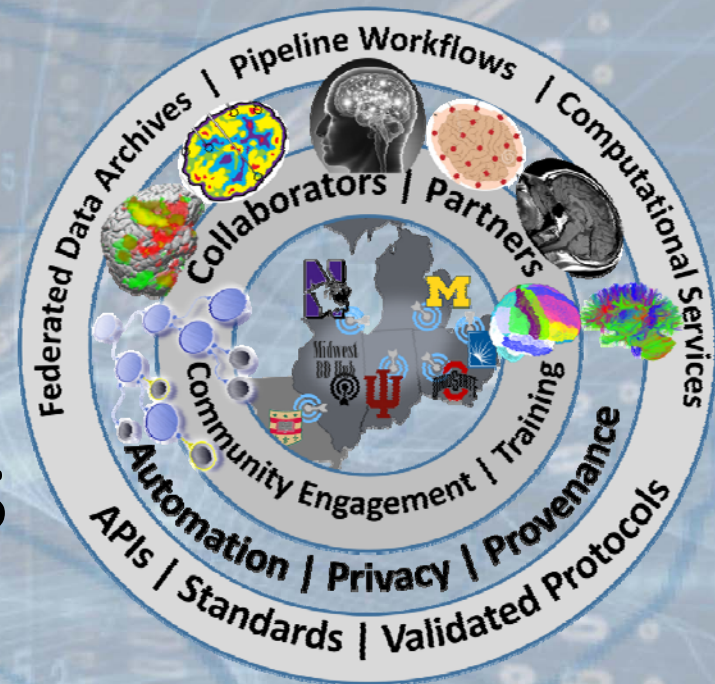


# Midwest Workshop on Big Neuroscience Data, Tools, Protocols & Services



Computational Neuroscience Network (ACNN)

[http://www.NeuroscienceNetwork.org/ACNN Workshop 2016.html](http://www.NeuroscienceNetwork.org/ACNN_Workshop_2016.html)

# Michigan Institute for Data Science

Ivo D Dinov

Statistics Online Computational Resource (SOCR)

Michigan Institute for Data Science (MIDAS)

