:* Jeffrey Immelt

*“Named entity” extraction*

# Broader View

Up to now we have been focused on segmentation and classification



## Steps 1 & 2: Hand Coded Rule Example: Conference Name

# These are subordinate patterns

```
$wordOrdinals="(?:first|second|third|fourth|fifth|sixth|seventh|eighth|ninth|tenth|eleventh|twelfth|thirteenth|fourteenth|fifteenth)";
```

```
my $numberOrdinals="(?:\d?(?:1st|2nd|3rd|1th|2th|3th|4th|5th|6th|7th|8th|9th|0th))";
```

```
my $ordinals="(?:$wordOrdinals|$numberOrdinals)";
```

```
my $confTypes="(?:Conference|Workshop|Symposium)";
```

```
my $words="(?:[A-Z]\w+\s*)"; # A word starting with a capital letter and ending with 0 or more spaces
```

```
my $confDescriptors="(?:international\s+|[A-Z]+\s+)"; # .e.g "International Conference ..." or the conference name for workshops (e.g. "VLDB Workshop ...")
```

```
my $connectors="(?:on|of)";
```

```
my $abbreviations="(?:\([A-Z]\w\w+[\W\s]*?(?:\d\d+)?\))"; # Conference abbreviations like "(SIGMOD'06)"
```

# The actual pattern we search for. A typical conference name this pattern will find is

# "3rd International Conference on Blah Blah Blah (ICBBB-05)"

```
my $fullNamePattern="((?:$ordinals\s+$words*|$confDescriptors)?$confTypes(?:\s+$connectors\s+.*?)|\s+$abbreviations?)(?:\n|\r|\s|<|>)";
```

```
#####
```

# Given a <dbworldMessage>, look for the conference pattern

```
#####
```

```
lookForPattern($dbworldMessage, $fullNamePattern);
```

```
#####
```

# In a given <file>, look for occurrences of <pattern>

# <pattern> is a regular expression

```
#####
```

```
sub lookForPattern {
```

```
    my ($file,$pattern) = @_;
```

# Machine Learning Methods

- Sequence models: HMMs, CMMs/MEMMs, CRFs
- Can work well when training data is easy to construct and is plentiful
- Can capture complex patterns that are hard to encode with hand-crafted rules
  - e.g., determine whether a review is positive or negative
  - extract long complex gene names

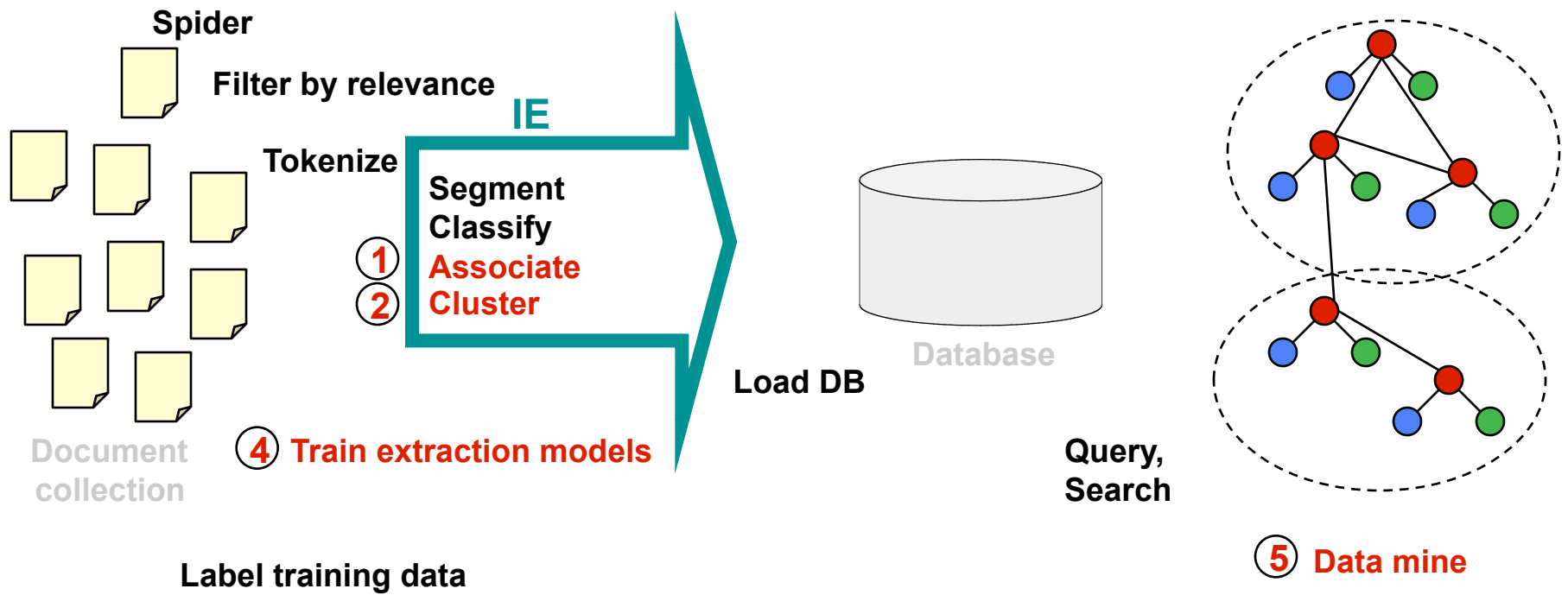
*The **human T cell leukemia lymphotropic virus type 1 Tax protein** represses MyoD-dependent transcription by inhibiting MyoD-binding to the KIX domain of p300.“*

