

The perfect neuroimaging-genetics-computation storm: collision of petabytes of data, millions of hardware devices and thousands of software tools

Ivo D. Dinov · Petros Petrosyan · Zhizhong Liu · Paul Eggert · Alen Zamanyan · Federica Torri · Fabio Macciardi · Sam Hobel · Seok Woo Moon · Young Hee Sung · Zhiguo Jiang · Jennifer Labus · Florian Kurth · Cody Ashe-McNalley · Emeran Mayer · Paul M. Vespa · John D. Van Horn · Arthur W. Toga ·
for the Alzheimer's Disease Neuroimaging Initiative

Published online: 23 August 2013
© Springer Science+Business Media New York 2013

Abstract The volume, diversity and velocity of biomedical data are exponentially increasing providing petabytes of new neuroimaging and genetics data every year. At the same time, tens-of-thousands of computational algorithms are developed and reported in the literature along with thousands of software tools and services. Users demand intuitive, quick and

platform-agnostic access to data, software tools, and infrastructure from millions of hardware devices. This explosion of information, scientific techniques, computational models, and technological advances leads to enormous challenges in data analysis, evidence-based biomedical inference and reproducibility of findings. The Pipeline workflow environment

Some of the data used in preparation of this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database (adni.loni.ucla.edu). As such, the investigators within the ADNI contributed to the design and implementation of ADNI and/or provided data but did not participate in analysis or writing of this report. A complete listing of ADNI investigators can be found at: http://adni.loni.ucla.edu/wp-content/uploads/how_to_apply/ADNI_Acknowledgement_List.pdf.

I. D. Dinov (✉) · P. Petrosyan · Z. Liu · P. Eggert · A. Zamanyan · S. Hobel · J. D. Van Horn · A. W. Toga
Laboratory of Neuro Imaging (LONI), David Geffen School of Medicine at UCLA, University of California, Los Angeles, 635 S. Charles Young Drive, Suite 225, Los Angeles, CA 90095-7334, USA
e-mail: dinov@ucla.edu
URL: <http://www.loni.ucla.edu>

I. D. Dinov · F. Torri · F. Macciardi · A. W. Toga
Biomedical Informatics Research Network (BIRN), Information Sciences Institute, University of Southern California, Los Angeles, CA 90292, USA

F. Torri · F. Macciardi
Department of Psychiatry and Human Behavior, University of California, Irvine, Irvine, CA 92617, USA

P. Eggert
Department of Computer Science, University of California, Los Angeles, Los Angeles, CA 90095, USA

S. W. Moon
Department of Psychiatry, Konkuk University Chungju Hospital, Seoul, South Korea

Y. H. Sung
Department of Neurology, Gachon University, Gil Hospital, Incheon, South Korea

J. Labus · F. Kurth · C. Ashe-McNalley · E. Mayer
Center for Neurobiology of Stress, University of California, Los Angeles, Los Angeles, CA 90095, USA

P. M. Vespa
Brain Injury Research Center, Ronald Reagan UCLA Medical Center, Los Angeles, CA 90095, USA

Z. Jiang
Human Performance and Engineering Laboratory, Kessler Foundation Research Center, West Orange, NJ 07051, USA

Z. Jiang
Department of Biomedical Engineering, New Jersey Institute of Technology, Newark, NJ 07102, USA

provides a crowd-based distributed solution for consistent management of these heterogeneous resources. The Pipeline allows multiple (local) clients and (remote) servers to connect, exchange protocols, control the execution, monitor the states of different tools or hardware, and share complete protocols as portable XML workflows. In this paper, we demonstrate several advanced computational neuroimaging and genetics case-studies, and end-to-end pipeline solutions. These are implemented as graphical workflow protocols in the context of analyzing imaging (sMRI, fMRI, DTI), phenotypic (demographic, clinical), and genetic (SNP) data.

Keywords Aging · Pipeline · Neuroimaging · Genetics · Computation solutions · Workflows · IBS · Pain · Parkinson's disease · Alzheimer's disease · Shape · Volume · Analysis · Big data · Visualization

Introduction

Process understanding is frequently the core research question in many biomedical, health and environmental applications. As we rarely know the exact process characteristics, we collect data (observations) which is used as proxy of the underlying physiological, physical or environmental phenomena. As such, the observed information (data) becomes the pivotal aspect of the scientific inquiry. The data variability, complexity and heterogeneity directly affect the scientific inference, accuracy of the results and reproducibility of findings.

Three data characteristics make contemporary biomedical data different, challenging and powerful. These are the data *volume* (size), typically in the petabyte range (1 PB=10¹⁵ bytes), data *heterogeneity*, including (un)formatted, ASCII/Binary, (un)structured, and the data *velocity*, or data *derivative*, which captures the change, transfer, and discovery of raw and derived data (Foster et al. 2011; Howe et al. 2008; Lynch 2008).

Table 1 illustrates the Kryder's law for exponential increase of the volume of data (Walter 2005). Using two decades of data, this law predicts that the density of information on hard drives, areal density, increases by a factor of 1,000 every 10–11 years. This storage rate increase is driven by the rapid expansion of data volume and velocity and translates into doubling of data size each 12–13 months. Both Moore's and Kryder's laws indicate similar exponential increase (of computational power and data storage, respectively) over time (Sood et al. 2012).

There are thousands of software tools for acquisition, processing, storage/databasing, service, migration, mining, analysis, visualization, annotation, and *data-driven process understanding*. For example, the field of biomedical imaging includes hundreds of different types of image processing algorithms and filters. For each type of process, there may be dozens

of concrete software products (instance implementations). More specifically, the Neuroimaging Informatics Tools and Resources Clearing House (NITRC) (Luo et al. 2009) lists over 500 openly shared neuroimaging software tools. For each openly shared tool, there may be dozens of proprietary or less commonly used analogues. Similarly, in genomics and bioinformatics there are over 200 data and cloud computing service providers, and hundreds of public, private and non-profit organizations that provide thousands of stand-alone tools (Eliceiri et al. 2012). Resource organization, classification, discovery, traversal and utilization of these software products require flexible human and machine interfaces (Dinov et al. 2008).

Another computational challenge is the proliferation of millions of hardware devices. According to Cisco (Cisco Systems Inc 2012), by the end of 2012, the number of mobile-connected devices will exceed the number of people on Earth and there will be over 10 billion mobile-connected devices in 2016; i.e., there will be more than 1.3 mobile devices per capita worldwide. These include phones, tablets, laptops, handheld gaming consoles, e-readers, in-car entertainment systems, digital cameras, and “machine-to-machine modules.” There is a clear need for bridges between these mobile devices and for efficient connections to distributed databases, clients, servers, compute-nodes, web-services, variety of interfaces.

Methods

The LONI Pipeline environment (<http://Pipeline.loni.ucla.edu>) (Torri et al. 2012; Dinov et al. 2011) is a graphical workflow middleware providing an interface to computational libraries, informatics resources, computational expertise and cloud services (e.g., cloud data storage, cloud computing services). The Pipeline facilitates the design, validation, execution, monitoring and sharing of advanced heterogeneous computational protocols as graphical workflows. It also mediates the tool discovery and interoperability and provides distributed computing infrastructure for *en masse* data processing. The Pipeline's user-friendly interface enables access to disparate data, services, hardware infrastructure, computational expertise and cloud computing services (Torri et al. 2012).

Alternative infrastructures to the Pipeline environment that also facilitate visual informatics and computational genomics include Taverna (Oinn et al. 2004), Kepler (Ludäscher et al. 2006), Triana (Taylor et al. 2006), Galaxy (Goecks et al. 2010), AVS (Lord 1995), VisTrails (Freire et al. 2006), Bioclipse (Spjuth et al. 2007), KNIME (Berthold et al. 2008), NyPipe (Gorgolewski et al. 2011), PSOM (Bellec et al. 2012) and others. The choice of a workflow environment depends on the specific research domain, scientific application and computational need. The Pipeline environment provides

