SystemT: an Algebraic Approach to Declarative Information Extraction

Yunyao Li

Scalable Natural Language Processing IBM Research | Almaden

Outline

- Emerging Applications: Case Studies
- Enterprise Requirements
- SystemT : an Algebraic approach
 - Support Enterprise Information
 - Multilingual Support
 - Machine Learning
 - Tooling
- SystemT Courses
- Summary

Case Study 1: Sentiment Analysis





100M+ consumers, ...

Case Study 2: Machine Analysis



Case Study 2: Machine Analysis



Outline

- Emerging Applications: Case Studies
- Enterprise Requirements
- SystemT : an Algebraic approach
 - Support Enterprise Information
 - Multilingual Support
 - Machine Learning
 - Tooling
- SystemT Courses
- Summary

Text Analytics vs Information Extraction



7

Enterprise Requirements

• Expressivity

- Need to express complex NLP algorithms for a variety of tasks and data sources

Scalability

- Large data volumes, often orders of magnitude larger than classical NLP corpora
 - Social Media: Twitter alone has 500M+ messages / day; 1TB+ per day
 - Financial Data: SEC alone has 20M+ filings, several TBs of data, with documents range from few KBs to few MBs
 - Machine Data: One application server under moderate load at medium logging level → 1GB of logs per day

• Transparency

- Every customer's data and problems are unique in some way
- Need to easily **comprehend**, **debug** and **enhance** extractors

Expressivity Example: Different Kinds of Parses

Natural Language

We are raising our tablet forecast.



Machine Log

Oct 1 04:12:24 9.1.1.3 41865: %PLATFORM_ENV-1-DUAL_PWR: Faulty internal power supply B detected

Time	Oct 1 04:12:24
Host	9.1.1.3
Process	41865
Category	%PLATFORM_ENV-1- DUAL_PWR
Message	Faulty internal power supply B detected

EXPRESSIVITY EXAMPLE: FACT EXTRACTION (TABLES)

OPERATING EXPENSES

PUBLIC UTILITIES BOARD AND ITS SUBSDAMIES

STATEMENTS OF COMPREHENSIVE INCOME

That studed 31 March 2012

	GROUP		BOARD		
-	Note	js Warth 2012 S\$'200	32 Merth 2021 S\$1000	31 Warch 2012 S\$1000	yı Mərch zosa SS'scon
Operating income	1	4/532-549	1,010,737	1,230,760	1,005,434
Operating expenses	5.4	11,037,0960	(998,773)	0.030/6671	(993,502)
pea discard perma-		493	11,968	33	13,825
Non-operating income	. 5	- 101,0050:	19,748	22,001	19,771
Bet Income before Disarched expression and speciality grants		36,493	31.735	37.974	33.695
Financing expenses	- 6	008,0301	0.03/60/8	0.08,030	0.03.608
Hel Lass Arfan contains proto.		(01437)	171,0743	(80,056)	(75,819)
Operating grants from government	-18.1	119.035	185.218	189/615	185,718
Met bicome after granty and before contribution to government, consolidated fund and taxation		112,498	1131347	118,979	03.305
Contribution to government canacilitated fund and taxation	7	[20,230]	110,2602	[20,231]	bace and
Matt Internet office grants and office condellation to prevenent. Associations fluid and to eather	1.1.1	\$7,768	94,073	98,748	94.054
Other comprehensive income					÷
Total contemportement income for the year		37,365	\$4,973	98,748	us ma
Attributation to:	20.3	11.05	64,073	98.748	95,055

Identify line item for Operating expenses from Income statement (financial table in pdf document)



Singapore 2012 Annual Report (136 pages PDF)

CROUR

Identify note breaking down Operating expenses line item, and extract opex components

	GROUP		BOARD		
	Note	31 March 2012 S\$'000	11 March 2011 S\$'eee	31 March 2012 S\$'000	31 March 2011 S\$'000
Direct operating expenses					
- electricity		147,427	126,539	147,427	126,539
- manpower		177,901	185,272	177,852	185,128
 depreciation 		264,431	254,436	264,431	253,753
- plant rental		10,071	24,801	10,071	24,801
 property tax 		15,014	14,365	15,014	14,365
 maintenance and others 	4.1	293,002	266,880	286,642	262,436
Indirect operating expenses					
 service departments' costs 	4.2	129,210	126,480	129,210	126,480
0- (M)	4-3	1,037,056	998,773	1,030,647	993,502

Expressivity Example: Sentiment Analysis

	Intel's 2013 capex is elevated at 23% of sales, above average of 16%
	IBM announced 4Q2012 earnings of \$5.13 per share, compared with 4Q2011 earnings of \$4.62 per share, an increase of 11 percent
	We continue to rate shares of MSFT neutral.
	FHLMC reported \$4.4bn net loss and requested \$6bn in capital from Treasury.
Analyst Research Reports	Sell EUR/CHF at market for a decline to 1.31000
	Not a pleasant client experience. Please fix ASAP.
	I'm still hearing from clients that Company A's website is better.
Customer Surveys	X fixing something that wasn't broken
	Makin chicken fries at home bc everyone sucks!
	Bank X got me ****ed up today!
	Mcdonalds mcnuggets are fake as shit but they so delicious.
- (+	You are never too old for Disney movies.
🔍 Social Media 🛛 🧲	We should do something cool like go to <mark>Z</mark> (kidding).

