

Social Informatics Data Grid

Cyberinfrastructure for Collaborative Research in the Neural, Social and Behavioral Sciences

Bennett I. Bertenthal Indiana University bbertent@indiana.edu







Infrastructure for Social and Behavioral Sciences

Goal:

Compare, measure and search for patterns in structured, semistructured, and heterogeneous data sets.

Challenge:

Integrate information over time, place, and types of data

Needs:

- (1) Data interface (shared datasets & databases)
- (2) Service interface (shared tools for analysis)
- (3) Intellectual interface (shared problems & theories)

An Example: What are S1 and S2 talking about?

- S1 you know like those ¿fireworks?
- S2 well if we're trying to drive'em / out her<r>e # we need to put'em up her<r>e
- S1 yeah well what I'm saying is we should*

S2 in front

S1 we should do it* we should make it a lin<n>e through the room<m>s / so that they explode like here then here then here then here



Embodied Cognitive and Social Behavior

- Social communication is not organized by speech alone
- Multiple forms of embodiment – speech, prosody, gesture, eye gaze, posture, facial expression



	desture	Farticipation Framewor
Referential	Matter under	Orientation of Participan
Content	Discussion	
Temporal Scope	Limited Topical	Extended Strips of Talk
	Items	Charles Goodwin

- Opportunity: Evaluate actions by which we construct in concert with each other the cognitive, social and cultural world
- Challenge: Lack tools to assess multiple measures at multiple levels simultaneously and to store and analyze these measures in a common database

Why SIDGrid?

- Data-sharing
 - Understanding complex behavior requires sharing across labs and across disciplines
 - Sciences need an established empirical base
 - Sharing maximizes coverage
- Protocol, coding standardization
- Tool creation
- Multimedia window is crucial
 - Multiple time scales
 - Embodied cognition and linkage
 - Social embedding and linkage

Testbed for Research Multiple Measures at Multiple Time Scales

- Multimodal Communication in Humans and Machines
- Cognitive and Social Neuroscience
- Neurobiology of Social Behavior in Humans and Animals

Primary Objectives

- Develop prototype of core facility for collecting multiple measures of time-synchronized data
- Develop integrated tools for storage, retrieval, annotation, and analyses of multiple data sets at different time scales
- Develop scripts for parallelizing code to run on grid clusters

SuperLab

Objectives of SuperLab

- To expand behavioral and biobehavioral models by simultaneous monitoring of physiological, psychological, and behavioral parameters
- To apply grid computing technologies to study multiple response system in behavioral and biobehavioral research
- To exploit the research potential of grid computing by facilitating "remote" and "virtual" collaborations
- To interface data collection and analysis with the SIDGrid

Sensor and Stimulus Systems

- Stimulus Display
- Audio Stimuli
- Video Recording
- Audio Recording
- Eye Tracking
- Motion Analysis
- IR Thermography
- Physiological Recording
 - electrocardiogram
 electromyogram
 galvanic skin response
 respiration
 blood pressure

•Scene

•Face •Full body

Ambient arrayPersonal boom mic

Multi-projector wall

Transparent display

High density display

- Audio feedback
- Spatialization

•Head mounted •"Remote optics"



•Visualeyez

•512 active markers

•Mid-infrared 320x240 NTSC

, andra

High density electroencephalogram























SuperLab





Project-Sensor Matrix

Features of the SuperLab

- Capable of simultaneous and dynamic monitoring of several channels of behavioral (e.g., audio, video) and physiological (e.g., autonomic activity, EMG, EEG, eye gaze, facial thermography) variables.
- 2. Synchronizing stimuli and with multimodal response systems.
- 3. The ability to monitor variables on at least two individuals simultaneously.
- 4. The ability for collaborators to contribute, observe, and collect data from remote sites connected to the SIDGrid.

Access Grid



ABOUT THE SIDGRid PROJECT

SIDGrid is a computing infrastructure that provides integrated computational resources for the processing of multimodal data in the social sciences.

The Social Informatics Data (SID) Grid enables researchers to collect real-time multimodal behavior at multiple time scales. Multimedia data (voice, video, images, text, numerical) is stored in a distributed data warehouse that employs Web and Grid services to support data storage, access, exploration, annotation, integration, analysis, and mining of individual and combined data sets. The SIDGrid includes the development of a user interface that supports access to the data regardless of the user's geographic location and can be coupled with existing Grid technology.



