

# Structure in Text: Extraction and Exploitation

Prabhakar Raghavan



# Agenda

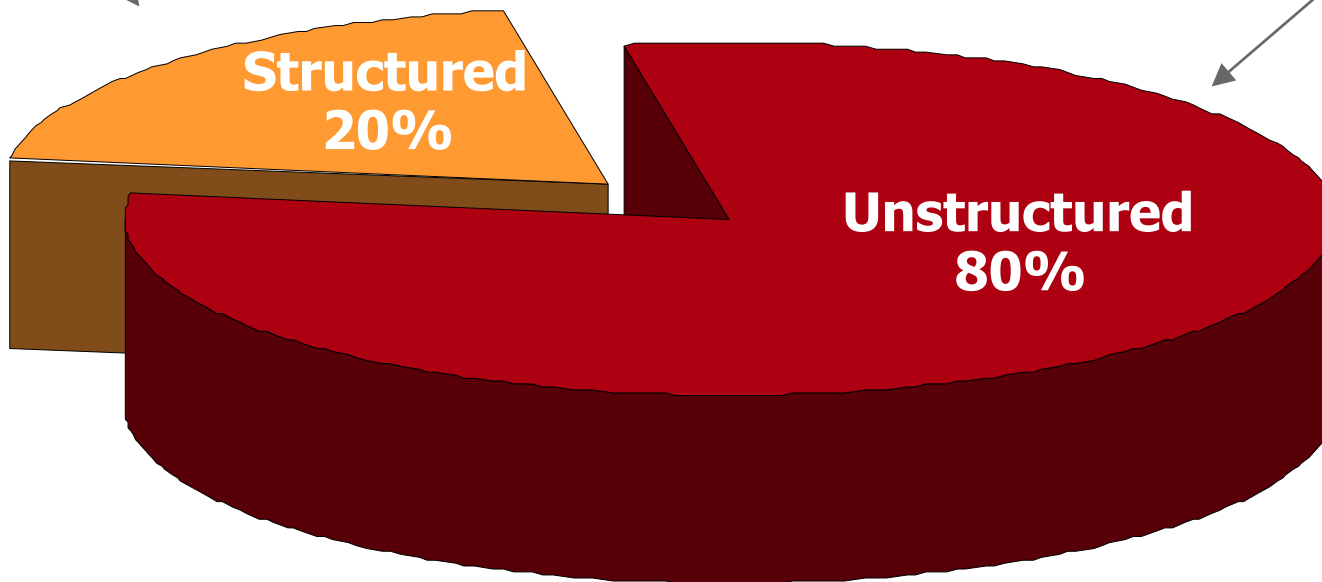
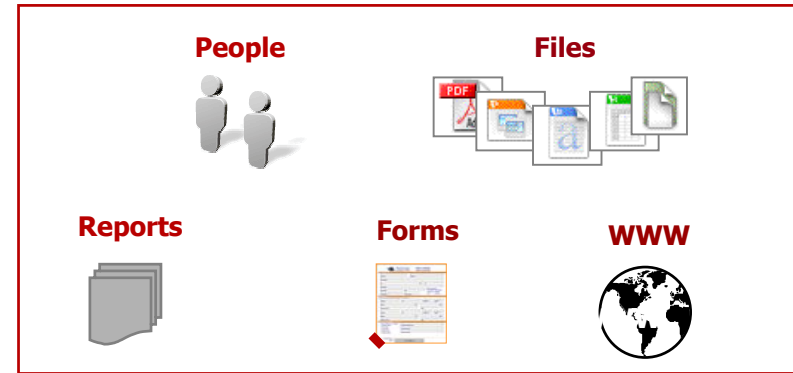
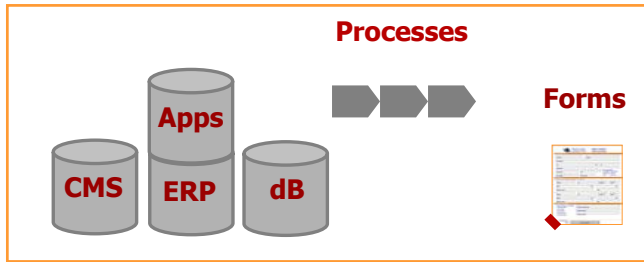
- **Market overview**
  - Structured and unstructured search – what convergence?
- **Application scenarios**
  - What my customers (think they) want
  - What it will take
- **Technical tofu**
  - Distance metrics for XML documents
  - Linear algebra as a foundation
- **Summary/research agenda**