- Can be labor intensive to construct training data
  - Question: how much training data is sufficient?

# Broader View

Now touch on some other issues

## ③ Create ontology



# Relation Extraction: Disease Outbreaks

- Extract structured relations from text

**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...

## Disease Outbreaks in *The New York Times*

<i>Date</i>	<i>Disease Name</i>	<i>Location</i>
Jan. 1995	Malaria	Ethiopia
July 1995	Mad Cow Disease	U.K.
Feb. 1995	Pneumonia	U.S.

**Information  
Extraction System  
(e.g., NYU's Proteus)**

## Example: Protein Interactions

„We show that CBF-A and CBF-C interact with each other to form a CBF-A-CBF-C complex and that CBF-B does not interact with CBF-A or CBF-C individually but that it associates with the CBF-A-CBF-C complex.“



# Relation Extraction

- Typically requires Entity Tagging as preprocessing
- Knowledge Engineering
  - Rules defined over lexical items
    - “<company> located in <location>”
  - Rules defined over parsed text
    - “((Obj <company>) (Verb located) (\*) (Subj <location>))”
  - Proteus, GATE, ...
- Machine Learning-based
  - Learn rules/patterns from examples  
Dan Roth 2005, Cardie 2006, Mooney 2005, ...
  - Partially-supervised: bootstrap from “seed” examples  
Agichtein & Gravano 2000, Etzioni et al., 2004, ...
- Recently, hybrid models [Feldman2004, 2006]



## Example Extraction Rule [NYU Proteus]

```
;;; For <company> appoints <person> <position>

(defpattern appoint
  "np-sem(C-company)? rn? sa? vg(C-appoint) np-sem(C-person) ', '?
  to-be? np(C-position) to-succeed?:
  company-at=1.attributes, sa=3.span, lv=4.span, person-at=5.attributes
  position-at=8.attributes |
  ...

(defun when-appoint (phrase-type)
  (let ((person-at (binding 'person-at))
        (company-entity (entity-bound 'company-at))
        (person-entity (essential-entity-bound 'person-at 'C-person))
        (position-entity (entity-bound 'position-at))
        (predecessor-entity (entity-bound 'predecessor-at))
        new-event)
    (not-an-antecedent position-entity)
    ;; if no company is specified for position, use agent
    ...
```

# Example Extraction Patterns: Snowball [AG2000]

<i>ORGANIZATION</i>	{<'s 0.7> <in 0.7> <headquarters 0.7>}	<i>LOCATION</i>
---------------------	---	-----------------

<i>LOCATION</i>	{<- 0.75> <based 0.75>}	<i>ORGANIZATION</i>
-----------------	----------------------------	---------------------

# (1) Association as Binary Classification

**Christos Faloutsos** conferred with **Ted Senator**, the **KDD 2003 General Chair**.

Person

Person

Role

Person-Role (**Christos Faloutsos**, **KDD 2003 General Chair**) → NO

Person-Role ( **Ted Senator**, **KDD 2003 General Chair**) → YES

Do this with SVMs and tree kernels over parse trees.