User Login

Username: *

- ----

Password: *

Log in

Create new account Request new password

USING SIDGrid

> Go to the getting started page to find out more about using SIDGrid and how to obtain a login account for the SIDGrid portal.

ABOUT THE TERAGRID



TeraGrid is an open scientific discovery infrastructure combining leadership class resources at nine partner sites to create an integrated, persistent computational resource.

[TeraGrid website][State of TeraGrid pdf]

SITE ADMINISTRATION

Social Informatics Data Grid



- A general purpose architecture for streaming data applications (e.g., video, audio, time series)
- Built on well established database, multimedia and web and grid services standards
- Time alignment in distributed heterogeneous datasets
 - Software and hardware based
 - Integrated with existing laboratory time stamping and registration techniques

Scalable

- Number of datasets
- Types of data
- Multiple end user applications



Client Side

000

Elan - DemoProject.eaf

File Edit Search View Options Help



000

Elan - DemoProject.eaf

Open from SidGrid Save to SidGrid Save As Save Selection As .eaf Merge Transcriptions Automatic Backup Page Setup Print Preview Print Print Save Save Save Save Save Selection Automatic Backup Print Preview Print Super As Save Save Save Save Save Save Save Save Print Save	New #N Open #O	<u>prions</u>		Grid T	ext Subtitles Controls	
Save \$\frac{\}\ }\\!\!}}}}}} & & & & & & & & & & & & & & & &	Open from SidGrid Save to SidGrid		Volume:			
Save As Save As Template Save Selection As .eaf Automatic Backup Print Preview Print merei Print Preview Print merei Selection: 00:00 Selection: 00:00	Save #S.		100 _			
Save As Template Save Selection As .eaf Merge Transcriptions Automatic Backup Print Preview Print Breview Print	Save As			· ·	50	100
Save Selection As lear Merge Transcriptions Automatic Backup Page Setup Print Preview Print ##P 587 Selection: 00:00 Selection: 00:00 O O O Open 20/files.eaf 66 GB 5 mov 2 wav 368 GB 23 mov 6 wav 21/Laursen.eaf 368 GB 23 mov 6 wav 21 GB 3 mov 12 wav 24/1995.eaf 4 GB 9 mov 1 wav 1 GB 0 mov 2 wav 1 GB 0 mov 2 wav	Save As Template		Ĭ	(50	100
Merge Transcriptions Automatic Backup Page Setup Print Preview Print #Preview Print #Preview Selection: 00:00 Selection: 00:00 Selection: 00:00	Save Selection As .eat			000 Op	en	
Automatic Backup Image Setup Print Preview Image Setup Print Preview Image Setup Print Breview Image Setup Selection: 00:00 Selection: 00:00 Selection: 00:00 Selection: 00:00 Selection: 00:00 Selection: 00:00	Merge Transcriptions		Rate:	20/files.eaf	66 GB 5 mov 2 way	6
Page Setup Print Preview Print #P Selection: 00:00 22/Cassell.eaf 5 GB 1 mov 0 wav 23/Pres2.eaf 4 GB 9 mov 1 wav 24/1995.eaf 4 GB 4 mov 1 wav	Automatic Backup		100 .	21/Laursen.eaf	368 GB 23 mov 6 wav	
Print Preview Print #P 23/Pres2.eaf 23/Pres2.eaf 23/Pres2.eaf 200 Selection: 00:00 25/ISL eaf 1 GB 0 mov 1 wav 200	Page Setup		1	22/Cassell.eaf	5 GB 1 mov 0 wav	200
Print %P 24/1995.eaf 4 GB 4 mov 1 wav Export Ac 587 Selection: 00:00 25 //SL eaf 1 GB 0 mov 2 way	Print Preview			23/Pres2.eaf	21 GB 3 mov 12 wav	200
Export Ac Selection: 00:00 25 /ISL eaf 1 GB 0 mov 2 way	Print %P			24/1995.eaf	$4 \text{ GB} 9 \text{ mov} 1 \text{ wav} \dots$ 4 GB 4 mov 1 wav	
	Export As	587	Selection: 00:00	25/ISL eaf	1 GB 0 mov 2 wav	
SidGrid Transformation + F 1 F F 5 8 8 26/1985 eaf 945 GB 1 mov 66 wav	SidGrid Transformation		DS 8 K	26/1985 eaf	945 GB 1 mov 66 wav	¥.
New from SidGrid 27/1991 eaf 20 CP 12 more 2 more	New from SidGrid			27/1991 eaf	8 GB 3 mov 0 wav	Ţ
				27/1591.cai	20 GB 13 mov 2 wav	
Import	Import	••••••••••••••		1		
Exit) 00:00:24.000 00:00:25.000 00:0	Exit) 00:00:24.000 00:00:	25.000 00:0		OK Can	cel 0:30.00
K-Spch	K-Spch	<u>a dnother</u>	L	L	· · · · · · · · · · · · · · · · · · ·	
W-Spch /n t rhine eh valley yeah that's another eh kind of rotunde and then you for	W-Spch /n t	rhine eh valley	yeah that's ar	nother eh	kind of rotunde	and then you fo
he Rhine eh valley yeah th another eh kin o rotunde an th you foll	he he	Rhine eh valley	yeah th and	other eh	kin o rotunde	an th you foll
w-words	w-words					
W-POS IT n pon int de adj part n p n co ad pro v	W-POS IT	n po n	int de ad	j part	n p n	co ad pro v i
W-IPA to ð ram a væli j£: ðæts anaða a kaind av rotond and ðen ju: fola	W-IPA to a	raın ə vælı	jɛ: ðæts ənað	9 9	kaind av rotund	ənd ðen ju: fola
W-RGU	W-RGU					

