

Open Source + Open Data = Open Science

Terry S. Yoo, PhD

OHPCC / LHNCBC / National Library of Medicine
National Institutes of Health



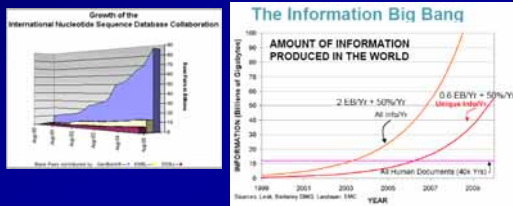
Open Science

What do you do when your customers are drowning in data?
What do you do when your engineers are starving for data?
What do you do when your scientists are reinventing the wheel?

- Infrastructure building.
- Validation, reproducible science.
- Open Source. Open Data. Transparency.

2

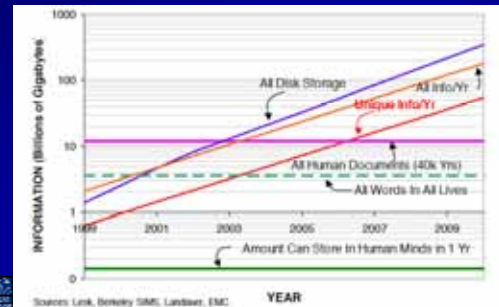
Do you recognize this problem?



- European Molecular Biology Laboratory's European Bioinformatics Institute (EMBL-Bank in Hinxton, UK)
- DNA Data Bank of Japan, had been launched at the National Institute of Genetics in Mishima

3

The Challenge of the Future: Information Big Bang



Sources: Link, Berkeley IRI, Laidlaw, EMC

A Poverty of Attention



"A wealth of information creates a poverty of attention and a need to allocate it efficiently."

-Herbert Simon,
Mathematician
Noble Laureate, Economics

5

Current State of most Three-D Informatics Research

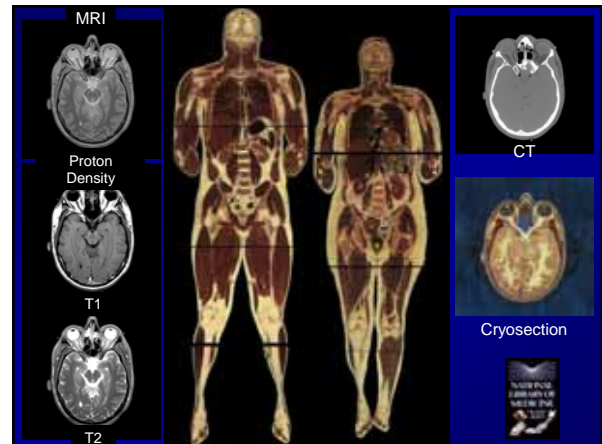
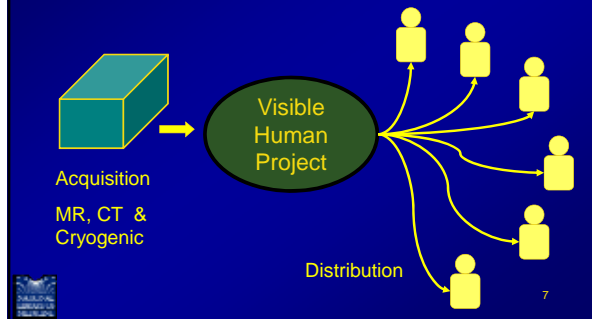
- Closed source code environment
- Protected proprietary data

This is a two way street. If you beg, borrow, or steal someone's data, you should be willing to donate your code or data back to the community!

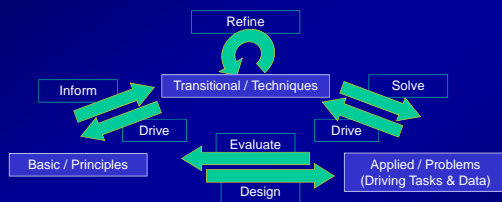
Open Source + Open Data = Open Science

6

Where did we start ?



Research Cycles



- Moving beyond our current development trends.
- Balancing our concentrations of research effort.
- Closing the loop.
- Inherently interdisciplinary.

Open Source Initiatives

- Encourage high-level technical communication
- Provide conventions (vs. standards) for interoperable software development
- Establish a baseline for improvement
- Opens the field to “beginners”
- Creates common ground for product development
 - example: the creation of HTML enabled Web-based internet development

Imaging Tools - “Insight”

The purpose of computing is insight, not numbers - Hamming

- A common platform to encourage communication and dissemination of research results.
- This initiative specifically does not include the development of visualization techniques or a user interface.
- Original release: September 2000.
- Current release: ITK 3.20.
- 45 countries.

