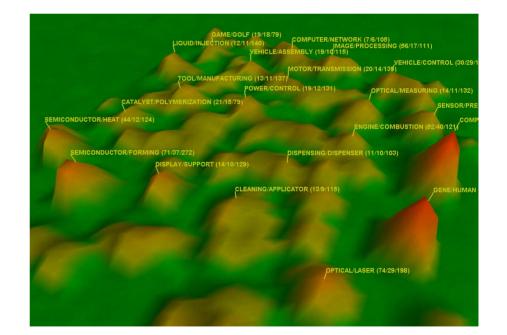
Applications of Data Mining and Information Visualization at SNL

Kevin W. Boyack

Sandia National Laboratories

SIAM International Data Mining Conference April 23, 2004





Sandia is a multiprogram laboratory operated by Sandia Corporation, a Lockheed Martin Company, for the United States Department of Energy under Contract DE-AC04-94AL85000.





- Lots of different data sources
- Lots of different format
- Lots of missing or dirty data
- Different levels of detail, styles, writers, objectives
- Different data types (text, signal, numeric, multimedia)
- Everyone has the similar issues, but different applications







- Lots of tools available
 - Commercial
 - Extraction, indexing, retrieval, portal, personal
 - Visualization
 - Home-grown, particularly in bioinformatics
 - Mostly algorithmic for a single purpose
 - Lots of one-off approaches, no silver bullet
- SNL uses some of our own, some commercial for different projects
 - Use fragmented within company







Outline

- Knowledge domain visualization
 - Overview of process
 - Recent PNAS special issue
- Overview of applications
 - Each has different gotchas







Domain Visualization

- Domain
 - Scientific or technical field
 - Corporate database (projects, requirements, IP ...)
 - Both of the above imply a library of documents
 - Genome
 - Network (social, email, protein ...)
- Visualization
 - Organized visual picture of the domain with both context and content
 - Best ones support interaction and retrieval





Why Use Domain Visualization?

- Reduces information overload
 - Ensure better coverage of related material
 - Less refined initial searches are sufficient
 - Takes advantage of human cognition and visual perception
 - Effectively screen out unrelated material
- Useful for experts and non-experts
 - Expert validate paradigm, easy access, surprises
 - Non-expert quick overview of a field





Literature Domain Analyses Uses

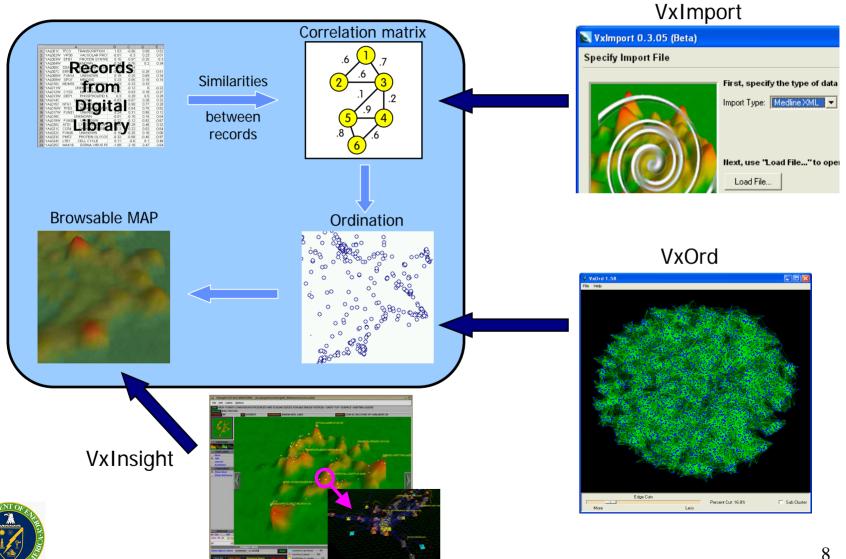
QUESTIONS RELATED TO

		Fields and paradigms	Communities and networks	Research performance or competitive advantage	Commonly used algorithms
UNIT OF ANALYSIS	Authors		Social structure, intellectual structure, some dynamics	Use network characteristics as indicators	Social network packages, MDS, factor analysis, Pathfinder networks
	Documents	Field structure, dynamics, paradigm development		Use field mapping with	Co-citation, co-term, vector space, LSA, PCA, various clustering methods
	Journals	Science structure, dynamics, classification, diffusion between fields			Co-citation, intercitation
	Words		Cognitive structure, dynamics		Vector space, LSA, LDA (20)
	Indicators and metrics			Comparisons of fields, institutions, countries, etc., input-output	Counts, correlations





