# Building community developed open source infrastructure to support large-scale biology research

Brad Chapman

Bioinformatics Core, Harvard School of Public Health

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

http://j.mp/bcbiolinks

12 September 2014

## Buffering of crucial functions by paleologous duplicated genes may contribute cyclicality to angiosperm genome duplication

Brad A. Chapman [*],[†], John E. Bowers [*], Frank A. Feltus [*], and Andrew H. Paterson [*],[†],[‡],[§],[¶]

Author Affiliations ⌄

[*]Plant Genome Mapping Laboratory and Departments of

[†]Plant Biology,

[‡]Genetics, and

[§]Crop and Soil Science, University of Georgia, Athens, GA 30602
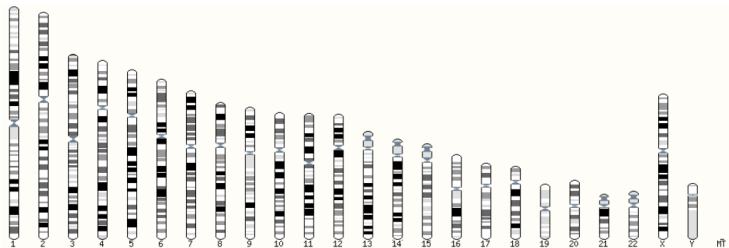
# Synthetic biology startup (2004-2009)



http://www.synthesis.cc/2009/04/on-the-demise-of-condon-devices.html

**HARVARD**
**SCHOOL OF PUBLIC HEALTH**

Powerful ideas for a healthier world

http://compbio.sph.harvard.edu/chb/

- Community developed variant calling analyses
- Validation enables science
- Science at scale: 50 to 1500 genomes
- Supporting a community of users
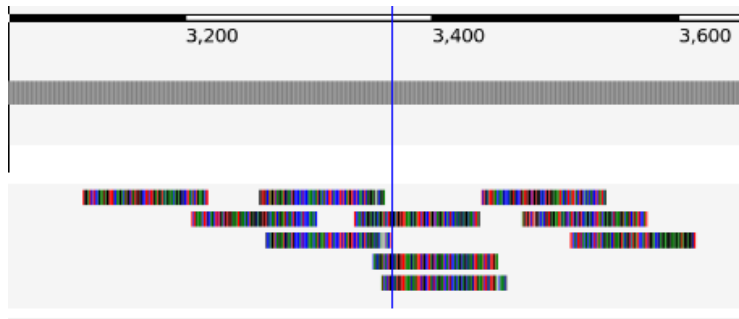- Software development and science

# Human whole genome sequencing



**Click on the image above to jump to a chromosome, or click and drag to select a region**

### Summary

| | |
|---|---|
| **Assembly** | GRCh37.p13 (Genome Reference Consortium Human Reference 37), INSDC Assembly GCA_000001405.14, Feb 2009 |
| **Database version** | 75.37 |
| **Base Pairs** | 3,326,743,047 |

http://ensembl.org/Homo_sapiens/Location/Genome

# High throughput sequencing

# Variant calling



Aligned Reads

Reference

http://en.wikipedia.org/wiki/SNV_calling_from_NGS_data

# Scale: exome to whole genome



**The haploid human genome sequence**

https://www.flickr.com/photos/119980645@N06/

# White box software

# Overview



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

# Uses

- Aligners: bwa-mem, novoalign, bowtie2
- Variantion: FreeBayes, GATK, Platypus, MuTecT, scalpel, SnpEff, VEP, GEMINI, Lumpy, Delly
- RNA-seq: Tophat, STAR, cufflinks, HTSeq
- Quality control: fastqc, bamtools, RNA-SeQC
- Manipulation: bedtools, bcftools, biobambam, sambamba, samblaster, samtools, vcflib

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

# Complex, rapidly changing pipelines

# Large number of specialized dependencies

```
######################################
# HugeSeq                            #
# The Variant Detection Pipeline     #
######################################

-- DEPENDENCIES

+ ANNOVAR version 20110506
+ BEDtools version 2.16.2
+ BreakDancer version 1.1
+ BreakSeq Lite version 1.3
+ BWA version 0.6.1
+ CNVnator version 0.2.2
+ GATK version 1.6-9
+ JDK version 1.6.0_21
+ Modules Release 3.2.8
+ Perl
+ Picard Tools version 1.64
+ Pindel version 0.2.2
+ Plantation version 2
+ pysam version 0.6
+ Python version 2.7
+ Simple Job Manager version 1.0
+ Tabix version 0.1.5
+ VCFtools version 0.1.5
```

https://github.com/StanfordBioinformatics/HugeSeq

# Solution



http://www.amazon.com/Community-Structure-Belonging-Peter-Block/dp/
1605092770

# Community: contribution



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

# Community: documentation



https://bcbio-nextgen.readthedocs.org

Tests for implementation and methods

- Family/population calling
- RNA-seq differential expression
- Structural variations
- Cancer tumor/normal
  http://j.mp/cancer-var-chal
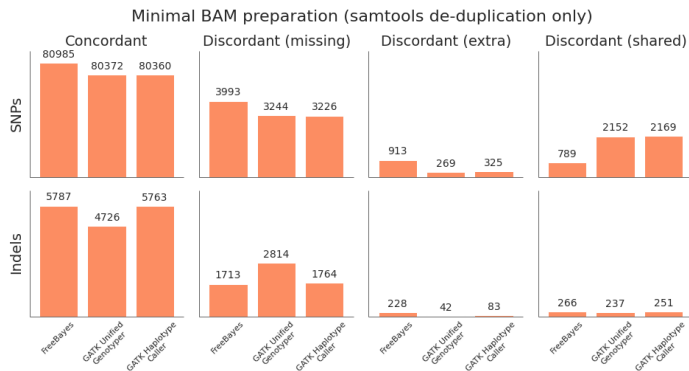
# Example evaluation

