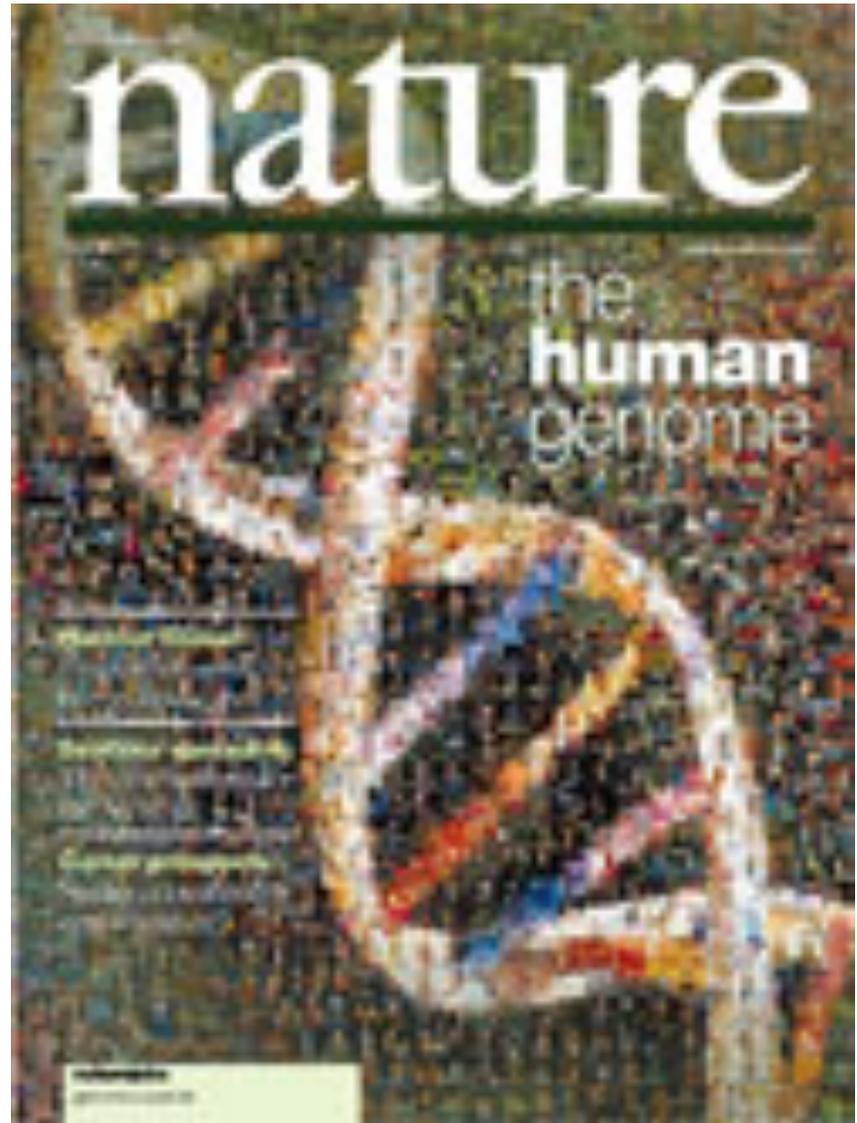


Functional Big-Data Genomics

Ashish Agarwal, Sebastien Mondet,
Paul Scheid, Aviv Madar, Richard Bonneau, Jane Carlton, Kris Gunsalus

*Center for Genomics and Systems Biology
Department of Biology
New York University*

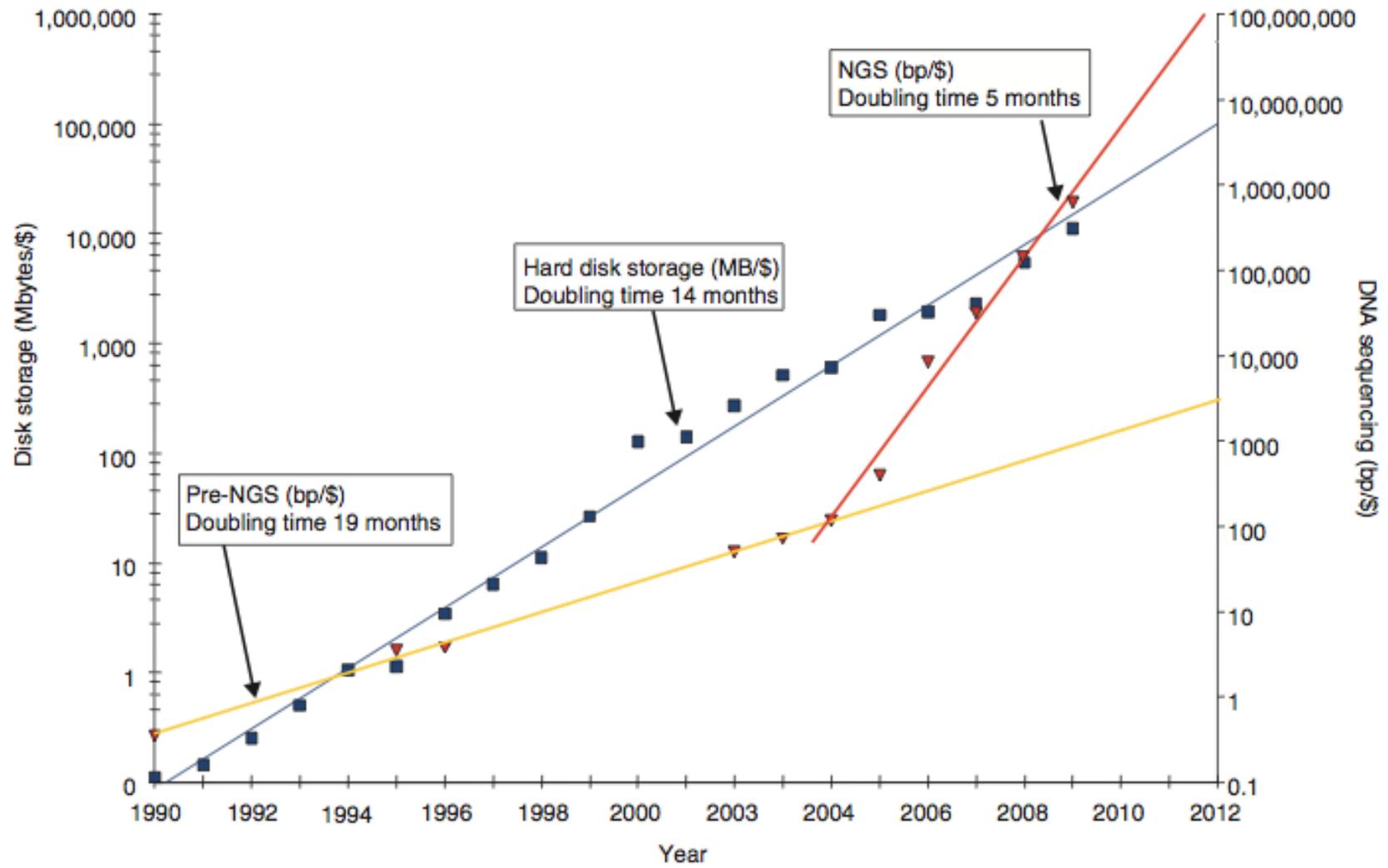
Commercial Users of Functional Programming
Copenhagen, Denmark
Sep 15, 2012

