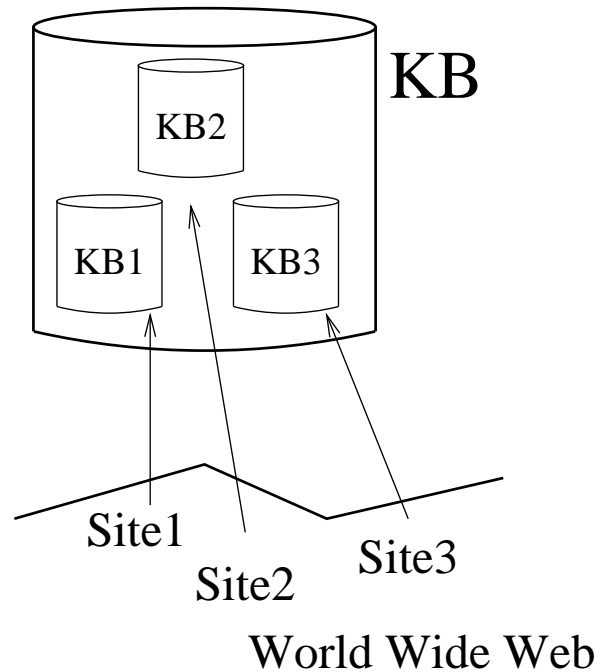


# How I Learned to Stop Worrying and Love Information Integration

William Cohen  
CALD, CMU

## What's the research problem?

We **don't know** how to **reason with** information that comes from many **different, autonomous** sources.



all mallards      duck.jpg is      duck.jpg is  
 are waterfowl    +    a picture of    =    a picture of  
                                          a mallard                                    a waterfowl

*Taxonomy*

Order	Species
waterfowl	mallard
waterfowl	bufflehead
raptor	osprey
raptor	bald eagle
...	...

*Images*

Species	File
robin	robin.jpg
mallard	duck.jpg
osprey	hawk.jpg
penguin	tweety.jpg
...	...

+

=

Order	Species	File
waterfowl	mallard	duck.jpg
raptor	osprey	hawk.jpg
...	...	...

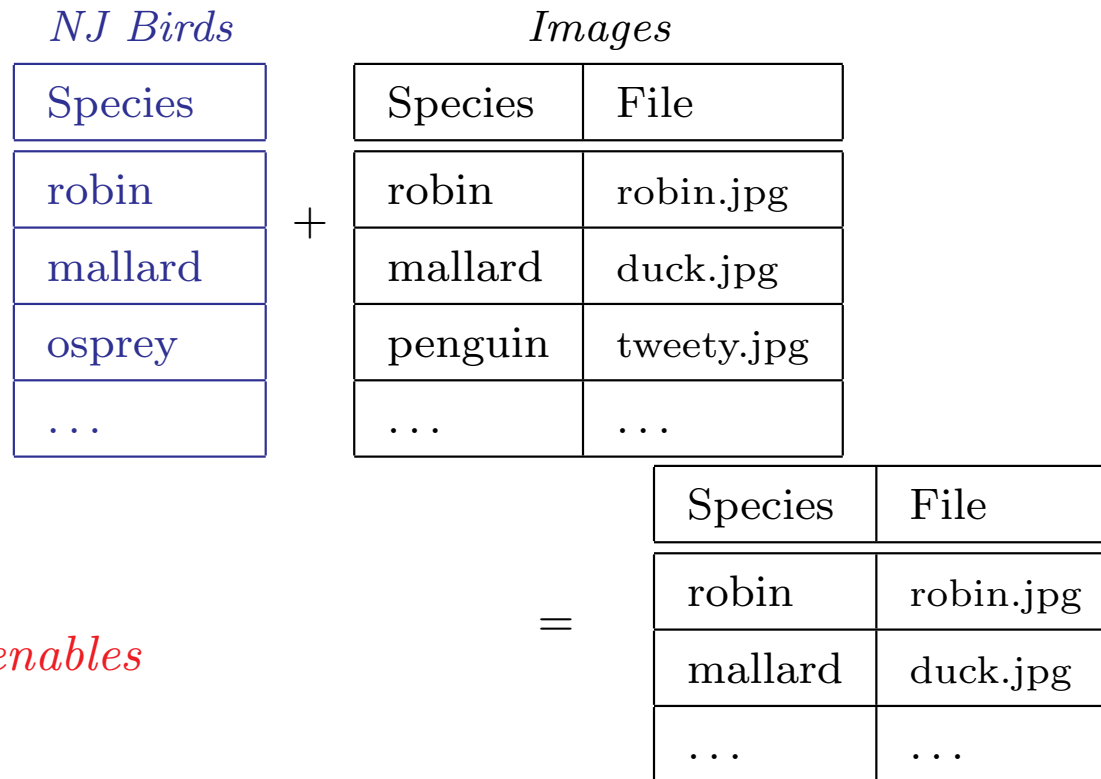
mallards are  
 found in  
 New Jersey

+

duck.jpg is  
 a picture of  
 a mallard

=

duck.jpg is a  
 picture of something  
 found in New Jersey



*Deduction enables modularity.*

## Why deduction requires co-operation

```
-? nj_bird(X),image(X,File).  
nj_bird(mallard). nj_bird(robin). ...  
image(mallard,'duck.jpg'). image(american_robin,'robin.jpg'). ...
```

The providers of the `nj_bird` and `image` facts have to agree on:

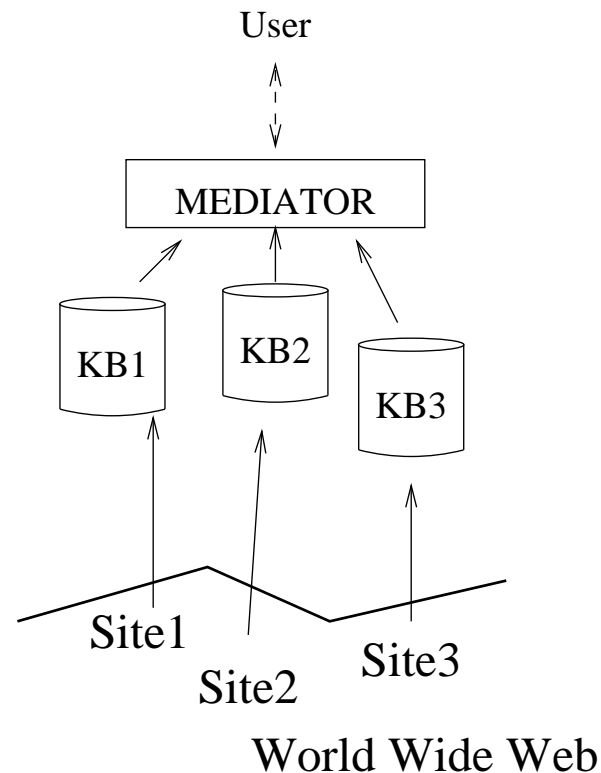
- predicate names and argument positions (schema);
- taxonomic information;
- **formal names** (OIDs) for **every entity they describe**;
- ...

## Deduction without co-operation

If information providers  
don't co-operate, then a  
“mediator” program must  
translate:

'robin' → 'american\_robin'

How hard is it to  
determine if two names  
refer to the same thing?



Humongous

Humongous  
Entertainment

Microsoft

Microsoft Kids  
Microsoft/Scholastic

Headbone

Headbone  
Interactive

The Lion King:  
Storybook

Lion King  
Animated  
StoryBook

Kestrel

American Kestrel  
Eurasian Kestrel

Disney's Activity  
Center, The  
Lion King

The Lion King  
Activity Center

Canada Goose

Goose,  
Aleutian Canada

Mallard

Mallard, Mariana

Bell Labs

AT&T Bell Labs

AT&T Research

AT&T Labs

Bell Telephone Labs

AT&T Labs—Research

AT&T Labs—Research,

Lucent Innovations

Shannon Laboratory

Bell Labs Technology

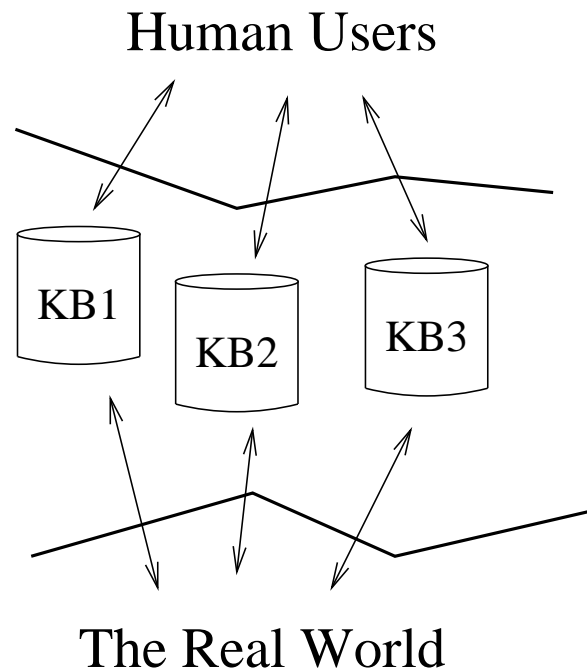
**Conclusion:** name-coreference is an AI-complete problem.



