

Categorization & Visualization: The Basics

- Ramana Rao, CTO & SVP
- Inxight Software, Inc

www.ramanarao.com

www.inxight.com

Beyond Search & Browse

- ◆ Search
 - Precise, but brittle ... leaves users searching, not finding ...
- ◆ Browse
 - Robust, but vague ... leaves users wandering & lost, not found ...
- ◆ Opportunity is to blend Browsing & Search
 - Categorization
 - Information Extraction
 - Information Visualization

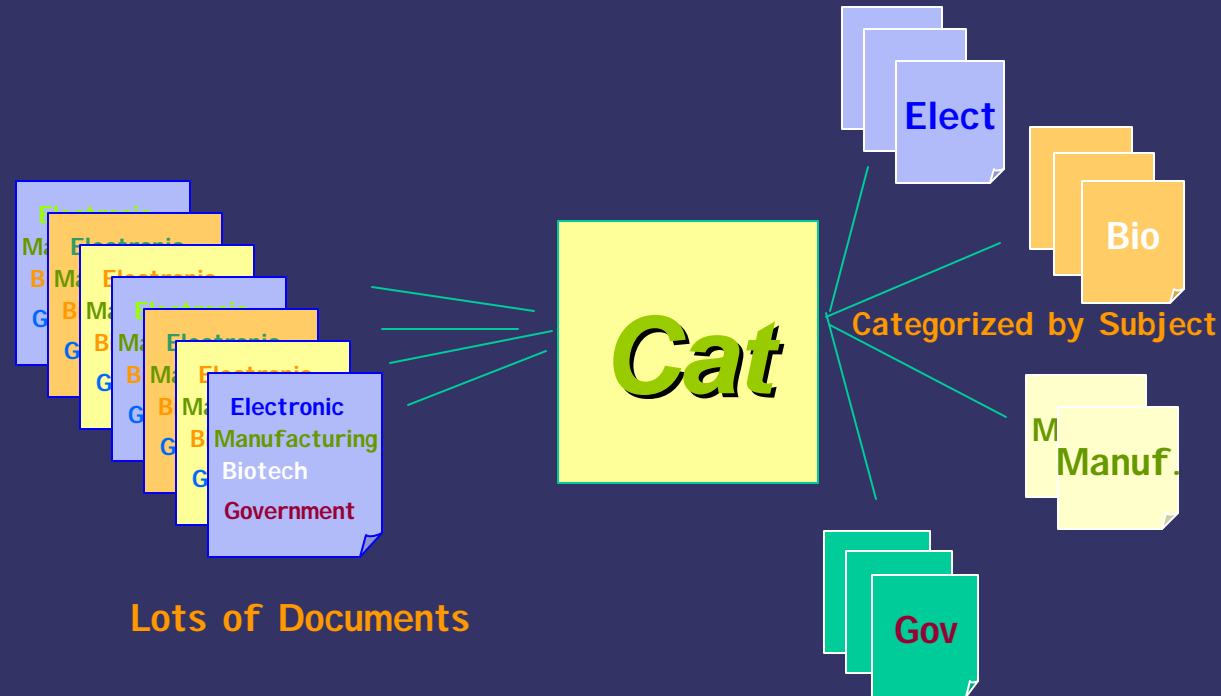
Automatic Categorization

What:

Subject Categorization classifies textual documents into categories based on what they are about.

Why:

Increases Efficiency & Effectiveness of Systems and People that utilize the content



Information Extraction

- ◆ Information extraction is about pulling elements out of documents and collections that guides the more intelligent use of content
- ◆ Often characterized as metadata that provides context
- ◆ Types of metadata include:
 - Noun phrases
 - Named entities (e.g., people, companies, places, products)
 - Key sentences
 - Concepts and topic relationships
 - Similarity between documents, paragraphs and phrases

MetaData from Information Extraction ...

...and Search
Categorization,
Clustering Etc

Summary

- Wall Street is optimistic as Fed cuts rates.
- Stocks Soar with Dow up 130 points.
- NASDAQ gains 2 %.

Similar Docs

- Document 1
- Document 176
- Document 3456

Embedded Entities

- Companies
 - IBM
 - Aventis
 - Goldman Sachs
- People
 - Alan Greenspan
 - George Bush

Optimism that Wall Street is indeed emerging from its slump sent technology stocks higher Thursday, adding to the previous session's triple-digit surge. Blue chips struggled to keep up, fluctuating in light profit-taking. "The market's beginning to buy the scenario that the interest rate cuts by the Federal Reserve are going to help," said Gregory Nie, technical analyst at First Union Securities.

Topical categories

- Financial reports
- FDA Approvals

Embedded Concepts

- "...White House source..."
- "...hot and cold running water..."
- "...20 Gb hard drive..."