Open Source



- Source code will be openly distributed.
- Consistent code style and data flow model.
- Regression tests will be included with all software.
- Free licensing for all software.

Software Development Program: Promoting Research

- Image Segmentation
 - multivalued (multimodal) data
- Image Registration
 - rigid and deformable registration
- Validation
 - Generation of mathematical models as test data
 - Acquisition of validation datasets from medical scanners



13

ITK Charter Sponsors



The National Institute of Neurological Disorders and Stroke



14

The Insight Team (Charter Members)

"The purpose of computing is insight, not numbers." - Hamming



The Technology of the Insight Toolkit

Segmentation

- Statistical, Fuzzy Logic, Markov Random Fields, Mixture Modeling, Parzen Windows, Nearest Neighbor, K-Means, ...
- Level Set, Finite Element, Region Growing, Hybrid, Watershed, Connected Components, Parameterized Models, ...

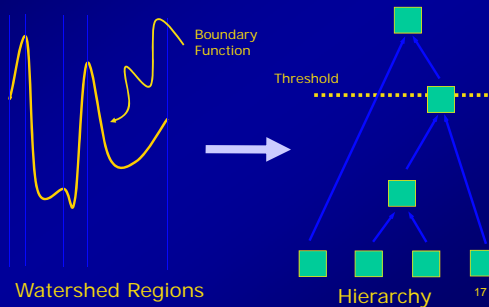
Registration

- Rigid, Similarity, Affine, Vector Field, Hierarchical, Quaternion, Versor, Parameterized Deformation, Euler, Perspective, ND/MD, ...
- Mutual Information, Normalized Correlation, Demons, Mean Squared, Landmark, ...



16

Watershed Algorithm

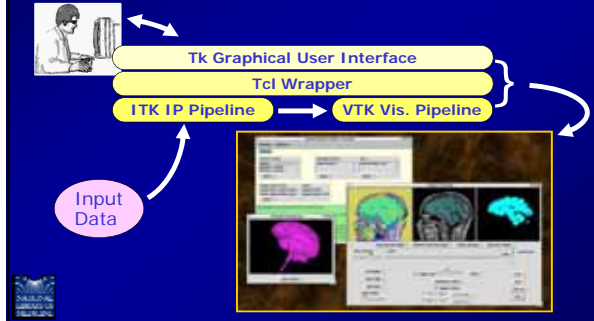


17

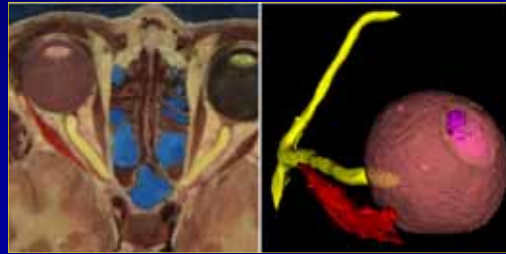
Watershed Segmentation Results



Watershed Segmentation User Interaction



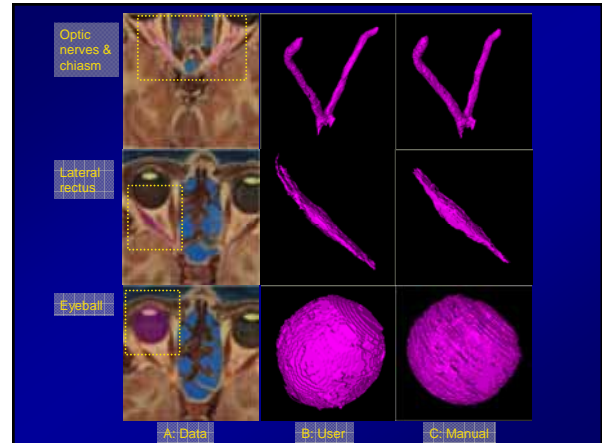
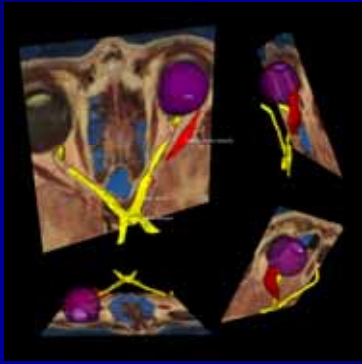
Watershed Segmentation