University of Michigan

<http://www.umich.edu/~dinov>



SCHOOL OF NURSING  
STATISTICS ONLINE  
COMPUTATIONAL RESOURCE (SOCR)  
UNIVERSITY OF MICHIGAN



# Big Neuroscience Data

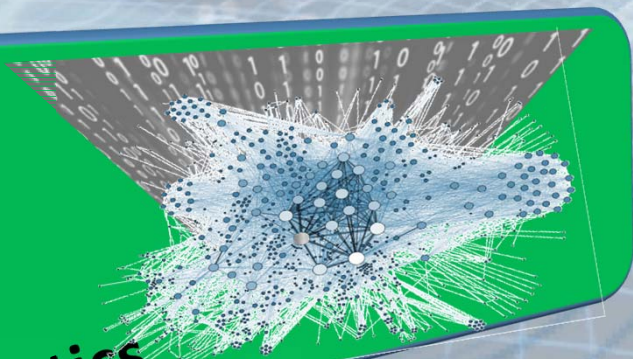
**Data  
Wrangling**



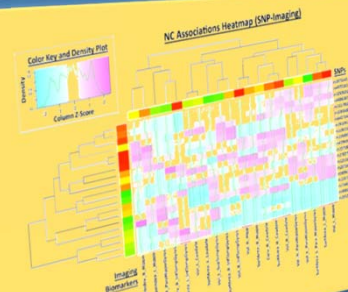
**Data  
Modeling**



**Data  
Analytics**



**Data  
Visualization**



# Big Neuroscience Data

Case-Studies	Sample-size	Data Elements	Description
<b>ALS</b>	Longitudinal data for 8K patients (data comes from 2 independent EU registries)	>300 variables including incomplete clinical, physiological, and cognitive data	This case-study needs innovative machine learning methods for automated diagnostic classification (e.g., patients vs. controls) and unsupervised prediction of cognitive and behavioral decline in ALS patients. The multi-source data will be partitioned into training (estimation) and testing (validation) sets. We plan to fit different models and estimate model-free classifiers to either cluster the participants into groups or hierarchies, or to forecast the progression of the disease over time
<b>Depression / SZ</b>	Extremely large longitudinal data (ms samples), 150 subjects	Study design includes 20 (2*2*5) stimuli types, 12 electrode locations, 1,000's of ms measurements, 4 summary measures	This case-study is focused on identifying a set of P300 biomarkers (using the incomplete high frequency data) that can classify the core cohorts (Bipolar, Schizophrenia, Depression, and healthy controls). We also aim to investigate if P300 responses to emotional stimuli classify the groups better than those to non-emotional stimuli (from standard Go/NoGo tasks). Mean and peak amplitude and latency are candidates based on previously reported results
<b>PD</b>	550 PPMI subjects with 1-10 observation time points	Demographics, clinical tests, vital signs, MDS-UPDRS scores, ADL, MoCA, sMRI, ESS Sleepiness Scale, GDS-15, genetics	Using heterogeneous data of Parkinson's Disease (PD), the study aims to develop a comprehensive end-to-end protocol for data characterization, manipulation, processing, cleaning, analysis and validation. Specifically: (1) introduce methods for rebalancing imbalanced cohorts, (2) utilize a wide spectrum of classification methods to generate consistent and powerful phenotypic predictions, and (iii) generate reproducible machine-learning based classification that enables the reporting of model parameters and diagnostic forecasting based on new data
<b>TBI / Trauma</b>	Over 2,000 patients and controls, acute and multiple chronic times	Dozens of structured (imaging, phenotypic, clinical) and unstructured (injury type, notes) data elements	Three major data sets are included: Volumetric in Brain Trauma (VBT), HeadSmat, and PROTECT II. Each of these datasets include patient's brain CT and/or MRI scans at time of admission, and in some cases during ICU stay, and even during long-term follow up after hospital discharge. These datasets also include time-course data on some physiological measures and blood factors, captured throughout the course of treatment. The UMich/Massey Foundation Grand Challenge provides additional motivation and testing data. The PIs are involved in research funded by this foundation to integrate and analyze these datasets

