# Validated variant calling: structural, joint, somatic

Brad Chapman

Bioinformatics Core, Harvard School of Public Health

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

http://j.mp/bcbiolinks

15 October 2014

# Variant calling



Aligned Reads

Reference

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

# Ok, VCFs

```
##FORMAT=<version 1.1 reads-3
#CHROM  POS     ID      REF     ALT     QUAL    FILTER  INFO    FORMAT  Test1
chrM    150     .       T       C       228     .       AC=2;AN=2;DP=250;DP4=0,0,430,504;FS=0.000;GC=44.55;HRun=1;Hapl
otypeScore=0.0000;MQ=60.00;MQ0=0;MQ0F=0;MQSB=1;QD=0.91;SGB=-0.693147;VDB=0.000445395      GT:AD:DP:PL     1/1:0,250:93
4:255,255,0
chrM    152     rs117135796     T       C       228     .       AC=2;AN=2;BQB=1;DB;DP=250;DP4=1,0,422,493;FS=0.000;GC=
42.57;HRun=1;HaplotypeScore=0.7340;MQ=60.00;MQ0=0;MQ0F=0;MQB=1;MQSB=1;QD=0.91;RPB=1;SGB=-0.693147;VDB=6.39968e-05 GT:A
D:DP:PL   1/1:0,250:916:255,255,0
chrM    195     .       C       T       228     .       AC=2;AN=2;DP=250;DP4=0,0,340,405;FS=0.000;GC=29.70;HRun=1;Hapl
otypeScore=0.9665;MQ=60.00;MQ0=0;MQ0F=0;MQSB=1;QD=0.91;SGB=-0.693147;VDB=0.0355506        GT:AD:DP:PL     1/1:0,250:82
5:255,255,0
chrM    410     .       A       T       228     .       AC=2;AN=2;DP=249;DP4=0,0,434,171;FS=0.000;GC=38.61;HRun=3;Hapl
otypeScore=1.9573;MQ=60.00;MQ0=0;MQ0F=0;MQSB=1;QD=0.92;SGB=-0.693147;VDB=0.0003518        GT:AD:DP:PL     1/1:0,249:60
5:255,255,0
chrM    2261    .       C       T       228     .       AC=2;AN=2;DP=250;DP4=0,0,461,357;FS=0.000;GC=37.62;HRun=0;Hapl
otypeScore=0.8667;MQ=60.00;MQ0=0;MQ0F=0;MQSB=1;QD=0.91;SGB=-0.693147;VDB=0.141436 GT:AD:DP:PL     1/1:0,250:818:255,25
5,0
chrM    2354    .       C       T       228     .       AC=2;AN=2;DP=250;DP4=0,0,503,367;FS=0.000;GC=37.62;HRun=1;Hapl
otypeScore=0.9995;MQ=60.00;MQ0=0;MQ0F=0;MQSB=1;QD=0.91;SGB=-0.693147;VDB=0.357219 GT:AD:DP:PL     1/1:0,250:870:255,25
5,0
chrM    2485    .       C       T       228     .       AC=2;AN=2;DP=250;DP4=0,0,306,224;FS=0.000;GC=41.58;HRun=0;Hapl
otypeScore=4.0391;MQ=60.00;MQ0=0;MQ0F=0;MQSB=1;QD=0.91;SGB=-0.693147;VDB=0.186858 GT:AD:DP:PL     1/1:0,250:530:255,2
55,0
chrM    2708    .       G       A       228     .       AC=2;AN=2;DP=218;DP4=0,0,137,23;FS=0.000;GC=43.56;HRun=1;Haplo
typeScore=0.9989;MQ=59.95;MQ0=0;MQ0F=0;MQSB=1;QD=1.05;SGB=-0.693147;VDB=0.00107048        GT:AD:DP:PL     1/1:0,218:16
0:255,255,0
chrM    4746    .       A       G       228     .       AC=2;AN=2;BQB=0.941685;BaseQRankSum=0.619;DP=250;DP4=1,1,514,4
12;FS=0.000;GC=31.68;HRun=0;HaplotypeScore=1.6385;MQ=60.00;MQ0=0;MQ0F=0;MQB=1;MQRankSum=-1.482;MQSB=1;QD=0.91;RPB=0.81
0035;ReadPosRankSum=-1.537;SGB=-0.693147;VDB=0.225191     GT:AD:DP:PL     1/1:1,249:928:255,255,0
chr22   14257   .       CTG     C       5.61    .       AC=2;AN=2;DP=3;DP4=0,0,0,1;GC=53.47;HRun=0;IDV=2;IMF=0.666667;
INDEL;MQ=60.00;MQ0F=0;OLD_VARIANT=chr22:14257:CTGTGTGTGTGTGTG/CTGTGTGTGTGTG;QD=1.87;SGB=-0.379885 GT:DP:PL     1/1:
1:32,3,0
chr22   14259   .       G       C       10.20   .       AC=1;AN=2;DP=3;DP4=0,0,0,1;FS=0.000;GC=53.47;HOB=0.5;HRun=0;Ha
plotypeScore=10.6316;ICB=1;MQ=60.00;MQ0=0;MQ0F=0;QD=3.40;SGB=-0.379885    GT:AD:DP:PL     0/1:0,1:1:38,3,0
chr22   14424   .       C       T       131     .       AC=2;AN=2;DP=5;DP4=0,0,3,2;FS=0.000;GC=52.48;HRun=0;HaplotypeS
core=0.0000;MQ=60.00;MQ0=0;MQ0F=0;MQSB=1;QD=26.20;SGB=-0.590765;VDB=0.125998      GT:AD:DP:PL     1/1:0,5:5:159,15,0
chr22   15494   .       C       T       6.95    .       AC=1;AN=2;BQB=0;BaseQRankSum=-1.231;DP=5;DP4=0,2,0,3;FS=0.000;
```
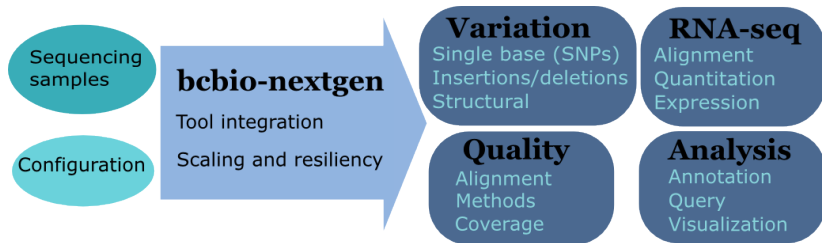
# Summary