Interactive results (Visible Human Project female head)

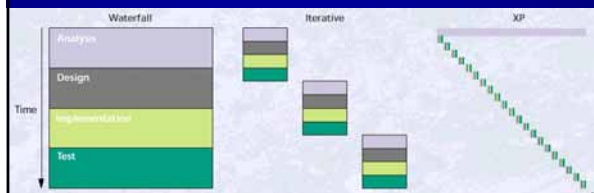
20

Segmentation Results



Extreme Programming

Compresses the standard analyze, design, implement, test cycle into a continuous process



Insight - Open Source Products

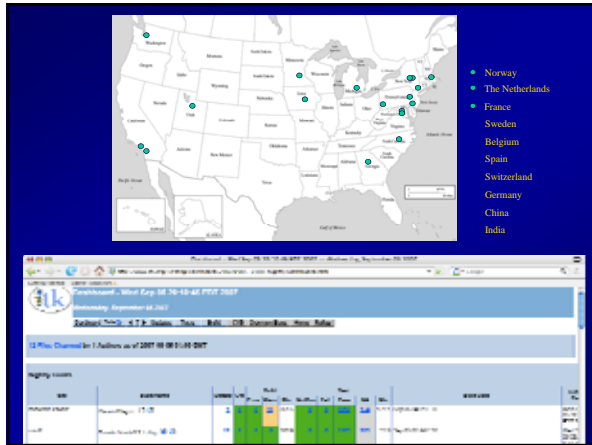
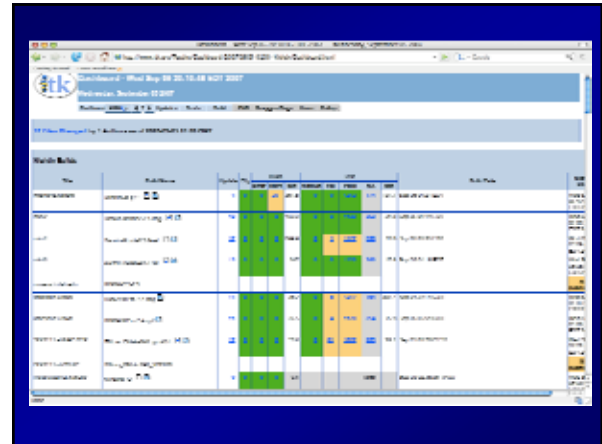


Extreme Programming Daily Testing Is The Key

- Testing anchors and drives the development process (Dart)
- Opens up the development process to everyone
- Developers monitor the testing dashboard constantly
- Problems are identified and fixed immediately
- Developers receive e-mail if they "Break the Build"