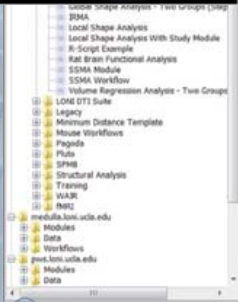
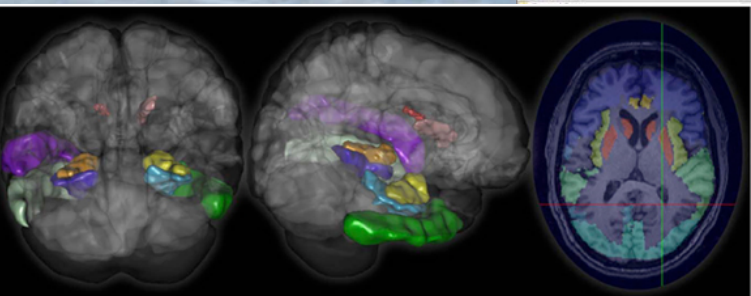
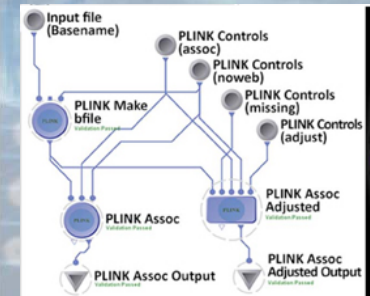
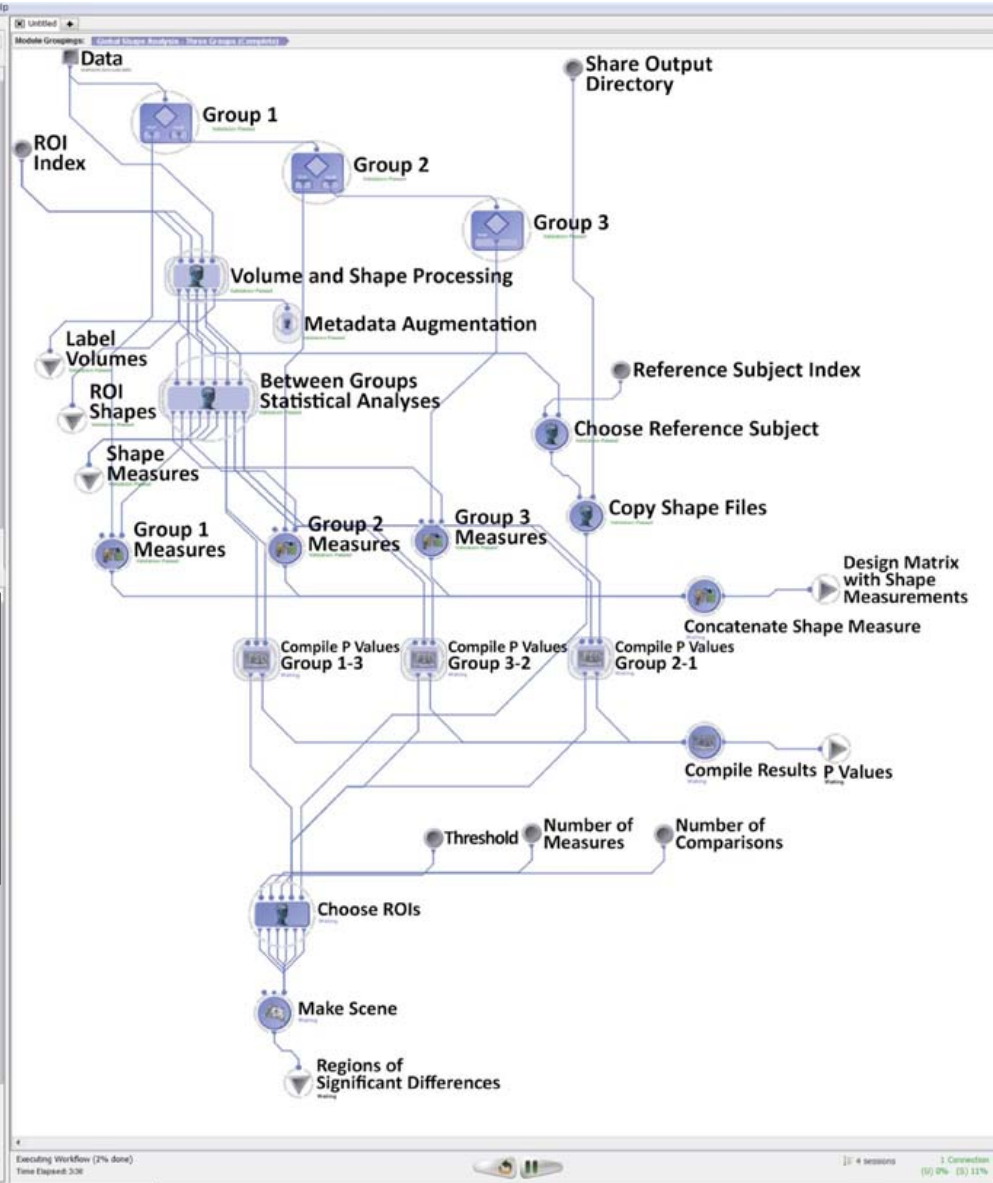
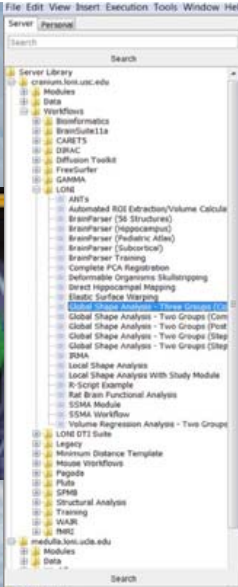


Big Data	Information	Knowledge	Action
Raw Observations	Processed Data	Maps, Models	Actionable Decisions
Data Aggregation	Data Fusion	Causal Inference	Treatment Regimens
Data Scrubbing	Summary Stats	Networks, Analytics	Predictions
Semantic-Mapping	Derived Biomarkers	Linkages, Associations	