Client Side

- Leveraging efforts for annotation and analysis of multimodal data
 - Familiarity and Interoperability
 - Elan (Max Planck Institute for Psycholinguistics, The Netherlands)
 - Talkbank (Carnegie Mellon University, US)
 - Digital Replay System (Nottingham University, UK)
 - XML, Java
 - Cross platform
- Adding SIDGrid functionality to Elan
 - Minimally intrusive
 - Avoid complicated co-development w/ELAN team
 - Browsing SIDGrid data
 - Additional data types
 - Upload / Download to SIDGrid server

Current Status of CHILDES and TalkBank

- 5000 members of CHILDES, 1200 of TalkBank
- Over 2000 articles based on use of CHILDES, 100 based on TalkBank
- 3 million utterances in CHILDES
- 2 million utterances in TalkBank
- TalkBank consists mostly of linked media; CHILDES only has some (mostly transcriptions)

TalkBank Groups, Areas, Topics

- Child Language (CHILDES, PhonBank)
- Conversation Analysis (MOVIN, CA)
- SLA, Bilingualism, LIDES
- Legal (Supreme Court)
- AphasiaBank
- Classroom Discourse
- Linguistic Exploration
- Sociolinguistics (SLX)

Server Side

							contact us
						1116	SIDGrid
	home about p	people	news 8	even	ts partners por	rtal log out	
Welcome hereld					projects vid	eos administra	te
	Search for Keywords	\$	like			Search	
my projects	ALL PROJECTS						
all projects		.mov	.wav	.ea	f GB		NEXT ►
GROUPS:	F 📄 🦳 fi es	10	0	0	45		2005-11-30 17:13:35
talkbank	E 📄 🧰 Cassell	4	30	0	20		2005-11-30 17:13:35
	► Pres2	2	2	1	3		2005-11-30 17:18:35
	1995	12	100	9	200		2005-11-30 17:13:35
	🕨 🖂 🛄 ISL	1	1	1	1		2005-11-30 17:13:35
	► 🖂 🤐 1985	6	2	0	12		2005-11-30 17:13:35
	Þ 🖂 🦳 1991	400	0	1	1001		2005-11-30 17:18:35
	1989	0	666	1	312		2005-11-30 17:18:35
	Þ 🖂 🧰 1992	0	0	13	0.1		2005-11-30 17:13:35
	Þ 📄 🧰 1986	0	0	0	0.0		2005-11-30 17:19:35
		18	4	0	66		2005-11-30 17:13:35

				contact us
	home about pe	ople news & events partr	ners portal log out	SIDGrid
Welcome hereld		ргој	ects videos administra	te
	Search for Keywords	‡) like	Search	
CLICK ON A THUMBNAIL TO	D DOWNLOAD THE CORRE	SPONDING VIDE0		
group5.mov	2010CZ304.mov	dec19f.mov	7004JP405.mov	NEXT ► Solution 8-5.mov
clip18.mov	dec21a.mov	site_administration_issues.mov	2011CZ406.mov	dec19e.mov
	T-1-0-		and the second	3 0 5 5 4

