# Quantifiable and scalable detection of genomic variants
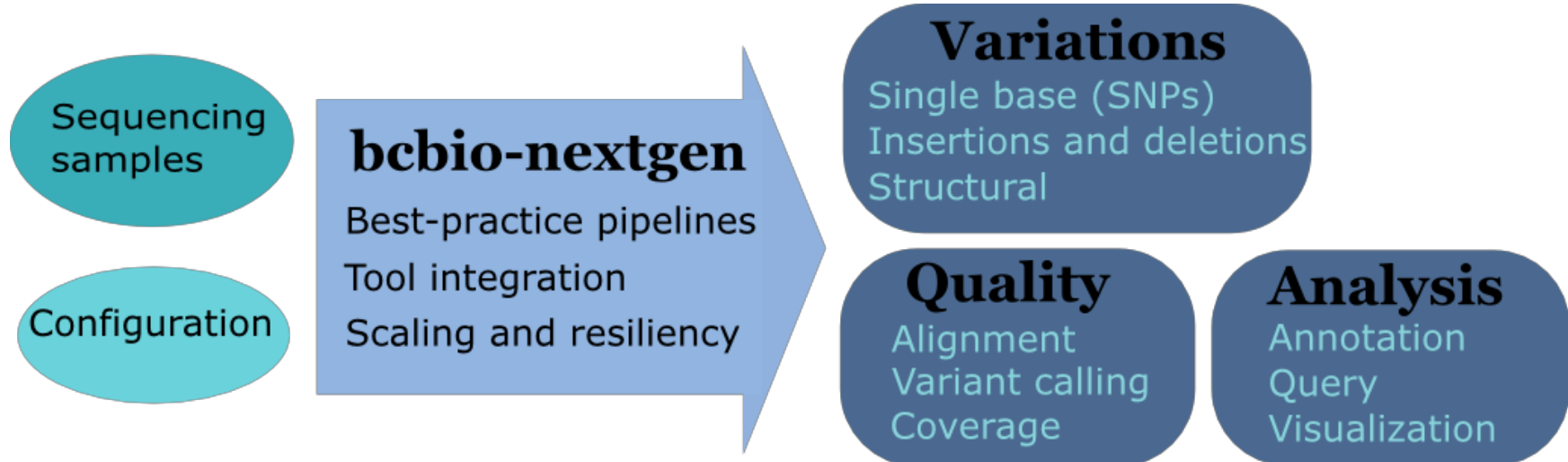
Brad Chapman

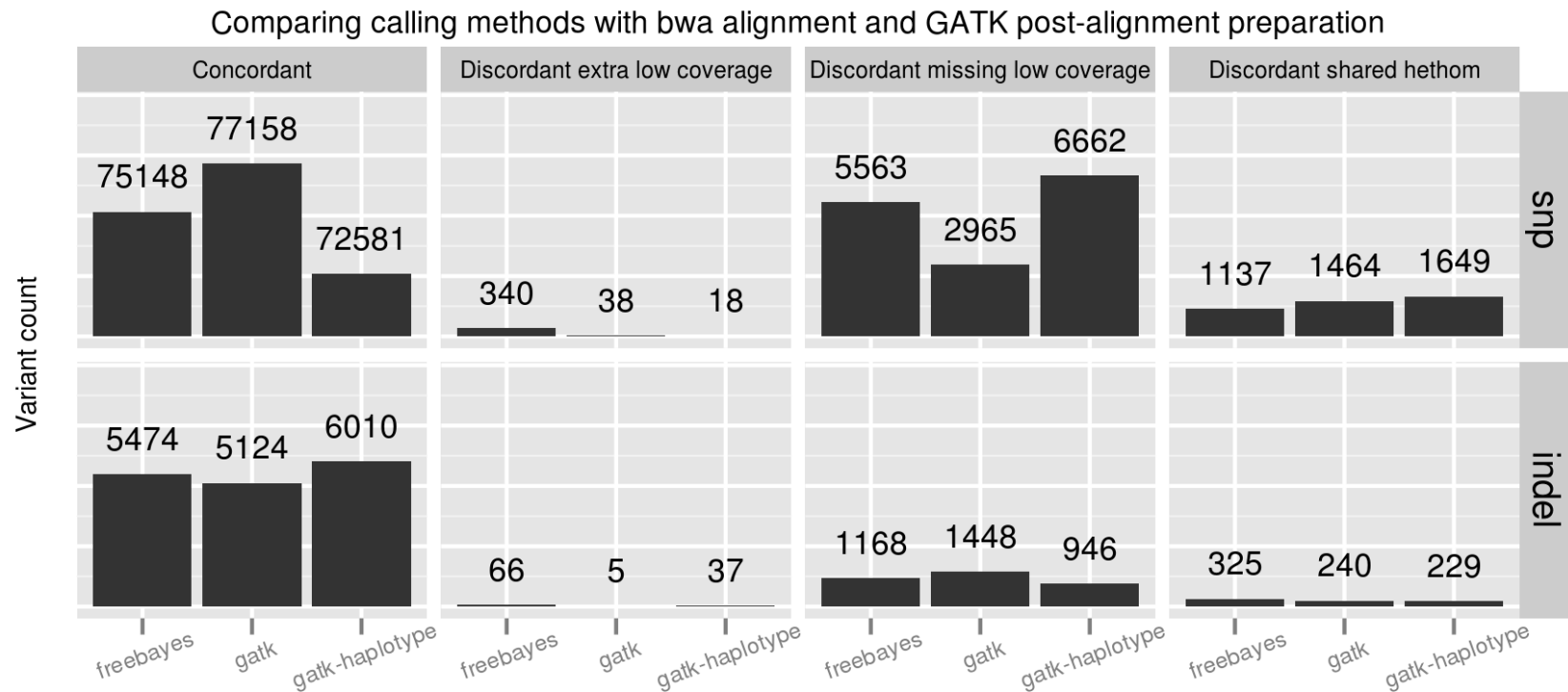Bioinformatics Core, Harvard School of Public Health

