

Interoperable community developed variant calling with bcbio and the Common Workflow Language

Brad Chapman
Bioinformatics Core, Harvard Chan School


<https://bcb.io>
<http://j.mp/bcbiolinks>

3 November 2016

- Open source communities
- bcbio: community developed analyses
- Value of variant validation
- Interoperable infrastructure

Supporting bioinformatics

LATEST OPEN RNA-SEQ CHIP-SEQ SNP ASSEMBLY TUTORIALS TOOLS JOBS FORUM PLANET ALL »

 **Biostars**
— BIOINFORMATICS EXPLAINED —

Brad Chapman • 9.0k | Logout about • faq • rss

Community Messages Votes (1) My Posts My Tags Following Bookmarks New Post

Live search: start typing... or Classic search

Limit to: all time < prev • 41,829 results • page 1 of 1046 • next > Sort by: update

0 votes **1** answer **130** views **Mapping drug targets to KEGG pathways**
pathway drug target kegg
written 8 weeks ago by ayanava18 • 10 • updated 6 minutes ago by Ron • 240

0 votes **0** answers **3** views **Normalized read density**
Insert size atac-seq
written 7 minutes ago by Kramdi • 0

This ad appears once a day.
Support the site by visiting the sponsor.

Recent Votes

- A: Mapping /Need software
- C: Extract some features from a combined gene bank file in bash
- C: RNA-Seq Best way to handle multimapping

<https://www.biostars.org/>

O|B|F

- Main Page
- Projects
- News
- BOSC
- OBF Board
- Join

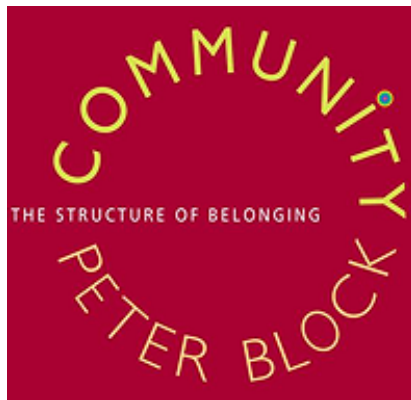
Page **Discussion**

BOSC 2016



http://www.open-bio.org/wiki/BOSC_2016

Build communities



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

Challenge: many communities

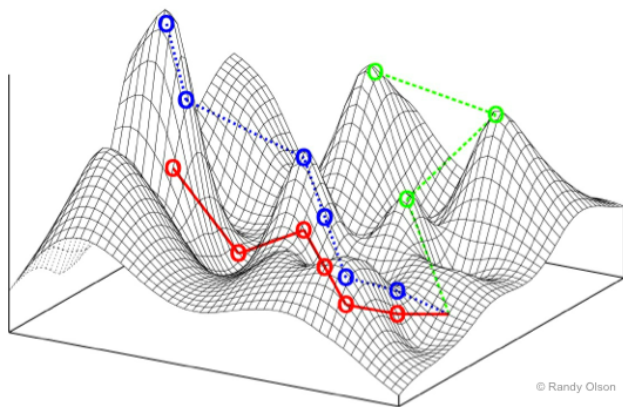


<https://galaxyproject.org/>

<http://www.cbioportal.org/>

<https://www.synapse.org/>

Challenge: open source communities not yet optimal



https://en.wikipedia.org/wiki/Fitness_landscape

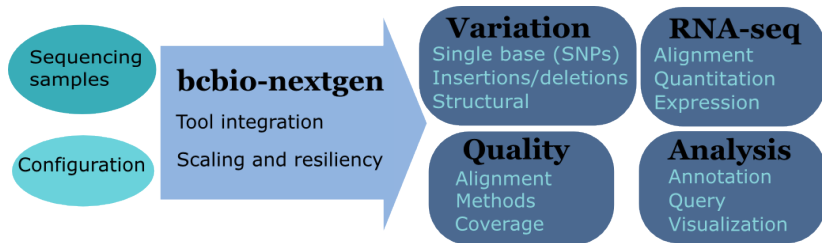
Large scale infrastructure development

- Shared problems – academic, industry, startups
- Community developed analyses
- Validation
- Scaling
- Supporting a community of users

