

Searching for Statistical Diagrams

Michael Cafarella
University of Michigan

Joint work with
Shirley Zhe Chen and Eytan Adar

U.S. Frontiers of Engineering Symposium
2011



Statistical Diagrams

- n Everywhere in serious academic, governmental, scientific documents
- n Our only peek into data behind docs
- n Previously rare and precious, Web gives us a *flood*
 - n In small Web crawl, found 319K diagrams in 153K academic papers
- n Google makes it easy to find docs, images; very hard to find diagrams
- n Searching for diagrams part of larger semantic processing trends

U.S. Population 1900 - 2100

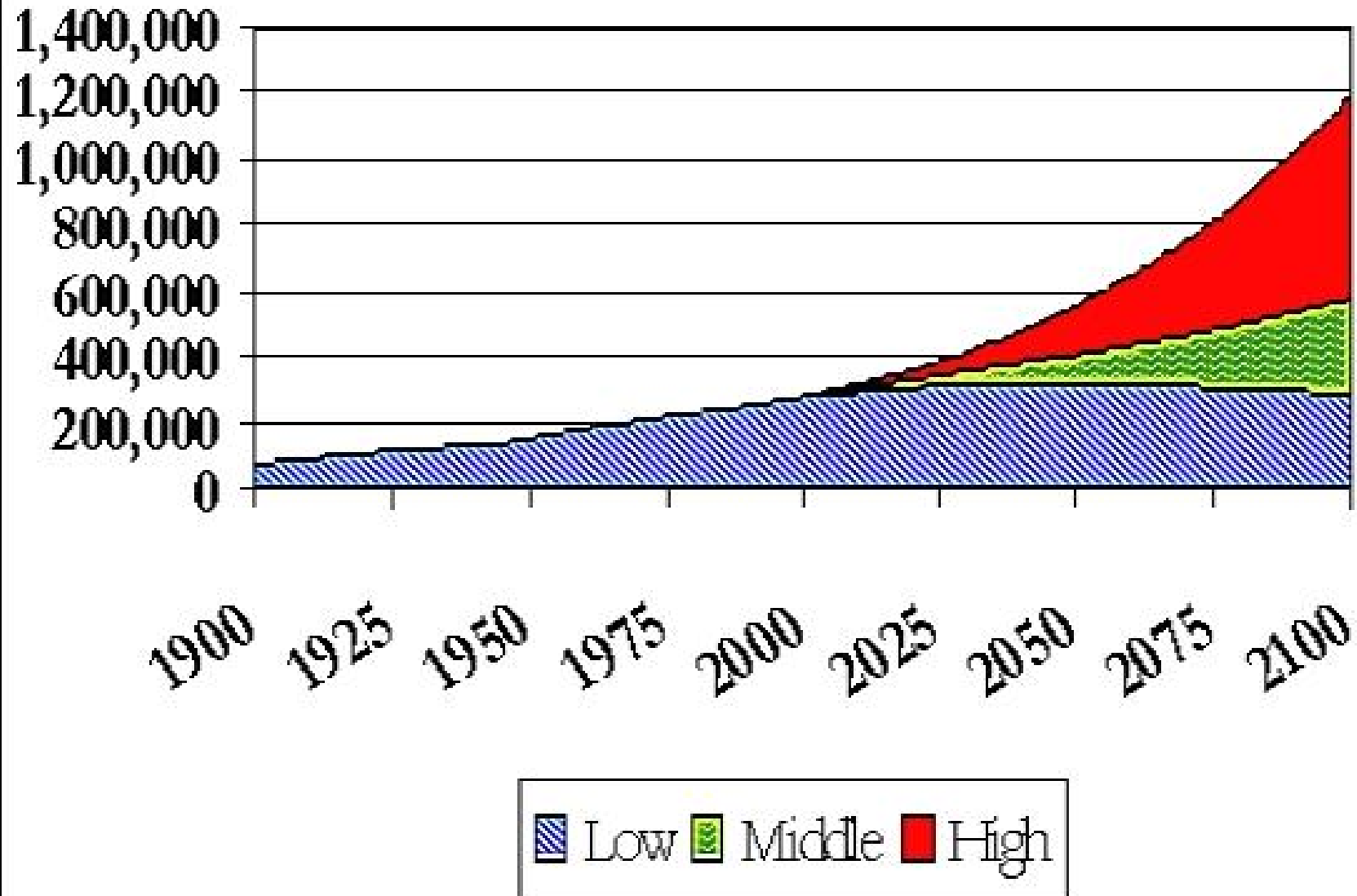
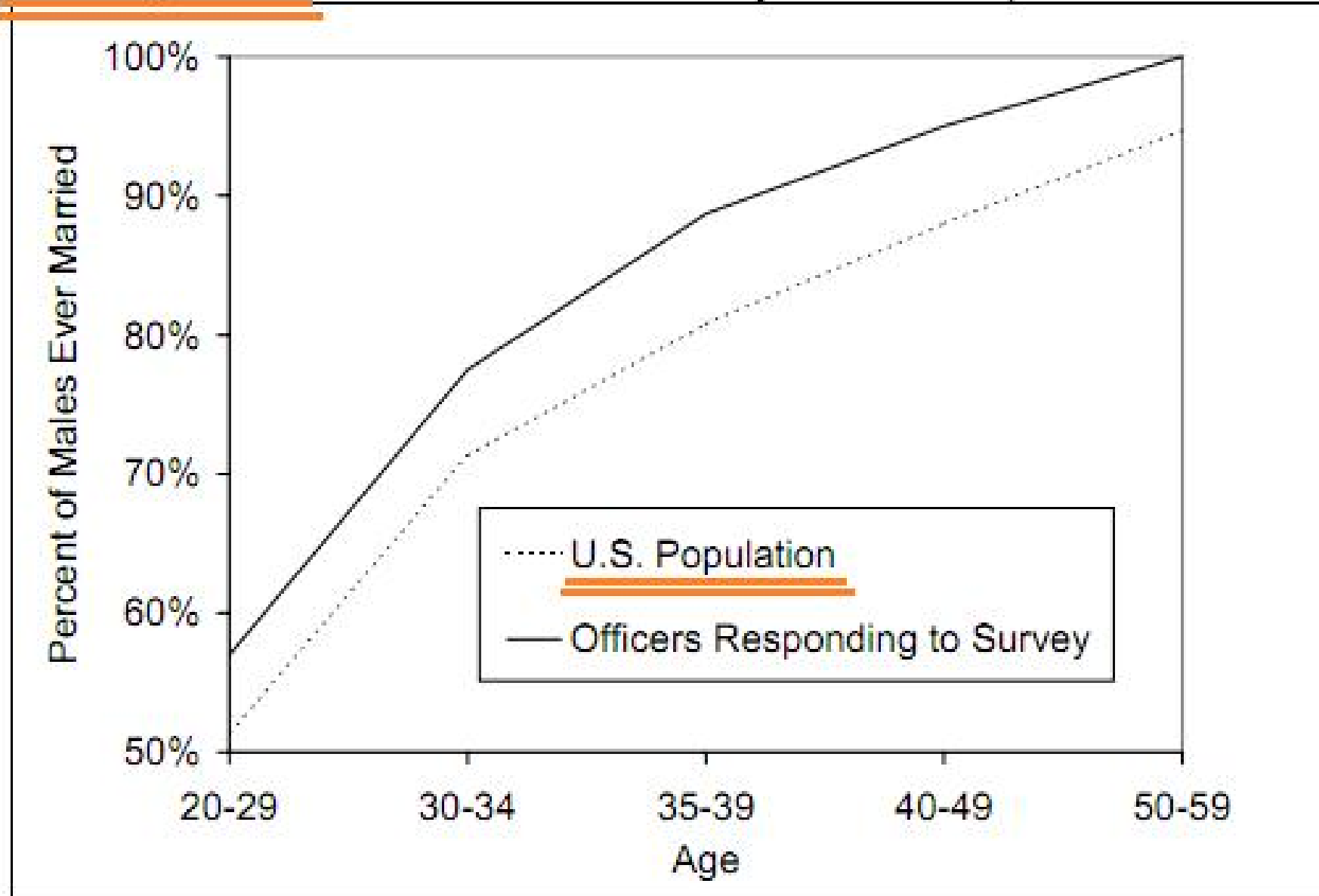
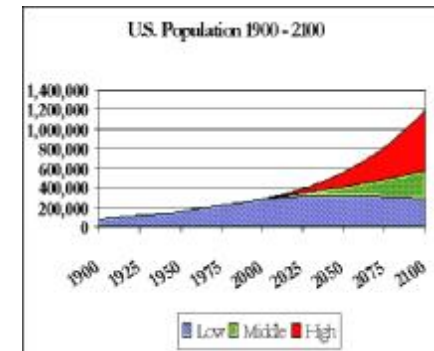
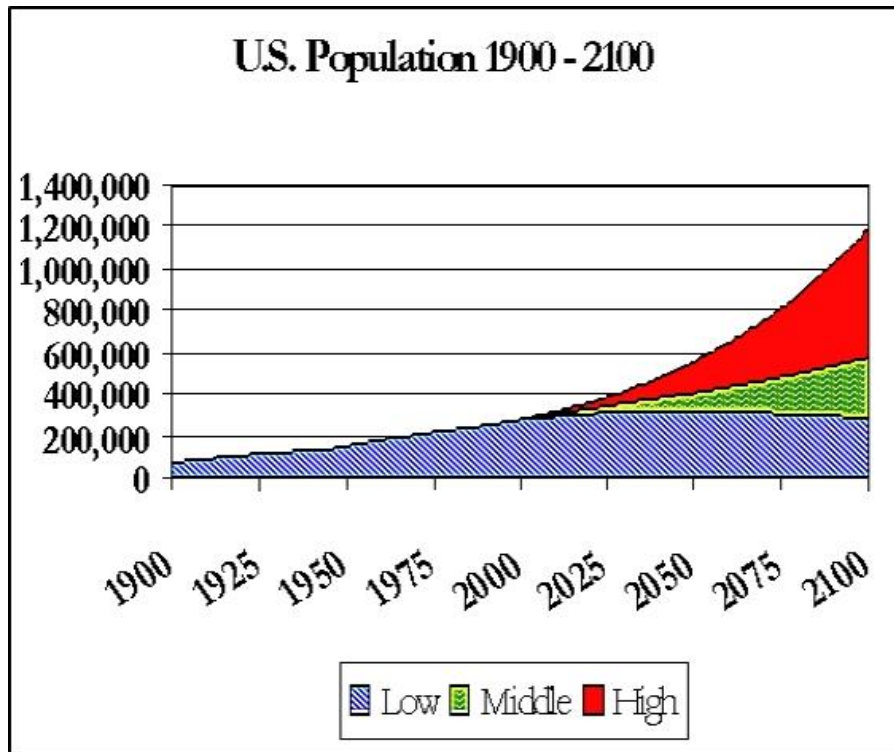