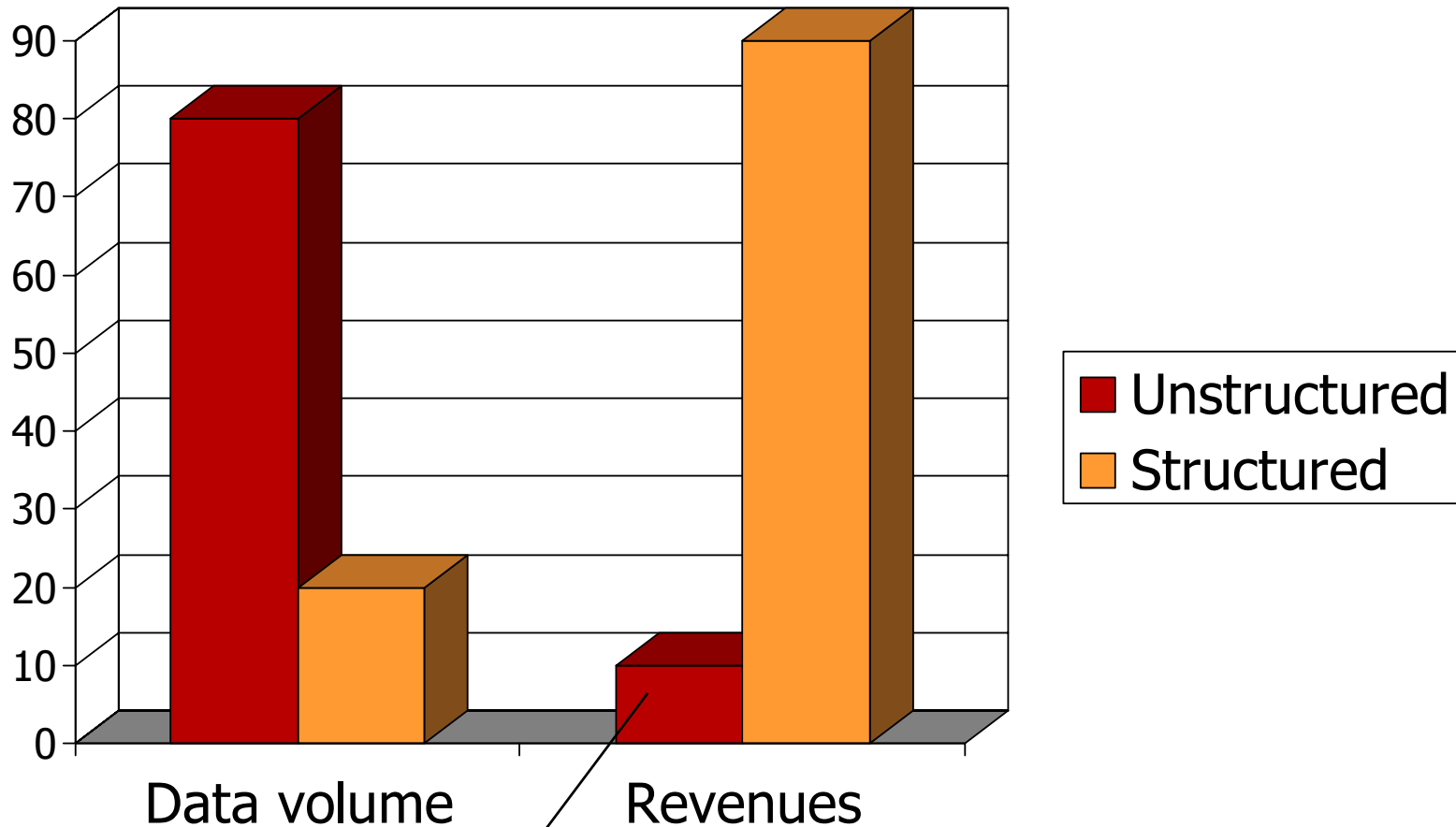
# Market overview



# Unstructured vs. Structured Data in the Enterprise



# Data: Volume vs. Revenues



~\$500M LFR



# Consequences



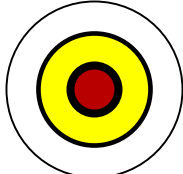



- **Value creation demands:**
- **Extract structure**
  - Document capture/conversion
  - Classification, linguistic tagging
  - Entity/relation extraction
  - The results will be noisy
- **Exploit structure**
  - Ad-words
  - Structured navigation
  - XML querying



