# Connecting Galaxy to tools with alternative storage and compute models

Brad Chapman

Bioinformatics Core, Harvard School of Public Health

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

http://j.mp/bcbiolinks

1 July 2014

# Community > Implementation

Galaxy
Biopython: http://biopython.org

OpenBio: http://www.open-bio.org

# Validation > Replication

Genome in a Bottle: http://www.genomeinabottle.org/
ICGC-TCGA DREAM: https://www.synapse.org/#!Synapse:syn312572

SMaSH: http://smash.cs.berkeley.edu/

Scaling > Configurability

bcbio scaling: http://j.mp/bcbioscale

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

# Uses

- Aligners: bwa-mem, novoalign, bowtie2
- Variantion: FreeBayes, GATK, MuTecT, 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
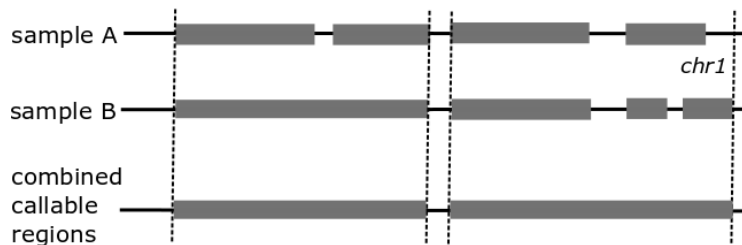
- Community – collected set of expertise
- Tool integration
- Validation – outputs + automated evaluation
- Installation of tools and data
- Scaling

```
("{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")
```

Selection of genome regions for parallel processing

- Intermediates – 6x final

```
$ du -sh *
353G final
2.2T work
```

- 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
```

# bcbio as a Galaxy tool

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

# Summary

- Focus: Community, Validation, Scaling
- bcbio-nextgen

  `https://github.com/chapmanb/bcbio-nextgen`
- Challenges: parallelization, scaling and storage
- Galaxy integration: Simple tool with Docker installation