**@chapmanb**

26 June 2013

# Development goals

- Quantifiable: assess variant quality

- Scalable: 1500 whole genome samples

- Reproducible: text-configurable, provenance, version tracking

- Community developed: open source, documented and widely deployable
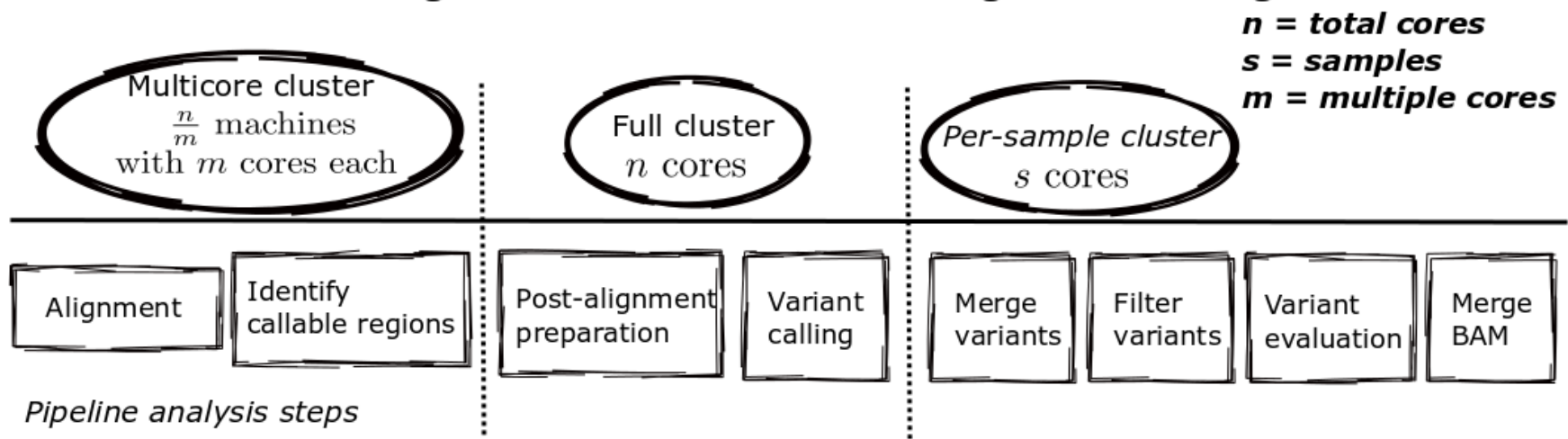
# Quantify quality



Comparing calling methods with bwa alignment and GATK post-alignment preparation

Reference materials:
**http://www.genomeinabottle.org/**

# Parallel scaling



Heterogeneous cluster creation during variant calling

$n$ = total cores
$s$ = samples
$m$ = multiple cores

Multicore cluster
$\frac{n}{m}$ machines
with $m$ cores each

Full cluster
$n$ cores

Per-sample cluster
$s$ cores

Pipeline analysis steps

| Alignment | Identify callable regions | Post-alignment preparation | Variant calling | Merge variants | Filter variants | Variant evaluation | Merge BAM |

Infrastructure:
**http://ipython.org/ipython-doc/dev/parallel/index.html**

# Reproducible configuration

```
- files: [NA12878-NGv3-LAB1360-A_1.fastq.gz, NA12878-NGv3-LAB1
360-A_2.fastq.gz]
  description: NA12878
  analysis: variant2
  genome_build: GRCh37
  algorithm:
    aligner: bwa
    recalibrate: gatk
    realign: gatk
    variantcaller: [gatk, freebayes, gatk-haplotype]
    coverage_interval: exome
    coverage_depth: high
    platform: illumina
    quality_format: Standard
    validate: NA12878-nist-v2_13-NGv3-pass.vcf
```

# Community developed

- Fully automated installation: CloudBioLinux
- Deployable on multiple clusters (LSF, SGE, Torque)
- Integrated with web platforms (Galaxy, STORMSeq)
- Open source and documented

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