

# Midas: Scalable Entity Integration for Unstructured Data Sources

Rajasekar Krishnamurthy IBM Research - Almaden

NFIC 2010

8/14/2010

© 2010 IBM Corporation



- Why Entity Integration for Unstructured Data Sources?
  Challenges in Scalable Entity Integration
- Midas Financial Insights Demo



# Entity View of the World

- Data is prevalent
  - Business Data:
    - Company filings to regulatory bodies
    - Security market (e.g., stock, fund, option) trading data
    - News articles, analyst reports, ...
  - Government Data:
    - US federal government spending data, earmarks data
    - Congress data (voting, members, ...)
- Users and applications prefer an entity view of the underlying data
  - Entities (Companies, People, Securities, ...)
  - Relationships (Employment, Investment, Ownership, ...)
  - Events (Mergers, Acquisitions, Bankruptcy, Appointment, ...)



# Sample questions posed over the entity view

## Questions posed over Business Data

- Which public companies currently share one or more board member?
- Which high-level federal government officials moved between federal government and industry recently?
- How has Berkshire Hathaway's investment profile changed recently?
- How has bank lending to small businesses changed over time?
- Which companies do business together a lot (e.g., banks making joint loans to other large institutions) ?

#### Questions posed over Government Data

- What is spending by each government department for each geographic region ?
- How many (or total value) earmarks in 2009 were solely sponsored by Republican (or Democrat) congress members ?
- Who are the Top k congress members with the most number of earmarks tied to the Department of Defense ?



# Why do we need entity integration ?

- A significant portion of business data is in unstructured format
- Financial service firms use manual methods to analyze regulatory files, news articles etc.
  - Error-prone, cost ineffective, not scalable

Errors made while manually converting corporate actions information into electronic format annually cost financial services firms from \$400 million to \$900 million a year, according to U.K. consulting firm Oxera.

However, a company's responsibility ends once it releases its announcement and makes any required filings. "The issuer is obviously concerned if there is an error in the redistribution of that information, but the fact is, any redistribution is not the issuer's responsibility," Morgan said, adding that company press releases are not regulated so companies won't be required to put them in XBRL.

 $\textbf{Source: http://www.treasuryandrisk.com/News/Pages/Corporate-Actions-Reporting-in-XBRL-Crosshairs.aspx}}$ 

- Users have to "search" data sources to obtain answers
  - Hard to answer questions that need to
    - combine facts from multiple places of the same data source or multiple data sources, or
    - aggregate all data with certain properties
  - For each question asked, manual post-processing of related facts is needed