# Round pegs and square holes



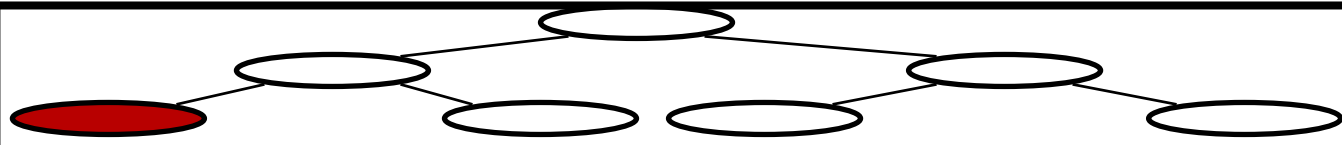
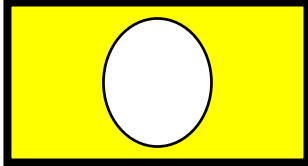

	Unstructured	Structured
Unit of Retrieval		
Selection for Retrieval	<p>Distance Metric</p> 	<p>Hard Selection</p> 
Results	<p>Consumed by humans: Limited Number. “Sort of” right; OK to get it wrong – Human “ratification”.</p>	<p>Consumed by machines: Composition: Need correct answers. Nuclear powerplant syndrome.</p>



# Semi-structured text: Round pegs meet square holes



## Semi-structured

<b>Unit of Retrieval</b>	
<b>Selection for Retrieval</b>	Select then rank  ??  <u>Semi-soft measures</u>
<b>Results</b>	Will still get it wrong sometimes. Want to dramatically limit human ratification.





# Example – relational taxonomies

**Shop. Learn. Improve.**
homedepot.com ...we're always open.

Home | Gift Cards | Expo | Store Locator | Project Index | Credit Center

Choose a Department |  |  | YOUR SHOPPING CART

NUMHITS 1984 - NUMPROCS 46341

**Aisles/Tools**

[Air Tools-Air Compressors-Air Tool Fasteners \(9\)](#)

[Hand Tools \(48\)](#)

[Power Tool Accessories \(1725\)](#)

[Power Tools \(103\)](#)

[Saw Horses and Workbenches \(1\)](#)

[Shop Vacuums-Air Filtration \(1\)](#)

[The Hilti Store \(93\)](#)

[Welding and Soldering \(4\)](#)

**Brand**

[Norton \(221\)](#)

[Oldham \(191\)](#)

[Blu-Mol \(166\)](#)

[Dewalt \(150\)](#)

[Vermont American \(136\)](#)

[Bosch \(111\)](#)

[more choices...](#)

**Price**

[< \\$10.00 \(1140\)](#)

[Home](#) | [Tools/](#)

Type

Use

Brand/Model Compatibility

[more choices...](#)

1 2 3 4 5 6 7 8 9 10 11 [Next](#)

Sort results by:				
	<a href="#">Brand</a>	<a href="#">Detail/Category</a>	<a href="#">Country</a>	<a href="#">Price</a>
	<a href="#">52-Piece Power Drill and Bit Set</a> SKU: 793610			
	Workforce	Drill Bit Kits	China	19.97
	<a href="#">3/8" Air Powered Drill</a> SKU: 997919			
	Campbell Hausfeld	Air Drills	Taiwan	29.84
	<a href="#">7/16" Masonry Drill Bit</a> SKU: 508650			
	Vermont American	Masonry Drill Bits	United States	3.96
	<a href="#">3/8" Masonry Drill Bit</a> SKU: 232557			
	Vermont American	Masonry Drill Bits	United States	3.29
	<a href="#">1/4" Masonry Drill Bit</a> SKU: 232526			
	Vermont American	Masonry Drill Bits	United States	2.79
	<a href="#">3/16" Masonry Drill Bit</a>			

**Homer Recommends**

[Tools/Power Tool Accessories/Portable Power Tool Accessories/Rotary Hammer Accessories](#)

[Tools/Power Tool Accessories/Portable Power Tool Accessories/Drilling and Screwdriver Accessories/Masonry Drill Bits](#)

[Tools/Power Tool Accessories/Portable Power Tool Accessories/Demolition Hammer Accessories](#)

[Bosch](#)

Worm Drive Power Cord  
SKU: 499255

15-Piece Titanium Drill Bit Set  
SKU: 152794

3/16" x 4" x 6" Blue Granite Hammer Drill Bit  
SKU: 664522

©2002 Verity, Inc.