TRANSPARENCY EXAMPLE: NEED TO INCORPORATE IN-SITU DATA

FHLMC reported \$4 4bn net loss and requested \$6bn in capital from Treasury
Entity of interest
Good or had?
Intel's 2018 capey is elevated at 23% of sales above evenage of 16%
The s 2013 capex is elevated at 25% of sales, above average of 10%
I'm still hearing from clients that Merrill's website is better.
Customer or
I need to go back to Walmart, Toys R Us has the same competitor?
toy \$10 cheaper!

Outline

- Emerging Applications: Case Studies
- Enterprise Requirements
- SystemT : an Algebraic approach
 - Support Enterprise Information
 - Multilingual Support
 - Machine Learning
 - Tooling
- SystemT Courses
- Summary

SystemT: IBM Research project started in 2006



Fundamental Results & Theorems

- **Expressivity**: The class of extraction tasks expressible in AQL is a strict superset of that expressible through cascaded regular automata.
- **Performance:** For any acyclic token-based finite state transducer *T*, there exists an operator graph *G* such that evaluating *T* and *G* has the same computational complexity.

EXPRESSIVITY: RICH SET OF OPERATORS



EXPRESSIVITY: AQL VS. CUSTOM JAVA CODE



Scalability: Class of Optimizations in SystemT

- Rewrite-based: rewrite algebraic operator graph
 - Shared Dictionary Matching
 - Shared Regular Expression Evaluation
 - On-demand tokenization
- Cost-based: relies on novel selectivity estimation for text-specific operators
 - Standard transformations
 - E.g., push down selections
 - Restricted Span Evaluation
 - Evaluate expensive operators on restricted regions of the document





Performance Comparison (with ANNIE)

Task:Named EntityDataset : Different document collections from the Enron corpus obtainedby randomly sampling 1000 documents for each size



[Chiticariu et al., ACL'10] ¹⁸

Performance Comparison on Larger Documents

Datasets : Web crawl and filings from the Securities and Exchanges Commission (SEC)

Date	aset	set Document Size		Throughput (KB/sec)		Average Memory (MB)	
		Range	Average	ANNIE	SystemT	ANNIE	SystemT
Web	Crawl	68 B – <u>388 KB</u>	8.8 KB	42.8	498.8	201.8	77.2
Med	Theorem: For any acyclic token-based FST T,						
SEC I La	tł	nere exists an	operato	perator graph G such that evaluating			
SEC	T and G has the same computational complexity						



Extraction Task: Named-entity extraction

Systems compared: SystemT (customized) vs. [Florian et al.' 03] [Minkov et al.' 05]

Dataset	Entity Type	System	Precision	Recall	F-measure
		SystemT	93.11	91.61	92.35
	Location	Florian	90.59	91.73	91.15
CoNLL 2003	Organization	SystemT	92.25	85.31	88.65
	Organization	Florian	85.93	83.44	84.67
	Democra	SystemT	96.32	92.39	94.32
	Person	Florian	92.49	95.24	93.85
Гакар	Devee	SystemT	87.27	81.82	84.46
Enron	Person	Minkov	81.1	74.9	77.9

Transparency without machine learning outperforms machine learning without transparency.

[Chiticariu et al., EMNLP' 10]

Outline

- Emerging Applications: Case Studies
- Enterprise Requirements
- SystemT : an Algebraic approach
 - Support Enterprise Information
 - Multilingual Support
 - Machine Learning
 - Tooling
- SystemT Courses
- Summary

Tokenization	 26+ languages All other languages supported via whitespace/punctuation tokenizer
Part of Speech and Lemmatization	 18+ languages, including English and Western languages, Arabic, Russian, CJK
Semantic Role Labeling	 6+ major languages Ongoing work on Multilingual SRL [Akbik et al, ACL'15]
Annotator Libraries	 Out-of-the-box customizable libraries for multiple applications and data sources in multiple languages