Table 1 Rapid increase of the volume of neuroimaging and genetics data

| Volume of data | | B. Neuroimaging (annually) | | C. Genomics (BP/year) | | D. Computational Power (CPU transistor counts) Moore's Law | | Years |
|--|--|----------------------------|-----------|-----------------------|--------|--|--------------------|-----------------------|
| MB = megabyte = 10^6 bytes, GB = gigabyte = 10^9 bytes, TB = terabyte = 10^{12} bytes, PB = petabyte = 10^{15} bytes | | | | | | | | |
| Single Cryo brain volume 1600 cm ² | | | | | | | | |
| A. Voxel resolution | | Gray scale | RGB color | 200 GB | 10 MB | | | |
| Size | Count | 8bits | 16bits | 24bits | 1 TB | 100 MB | | |
| 1 cm | $12 \times 15 \times 9$ | 1620 | 3000 | 4860 | 50 TB | 10 GB | | |
| 1 mm | $120 \times 150 \times 90$ | 1.62 MB | 3.24 MB | 4.86 MB | 250 TB | 1 TB | | |
| 100 μ m | $12000 \times 15000 \times 9000$ | 1.62 GB | 3.24 GB | 4.86 GB | 1 PB | 30 TB | | |
| 10 μ m | $120000 \times 150000 \times 90000$ | 1.62 TB | 3.24 TB | 4.86 TB | 5 PB | 1 PB | | |
| 1 μ m | $1200000 \times 1500000 \times 900000$ | 1.62 PB | 3.24 PB | 4.86 PB | 10+ PB | 20+ PB | | |
| | | | | | | | 1×10^5 | 1985–1989 |
| | | | | | | | 1×10^6 | 1990–1994 |
| | | | | | | | 5×10^6 | 1995–1999 |
| | | | | | | | 1×10^7 | 2000–2004 |
| | | | | | | | 8×10^6 | 2005–2009 |
| | | | | | | | 1×10^9 | 2010–2014 |
| | | | | | | | 1×10^{11} | 2015–2019 (estimated) |

^A Recent technological advances enable significant increases of the level of detail of optical imaging (e.g., cryotomographic brain images) into the micron (μ m) resolution (Niziachristos 2010; Roy et al. 2009; Scholl et al. 2011)

^B By 2012, there were 55 PBs of neuroimaging data (Breeze et al. 2012; Mennes et al. 2013), which may exaggerate the volume of neuroimaging data due to different publications sharing the same datasets. As of 2010, the Imaging Data Archive, a Laboratory of Neuro Imaging brain database, stored about 5×10^{15} B = 5 PBs data. Recent neuroimaging studies may generate 1.5 TB of data each week (Olabarriaga et al. 2010)

^C In 2011, the size of the genetics data is estimated to be 30 TBs (based on 10,000 human genomes) (Grossman and White 2012; Manusina 2012). As the total number of complete human genomes sequenced by the end of 2011 worldwide was >10,000, this figure may be orders of magnitude smaller than the real genomics data size. Furthermore, data derived from genome sequencing of other species and 'partial genomes' (e.g., exome capture sequencing, RNA sequencing and chromatin immunoprecipitation sequencing) is not included in this estimate. By 2015 more than a 10^6 human genomes will be sequenced (Grossman and White 2012). Assuming each genome takes about 10^{11} B (100 GB) this translates into a total data volume of 10^{17} B (100 PB). Some of the sequences may be whole-genome $100 \times$ depth/coverage acquisitions, and some may be acquired at lower depth

^D Data volume may be increasing at a faster pace compared to the well-established growth of computational power, Moore's law (Fuller and Millett 2011a; Rupp and Selberherr 2011)

some advantages over the alternative architectures. These include distributed client–server architecture, an array of scheduler grid plug-ins, external lightweight data manager, easy incorporation of new software tools and libraries, and dynamic workflow design, validation, execution, monitoring and dissemination of complete end-to-end computational solutions (Dinov et al. 2010).

The main types of computational tools available in the Pipeline library include software for neuroimaging and genetics data processing and visualization. For each of these types there are 3 categories of resources – data, atomic modules, and workflows. These resources can be explored via the Pipeline Navigator (<http://pipeline.loni.ucla.edu/explore/library-navigator/>) and can be tested via the guest-access Pipeline Web-Start server (<http://pipeline.loni.ucla.edu/PWS>). Many interesting end-to-end computational workflow solutions (pipelines) are documented online (<http://pipeline.loni.ucla.edu/explore/pipeline-workflows/>). There are also many video tutorials, screencasts, and training materials (<http://pipeline.loni.ucla.edu/learn/basic-videos/>), which illustrate the basic and advanced features of the pipeline client–server architecture, and the protocols for workflow design, execution and management.

Neuroimaging processing tools There are several hundred atomic neuroimage processing tools, from a variety of software suites available in the LONI pipeline library, Fig. 1a. These tools may be used for analysis of structural brain images (e.g., AFNI (Cox 1996), ROBEX (Iglesias et al. 2011), MDT Atlasing (Wang et al. 2005; Tang et al. 2010), BrainParser

(Tu et al. 2008), SVPASEG (Tohka et al. 2007, 2010), AIR (Woods et al. 1999), FSL (Smith et al. 2004), BrainSuite (Shattuck and Leahy 2000), SSMA (Leung 2011; Leung et al. 2008), ANTS (Avants et al. 2009), ITK (Pieper et al. 2006), MINC (Evans 2002)), functional brain data (e.g., FLIRT (Smith et al. 2005), AFNI (Cox 1996), WAIR (Dinov et al. 2002), Matlab (Hanselman and Littlefield 1997)), diffusion data (e.g., DTK (Wang et al. 2007a), DIRAC (Patel et al. 2010a), MiND (Patel et al. 2010b)), statistical analyses (e.g., R (Ihaka and Gentleman 1996), GAMMA (Chen and Herskovits 2005), SOCR (Che et al. 2009; Dinov 2006), SPM (Hu et al. 2005; Friston et al. 2011)), shape and surface modeling (e.g., sulcal analysis (Joshi et al. 2012), local and global shape analyses (Dinov et al. 2010), shape mapping DHM (Shi et al. 2007), FreeSurfer surface extraction, and cortical thickness (Fennema-Notestine et al. 2006; Fischl and Dale 2000a)).

Informatics and genomics computational library The breadth of genomics tools available as pipeline modules and workflows is illustrated by the variety of sequence alignment solutions (Torri et al. 2012), Fig. 1b. Some different categories of informatics and genomics computing software tools available in the Pipeline library include: sequence alignment (Mosaik (Smith et al. 2008), MAQ (Li et al. 2008), PERM (Chen et al. 2009), BWA/BWA-SW (Li and Durbin 2009, 2010), Bowtie (Langmead and Salzberg 2012), Novoalign (Li and Homer 2010), SOAPv2 (Li et al. 2009a), BLAST (Kent 2002)), indexing (mrFAST/mrsFAST (Hach et al. 2010)), genome-wide association studies (GWASS (Marchini et al. 2007), PLINK (Purcell et al. 2007)), basic

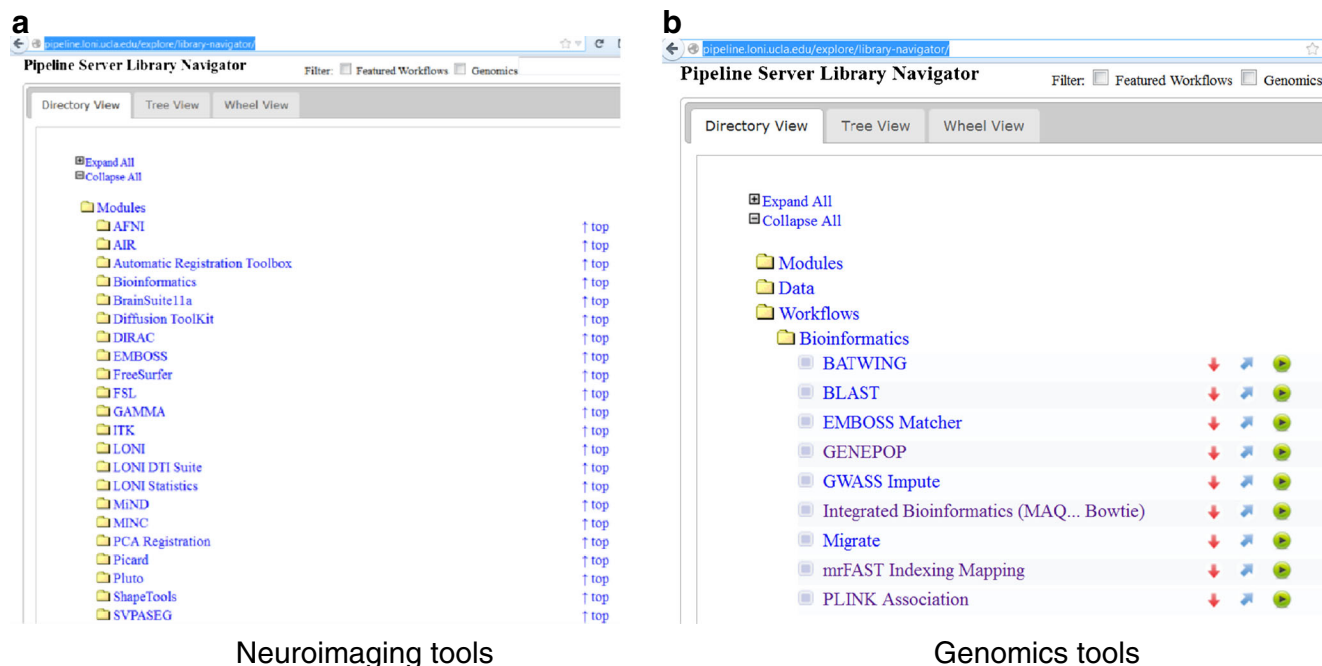


Fig. 1 Examples of classes of tools available in the Pipeline computational library

and advanced quality control (SAMTools (Li et al. 2009b), GATK (McKenna et al. 2010)), CNV calling (CNV/CNVR (Stranger et al. 2007; Wang et al. 2007b)), annotation (Artemis (Rutherford et al. 2000)), de novo assembly (Trinity (Grabherr et al. 2011), Velvet (Zerbino and Birney 2008)), molecular biology (EMBOSS (Olson 2002)), population genetics (GENEPOP (Raymond and Rousset 1995)), and many others.