## What's the research problem?

We need a **general** means for integrating **formally unconnected** knowledge bases.

We must exploit these facts: the individual KB's model the **same real world**, and communicate with the **same users**.



## The WHIRL approach

Key points:

- Use **informal** names and descriptions as object identifiers.
- Use techniques from **information retrieval** (IR) to guess if two descriptions refer to the **same object**.
- Use **soft** ( $\approx$  probabilistic) **reasoning** for deduction.

**Formal** reasoning methods over **informally** identified objects.

## Overview of WHIRL

- WHIRL (Word-based Heterogeneous Information Representation Language) is somewhere **between** IR systems (document delivery) and KR systems (deduction).
- Outline:
  - Data model: how information is stored.
  - WHIRL query language
  - Accuracy results
  - Key ideas for implementation
  - Efficiency results
  - More results and conclusions

## Background: Information retrieval

**Ranked retrieval:** (e.g., Altavista, Infoseek, ...) given a query  $Q$ , find the documents  $d_1, \dots, d_r$  that are **most similar** to  $Q$ .

**Similarity** of  $d_i$  and  $d_j$  is measured using set of terms  $T_{ij}$  common to  $d_i$  and  $d_j$ :

$$SIM(d_i, d_j) = \sum_{t \in T_{ij}} weight(t, d_i) \cdot weight(t, d_j)$$

- A **term** is a single word (modulo stemming, ...)
- Heuristic: make  $weight(t, d)$  large if  $t$  is frequent in  $d$ , or if  $t$  is rare in the corpus of which  $d$  is an element.

## Background: Information retrieval

Similarity of  $d_i$  and  $d_j$  is measured using set of terms  $T_{ij}$  common to  $d_i$  and  $d_j$ :

$$SIM(d_i, d_j) = \sum_{t \in T_{ij}} weight(t, d_i) \cdot weight(t, d_j)$$

- Heuristic: make  $weight(t, d)$  large if  $t$  is frequent in  $d$  (TF), or if  $t$  is rare in the corpus of which  $d$  is an element (IDF).
- Example: if the corpus is a list of company names:
  - Low weight: “Inc”, “Corp”, ...
  - High weight: “Microsoft”, “Lucent”, ...
  - Medium weight: “Acme”, “American”, ...

## Background: Information retrieval

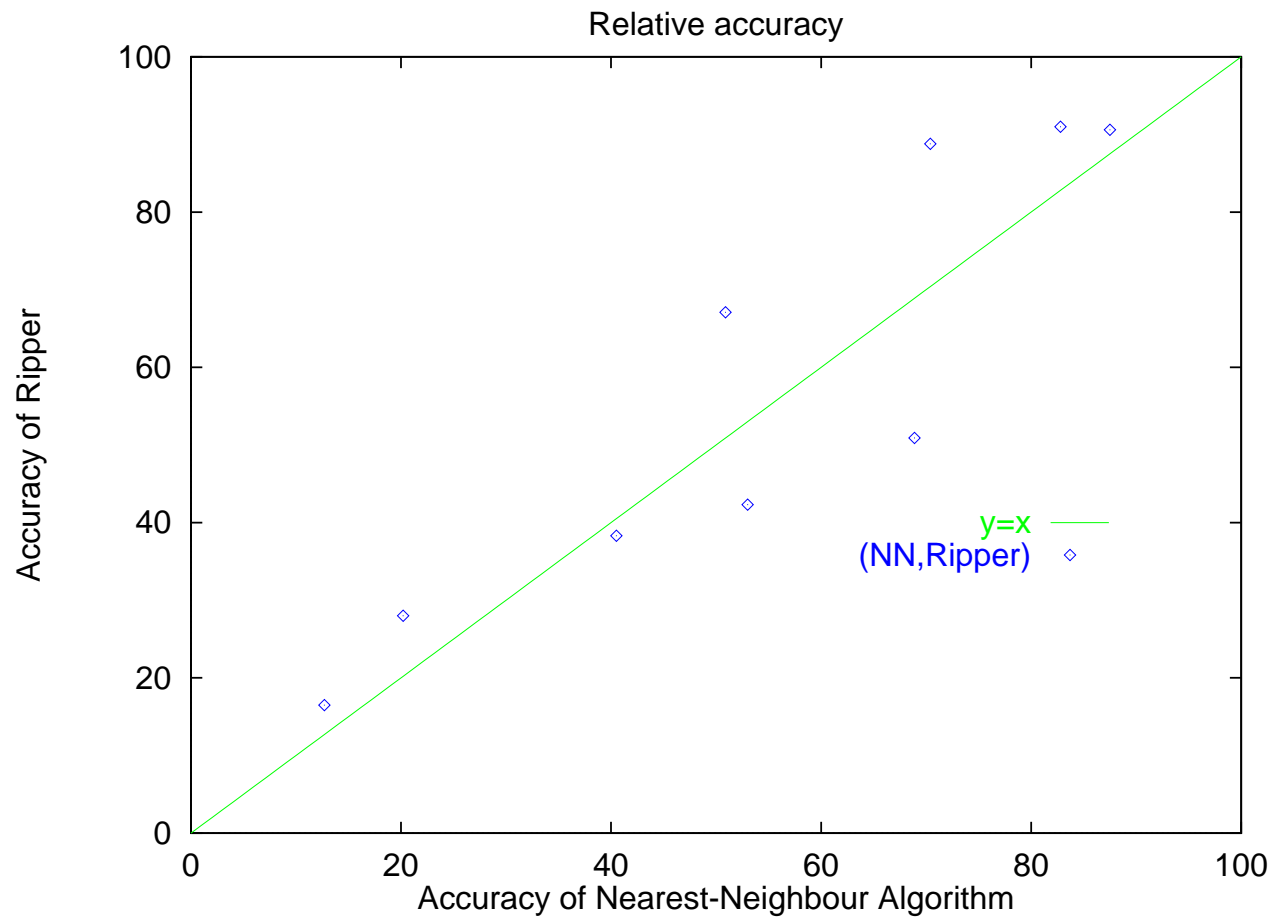
It's notationally convenient to think of a document  $d_i$  as a long, sparse **vector**,  $v_i$ .

If  $\vec{v}_i = \langle v_{i,1}, \dots, v_{i,|T|} \rangle$ ,  $v_{i,t} = \text{weight}(t, d_i)$ , and  $\|v_i\| = 1$ :

$$\begin{aligned} \text{SIM}(d_i, d_j) &= \sum_{t \in T} \text{weight}(t, d_i) \cdot \text{weight}(t, d_j) \\ &= \vec{v}_i \cdot \vec{v}_j \end{aligned}$$

Also,  $0 \leq \text{SIM}(d_i, d_j) \leq 1$ .

## Effectiveness of the TF-IDF “vector space” representation



Cinema	Movie	Show Times
Roberts Theaters Chatham	Brassed Off	7:15 - 9:10
Berkeley Cinema	Hercules	4:15 - 7:30
Sony Mountainside Theater	Men In Black	7:40 - 8:40 - 9:30 - 10:10

listing( $\vec{v}_{RTC}, \vec{v}_{BO}, \vec{v}_{T79}$ ), 1.  
listing( $\vec{v}_{BC}, \vec{v}_H, \vec{v}_{T47}$ ), 1.  
listing( $\vec{v}_{SMT}, \vec{v}_{MIB}, \vec{v}_{T789}$ ), 1.

review( $\vec{w}_{MIB97}, \vec{w}_{R1}$ ), 1.  
review( $\vec{w}_{FO}, \vec{w}_{R2}$ ), 1.  
review( $\vec{w}_{SB}, \vec{w}_{R3}$ ), 1.

Each  $\vec{v}_i, \vec{w}_i$  is a **document vector**.  
Each fact has a **score**  $s \in [0, 1]$ .

Movie	Review
Men in Black, 1997	(* * *) One of the biggest hits of ...
Face/Off, 1997	(* * $\frac{1}{2}$ ) After a slow start, ...
Space Balls, 1987	(* $\frac{1}{2}$ ) Not one of Mel Brooks' best efforts, this spoof ...

$$\vec{v}_{MIB} = \langle \dots, v_{black}, \dots, v_{in}, \dots, v_{men}, \dots \rangle$$

$$\vec{w}_{MIB97} = \langle \dots, w_{black}, \dots, w_{in}, \dots, w_{men}, \dots, w_{1997}, \dots \rangle$$

$$w_{1997} \approx 0 \implies sim(\vec{v}_{MIB}, \vec{w}_{MIB97}) \approx 1$$



## Queries in WHIRL

- **Syntax:** WHIRL = (similarity)  
Prolog – function symbols – recursion – negation +  $X \sim Y$
- **Semantics** (details in Cohen, SIGMOD98):
  - A ground formula gets a **score**  $s \in [0, 1]$
  - $\text{Score}(p(a_1, \dots, a_k)) = s$  for DB literals.
  - $\text{Score}(a \sim b) = \text{SIM}(a, b)$  for similarity literals.
  - $\text{Score}(\phi \wedge \psi) = \text{Score}(\phi) \cdot \text{Score}(\psi)$ .
  - $\text{Score}(\phi \vee \psi) = 1 - (1 - \text{Score}(\phi))(1 - \text{Score}(\psi))$
  - Answer to a query  $Q$  is an **ordered** list of the  $r$  substitutions  $\theta_1, \dots, \theta_r$  that give  $Q\theta_i$  the **highest scores**.  
(User provides  $r$ ).

