# Scalable and Validated Variant Calling Work in the Bioinformatics Core

Brad Chapman

Bioinformatics Core, Harvard Chan School

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

http://bcb.io

http://j.mp/bcbiolinks

6 February 2015

Oliver Hofmann     Shannan Ho Sui     John Hutchinson     Lorena Pantano

Meeta Mistry     John Morrissey     Rory Kirchner     Brad Chapman

Radhika Khetani     Mary Piper     Andreas Sjödin     Peter Kraft

- Project design
- Analysis and consulting
- Teaching and training
- Infrastructure
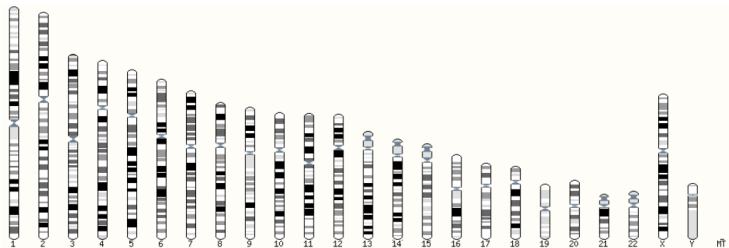
FXB 202B

http://hsphbio.ghost.io
http://bioinformatics.hms.harvard.edu

- What is bcbio?
- Validation
- Support
- Scaling

# Human whole genome sequencing



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

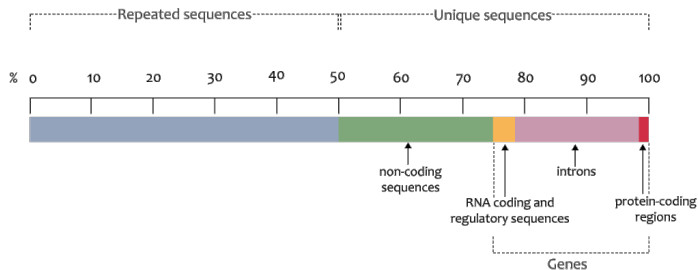# High throughput sequencing

# 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

# White box software

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

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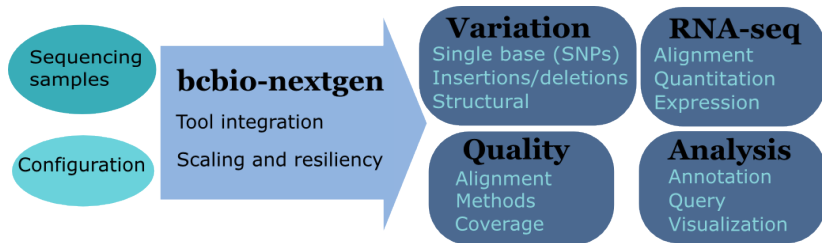
- Aligners: bwa-mem, novoalign, bowtie2
- Variantion: FreeBayes, GATK, Platypus, MuTecT, scalpel, SnpEff, VEP, GEMINI, Lumpy, Delly
- RNA-seq: Tophat, STAR, cufflinks, HTSeq
- Quality control: fastqc, bamtools, RNA-SeQC
- Manipulation: bedtools, bcftools, biobambam, sambamba, samblaster, samtools, vcflib

- Community – collected set of expertise
- Validation
- Scaling
- Multi-architecture parallel processing

# Complex, rapidly changing pipelines



Whole genome, deep coverage v1

Warning: the material on this page is considered out of date by the GSA team.

Best Practice Variant Detection with the GATK v2

Warning: the material on this page is considered out of date by the GSA team.

RETIRED: Best Practice Variant Detection with the GATK v3

Best Practice Variant Detection with the GATK v4, for release 2.0 [RETIRED]

Mark_DePristo Posts: 153
July 2012   edited February 4

The Best Practices have been updated for GATK version 3. If you are running an older version, you should seriously consider upgrading. For more details
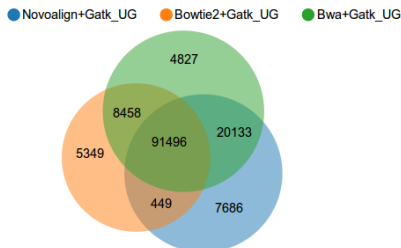
http://www.bioplanet.com/gcat

# Benefits of improved filtering



FreeBayes somatic calling: filter improvements

http://j.mp/cancervalpre

# Solution

# Community: contribution



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

# Community: documentation



https://bcbio-nextgen.readthedocs.org

# Community

## Contributors

- Miika Ahdesmaki, AstraZeneca
- Luca Beltrame, IRCCS "Mario Negri" Institute for Pharmacological Research, Milan, Italy
- Alla Bushoy, AstraZeneca
- Guillermo Carrasco, Science for Life Laboratory, Stockholm
- Nick Carriero, Simons Foundation
- Brad Chapman, Harvard Chan Bioinformatics Core
- Saket Choudhary, University Of Southern California
- Peter Cock, The James Hutton Institute
- Matt Edwards, MIT
- Mario Giovacchini, Science for Life Laboratory, Stockholm
- Karl Gutwin, Biogen
- Jeff Hammerbacher, Icahn School of Medicine at Mount Sinai
- John Kern
- Rory Kirchner, Harvard Chan Bioinformatics Core
- Jakub Nowacki, AstraZeneca
- John Morrissey, Harvard Chan Bioinformatics Core
- Lorena Pantano, Harvard Chan Bioinformatics Core
- Brent Pedersen, University of Colorado Denver
- James Porter, The University of Chicago
- Valentine Svensson, Science for Life Laboratory, Stockholm
- Paul Tang, UCSF
- Roman Valls, Science for Life Laboratory, Stockholm
- Kevin Ying, Garvan Institute of Medical Research, Sydney, Australia