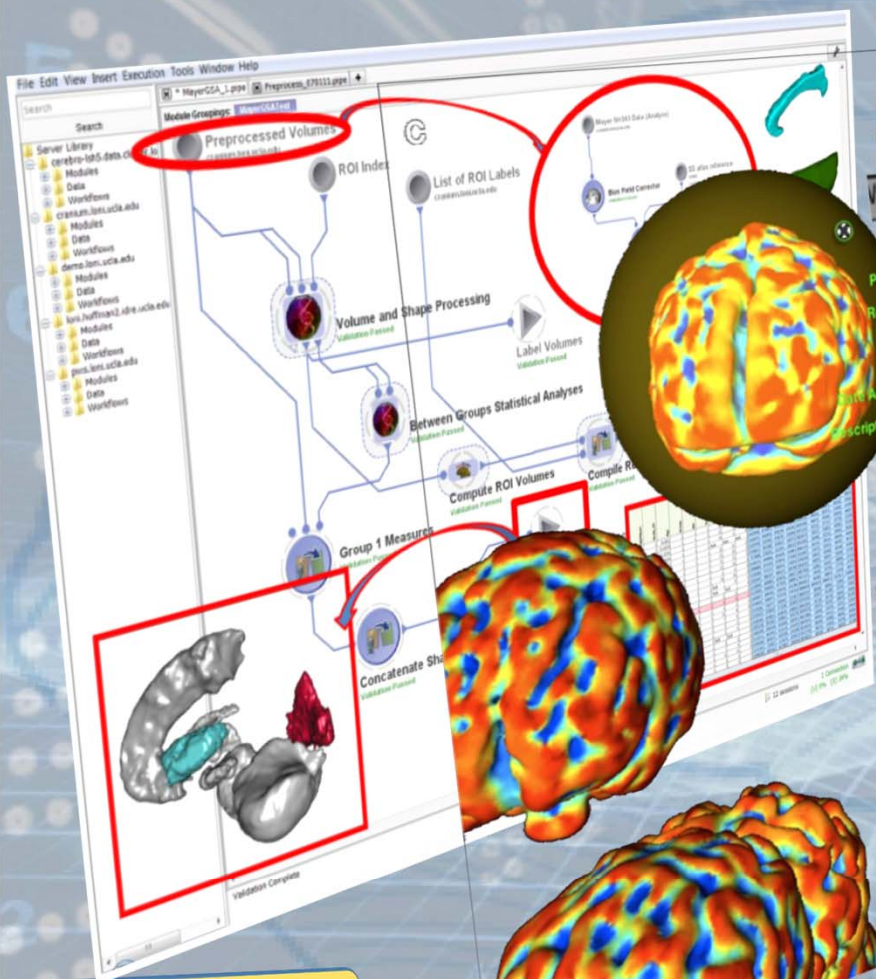


## Normal & Schizophrenia Pediatric Neuroimaging Study

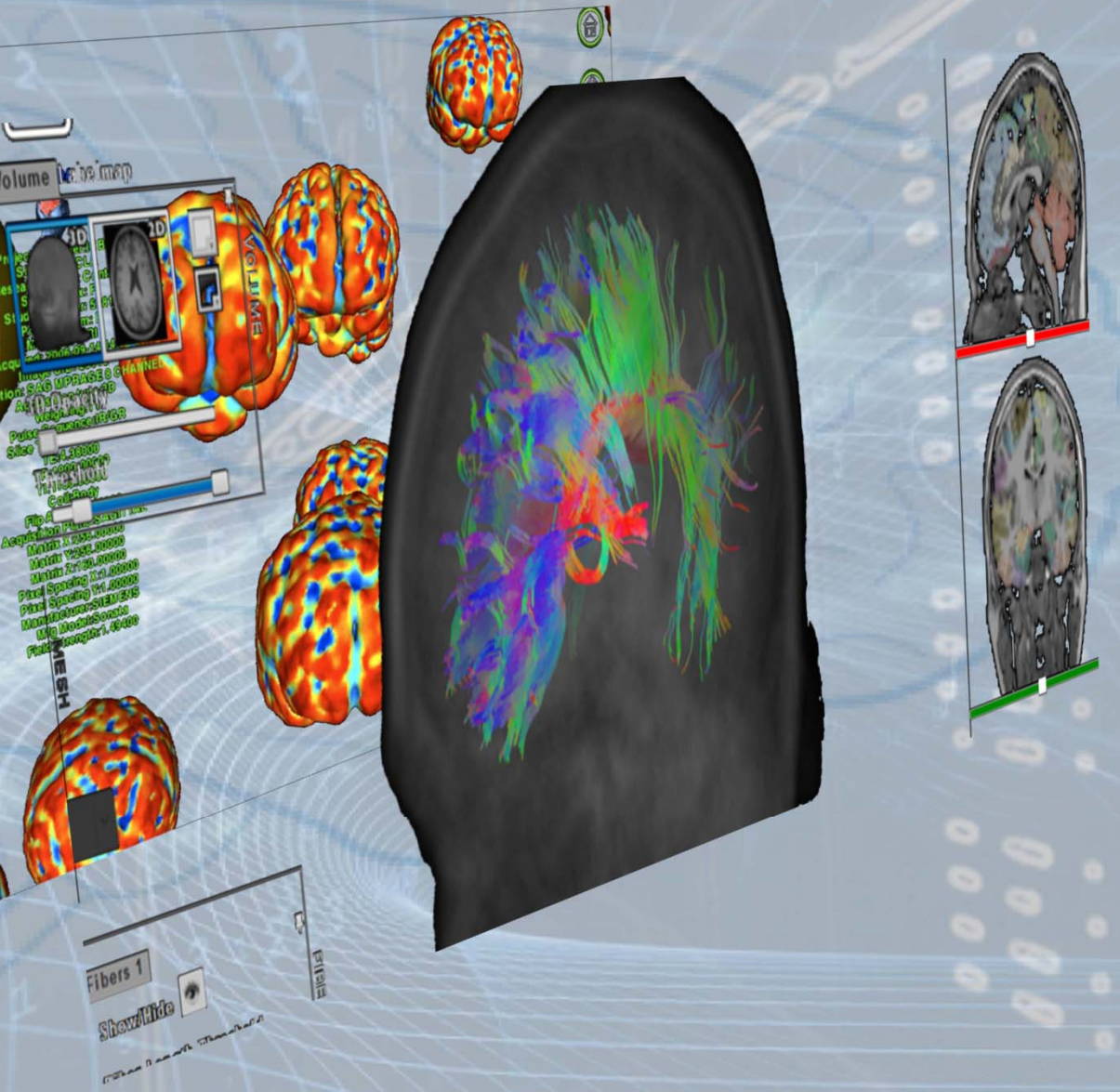
	Sex	FS_IQ	TBV	GMV	WMV	CSF	Background	L_superior frontal gyrus	R_superior frontal gyrus
1	1	118	1663407	654981					
1	1	116	1583940	527918	840874	167553			
1	2	107	1299470	550035	851560	204463	8362251	67276	
1	1	107	1535137	599372	614598	134837	8586107	19866	66079
1	2	107	1431890	593306	780127	155638	8627510	15450	17346
1	1	100	1578698	822056	695296	143288	8607550	15483	14219
1	1	94	1453510	534556	546395	210247	8617899	18538	13942
1	2	114	1650348	608916	730155	188800	8593510	16515	17094
1	1	107	1288971	479707	860876	180555	8624978	17285	15956
1	2	95	1366346	699476	668576	140688	8631722	15062	20009
1	1	109	1326402	438000	479900			15918	15280
1	1	125	1500000	438000	479900			13257	16916
1	1	125	1500000	438000	479900			8627	6539