## Sample WHIRL queries

Standard ranked retrieval:

*“find reviews of sci-fi comedies”.*

?- review(Title,Rev)  $\wedge$  Rev $\sim$ “sci-fi comedy”

(score 0.22):  $\theta_1 = \{\text{Title}/\vec{w}_{MIB97}, \text{Rev}/\vec{w}_{R1}\}$

(score 0.19):  $\theta_2 = \{\text{Title}/\vec{w}_{SB}, \text{Rev}/\vec{w}_{R4}\}$

(score 0.13):  $\theta_2 = \dots$

## Sample WHIRL queries

Standard DB queries: “*find reviews for movies playing in Mountainside*” (assume single-term “movie IDs” in DB)

?- review(Id1,T1,Rev)  $\wedge$  listing(C,Id2,T2,Time)  
 $\wedge$  Id1 $\sim$ Id2  $\wedge$  C $\sim$ “Sony Mountainside Theater”

(score 1.00):  $\theta_1 = \{\text{Id1}/\vec{v}_{93}, \text{Id2}/\vec{w}_{93}, \text{Rev}/\vec{w}_{R1}, \dots\}$

(score 1.00):  $\theta_2 = \dots$

Cinema	Id	Movie	Time
...	21	Brassed Off	...
Sony ...	93	Men In Black	...

Id	Movie	Review
93	Men in Black, 1997	...
44	Face/Off, 1997	...

## Sample WHIRL queries

Mixed queries: “*where is [Men in Black] playing?*”

?- review(Id1,T1,Rev)  $\wedge$  listing(C,Id2,T2,Time)  
 $\wedge$  Id1 $\sim$ Id2  $\wedge$  Rev $\sim$ “sci-fi comedy with Will Smith”

(score 0.22):  $\theta_1 = \{Id1/\vec{v}_{93}, Id2/\vec{w}_{93}, Rev/\vec{w}_{R1}, \dots\}$

(score 0.13):  $\theta_2 = \dots$

Cinema	Id	Movie	Time
...	21	Brassed Off	...
Sony ...	93	Men In Black	...

Id	Movie	Review
93	Men in Black, 1997	...
44	Face/Off, 1997	...

## A realistic situation

Cinema	Movie	Show Times
Roberts Theaters Chatham	Brassed Off	7:15 - 9:10
Berkeley Cinema	Hercules	4:15 - 7:30
Sony Mountainside Theater	Men In Black	7:40 - 8:40 - 9:30 - 10:10

With real Web data, there will be no common **ID** fields, only **informal names**.

Movie	Review
Men in Black, 1997	(* * *) One of the biggest hits of ...
Face/Off, 1997	(* * $\frac{1}{2}$ ) After a slow start, ...
Space Balls, 1987	(* $\frac{1}{2}$ ) Not one of Mel Brooks' best efforts, this spoof ...

## Sample WHIRL queries

“Similarity” joins: *“find reviews of movies currently playing”*

?- review(Title1,Rev)  $\wedge$  listing(-,Title2,Time)  $\wedge$  Title1~Title2

(score 0.97):  $\theta_1 = \{ \text{Title1}/\vec{v}_{MIB}, \text{Title2}/\vec{w}_{MIB97}, \dots \}$   
(Men in Black) (Men in Black, 1997)

...

(score 0.41):  $\theta_2 = \{ \text{Title1}/\vec{v}_{BO}, \text{Title2}/\vec{w}_{FO}, \dots \}$   
(Brassed Off) (Face/Off)

...

## How well do similarity joins work?

?- top500(X), hiTech(Y), X~Y

*top500:*

Abbott Laboratories  
Able Telcom Holding Corp.  
Access Health, Inc.  
Acclaim Entertainment, Inc.  
Ace Hardware Corporation  
ACS Communications, Inc.  
ACT Manufacturing, Inc.  
Active Voice Corporation  
Adams Media Corporation  
Adolph Coors Company  
...

*hiTech:*

ACC CORP  
ADC TELECOMMUNICATION INC  
ADELPHIA COMMUNICATIONS CORP  
ADT LTD  
ADTRAN INC  
AIRTOUCH COMMUNICATIONS  
AMATI COMMUNICATIONS CORP  
AMERITECH CORP  
APERTUS TECHNOLOGIES INC  
APPLIED DIGITAL ACCESS INC  
APPLIED INNOVATION INC  
...

Sample company-name pairings

*WHIRL output on business.html*



## Evaluating similarity joins

- **Input:** query
- **Output:** ordered list of documents

1	✓	$a_1$	$b_1$
2	✓	$a_2$	$b_2$
3	✗	$a_3$	$b_3$
4	✓	$a_4$	$b_4$
5	✓	$a_5$	$b_5$
6	✓	$a_6$	$b_6$
7	✗	$a_7$	$b_7$
8	✓	$a_8$	$b_8$
9	✓	$a_9$	$b_9$
10	✗	$a_{10}$	$b_{10}$
11	✗	$a_{11}$	$b_{11}$
12	✓	$a_{12}$	$b_{12}$

Precision at  $K$ :  $G_K/K$

Recall at  $K$ :  $G_K/G$

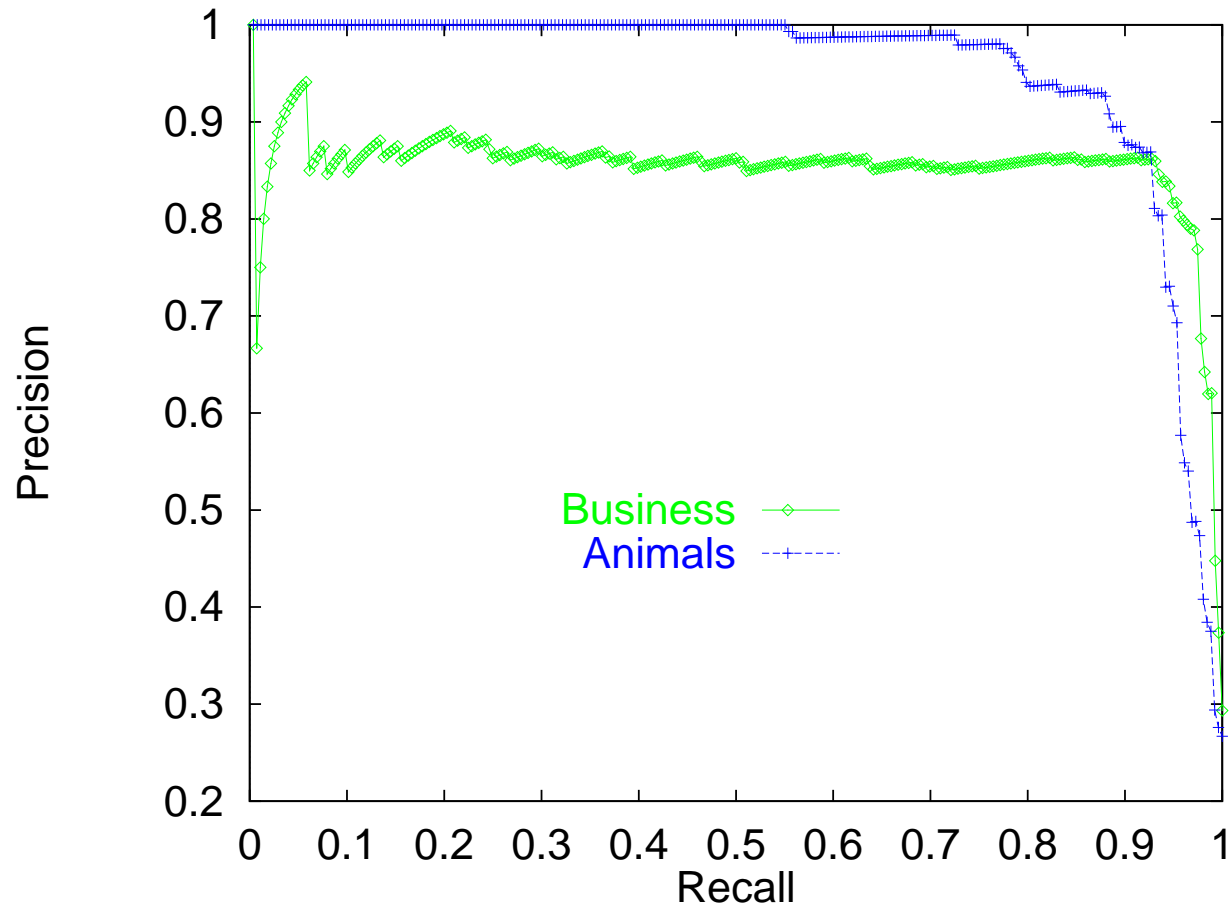
$G$ : # good pairings

$G_K$ : # good pairings in first  $K$

## Evaluating similarity joins

- Pick relations  $p, q$  with  $> 2$  plausible keys
- Perform “similarity join” using first key field
- Mark a pairing correct (“relevant”) if secondary key matches
- Compute precision and recall over first 1000 rankings
- Examples
  - Business: company name, web site
  - Animals: common name, scientific name
  - etc