- What is bcbio?
- Community software development
- Variation validation
- Software support

# Overview

# 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

# Quality differences between methods



**Variant Calling Test**

Discuss

We compare combinations of variant calling pipelines across different data sets. Browse our public facing reports to see how various aligner + variant caller combinations perform against each other. Test your own combination of tools by creating your own report. Below is a sample conconcordance view on our "Illumina 100bp Paired End 30x Coverage" data set.

**Variant Concordance - "illumina-100bp-pe-exome-30x"**

● Novoalign+Gatk_UG    ● Bowtie2+Gatk_UG    ● Bwa+Gatk_UG

4827

8458

5349     91496     20133

449     7686

http://www.bioplanet.com/gcat

# Solution

# Community: contribution



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

# Community: documentation

Introduction
Installation
**Getting started**
Overview
Logging
Example pipelines
Test suite
Configuration
Parallel execution
Pipelines
Outputs
Code
Internals
Presentations
Deprecated functionality

Docs » Getting started                                    ⌂ Edit on GitHub

## Getting started

### Overview

1. Create a sample configuration file for your project (substitute the example BAM and fastq names below with the full path to your sample files):

```
bcbio_nextgen.py -w template gatk-variant project1 sample1.bam sample2_1.fq sample2_2.fq
```

This uses a standard template (GATK best practice variant calling) to automate creation of a full configuration for all samples. See *Automated sample configuration* for more details on running the script, and manually edit the base template or final output file to incorporate project specific configuration. The example pipelines provide a good starting point and the *Sample information* documentation has full details on available options.