White box software



Overview



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

High level configuration

```
- analysis: variant2
  genome_build: hg38
  algorithm:
    aligner: bwa
    mark_duplicates: true
    recalibrate: false
    realign: false
    variantcaller: [gatk-haplotype, freebayes, vardict]
    ensemble:
      numpass: 2
    svcaller: [lumpy, manta]
```

[https://bcbio-nextgen.readthedocs.org/en/latest/contents/
configuration.html](https://bcbio-nextgen.readthedocs.org/en/latest/contents/configuration.html)

- Aligners: bwa, novoalign, bowtie2, HISAT2
- Variantion: FreeBayes, GATK, VarDict, MuTect2, Scalpel, SnpEff, VEP, GEMINI, Lumpy, Manta, CNVkit, WHAM
- RNA-seq: Tophat, STAR, Cufflinks, Sailfish
- Quality control: FastQC, samtools, Qualimap, MultiQC
- Manipulation: bedtools, bcftools, biobambam, picard, sambamba, samblaster, samtools, vcflib, vt

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

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>

Complex, rapidly changing baseline functionality

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

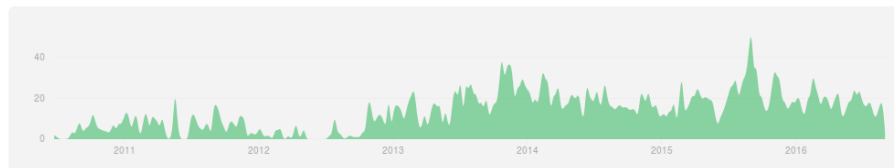
The [Best Practices](#) have been updated for GATK version 3. If you are running an older version, you should seriously consider upgrading. For more details

Community: sustainability

Jul 18, 2010 – Sep 27, 2016

Contributions: **Commits** ▾

Contributions to master, excluding merge commits



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

Community: support

<input type="checkbox"/>	95 Open ✓ 1,215 Closed	Author ▾	Labels ▾	Milestones ▾	Assignee ▾	Sort ▾
<input type="checkbox"/>	update yaml templates #1575 opened 3 minutes ago by saboswell					
<input type="checkbox"/>	HG38 and Gemini #1573 opened a day ago by matthdsm					7
<input type="checkbox"/>	Test run error #1572 opened 4 days ago by firatuyulur					2
<input type="checkbox"/>	vep annotation fields + hgvs #1571 opened 4 days ago by matthdsm					7
<input type="checkbox"/>	how to force bam to stream directly to bwa? #1567 opened 6 days ago by brentp					2
<input type="checkbox"/>	Would it be possible to run the QC stage in parallel? #1556 opened 14 days ago by NeillGibson					14
<input type="checkbox"/>	consider samtools depth to replace sambamba bedtools in callable #1549 opened 18 days ago by brentp					21

<https://bcbio-nextgen.readthedocs.org>

Community: contribution

The screenshot shows the GitHub repository page for `chapmanb/bcbio-nextgen`. At the top, there are buttons for 'Unwatch' (74), 'Unstar' (342), and 'Fork' (172). Below this is a navigation bar with links to 'Code', 'Issues' (95), 'Pull requests' (4), 'Projects' (0), 'Pulse', 'Graphs', and 'Settings'. A description of the repository follows: 'Validated, scalable, community developed variant calling, RNA-seq and small RNA analysis <https://bcbio-nextgen.readthedocs.org> — Edit'. Below the description is a summary bar showing '5,060 commits', '2 branches', '35 releases', '43 contributors', and the license 'MIT'. Further down, there are buttons for 'Branch: master', 'New pull request', 'Create new file', 'Upload files', 'Find file', and 'Clone or download'. The main content area shows a list of recent commits, with the most recent being 'Docs: how to change logging directory location' by 'chapmanb' 3 hours ago. Below this is a table of files in the repository.