## Evaluating similarity joins



## Evaluating WHIRL queries

Additional experiments:

- Repeat with **more datasets** from more domains.
  - Average precision ( $\approx$  area under precision-recall curve) ranges from 85% to 100% over 13 joins in 6 domains.
- Repeat for **more complex** join queries.
  - Average precision drops from 94% for 2-way joins to 90% for 5-way joins (averaged over many queries in one domain).
- Evaluate other things to do with WHIRL.
- **How can you implement WHIRL efficiently?**

## An efficient implementation

Key ideas for current implementation:

- Central problem: given  $Q$ , find **best** substitution.
  - Currently, using  **$A^*$  search**.
- Search space: partial substitutions.  
e.g., for “?- r(X),s(Y),X~Y”, possible state is  $\{X = \vec{x}\}$ .
- Key operator: when  $Q$  contains “ $\vec{x} \sim Y$ ”, find good candidate bindings for  $Y$  quickly.
  - Use **inverted indices**.

## An efficient implementation

- Key step: state is a substitution  $\theta$ ,  $Q\theta$  contains “ $s(Y), \vec{x} \sim Y$ ”.  
Need to find good candidate bindings for  $Y$  quickly.
  1. Pick some term  $t$  with large weight in  $\vec{x}$ .
  2. Use inverted index to get

$$I_{t,s,1} = \{\vec{y} : s(\vec{y}) \in \text{DB and } y_t > 0\}$$

- To compute heuristic value of state, use fact that

$$\text{score}(\vec{x} \sim Y) \leq \max_{\vec{z} \in I_{t,s,1}} \left( \sum_t x_t \cdot z_t \right) \leq \sum_t x_t \cdot \left( \max_{\vec{z} \in I_{t,s,1}} z_t \right)$$

- Indexing and bounds well-known in IR  
(Buckley-Lewitt, Turtle-Flood’s *maxscore* alg)

## An efficient implementation

- **For instance:** I used WHIRL as the DBMS for two real-life integration systems:
  - Birds of North America:  $\approx 35$  sites
  - Computer Games for Kids:  $\approx 15$  sites
- Both were made available on the Web, and queries were logged.

## Results on real-world queries

	Domain	
	Games	Birds
# sites indexed	15	34
# facts stored in DB	23,435	143,666
# queries in sample	100	91
avg time/query (sec)	0.3	0.2
max time/query (sec)	5.2	5.4



Domain	$k$	# $k$ -way Joins	Avg #Sim Literals	Average Time
Birds	$\leq 2$	47	2.0	0.02
	3	22	3.3	0.03
	4	14	3.8	0.35
	5	4	3.8	1.90
	6	4	5.0	0.22
Games	$\leq 2$	35	1.4	0.06
	3	20	3.9	0.08
	4	16	4.1	0.50
	5	23	5.3	0.26
	6	6	6.0	1.61

## The extraction problem

Sometimes it's difficult to extract even an **informal** name from its context:

- Fox Interactive has a fully working demo version of the Simpsons Cartoon Studio. (Win and Mac)
- Vividus Software has a free 30 day demo of Web Workshop (web authoring package for kids!) Win 95 and Mac
- Scarlet Tanager (58kB) *Piranga olivacea*. New Paltz, June 1997.  
“...Robin-like but hoarse (suggesting a Robin with a sore throat).”  
(Peterson) “..a double-tone which can only be imitated by strongly humming and whistling at the same time.” (Mathews)

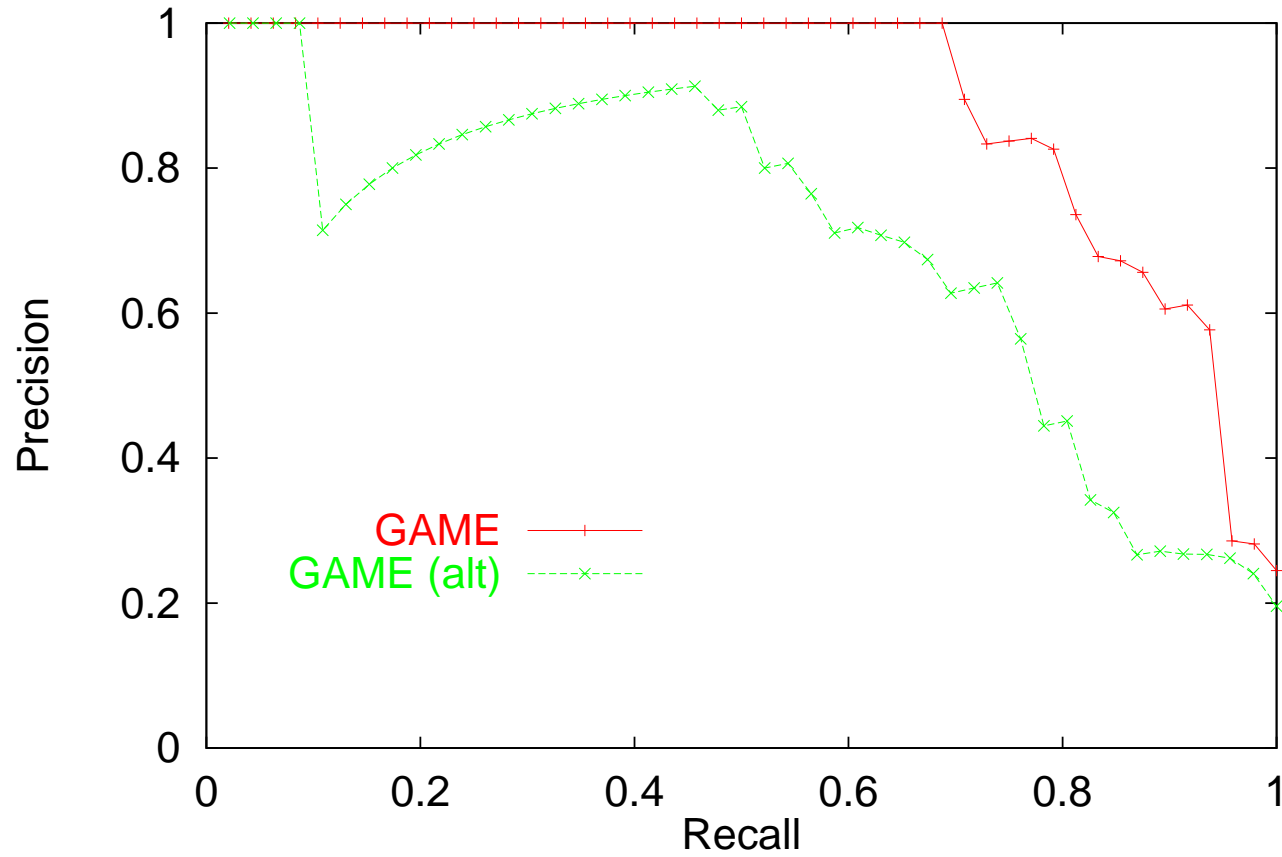
## The extraction problem

Idea: use text **without** trying to extract names.

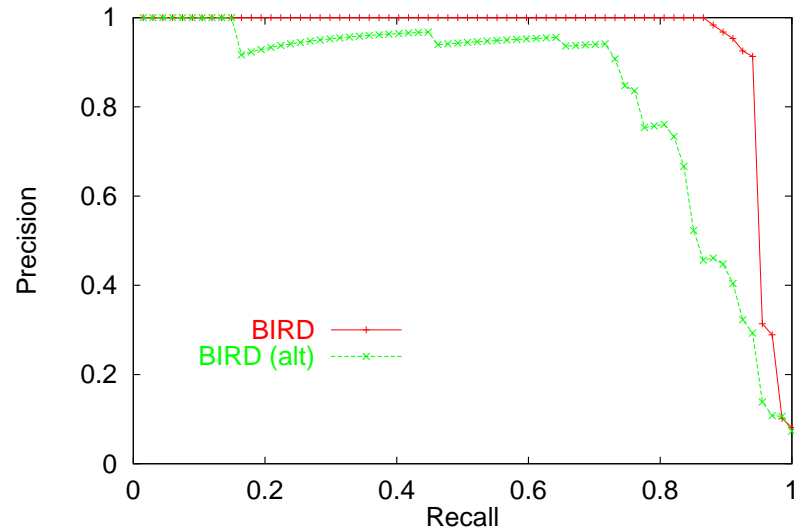
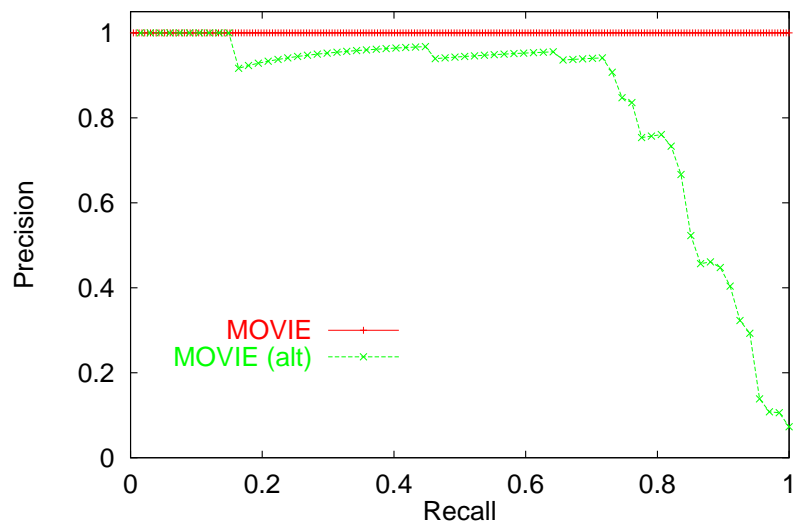
?- paragraph(X),name(Y),X~Y