Example Process with SNL Tools







- Research indicators
- Corporate information
- Bioinformatics, networks







- Models of science
 - Identify recent relevant high impact work
 - Identify possible collaborators
 - Focus new ideas (proposal writing)
 - Evaluate internal work
 - Test theories of scientific progress and innovation

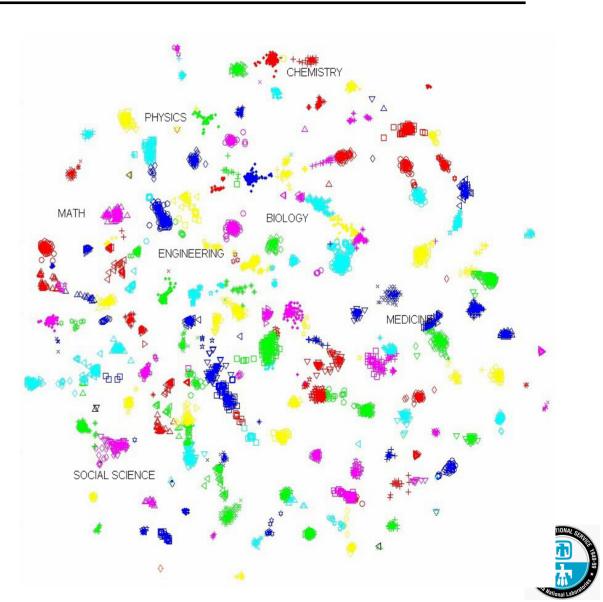




10

Macromodel (Structure) of Science

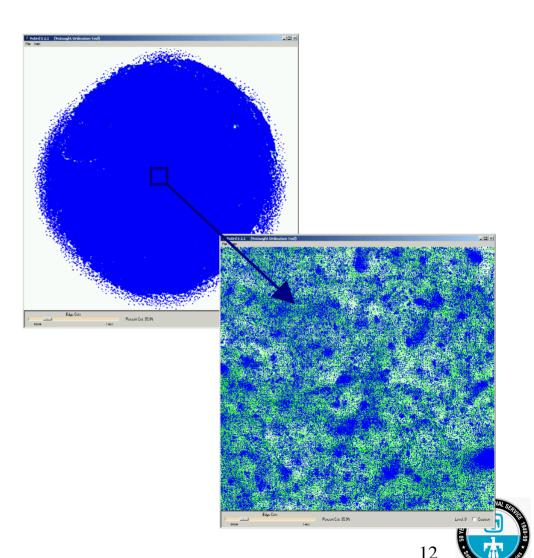
- Map of 7000 journals for the year 2000
- Based on citation statistics
- Validates methods so we can construct micromodel





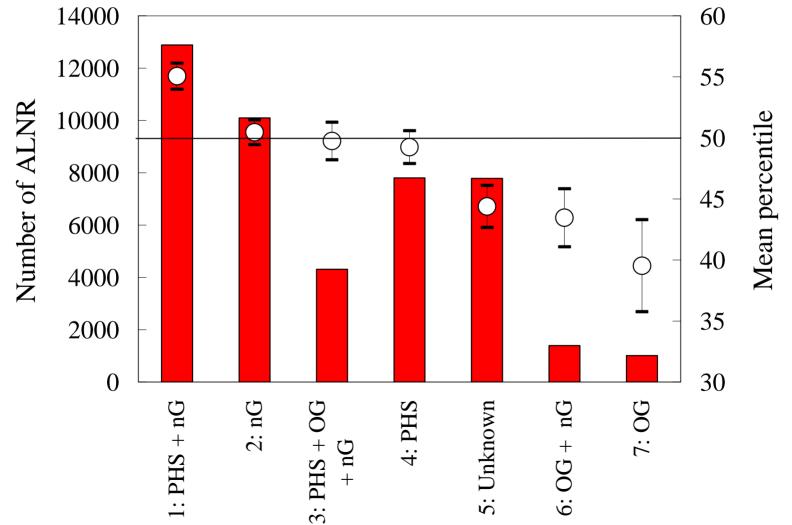
Micromodel of Science

- Models science on scale that allows metric-based assessment to promote innovation
 - 833k nodes (papers),5.6M edges
 - Local structure preserved
 - Develop local indicators

