# Improving support and distribution of validated analysis tools

Brad Chapman

Bioinformatics Core, Harvard Chan School

https://bcb.io

http://j.mp/bcbiolinks

9 November 2015

# We need to do science faster



**Karyn MeltzSteinberg**
@KMS_Meltzy
Following

My heart is breaking for friend whose 1 wk old son has been diagnosed w a rare genetic disorder w/o a cure. Motivation to work harder.

FAVORITE
1

9:39 AM - 2 Nov 2015

https://twitter.com/KMS_Meltzy/status/661206070308794368

# We need to incorporate improvements faster

**New human genome assembly (GRCh38) released!**

Tuesday, December 24, 2013

On December 24th, the Genome Reference Consortium (GRC) submitted a new assembly for the human genome (GRCh38) to GenBank. These data are now available in the Assembly database

8

Switch from hg19/build37 to hg20/build38?

(self.genome)

submitted 4 months ago by coopergm

I am curious to what extent there is interest among people that routinely use the reference assembly and associated data (variant datasets, functional genomic annotations, conservation, what-have-you) to change from hg19 to hg20.

https://www.reddit.com/r/genome/comments/3b3s3t/switch_from_hg19build37_to_hg20build38/

- Install tools
- Put tools together
- Test and validate
- Improve algorithms
- Scale
- Read literature
- Do biology

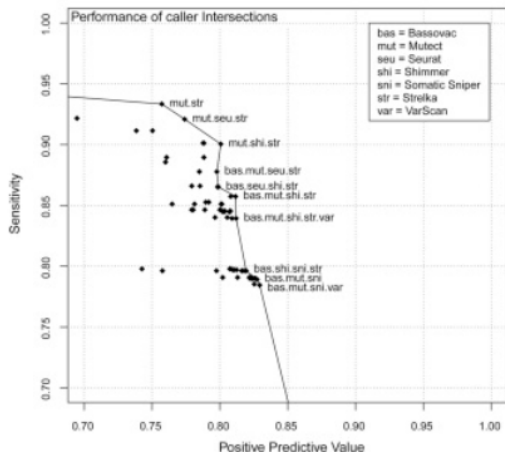*Four major genome centers predicted single-nucleotide variants (SNVs) for The Cancer Genome Atlas (TCGA) lung cancer samples, but only 31.0% (1,667/5,380) of SNVs were identified by all four.*

http://www.nature.com/nmeth/journal/vaop/ncurrent/full/nmeth.3407.html

# Combining analyses



D  **Multiple variant callers**

Performance of caller Intersections

bas = Bassovac
mut = Mutect
seu = Seurat
shi = Shimmer
sni = Somatic Sniper
str = Strelka
var = VarScan

# Working together produces great things

**ExAC Principal Investigators**
- Daniel MacArthur
- David Altshuler
- Diego Ardissino
- Michael Boehnke
- Mark Daly
- John Danesh
- Roberto Elosua
- Jose Florez
- Gad Getz
- Christina Hultman
- Sekar Kathiresan
- Markku Laakso
- Steven McCarroll
- Mark McCarthy
- Dermot McGovern
- Ruth McPherson
- Benjamin Neale
- Aarno Palotie
- Shaun Purcell
- Danish Saleheen
- Jeremiah Scharf
- Pamela Sklar
- Patrick Sullivan
- Jaakko Tuomilehto
- Hugh Watkins
- James Wilson

**Contributing projects**
- 1000 Genomes
- Bulgarian Trios
- Finland-United States Investigation of NIDDM Genetics (FUSION)
- GoT2D
- Inflammatory Bowel Disease
- METabolic Syndrome In Men (METSIM)
- Jackson Heart Study
- Myocardial Infarction Genetics Consortium:
  - Italian Atherosclerosis, Thrombosis, and Vascular Biology Working Group
  - Ottawa Genomics Heart Study
  - Pakistan Risk of Myocardial Infarction Study (PROMIS)
  - Precocious Coronary Artery Disease Study (PROCARDIS)
  - Registre Gironi del COR (REGICOR)
- NHLBI-GO Exome Sequencing Project (ESP)
- National Institute of Mental Health (NIMH) Controls
- SIGMA-T2D
- Sequencing in Suomi (SISu)
- Swedish Schizophrenia & Bipolar Studies
- T2D-GENES
- Schizophrenia Trios from Taiwan
- The Cancer Genome Atlas (TCGA)
- Tourette Syndrome Association International Consortium for Genomics (TSAICG)

**Production team**
- Monkol Lek
- Fengmei Zhao
- Ryan Poplin
- Eric Banks
- Timothy Fennell

**Analysis team**
- Monkol Lek
- Kaitlin Samocha
- Konrad Karczewski
- Eric Minikel
- James Ware
- Anne O'Donnell Luria
- Andrew Hill
- Beryl Cummings
- Daniel Birnbaum
- Taru Tukiainen
- Laramie Duncan
- Karol Estrada
- Menachem Fromer
- Adam Kiezun
- Mitja Kurki
- Ron Do
- Pradeep Natarajan
- Gina Peloso
- Hong-Hee Won