				contact us		
			SID	Grid		
home abo	ut people news & events p	artners poi	tal log out			
Welcome hereId		projects vid	eos administrate	000	SIDGrid Preview	
Search for Key	vords 🗘 like		Search		E.	
projects > DemoProject						
DemoProject		l	created by: susanne updated			
				- 4		
KEYWORDS:						
Actions						
demo yoho testkey funfun						
					the state of the s	
PROJECT FILES:					N TON	
					ja	
elan-example 1 mpg	🔲 elan-e	xample1.wav		4.1	O - 4 P ×	
				Color Annotatio	on Tier	
ANNOTATION TIERS:				K-Spch		
K-Spch	W-Spch	-	W-Words	w-r03		
W-POS	W-IPA	6] W-RGU			
W-RGph	W-RGMe	6	K-RGU			
K-RGph	K-RGMe					
(Preview) (Export Selection) (Export All)	Export range from	to	(Ranges from -1 to 3618	0 milliseconds)		
Process Data (select script and new experiment nam	e)					
New Experiment Name:	Script: recon-all.fre	esurfer 🛟	View Script	(Run Script		

Search and Query (4,000 projects)

- Data Files
 - Names
 - Keywords
 - Attributes (keyword-value)
 - Date
 - Type (Elan, Chat)
- Contents of Files
 - Metadata
 - Tier
 - Annotations

Server Side

- Web services
 - Query
 - Data download / upload
- Portal interface
 - Security
 - Data and metadata browsing
 - Preview
 - Tags, attributes
 - Projects
 - Groups
 - Search
 - Data transformation using grid resources



Science Gateway



What Is The TeraGrid?





				contact us
home about	people news & events	s partners portal	SID(Job id status
Welcome hereld		projects videos	administrate	in queue in queue in queue
Search for Keywor	rds 🛟 like		Search	in queue
projects > waverunner	30 200	ant objects	rgne World Select Pan Font Help	
waverunner	Contra Sound	CI New Read Write Help Figev01a Sound help File Edit Query View Select Spectrum Pitch II	ntensity Formant Pulses Help	DOI IDCOCCO2 IDCOCCO3 IDC r:praat: TR prsat:pra at TR prsat:pra at TR pr
Actions	Rena		Sound editor	
PROJECT FILES:		Praat 4.	1.28 Feb 2004	LDODODO
 ✓ his1.wav ✓ his4.wav 	☑ his2.wav	🗹 hisi	3.wav	TR upload::java
Preview Export Selection Export All	Export range from	to	(Ranges from None to None	e milliseconds)
Process Data (select script and new experiment name) New Experiment Name: WACE2006test2	Script: pitch.p	raat 🗘 View	Script (R	un Script

Grid-enabling Cognitive Neuroscience

- Work done in collaboration with Jed Dobson, Dartmouth Brain Imaging Center
- Test problem: spatially normalize functional MRI volumes prepare studies for analysis
- Goals:
 - Provide significant speedup of analysis runs
 - Reduce the manual effort (and errors) of ad-hoc methods
 - Benefits of provenance: share data and methods with remote researchers
 - Grid needs to be transparent to the users: Literally, "Grid as a Workstation"

fMRI Preprocessing









Workflow courtesy James Dobson, Dartmouth Brain Imaging Center



Pessoa Lab, IU

To create 2D cortical surface renderings from 3D anatomical images, typical modern cpu consumes ~40 hours to process one subject

Typical neuroimaging studies have 15-20 subjects: implies more than **30 days** to get surface renderings!!

Longitudinal studies: hundreds of brain scans!!!



Freesurfer

High resolution anatomical slice



Segmented grey (pink line) and white (green line) matter



View: Inflated surface



View: Pial surface

SIDGrid

Generate cortical surface representations of whole study in 2 days!



Download output files









Modulation of Oscillatory Neural Synchronization



Electroencephalography

- Electromagnetic oscillations picked up by EEG electrodes are created by the somato-dendritic potentials of very large groups of 'open field' neurons.
 - These are neurons arranged in a parallel and structured fashion such that their membrane potentials superimpose
- Signal is filtered by the varying density and thickness of the skull and meninges.
- Raw signal is often analyzed as a superposition of frequency bands:
 - delta (.1-4hz), theta (4-8hz), alpha (8-12hz), beta (12-26hz), and gamma (26 hz+)



Functional Significance of Mu Rhythms (Pineda, 2005, Brain Research Reviews)

- Mu rhythm is an EEG oscillation with dominant frequencies between 8-13 and 15-25 Hz
- Recent studies suggest mu rhythms reflect downstream modulation of motor cortex by prefrontal mirror neurons
 - Mu rhythms represent S-M transformations, specifically "seeing" and "hearing" into "doing"
- Mu power reduced in normal adults by self-initiated movement, imagined movement, and observed movement

Time-Frequency Analysis



Figure 11.7

Three possibilities to extract frequency information from ERP data: two 35–45 Hz filtered ERPs (*left*), two FFT spectra of the ERPs (*middle*) and the wavelet transforms of the ERPs (*right*). Note that only the filtered signal and the wavelet transform still represent changes over time. The FFT spectra show the entire frequency range but no temporal information.