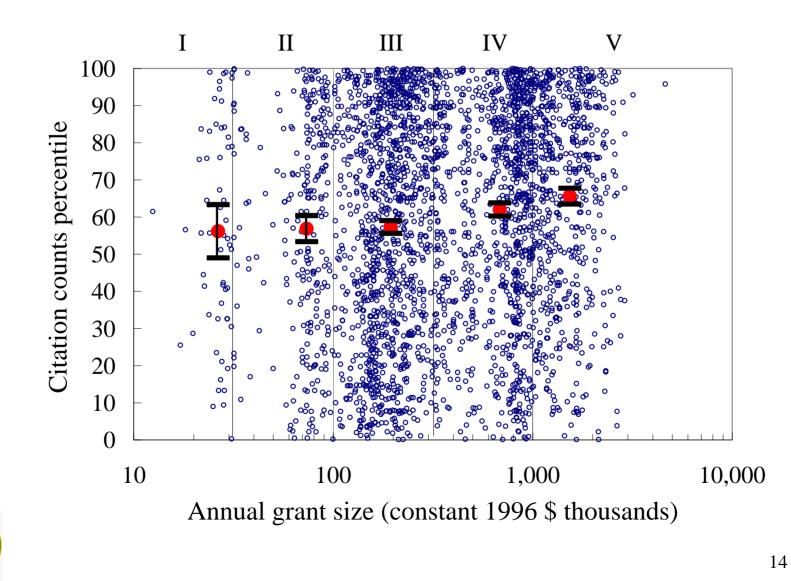








Impact by Grant Amount





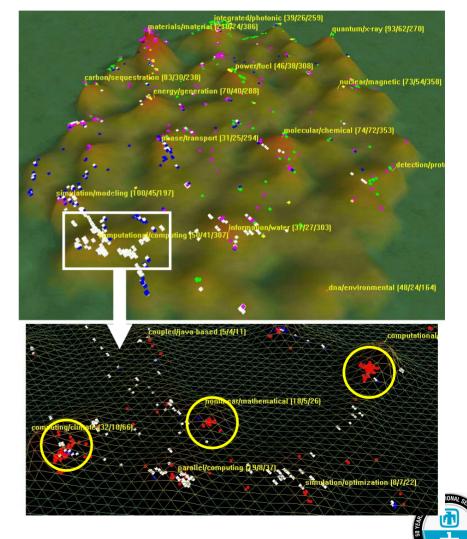
- Internal investment overlaps
- Investment opportunities
- Retrieval
- Parallel computer loganalysis





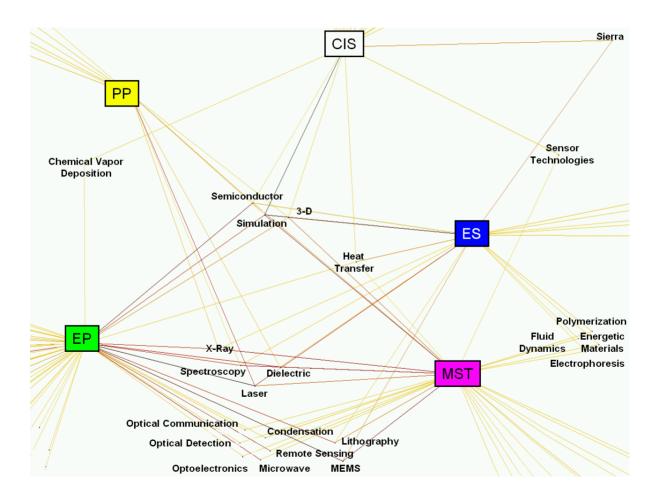
LDRD Investment Area Analysis

- Maps of Sandia LDTD IA's
 - As good as the mental models of IA leads
 - Trends, overlaps between IA's, identifies internal leverage points
- Maps of all DOE LDRD projects
 - Identify potential opportunities for IA's
 - Identify trends at other labs where they found a new opportunity