http://gemini.readthedocs.org

Tests for implementation and methods

- Family/population calling
- Structural variations
- Cancer tumor/normal

http://www.genomeinabottle.org/

- Single sample calling
- Pooled calling
- Joint calling
- Squaring off/backfilling

http://j.mp/bcbiojoint

# Squared off VCF

- Parallelize: call samples individually
- Add single new sample to analysis
- Combine existing populations
- Inform calls based on previously known variants

- GATK HaplotypeCaller – gVCFs
- FreeBayes – recalling
- Platypus – recalling
- samtools 1.x – recalling

https://github.com/chapmanb/bcbio.variation.recall

# Multiple approaches work well



Incremental joint calling: GATK HaplotypeCaller, FreeBayes, Platypus and samtools

# Joint vs batch vs single



single, pooled and joint: GATK HaplotypeCaller

- Goal: identify regions with potential issues
- Rough boundaries for additional analysis
- Ensemble: union of all calls
- Understand sensitivity and precision

http://j.mp/bcbiosv

# Structural variant callers

- LUMPY https://github.com/arq5x/lumpy-sv

- Delly https://github.com/tobiasrausch/delly

- cn.mops http://www.bioconductor.org/packages/release/bioc/html/cn.mops.html

- CNVkit http://cnvkit.readthedocs.org/

- WHAM https://github.com/jewmanchue/wham

# Structural variant evaluation

# Making bcbio easy to use



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?

# Need a consistent support environment

http://docker.com

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

- Ready to run
- Easy interface to start/stop clusters
- Pull/push data from encrypted S3
- Lustre and encrypted NFS filesystems

http://bcb.io/2014/12/19/awsbench/

- Odyssey at FAS
  https://rc.fas.harvard.edu/

- Orchestra at HMS
  https://rc.hms.harvard.edu/

- 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
```

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 |

James Cuff, John Morrissey (FAS)
Kristina Kermanshahche (Intel)

# Scaling: avoid intermediates

```
("{bwa} mem -M -t {num_cores} -R '{rg_info}' -v 1 "
 "  {ref_file} {fastq_file} {pair_file} "
 "| {samblaster} "
 "| {samtools} sort -@ {cores} -m {mem} -T {tmp_file}"
 "   -o {tx_out_file} /dev/stdin")
```
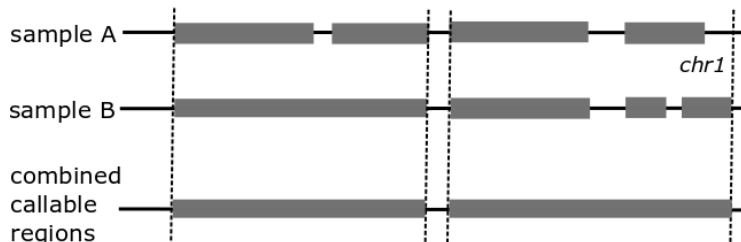
Selection of genome regions for parallel processing

# Scaling: AWS benchmarking
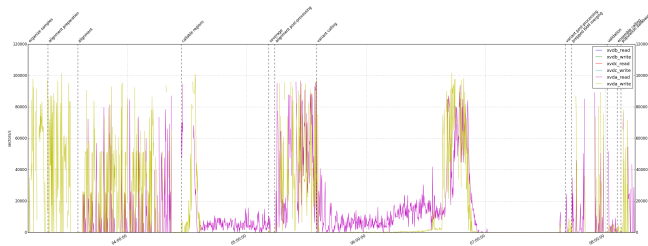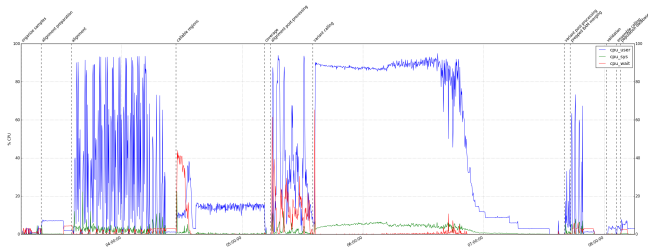
|                          | AWS (Lustre) |
|--------------------------|--------------|
| Total                    | 4:42         |
| genome data preparation  | 0:04         |
| alignment preparation    | 0:12         |
| alignment                | 0:29         |
| callable regions         | 0:44         |
| alignment post-processing| 0:13         |
| variant calling          | 2:35         |
| variant post-processing  | 0:05         |
| prepped BAM merging      | 0:03         |
| validation               | 0:05         |

100X cancer tumor/normal exome on 64 cores (2 c3.8xlarge)

# Scaling: Resource usage plots

# Summary

- bcbio – community built variant calling and RNA-seq analyses
- Validation – measure quality = good science
- Support – AWS and local HPC
- Scaling – diverse teams

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