**Data Visualization**





# Gaps, Barriers & Opportunities

- There is no analytical foundation for systematic representation of Big Data that facilitates the handling of data complexities and at the same time enables joint modeling, information extraction, high-throughput and adaptive scientific inference (cf. CBDA, PMID: PMC479548)
- Kryder's Law  $\gg$  Moore's Law (more data than we can possibly manage with projected increase of computational power) PMID: PMC3933453
- Enormous opportunities for algorithm development, trans-disciplinary data-science training, collaborative research using Big Neuroscience Data
- Advance "Open-Science"

# Big Data Analytics Resourceome

The image displays a comprehensive grid of logos for various Big Data Analytics tools and services, organized into several categories:

- Data Analysis & Platforms:** Includes logos for Hadoop, PARACCEL, Storm, HPC Systems, GridGain, Dremel, Hortonworks, Apache Drill, calpont, ORACLE, HD, Zettaset, and Spark.
- Databases / Data warehousing:** Includes INFOBRIGHT, Cassandra, HBASE, Hibari, riak, Infinispan, Bigdata, OrientDB, Neo4j, HYPERTABLE, socr, STAILISTICS, HIVE, redis, and Globals.
- Workflows:** Includes Pipeline, Galaxy, and tranSMART Foundation.
- Multivalue database:** Includes Rocket, U2-REVELATION, northgate, QM, and jBASE INTERNATIONAL.
- Big Data search:** Includes Lucene, Apache Solr, and elasticsearch.
- Data aggregation:** Includes OQOOP, Ocean, and chubex.
- Big Data to Knowledge (BD2K):** Includes talend, pentaho, python, BIRT, Exchange, ACTUATE, and KNIME.
- Data Mining:** Includes RAPID MINER, orange, RAPID ANALYTICS, WEKA, KEEL, togaware, and SPINF.
- Social:** Includes Apache Kafka, ThinkUp, and Corona.
- Key Value:** Includes AEROSPIKE, leveldb, GENIEDB, Chordless, Tokyo Cabinet, Scalaris, SCALIEN, Project Voldemort, hamsterdb, RAPTORDB, FairCom, STSDB, HyperDex, IQLECT, OpenLDAP, and ioremap.net.
- Document Store:** Includes mongoDB, Couchbase, CouchDB, Raven DB, CLUSTERPOINT, RaptorDB, EJDB, djon, JasDB, SchemafreeDB, sisodb, and denso db.
- Object databases:** Includes db4objects, ZOPE, NEOPPOD, STARCOUNTER, Magma, Sterling, EyeDB, Picolisp, siaqodb, HSS Database, and NDatabase.
- Graphs:** Includes Gephi, InfiniteGraph, AllegroGraph 4.9, FlockDB, GraphBuilder, Gremlin, INFO GRID, HYPERGRAPH-DB, dex, meronymy, GraphBase, and BrightstarDB.
- Multidimensional:** Includes GT.M, SciDB, and rasdaman.
- Grid Solutions:** Includes BIGASPACE, HAZELCAST, and Galaxy.
- Multimodel:** Includes ArangoDB and alchemydatabase.
- XML Databases:** Includes existdb, BASE, Qizx, sedna, and xindice.