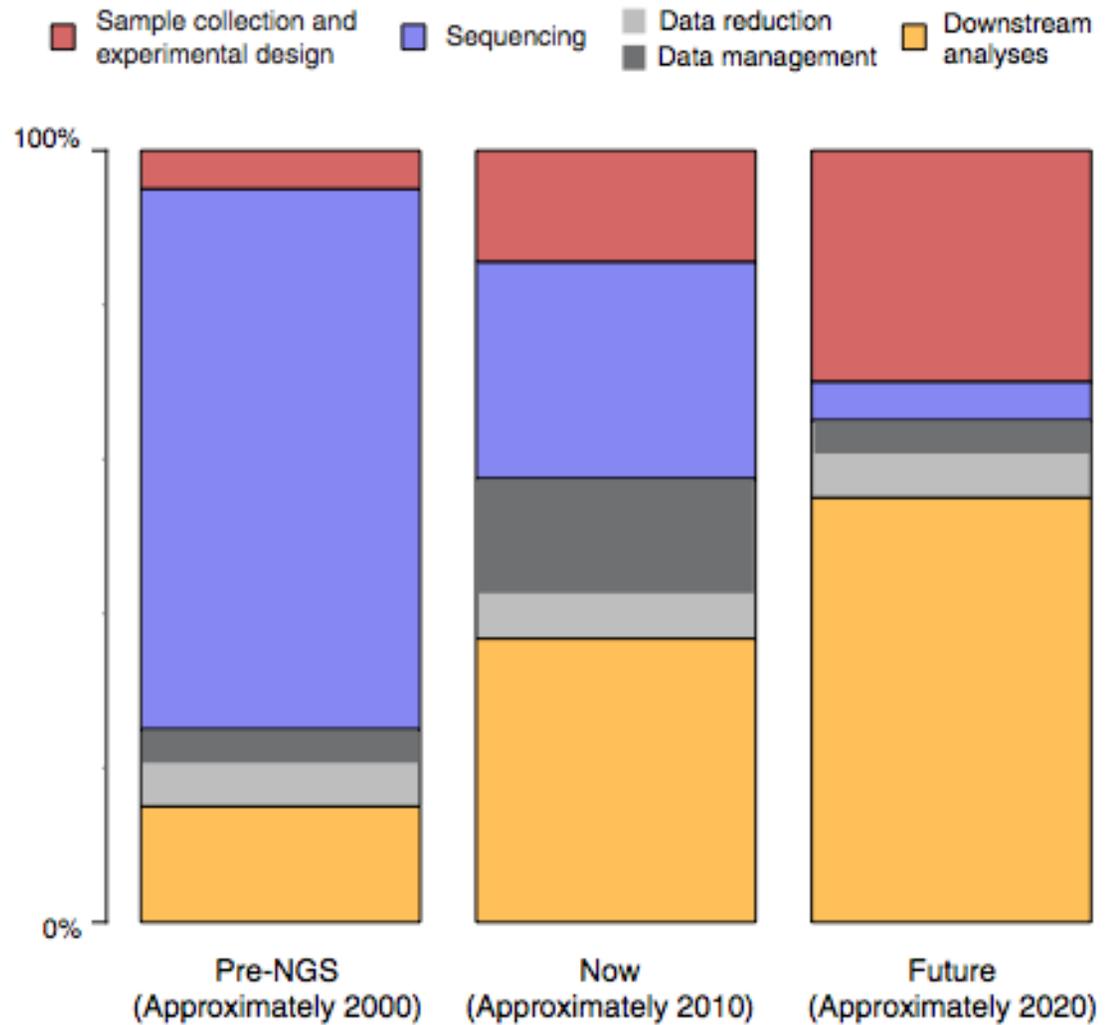
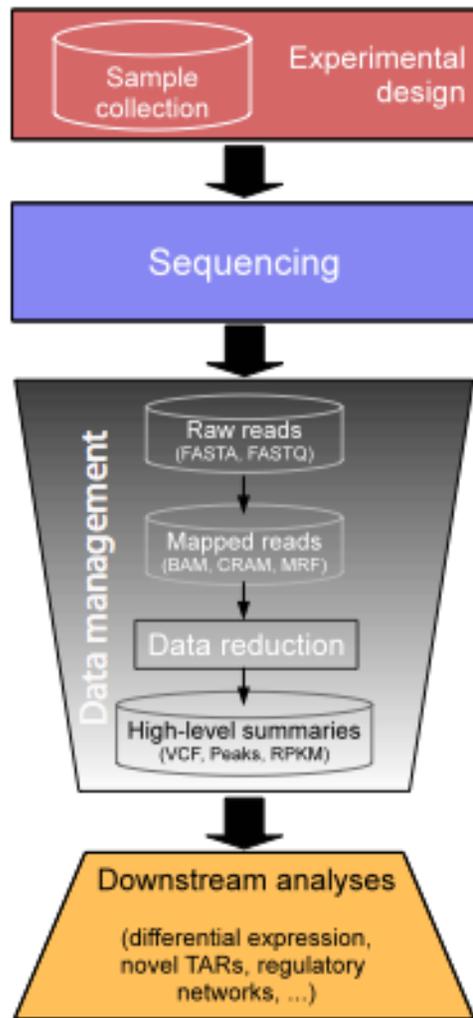