<b>80.26</b>	Ubi Software has a demo of Amazing Learning Games with Rayman.	Amazing Learning Games with Rayman	✓
<b>78.25</b>	Interplay has a demo of Mario Teaches Typing. (PC)	Mario Teaches Typing	✓
<b>75.91</b>	Warner Active has a small interactive demo for Where's Waldo at the Circus and Where's Waldo? Exploring Geography (Mac and Win)	Where's Waldo? Exploring Geography	✓
<b>74.94</b>	MacPlay has demos of Marios Game Gallery and Mario Teaches Typing. (Mac)	Mario Teaches Typing	✓
<b>71.56</b>	Interplay has a demo of Mario Teaches Typing. (PC)	Mario Teaches Typing 2	✗

## Deduction without extraction



## Deduction without extraction



Movie 1: full review (no extraction).

Movie 2: movie name, cinema name & address, showtimes.

## More uses of WHIRL: Classification?

review(“Putt-Putt Travels Through Time”, url1).  
category(“Putt-Putt’s Fun Pack”, “adventure”).  
category(“Time Traveler CD”, “history”).

...

*“find me reviews of adventure games”*

$v(\text{Url}) \leftarrow$

$\text{review}(\text{Game1}, \text{Url}) \wedge \text{category}(\text{Game2}, \text{Cat})$

$\wedge \text{Game1} \sim \text{Game2} \wedge \text{Cat} \sim \text{“adventure”}$

To answer this query, WHIRL **guesses** the class “adventure” based on **similarities** between names.

## More uses of WHIRL: Classification

$$\text{category}(\text{Cat}) \leftarrow \text{test}(\text{X}) \wedge \text{train}(\text{Y}, \text{Cat}) \wedge \text{X} \sim \text{Y}$$

- Here **train** contains a **single** unclassified example, and **test** contains a **set** of training examples with known **categories**.  
(from Cohen&Hirsh, KDD-98)
- WHIRL here performs a sort of  $K$ -NN classification.
  1. Find  $r$  best bindings for **X,Y,Cat**
  2. Combine evidence using noisy-or:  
$$\text{Score}(\phi \wedge \psi) = \text{Score}(\phi) \cdot \text{Score}(\psi)$$

## Using WHIRL for Classification

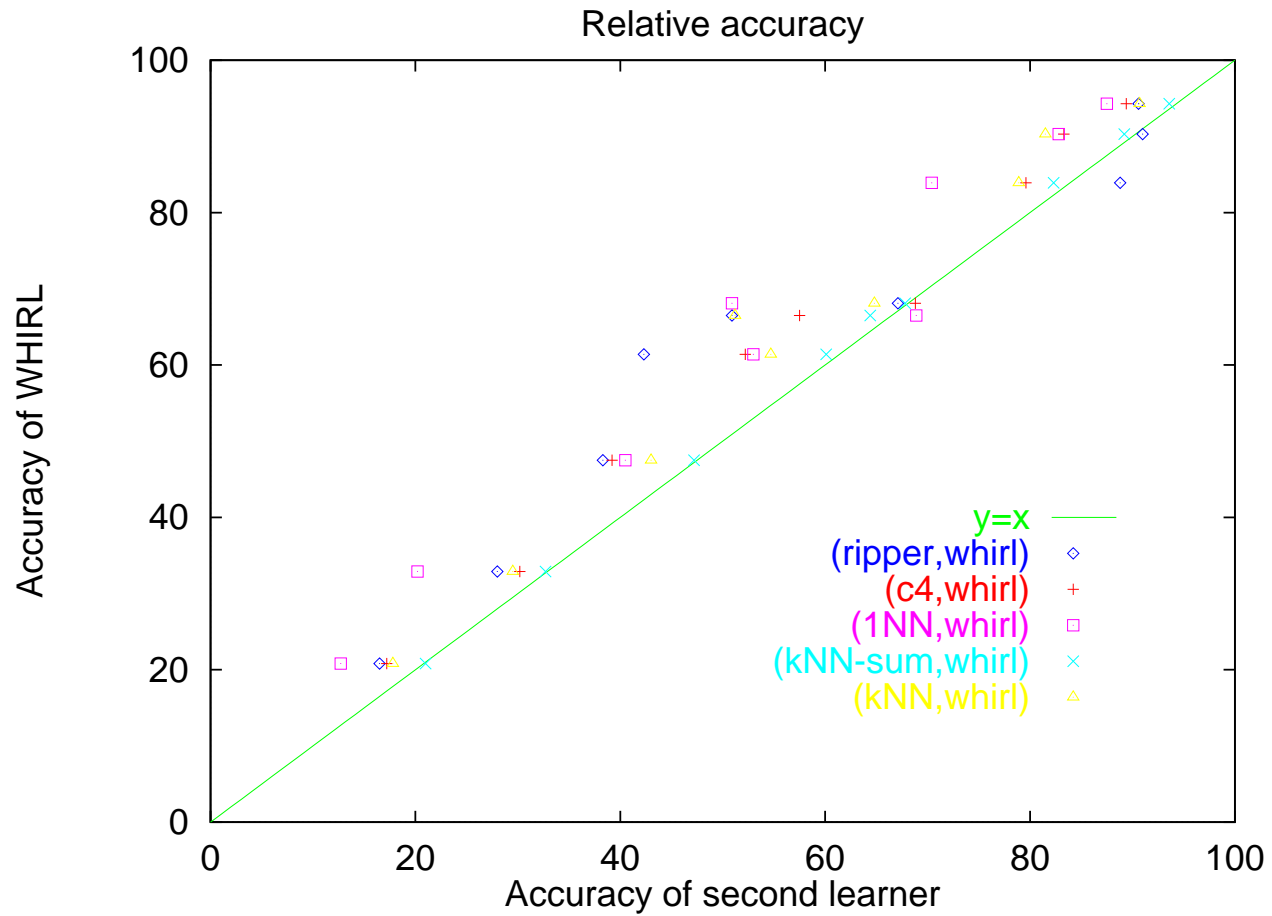
- Created nine representative datasets using data from Web.
- All instances were short “names”
  - *book title*: inst=“The Humbugs of the World by P. T. Barnum (page images at MOA)”, class=“General Works”
  - *company name*: inst=“National City Corporation”, class=“Banks–Midwest”
  - Also bird names, Web page titles, ...
- # classes ranged from 6 to 228, #instances ranged from  $\approx 300$  to  $\approx 3000$ .



## Benchmark classification problems

problem	#train/ #test	#classes/ #terms	text-valued field/label
memos	334/10cv	11/1014	document title/category
cdroms	798/10cv	6/1133	CDRom game name/category
birdcom	914/10cv	22/674	common name of bird/phylogenic order
birdsci	914/10cv	22/1738	common+sci name/phylogenic order
hcoarse	1875/600	126/2098	company name/industry (coarse grain)
hfine	1875/600	228/2098	company name/industry (fine grain)
books	3501/1800	63/7019	book title/subject heading
species	3119/1600	6/7231	animal name/phylum
netvet	3596/2000	14/5460	URL title/category

## Using WHIRL for Classification



## Using WHIRL for Classification

Later work by Zelikovitz & Hirsh:

- Slightly more complex WHIRL queries (2-way chain join)
- Linked **test** and **train** documents via a set of “similar” unlabeled documents
- Showed **improved** classification performance for **short** examples or **small** training sets.

## Classification with “side information”

Consider classification...

- **Observation:** Performance can often be improved by obtaining additional features about the entities involved.
- **Question:** Can performance be improved using weaker “side information”—like additional features that **might or might not** be about the entities involved in the classification task?

Instance		Label
Itzak Perlman	BMG	classic
Billy Joel	RCA	pop
Metallica	...	pop
...	...	...