Channel 2 FB

Channel 2 SB



Channel 14 FB





Neural Synchronization at 10Hz

Channel 2, SB v. FB, 10 Hz



Channel 15, SB v. FB, 10 Hz



Channel 14, SB v. FB, 10 Hz



Channel 16, SB v. FB, 10 Hz



Scripts for Running Jobs on Grid

- Matlab (high-level language and interactive environment for peforming computationally intensive tasks)
- R (software environment for statistical computing and graphics)
- **Praat** (software for acoustic analysis)
- Free Surfer (automated tools for reconstruction of the brain's cortical surface from structural MRI data)
- **AFNI** (programs for processing, analyzing, and displaying FMRI data)
- SUMA (adds cortical surface based functional imaging analysis to the AFNI suite of programs)

Advantages of Grid Computing

- Vastly expanded computing and storage
- Reduced effort as needs scale up
- Improved resource utilization, lower costs
- Facilities and models for collaboration
- Sharing of tools, data, and procedures and protocols
- Recording, assessment and reuse of complex tasks



Challenges and Opportunities

	Today	Tomorrow with SID Grid	Milestones
Theories & Models	Static Single cause Linear Component processes Specific time scale Single level of analysis (e.g., neural)	Dynamic Multiple causes Linear and Nonlinear Systems or networks Multiple time scales Multiple levels of analysis (e.g., neural, social)	Empirical tests of theories & models
Collaboration	Single labs Annotations by single investigators Local access only	Community of collaborators Collaborative annotation Remote & distributed access	Collaborative annotation tool
Query and Analysis	Standard statistical analyses Single stream Non-standard formats & coarse alignment of data streams Single location Stand alone application	Automated query, exploration, and analysis Multiple streams Tools to acquire, transform & align multiple data streams Multiple locations Extensible SID Grid application	Query & analysis services
Measurement & Annotation	Single measure Uni-modal Single time scale Manual coding	Multiple measures Multimodal Multiple time scales Automated coding	Multimodal data stream tool
Data Collection	Single investigator populating database on single workstation	Community of collaborators creating SID Grid data resources on grid	SuperLab Legacy datasets

Research Team

University of Chicago:

Bennett Bertenthal:	Professor, Psychological and Brain Sciences, Indiana University
David McNeill:	Professor Emeritus, Department of Psychology
Howard Nusbaum:	Chair, Department of Psychology, co-Director, Center for Cognitive & Social Neuroscience
Jean Decety:	Professor of Psychology, former Director of Research, CNRS, Lyons.
Gina Levow:	Assistant Professor of Computer Science

Argonne National Laboratory:

lan Foster:	Distinguished Service Professor of Computer Science, Head, Distributed Systems Lab, Associate Division
	Director, Mathematics & Computer Science Division; Director, Computation Institute
Mark Hereld:	Experimental Systems Engineer, ANL, Fellow, Computation Institute
Michael Papka:	Acting Deputy Division Director, Mathematics & Computer Science Division, Fellow, Computation Institute
Rick Stevens:	Acting Chief Scientist, ANL, Professor of Computer Science, Director, Computation Institute, Director, Mathematics & Computer Science Division, ANL
Michael Wilde:	Software Architect, Mathematics & Computer Science Division, Fellow, Computation Institute

University of Illinois at Chicago:

Robert Grossman	Professor of Mathematics, Statistics and Computer Science, Director, National Center for Data Mining
Stephen Porges:	Professor of Psychology and Psychiatry, Director, Brain-Body Center
Sue Carter:	Professor of Psychiatry, Physiology, and Biophysics, co-Director, Brain-Body Center

Consultants:

Brian MacWhinney:	Professor of Psychology, CMU, Developer of CHILDES & TalkBank
Francis Quek:	Professor of Computer Science, VPI, Director, Center for Human-Computer Interaction

Programmers

- David Hanley
- Sarah Kenny
- Kavithaa Rajavenkateshwaran
- Thomas D. Uram
- Wenjun Wu

Questions