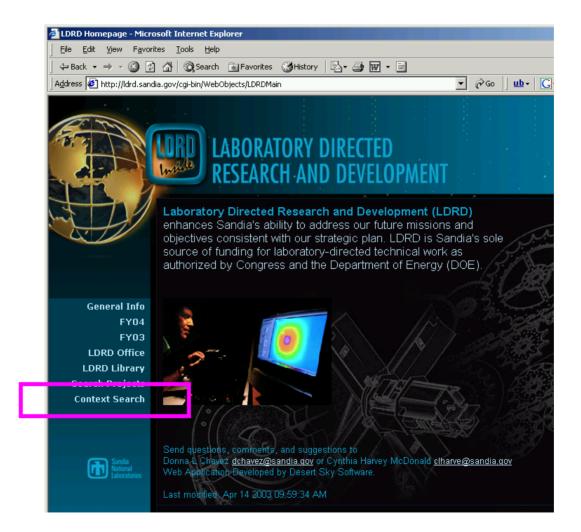
LDRD Investment Area Analysis







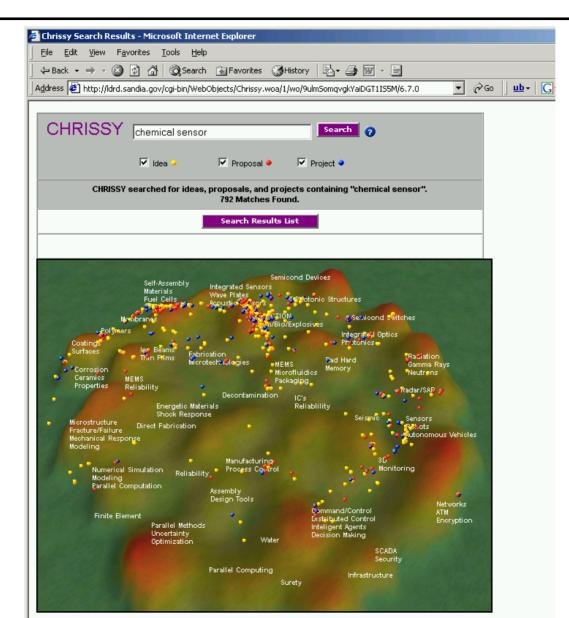
CHRISSY – LDRD Context Search







CHRISSY – LDRD Context Search

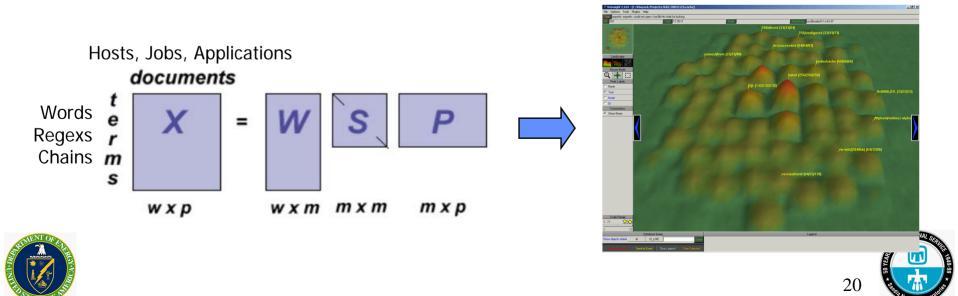






Loganalysis

- Established expert-generated annotation of known error messages
- Established log archive of CPlant and ICC logs, in a format that supports precise and flexible subset selection using SQL
- Using multiple analysis types (Teiresias, LSA)



Loganalysis

- Implemented automatic regex generation using Teiresias
- Successes
 - Cplant TSUNAMI hardware faults on Ross
 - ICC PBS scheduler problems on Liberty
 - LosLobos (UNM AHPCC) version mismatch and portscan