Figure 7.1 Comparison of the Percent of Males Ever Married, U.S. Population and Sacramento County Sheriff's Department Officers

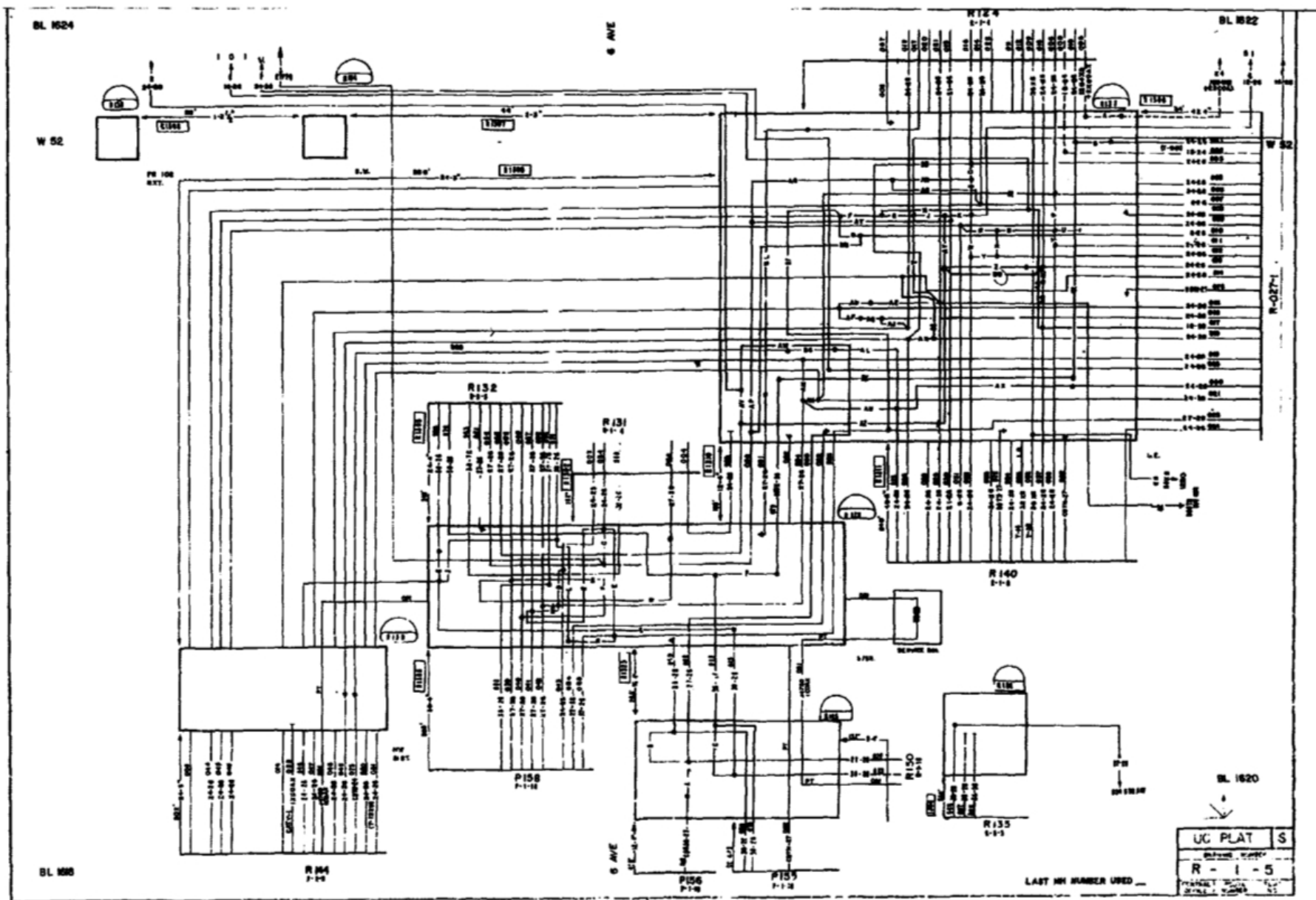




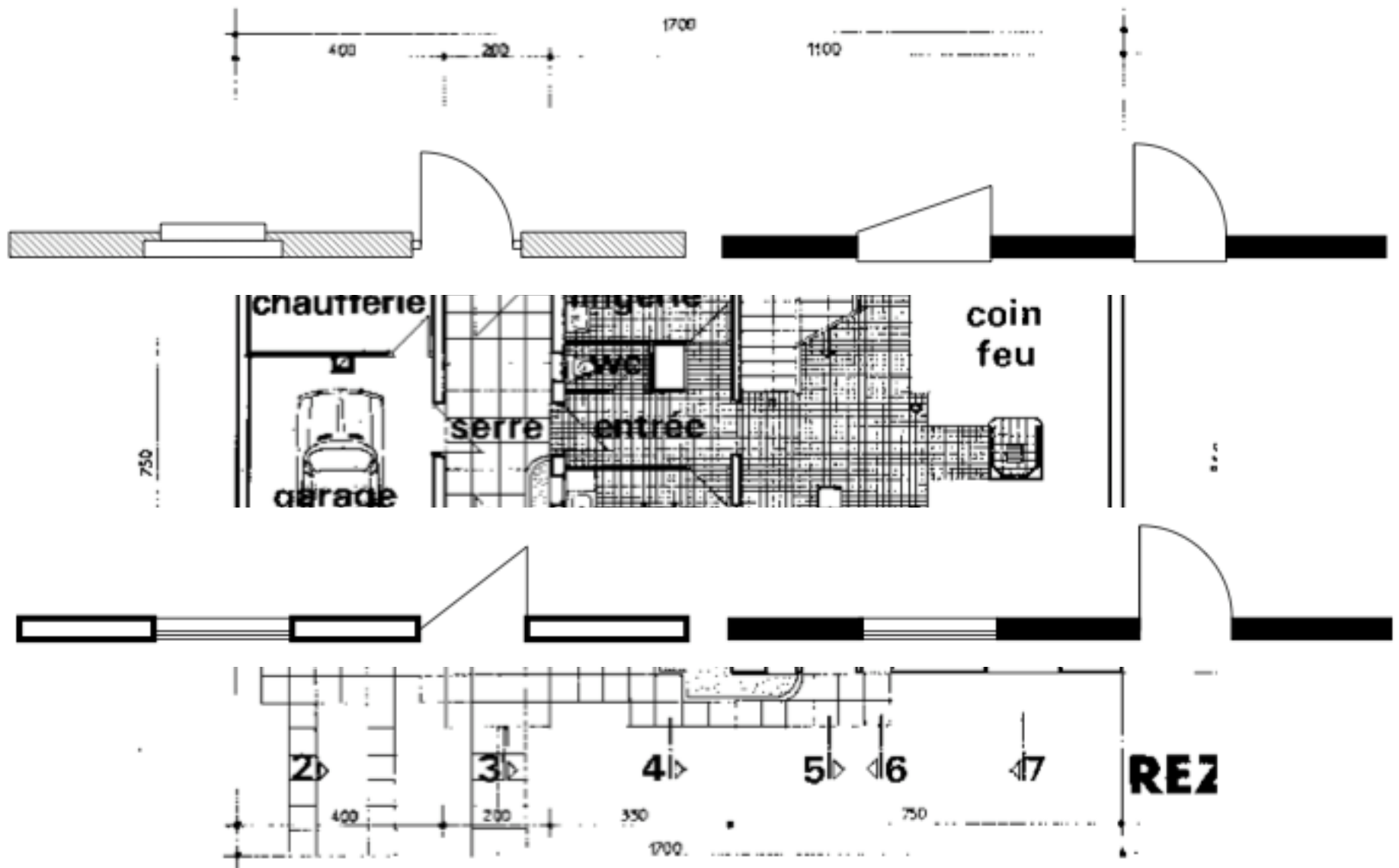


Previous Work

- n Searching for diagrams requires some amount of understanding
- n Lots of work in image search, most inapplicable to diagrams
- n But even understanding diagrams isn't new