J. Craig Venter^{1,2}, Mark D. Adams¹, Eugene W. Myers¹, Peter W. Li¹, Richard J. Mural¹, Granger G. Sutton¹, Hamilton O. Smith¹, Mark Yandell¹, Cheryl A. Evans¹, Robert A. Holt¹, Jeannine D. Gocayne¹, Peter Amanatides¹, Richard M. Ballew¹, Daniel H. Huson¹, Jennifer Russo Wortman¹, Qing Zhang¹, Chinnappa D. Kodira¹, Xiangqun H. Zheng¹, Lin Chen¹, Marian Skupski¹, Gangadharan Subramanian¹, Paul D. Thomas¹, Jinghui Zhang¹, George L. Gabor Miklos², Catherine Nelson³, Samuel Broder¹, Andrew G. Clark⁴, Joe Nadeau⁵, Victor A. McKusick⁶, Norton Zinder⁷, Arnold J. Levine⁷, Richard J. Roberts⁸, Mel Simon⁹, Carolyn Slayman¹⁰, Michael Hunkapiller¹¹, Randall Bolanos¹, Arthur Delcher¹, Ian Dew¹, Daniel Fasulo¹, Michael Flanagan¹, Liliana Florea¹, Aaron Halpern¹, Sridhar Hannenhalli¹, Saul Kravitz¹, Samuel Levy¹, Clark Mobarry¹, Knut Reinert¹, Karin Remington¹, Jane Abu-Threideh¹, Ellen Beasley¹, Kendra Biddick¹, Vivien Bonazzi¹, Rhonda Brandon¹, Michele Cargill¹, Ishwar Chandramouliswaran¹, Rosane Charlab¹, Kabir Chaturvedi¹, Zuoming Deng¹, Valentina Di Francesco¹, Patrick Dunn¹, Karen Eilbeck¹, Carlos Evangelista¹, Andrei E. Gabrielian¹, Weiniu Gan¹, Wangmao Ge¹, Fangcheng Gong¹, Zhiping Gu¹, Ping Guan¹, Thomas J. Heiman¹, Maureen E. Higgins¹, Rui-Ru Ji¹, Zhaoxi Ke¹, Karen A. Ketchum¹, Zhongwu Lai¹, Yiding Lei¹, Zhenya Li¹, Jiayin Li¹, Yong Liang¹, Xiaoying Lin¹, Fu Lu¹, Gennady V. Merkulov¹, Natalia Milshina¹, Helen M. Moore¹, Ashwini Kumar K Naik¹, Vaibhav A. Narayan¹, Beena Neelam¹, Deborah Nusskern¹, Douglas B. Rusch¹, Steven Salzberg¹², Wei Shao¹, Bixiong Shue¹, Jingtao Sun¹, Zhen Yuan Wang¹, Aihui Wang¹, Xin Wang¹, Jian Wang¹, Ming-Hui Wei¹, Ron Wides¹³, Chunlin Xiao¹, Chunhua Yan¹, Alison Yao¹, Jane Ye¹, Ming Zhan¹, Weiqing Zhang¹, Hongyu Zhang¹, Qi Zhao¹, Liansheng Zheng¹, Fei Zhong¹, Wenyang Zhong¹, Shiaoqing C. Zhu¹, Shaying Zhao¹², Dennis Gilbert¹, Suzanna Baumhueter¹, Gene Spier¹, Christine Carter¹, Anibal Cravchik¹, Trevor Woodage¹, Feroze Ali¹, Huijin An¹, Aderonke Awe¹, Danita Baldwin¹, Holly Baden¹, Mary Barnstead¹, Ian Barrow¹, Karen Beeson¹, Dana Busam¹, Amy Carver¹, Angela Center¹, Ming Lai Cheng¹, Liz Curry¹, Steve Danaher¹, Lionel Davenport¹, Raymond Desilet¹, Susanne Dietz¹, Kristina Dodson¹, Lisa Doup¹, Steven Ferreira¹, Neha Garg¹, Andres Gluecksmann¹, Brit Hart¹, Jason Haynes¹, Charles Haynes¹, Cheryl Heiner¹, Suzanne Hladun¹, Damon Hostin¹, Jarrett Houck¹, Timothy Howland¹, Chinyere Ibegwam¹, Jeffery Johnson¹, Francis Kalush¹, Lesley Kline¹, Shashi Koduru¹, Amy Love¹, Felecia Mann¹, David May¹, Steven McCawley¹, Tina McIntosh¹, Ivy McMullen¹, Mee Moy¹, Linda Moy¹, Brian Murphy¹, Keith Nelson¹, Cynthia Pfannkoch¹, Eric Pratts¹, Vinita Puri¹, Hina Qureshi¹, Matthew Reardon¹, Robert Rodriguez¹, Yu-Hui Rogers¹, Deanna Romblad¹, Bob Ruhfel¹, Richard Scott¹, Cynthia Sitter¹, Michelle Smallwood¹, Erin Stewart¹, Renee Strong¹, Ellen Suh¹, Reginald Thomas¹, Ni Ni Tint¹, Sukyee Tse¹, Claire Vech¹, Gary Wang¹, Jeremy Wetter¹, Sherita Williams¹, Monica Williams¹, Sandra Windsor¹, Emily Winn-Deen¹, Keriellen Wolfe¹, Jaysree Zaveri¹, Karen Zaveri¹, Josep F. Abril¹⁴, Roderic Guigó¹⁴, Michael J. Campbell¹, Kimmen V. Sjolander¹, Brian Karlak¹, Anish Kejarival¹, Huaiyu Mi¹, Betty Lazareva¹, Thomas Hatton¹, Apurva Narechania¹, Karen Diemer¹, Anushya Muruganujan¹, Nan Guo¹, Shinji Sato¹, Vineet Bafna¹, Sorin Istrail¹, Ross Lippert¹, Russell Schwartz¹, Brian Walenz¹, Shibu Yooseph¹, David Allen¹, Anand Basu¹, James Baxendale¹, Louis Blick¹, Marcelo Caminha¹, John Carnes-Stine¹, Parris Caulk¹, Yen-Hui Chiang¹, My Coyne¹, Carl Dahlke¹, Anne Deslattes Mays¹, Maria Dombroski¹, Michael Donnelly¹, Dale Ely¹, Shiva Esparham¹, Carl Foster¹, Harold Gire¹, Stephen Glanowski¹, Kenneth Glasser¹, Anna Glodek¹, Mark Gorokhov¹, Ken Graham¹, Barry Gropman¹, Michael Harris¹, Jeremy Heil¹, Scott Henderson¹, Jeffrey Hoover¹, Donald Jennings¹, Catherine Jordan¹, James Jordan¹, John Kasha¹, Leonid Kagan¹, Cheryl Kraft¹, Alexander Levitsky¹, Mark Lewis¹, Xiangjun Liu¹, John Lopez¹, Daniel Ma¹, William Majoros¹, Joe McDaniel¹, Sean Murphy¹, Matthew Newman¹, Trung Nguyen¹, Ngoc Nguyen¹, Marc Nodell¹, Sue Pan¹, Jim Peck¹, Marshall Peterson¹, William Rowe¹, Robert Sanders¹, John Scott¹, Michael Simpson¹, Thomas Smith¹, Arlan Sprague¹, Timothy Stockwell¹, Russell Turner¹, Eli Venter¹, Mei Wang¹, Meiyuan Wen¹, David Wu¹, Mitchell Wu¹, Ashley Xia¹, Ali Zandieh¹ and Xiaohong Zhu¹