# labe		efinition		
		period		motif
r0	27	0	0	NOCLASS
L137	64	0	1	rte: succeeded
L47	33	1	180	RAM disk driver initialized: 16 RAM disks of 32768K
L48	33	1	180	eth0: Digital DS21143 Tulip rev 65 at 0x8000, * IF
L53	33	1	180	eth1: Digital DS21143 Tulip rev 65 at 0x8800, * IF
L35	32	1	3	Dentry hash table entries: 131072 (order 8, 2048k)
L56	32	1	3	Booting on Tsunami variation Webbrick using machine
L57	32	1	3	IP-Config: eth0: Got DHCP answer from 192.168.37.2
L95	32	1	3	HWRPB cycle frequency (462962962) seems inaccurate -
L122	32	1	3	if=eth0, addr=* mask=255.255.255.0, gw=255.255.255.255,
L125	32	1	3	bootserver=192.168.37.2, rootserver=192.168.37.2, rootpat
L146	32	1	1	rte-init: /cplant/init.d/enfs_client running: mount_nfs
L144	32	0	1	rte-init: 1816 routes read (1816 valid). Max 15, av
L80	31	0	2	init.c:118 Using Alpha PCI_OFFSET_TSUNAMI (0xfffffd0000
L118	28	1	196	Sending DHCP requests, OK
L136	28	1	3	Memory: 1033592k available
L61	26	1	1	rte-init: Found a LANai type 7.2 with 2097152 bytes
L7	14	17	11686	PAM_pwdb: (rsh) session * for user root
L6	14	1	11688	root@admin-0 as
L4	7	4806	16024	connect from
L62	6	3	6	rte-init: Found a LANai type * with 2097152 bytes
L119	4	6	2	Sending DHCP requests * OK
L177	4	2	27	Memory: * available
L131	3	0	0	startup succeeded
L84	1	0	0	Looking up port of RPC 100003/2 on
L86	1	0	0	Looking up port of RPC 100005/1 on
L94	1	0	0	IP-Config: * Got DHCP answer from
L120	1	0	0	Sending DHCP
\triangleleft				
				View lines in original (ungrouped) order
110 1	OF	17 52 07		$-\frac{1}{2}$
L118 N	iov 25	17:53:07	src@nod	<pre>e/if-0.n-4.t-37/if-1.n-0.t-37 Sending DHCP requests, OK e/if-0.n-28.t-37/if-1.n-0.t-37 Sending DHCP requests, OK</pre>
L118 N	iov 25	17:53:10	src@nod	e/if-0.n-32.t-37/if-1.n-0.t-37 Sending DHCP requests, OK
L118 N	iov 25	17:53:10	src@nod	e/if-0.n-7.t-37/if-1.n-0.t-37 Sending DHCP requests, OK
				e/if-0.n-25.t-37/if-1.n-0.t-37 Sending DHCP requests, OK e/if-0.n-11.t-37/if-1.n-0.t-37 Sending DHCP requests, OK
L118 N	iov 25	17:53:14	srcenco	e/if-0.n-14.t-37/if-1.n-0.t-37 Sending DHCP requests, 0K
				e/if-0.n-23.t-37/if-1.n-0.t-37 Sending DHCP requests OK
L118 N	iov 25	17:53:15	src0nod	e/if-0.n-18.t-37/if-1.n-0.t-37 Sending DHCP requests, OK
L118 N	ov 25	17:53:16	src@nod	e/if-0.n-17.t-37/if-1.n-0.t-37 Sending DHCP requests, OK
1118 N	ov 25 oπ 25	17:53:16	srconod	e/if-0.n-22.t-37/if-1.n-0.t-37 Sending DHCP requests, OK e/if-0.n-29.t-37/if-1.n-0.t-37 Sending DHCP requests, OK
L118 N	ov 25	17:53:19	srcenod	e/if-0.n-9.t-37/if-1.n-0.t-37 Sending DHCP requests, OK
L118 N	iov 25	17:53:20	src0nod	e/if-0.n-2.t-37/if-1.n-0.t-37 Sending DHCP requests, OK
T118 M				e/if-0.n-5.t-37/if-1.n-0.t-37 Sending DHCP requests, OK
	iov 25	17:53:21	src@nod	e/if-0.n-8.t-37/if-1.n-0.t-37 Sending DHCP requests, OK
L118 N				
L118 N L119	iov 25	17:52:22	src@nod	e/if-0.n-13.t-37/if-1.n-0.t-37 Sending DHCP requests .?., OK
L118 N L119 L119 N L119 N	iov 25	17:52:28	src@nod	e/if-0.n-13.t-37/if-1.n-0.t-37 Sending DHCP requests .?., OK e/if-0.n-1.t-37/if-1.n-0.t-37 Sending DHCP requests .?., OK e/if-0.n-10.t-37/if-1.n-0.t-37 Sending DHCP requests .?., OK