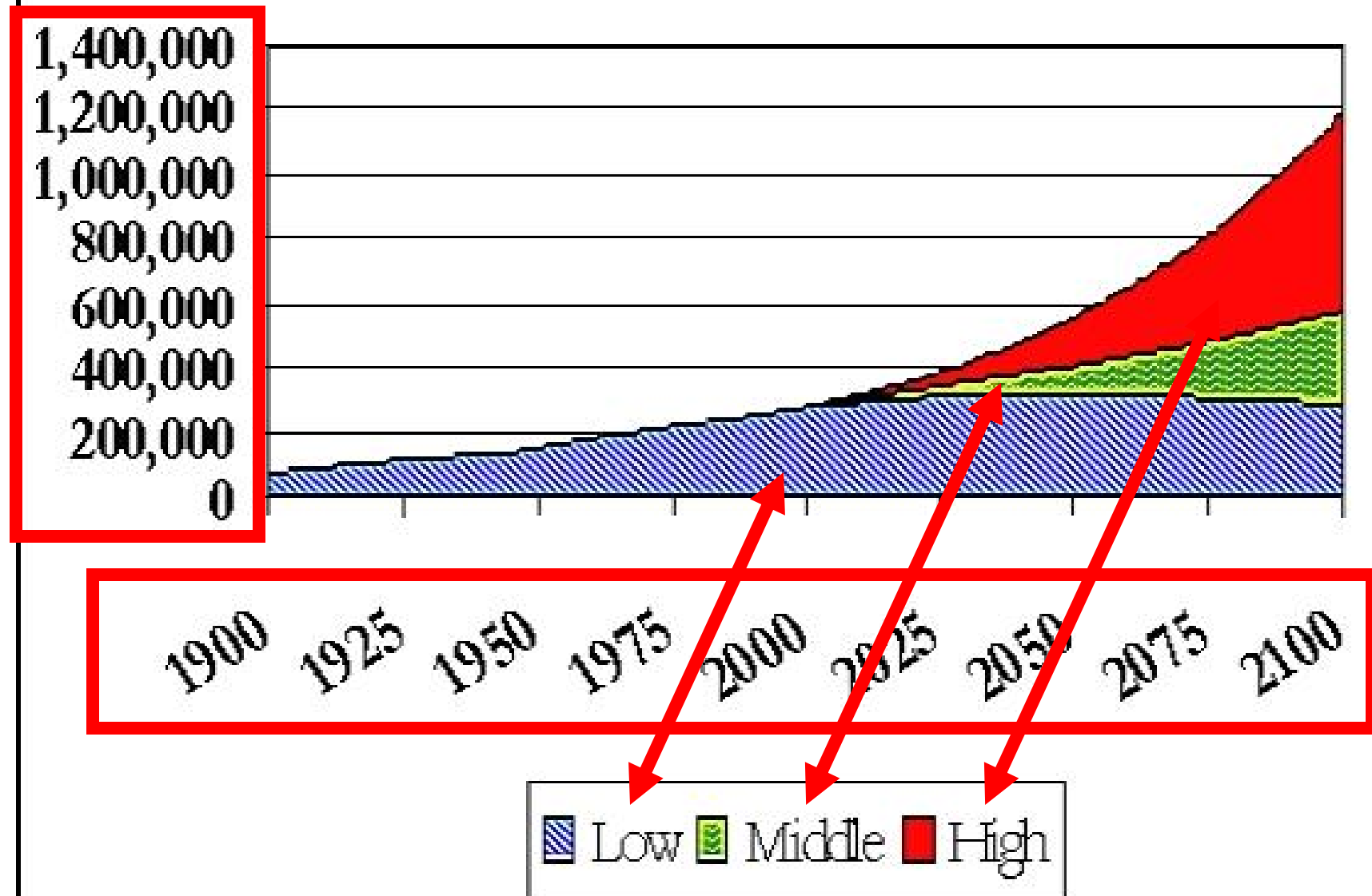


Telephone System Manhole Drawing, from
Arias, et al, Pattern Recognition Letters 16, 1995



Sample Architectural Drawing, from [Ah-Soon and Tombre](#),
[Proc'ds of Fourth Int'l Conf on Document Analysis and Recognition, 1997](#).

U.S. Population 1900 - 2100





Previous Work

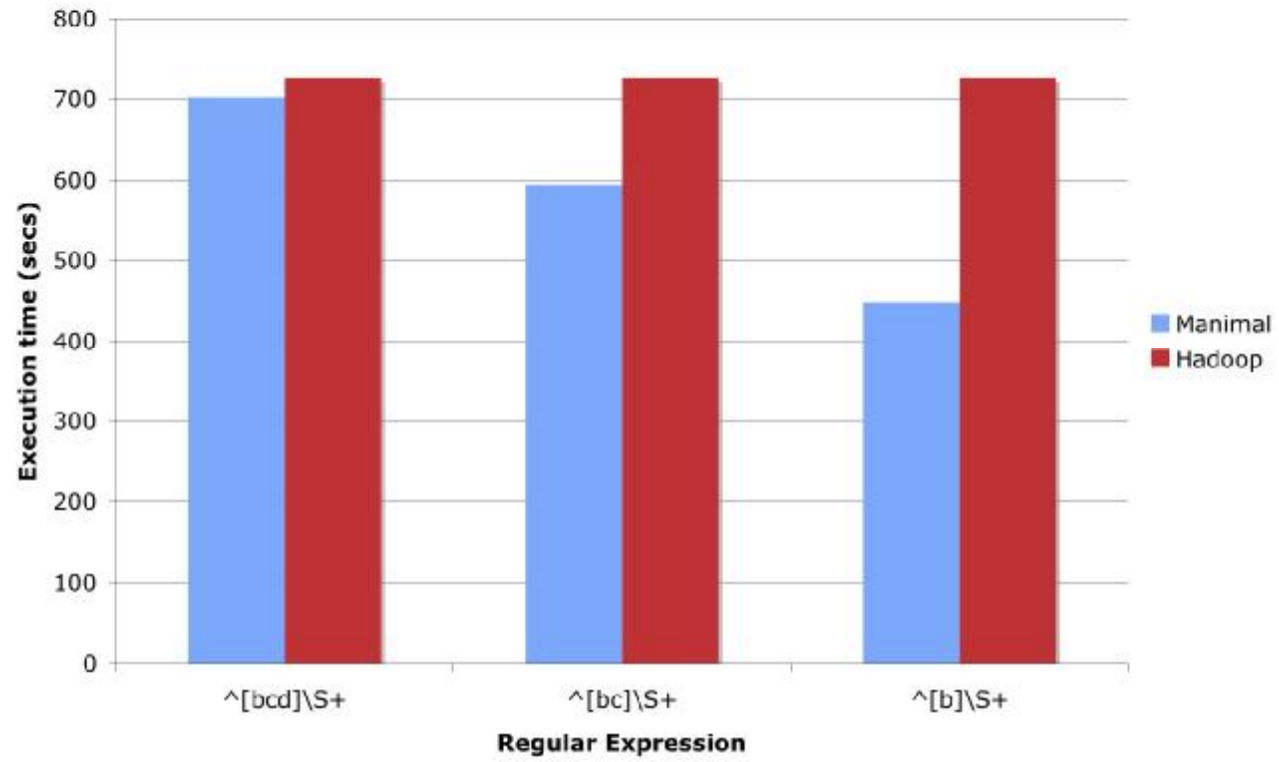
- n Understanding diagrams isn't new
- n Understanding a *Web's worth* of diagrams is new
 - n Need to search statistical diagrams in medicine, economics, biology, physics, *etc*
- n The phone company can afford a system tailored for manhole diagrams, but we can't
- n Effective scaling with # of topics is central goal of *topic-independent information extraction*