- Variant calling
  - GATK UnifiedGenotyper
  - GATK HaplotypeCaller
  - FreeBayes
- Two preparation methods
  - Full (de-duplication, recalibration, realignment)
  - Minimal (only de-duplication)

http://www.genomeinabottle.org/

# Quantify quality



Minimal BAM preparation (samtools de-duplication only)

- Quantification details: http://j.mp/bcbioeval2

- Little value in realignment when using haplotype aware caller
- Little value in recalibration when using high quality reads
- Streaming de-duplication approaches provide same quality without disk IO

- Initial pipeline scales with exomes
- 50 whole genomes = 3 months
- Next project: 1500 whole genomes

1500 whole genome scale – 110Tb

```
$ du -sh alz-p3f_2-g5/final
3.4T  alz-p3f_2-g5/final
$ ls -lhd *alz* | wc -l
31
```
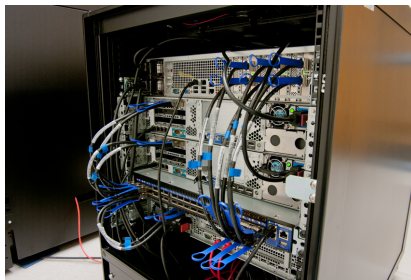
- Network bandwidth
- Better shared filesystems: Lustre
- Avoid file intermediates
- Parallel alignment
- Parallel genome processing

1 GigE to Infiniband



Dell Genomic Data Analysis Platform; Glen Otero

http://www.dell.com/learn/us/en/555/hpcc/

high-performance-computing-life-sciences?c=us&l=en&s=biz&cs=555

480 cores, 30 samples

| Step | Lustre | NFS |
|------|--------|-----|
| alignment | 4.5h | 6.1h |
| alignment post-processing | 7.0h | 20.7h |

```
("{bwa} mem -M -t {num_cores} -R '{rg_info}' -v 1 "
 "  {ref_file} {fastq_file} {pair_file} "
 "| {samblaster} "
 "| {samtools} view -S -u /dev/stdin "
 "| {sambamba} sort -t {cores} -m {mem} --tmpdir {tmpdir}"
 "   -o {tx_out_file} /dev/stdin")
```
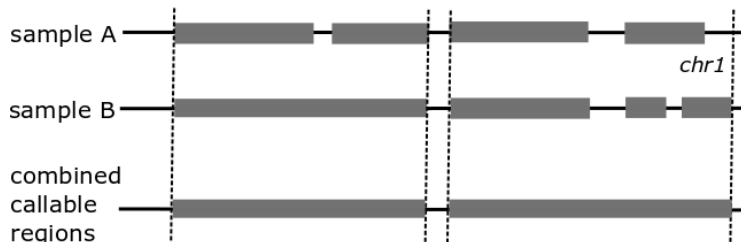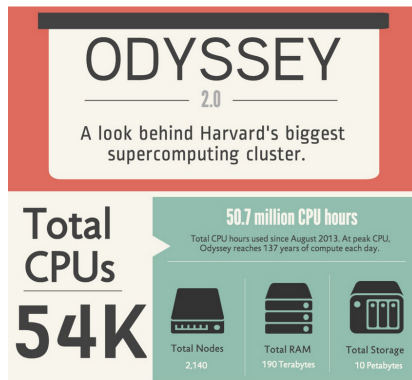
# Scaling: Parallel alignment



https://github.com/arq5x/grabix

Selection of genome regions for parallel processing

James Cuff, John Morrissey, Kristina Kermanshahche
https://rc.fas.harvard.edu/

# Make installation easy



**John Davey**
@johnomics

⚙ Following

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?

# Need a consistent support environment

- Fully isolated
- Reproducible – store full environment with analysis (1Gb)
- Improved installation – single download + data

- External Python wrapper
  - Installation
  - Start and run containers
  - Mount external data into containers
  - Parallelize
- All analysis tools inside Docker

https://github.com/chapmanb/bcbio-nextgen-vm
http://j.mp/bcbiodocker

http://software-carpentry.org
http://mozillascience.org

http://github.com
https://bitbucket.org

http://ipython.org
http://www.rstudio.com/

http://www.open-bio.org
http://www.open-bio.org/wiki/BOSC_2014
http://usegalaxy.org
https://wiki.galaxyproject.org/Events/GCC2014

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

- Wide range of projects
- Collaboration
- Respected
- Help others
- Grow and learn

# Summary

- Community developed variant calling analyses
- Validation enables science
- Science at scale: 50 to 1500 genomes
- Supporting a community of users
- Software development and science

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