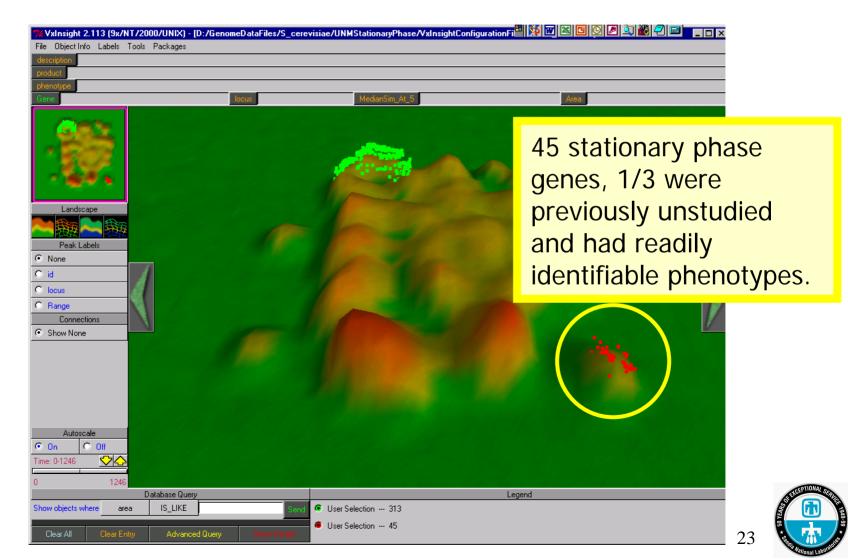
- Yeast
- C. elegans
- Infant leukemia
- How do you enable junior researchers to perform at expert level?





Yeast Stationary Phase Data

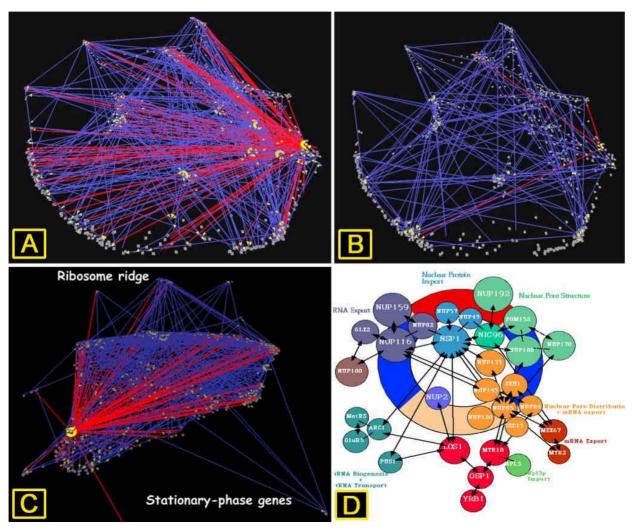
(Werner-Washburne lab)