**Website team**
- Konrad Karczewski
- Brett Thomas
- Daniel Birnbaum
- Ben Weisburd

**Ethics team**
- Stacey Donnelly
- Andrea Saltzman
- Namrata Gupta
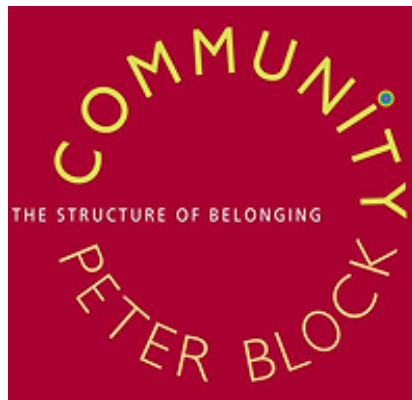
**Broad Genomics Platform**
- Stacey Gabriel

Many thanks to the Genomics Platform both for generating much of the exome data displayed here and for providing the computing resources required for this analysis.

**Funding**
- NIGMS R01 GM104371 (PI: MacArthur)
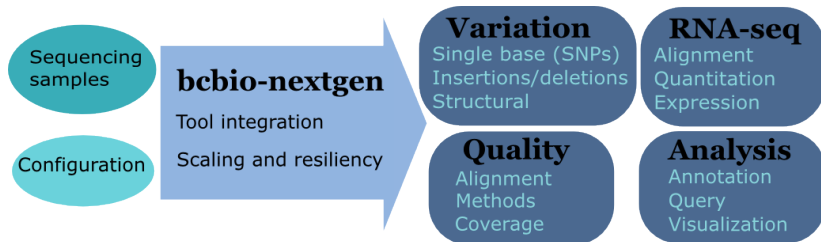- NIDDK U54 DK105566 (PIs: MacArthur and Neale)

http://exac.broadinstitute.org/about

# Solution

- Find shared problems
- Community developed analyses
- Validation
- Scaling
- Supporting a community of users

https://github.com/chapmanb/bcbio-nextgen

- Aligners: bwa, novoalign, bowtie2
- Variantion: FreeBayes, GATK, VarDict, MuTecT, Scalpel, SnpEff, VEP, GEMINI, Lumpy, Manta, CNVkit
- RNA-seq: Tophat, STAR, cufflinks, Sailfish, HISAT2
- Quality control: fastqc, bamtools, RNA-SeQC
- Manipulation: bedtools, bcftools, biobambam, sambamba, samblaster, samtools, vcflib, vt

- Community – collected set of expertise
- Tool integration
- Validation – outputs + automated evaluation
- Scaling
- Installation of tools and data

# We made a pipeline – so what?

*There have been a number of previous efforts to create publicly available analysis pipelines for high throughput sequencing data. Examples include Omics-Pipe, bcbio-nextgen, TREVA and NGSane. These pipelines offer a comprehensive, automated process that can analyse raw sequencing reads and produce annotated variant calls. However, the main audience for these pipelines is the research community. Consequently, there are many features required by clinical pipelines that these examples do not fully address. Other groups have focused on improving specific features of clinical pipelines. The Churchill pipeline uses specialised techniques to achieve high performance, while maintaining reproducibility and accuracy. However it is not freely available to clinical centres and it does not try to improve broader clinical aspects such as detailed quality assurance reports, robustness, reports and specialised variant filtering. The Mercury pipeline offers a comprehensive system that addresses many clinical needs: it uses an automated workflow system (Valence) to ensure robustness, abstract computational resources and simplify customisation of the pipeline. Mercury also includes detailed coverage reports provided by ExCID, and supports compliance with US privacy laws (HIPAA) when run on DNANexus, a cloud computing platform specialised for biomedical users. Mercury offers a comprehensive solution for clinical users, however it does not achieve our desired level of transparency, modularity and simplicity in the pipeline specification and design. Further, Mercury does not perform specialised variant filtering and prioritisation that is specifically tuned to the needs of clinical users.*

http://www.genomemedicine.com/content/7/1/68

A piece of software is being sustained if people are using it, fixing it, and improving it rather than replacing it.
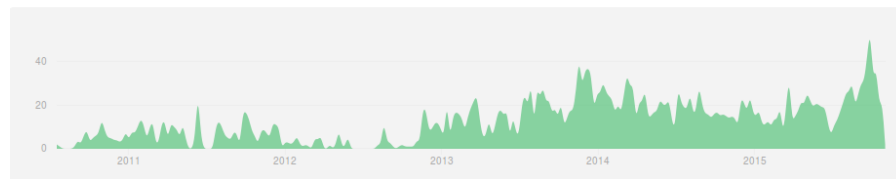
http://software-carpentry.org/blog/2014/08/
sustainability.html

# Community: sustainability



Jul 18, 2010 – Nov 2, 2015

Contributions to master, excluding merge commits

Contributions: **Commits**

https://github.com/chapmanb/bcbio-nextgen

# Community: support



https://bcbio-nextgen.readthedocs.org

## What can we replace?

- Installation
- Infrastructure – runs on your cluster
- Tool integration
- Validation – stability
- Rapid development – new improvements

# Installation



**John Davey**
@johnomics

The trepidation of opening an INSTALL file.
"Please say ./configure; make; make
install… please say ./configure; make; make
install…"

↩ Reply  ♺ Retweet  ★ Favorite  ••• More

## Automated Install

We made it easy to install a large number of biological tools.
Good or bad idea?