<http://socr.umich.edu/docs/BD2K/BigDataResourceome.html>

# Examples of Available Resources

- Source Code: <https://github.com/SOCR>
- End-to-End Pipeline Workflows:
  - Docs: <https://wiki.loni.usc.edu>
  - Protocols: <http://pipeline.loni.usc.edu/explore/library-navigator>
- Pubs: <http://www.socr.umich.edu/people/dinov/publications.html>
- Training/Learning Resources:  
<http://wiki.socr.umich.edu/index.php/SMHS>
- Data:
  - Classical: [http://wiki.socr.umich.edu/index.php/SOCR\\_Data](http://wiki.socr.umich.edu/index.php/SOCR_Data)
  - Case-Studies:  
[https://umich.instructure.com/courses/38100/files/folder/Case\\_Studies](https://umich.instructure.com/courses/38100/files/folder/Case_Studies)

# Distributed Services

## Processing

- Pipeline Try-It-Now Graphical Workflow Environment  
(Guest access): <http://pipeline.loni.usc.edu/products-services/pws/>
- socr-pipeline.nursing.umich.edu
- AWS/Galaxy: <http://GalaxyProject.org>
- tranSMART: <https://github.com/transmart>

## Data

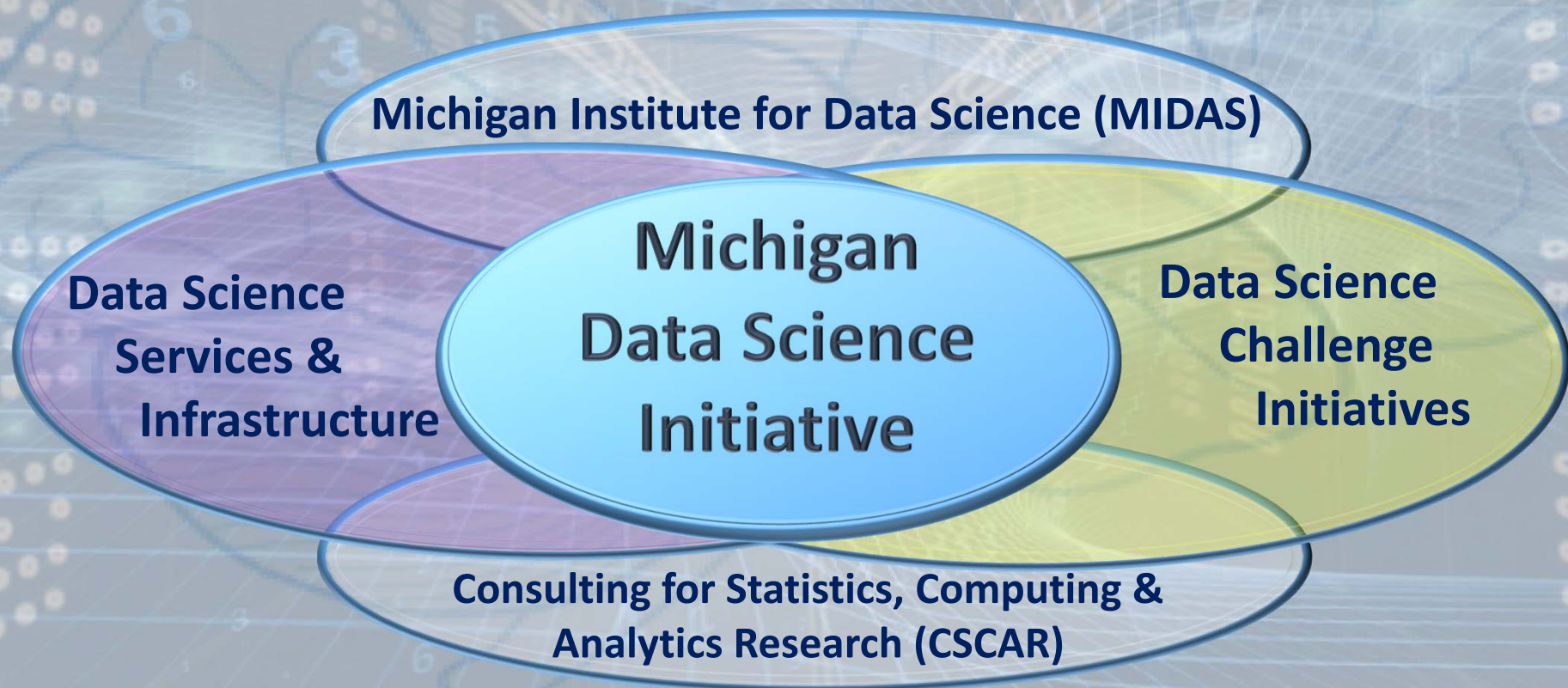
- dbGaP <http://dbgap.ncbi.nlm.nih.gov>
- Neuroimaging <http://IDA.loni.usc.edu>
- XNAT: <https://central.xnat.org>