Yeast Protein Interactions on Gene Map

(Werner-Washburne lab) Genome Research (**12**) 1564-1573, Oct. 2002



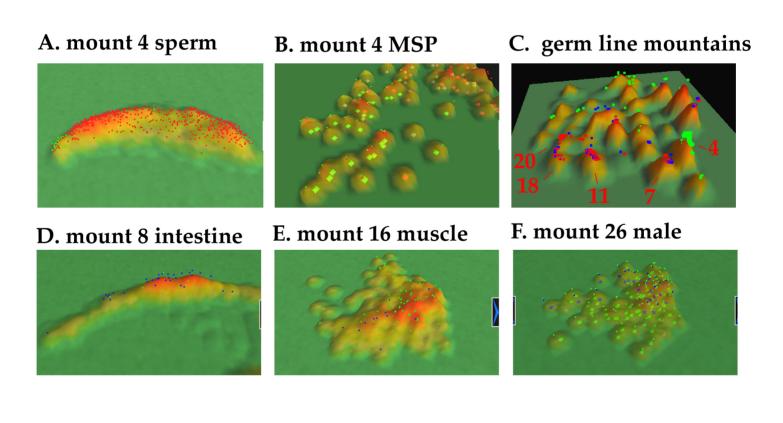




24

Topo Map of *C. elegans* genes

(Stuart Kim lab) Science (**293**) 2087-2092, 14 Sept. 2001



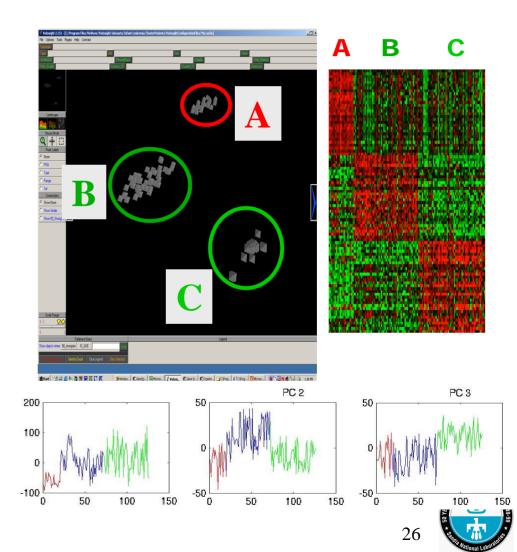




Infant Leukemia – New Results

Ç,

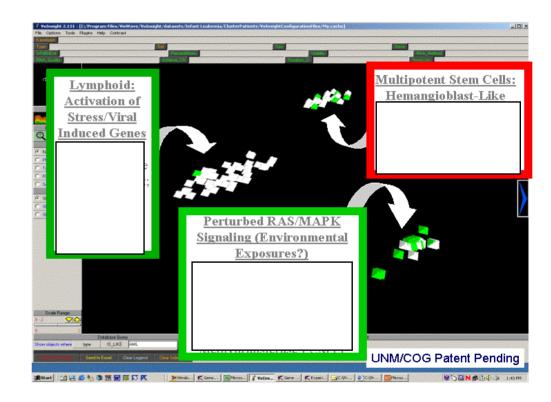
- Vx finds structure in the expression patterns across the patients (A,B,C)
- All of our prediction methods improve when we condition on these patient subgroups
- What is going on biologically in these groups? (back to our fundamental question)





What the Genes Reveal

- <u>A</u> is a new and apparently very significant *stem cell leukemia*
- <u>B</u> is mostly ALL and, possibly, *virally induced*
- <u>C</u> is, possibly, *induced by* environmental exposure

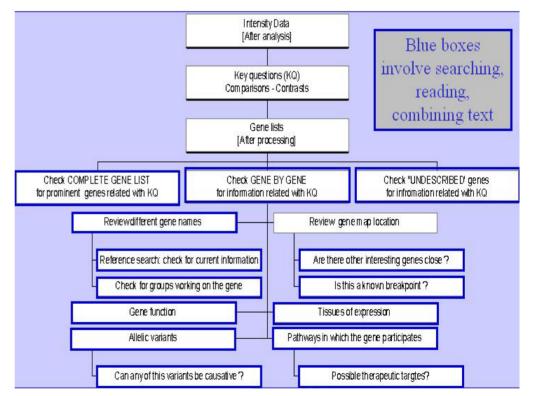






Process to Enable Discovery

- Cluster by genes
- Cluster by arrays
- Contrast between clusters of patients to generate differentiating gene lists
- Used these gene lists as entries into the online databases and back to the original literature
- <u>Read, read, and read some</u> <u>more</u>.