Backend pipeline servers Pipeline web-start server (PWS) uses Java Web-Start technology enabling guest users to test the LONI Pipeline application from a web browser without the installation of either a pipeline client or a server. The PWS server provides access to all of the functions and features included in the downloadable version. PWS is accessible via an anonymous guest login or user-authentication to connect to remote Pipeline servers, e.g., <http://ucla.in/GRSc8a>. Several alternative Pipeline servers provide secure access-controlled connections to independent computational infrastructures. Examples include LONI Genomics Server (Genomics.loni.ucla.edu, 1 TB RAM/40-core), Cranium Server (Cranium.loni.ucla.edu, 16GB RAM/core, 1,200 cores) and Medulla Server (Medulla.loni.ucla.edu, 24GB RAM/core, 4,300 slots). The Distributed Pipeline Server infrastructure (<http://pipeline.loni.ucla.edu/DPS>) facilitates the deployment of independent disparate Pipeline services on available hardware resources, including Amazon EC2 (<http://pipeline.loni.ucla.edu/products-services/pipeline-server-on-ec2/>).

Big data Modern protocols for imaging and genetics data collection generate enormous amounts of data. Table 2 illustrates some of the data-management, storage and processing challenges associated with common neuroimaging and genetics analysis protocols. Figure 2 shows an example of the multi-channel imaging brain data typically acquired in traumatic brain injury studies.

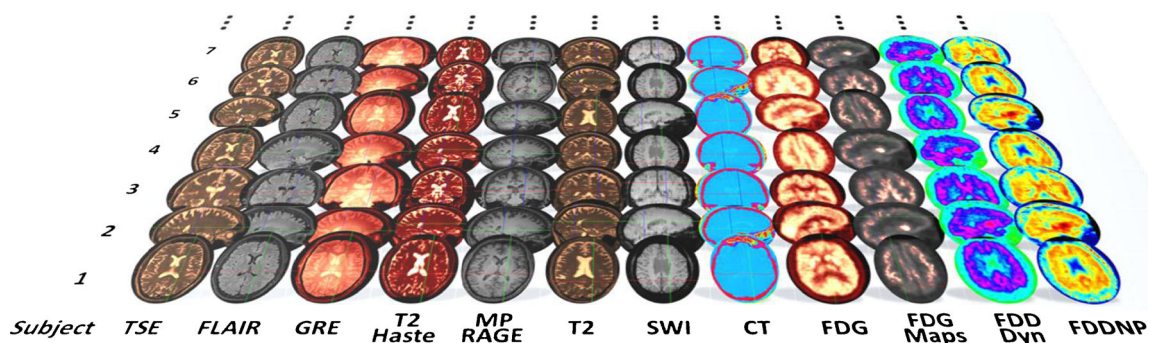


Fig. 2 Traumatic brain injury (TBI) studies demonstrate the diversity of the neuroimaging data in clinical applications. Imaging modalities included in many TBI studies include: *TSE* Turbo-Spin-Echo magnetic resonance imaging (MRI); *FLAIR* Fluid Attenuated Inversion Recovery MRI; *GRE* Gradient-Recalled-Echo (MRI); *T2 Haste* Half-Fourier Acquisition Single-Shot Turbo Spin-Echo MRI; *MP RAGE* Magnetization-Prepared

Applications and results

To demonstrate the Pipeline management of heterogeneous neuroimaging, genetics, phenotypic and clinical data, and the diversity of computational data processing tools available through the Pipeline library, we have chosen three complementary applications. These include studies of imaging-based genome-wide association, hippocampal morphometry, persistent pain and irritable bowel syndrome. Each of these three applications demonstrates exemplary solutions to the resource-scalability and processing-efficiency challenges related to the data complexity (size, heterogeneity and velocity), software tools interoperability and diversity of hardware devices. Specifically, these case-studies demonstrate (1) how seemingly incongruent imaging, phenotypic and clinical data can be jointly processed and analyzed in an integrated computational workflow protocol; (2) how pipeline workflows can wrap independent software tools to make them interoperate; and (3) how these data and computational resources (tools and services) can be accessed via different client devices (e.g., desktop or laptop computers or mobile devices running different operating systems and browser configurations).

ADNI imaging-genetics GWAS study The Alzheimer's disease data used in this article were obtained from the Alzheimer's Disease Neuroimaging Initiative (ADNI) database (adni.loni.ucla.edu). ADNI is the result of efforts of many co-investigators from a broad range of academic institutions and private corporations, and subjects have been recruited from over 50 sites across the U.S. and Canada. The initial goal of ADNI was to recruit 800 subjects but ADNI has been followed by ADNI-GO and ADNI-2. For up-to-date information, see www.adni-info.org.

The Alzheimer's Disease Neuroimaging Initiative (ADNI) (Jack et al. 2008; Mueller et al. 2005; Shen et al. 2010) data

Rapid Acquisition with Gradient Echo (MRI); *T2* T₂-weighted MRI; *SWI* Susceptibility Weighted Imaging (MRI); *CT* Computed Tomography; *FDG* Fludeoxyglucose Positron Emission Tomography (PET); *FDG Maps* Statistical maps of Fludeoxyglucose; *FDDNP* 2-(1-{6-[(2-[F-18]fluoroethyl)(methyl)amino]-2-naphthyl] ethylidene}malononitrile PET imaging

Table 2 Storage and processing of Big neuroimaging and genetics data

| Types of studies (Number of Subjects) | Neuroimaging Depending on acquisition & analysis protocols | | C. Genetics Depending on coverage(X) | |
|---|---|-----------------------------|---|---------------------------|
| | A. Storage | B. Processing | Storage | Processing |
| $N=1$ | Raw data: 10 GB (e.g., 512 directional diffusion data) Derived: 100 GB | 100+ GB RAM 70+ h CPU | 320 GB (at 80X) | 2+ TB RAM 100+ h CPU |
| D. Cohort Studies ($N\sim 100$) | 100 GB – 1 TB | 1 TB RAM 100's h CPU | 3+ TB | 2+ TB RAM 100's h CPU |
| Multi-site population wide studies ($N>1,000$) | 1–10 TB | 1+ TB RAM 1000's h CPU | 30+ TB | 2+ TB RAM 1000's h CPU |
| Longitudinal (Time ≥ 2) | > 5 TB | > 2 TB RAM > 5,000 h CPU | – | – |

^A Relative to the mouse brain, the field of view of human brain imaging data is several orders of magnitude larger (Sultan and Braitenberg 1993). Diffusion imaging of mouse brain may reach 1.9 GB ($7 \times 512 \times 256 \times 256$ points with real and imaginary parts, represented as 4 bits float numbers) (Jiang and Johnson 2010), and correspondingly diffusion spectral or high-angular resolution images may exceed 10 GB per human subject and session (Scholl et al. 2011). The Global Shape Analysis pipeline workflow (Dinov et al. 2010) includes about 100 processing steps and depending on the server load and the number of subjects provided as input may take 7 days to complete on the LONI Pipeline Medulla cluster (4 TB RAM, 3,000 slots)

^B Many computationally intensive neuroimaging processing tools require significant hardware resources including storage, memory and CPU cycles (Dinov et al. 2011)

^C In 2011, many alternative commercial DNA sequencing platforms generated whole genome sequences of size 100–600 GB (Glenn 2011), which require days of computations on powerful grid systems. For example, our experience shows that Trinity whole-genome de novo assembly (Zhang et al. 2011; Li et al. 2010) takes over 14 days of calculations on the LONI Pipeline Genomics server (1.4 TB RAM, 40-core)