# Information retrieval platforms today

- Text and parametric search
- **Classification**
- Clustering and taxonomy induction
  - A modicum of “text mining”
- **Index only, no repository**
  - No content retrieval
    - (Exception: web search caches)
  - Affects “Return” portion of XML queries



## A step back: what's all this in aid of?

- End users don't want to
  - Search
  - Browse
  - Solve relevance problems
  - Deal with performance problems
  - Hear of your magic panacea for the above
- ***Users want to solve their business problems***



# Application scenarios



## Parental advisory

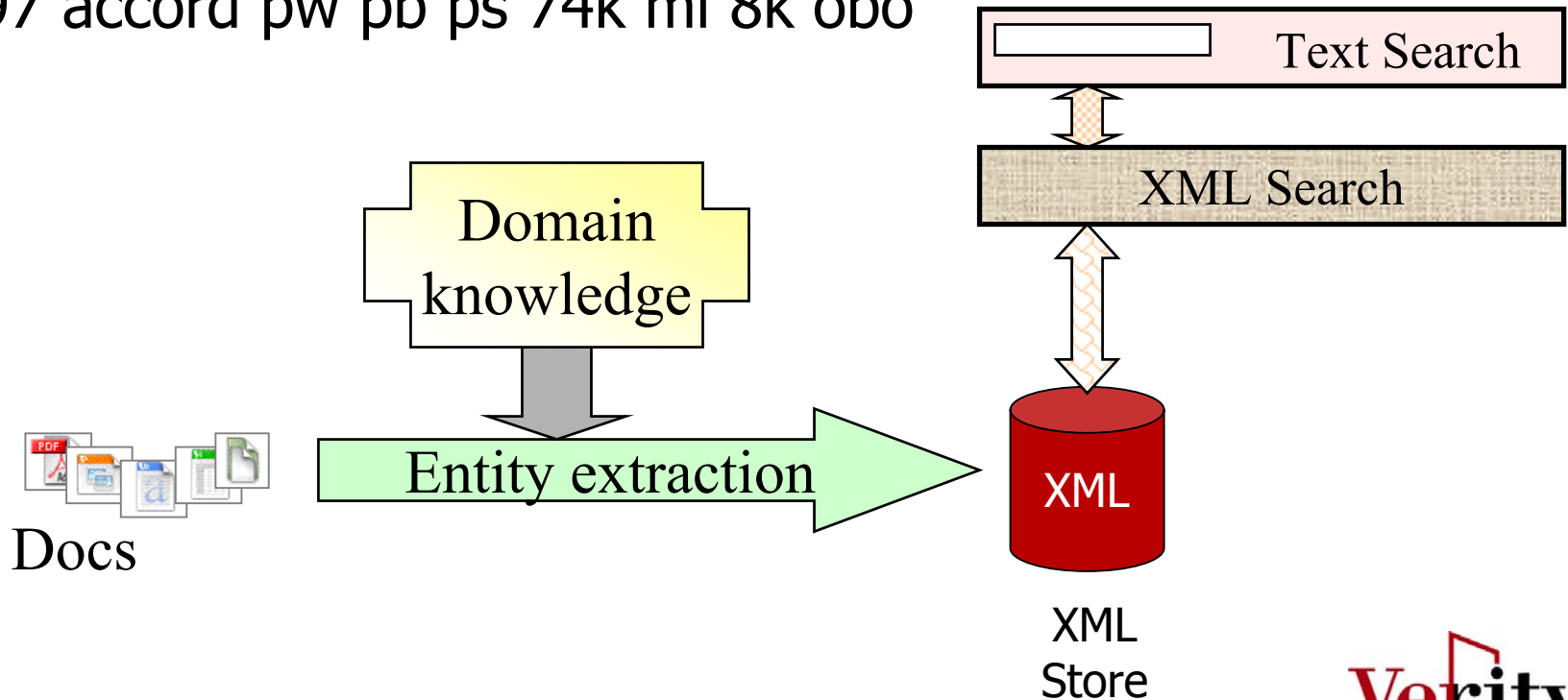


**Warning:** Some of the following slides have been produced by Corporate Marketing. The graphic(s) nature of these slides may overwhelm innocent researchers.



# Domain knowledge in search

- Giants tix pair Saturday game \$60
- Sunny 3br 2ba top floor vu 1600+util
- 97 accord pw pb ps 74k mi 8k obo



# Mining semi-structured text

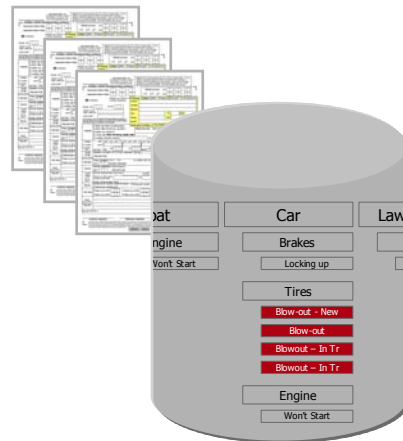


Product recalls can cost company's millions in profits and loss of brand value.