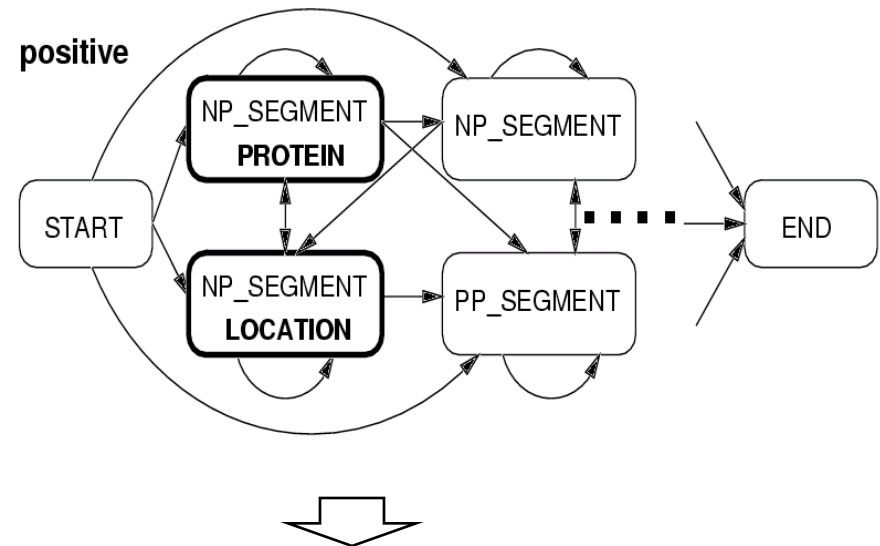
*[Zelenko et al, 2002]*

# (1) Association with Finite State Machines

[Ray & Craven, 2001]

... This enzyme, UBC6,  
localizes to the endoplasmic  
reticulum, with the catalytic  
domain facing the cytosol. ...

DET	this
N	enzyme
N	ubc6
V	localizes
PREP	to
ART	the
ADJ	endoplasmic
N	reticulum
PREP	with
ART	the
ADJ	catalytic
N	domain
V	facing
ART	the
N	cytosol