- bcbio tools + code in Docker containers
- Bootstrap from plain AMIs to cluster
- Pull/push data from S3
- Lustre and encrypted NFS filesystems

http://bcb.io/2014/12/19/awsbench/
https://github.com/chapmanb/bcbio-nextgen-vm

# Common Workflow Language

- Standard way to describe workflows
- Explicit markup of input/output files
- Automatically generated by bcbio
- Run on multiple infrastructures
- Community

https://github.com/chapmanb/bcbio-nextgen/tree/master/cwl

## Arvados Core Platform

The Arvados core is a platform for production data science with very large data sets. It is made up of two major systems and a number of related services and components including APIs, SDKs, and visual tools.

### Keep

Keep is a content-addressable storage system for managing and storing large collections of files with durable, cryptographically verifiable references and high-throughput processing. Keep works on a wide range of underlying file systems. Learn More >

### Crunch

Is a containerized workflow engine for running complex, multi-part pipelines or workflows in a way that is flexible, scalable, and supports versioning, reproducibility, and provenance. Crunch runs in virtualized computing environments.

| Interfaces | Arvados Workbench | Command Line Interface | Tools & Pipelines | 3rd Party Web Apps |
|---|---|---|---|---|
| | SDKs | | | |
| | API & Access Control | | | |
| Core Services | Keep | | Crunch | |
| Elastic Computing Foundation | | | | |

https://arvados.org/

# Infrastructure



https://github.com/galaxyproject/planemo
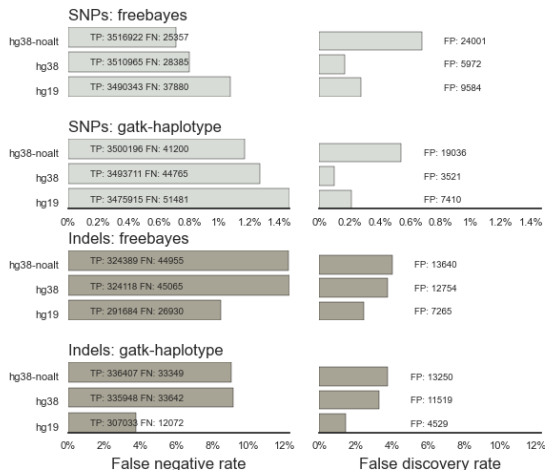
- CWL description + platforms for run
- Docker containers with tools + code
- Mix and match implementations
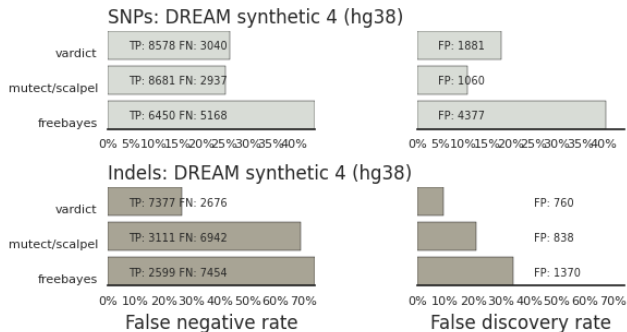- Do research and development and production in same environment

# Practical: Human build 38 validation



hg19/hg38 comparison: NA12878 Platinum Genomes
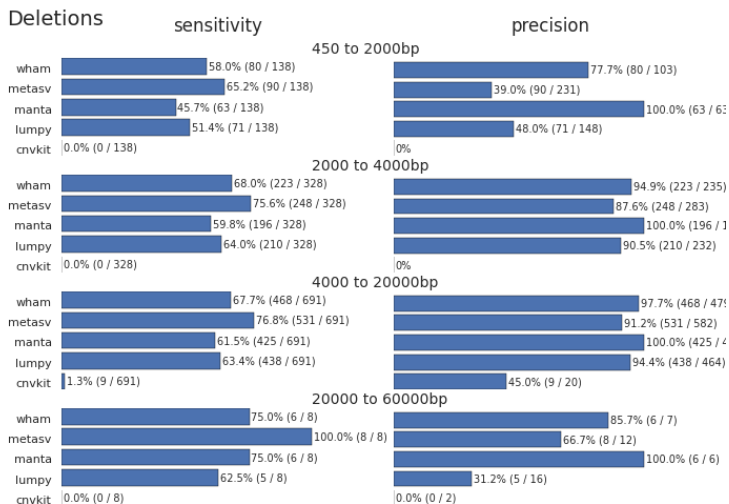
# Practical: cancer validation



SNPs: DREAM synthetic 4 (hg38)

Indels: DREAM synthetic 4 (hg38)

# Practical: structural variant calling

# Summary

- Do more science faster
- Community – integrate, not re-implement
- Docker + CWL enables integration
- Let's talk about ways to work together