1. Product "problem tickets" completed at point of origin



2. "Tickets" converted real-time to dynamic, categorized information



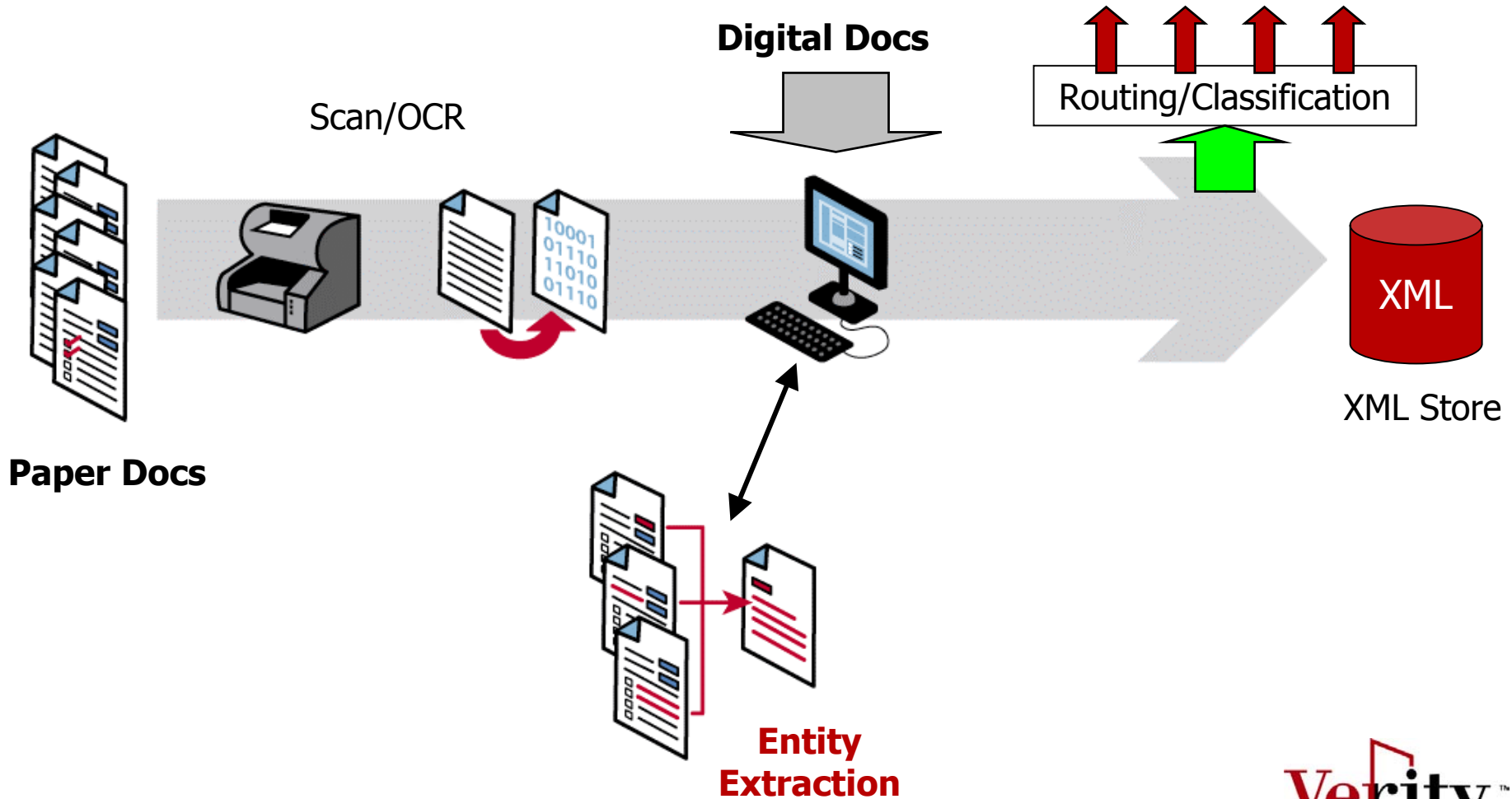
3. Smart monitoring identifies critical problems and send alert notice