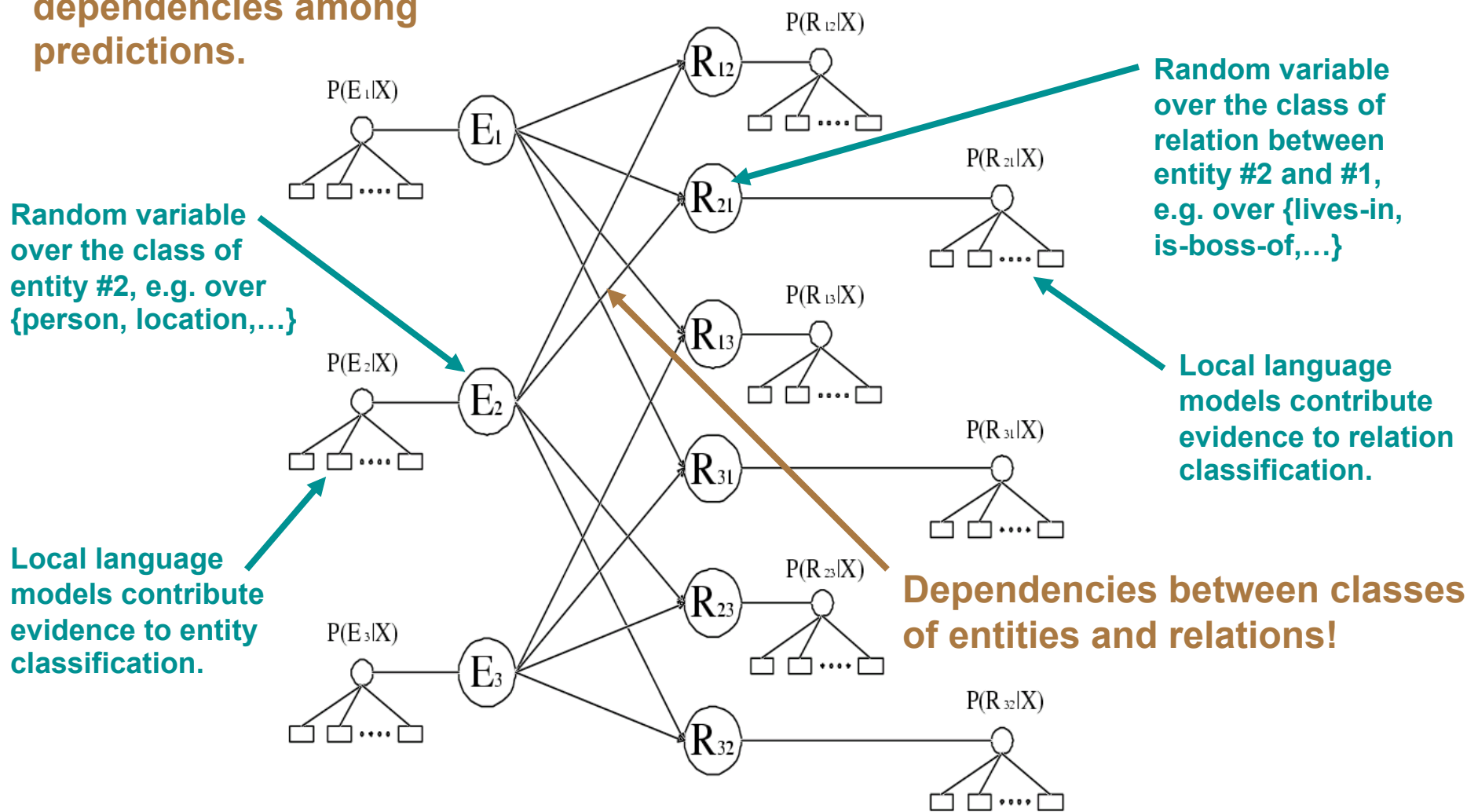


Subcellular-localization (UBC6, endoplasmic reticulum)

# (1) Association with Graphical Models

[Roth & Yih 2002]

Capture arbitrary-distance dependencies among predictions.



Random variable over the class of entity #2, e.g. over {person, location,...}

Random variable over the class of relation between entity #2 and #1, e.g. over {lives-in, is-boss-of,...}

Local language models contribute evidence to relation classification.

Local language models contribute evidence to entity classification.