Topic-Independent IE

- n Information extraction topic since early 1990s
- n Goal is to obtain structured information from unstructured raw documents
 - n [Title, Price] from online bookstores
 - n [Director , Film] from discussion boards
 - n [Scientist, Birthday] from biographies
- n Traditional solutions require topic-specific code, features, data
- n Costs of TI IE do not grow with # topics

Grep Task



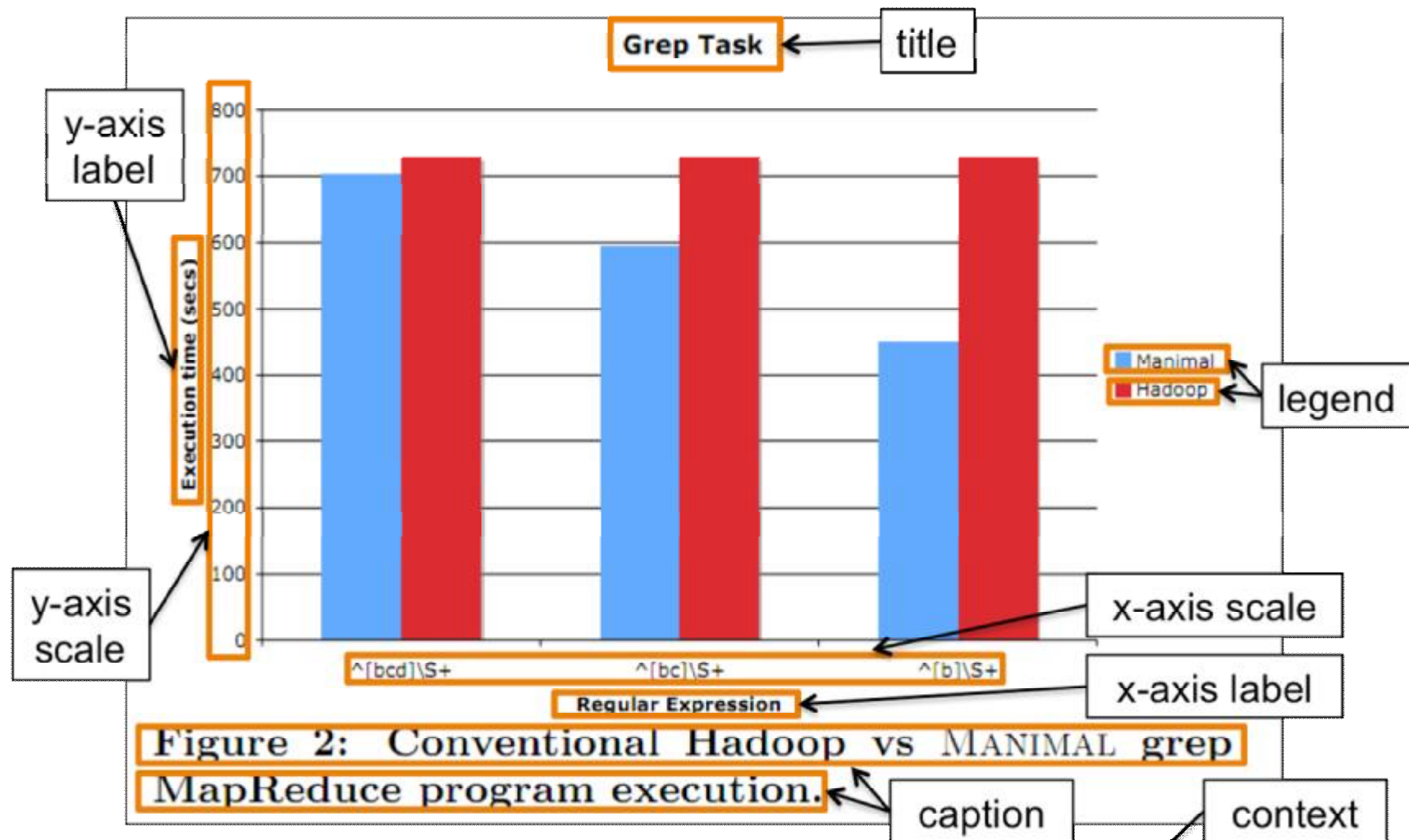
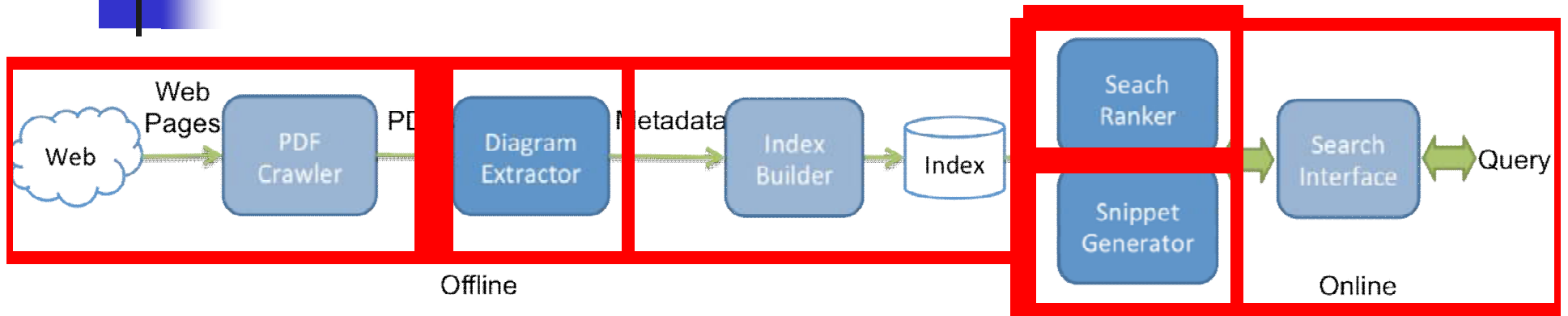