# Business Data is exchanged frequently in unstructured format

Example Scenario : Corporate Action flow in the U.S. market



**Source**: XBRL Pacific Rim Workshop, 2009 http://xbrl.us/events/Documents/PacificRim/CorporateActions\_Hands.pdf



## Complications in understanding multiple "related" facts

Oracle Corp. (OF	RCL)					On Aug 13: 22.66 4	0.28 (1.22%)	
MORE ON ORCL Quotes Summary Real Time Options	S7 Fide	95 MUNE May ansactions	Scot	trade Online Trades			ERITRADE ve guidance	
Historical Prices	insider in	111340110113					(00)	
Cherts	NET SHARE P	URCHASE ACTI	VITY					
Interactive Regio Charf	Insider Purol	lases - Lost & Mo	ontha					
Basic Toch, Analysis	Descharges		anare	s Trans				
News & Info	Sales		6 862 1	100 1				
Headlines Financial Blogs	Net Shares P (Sold)	urchased	(6,862,1	00) 1	9			
Company Events	Total Insider Shares Held		1.1	168 N/A				
Message Boards Company	% Net Shares Purchased (Sold)		(0.6	96) N	•		Pur	cha
Profile	Net Institution	al Purchases -	Prior Qtr	to Latest Qt				011
Key Statistics				Shares			acro	
Compositora	Net Shares P	urchased		(39,336,80	0)		aure	133
Industry	(Solid) Sh Change In	Institutional Sha	ves Held	/1.285	63	L. L		
Components	Dete	provided by The		andal	~			
Analyst Coverage								
Analyst Opinion Analyst Palmatas	IN SIDER TRA	NSACTIONS RE	PORTED	- LAST TW	D YEAR	8		
Research Reports	Date	Insider		Shares	Type	Transaction	Value*	
Ster Analysts	Aug 4, 2010	SPLAIN MICHA Officer	<u>VEL E</u>	33	Direct	Acquisition (Non Open Market) at \$24.49 per share.	\$808	
Major Holders	Aug 4, 2010	FOWLER JOH	NE	52	Direct	Acquisition (Non Open Market) at \$24.49 per share.	\$1,273	
Insider Transactions	Jul 30, 2010	SPLAIN MICHA Officer	VEL E	466	Direct	Disposition (Non Open Market) at \$23.64 per share.	\$11,016	/_
Pinenciela Income Statement	Jul 30, 2010	FOWLER JOHN	NF	3,638	Direct	Disposition (Non Open Market) at \$23.64 per share.	\$86,002	
Selence Sheet Cesh Flow	Jul 27, 2010	FOWLER JOHN	NE	1,323	Direct	Disposition (Non Open Market) at §24.57 per share.	\$32,50	
	Jul 26, 2010	ROTTLER JUE Officer	RGEN	250,000	Direct	Option Exercise at \$20.73 per share.	\$5,182,500	
	Jul 26, 2010	ROTTLER JUE Officer	RGEN	250,000	Direct	Sale at \$24.55 per share.	\$6,137,500	
	Jul 19, 2010	ROTTLER JUE Officer	RGEN	250,000	Direct	Option Exercise at \$21.04 per share.	\$1,250,000	
	Jul 19, 2010	ROTTLER JUE Officer	RGEN	250,000	Direct	Sale at \$23.59 per share.	\$5,897,500	
	Jul 15, 2010	PHILLIPS JR CHARLES E Officer		3,500,000	Direct	Option Exercise at \$14.57 - \$21.04 per share.	N/A	
	Jul 15, 2010	PHILLIPS JR CHARLES E		3,500,000	Direct	Sale at \$23.68 per share	\$82,880,000	
	Jul 8, 2010	MEISLER LUIZ		3,500	Direct	Purchase at §23.14 per share.	\$80,990	
	Jul 2, 2010	AU YEUNG ST Officer	EVE	13,356	Direct	Statement of Ownership	N/A	
	Jun 29, 2010	CATZ SAFRA Officer		750,000	Direct	Automatic Sale at §21.98 per share.	\$16,485,000	
	Jun 29, 2010	DALEY DORIA	N	32,500	Direct	Option Exercise at \$12.34 per share.	\$401,050	
	Jun 29, 2010	DALEY DORIA	N	32,500	Direct	Automatic Sale at \$22.06 per share.	\$716,950	

Purchase numbers and number of transactions differ across aggregate report and individual transactions list

- Aggregate Data and individual transactions list provided by different data providers!!
- Possible semantic differences on what is a purchase across the data providers
- Understanding multiple facts can be complicated even on a single document !!

Source: http://finance.yahoo.com/q/it?s=ORCL+Insider+Transactions Aug 13, 2010

Data provided by EDGAR Online

8/14/2010



# The Value of Entity Integration



- Obtain an entity view of the world
- Entities, relationships and events are represented as structured objects
- Can answer complex questions over the Integrated Entity Data:
  - Which public companies currently share one or more board member?
  - Which high-level federal government officials moved between federal government and industry recently?
  - How has Berkshire Hathaway's investment profile changed recently?



# **Example: Government and Corporate Positions**

#### Questions we can answer by extracting the employment history of key financial officers:

- What are the (past) government positions held by directors of different companies?
- What is the employment history of key government officers?
- How significant is the interlock between companies receiving TARP funding and government officers?

Citigroup's Proxy Statement filed on April 22 2008

Citigroup Inc.

**Chairman of the Executive Committee** 

 Management Committee — 1980 General Partner — 1971

Joined Goldman, Sachs & Co. — 1966

 Chairman of the Executive Committee, Citigroup Inc. — 1999 to present Chairman of the Board, Citigroup Inc. — November 2007 to December 2007

Co-Senior Partner and Co-Chairman, Goldman, Sachs & Co. — 1990 to 1992

Secretary of the Treasury of the United States — 1995 to 1999

Assistant to the President for Economic Policy — 1993 to 1995

Vice-Chairman and Co-Chief Operating Officer — 1987 to 1990

#### **Technical Challenges:**

- Extraction of current/past positions from biographies and appointment/resignation of special officers
- Entity resolution of person and companies mentioned in the biographies (and other locations)
- Fusion of employment history,
- Selection of government positions from the employment records



#### Examples:

- Robert E. Rubin: Former Secretary of the Treasury and Former officer in both Citigroup and Goldman Sachs.
- Arthur Levitt: Former Chairman of the SEC and officer in AIG.