International Human Genome Sequencing Consortium Eric S. Lander¹, Lauren M. Linton¹, Bruce Birren¹, Chad Nusbaum¹, Michael C. Zody¹, Jennifer Bateman¹, Ken Devon¹, Ken Dewar¹, Michael Doyke¹, William Fitzhugh¹, Rolf Funke¹, Diane Gage¹, Katrina Harris¹, Andrew Hesford¹, John Howland¹, Lisa Kann¹, Jessica Lehoczyk¹, Rosie Levine¹, Paul McEwan¹, Kevin McKernan¹, James Meldrim¹, Jill P. Mesirov¹, Cher Miranda¹, William Morris¹, Jerome Naylor¹, Christina Raymond¹, Mark Rosetti¹, Ralph Santos¹, Andrew Sheridan¹, Carrie Sougnez¹, Nicole Stange-Thomann¹, Nikola Stojanovic¹, Aravind Subramanian¹ & Dudley Wyman¹ for Whitehead Institute for Biomedical Research, Center for Genome Research, Jane Rogers¹, John Suton¹, Rachael Ainscough², Stephan Beck², David Bentley², John Burton², Christopher Clew², Nigel Carter², Alan Coulson², Rebecca Deadman², Panos Deloukas², Andrew Dunham², Ian Dunham², Richard Durbin², Lisa French², Darren Grafham², Simon Gregory², Tim Hubbard², Sean Humphray², Adrienne Hunt², Matthew Jones², Christine Lloyd², Amanda McMurray², Lucy Matthews², Simon Mercer², Sarah Milne², James C. Mullikin², Andrew Mungall², Robert Plumb², Mark Ross², Ratna Showkneen² & Sarah Sims² for The Sanger Centre, Robert H. Waterston³, Richard K. Wilson³, LaDeana W. Hillier³, John D. McPherson³, Marco A. Marra³, Elaine R. Mardis³, Lucinda A. Fulton³, Asif T. Chinwalla³, Kymberlie H. Pepin³, Warren R. Gish³, Stephanie L. Chissole³, Michael C. Wendi³, Kim D. Delehaunty³, Tracie L. Miner³, Andrew Delehaunty³, Jason B. Kramer³, Lisa L. Cook³, Robert S. Fulton³, Douglas L. Johnson³, Patrick J. Minx³ & Sandra W. Clifton³ for Washington University Genome Sequencing Center, Trevor Hawkins⁴, Elbert Branscomb⁴, Paul Predki⁴, Paul Richardson⁴, Sarah Wenning⁴, Tom Slezacek⁴, Norman Doggett⁴, Jan-Fang Cheng⁵, Anne Olsen⁵, Susan Lucas⁵, Christopher Eklund⁵, Edward Uberbacher⁵ & Marvin Frazer⁵ for US DOE Joint Genome Institute, Richard A. Gibbs⁶, Donna M. Muzny⁶, Steven E. Scherer⁶, John B. Bouck⁶, Erica J. Sodergren⁶, Kim C. Worley⁶, Catherine M. Rives⁶, James H. Gorrell⁶, Michael L. Metzker⁶, Susan L. Naylor⁶, Raju S. Kucherlapati⁶, David L. Nelson & George M. Weinstock⁶ for Baylor College of Medicine Human Genome Sequencing Center, Yoshiyuki Sakaki⁷, Asao Fujiyama⁷, Masahira Hattori⁷, Tetsushi Yada⁷, Atsushi Toyoda⁷, Takehiko Itoh⁷, Chiharu Kawagoe⁷, Hideomi Watanabe⁷, Yasushi Totoki⁷ & Todd Taylor⁷ for RIKEN Genomic Sciences Center, Jean Weissenbach⁸, Roland Heilig⁸, William Saunier⁸, Francis Artiguenave⁸, Philippe Brottier⁸, Thomas Bruggé⁸, Eric Pelletier⁸, Catherine Robert⁸ & Patrick Wincker⁸ for Genoscope and CNRS UMR-8030, André Rosenthal⁸, Matthias Platzer⁸, Gerald Nyakatura⁸, Stefan Taudien⁸ & Andreas Rump⁸ for Department of Genome Analysis, Institute of Molecular Biotechnology, Douglas R. Smith⁹, Lynn Doucette-Stamm⁹, Marc Rubinfeld⁹, Keith Weinstock⁹, Hong Mei Lee⁹ & JoAnn Dubois⁹ for GTC Sequencing Center, Huanming Yang¹⁰, Jun Yu¹⁰, Jian Wang¹⁰, Guyang Huang¹⁰ & Jun Gu¹⁰ for Beijing Genomics Institute/Human Genome Center, Leroy Hood¹¹, Lee Rowen¹¹, Anup Madan¹¹ & Shizen Qin¹¹ for Multimegabase Sequencing Center, The Institute for Systems Biology, Ronald W. Davis¹², Nancy A. Federspiel¹², A. Pia Abola¹² & Michael J. Proctor¹² for Stanford Genome Technology Center, Bruce A. Roe¹², Feng Chen¹² & Huaqin Pan¹² for University of Oklahoma's Advanced Center for Genome Technology, Juliane Rasmussen¹³, Hans Lehrach¹³ & Richard Reinhardt¹³ for Max Planck Institute for Molecular Genetics, W. Richard McCombs¹⁴, Melissa de la Bastide¹⁴ & Nellya Dethia¹⁴ for Cold Spring Harbor Laboratory, Liza Annenberg Hazen Genome Center, Helmut Blocker¹⁵, Klaus Hornischer¹⁵ & Gabriele Nordtsiek¹⁵ for GBF—German Research Center for Biotechnology, Richa Agarwal¹⁶, L. Aravind¹⁶, Jeffrey A. Bailey¹⁶, Alex Bateman¹⁶, Serafim Batzoglou¹⁶, Ewan Birney¹⁶, Peer Bork^{16,20}, Daniel G. Brown¹⁶, Christopher B. Burge¹⁶, Lorenzo Cenutti¹⁶, Hsiu-Chuan Chen¹⁶, Deanna Church¹⁶, Michele Clamp¹⁶, Richard R. Copley¹⁶, Tobias Doering^{16,20}, Sean R. Eddy¹⁶, Evan E. Eichler¹⁶, Terrence S. Furey¹⁶, James Galagan¹⁶, James G. R. Gilbert¹⁶, Cyrus Harmon¹⁶, Yoshihide Hayashizaki¹⁶, David Haussler¹⁶, Henning Hermjakob¹⁶, Karsten Hokamp¹⁶, Wonhee Jang¹⁶, L. Steven Johnson¹⁶, Thomas A. Jones¹⁶, Simon Kasir¹⁶, Arek Kasprzyk¹⁶, Scot Kennedy^{16,20}, W. James Kent¹⁶, Paul Kitts¹⁶, Eugene V. Koonin¹⁶, Ian Korf¹⁶, David Kulpa¹⁶, Doron Lancet¹⁶, Todd M. Lowe¹⁶, Aoife McLysaght¹⁶, Tarjei Mikkelson¹⁶, John V. Moran¹⁶, Nicola Mulder¹⁶, Victor J. Pollar¹⁶, Chris P. Ponting¹⁶, Greg Schuler¹⁶, Jörg Schultz¹⁶, Guy Suter¹⁶, Arian F. A. Smit¹⁶, Ella Stupka¹⁶, Joseph Szustakowski¹⁶, Danielle Thierry-Mieg¹⁶, Jean Thierry-Mieg¹⁶, Lukas Wagner¹⁶, John Wallis¹⁶, Raymond Wheeler¹⁶, Alan Williams¹⁶, Yuri I. Wolf¹⁶, Kenneth H. Wolfe¹⁶, Shlaw-Pyng Yang¹⁶ & Ru-Fang Yeh¹⁶ for *Genome Analysis Group (listed in alphabetical order, also includes individuals listed under other headings), Francis Collins¹⁷, Mark S. Guyer¹⁷, Jane Peterson¹⁷, Adam Feisenfeld¹⁷ & Kris A. Wetterstrand¹⁷ for Scientific management: National Human Genome Research Institute, US National Institutes of Health, Richard M. Myers¹⁸, Jeremy Schmutz¹⁸, Mark Dickson¹⁸, Jane Grimwood¹⁸ & David R. Cox¹⁸ for Stanford Human Genome Center, Maynard V. Olson¹⁹, Rajinder Kaul¹⁹ & Christopher Raymond¹⁹ for University of Washington Genome Center, Nobuyoshi Shimizu²⁰, Kazuhiko Kawasaki²⁰ & Shinsai Minoshima²⁰ for Department of Molecular Biology, Keio University School of Medicine, Glen A. Evans²¹, Maria Athanasiou²¹ & Roger Schultz²¹ for University of Texas Southwestern Medical Center at Dallas, Aristides Patrino²² for Office of Science, US Department of Energy, & Michael J. Morgan²² for The Wellcome Trust.