**Goal:** from the data above, learn to classify musical artists as classical *vs.* popular.

**Basic ideas:** introduce **new features** for artist names that

- appear in certain lists or tables; (e.g., italicized names in the ‘Guest Artist’ page)
- are modified by certain words (e.g., “KØØL”)

### Guest Artists: Spring 2000

- Apr 9, *Itzak Perlman*
- May 3, *Yo Yo Ma*
- May 17, *The Guanari Quartet*
- ...

### Biff’s KØØL Band Links

- Nine Inch Nails (new!)
- Metallica!! Rockin’! Anyone know where can I find some MP3s?

• ...

...

## The extraction algorithm

1. Parse the **HTML** markup
2. Associate each short marked-up section with its “**tag-path position**”  $(x_1, p_1), (x_2, p_2), \dots$
3. Find all triples  $(a_j, x_i, p_i)$  such that instance  $a_j$ 's name is **highly similar** to  $x$  (with a WHIRL query.)
4. Define  $g_p(a) = 1$  iff  $\exists x : (a, x, p)$  is a triple.
5. Determine the “**scope**” of each HTML header (e.g., **h1, h2, ...**)
6. Define  $g_w(a) = 1$  iff  $\exists x, h : (a, x, p)$  is a triple,  $h$  is a header,  $x$  is in the scope of  $h$ , and  $w$  is a word  $h$ .

## Feature construction: an example

```
<html><head>Biff's Home Page</head>
<body>
<h2>KØØL Band Links</h2>
<table> <tr>
  <td>Metallica
  <td>Nine Inch Nails (new!)
</tr><tr>
  <td>Barry Manilow
  ...
```

```
html(head(...),
      body(
        h2(KØØL Band Links),
        table(
          tr(td(Metallica),
             td(Nine Inch Nails (new!))),
          tr(td(Barry Manilow),
             ...
            )
        )
      )
    )
```

*Instances:*

...
Metallica
Nine Inch Nails
Itzak Perlman
...

(“KØØL Band Links”, [www.biff.com/html\\_body\\_h1](http://www.biff.com/html_body_h1))

(“Metallica”, [www.biff.com/html\\_body\\_table\\_tr\\_td](http://www.biff.com/html_body_table_tr_td))

(“Nine Inch Nails (new!)”, [www.biff.com/html\\_body\\_table\\_tr\\_td](http://www.biff.com/html_body_table_tr_td))

(“Barry Manilow”, [www.biff.com/html\\_body\\_table\\_tr\\_td](http://www.biff.com/html_body_table_tr_td))



```
html(head(...),
      body(
        h2(KØØL Band Links),
        table(
          tr(td(Metallica),
             td(Nine Inch Nails (new!))),
          tr(td(Barry Manilow),
             ...
```

*(instance-name, instance-mention, position)*

(“Metallica”, “Metallica”, table\_tr\_td)

(“Nine Inch Nails”, “Nine Inch Nails (new!)”, table\_tr\_td)

(“Barry Manilow”, “Barry Manilow”, table\_tr\_td)

```

html(head(...),
      body(
        h2(KØØL Band Links),
        | table(
        |   tr(td(Metallica),
        |     td(Nine Inch Nails (new!))),
        |   tr(td(Barry Manilow),
        |     h1(...),
        |     ...

```

$g_{\text{table\_tr\_td}}(\text{"Metallica"}) = 1$

$g_{\text{KØØL}}(\text{"Metallica"}) = 1$

$g_{\text{table\_tr\_td}}(\text{"Nine Inch Nails"}) = 1$

$g_{\text{band}}(\text{"Metallica"}) = 1$

$g_{\text{table\_tr\_td}}(\text{"Barry Manilow"}) = 1$

$g_{\text{links}}(\text{"Metallica"}) = 1$

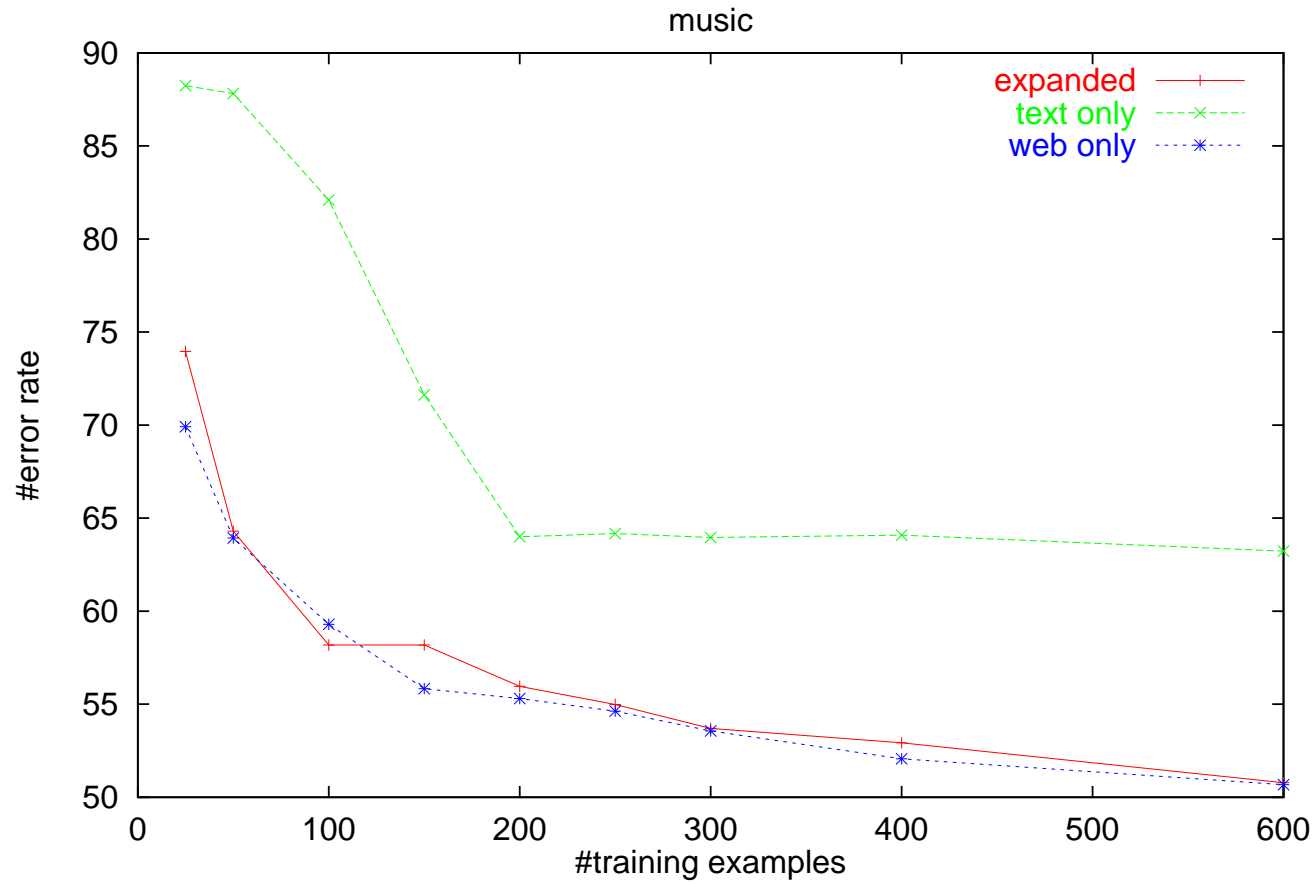
...