2. Run analysis, distributed across 8 local cores:

```
bcbio_nextgen.py bcbio_sample.yaml -n 8
```

https://bcbio-nextgen.readthedocs.org

Tests for implementation and methods

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

http://www.genomeinabottle.org/

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

http://j.mp/bcbiojoint

# Squared off VCF

# Implementation

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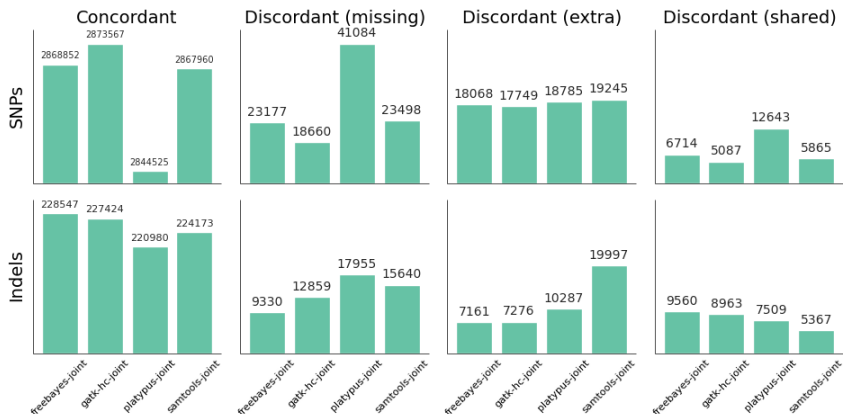- Platypus – recalling
- samtools 1.x – recalling

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

- Parallelize: call samples individually
- Add single new sample to analysis
- Combine existing populations

# Multiple approaches



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

- Truth calls: synthetic data from DREAM challenge
- Mixed population of subclones
- Need additional complexity: mixed cellularity

http://j.mp/dreamsyn3

# Community built



- Luca Beltrame – IRCCS, Italy
- Miika Ahdesmaki – AstraZeneca
- Mario Giovacchini – SciLifeLab, Sweden
- Lorena Pantano – HSPH

# Callers available

- MuTect
  https://www.broadinstitute.org/cancer/cga/mutect
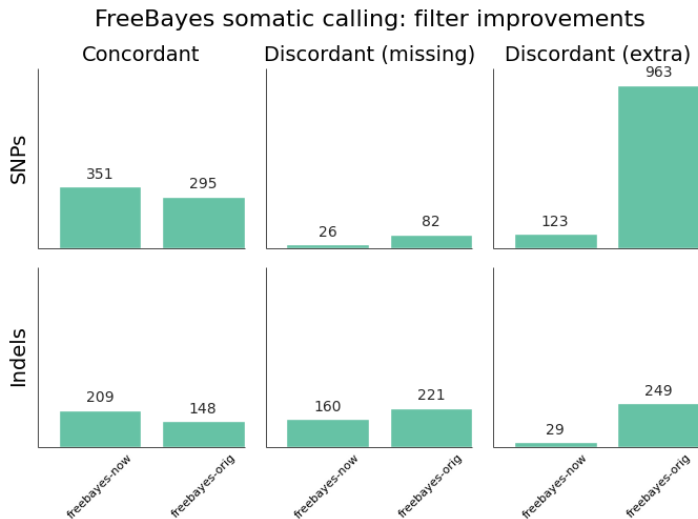
- FreeBayes https://github.com/ekg/freebayes

- VarScan http://varscan.sourceforge.net/

- VarDict https://github.com/AstraZeneca-NGS/VarDict

- Ensemble

# Somatic evaluation



DREAM syn3 exomes: Ensemble, MuTect/Scalpel, FreeBayes, VarDict, VarScan

# Benefits of improved filtering



FreeBayes somatic calling: filter improvements

# Validation enables scaling

- 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

http://j.mp/bcbioeval2

# 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

# Docker lightweight containers

- 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

# Sustainability

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

# Summary

- What is bcbio?
- Community software development
- Variation validation
- Software support

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