2000:	8 YEARS	\$3,000,000,000
Now:	10 days	\$10,000

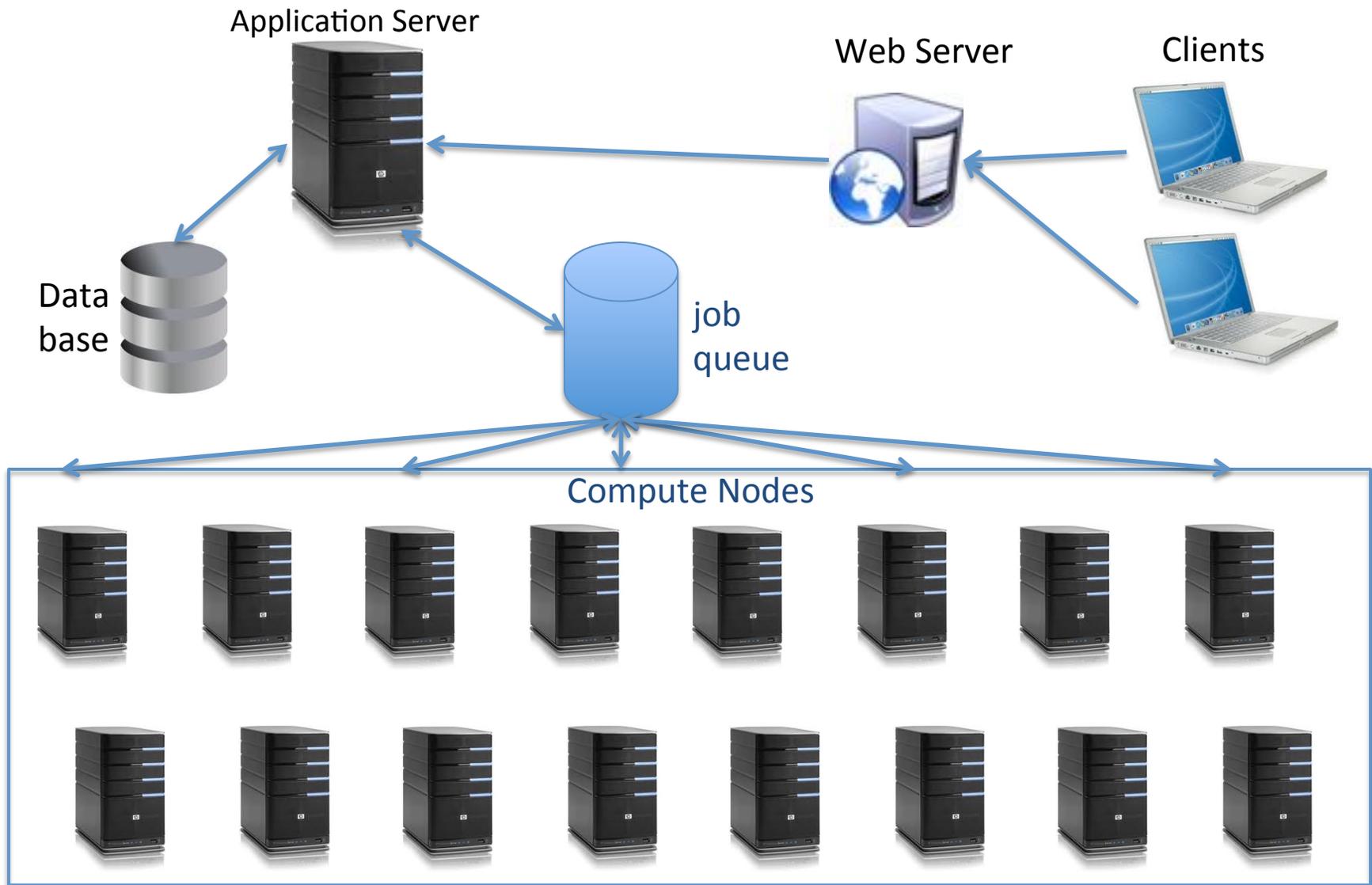




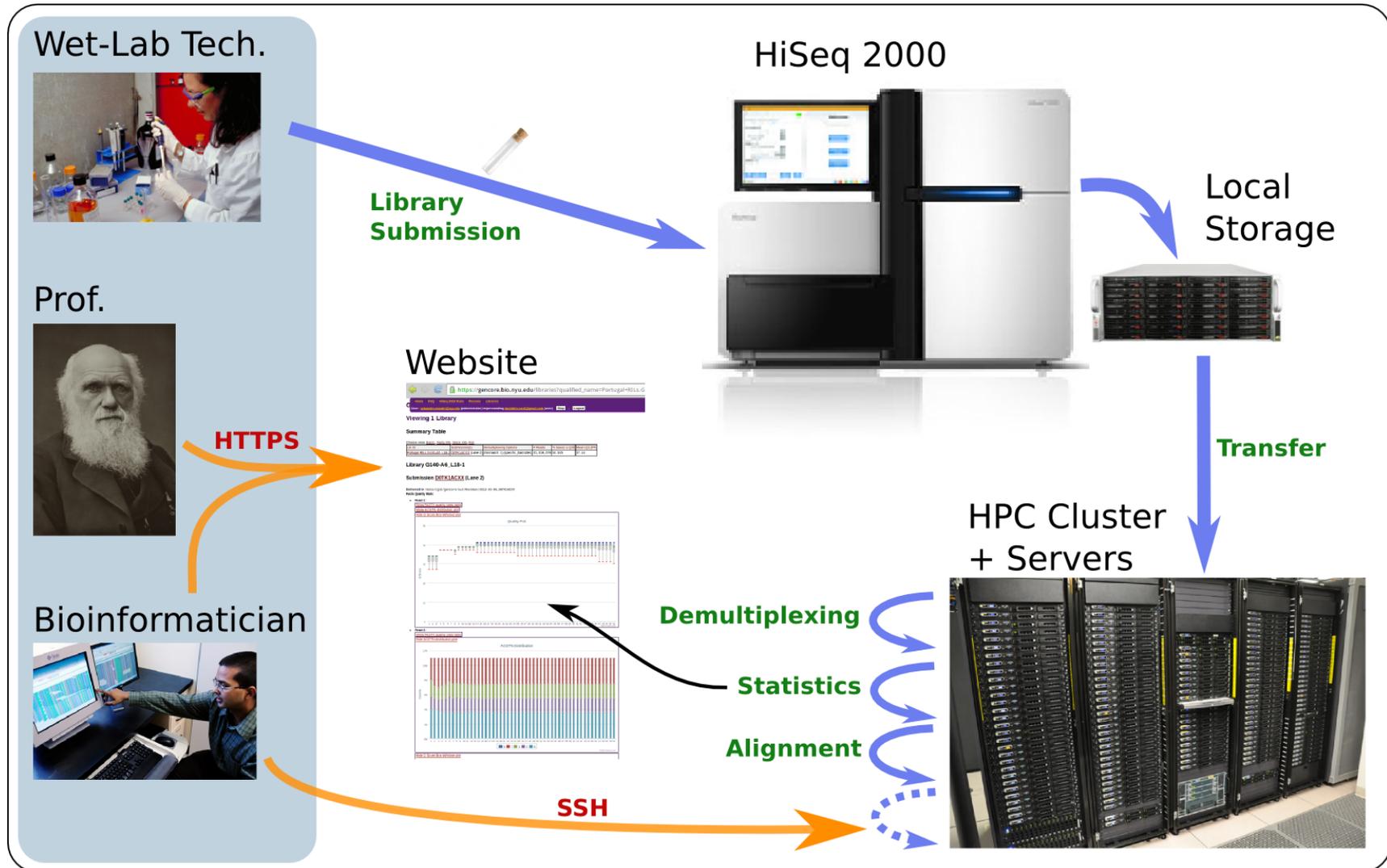
Computational Stack

- Website
 - data visualization, upload/download
- Analysis Pipeline
 - run 3rd party tools, distribute jobs, track results
- LIMS (Laboratory Information Management System)
 - metadata tracking of samples, libraries, protocols, etc
- Systems Infrastructure
 - operating systems, packages, virtual machines, network

Software Architecture



The Whole System



Big-Data is Not Just Big

- **V**olume
- **V**ariety
- **V**elocity

The **Types** of our DSL

type =

| Bool

| Timestamp

| Int

| Real

| String

| Option of type

| Array of type

| Record of (string * type) list

| Enumeration of string list

| Function of type * type

| Volume of volume

- More expressive than SQL
- More closely models the schemas we need
- We have full control over DSL objects

DSL Compiler

- From a given program in the DSL, we *automatically* generate
 - SQL scripts to initialize and reset the database
 - About ~10k lines (and growing) of OCaml strongly typed database reads/inserts
 - Web widgets
 - Figures in this talk