chapmanb / **bcbio-nextgen** Unwatch 74 Unstar 342 Fork 172

<> Code Issues 95 Pull requests 4 Projects 0 Pulse Graphs Settings

Validated, scalable, community developed variant calling, RNA-seq and small RNA analysis <https://bcbio-nextgen.readthedocs.org> — Edit

5,060 commits 2 branches 35 releases 43 contributors MIT

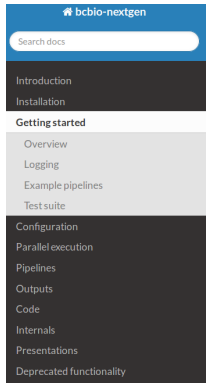
Branch: master New pull request Create new file Upload files Find file Clone or download

chapmanb Docs: how to change logging directory location Latest commit `791a5e8` 3 hours ago

artwork	add logo to README and docs	a year ago
bcbio	VEP: add support for HGVS with clinical_reporting	8 hours ago
config	Update resources to match new bgzipped BED files	9 days ago
docs	Docs: how to change logging directory location	3 hours ago

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

Community: documentation



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>

Supported analysis types

⊞ Pipelines

☐ Germline variant calling

Basic germline calling

Population calling

Cancer variant calling

Structural variant calling

RNA-seq

single-cell RNA-seq

smallRNA-seq

ChIP-seq

<https://bcbio-nextgen.readthedocs.org/en/latest/contents/pipelines.html>

- Integration tests for pipelines
- Unbiased algorithm comparisons
- Baseline for improving methods



Genome in a Bottle
Consortium



Global Alliance
for Genomics & Health

ICGC-TCGA DREAM Mutation Calling challenge

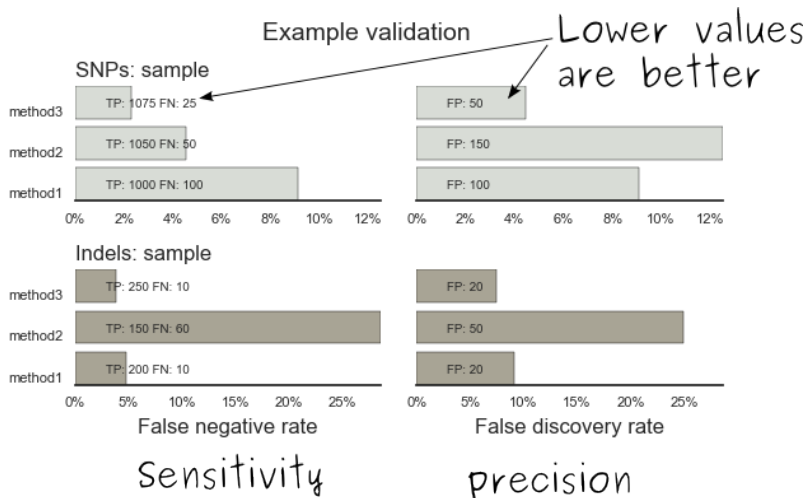
<http://www.genomeinabottle.org/>

<http://ga4gh.org/#/benchmarking-team>

<https://www.synapse.org/#!Synapse:syn312572>

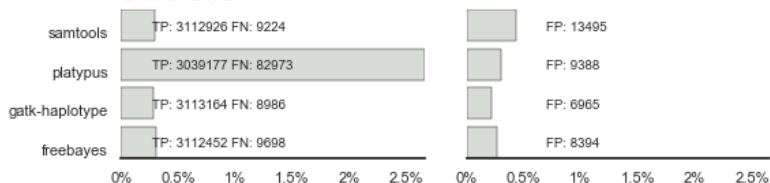
- Collaboration with GATK methods development
- Compare HaplotypeCaller to other methods
- Germline validation
- Genome in a Bottle reference materials
 - NA12878 – Caucasian
 - NA24385 – Ashkenazim Jewish
 - NA24631 – Chinese