^D The infrastructure needs of cohort-based and multi-institutional studies increase linearly with the increase of the number of cases that require processing. Thus, a brain study of 1,000 subjects (e.g., Chinese Probabilistic Brain Atlas (Xing et al. 2013), vGWAS (Hibar et al. 2012)) or a computational genetics study of 1,000 whole-genome sequences (e.g., prostate cancer (Gudmundsson et al. 2012), autism spectrum disorder (Buxbaum et al. 2012)) may require Terabytes of storage and extensive infrastructure for data management, processing and interrogation. Longitudinal neuroimaging studies add another layer of complexity, as these typically require baseline as well as several (1+) follow up scans, which increases proportionately the volume of the imaging data

was screened and from 589 study participants, 188 qualified for an Alzheimer's Disease (AD) diagnosis at baseline, 401 had mild cognitive impairment (MCI). Among them, 9 were early-onset (EO) AD (Male: 4, Female: 5) and 27 were early-onset MCI (Male: 15, Female: 12). Subjects (ages 55 to 65) were divided into two groups: EO-AD and EO-MCI. Individual ADNI genotype and imaging data were downloaded and merged to form a single dataset containing genome-wide information for 36 individuals. Genetic analysis, including quality control, were performed using PLINK version 1.09. All the genetic processing was done via the LONI Pipeline environment. The 20 most significant single nucleotide polymorphisms (SNPs) were chosen by Manhattan plot and were associated with specific neuroimaging biomarkers. The structural ADNI data (1.5 T MRI) was parcellated using BrainParser, and the 15 most important neuroimaging markers were extracted by the Global Shape Analysis (GSA) Pipeline workflow.

The goal of this application is to demonstrate the use of the pipeline environment for genome-wide association study (GWAS) using early-onset ADNI data including cognitive impairment measures, neuroimaging and genetics biomarkers. After standard SNP quality control (Hibar et al. 2011; Rimol et al. 2010), the raw SNP data (630K SNPs) was reduced to

360K SNPs. A new pipeline workflow was designed to integrate the global shape analysis, tensor-based morphometry and SOCR multivariate regression analyses. The results of the automated pipeline workflows included significant correlations between SNPs and various neuroimaging biomarkers in the EO subjects and discriminated between EO-AD and EO-MCI cohorts, Fig. 3. A connectomics diagram can be used to illustrate the strength of the associations between the 15 derived neuroimaging biomarkers and the top 20 SNP genetic markers. In this case-study, the small sample-size ($N=36$) has a negative effect on the (statistical) power to detect significant associations between the biomedical imaging markers (e.g., regional volume and shape metrics) and the genetic traits (SNPs/chromosomes). However, the same computational pipeline workflows can be used to analyze similarly larger cohorts (e.g., $N>700$), where sufficient power may be available to detect interactions between imaging and genetics effects (after Bonferonni correction for multiple testing). The imaging, genetics and clinical data used in this example were directly imported into the Pipeline workflow environment from the ADNI database using the Pipeline's IDAGet module. This pipeline workflow protocol can be designed on one client, and execution may be initiated on a user-specified pipeline server from another pipeline client, and the workflow

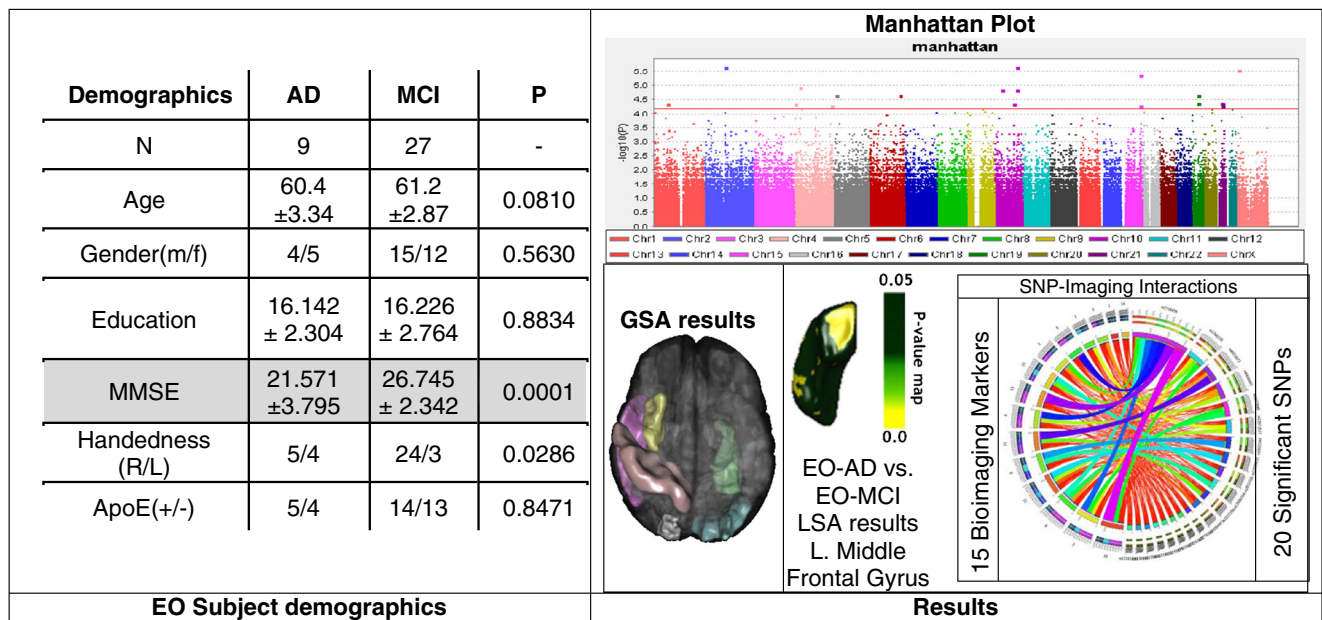


Fig. 3 Early Onset (EO) ADNI Imaging-Genetics GWAS Study using the pipeline environment

progress, and final result inspection, may be monitored or examined on a different client device.

Genetic associations with hippocampal function and shape A recent study investigated the genetics effects (single-nucleotide polymorphisms, SNP, associated with FKBP5 gene regulation, rs1360780) related to attention, behavioral, and hippocampal morphometrics (Fani et al. 2013). The FKBP5 gene regulates glucocorticoid receptor sensitivity and is associated with hypothalamic-pituitary-adrenal axis functioning and stress-related psychiatric disorders (Binder 2009). In this cross-sectional study using fMRI/MRI, African American cohort of adults ($N=103$) separated into 2 groups by genotype: Group 1 included carriers of the rs1360780 T allele, associated

with increased risk for posttraumatic stress disorder; Group 2 included non-carriers. The study used the local shape analysis pipeline workflow to identify attention bias toward threat ($F_{1,90}=5.19, p=0.02$), and revealed alterations in the hippocampal shape for TT/TC compared with the CC genotype groups. Figure 4 shows part of the computational protocol implemented as a pipeline workflow and the exemplary result from this morphometric analysis.

Persistent pain and irritable bowel syndrome (IBS) A UCLA IRB approved study recruited 328 female normal controls (NC) and IBS subjects. A diagnosis of IBS was made using the ROME III symptom criteria (Sperber et al. 2007; Drossman and Dumitrascu 2006) based on the assessment

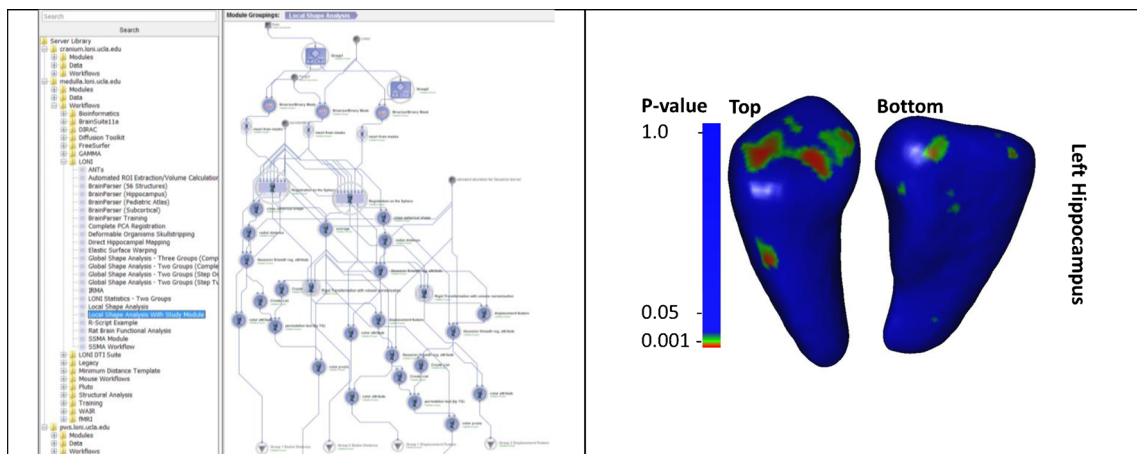


Fig. 4 Example of using the pipeline environment to complete a neuroimaging genetics study of FKBP5 gene (rs1360780) association with attention, measured through behavioral response (dot probe task) and hippocampal morphometrics. The superior and inferior views of the

hippocampal surface map illustrate the vertex locations, on the mean left hippocampus, where FKBP5 carriers (group 1) and non-carriers (group 2) showed significant shape differences