AIG's Proxy Statement filed on April 05 2006

In July 2005 Arthur Levitt, a former Chairman of the Securities and Exchange Commission named a special advisor to the Board of Directors and the Committee. Mr. Levitt has worked closely with the Committee advising on corporate governance initiatives and possible independent director nominees.



8/14/2010

Robert E. Rubin

## IBM®

## How to Bridge the Gap? Multiple Raw Unstructured Datasets → Consolidated Entities?

## Core Technology Requirements for Understanding Unstructured Data

#### Information Extraction

- Information is present in multiple formats (e.g., Text, XML, HTML)
- Extract entities, events, relationships from unstructured documents
- Unstructured data → Structured data

#### Entity Integration

- Resolving mentions to same real-world entity across filings
- Normalize and cleanse extracted values
- Aggregate related facts extracted from multiple filings

#### Scalable Architecture

- Millions of documents of varying size and format
- New documents arrive daily

Above challenges to handle unstructured data sources are complementary

to issues discussed in earlier talk by Prof. Felix Naumann on "Web Data Integration"



# Why Entity Integration for Unstructured Data Sources ? Challenges in Scalable Entity Integration Midas Financial Insights Demo



# Midas Architecture: A Detailed View





# Use case : Midas Financial Insight Entity Integration Over Financial Data







#### □ Information Extraction





## Example: Extraction of loan information data

#### Extract and cleanse information from headers, tables main content and signatures



Notes: Documents filed by Charles Schwab Corporation On Aug 6, 2009

#### Loan Company Information



## Example: Extraction of person information across documents

Do these filings refer to the same person ?				Who Is James Dimon?				_
variability in the person	n's name			/	i		ommitte	e
- look of a key identifier			posi	tion /	;	\ <mark>m</mark>	embersh	nip
ack of a key identifier				;	histo	ry 🔪		
supporting attributes v	ary depending on the context (form ty	vpe)		/	:		100	
				i i	Compensation	Corporate		
Sincerely,				;	Management	Governance &	Public	Risk
	Director Crandall C. Danias			Audit	Development	Nominating	Responsibility	Policy
$\langle \cap$	Stephen B. Burke			Iviende	Member	Member	\`	
(1)	David M. Cote				1		Member	Member
mus	James S. Crown		,	;	/		Member	Chair
Juryma	ames Dimon		;	ļ	1			
James Dimon	Ellen v. Futter		!	į			Member	Member
Chairman and Chief Executive Off			;	Stock awards	Option awards			- 12
			!	shares of	Number of		C	
			Approval	stock or	underlying	price	Grant aate fair	
	Name	Grant date	date	264 049	options (#)	(\$/Sh)	value (\$)	
	James Dinon	1/22/2008	1/13/2008	504,048	2 222 222 (3)		314,300,000	
		1/22/2008 **	1/15/2008		2,000,000	\$39.830	19,868,000	1
11 31	Michael J. Cavanagh (7)	1/22/2008	1/15/2008	94,151	200.000	20.920	3,/50,000	
1200		1/30/2008	N/A	;	54 271	47.835	360 577	_
Turber 4		1/30/2008	N/A	/	10,625	47.835	91,821	_
	Frank J. Bisignano	1/22/2008	1/15/2008	94,151			3,750,000	
Alle		1/22/2008	1/15/2008	!	300,000	39.830	2,980,200	
		ļ.		į				
		¥		;				
James Dimon, 53, Chairman an	nd Chief Executive Officer of JPMorgan Chase. D	irector since 200	0.	!				
Mr. Dimon became Chairman of	the Board on December 31, 2006, and has been Ch	ief Executive Offic	cer and Pres	ident since D	ecember 31 2	005. He had	d been Presid	ent

and Chief Operating Officer since JPMorgan Chase's merger with Bank One Corporation in July 2004. At Bank One he had been Chief Executive Officer since JPMorgan Chase's merger with Bank One Corporation in July 2004. At Bank One he had been Chairman and Chief Executive Officer since March 2000. Mr. Dimon is a graduate of Tufts University and received an MBA from Harvard Business School. He is a director of The College Fund/UNCF and serves on the Board of Directors of The Federal Reserve Bank of New York, The National Center on Addiction and Substance Abuse, Harvard Business School and Catalyst. He is on the Board of Trustees of New York University School of Medicine.