Figure 2 shows the results. As expected, Hadoop's conventional MapReduce execution time is almost wholly insensitive to the selectivity of the program. MANIMAL execution time, in contrast, decreases as the conditional test becomes more restrictive, dropping to 63% of Hadoop in the case of $^[b]\\S+$. There is nothing about the MANIMAL approach

Our Approach



- n Typical Web search pipeline
 - n Crawl Web for documents
 - n Obtain and index text
 - n Make index queryable
- n Our novel components
 - n Diagram metadata extraction
 - n Custom search ranker
 - n Snippet generator



Metadata Extraction

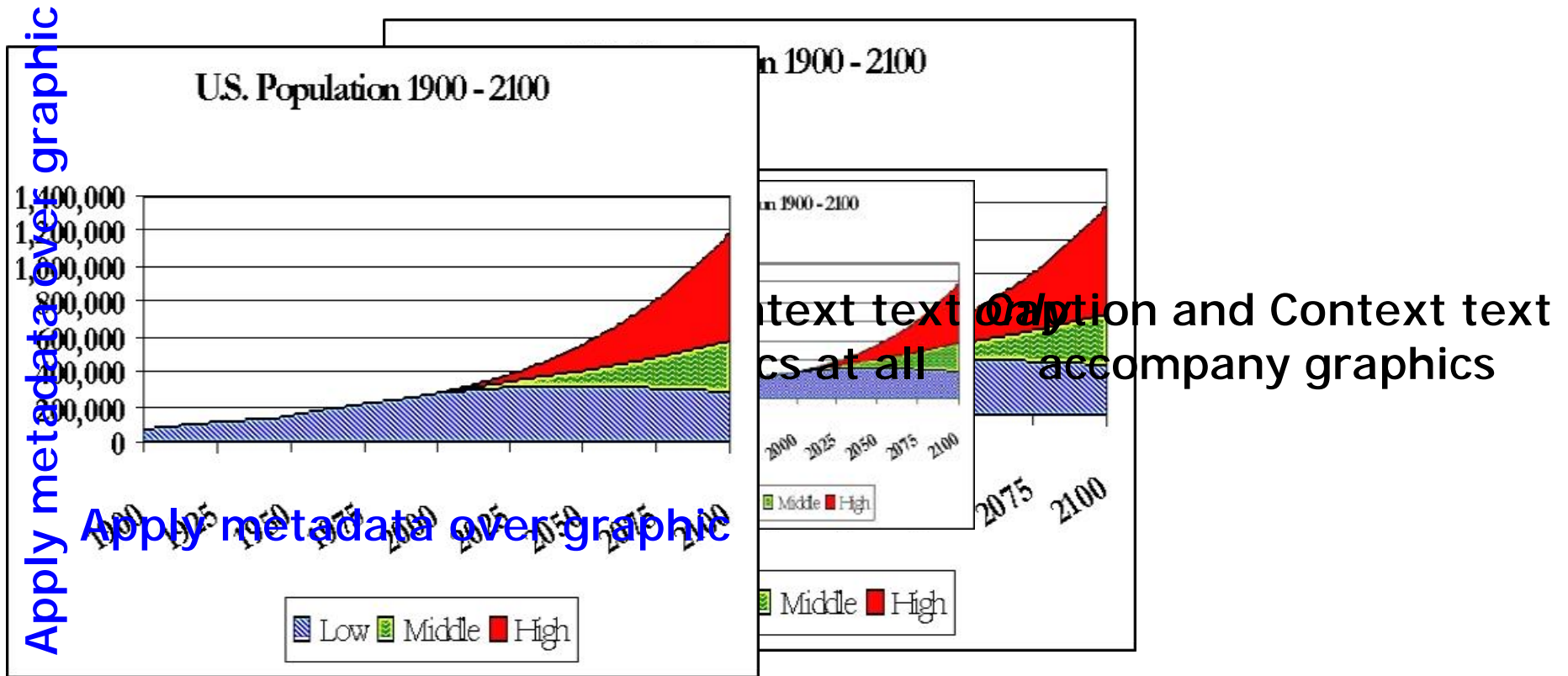
1. Recover good (text, x, y) from PDFs
2. Apply simple role label: *title*, *legend*, etc

Figure 1: Frequency histogram for 2925 syllables derived from the Switchboard corpus (upper illustration).

3. Group texts into “model diagram” candidates, throw away unlikely ones
 - n E.g., must include *something* on x scale
4. Relabel text using geometric relationships
 - n Distance, angle to diagram’s origin?
 - n Leftmost in diagram? Under a caption?

Snippet Generation

Tested five versions



3. Only generate snippet



Experiments

- n Crawled Web for scientific papers
 - n From ClueWeb09
 - n Any URL ending in **.pdf** from **.edu** URL
 - n 319K diagrams
- n Fed data to prototype search engine
- n Evaluated
 - n Metadata extraction
 - n Rank quality
 - n Snippet effectiveness
- n All results compared against human judgments



population

Search

Human Population Billions

compared with ~1Tg only 50 years ago (The Fertilizer Institute, 2000; International Fertilizer Industry Association (IFA), 2004).