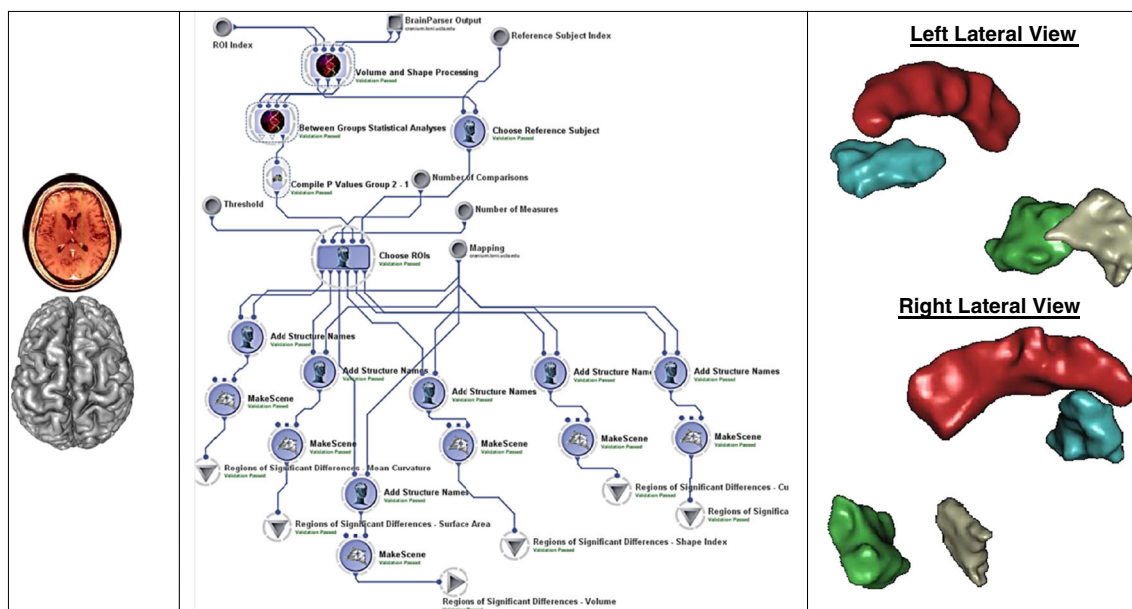


Fig. 5 Analyzing IBS/NC regional differences: (Left) raw sMRI data, (Middle) GSA workflow including data processing, surface reconstruction, 3D parcellation, and statistical analysis, (Right) Statistically signif-

icant ROI between-differences rendered as 3D scenes (left cuneus is green, and right angular gyrus is gray; the red cingulate gyrus and the blue insula are shown for orientation only)

by one of 4 gastroenterologists experienced in the diagnosis of functional bowel disease and the exclusion of organic disease. A subject's medical history and physical examination were obtained by a gastroenterologist. IBS patients with all types of predominant bowel habit were included. Subjects with a history of any chronic functional symptom or syndrome, or symptoms suggestive of disordered mood or affect, by history or by questionnaire, were excluded. In addition, potential subjects are excluded if by either history or questionnaire they a) have a serious medical condition or are taking medications which may interfere with interpretation of the brain imaging or physiological measures (other than IBS); b) have an ongoing major psychiatric diagnosis or psychotropic medication use over the past 6 months (subjects are not excluded for lifetime incidence of psychiatric disorder, or for intake of low dose tricyclic antidepressant for non-psychiatric indication); c) have a positive symptom score on the Hospital Anxiety and Depression Scale consistent with depression or anxiety d) do excessive physical exercise (i.e., marathon runners).

Brain images were obtained from all 328 subjects (107 IBS, 221 NC) using 1.5 and 3 T MRI scanners (Jiang et al. 2013). We collected phenotyping data on catastrophizing (Coping Strategies Questionnaire) (Geisser et al. 1994), early life trauma (Early Trauma Inventory) (Bremner et al. 2000), state anxiety and depression (Hospital Anxiety and Depression Scale) (Zigmond and Snaith 1983), health status (12-Item Short-Form Health Survey) (Ware et al. 1996), trait anxiety scores (State Trait Anxiety Inventory) (Spielberger 2005) and IBS symptom severity and duration (Bowel Symptoms Questionnaire) (Talley et al. 1995).

As a first step for shape-based neuroimage analysis, we reconstruct surface representation of anatomical structures of interest. Then, we analyze both cortical and subcortical structures. The cortical surfaces, including both white matter and pial surfaces, are reconstructed from T1-weighted MR images using FreeSurfer (Fischl and Dale 2000b). For sub-cortical structures, we applied the LONI BrainParser (Tu et al. 2008) to automatically segment the T1-weighted MR image into fifty-six regions. Using masks generated by BrainParser, accurate surface representations of the segmented regions are reconstructed with a novel algorithm we developed recently. This tool can remove segmentation artifacts without volume shrinkage and guarantees all surfaces guaranteed have the correct topology. All surfaces are represented as triangular meshes with spherical topology. The global shape analysis (GSA) pipeline workflow was used to identify regional differences between the NC and IBS subjects using the 56 regions of interest (ROIs) on 6 different volumetric and shape metrics (average mean curvature, surface area, volume, shape index, curvedness, and fractal dimension). Figure 5 shows the 3 steps in this analysis (data inputs, pipeline workflow and results of regional group differences).

Conclusions

Although there are a number of useful software discovery and navigation frameworks (Dinov et al. 2008; Kennedy 2006; Tenenbaum et al. 2011), the protocols for tool interoperability

continue to present significant biomedical computing challenges. There are considerable design differences between independent software suites. Furthermore, the varieties of computer programming languages for algorithm implementation, the substantial diversity of compilers and optimization strategies, and the gamut of hardware resources present additional hurdles in biomedical computing. Mediating these computational issues, coping with the enormous amounts of incongruent data, and handling a wide spectrum of devices require a paradigm shift of how we manage, process, interrogate and utilize biomedical and health related data.

The evidence is clear that we are in the front of an enormous storm of exponentially increasing wave of data, processing power and resource diversity. Multidisciplinary science efforts, technologies like Hadoop (White 2012), OpenStack (Wen et al. 2012), Elastic Cloud Computing (Ostermann et al. 2010), Pipeline workflow systems (Dinov et al. 2010; Heinis 2010) and super high-bandwidth networking (Chowdhury et al. 2010; Wang and Guo 2012) will be critical for riding this storm and uncovering novel biomedical knowledge. Embracing the *science interactome* (the multidisciplinary interactions between biomedical, computational and basic scientific areas, which often lead to new discoveries) will also be essential for establishing, maintaining and expanding the cyclical flow from Biomedical Challenges ↔ Scientific Models ↔ Data Analysis ↔ Computational Infrastructure ↔ Sustainable Education.

In this manuscript, we presented evidence of the rapid increase of the volume, diversity and velocity of biomedical data (e.g., neuroimaging and genetics (Novak et al. 2012; Van Essen et al. 2012; Thompson et al. 2013)), and the growth of computational models, algorithms, software tools, services and electronic devices that manipulate these data (Toga et al. 2012; Berger et al. 2013; Meir and Rubinsky 2009). There is evidence that software tool expansion always occurs within the limits of the available hardware infrastructure (Fuller and Millett 2011b). This close connection between the Moore's law for increase of computational power facilitates the observed expansion of new and more powerful software tools (e.g., Software as a Service (SaaS) (Hashizume et al. 2012), Platform-as-a-Service (PaaS) (Truong and Dustdar 2012)). For example, in 1993, Windows NT OS 3 consisted of 5-million lines of code, which 10-years later grew 10-fold to 50-million lines in Windows, Server OS 2003 (Maraia 2005). Similarly, from 2000 to 2007, the Linux Debian OS grew from 59-million to 280-million lines of code (Matellán Olivera 2012). Web and mobile applications, or webapps, are software systems running on portable devices, which have significantly grown since 2005 into a multi-billion dollar business (Minelli and Lanza 2013). The explosion of webapp software development can be measured in terms of pure source code, usage of third-party APIs, and historical data. Studies of lines of code in specific areas indicate that over the past few decades

there is an exponential increase of software development efforts (Knobloch 2013; German et al. 2013). This advancement of the software tool capabilities in turn pushes the introduction of more efficient and omnipotent hardware devices (e.g., Infrastructure as a Service (IaaS) and Virtual Machines (VMs) (Alarifi and Wolthusen 2013)).