4. Management can research, isolate and rectify problems quickly

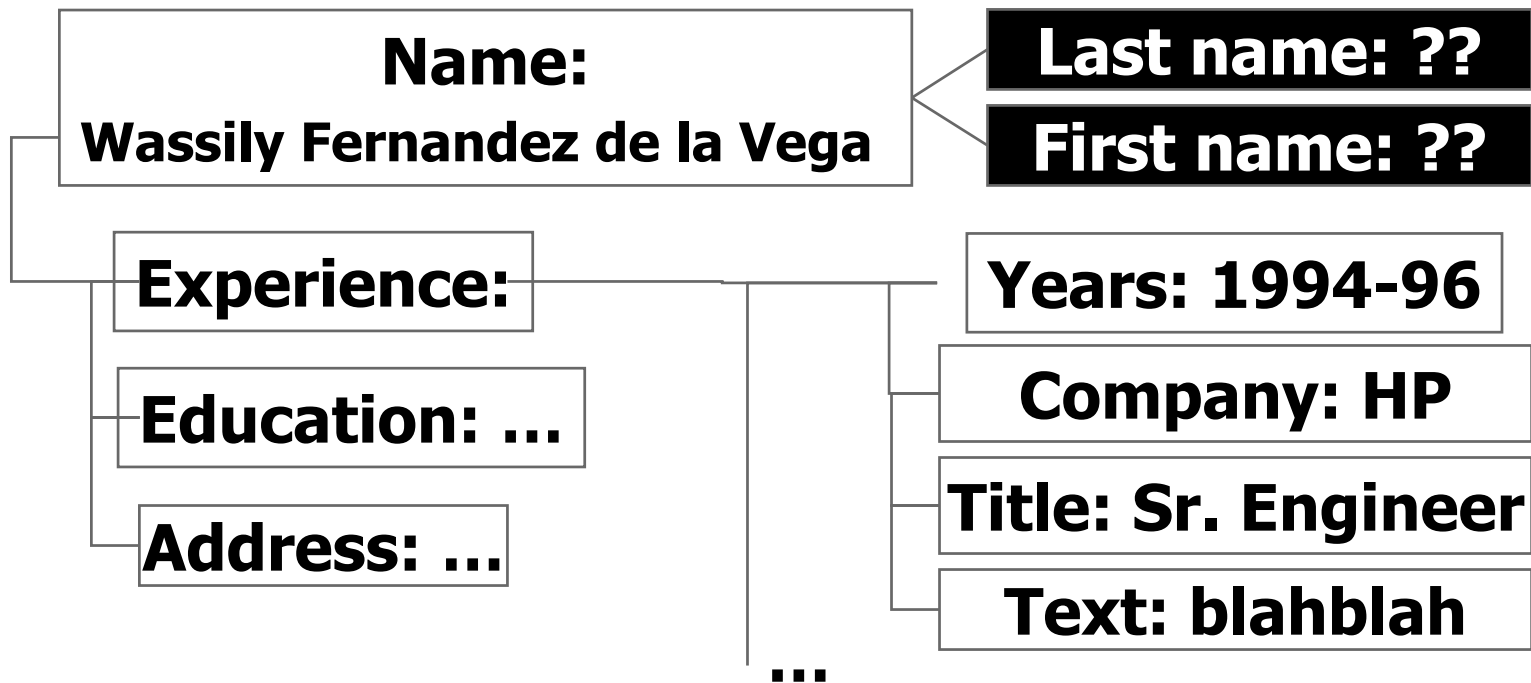


# Resume Extraction/routing





# Extracted Resume



# More semi-structured user needs

- Find loss-making software companies with revenues of 50-200M last year, and repeated mutual deals with other software companies within 90-day intervals.
  - *Input: SEC Filings, Press Releases*
- Want a family vacation to a warm place with outdoor activities, a few museums and churches, budget 3000 and I prefer to fly Singapore Airlines.
  - *Input: Car, Air, Hotel, Travel info/opinion websites*



# What will it take?

- User need parsing
  - Query expression and semantics
- Entity Extraction
  - (Semi-)structure from unstructured text
- Retrieval
  - XML querying on unreliable structure
- Integration
  - Heterogeneous data sources
- It's ok to get it wrong! (occasionally)



# Entity extraction and noisy XML



- **Static extraction:** \$60 is a price
  - (for a pair!)
- Relationship extraction
  - Microsoft acquires SAP
- **Dynamic extraction**
  - Saturday's game is probably 

--	--	--	--	--
- **Net:** you won't get a clean table;
  - Noisy, incomplete XML

Vendors can do the easy portions well.



# Retrieval engine

- Integrate results from heterogeneous sources
  - No presumption of a common DTD
- Need a semi-soft notion of query-doc proximity
  - *Camembert distance measure*
  - Ideally with a time-tested mathematical foundation
    - Linear algebra, probability, ...



