# Supporting dynamic community developed biological pipelines

Brad Chapman

Bioinformatics Core, Harvard School of Public Health

https://github.com/chapmanb

http://j.mp/bcbiolinks

17 April 2014

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Warning: the material on this page is considered out of date by the GSA team.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Warning: the material on this page is considered out of date by the GSA team.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

# 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



http://www.bioplanet.com/gcat

# Solution

Sequencing samples

Configuration

**bcbio-nextgen**

Best-practice pipelines
Tool integration
Scaling and resiliency

**Variation**
Single base (SNPs)
Insertions/deletions
Structural

**RNA-seq**
Alignment
Quantitation
Expression

**Quality**
Alignment
Methods
Coverage

**Analysis**
Annotation
Query
Visualization

# Uses

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

# Provides

- Best practice analysis pipelines
- Tool integration
- Multi-platform support
- Scaling

- Community developed

- Quantifiable

- Scalable

- Reproducible

**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
Bare machine to ready-to-run pipeline, 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

# Example evaluation

- 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

- 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

# Scaling overview



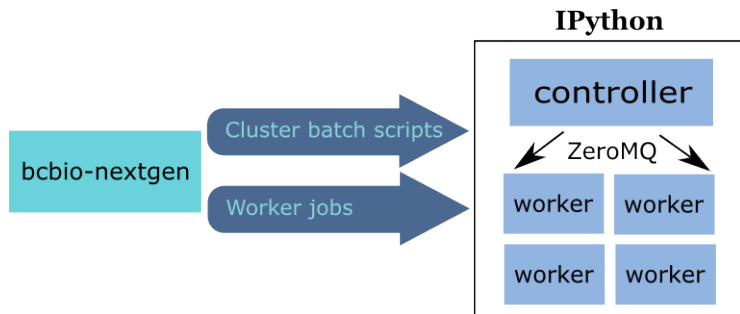- Infrastructure details: http://j.mp/bcbioscale
- IPython: http://ipython.org/ipython-doc/dev/parallel/index.html
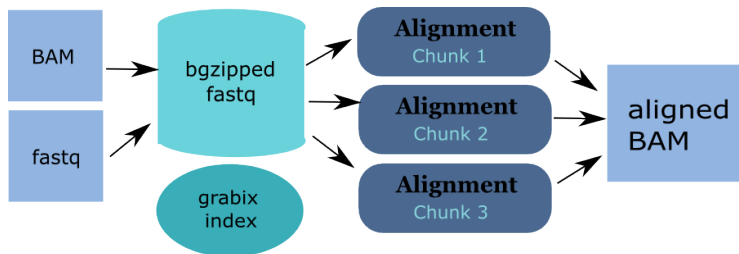
- Cluster scheduler
  - SLURM
  - Torque
  - SGE
  - LSF
- Shared filesystem
  - NFS
  - Lustre
- Local temporary disk
  - SSD

- Split alignments
- Split by genome regions
- Manage memory
- Avoid IO

# Alignment parallelization



https://github.com/arq5x/grabix
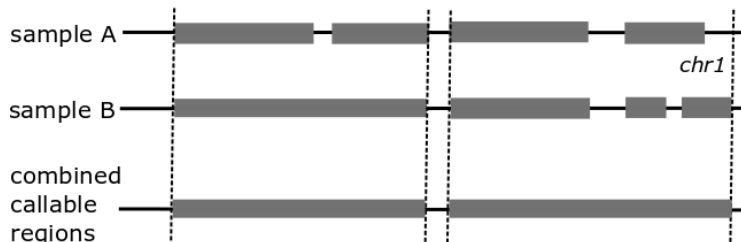
# Variant calling parallelization



Selection of genome regions for parallel processing

# Memory usage

*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
```

# Avoid filesystem IO

Pipes and streaming algorithms

```
("{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")
```

# Dell System



## Dell Active Infrastructure for HPC Life Sciences

**High Performance Computing**

› Dell Advantage
› Strategy
› Products & Services
› Resource Library

"With diseases like neuroblastoma, hours matter. Our new Dell HPC cluster allows us to do the processing we need to get a meaningful result in a clinically relevant amount of time."
— Jason Corneveaux, Bioinformatician, Neurogenomics Division, the Translational Genomics Research Institute [1]

**High performance for high-volume genomics research**

Processing complex genomic data sets requires massive compute power, storage and network capabilities. Getting the balance right is critical to success, but without proper support and expertise, it can take months to integrate the necessary computing components and tune them for maximum performance and efficiency.

Glen Otero, Will Cottay
http://dell.com/ai-hpc-lifesciences

System

- 400 cores
- 3Gb RAM/core
- Lustre filesystem
- Infiniband network

Samples

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

# Timing: Alignment

| Step | Time | Processes |
|---|---|---|
| Alignment preparation | 13 hours | BAM to fastq; bgzip; grabix index |
| Alignment | 30 hours | bwa-mem alignment |
| BAM merge | 7 hours | Merge alignment parts |
| Alignment post-processing | 6 hours | Calculate callable regions |

# Timing: Variant calling

| Step | Time | Processes |
| --- | --- | --- |
| Post-alignment BAM preparation | 6 hours | De-duplication |
| Variant calling | 18 hours | FreeBayes |
| Variant post-processing | 2 hours | Combine variant files; annotate: GATK and snpEff |

# Timing: Analysis and QC

| Step | Time | Processes |
| --- | --- | --- |
| BAM merging | 6 hours | Combine post-processed BAM file sections |
| GEMINI | 3 hours | Create GEMINI SQLite database |
| Quality Control | 5 hours | FastQC, alignment and variant statistics |

- 4 days for 60 samples
- ~2 hours per sample at 400 cores
- In progress: optimize for single samples

http://docker.io

# 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/

Arvados is a free and open source bioinformatics platform for genomic and biomedical data.
Store | Organize | Compute | Share

https://arvados.org/
https://curoverse.com/

# Integrated



https://usegalaxy.org/

# Summary

- Community developed pipelines > challenges
- Focus
    - Community: easy to install and contribute
    - Validation of methods
    - Scalability
    - Reproducibility and virtualization
- Widely accessible

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