The Pipeline workflow environment is one of many solutions that provide a distributed and platform-independent management of heterogeneous resources using dispersed clients and servers, elaborate exchange protocols, and flexible mechanisms for control, execution, monitoring and sharing of complete computational protocols. We demonstrated three advanced end-to-end computational pipeline solutions for neuroimaging, genetics and computational morphometry. These solutions are implemented as graphical workflow protocols in the context of analyzing imaging (sMRI, fMRI, DTI), phenotypic (demographic, clinical), and genetic (SNP) data.

Acknowledgments This study was partially supported by National Institutes of Health grants NIA P 50 AG16570, NIBIB EB01651, NLM LM05639, NIMH R01 MH071940, NIBIB P41EB015922, U24-RR025736, U24-RR021992, U24 GM104203, as well as National Science Foundation grants 0716055 and 1023115. The authors are also indebted to the faculty, staff and students in the Laboratory of Neuro Imaging (LONI) for their support and dedication.

Some of the data processed in this study was partly funded by the Alzheimer's Disease Neuroimaging Initiative (ADNI) (National Institutes of Health Grant U01 AG024904). ADNI is funded by the National Institute on Aging, the National Institute of Biomedical Imaging and Bioengineering, and through generous contributions from the following: Alzheimer's Association; Alzheimer's Drug Discovery Foundation; BioClinica, Inc.; Biogen Idec Inc.; Bristol-Myers Squibb Company; Eisai Inc.; Elan Pharmaceuticals, Inc.; Eli Lilly and Company; F. Hoffmann-La Roche Ltd and its affiliated company Genentech, Inc.; GE Healthcare; Innogenetics, N.V.; IXICO Ltd.; Janssen Alzheimer Immunotherapy Research & Development, LLC.; Johnson & Johnson Pharmaceutical Research & Development LLC.; Medpace, Inc.; Merck & Co., Inc.; Meso Scale Diagnostics, LLC.; NeuroRx Research; Novartis Pharmaceuticals Corporation; Pfizer Inc.; Piramal Imaging; Servier; Synarc Inc.; and Takeda Pharmaceutical Company. ADNI data are disseminated by the Laboratory for Neuro Imaging at the University of California, Los Angeles.

As of September 2013, the Laboratory of Neuro Imaging (LONI) will be relocated to the University of Southern California (USC). Thus, some of the URL links, web-page references, and internet resources cited throughout this manuscript may be relocated to appropriate subdomains under <http://www.loni.usc.edu>. If you find broken links or defunct URLs please contact help@loni.usc.edu.

References

- Alarifi, S., & Wolthusen S. (2013). *Anomaly detection for ephemeral cloud IaaS virtual machines*, In *Network and system security*, Springer. p. 321–335.
- Avants, B. B., Tustison, N., & Song, G. (2009). Advanced Normalization Tools (ANTS). *Insight J.*

- Bellec, P., et al. (2012). The pipeline system for Octave and Matlab (PSOM): a lightweight scripting framework and execution engine for scientific workflows. *Frontiers in Neuroinformatics*, 6.
- Berger, B., Peng, J., & Singh, M. (2013). Computational solutions for omics data. *Nature Reviews Genetics*, 14(5), 333–346.
- Berthold, M. R., et al. (2008). *KNIME: The Konstanz information miner*, in *Data analysis, machine learning and applications*. C. Preisach, et al., (Eds.), Springer Berlin Heidelberg. p. 319–326.
- Binder, E. B. (2009). The role of FKBP5, a co-chaperone of the glucocorticoid receptor in the pathogenesis and therapy of affective and anxiety disorders. *Psychoneuroendocrinology*, 34, S186–S195.
- Breeze, J. L., Poline, J.-B., & Kennedy, D. N. (2012). Data sharing and publishing in the field of neuroimaging. *Giga Science*, 1(1), 1–3.
- Bremner, J. D., Vermetten, E., & Mazure, C. M. (2000). Development and preliminary psychometric properties of an instrument for the measurement of childhood trauma: the Early Trauma Inventory. *Depression and Anxiety*, 12(1), 1–12.
- Buxbaum, J. D., et al. (2012). The autism sequencing consortium: large-scale, high-throughput sequencing in autism spectrum disorders. *Neuron*, 76(6), 1052–1056.
- Che, A., Cui, J., & Dinov, I. (2009). SOCR analyses: implementation and demonstration of a New graphical statistics educational toolkit. *JSS*, 30(3), 1–19.
- Chen, R., & Herskovits, E. H. (2005). Graphical-model-based morphometric analysis. *Medical Imaging, IEEE Transactions on*, 24(10), 1237–1248.
- Chen, Y., Souaiaia, T., & Chen, T. (2009). PerM: efficient mapping of short sequencing reads with periodic full sensitive spaced seeds. *Bioinformatics*, 25(19), 2514–2521.
- Chowdhury, A., et al. (2010). Next-generation E-health communication infrastructure using converged super-broadband optical and wireless access system. In *World of Wireless Mobile and Multimedia Networks (WoWMoM)*, 2010 I.E. International Symposium on. *IEEE*.
- Cisco Systems Inc. *Cisco Visual Networking Index: Global Mobile Data Traffic Forecast Update, 2012–2017*. Cisco 2012; Available from: http://www.cisco.com/en/US/solutions/collateral/ns341/ns525/ns537/ns705/ns827/white_paper_c11-520862.pdf.
- Cox, R. W. (1996). AFNI: software for analysis and visualization of functional magnetic resonance neuroimages. *Computers and Biomedical Research*, 29(3), 162–173.
- Dinov, I. (2006). Statistics online computational resource. *Journal of Statistical Software*, 16(1), 1–16.
- Dinov, I. D., et al. (2002). Quantitative comparison and analysis of brain image registration using frequency-adaptive wavelet shrinkage. *IEEE Transactions on Information Technology in Biomedicine*, 6(1), 73–85.
- Dinov, I., et al. (2008). iTools: a framework for classification, categorization and integration of computational biology resources. *PLoS One*, 3(5), e2265.
- Dinov, I., et al. (2010). Neuroimaging study designs, computational analyses and data provenance using the LONI pipeline. *PLoS One*, 5(9), e13070. doi:10.1371/journal.pone.0013070.
- Dinov, I., et al. (2011). Applications of the pipeline environment for visual informatics and genomics computations. *BMC Bioinformatics*, 12(1), 304.
- Drossman, D., & Dumitrescu, D. (2006). Rome III: new standard for functional gastrointestinal disorders. *Journal of Gastrointestinal and Liver Diseases: JGLD*, 15(3), 237.
- Eliceiri, K. W., et al. (2012). Biological imaging software tools. *Nature Methods*, 9(7), 697–710.
- Evans, A. (2002). *Automated 3D analysis of large brain MRI databases*. Neuropsychopharmacology: The Fifth Generation of Progress: American College of Neuropsychopharmacology. Nature Publishing, London: p. 301–313.
- Fani, N., G.D.T.E.B., et al. (2013). Fkbp5 and attention bias for threat: associations with hippocampal function and shape. *JAMA Psychiatry*: p. 1–9.
- Fennema-Notestine, C., et al. (2006). Quantitative evaluation of automated skull-stripping methods applied to contemporary and legacy images: effects of diagnosis, bias correction, and slice location. *Human Brain Mapping*, 27(2), 99–113.
- Fischl, B., & Dale, A. M. (2000). Measuring the thickness of the human cerebral cortex from magnetic resonance images. *Proceedings of the National Academy of Sciences of the United States of America*, 97(20), 11050–11055.
- Foster, K., Spicer, M., & Nathan, S. (2011). *IBM infosphere streams: Assembling continuous insight in the information revolution*. San Jose: International Technical Support Organization.
- Freire, J., et al. (2006). *Managing rapidly-evolving scientific workflows*, in *IPAW 2006*, L.M.a.I.F. (Eds.), Springer-Verlag: Berlin Heidelberg. p. 10–18.
- Friston, K. J., et al. (2011). *Statistical parametric mapping: The analysis of functional brain images: The analysis of functional brain images*: Academic Press.
- Fuller, S. H., & Millett, L. I. (2011). Computing performance: game over or next level? *Computer*, 44(1), 31–38.
- Fuller, S. H., & Millett, L. I. (2011). *The future of computing performance: game over or next level?*: The National Academies Press.
- Geisser, M. E., Robinson, M. E., & Henson, C. D. (1994). The Coping Strategies Questionnaire and chronic pain adjustment: a conceptual and empirical reanalysis. *The Clinical Journal of Pain*.
- German, D. M., Adams, B., & Hassan, A. E. (2013). The Evolution of the R Software Ecosystem. In *Software Maintenance and Reengineering (CSMR)*, 2013 17th European Conference on. *IEEE*.
- Glenn, T. C. (2011). Field guide to next-generation DNA sequencers. *Molecular Ecology Resources*, 11(5), 759–769.
- Goecks, J., et al. (2010). Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences. *Genome Biology*, 11(8), R86.
- Gorgolewski, K., et al. (2011). Nipype: a flexible, lightweight and extensible neuroimaging data processing framework in python. *Frontiers in Neuroinformatics*. 5.
- Grabherr, M. G., et al. (2011). Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nature Biotechnology*, 29(7), 644–652.
- Grossman, R., & White, K. (2012). A vision for a biomedical cloud. *Journal of Internal Medicine*, 271(2), 122–130.
- Gudmundsson, J., et al. (2012). A study based on whole-genome sequencing yields a rare variant at 8q24 associated with prostate cancer. *Nature Genetics*.
- Hach, F., et al. (2010). mrsFAST: a cache-oblivious algorithm for short-read mapping. *Nature Methods*, 7(8), 576–577.
- Hanselman, D., & Littlefield, B. C. (1997). *Mastering MATLAB 5: A comprehensive tutorial and reference*: Prentice Hall PTR.
- Hashizume, K., Fernandez, E. B., & Larrondo-Petrie M. M. (2012). A pattern for Software-as-a-Service in Clouds. In *BioMedical Computing (BioMedCom)*, 2012 ASE/IEEE International Conference on. *IEEE*.
- Heinis, T. (2010). *Workflow-based services: infrastructure for scientific applications*: Suedwestdeutscher Verlag fuer Hochschulschriften.
- Hibar, D. P., et al. (2011). Voxelwise gene-wide association study (vGeneWAS): multivariate gene-based association testing in 731 elderly subjects. *NeuroImage*, 56(4), 1875–1891.
- Hibar, D. P., et al. (2012). Genome-wide association identifies genetic variants associated with lentiform nucleus volume in N = 1345 young and elderly subjects. *Brain Imaging and Behavior* p. 1–14.
- Howe, D., et al. (2008). Big data: the future of biocuration. *Nature*, 455(7209), 47–50.
- Hu, D., et al. (2005). Unified SPM-ICA for fMRI analysis. *NeuroImage*, 25(3), 746–755.
- Iglesias, J. E., et al. (2011). Robust brain extraction across datasets and comparison with publicly available methods. *Medical Imaging, IEEE Transactions on*, 30(9), 1617–1634.