Outline

- Emerging Applications: Case Studies
- Enterprise Requirements
- SystemT : an Algebraic approach
 - Support Enterprise Information
 - Multilingual Support
 - Machine Learning
 - Tooling
- SystemT Courses
- Summary

Machine Learning in SystemT

- AQL provides a foundation of transparency
- Next step: Add machine learning without losing transparency
- Major machine learning efforts:
 - Embeddable Models in AQL
 - Learning using AQL as target language

Machine Learning in SystemT

- AQL provides a foundation of transparency
- Next step: Add machine learning without losing transparency
- Major machine learning efforts:
 - Embeddable Models in AQL
 - Model Training and Scoring \rightarrow integration with SystemML
 - Deep Parsing & Semantic Role Labeling [Akbik et al., ACL'15]
 - Text Normalization [Zhang et al. ACL'13, Li & Baldwin, NAACL '15]
 - Learning using AQL as target language

SIMPLIFY TRAINING AND APPLYING STATISTICAL PARSERS



SystemML in a Nutshell

- Provides a language for data scientists to implement machine learning algorithms
 - Declarative, high-level language with R-like syntax (also Python)
 - Also comes with approx. 20 algorithms pre-implemented
- Compiles execution plans ranging from single node (scale up multi threaded) to scale out (MapReduce, Spark)
 - Cost-based optimizer to generate execution plans, parallelize
 - Based on data and system characteristics
 - Operators for in-memory single node and cluster execution
- Runs in embeddable, standalone, and cluster mode
- Supports various APIs
- Apache SystemML Incubator project: <u>http://systemml.apache.org</u>
- Ongoing research effort at IBM Research Almaden

Machine Learning in SystemT

- AQL provides a foundation of transparency
- Next step: Add machine learning without losing transparency
- Major machine learning efforts:
 - Embeddable Models in AQL
 - Learning using AQL as target language
 - Low-level features: regular expressions [Li et al., EMNLP'08], dictionaries [Li et al., CIKM'11, Roy et al., SIGMOD'13]
 - Rule refinement [Liu et al., VLDB'10]
 - Rule induction [Nagesh et al., EMNLP'12]
 - Ongoing research

Outline

- Emerging Applications: Case Studies
- Enterprise Requirements
- SystemT : an Algebraic approach
 - Support Enterprise Information
 - Multilingual Support
 - Machine Learning
 - Tooling
- SystemT Courses
- Summary

TOOLING RESEARCH FOR PRODUCTIVITY





[Chiticariu et al., SIGMOD '11, Li et al. ACL'12]

Web Tools Overview

Projects	Extractors	Research Education History	🗷 🖪 🔤 🚟 🛛	K () V () -	Documents 🧄 🗙 🝷 📄 🗐 🧵 😹
 Type Generi Namer Financ Parta c Machie Sentim Sentim Sentim tauser 	e a string and press Enter ic d Entity Recognition se Actors of Speech ne Data Analytics nent Analysis - Surveys nent Analysis - General	Union 1 Education History degree 1-4 Major or Re 1-4 Institution Education History 2 degree 1-4 Institution		Higher Educa	Dan_Jurafsky.txt Dan Jurafsky is Professor and Chair of Linguistics and Professor of Computer Science at Stantistic University. He is the recipient of a 2002 MacArthur Fellowship, is the co- author with Jim Martin of the widely-used taxtbook "Speech and Language Processing", and co-created with Chris Manning one of the first massively open online courses, Stanford's course in Natural Language Processing. His new trade book "The Language of Pood: A Linguist Reads the Menu" just came out on September 15, 2014.
		Extractor Properties. Select an extractor or structure and format your output into columns. Learn	more.	General Settings Output	Dan received a B.A in Linguistics in 1983 and a Ph.D. in Computer Science in 1992 from the University of California at Berkeley, was a postdoc 1992-1995 at the International
		💠 • Education History • degree •	Major + Institution +		Diversity of Colorado, Bouldar until moving to Stanford in 2003.
		4 Span Span	Span Span		His research ranges widely across computational linguistics; special interests include natural language understanding, machine translation, spoken language and conversation, the relationship between human and machine processing, and
		Filters 🛛 🔶 New Filter 🦳 Manage overlapping	matches Output column: Education His	tory • Method: Contained Within •	the application of natural language processing to the social and behavioral sciences. He also works on the linguistics of food and the linguistics of Chinese. Dan was born in New York and grew up in California. He lives with his wife Janet in the Bernal Heights neighborhood of San Francisco.
			•		
		Results decision History (5) Education History 2 (3) Institution (30)	Major or Research Areas (21)	nion 1 (8) degree (13)	
		Document Education History (Span) degree	(Span) Major (Span)	Institution (Span)	
		Chuck, Fillmore.bt Ph.D. in 1961 from the Ph.D. University of Michigan		Univers Michigan	Ease of
		turafsky.txt Ph.D. in Computer Science Ph.D. in 1992 from the University of California	Computer Science	Univer	Programming
	Ease of	Concept catalog: share concepts		Canvas: Visual const Customization of exis	truction of extractors, ting extractors
Ę	Sharing	Project: share extractor development		Result Viewer: visua	lize/compare/evaluate
				[Li et al	., VLDB' 15]