Validation graphs

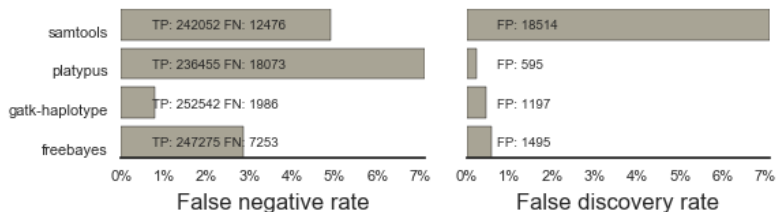


NA12878: Genome in a Bottle whole genome validation

SNPs: bwa

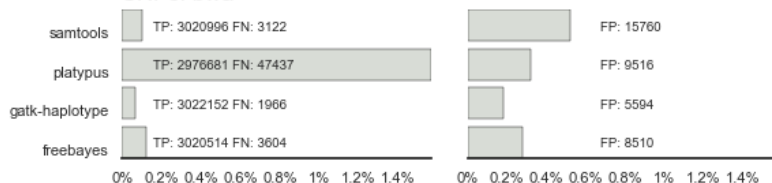


Indels: bwa

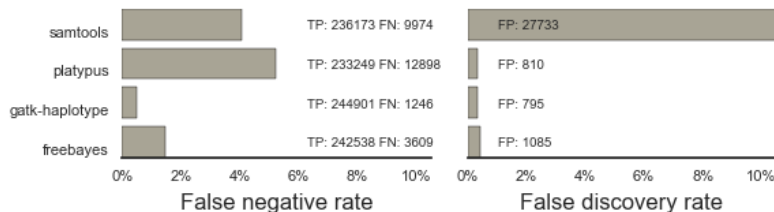


NA24385: Genome in a Bottle whole genome validation

SNPs: bwa



Indels: bwa



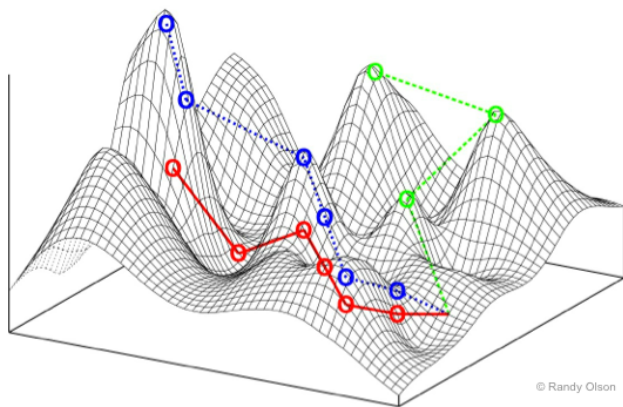
Validation results

- Good performance for GATK HaplotypeCaller
- Other good performing callers: FreeBayes
- Consistency across diverse samples
- Identify potential problem areas for tuning
 - samtools Indel false positive rates
 - Platypus SNP sensitivity
- PrecisionFDA: <https://precision.fda.gov/>

Infrastructure Goals

- Free, open source, community developed
- Welcoming to contributions
- Local machines
- Clusters: SLURM, SGE, Torque, PBS, LSF
- Clouds: Amazon, Google, Azure
- Clinical environments
- User interface for researchers
- Integrate with LIMS
- Accessible to the general public

Challenge: open source communities not yet optimal



https://en.wikipedia.org/wiki/Fitness_landscape

Better abstractions = more interoperability



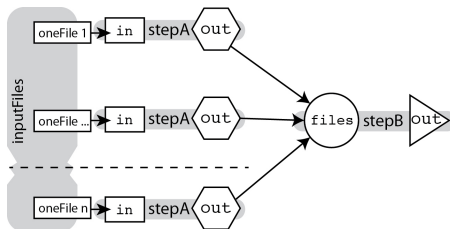
COMMON
WORKFLOW
LANGUAGE



<https://bcbio-nextgen.readthedocs.io/en/latest/contents/cwl.html>




Workflow Description Language (WDL)

```
workflow myWorkflowName {  
  call task_A  
  call task_B  
}  
task task_A {  
  ...  
}  
task task_B {  
  ...  
}
```



<https://software.broadinstitute.org/wdl/>

Common Workflow Language (CWL)

Workflow	pipeline-se-narrow.cwl		
Sub-workflow 1	01-qc-se.cwl		
Step 1	extract.cwl	extract.py	
Step 2	count.cwl	count.py	
Step 3	fastqc.cwl	fastqc	
Sub-workflow 2	02-trim.cwl		
...			