- Thakur, R., & Gentleman, R. (1996). R: a language for data analysis and graphics. *Journal of Computational and Graphical Statistics*, 5(3), 299–314.
- Jack, C. R., et al. (2008). The Alzheimer's disease neuroimaging initiative (ADNI): MRI methods. *Journal of Magnetic Resonance Imaging*, 27(4), 685–691.
- Jiang, Y., & Johnson, G. A. (2010). Microscopic diffusion tensor imaging of the mouse brain. *NeuroImage*, 50(2), 465–471.
- Jiang, Z., et al. (2013). Sex-related differences of cortical thickness in patients with chronic abdominal pain. in press.
- Joshi, S. H., et al. (2012). Diffeomorphic sulcal shape analysis on the cortex. *Medical Imaging, IEEE Transactions on*, PP(99): p. 1–1.
- Kennedy, D. N. (2006). The internet analysis tools registry: a public resource for image analysis. *Neuroinformatics*, 4, 263–270.
- Kent, W. J. (2002). BLAT—the BLAST-like alignment tool. *Genome Research*, 12(4), 656–664.
- Knobloch, J. (2013). *Four decades of computing in subnuclear physics—from bubble chamber to LHC*. arXiv preprint arXiv:1302.2974.
- Langmead, B., & Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nature Methods*, 9(4), 357–359.
- Leung, K. T. K. (2011). *Principal ranking meta-algorithms*. Los Angeles: University of California.
- Leung, K., et al. (2008). *IRMA: an image registration meta-algorithm - evaluating alternative algorithms with multiple metrics*. SSDBM 2008. Springer-Verlag.
- Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics*, 25(14), 1754–1760.
- Li, H., & Durbin, R. (2010). Fast and accurate long-read alignment with Burrows–Wheeler transform. *Bioinformatics*, 26(5), 589–595.
- Li, H., & Homer, N. (2010). A survey of sequence alignment algorithms for next-generation sequencing. *Briefings in Bioinformatics*, 11(5), 473–483.
- Li, H., Ruan, J., & Durbin, R. (2008). Mapping short DNA sequencing reads and calling variants using mapping quality scores. *Genome Research*, 18(11), 1851–1858.
- Li, R., et al. (2009a). SOAP2: an improved ultrafast tool for short read alignment. *Bioinformatics*, 25(15), 1966–1967.
- Li, H., et al. (2009b). The Sequence alignment/map format and SAMtools. *Bioinformatics*, 25, 2078–2079.
- Li, R., et al. (2010). De novo assembly of human genomes with massively parallel short read sequencing. *Genome Research*, 20(2), 265.
- Lord, H. D. (1995). Improving the application development process with modular visualization environments. *SIGGRAPH Computing Graph*, 29(2), 10–12.
- Ludäscher, B., Altintas, I., Berkley, C., Higgins, D., Jaeger, E., Jones, M., Lee, E. A., Tao, J., & Zhao, Y. (2006). Scientific workflow management and the Kepler system. *Concurrency and Computation: Practice and Experience*, 18(10), 1039–1065.
- Luo, X.-Z. J., Kennedy, D. N., & Cohen, Z. (2009). Neuroimaging informatics tools and resources clearinghouse (NITRC) resource announcement. *Neuroinformatics*, 7(1), 55–56.
- Lynch, C. (2008). Big data: how do your data grow? *Nature*, 455(7209), 28–29.
- Maraia, V. (2005). *The build master: Microsoft's software configuration management best practices*: Addison-Wesley Professional.
- Marchini, J., et al. (2007). A new multipoint method for genome-wide association studies by imputation of genotypes. *Nature Genetics*, 39(7), 906–913.
- Marusina, K. (2012). Big data requires big solutions. *Genetic Engineering & Biotechnology News* 32(15): p. 1, 34–40.
- Matellán Olivera, V. (2012). *Studying the evolution of libre software projects using publicly available data*; Available from: <https://buleria.unileon.es/handle/10612/1796>.
- McKenna, A., et al. (2010). The genome analysis toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Research*, 20(9), 1297–1303.
- Meir, A., & Rubinsky, B. (2009). Distributed network, wireless and cloud computing enabled 3-D ultrasound; a new medical technology paradigm. *PLoS One*, 4(11), e7974.
- Mennes, M., et al. (2013). Making data sharing work: the FCP/INDI experience. *NeuroImage*, (0).
- Minelli, R., & Lanza, M. (2013). Software analytics for mobile applications—insights & lessons learned. In Software maintenance and reengineering (CSMR), 2013 17th European Conference on. *IEEE*.
- Mueller, S. G., et al. (2005). Ways toward an early diagnosis in Alzheimer's disease: the Alzheimer's disease neuroimaging initiative (ADNI). *Alzheimer's and Dementia: The Journal of the Alzheimer's Association*, 1(1), 55–66.
- Novak, N. M., et al. (2012). EnigmaVis: online interactive visualization of genome-wide association studies of the Enhancing NeuroImaging Genetics through Meta-Analysis (ENIGMA) consortium. *Twin Research and Human Genetics: The Official Journal of the International Society for Twin Studies*, 15(3), 414.
- Ntzachristos, V. (2010). Going deeper than microscopy: the optical imaging frontier in biology. *Nature Methods*, 7(8), 603–614.
- Oinn, T., et al. (2004). Taverna: a tool for the composition and enactment of bioinformatics workflows. *Bioinformatics*, 20(17), 3045–3054.
- Olabarriaga, S. D., Glatard, T., & de Boer, P. T. (2010). A virtual laboratory for medical image analysis. *Information Technology in Biomedicine, IEEE Transactions on*, 14(4), 979–985.
- Olson, S. A. (2002). EMBOSS opens up sequence analysis. European molecular biology open software suite. *Briefings in Bioinformatics*, 3(1), 87.
- Ostermann, S., et al. (2010). A performance analysis of EC2 cloud computing services for scientific computing. *Cloud Computing*, p. 115–131.
- Patel, V., et al. (2010a). Mesh-based spherical deconvolution: a flexible approach to reconstruction of non-negative fiber orientation distributions. *NeuroImage*, 51(3), 1071–1081.
- Patel, V., et al. (2010b). LONI MiND: metadata in NIFTI for DWI. *NeuroImage*, 51(2), 665–676.
- Pieper, S., Lorensen, B., Schroeder, W., Kikinis, R. (2006). The NA-MIC Kit: ITK, VTK, pipelines, grids and 3D slicer as an open platform for the medical image computing community. In Biomedical Imaging: Nano to Macro, 2006. *3rd IEEE International Symposium on*
- Purcell, S., et al. (2007). PLINK: a tool set for whole-genome association and population-based linkage analyses. *The American Journal of Human Genetics*, 81(3), 559–575.
- Raymond, M., & Rousset, F. (1995). GENEPOP (Version 1.2): population genetics software for exact tests and ecumenicism. *Journal of Heredity*, 86(3), 248–249.
- Rimol, L. M., et al. (2010). Sex-dependent association of common variants of microcephaly genes with brain structure. *Proceedings of the National Academy of Sciences*, 107(1), 384–388.
- Roy, D., et al. (2009). 3D cryo-imaging: a very high-resolution view of the whole mouse. *The Anatomical Record*, 292(3), 342–351.
- Rupp, K., & Selberher, S. (2011). The economic limit to Moore's Law. *Semiconductor Manufacturing, IEEE Transactions on*, 24(1), 1–4.
- Rutherford, K., et al. (2000). Artemis: sequence visualization and annotation. *Bioinformatics*, 16(10), 944–945.
- Scholl, I., et al. (2011). Challenges of medical image processing. *Computer Science—Research and Development*, 26(1–2), 5–13.
- Shattuck, D., & Leahy R. (2000). BrainSuite: An automated cortical surface identification tool, in Medical image computing and computer-assisted intervention—MICCAI 2000, *Lecture Notes in Computer Science*. p. 50–61.
- Shen, L., et al. (2010). Whole genome association study of brain-wide imaging phenotypes for identifying quantitative trait loci in MCI and AD: a study of the ADNI cohort. *NeuroImage*, 53(3), 1051.
- Shi, Y., Thompson, P. M., Dinov, I. D., Osher, S., & Toga, A. W. (2007). Direct cortical mapping via solving partial differential equations on implicit surfaces. *Medical Image Analysis*, 11(3), 207–223.