## Entity Integration





#### **Resolving Person Names: An Example**

- 1 Build an authoritative list of insider names for each company based on insider filings
- 2. Compare extracted name references from other filings to entries in the list and merge data to the closest match

List of insiders for Bank of America & Merrill Lynch

BANKS, KEITH T. BRAMBLE, FRANK P. SR. COI BERT VIRGID GIFFORD, CHARLES K. HAMMONDS, BRUCE L. HANCE, JAMES H. JR. LEWIS, KENNETH D. Chairman of the Board and Chief Executive Officer of Merrill Lynch & Co. MONTAG, THOMAS K. MOYNIHAN, BRIAN T. PRUEHER, JOSEPH W. ROSSOTTI, CHARLES O. SARLES. H. JAY SLOAN, O. TEMPLE JR. TILLMAN. ROBERT L THAIN, JOHN A.



Upon completion of the merger, the board of directors of Bank of America will consist of those directors serving immedia completion of the merger and three directors to be mutually agreed upon by Bank of America and Merrill Lynch from among the peo erving as directors of Merrill Lynch immediately prior to the completion of the merger. It is anticipated that upon completion of the m Mr. Thain will become president of Global Banking, Securities and Wealth Management for the combined company. It is also anticipal

#### However, not all real-world cases are as simple...

Position, principal occupation

business experience and directorships

Name and age

John A. Thain (



## Entity Resolution: Context Is Important



How do we match the partial names with the corresponding correct directors ?

- We need to use additional attributes like position, gender and time period
- We capture these "matching" semantics as "rules". For example,

"IF the names partially match AND the dates of the position match,

THEN link the extracted data to the known director."



## Mapping, Temporal Analysis and Fusion : Creating Person Entities

Basic extracted information

Sample data records for John Thain extracted from various sources

Tror	eactione	hv	lohn	Thai	in
IIdi	1500110115	Dy	JUIII	IIIai	

personName:	Thain John A.
cik:	0001090355
filingDate:	2008-01-24
reportingDate:	2008-01-24
issuer:	BlackRock Inc.
isOfficer:	false
isDirector:	true
fillingType:	3
····	

cik: 0001090355 John A. Thain name: company: Bank of America ... Appointment announcement John Thain personName: cik: 0001090355 appointmentDate: 2008-01-16 filer: BlackRock Inc. appointedAs: Director

8-K



filingType:





## Mapping, Temporal Analysis and Fusion : Computing Current Holdings





# Midas Architecture: A Detailed View





# Handling Scalability

#### Scalability Challenges

- Large document corpora
  - Millions of documents of different formats and document types
  - Documents vary in size (10KB 10MB each)
  - New documents available daily
- Maintaining a complex analysis pipeline
  - Some document types require specialized analysis
  - New analysis needs to be incorporated incrementally
  - Semi-structured results
- Process data updates incrementally
  - Some analysis stages support incremental updates, while other stages may need to run over the entire data.
- Tolerance to errors
  - A failure when processing a document should not be fatal to the overall flow
- Scalable Platform on Cloud Infrastructure
  - Jaql : Declarative language for expressing transformations over semi-structured data
  - SystemT : High-performance declarative rule-based information extraction system



# Why Entity Integration for Unstructured Data Sources ?

Challenges in Scalable Entity Integration

# Midas Financial Insights Demo



# Midas: Financial Insight Scale of current running prototype





# Summary : Research Challenges

- Information Extraction from Text
  - Extracting entities, events, relationships from text and html documents
- Entity Integration
  - Resolving mentions to same real-world entity across filings
  - Normalize and cleanse extracted values
  - Aggregate related facts extracted from multiple filings
- Scalable architecture leveraging Cloud technology
  - Complex analysis over millions of documents in a scalable manner
- Tooling & Programmability
  - Enabling easier definition, deployment and customization of Entity Integration Flows