Linked Concepts

- "White House source" & "Environmental Policy"
- "20 Gb hard drive" & "Compaq Computer"

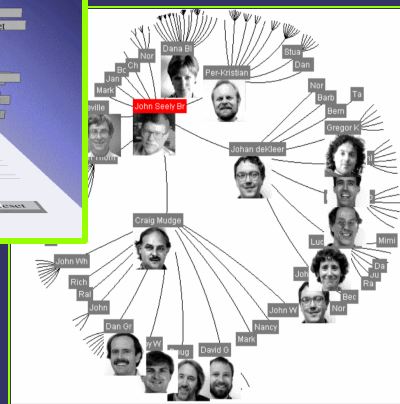
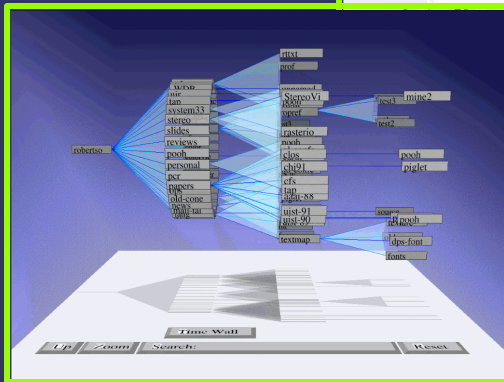
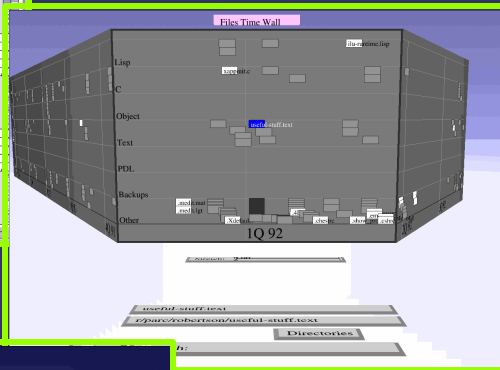
Information Visualization

- ◆ **The Role of Content Visualization**
 - Provides maps of large content spaces ... and also the means for getting to specific documents or items
 - Thus support early or organizing processes like orientation, assessing, survey, etc. ... as well as tune very focused processes like direct walk navigation
- ◆ **Nature of the Solution**
 - Leverage our visual/spatial skills
 - Like browsing, but shows much more, maps not just pages
 - Can eliminate mechanical overheads of browsing
 - Can integrate with searching more tightly
- ◆ **Two key types of Content Visualizations**
 - Content Terrain Maps
 - Wide Widgets

Wide Widgets

Table Lens: Baseball Player Statistics

Player	Year	Team	Pos	Salary	Age	Height	Weight	Bats	Throws	Games	Runs	Hits	Runs Batted In	Home Runs	RBI	Stolen Bases	Caught Stealing	Walks	Strikeouts	Errors	Fielding Percentage
Barry Bonds	2001	Pi	OF	\$32,000,000	32	6'5"	205	R	R	138	117	312	101	37	101	10	1	143	111	1	.989
Andre Braund	2001	Pi	1B	\$18,000,000	30	6'7"	255	L	R	152	124	343	115	27	115	12	0	156	107	0	.987
Jeffrey Lantini	2001	Pi	OF	\$18,000,000	29	6'4"	205	R	R	145	115	305	105	25	105	10	0	140	100	0	.986
Ken Griffey Jr.	2001	Pi	OF	\$18,000,000	28	6'2"	185	R	R	146	115	305	105	25	105	10	0	140	100	0	.986
Billy Beane	2001	Pi	OF	\$18,000,000	31	6'3"	190	R	R	145	115	305	105	25	105	10	0	140	100	0	.986
Bernie Johnson	2001	Pi	OF	\$18,000,000	31	6'7"	215	L	R	145	115	305	105	25	105	10	0	140	100	0	.986
Mike Shaver	2001	Pi	OF	\$18,000,000	31	6'7"	215	L	R	145	115	305	105	25	105	10	0	140	100	0	.986
Billy Beane	2001	Pi	OF	\$18,000,000	31	6'3"	190	R	R	145	115	305	105	25	105	10	0	140	100	0	.986
Paul Hase	2001	Pi	OF	\$18,000,000	31	6'7"	215	L	R	145	115	305	105	25	105	10	0	140	100	0	.986
Stanley Jons	2001	Pi	OF	\$18,000,000	31	6'7"	215	L	R	145	115	305	105	25	105	10	0	140	100	0	.986



- ◆ high bandwidth widgets for interacting w/ large collections
- ◆ arranged on a spine
 - *Hierarchical*
Cone Tree, Spiral Calendar, Hyperbolic Tree Browser
 - *Temporal*
Perspective Wall, Time Lens
 - *Pages* -
Document Lens, Web Books
 - *Calendars* -
Spiral Calendar
 - *Tabular* -
Table Lens, Time Lattice

DEMO

To be continued ...

- ◆ rao@inxight.com
 - Don't hesitate to write ...
- ◆ www.ramanarao.com
 - Papers from talks
 - Information Flow newsletter
- ◆ www.inxight.com
 - White papers
 - Demos & free downloads