# Information integration



- On-the-fly integration
  - Score/rank aggregation a first step
  - Classic schema normalization
- Query planning and optimization
- Ontology mapping and merging

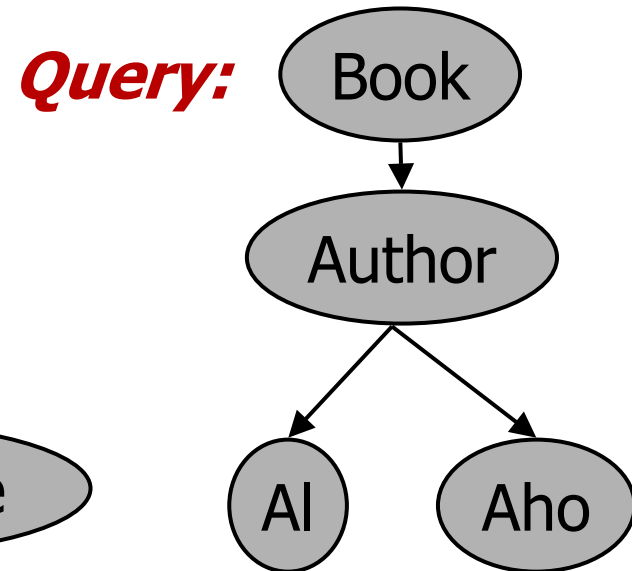
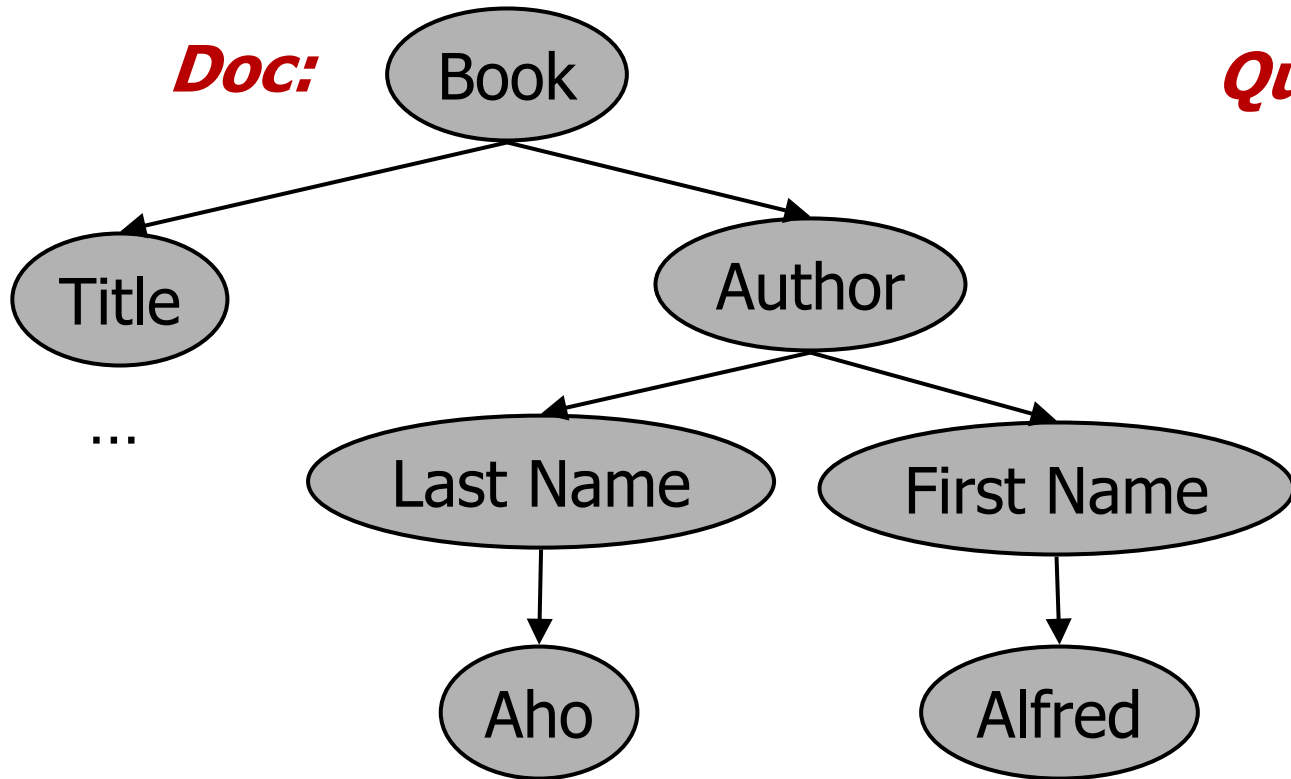
No enterprise vendor really does this well;  
Most stop at federation of results.  
Web companies beginning to.