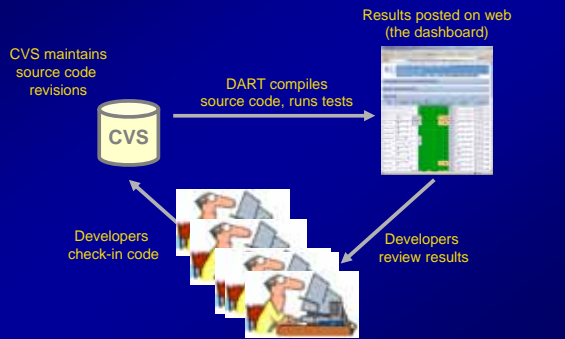
25



Multi-Platform Builds

Site	Build Name	Build Stamp	Status	Time	Detail
g4Cube SINTEF	Carves 6.0.cxx	20031112-0500-4grey	Failed	1:545189	CHILDSTATUS
linnors@stere	Carves 6.0.cxx	20031112-0500-4grey	Failed	1:206179	CHILDSTATUS
hw02_crd	Carves 6.0.cxx	20031112-0500-4grey	Failed	1:57116	CHILDSTATUS
hw02_crd	Carves 6.0.cxx	20031112-0500-4grey	Failed	2:352646	CHILDSTATUS
kraspe@uowa	Carves 6.0.cxx	20031112-0500-4grey	Failed	2:257906	CHILDSTATUS
hw02_crd	Carves 6.0.cxx	20031112-0500-4grey	Failed	2:054248	CHILDSTATUS
geeds_crd	Linux-2.4.19-24.8-1.0	20031112-1105-Continuous	Failed	0:12344	CHILDSTATUS
geeds_crd	Linux-2.4.19-24.8-1.0	20031112-1105-Continuous	Failed	0:114095	CHILDSTATUS
geeds_crd	Linux-2.4.19-24.8-1.0	20031112-1105-Continuous	Failed	0:140022	CHILDSTATUS
hw02_crd	Linux-2.4.19-24.8-1.0	20031112-1105-Continuous	Failed	0:419346	CHILDSTATUS
hw02_crd	Linux-2.4.19-24.8-1.0	20031112-1105-Continuous	Failed	0:356677	CHILDSTATUS
hw02_crd	Linux-2.4.19-24.8-1.0	20031112-1105-Continuous	Failed	0:664372	CHILDSTATUS
hwgen_crd	Linux-2.4.19-24.8-1.0	20031112-1105-Continuous	Failed	0:233957	CHILDSTATUS
hw02_crd	Linux-2.4.19-24.8-1.0	20031112-1105-Continuous	Failed	0:185997	CHILDSTATUS
hw02_crd	Linux-2.4.19-24.8-1.0	20031112-1105-Continuous	Failed	0:170946	CHILDSTATUS
hw02_crd	Linux-2.4.19-24.8-1.0	20031112-1105-Continuous	Failed	0:143543	CHILDSTATUS
hw02_crd	Linux-2.4.19-24.8-1.0	20031112-1105-Continuous	Failed	0:151034	CHILDSTATUS
hw02_crd	Linux-2.4.19-24.8-1.0	20031112-1105-Continuous	Failed	0:333872	CHILDSTATUS
hw02_crd	Linux-2.4.19-24.8-1.0	20031112-1105-Continuous	Failed	0:333353	CHILDSTATUS

How CDash Enables Collaboration



Achievements

- International collaboration of broad research interests and software developers.
- Approximately 4 companies, 7 medical schools, 5 engineering schools.
- Multi-agency \$18.5M Initiative
 - \$2.5M / year - 3 years.
 - \$3.26M - year 4.



30

Achievements (continued)



- Comparable to other Centers and Cooperative Agreements.
- Shared software - deliverables.
- Open Source - extends beyond the consortium; a living contribution to research community.
- Consortium planned for growth.
- An affirmation of our ideals came in 2004.



31

Success



- IGSTK.
- NA-MIC.
- OSIRIX
- Mayo Clinic - Analyze
- MITK - Heidelberg, SCIRun, MeVis, VolView
- Users: Pfizer, Allen Brain Institute
- Linux K Development Environment
- Mayo BIR Collection



32



NA-MIC

The National Alliance for Medical Image Computing



- Committed to open source.
- Shared vision for brain research through imaging.
- Fast start, building from existing tools.
- Likely to accomplish the goal of making lives better.

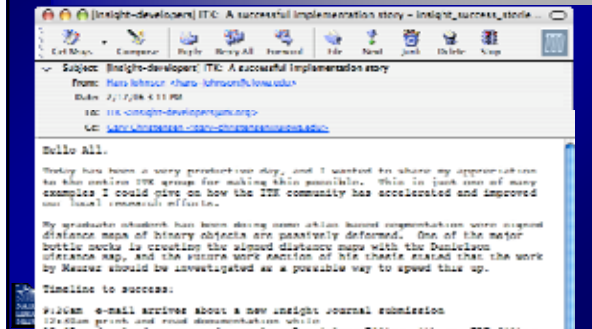


IGSTK

Image-Guided Surgery Toolkit

- Fast start, building from existing tools.
- STTR funded by NIBIB/NIH (Georgetown-Kitware) Grant 2R42EB000374-02
- The CADDLab at UNC joined the project and contributed SpatialObjects and the RF Ablation application
- 2005: Atamai Inc. joined the project and contributed the Tracker code
- March 2007: IGSTK 2.0.1

Success (continued) The Insight Journal



Changes in NIH Culture

- 1987 - NHLBI/Mayo Analyze. NIH Image.
 - Monolithic, proprietary projects
- 1996 - Visible Human Project - open data
- 1999 - ITK: The Insight Toolkit - open source
- 2001 - NCI Lung Image Database Consortium
- 2003 - NCI CAD Open Source Software
- 2003 - NIH Roadmap
- 2004 - NIH BECON/BISTI Symposium
 - Biomedical Informatics for Clinical Decision Support
- 2004 - NIH makes four (4) NCBC awards



36

Principles from the First International Strategy Meeting on Human Genome Sequencing (Bermuda, 25-28 February 1996)

- Primary Genomic Sequence Should be in the Public Domain.
- Primary Genomic Sequence Should be Rapidly Released.
 - Sequence assemblies should be released as soon as possible; in some centres, assemblies of greater than 1 Kb would be released automatically on a daily basis.
 - Finished annotated sequence should be submitted immediately to the public databases.