## Benchmark problems

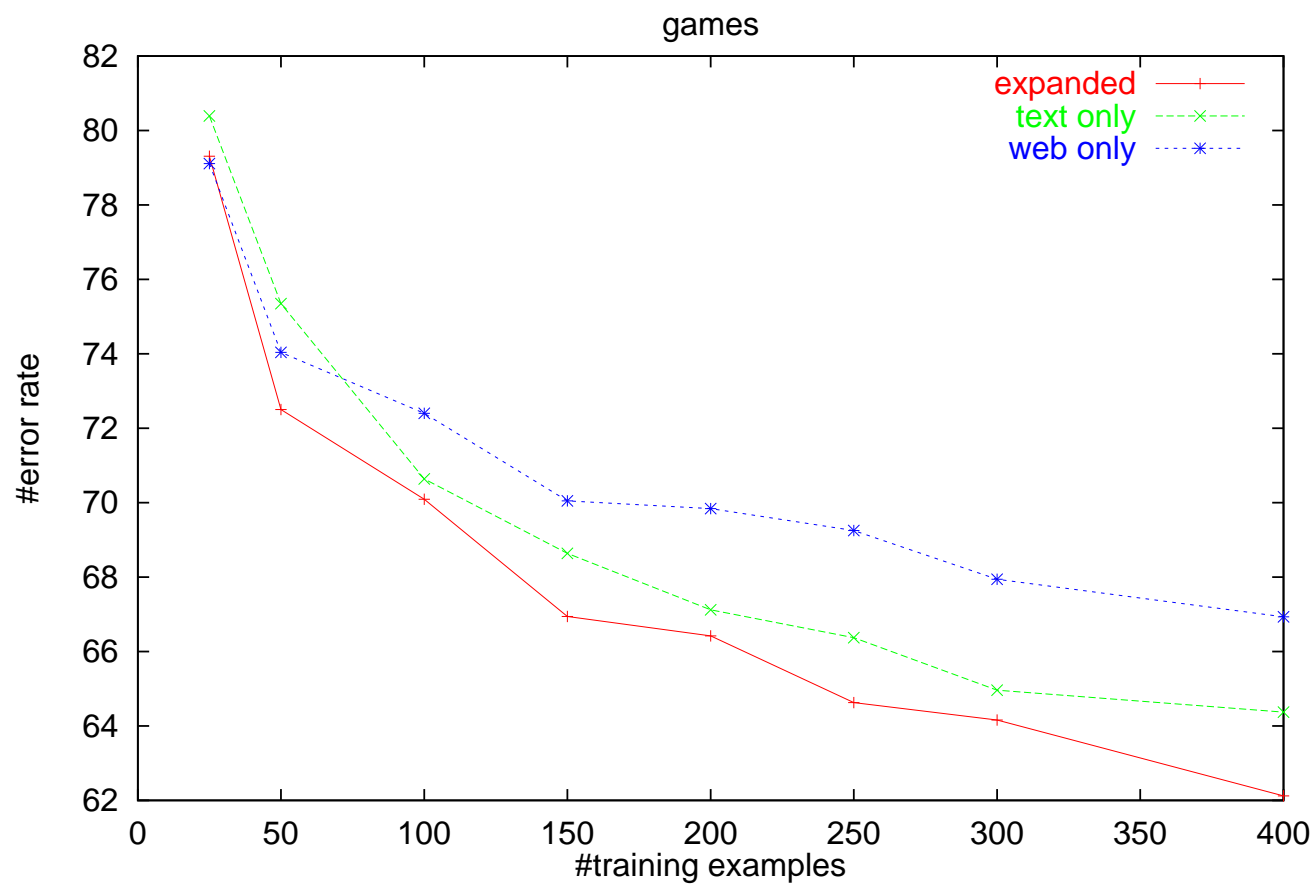
	#example	#class	#terms	#pages	#features added
music	1010	20	1600	217	1890
games	791	6	1133	177	1169
birdcom	915	22	674	83	918
birdsci	915	22	1738	83	533

- original data: names as bag-of-words
- music: (Cohen&Fan,WWW00) others: (Cohen&Hirsh,KDD98)
- note: test data must be processed as well (transduction).

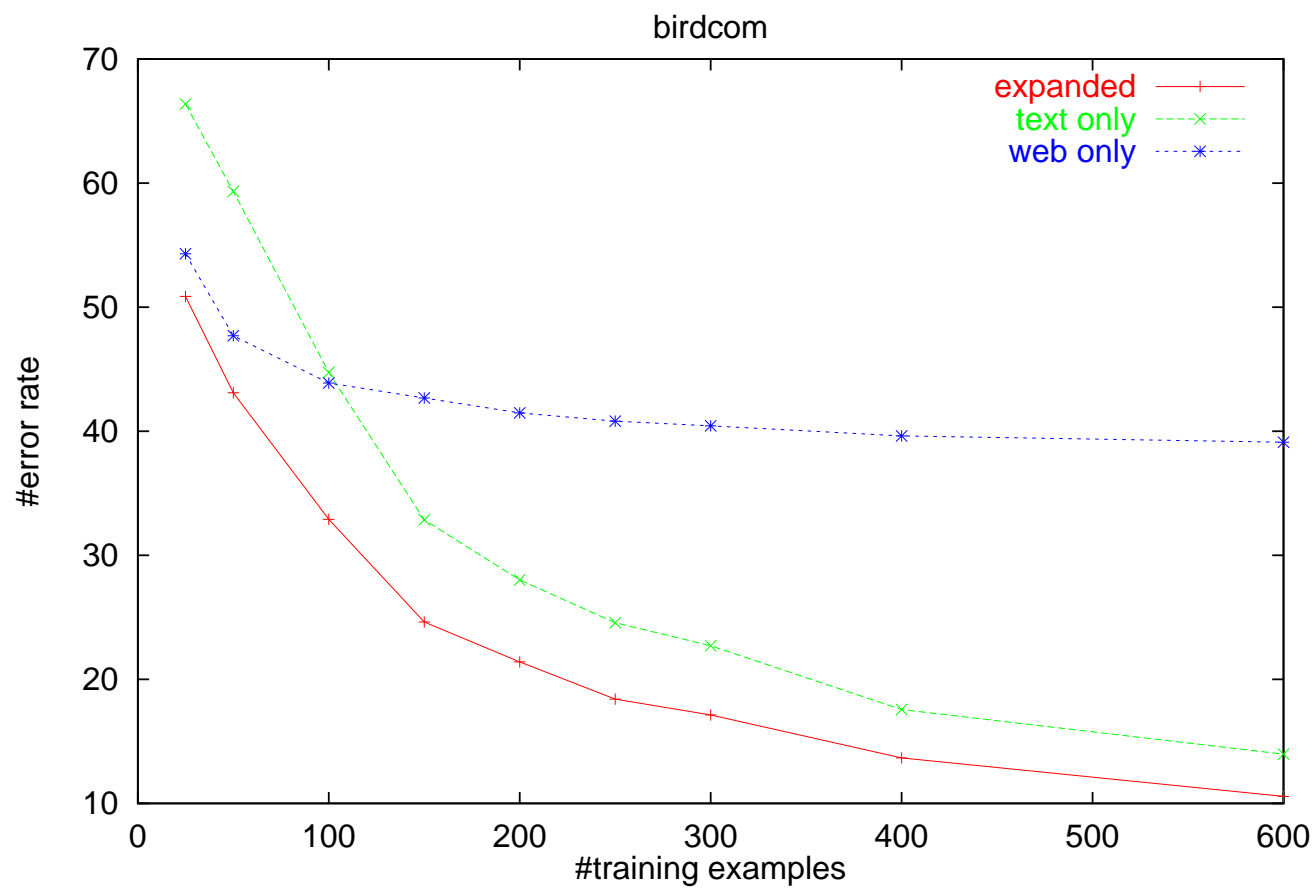
## Results (with RIPPER)