Outline

- Emerging Applications: Case Studies
- Enterprise Requirements
- SystemT : an Algebraic approach
 - Support Enterprise Information
 - Multilingual Support
 - Machine Learning
 - Tooling
- SystemT Courses
- Summary

System Internal & External Impact



Example SystemT Courses

University of Washington

- LING 575 Declarative Information Extraction
 - Audience: ~15 graduate students in CLMS program (Professional MS in Computational Linguistics)
 - Prerequisites:
 - LING 570
 - LING 572, Java, and Eclipse are a big plus, but not required

• University of Oregon

- CIS607 Natural Language Processing and Information Extraction
 - Audience: ~15 graduate students in CIS program (Computer and Information Science)
 - Prerequisites:
 - None. Basic knowledge of AI will be helpful

• Support from Almaden:

- All lecture materials (adapted by the professors for the classes):
- 1 lecture in person / by Telepresence / Skeype + 1 2 office hours
- Email supports available throughout the course (mainly needed during the first month)

• Coming soon:

SystemT MOOC (slides/video/story board done)

Outcome

University of Washington

- 9 group projects: 8 based on / related to SystemT
 - Extract of Personally Identifiable Information from Email Correspondence
 - Extract sentiment from movie reviews with Extracted Named Entities
 - Apply IE for analyzing learner corpora
 - Extract opinion from customer reviews
 - OrionsBelt: Integrating Sentiment Analysis augmented with Fuzzy Matching into SystemT
 - Extract poetic similes from the body of a poem and then related back to the poem title and author + Extract the poet's lifespan and gender
 - Game of Thrones: Character Relation Extraction
 - Extract family relation + post-processing (de-duplicate + anaphora resolution)
 - DIFFRENT: Differentiating Fiction From REality iN systemT
 - Extract (actor, character) pairs from Wikipedia
 - Extract grammatical constructions that indicate which candidate entities participate in the relations of interest.

University of Oregon

- 7 student projects + paper presentation by other students
 - Extract bios from Wikipedia
 - Extract information about the native habitat for the IUCN list of endangered mammals
 - The global location of the native habitat, typically in relation to a country or region
 - The terrain that the animal typically lives in, including altitude and the type of vegetation or specific setting the animal is found in or near
 - Extract financial information from financial documents (Company name / performance / profit/ loss / action / expenses / revenue / sales, etc.)
 - Extract financial information from news reports
 - Sentiment / Company actions / performance / profit / loss
 - Extract restaurant reviews that critiqued the service of the restaurant and the décor
 - Extract company name and # of patents granted for the company
 - Extract address and phone # for restaurants at a certain location

Summary

A declarative information extraction system with cost-based optimization, high-performance runtime and novel development tooling based on solid theoretical foundation [PODS'13, 14]. Shipping with 10+ IBM products, and installed on 400+ client's sites.



Find Out More about SystemT!





More Information

Documentation

Research Areas

- Computer Science
- Natural Language Processing
- Data Management

Awards

 SystemT
 Join/Edit Group

 Overview
 Publications
 Annotated Publications
 News
 Get SystemT
 Educators
 Demo

We are hiring! Multiple positions available. Email your resume to Laura Chiticariu (first 5 letters of last name {at} us.ibm.com)

Highlights

- State-of-the-art <u>AQL language</u> for expressing NLP algorithms, optimizer and runtime engine for execution at scale, and easy to use user interface (see a demo)
- Publications in top NLP, database systems, hardware and HCI conferences
- Currently taught in <u>multiple universities</u>
- · Winner of multiple IBM Corporate Awards for its contributions to IBM products and clients

Find Out More about SystemT!

https://ibm.biz/BdF4GQ



Τηανκ γου!

• For more information...

Visit our website:
 <u>http://ibm.co/1Cdm1Mj</u>

 Visit BigInsights Knowledge Center: <u>http://ibm.co/1DIouEv</u>

- Learning at your own pace: With the lab materials
- Contact me
 - yunyaoli@us.ibm.com

SystemT Team: Alan Akbik, Laura Chiticariu, Marina Danilevsky, Howard Ho, Rajasekar Krishnamurthy, Yunyao Li, Sriram Raghavan, Frederick Reiss, Shivakumar Vaithyanathan, Huaiyu Zhu