A Program in our DSL

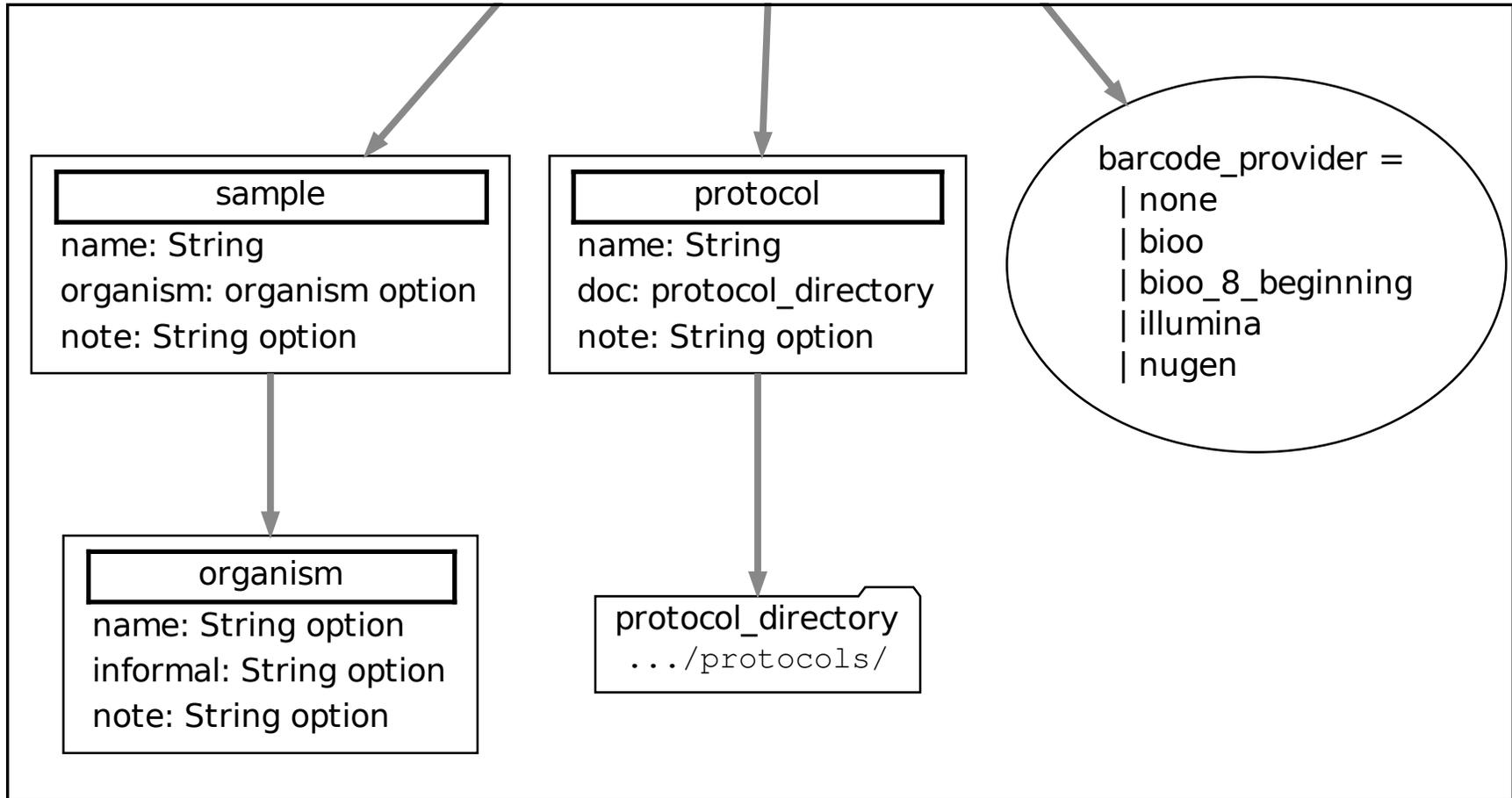
(record sample

 (name string)

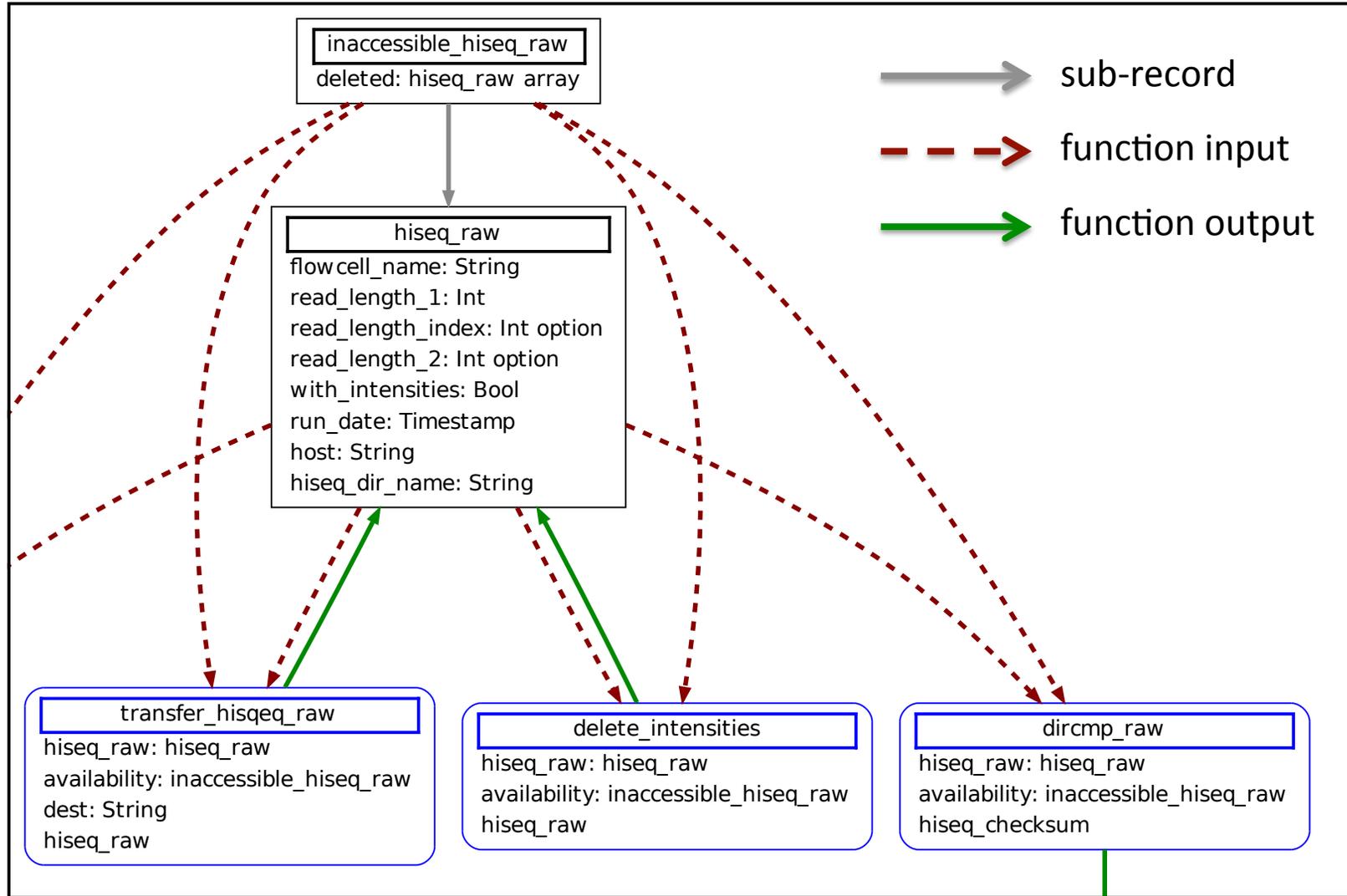
 (organism organism option)

 notable)