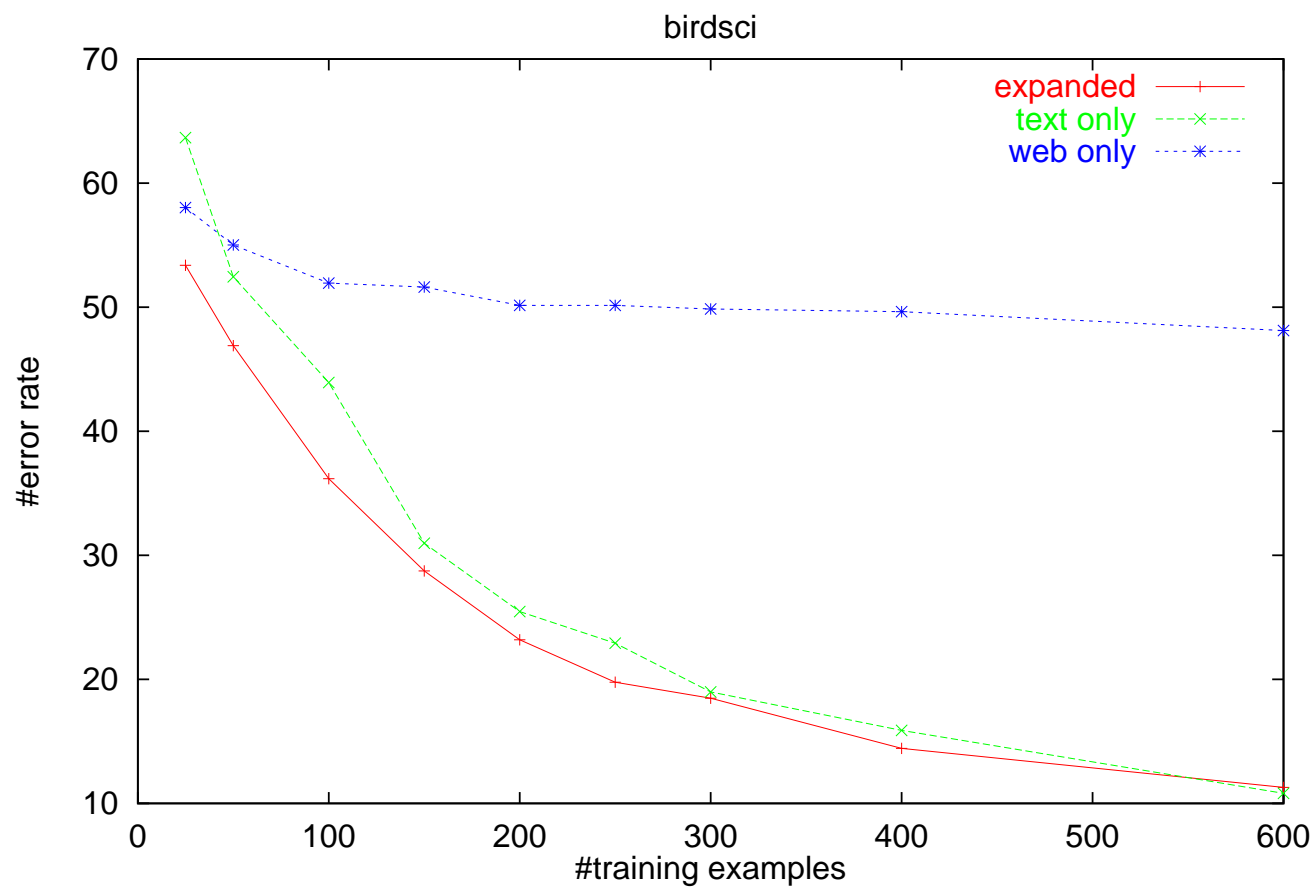
# Results

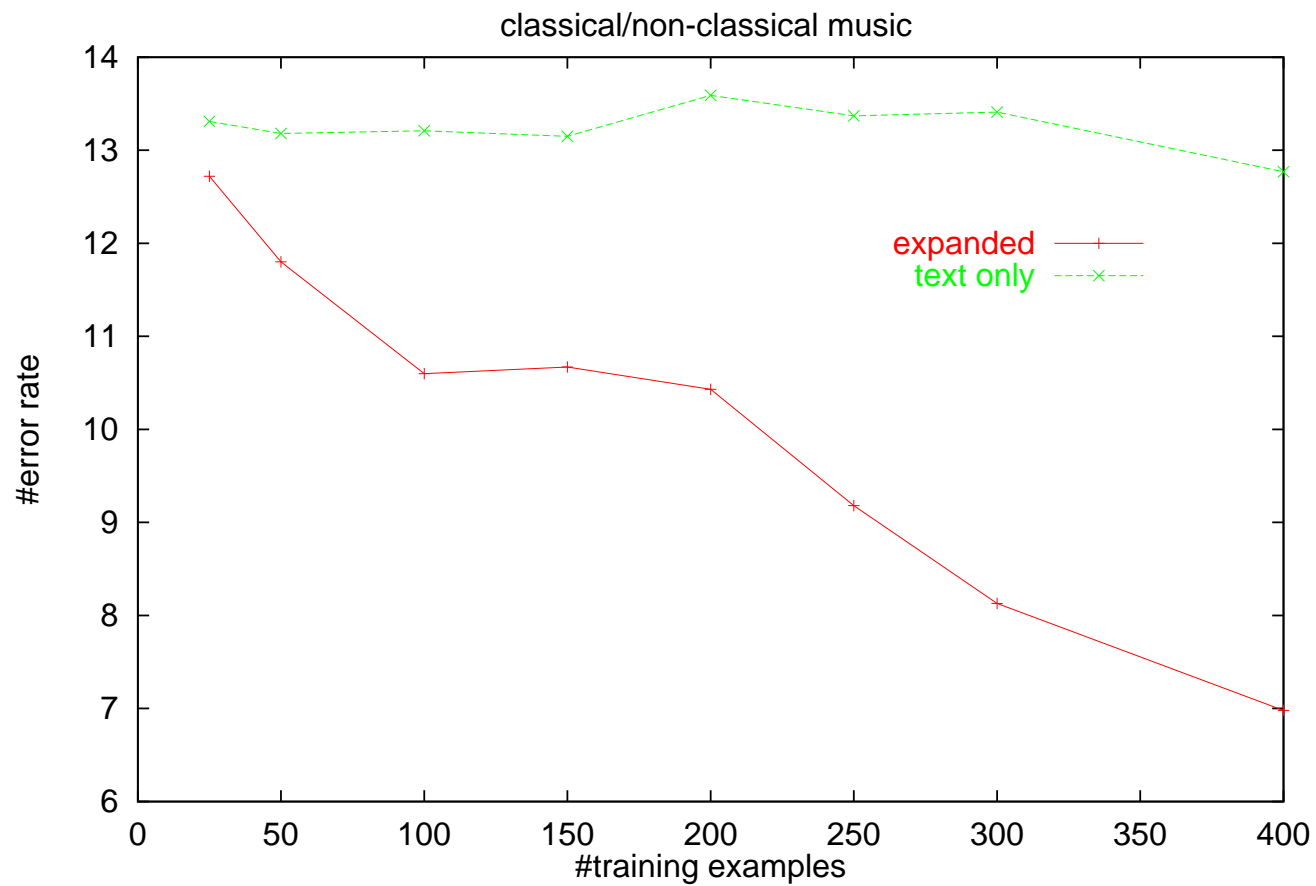


# Results



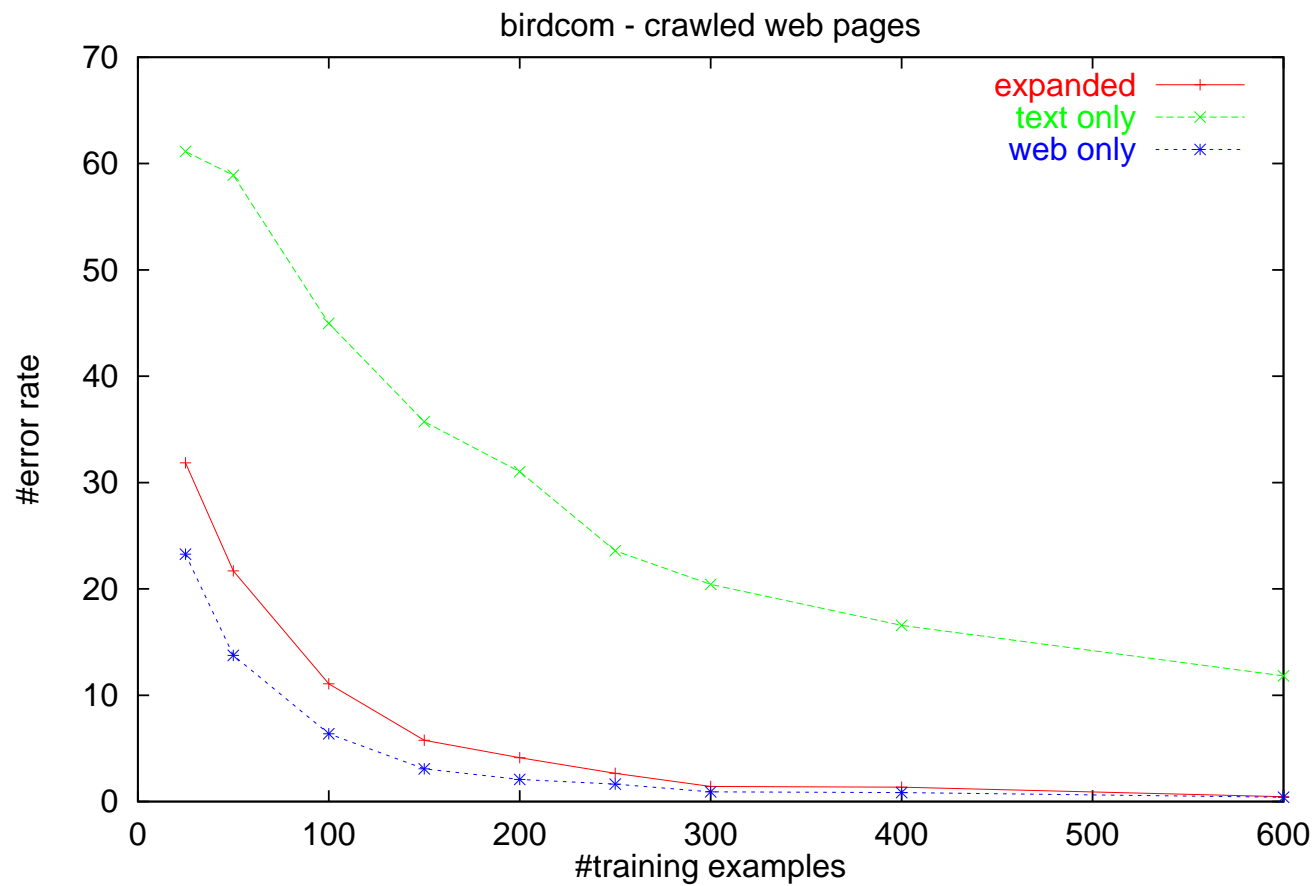
# Results





Distinguish classical/non-classical, restricting to artists for which some new features are constructed.





Web pages automatically crawled—not sampled from WHIRL DB on birds.

- Motivation: why this is the big problem.
- Data model: how information is stored.
- WHIRL query language
- Efficient implementation of WHIRL
- Results & applications
  - Queries without formal identifiers
  - Performance of a real query-answering system
  - Queries that generalize
  - Queries that don't require extraction
  - Queries that suggest extraction rules
  - Queries that automatically collect background knowledge for learning