# Validated, scalable, community developed variant calling and RNA-seq analysis

Brad Chapman
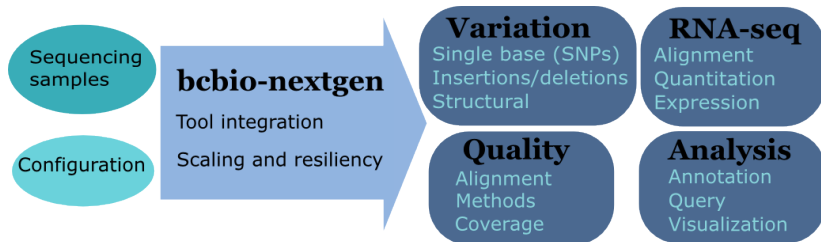
Bioinformatics Core, Harvard School of Public Health

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

http://j.mp/bcbiolinks

3 June 2014

# Development goals

- Community developed

- Quantifiable

- Configurable

- Scalable

- Reproducible

- Pipeline

- Best-practice

- Framework

# Uses

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

- Validation – outputs + automated evaluation
- Tool integration
- Multi-platform support
- Scaling

# Complex, rapidly changing tools

# 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
```
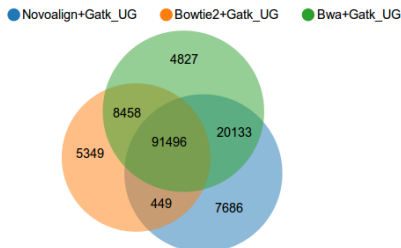
https://github.com/StanfordBioinformatics/HugeSeq

http://www.bioplanet.com/gcat

# Solution



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

**John Davey**
@johnomics

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

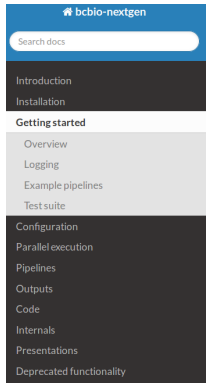↩ Reply   ♺ Retweet   ★ Favorite   ••• More

## Automated Install
Bare machine to ready-to-run with tools and data

- CloudBioLinux: http://cloudbiolinux.org
- Homebrew: https://github.com/Homebrew/homebrew-science
- Conda: http://j.mp/py-conda

## Easier install
Docker

# Community: documentation



https://bcbio-nextgen.readthedocs.org

# Community: contribution



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

Tests for implementation and methods

- Currently:
  - Family/population calling
  - RNA-seq differential expression
  - Structural variations
- Expand to:
  - Cancer tumor/normal
    http://j.mp/cancer-var-chal

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

# Reference materials



http://www.genomeinabottle.org/

# Quantify quality



Minimal BAM preparation (samtools de-duplication only)

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

# 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

- High level abstraction
- Adjust by intent, rather than command line
- Domain specific language
- YAML configuration file

# Getting started

- Start with examples
  https://bcbio-nextgen.readthedocs.org/en/latest/contents/testing.html#example-pipelines
- Automatically generate configuration
  https://bcbio-nextgen.readthedocs.org/en/latest/contents/configuration.html#automated-sample-configuration
- Parameter documentation
  https://bcbio-nextgen.readthedocs.org/en/latest/contents/configuration.html#algorithm-parameters

```
 5   details:
 6     - analysis: RNA-seq
 7       algorithm:
 8         aligner: star
 9         quality_format: Standard
10         trim_reads: read_through
11         adapters: [truseq, polya]
12       description: Test1
13       files: [1_110907_ERP000591_1_fastq.txt, 1_110907_ERP000591_2_fastq.txt]
14       genome_build: mm9
```

# Example – variant calling

```
11  details:
12    - files: [../input/NA12878_1.fastq.gz, ../input/NA12878_2.fastq.gz]
13      description: NA12878
14      metadata:
15        batch: ceph
16        sex: female
17      analysis: variant2
18      genome_build: GRCh37
19      algorithm:
20        aligner: bwa
21        align_split_size: 5000000
22        mark_duplicates: true
23        recalibrate: false
24        realign: false
25        variantcaller: [freebayes, gatk-haplotype]
26        quality_format: Standard
27        coverage_interval: genome
28        remove_lcr: true
29        validate: ../input/GiaB_NIST_RTG_v0_2.vcf.gz
30        validate_regions: ../input/GiaB_NIST_RTG_v0_2_regions.bed
```

# Scaling overview



- Infrastructure details: http://j.mp/bcbioscale
- IPython: http://ipython.org/ipython-doc/dev/parallel/index.html

- Cluster scheduler
  - SLURM
  - Torque
  - SGE
  - LSF
- Shared filesystem
  - NFS
  - Lustre
- Local temporary disk
  - SSD

# Configuration to batch scripts

*Configuration*

```
bwa:
  cmd: bwa
  cores: 16
samtools:
  cores: 16
  memory: 2G
gatk:
  jvm_opts: ["-Xms750m", "-Xmx2750m"]
```

*Batch file*

```
#PBS -l nodes=1:ppn=16
#PBS -l mem=45260mb
```

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

System

- 560 cores
- 4Gb RAM/core
- Lustre filesystem
- Infiniband network

Samples

- 75 samples
- 30x whole genome (100Gb)
- Illumina
- Family-based calling

# Timing: Alignment

| Step | Time | Processes |
| --- | --- | --- |
| Alignment preparation | 9.5 hours | BAM to fastq; bgzip; grabix index |
| Alignment | 31 hours | bwa-mem alignment samblaster deduplication |
| BAM merge | 5.5 hours | Merge alignment parts |
| Post-processing | 11 hours | Calculate callable regions |

# Timing: Variant calling

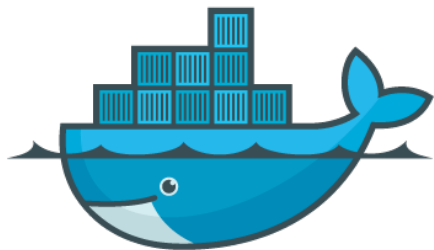| Step | Time | Processes |
|------|------|-----------|
| Variant calling | 30 hours | FreeBayes |
| Variant post-processing | 5 hours | Combine variant files; |
| | | annotate: GATK and snpEff |

# Timing: Analysis and QC

| Step | Time | Processes |
|------|------|-----------|
| GEMINI | 5 hours | Create GEMINI SQLite database |
| Quality Control | 2.5 hours | FastQC, alignment and variant statistics |

- 100 hours, ~4 days for 75 samples
- ~1 1/2 hours per sample at 560 cores
- In progress: optimize for single samples

# Reproducible environment

# 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

# Docker HPC parallelization

- Cost – spot instances
- Disk – local scratch, no EBS
- Organization – no shared filesystems, S3 push/pull
- Data – reconstitute on minimal machines
- Security – encryption at rest

Clusterk http://clusterk.com/

# Summary

- Community development > challenges
- Easy to install, learn and contribute
- Validated
- Configurable
- Scalability
- Reproducibility and virtualization

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