- Smith, S. M., et al. (2004). Advances in functional and structural MR image analysis and implementation as FSL. *NeuroImage*, 23(Supplement 1), S208–S219.
- Smith, S. M., et al. (2005). Variability in fMRI: a re-examination of inter-session differences. *Human Brain Mapping*, 24(3), 248–257.
- Smith, D. R., et al. (2008). Rapid whole-genome mutational profiling using next-generation sequencing technologies. *Genome Research*, 18(10), 1638–1642.
- Sood, A., et al. (2012). Predicting the path of technological innovation: SAW vs. Moore, bass, gompertz, and kryder. *Marketing Science*, 31(6), 964–979.
- Sperber, A. D., et al. (2007). A comparative reappraisal of the Rome II and Rome III diagnostic criteria: are we getting closer to the 'true' prevalence of irritable bowel syndrome? *European Journal of Gastroenterology & Hepatology*, 19(6), 441.
- Spielberger, C. D. (2005). *State-trait anxiety inventory*: Wiley Online Library.
- Spjuth, O., et al. (2007). Bioclipse: an open source workbench for chemo- and bioinformatics. *BMC Bioinformatics*, 8(1), 59.
- Stranger, B. E., et al. (2007). Relative impact of nucleotide and copy number variation on gene expression phenotypes. *Science*, 315(5813), 848–853.
- Sultan, F., & Braitenberg, V. (1993). Shapes and sizes of different mammalian cerebella. A study in quantitative comparative neuroanatomy. *Journal für Hirnforschung*, 34(1), 79.
- Talley, N., et al. (1995). Initial validation of a bowel symptom questionnaire* and measurement of chronic gastrointestinal symptoms in Australians. *Internal Medicine Journal*, 25(4), 302–308.
- Tang, Y., et al. (2010). The construction of a Chinese MRI brain atlas: a morphometric comparison study between Chinese and caucasian cohorts. *NeuroImage*, 51(1), 33–41.
- Taylor, I., Shields, M., Wang, I., & Harrison, A. (2006). Visual grid workflow in triana. *Journal of Grid Computing*, 3, 153–169.
- Tenenbaum, J. D., et al. (2011). The biomedical resource ontology (BRO) to enable resource discovery in clinical and translational research. *Journal of Biomedical Informatics*, 44(1), 137–145.
- Thompson, P. M., et al. (2013). Genetics of the connectome. *Neuroimage*.
- Toga, A. W., et al. (2012). The center for computational biology: resources, achievements, and challenges. *Journal of the American Medical Informatics Association*, 19(2), 202–206.
- Tohka, J., et al. (2007). Genetic algorithms for finite mixture model based voxel classification in neuroimaging. *Medical Imaging, IEEE Transactions on*, 26(5), 696–711.
- Tohka, J., et al. (2010). Brain MRI tissue classification based on local Markov random fields. *Magnetic Resonance Imaging*, 28(4), 557–573.
- Torri, F., et al. (2012). Next generation sequence analysis and computational genomics using graphical pipeline workflows. *Genes*, 3(3), 545–575.
- Truong, H.-L., & Dustdar, S. (2012). A survey on cloud-based sustainability governance systems. *International Journal of Web Information Systems*, 8(3), 278–295.
- Tu, Z., et al. (2008). Brain anatomical structure segmentation by hybrid discriminative/generative models. *IEEE Transactions on Medical Imaging*, 27(4), 495–508.
- Van Essen, D. C., et al. (2012). The human connectome project: a data acquisition perspective. *NeuroImage*, 62(4), 2222–2231.
- Walter, C. (2005). Kryder's law. *Scientific American*, 293(2), 32–33.
- Wang, W., & Guo, L. (2012). The Development and Applications of Wireless Streaming Media Technology. In Computer Science and Electronics Engineering (ICCSEE), 2012 International Conference on. *IEEE*.
- Wang, Q., et al. (2005). Construction and validation of mean shape atlas templates for atlas-based brain image segmentation. In *Information Processing in Medical Imaging*. Springer.
- Wang, R., et al. (2007). Diffusion toolkit: a software package for diffusion imaging data processing and tractography. In *Proc Intl Soc Mag Reson Med*.
- Wang, K., et al. (2007b). PennCNV: an integrated hidden Markov model designed for high-resolution copy number variation detection in whole-genome SNP genotyping data. *Genome Research*, 17(11), 1665–1674.
- Ware, J. E., Jr., Kosinski, M., & Keller, S. D. (1996). A 12-item short-form health survey: construction of scales and preliminary tests of reliability and validity. *Medical Care*, 34(3), 220.
- Wen, X., et al. (2012). Comparison of open-source cloud management platforms: OpenStack and OpenNebula. In Fuzzy Systems and Knowledge Discovery (FSKD), 2012 9th International Conference on. *IEEE*.
- White, T. (2012). *Hadoop: The definitive guide*: O'Reilly Media.
- Woods, R. P., Dapretto, M., Sicotte, N. L., Toga, A. W., & Mazziotta, J. C. (1999). Creation and use of a Talairach-compatible atlas for accurate, automated, nonlinear intersubject registration, and analysis of functional imaging data. *Human Brain Mapping*, 8(2–3), 73–79.
- Xing, W., et al. (2013). Probabilistic MRI brain anatomical atlases based on 1,000 Chinese subjects. *PLoS One*, 8(1), e50939.
- Zerbino, D. R., & Birney, E. (2008). Velvet: algorithms for de novo short read assembly using de Bruijn graphs. *Genome Research*, 18(5), 821–829.
- Zhang, W., et al. (2011). A practical comparison of de novo genome assembly software tools for next-generation sequencing technologies. *PLoS One*, 6(3), e17915.
- Zigmond, A. S., & Snaith, R. (1983). The hospital anxiety and depression scale. *Acta Psychiatrica Scandinavica*, 67(6), 361–370.