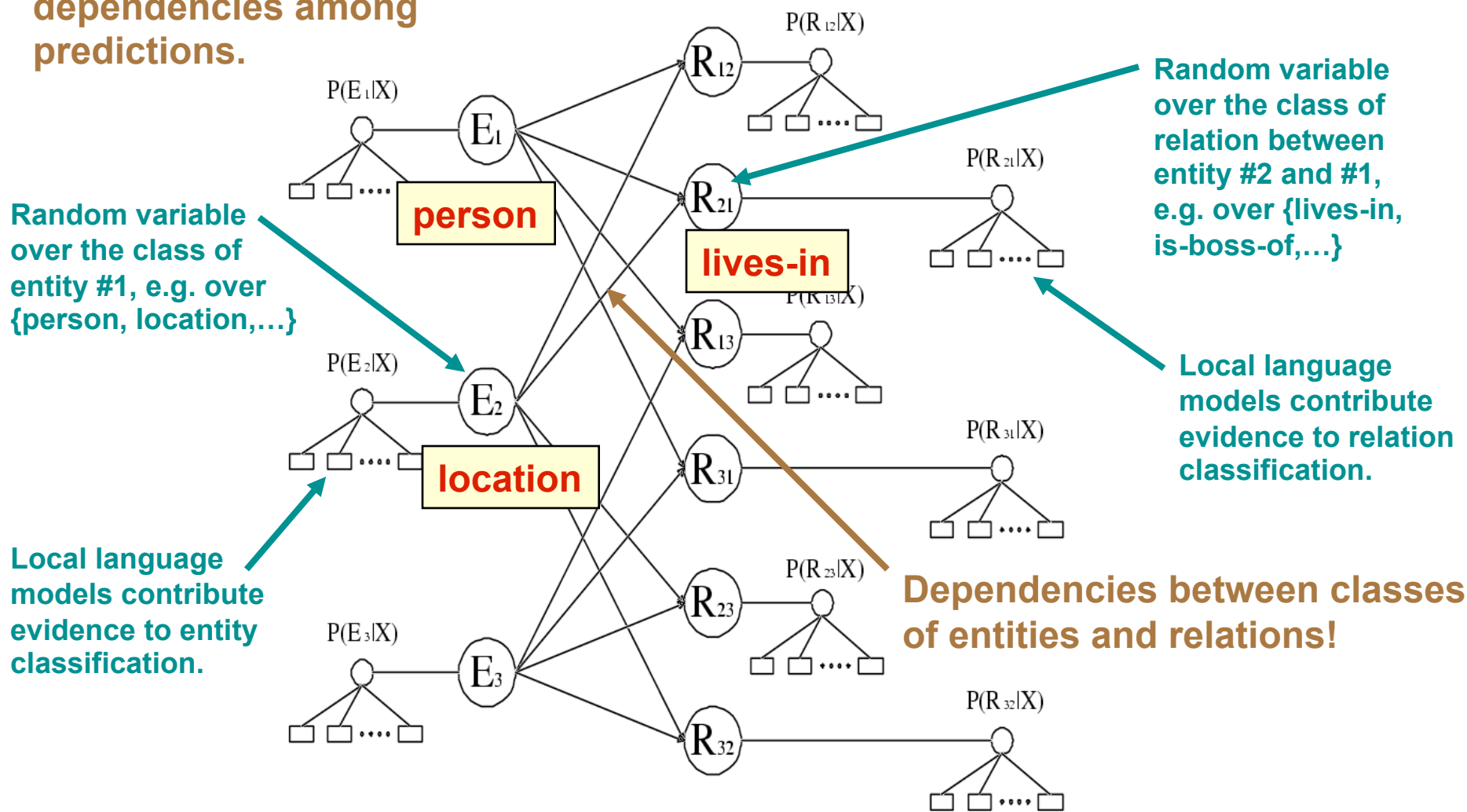
Dependencies between classes of entities and relations!

Inference with loopy belief propagation.

# (1) Association with Graphical Models

[Roth & Yih 2002]

Also capture long-distance dependencies among predictions.



Inference with loopy belief propagation.

# Accuracy of Information Extraction

Information Type	Accuracy
Entities	90-98%
Attributes	80%
Facts	60-70%
Events	50-60%