## Transfer

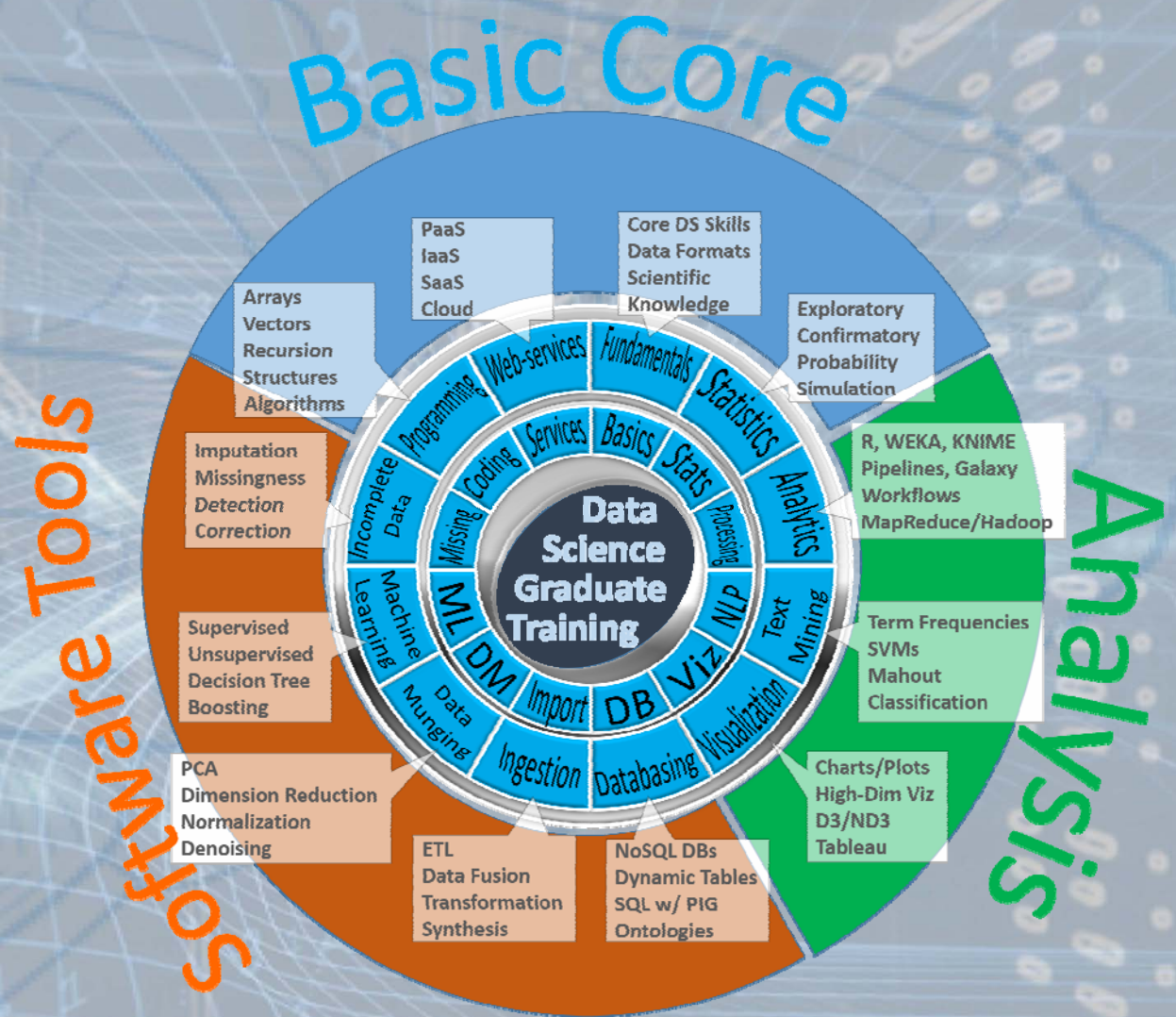
- Globus: <http://www.globusonline.org>
- GridFTP: <http://toolkit.globus.org/toolkit/docs/latest-stable/gridftp/>