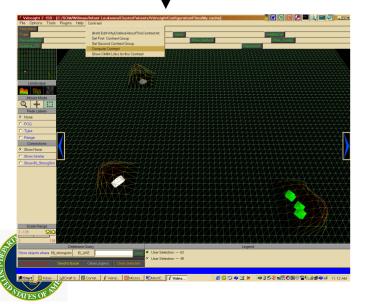
Multi-database knowledge mining

142 patients

12,625 genes

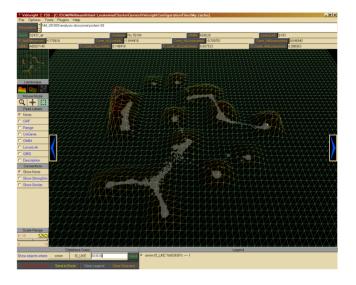
Cluster *by genes*

, Cluster *by arrays*



ANOVA & Contrast

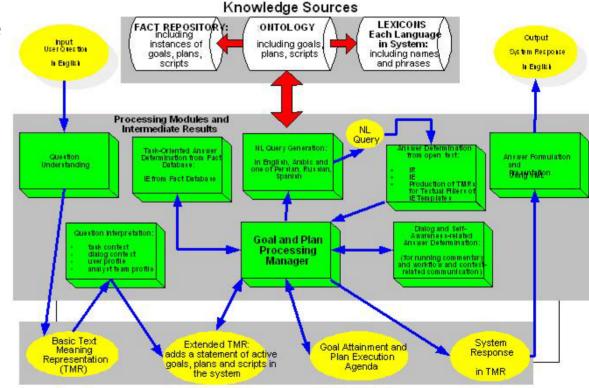
Contrast is Group1 - Group2									
ORF	_	Contrast	F	URL					
3		1.55	149.72						
3		1.40	149.38	Get OMIM details					
3		2.23	138.32	Get OMIM details					
3		1.87	135.59	Get OMIM details					
1		1.60	133.46	Get OMIM details					
3		1.24	126.09	Get OMIM details					
4		1.03	120.06	Get OMIM details					
2		1.08	119.77						
1		2.21	118.18	Get OMIM details					
3		1.36	117.75	Get OMIM details					
4		1.52	115.49						
3	at	2.96	112.13	Get OMIM details					
1		1.19	108.09	Get OMIM details					
3 5050_ a		1.02	107.84	Get OMIM details					





GLEE + +

- Proposing something more intelligent for a *much* larger text corpora.
- Want deeper understanding, and biologically useful reasoning.



System Working Memory







- Scale need to analyze and visualize extreme scale (>1M nodes) networks
- Data extraction, merging of different types and formats, relationship extraction, etc., to generate knowledge
- Understanding from NLP, deep semantics
- Validity





31

Recent Publications

- Boyack, K. W. (2004). Mapping knowledge domains: Characterizing PNAS. *Proceedings of the National Academy of Sciences*, 101(S1), 5192-5199.
- Boyack, K. W., & Börner, K. (2003). Indicator-assisted evaluation and funding of research: Visualizing the influence of grants on the number and quality of research papers. *Journal of the American Society for Information Science and Technology* 54(5), 447.
- Börner, K., Chen, C., & Boyack, K. W. (2003). Visualizing knowledge domains. *Annual Review of Information Science and Technology* 37, 179-255.
- Werner-Washburne, M., Wylie, B., Boyack, K., Fuge, E., Galbraith, J., Fleharty, M., Weber, J., & Davidson, G.S. (2002). Concurrent analysis of multiple genome-scale datasets. *Genome Research* 12(10), 1564-1573.
- Boyack, K. W., Wylie, B. N., & Davidson, G. S. (2002). Domain visualization using VxInsight for science and technology management. *Journal of the American Society for Information Science and Technology*, 53(9), 764-774.
- Davidson, G. S., Wylie, B. N., & Boyack, K. W. (2001). Cluster stability and the use of noise in interpretation of clustering. *Proc. IEEE Information Visualization 2001*, 23-30.
- Kim, S. K., Lund, J., Kiraly, M., Duke, K., Jiang, M., Stuart, J. M., Eizinger, A., Wylie, B. N., & Davidson, G. S. (2001). A Gene Expression Map for Caenorhabditis elegans. *Science*, *293*, 2087-2092.
- Boyack, K.W., Wylie, B.N., Davidson, G.S. & Johnson, D.K., *Analysis of patent databases using VxInsight*. Presented at New Paradigms in Information Visualization and Manipulation 2000, McLean, VA, Nov. 10, 2000.
- Beck, D.F., Boyack, K.W., Bray, O.H. & Siemens, W.D., "Landscapes, Games, and Maps for Technology Planning," *Chemtech* 29(6), 8-16, 1999.