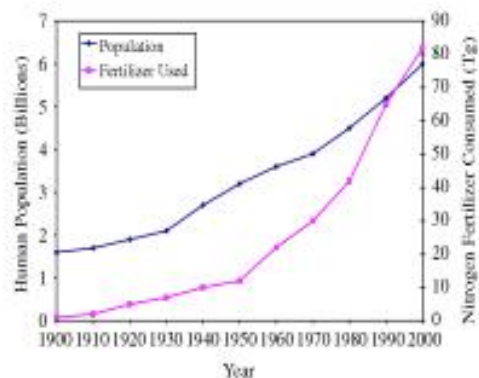


Fig. 1. Graph showing population increase and use of nitrogen fertilizer from 1900 to 2000.

While NH from agricul tural air poll concern, inc [e.g., nitrog nitrous oxid (VOCs) (e.g. organic aci particulates particle conv (e.g., hydrog emissions. o practices inc tions (CAFi manure and biomass burn

In many a managed cro parallelly to for food. In

Tags: fertilizer, century, population, billions, yr,

Caption: Fig 1 Graph showing **population** increase and use of nitrogen fertilizer from 1900 to 2000

X-axis Label: Year

Y-axis Label: Human **population** Billions Nitro g en Fertilizer Consumed T g

Legend: population Fertilizer Used

Title:

Context: Fig 1 shows the parallel increase in human **population** and fertilizer usage over the past century Currently the global production of fertilizer is more than 80 Tg of N yr 1 compared with 1 Tg only 50 y

Required Sample

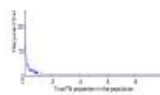


Figure 1: Minimal sample size for estimating within relative error bounds.

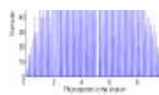


Figure 2: Minimal sample size for estimating within absolute error bounds.

Exam data with By discov strata, wh On the c estimate, be obtain

Y-axis Label: Required Sample Size [Percent]

Legend: Relative Error of 20 * **population** 100 000

Title: Sample size required to estimate the FN proportion in the **population**



1. Experiments - Extraction

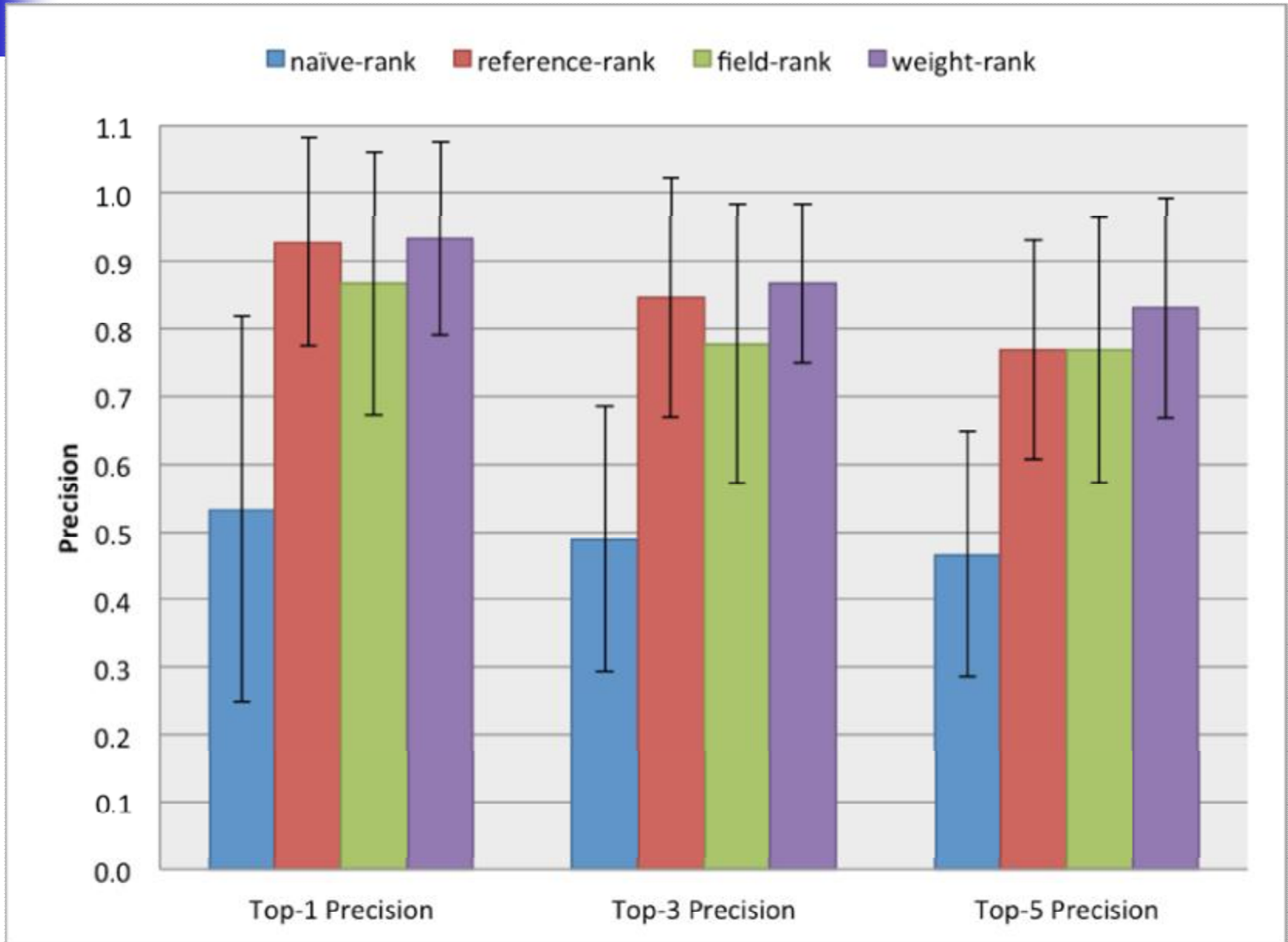
	Recall			Precision		
	<i>Text</i>	<i>All</i>	<i>Full</i>	<i>Text</i>	<i>All</i>	<i>Full</i>
title	0.256	0.651	0.674	0.344	0.609	0.617
Y-scale	0.782	0.796	0.754	0.899	0.843	0.900
Y-label	0.835	0.864	0.874	0.775	0.752	0.797
X-scale	0.903	0.835	0.835	0.616	0.915	0.896
X-label	0.241	0.681	0.681	0.340	0.842	0.835
legend	0.520	0.623	0.656	0.349	0.615	0.631
caption	0.952	0.887	0.839	0.450	0.887	0.929
nondiag	0.768	0.924	0.313	0.850	0.909	0.838