# Michigan Institute for Data Science (MIDAS)

MIDAS catalyzes data science at the University of Michigan through support for faculty, research, education and training, and industry engagement



# Fundamentals of Data Science Education



# Grad DS Curriculum: Prereqs & Competencies

Prerequisites	Skills	Rationale
<b>BS degree or equivalent</b>	Quantitative training and coding skills as described below	The DS certificate is a graduate program requiring a minimum level of quantitative skill
<b>Quantitative training</b>	Undergraduate calculus, linear algebra and intro to probability and statistics	These are the entry level skills required for most upper-level undergrad and grad courses in program
<b>Coding experience</b>	Exposure to software development or programming on the job or in the classroom	Most DS practitioners need substantial experience with Java, C/C++, HTML5, Python, PHP, SQL/DB
<b>Motivation</b>	Significant interest and motivation to pursue quantitative data analytic applications	Dedication for prolonged & sustained immersion into hands-on and methodological research

# Grad DS Curriculum: Prereqs & Competencies

Areas	Competency	Expectation	Notes
Algorithms & Applications	Tools	Working knowledge of basic software tools (command-line, GUI based, or web-services)	Familiarity with statistical programming languages, e.g., R or SciKit/Python, and database querying languages, e.g., SQL or NoSQL
	Algorithms	Knowledge of core principles of scientific computing, applications programming, API's, algorithm complexity, and data structures	Best practices for scientific and application programming, efficient implementation of matrix linear algebra and graphics, elementary notions of computational complexity, user-friendly interfaces, strings
	Application Domain	Data analysis experience from at least one application area, either through coursework, internship, research project, etc.	Applied domain examples include: computational social sciences, health sciences, business and marketing, learning sciences, transportation sciences, engineering and physical sciences
Data Management	Data validation & visualization	Curation, Exploratory Data Analysis (EDA) and visualization	Data provenance, validation, visualization via histograms, Q-Q plots, scatterplots (ggplot, Dashboard, D3.js)
	Data wrangling	Skills for data normalization, data cleaning, data aggregation, and data harmonization/registration	Data imperfections include missing values, inconsistent string formatting ('2016-01-01' vs. '01/01/2016', PC/Mac/Linux time vs. timestamps, structured vs. unstructured data
	Data infrastructure	Handling databases, web-services, Hadoop, multi-source data	Data structures, SOAP protocols, ontologies, XML, JSON, streaming
Analysis Methods	Statistical inference	Basic understanding of bias and variance, principles of (non)parametric statistical inference, and (linear) modeling	Biological variability vs. technological noise, parametric (likelihood) vs non-parametric (rank order statistics) procedures, point vs. interval estimation, hypothesis testing, regression
	Study design and diagnostics	Design of experiments, power calculations and sample sizing, strength of evidence, p-values, False Discovery Rates	Multistage testing, variance normalizing transforms, histogram equalization, goodness-of-fit tests, model overfitting, model reduction
	Machine Learning	Dimensionality reduction, k-nearest neighbors, random forests, AdaBoost, kernelization, SVM, ensemble methods, CNN	Empirical risk minimization. Supervised, semi-supervised, and unsupervised learning. Transfer learning, active learning, reinforcement learning, multiview learning, instance learning



# Vertical Integration of MIDAS Challenges and Analytical Methods

**Learning  
Analytics**

**Trans-  
portation**

**Social  
Sciences**

**Health  
Sciences**

**Analytics and Visualization of Complex Data**

**Machine Learning-enabled Analytics**

**Temporal, Multi-Scale and Statistical Models**

**Integration of Heterogeneous Data**

**Data Scrubbing, Wrangling and Provenance Tracking**

**Data Privacy and Cybersecurity**