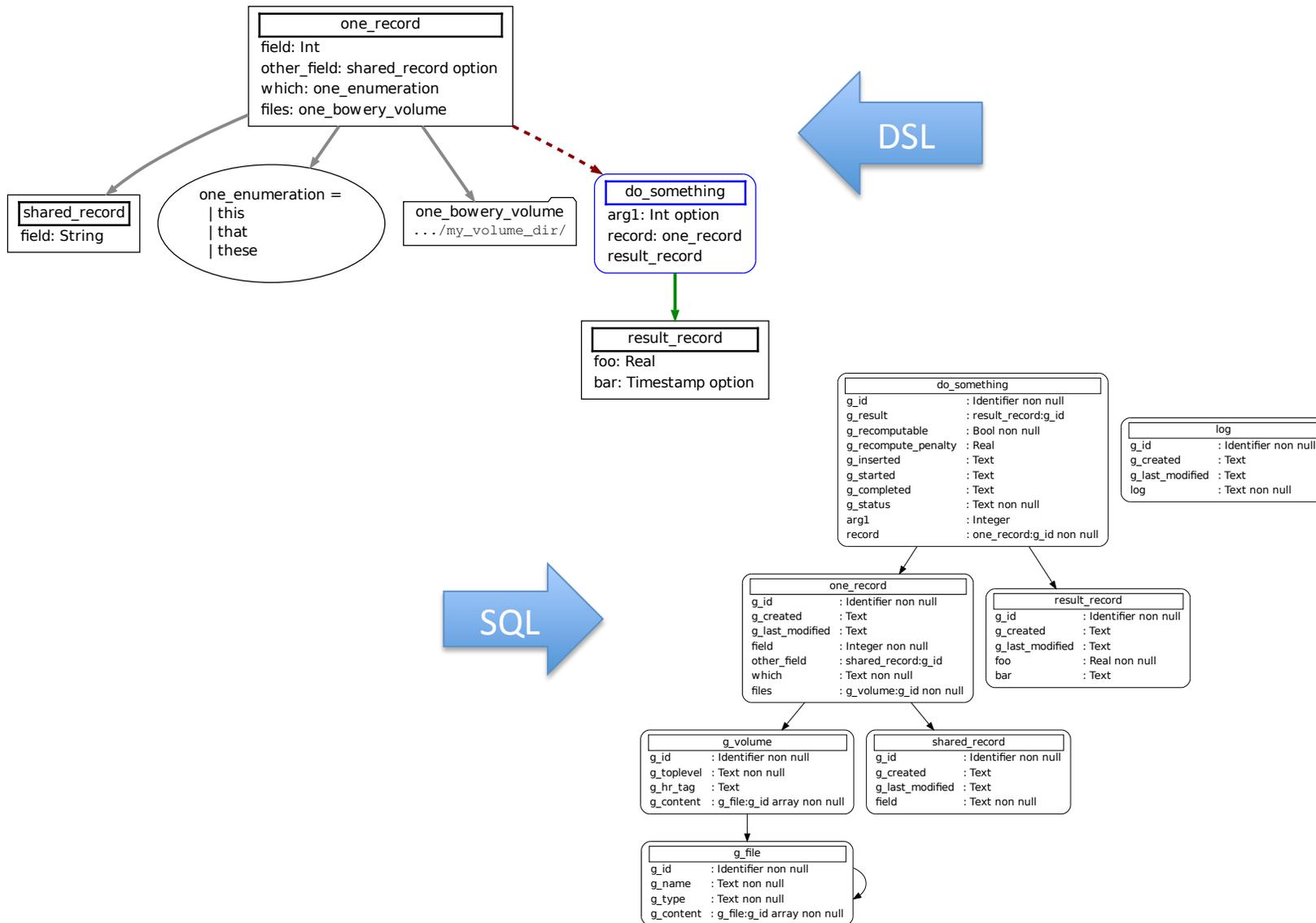
Samples, Organisms, and Protocols



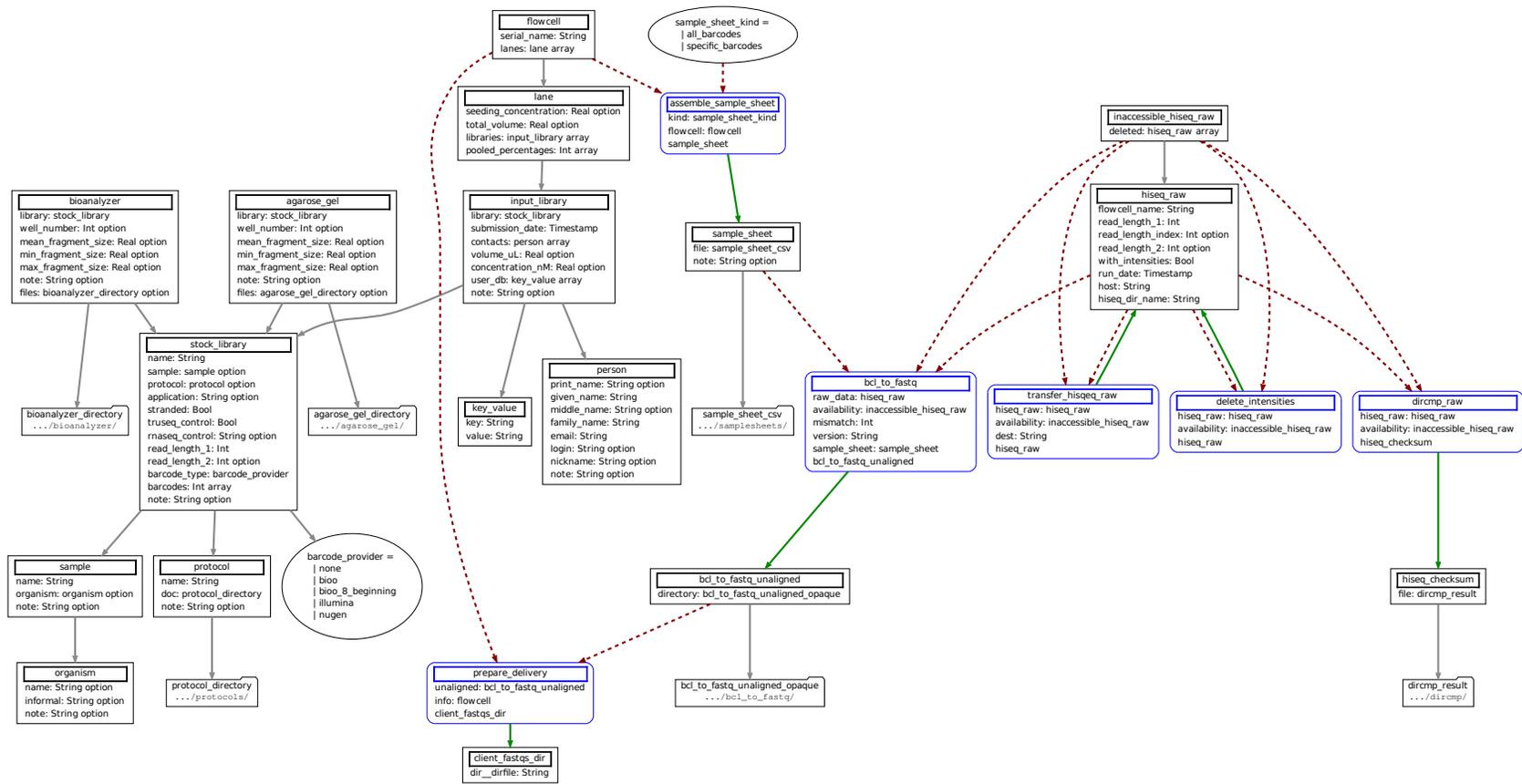
Managing HiSeq Data



Generating SQL



Entire NYU GenCore Data Schema (in our new DSL)