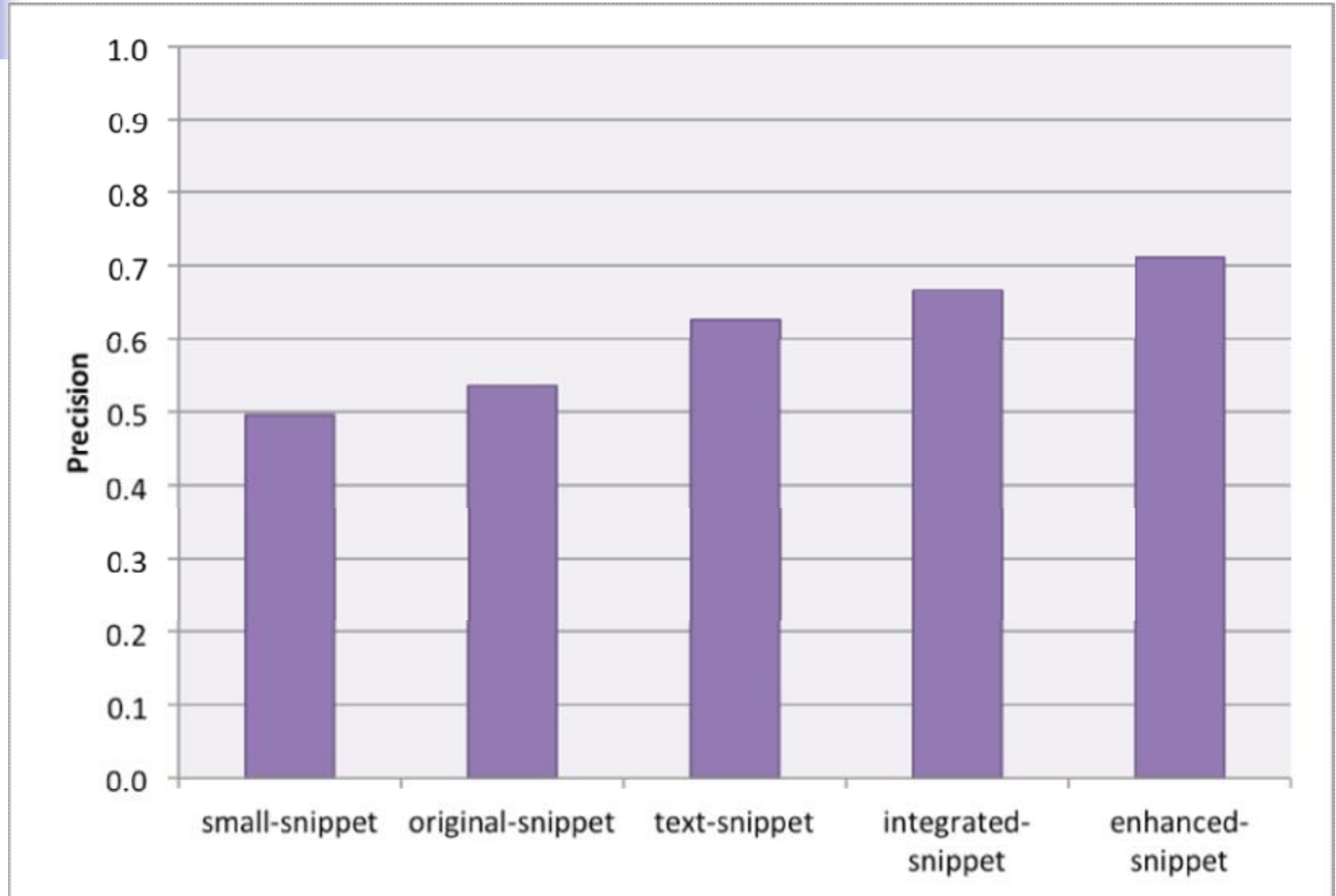
1. Experiments - Extraction

	Recall			Precision		
	<i>Text</i>	<i>All</i>	<i>Full</i>	<i>Text</i>	<i>All</i>	<i>Full</i>
title	0.256	0.651	0.674	0.344	0.609	0.617
Y-scale	0.782	0.796	0.754	0.899	0.843	0.900
Y-label	0.835	0.864	0.874	0.775	0.752	0.797
X-scale	0.903	0.835	0.835	0.616	0.915	0.896
X-label	0.241	0.681	0.681	0.340	0.842	0.835
legend	0.520	0.623	0.656	0.349	0.615	0.631
caption	0.952	0.887	0.839	0.450	0.887	0.929
nondiag	0.768	0.924	0.313	0.850	0.909	0.838

2. Experiments - Ranking



3. Experiments - Snippets





Other Applications

- n Working now

- n Search by axis label
 - n Search by range
 - n Given a query diagram (or paper), find related papers

- n In future:

- n Improved academic paper search
 - n Show plots that support my hypothesis



Future Work

- n Spreadsheets

- n Has experiment X ever been run before?
 - n WY GDP vs coal production in 2002
 - n Preemptively compute good diagrams

- n Deeper questions for messy data

- n HTML tables, data files, spreadsheets
 - n Lots of structured data lives outside DBMS

- n Structured search



Conclusions

- n Metadata extraction enables **52%** better search ranking
- n Extraction-enhanced snippets allow users to choose **33% more accurately**
- n We rely on open information extraction, but extracted data not the main product
 - n Can be successful even with imperfect extractors



Thanks

- n Academy of Engineering
- n FOE sponsors
- n Google
- n You!



Related Work

- n Suitable for Web search settings

- n Huang *et al*, "Associating text and graphics...", ICDAR 2005
- n Huang *et al*, "Model-based chart image recognition", GREC 2003
- n Kaiser *et al*, "Automatic extraction...", AAAI 2008
- n Liu *et al*, "Automated analysis...", IJDAR 2009

- n Diagram parsing

- n E.g., Futrelle, "Summarization...", 1999

- n Visually-impaired access

- n E.g., Demir *et al*, "Generating textual...", INLG 2008.