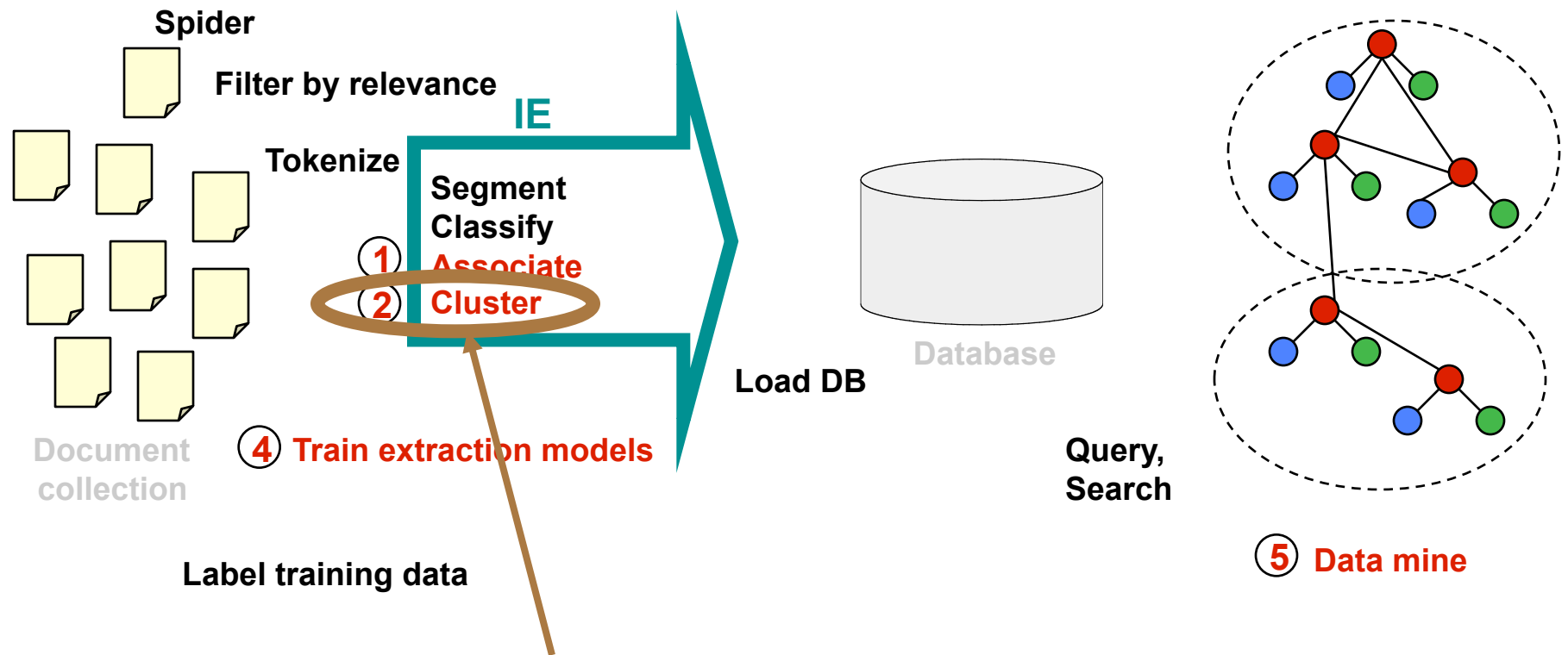
[Feldman, ICML 2006 tutorial]

- Errors cascade (error in entity tag → error in relation extraction)
- This estimate is optimistic:
  - Holds for well-established tasks
  - Many specific/novel IE tasks exhibit lower accuracy

# Broader View

Now touch on some other issues

## ③ Create ontology



When do two extracted strings refer to the same object?



# Extracted Entities: Resolving Duplicates



**Document 1:** *The Justice Department has officially ended its inquiry into the assassinations of **John F. Kennedy** and Martin Luther King Jr., finding "no persuasive evidence" to support conspiracy theories, according to department documents. The House Assassinations Committee concluded in 1978 that **Kennedy** was "probably" assassinated as the result of a conspiracy involving a second gunman, a finding that broke from the **Warren Commission**'s belief that Lee Harvey Oswald acted alone in **Dallas** on Nov. 22, 1963.*

**Document 2:** *In 1953, Massachusetts **Sen. John F. Kennedy** married Jacqueline Lee Bouvier in Newport, R.I. In 1960, Democratic presidential candidate **John F. Kennedy** confronted the issue of his Roman Catholic faith by telling a Protestant group in Houston, "I do not speak for my church on public matters, and the church does not speak for me."*

**Document 3:** ***David Kennedy** was born in Leicester, England in 1959. ...**Kennedy** co-edited *The New Poetry* (Bloodaxe Books 1993), and is the author of *New Relations: The Refashioning Of British Poetry 1980-1994* (Seren 1996).*

[From Li, Morie, & Roth, AI Magazine, 2005]

# Important Problem

- Appears in numerous real-world contexts
- Plagues many applications
  - Citeseer, DBLife, AliBaba, Rexa, etc.

## (2) Information Integration

*[Minton, Knoblock, et al 2001], [Doan, Domingos, Halevy 2001],  
[Richardson & Domingos 2003]*

Goal might be to **merge** results of two IE systems:

Name:	Introduction to Computer Science	→	Title:	Intro. to Comp. Sci.
Number:	CS 101	↘ ↗	Num:	101
Teacher:	M. A. Kludge	→	Dept:	Computer Science
Time:	9-11am		Teacher:	Dr. Klüdge
Name:	Data Structures in Java	↘	TA:	John Smith
Room:	5032 Wean Hall	→	Topic:	Java Programming
			Start time:	9:10 AM