# Technical Tofu



# Document and query trees





# Vector spaces for XML (w/V. Kakade, Stanford)



- Goal: Distance metric between XML queries and/or documents
- Approach: embed queries, docs in a vector space
- Benefits
  - Camembert distance metric
    - Query-by-example
  - Full power of linear algebra available
    - Dimension reduction/LSI
    - Advanced machine learning
      - Support Vector Machines for Classification
      - Mining – can cluster docs, elements etc.



# Background

- **Schlieder and Meuss:**
  - One axis for each possible sub-tree of any tree
  - Each query/doc becomes a vector in this space
  - Inner product (cosine) similarity measure
  - Exponentially many sub-trees
- **JuruXML (IBM Haifa)**
  - Use only root-leaf paths as axes
  - Enhanced to cope with sub-path matches
    - Book/Author/Aho vs. Book/Author/Last\_Name/Aho
  - Extra computation outside vector space
    - Good for query scoring
    - Loses other benefits commonplace in text retrieval



# Goals of new encoding

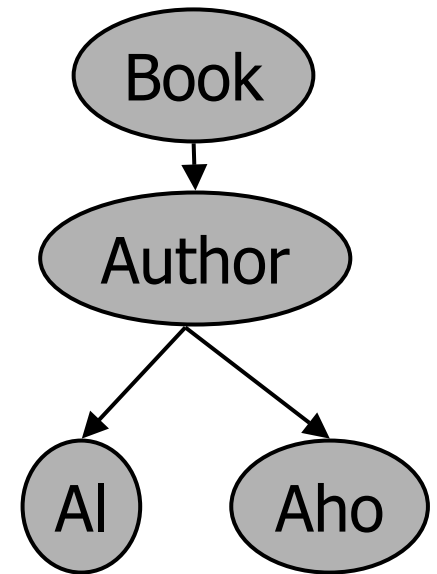


- Avoid exponential blowup
  - Control # axes = index complexity
  - Trade index size for retrieval quality (what's this?)
- **Camembert distance: a pure vector inner product**
  - Gives us LSI, SVM's and other TLA's
- Net – all text retrieval functionality for XML docs



# Main ideas

- A filter  $F$  selects a class of sub-trees
  - E.g.,  $F$  selects all sub-paths of length  $k$
  - Sub-trees thus selected are axes
- For  $k=2$ , this tree generates
  - Book/Author, Author/Al, Author/Aho
- By varying  $F$  we control
  - Size of index (# axes)
  - Quality of matches
- We almost have everything we wanted
  - Camembert distance measure
  - Full power of linear algebra



# What are we missing?

- Doc=Author/Last\_Name/Aho vs. Query=Author/Aho
- How do we capture this in a vector space?
- At index time, compute random substrings of root-leaf paths in docs
  - E.g., Author/Aho is a substring from the Doc
- Do the same for each query
- $\Rightarrow$  Randomized Camembert distance
  - not random docs or queries
  - Analysis possible – what does it say?



# Randomized index



- Reduces Cosine distance computation to a vector space inner product.
- **Analysis: Expected distance grows with similarity between docs/queries.**
  - Gives us the means for search, similar docs, clustering;
  - Use of linear algebra means we can invoke SVM's etc. for document classification/routing.



# Experiments

- INEX 2002 documents
    - 12000+ articles from IEEE transactions
  - Benchmark query suite (“Content Only – CO”)
    - Each query specifies a user need
    - Engine must return a sequence of doc elements
  - Known judgements for each element in the corpus
    - Relevance – how relevant to the user need
    - Coverage – too deep? too shallow?
- ⇒ Composite quality measure for each query/element pair



# Experiments



- Encouraging; perhaps not statistically significantly so.
  - Index blowup, retrieval quality compared
  - Quality can be improved using our techniques, relative to treating each doc as a “bag of words”.
- Classification and clustering “out of the box”
  - On some data, out-perform tailor-made XML clustering/classification algorithms.
- Kakade MS Thesis, Stanford CS Dept, June 2004.
- Much more detailed experimentation is needed.





# Summary



# Research challenges

- Step up on entity extraction – dynamic etc.
- **Retrieval platforms**
  - Unreliable, noisy structure
  - Semi-soft (“camembert”) distance measures
- **Information integration**
- Pilot vertical domain applications
  - Full semantic web too daunting
  - End-users won’t annotate – incentive structure?