<http://www.commonwl.org/>

<https://f1000research.com/slides/5-1617>

Abstraction > Implementation

$WDL \leftrightarrow CWL$

- Start with high level configuration file
- Generate CWL
- Run CWL:
 - Any infrastructure that supports CWL
 - Generated CWL
 - Docker or local bcbio installation
 - Genome data

<https://bcbio-nextgen.readthedocs.io/en/latest/contents/cwl.html>

Why use a workflow abstraction?

- Integrate with multiple platforms
 - Arvados
 - Toil
 - Cromwell
 - Galaxy
 - Nextflow
 - Seven Bridges
 - DNAnexus
- Stop maintaining bcbio specific infrastructure
- Focus on hard biological problems

Welcome to the Arvados Project

The Arvados community is dedicated to building a new generation of open source distributed computing software for bioinformatics, data science, and production analysis using massive data sets.



<https://arvados.org/>

<https://cloud.curoverse.com/>

- UCSC NIH Big Data to Knowledge Center for Translational Genomics
- Supports CWL via conversion to internal Python workflow description
- Local HPC support: SLURM, SGE
- Cloud: AWS + spot instances

<http://toil.readthedocs.io/en/latest/>

Cromwell via conversion to WDL

- bcbio workflow abstractions supported in WDL
 - Tasks, workflows, nested workflows
 - Scatter based parallelization
 - Grouping/batching of samples
- Work in progress CWL to WDL converter based on cwl2wdl
- Happy to collaborate

<https://github.com/broadinstitute/cromwell> <https://github.com/chapmanb/bcbio-nextgen/blob/master/scripts/utils/cwltool2wdl.py>

- CWL support in progress
- Supports subset of CWL – tool definitions
- Needs workflow support

<https://github.com/galaxyproject/planemo>

<https://f1000research.com/posters/5-2567>



Michael R. Crusoe

@biocrusoe



 Follow

I received permission to share:
[@CH_maria_CH](#) is working on a [#CommonWL](#)
to [@nextflowio](#) s2s converter! (No promises on
a release schedule :-)

RETWEETS

6

LIKES

4



4:24 PM - 21 Oct 2016

<https://www.nextflow.io/>

- Currently supports CWL v2 + extensions
- Moving to CWL v1.0
- External runner, Bunny, supports v1.0

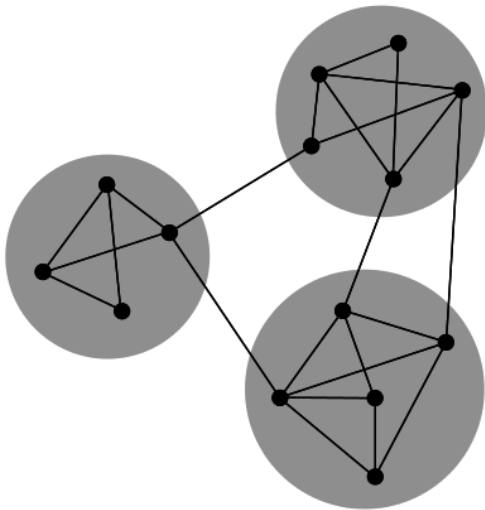
<https://github.com/rabix/bunny>

<http://docs.sevenbridges.com/docs/sdk-overview>

- Integration of bcbio Docker containers
- Single node runs
- Convert CWL to DNAnexus API for distributed

<https://www.dnanexus.com/developer-resources>

Connections



By jham3 - Own work, CC BY-SA 3.0,

<https://commons.wikimedia.org/w/index.php?curid=17125894>

Summary

- bcbio community developed resources
- Value of validation
 - Germline calling with Genome in a Bottle
- Interoperable infrastructure
 - Importance of abstractions
 - Common Workflow Language
 - Practical interoperability

<http://bcb.io>