Virtual Filesystem

- Database entries have to map to physical files
- Files in DSL are often *recomputable*
- But don't have to exist on physical filesystem, or can exist in compressed form
- File metadata is recorded
 - e.g. FASTQ quality-score encoding = 33 or 64
- Easier file restructuring/renaming
- Easier migration to new clusters

Function Values

- No lambdas
- Only hardcoded function constants
- Output values are “normal” data
- Implementation must account for
 - failures
 - time lags
 - serialization

Website – powered by Ocsigen

Home FAQ HiSeq 2000 Runs Persons Libraries Function evaluations Layout Navigaditor

User: ashish.agarwal@nyu.edu (administrator, user); Impersonate someone else: Start ;

Logout

GENCORE HOME

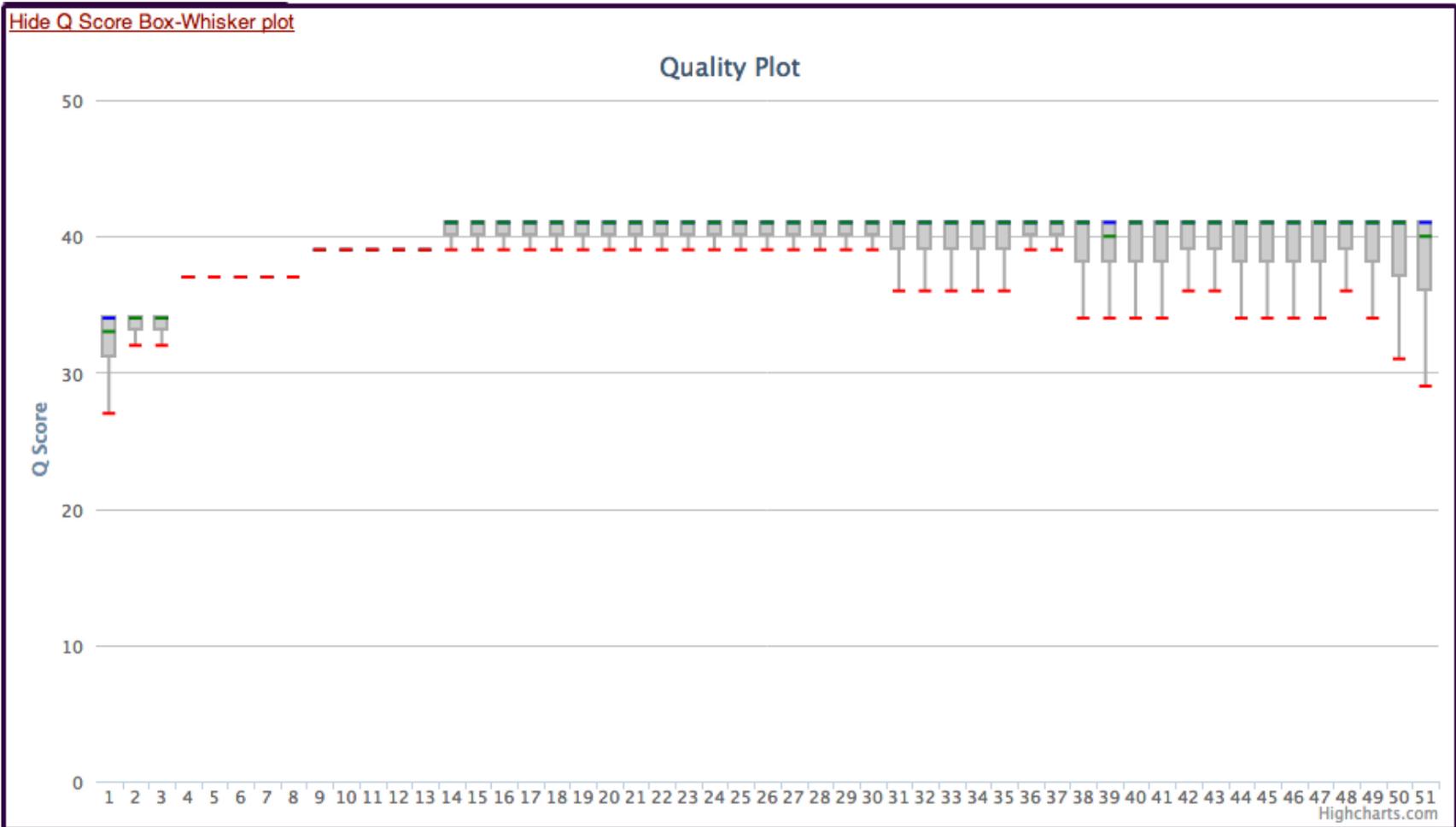
Welcome

This is Gencore's website; for library submission information see [the FAQ](#) or [GenCore's Google-Docs](#).

Menu

- [HiSeq 2000 Runs](#)
- [Persons](#)
- [Libraries](#)
- [Function evaluations](#)
- [Layout Navigaditor](#)

Hide Q Score Box-Whisker plot



Auto-generated Forms

Home FAQ HiSeq 2000 Runs Persons Libraries Function evaluations Layout Navigator

User: [sebastien.mondet@nyu.edu](#) (administrator); Impersonate someone else: Start ; Logout

- You may add a new person
- You may modify this person

Identifier	3179
g_created: Timestamp	2011-12-13
g_last_modified: Timestamp	2012-05-29
S-Exp: String	print_name ("Sebastien (family_name Mondet) (sebastien.mondet@nyu.edu) (sebastien.mondet.org) (log
print_name: String option	Sebastien Mondet
given_name: String	Sebastien
middle_name: String option	
family_name: String	Mondet
email: String	sebastien.mondet@nyu.edu
secondary_emails: String array	seb@mondet.org
login:	em4421

Actions:

- [You may add a new person](#)
- [You may modify this person](#)

g_id: Identifier 3179

Data Migrations

val migrator: S-Expression -> S-Expression

\$ hitscore dump-to-file backup_v42

\$./migrator backup_v42 backup_v43

\$ hitscore wipe-out-database

\$ hitscore init-database

\$ hitscore load-file backup_v43

\$ hitscore verify-layout

We have done over 25 data migrations!

OCaml Has Many Libraries

- Core & Batteries
- Lwt – lightweight threads
- Ocsigen – web programming framework
- Biocaml
- PG'OCaml
- Xmlm
- Ocamlnet
- ... and many more

Database -> Domain Logic -> Web Interface

Experience with OCaml

- The Good
 - Good libraries
 - Industrial strength
 - Hackable
 - Option to be unsafe feels safe
 - Excellent performance
- Could Be Better
 - Public relations
 - Build system
 - Blessed libraries
 - More libraries

Functional Programming in Biology

- “Functional programming” is becoming a recognized term
- Programmers are desperately needed
- Be sure to distinguish *software engineering* from *data analysis*
- Key to success:
 - acquire domain expertise
 - build software fast
- Discuss programming scientifically

Conclusions

- The Genomics Sequencing Core at NYU CGSB runs on OCaml
- Entire system built by ~1.3 programmers
- First version: in production within 2 months
- Biology needs you!