37

Current NHGRI Policy for Release and Database Deposition of Sequence Data (March 7, 1997)

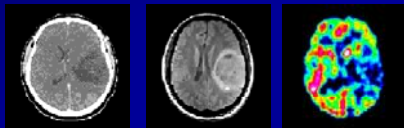
At the Second International Strategy Meeting on Human Genome Sequencing (Bermuda, 1997), attendees affirmed the principle that was set out at the First (1996) International Strategy meeting, that primary genomic sequence should be rapidly released. NHGRI has determined, therefore, that its grantees engaged in large-scale genomic DNA sequencing should now be automatically releasing sequence assemblies of 2 kb or larger within 24 hours of their generation. Any laboratory funded by NHGRI for large-scale human genomic sequencing must develop and submit to NHGRI a plan to implement such a data release program, which must be implemented within one month of its being approved by NHGRI. No non-competing or competing renewal will be funded until an acceptable plan has been approved. Mandatory data release as described above will be made a condition of the award for any grant funded by NHGRI for large-scale human sequencing.

Last Reviewed: August 2006



38

Retrospective Image Registration Evaluation Project (1995 – present) J. Michael Fitzpatrick, P.I.



- The SAT test for rigid, non-deformable multimodal brain registration.
 - Separate training and test data. Provide tests under careful control.
 - Only scores are returned... no solutions.
- Funded across NIH over several installments.
 - Originally NINDS, then NCI, then NIBIB.



39

The 2010 Insight Team “The purpose of computing is insight, not numbers.” - Hamming



ITK Version 4 - Sponsors



QuickTime™ and a decompressor are needed to see this picture.



41

ITK-v4 and ITK-A2D2-2010

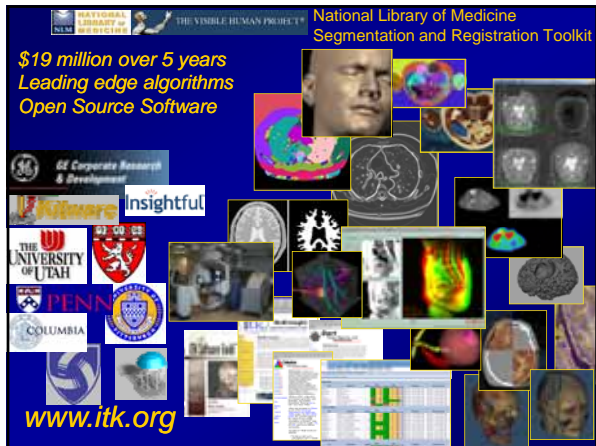
UPGRADE - v4
TEST and EVALUATE - A2D2
Requested to sequential
Grow the community
Lower entry bar for new users
seek new domains

Test - the architecture and the API

We know there's redundancy...
We're here to work it out.



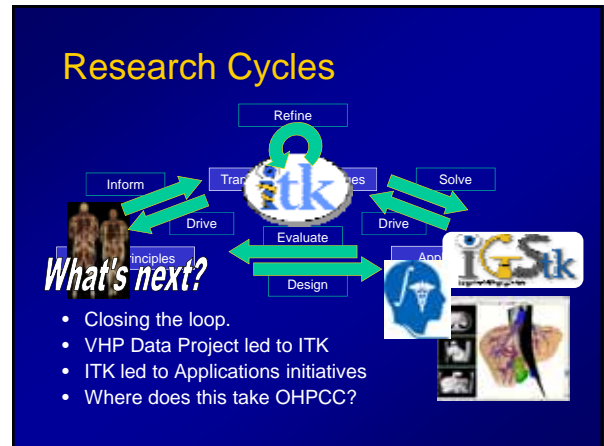
42


 National Library of Medicine
 Segmentation and Registration Toolkit

\$19 million over 5 years
 Leading edge algorithms
 Open Source Software

www.itk.org

Research Cycles



What's next?

- Closing the loop.
- VHP Data Project led to ITK
- ITK led to Applications initiatives
- Where does this take OHPCC?

Open Science

What do you do when your customers are drowning in data?
 What do you do when your engineers are starving for data?
 What do you do when your scientists are reinventing the wheel?

- Infrastructure building.
- Validation, reproducible science.
- Open Source. Open Data. Transparency.

45