# Variant calling with validated community developed tools

Brad Chapman
Bioinformatics Core, Harvard Chan School
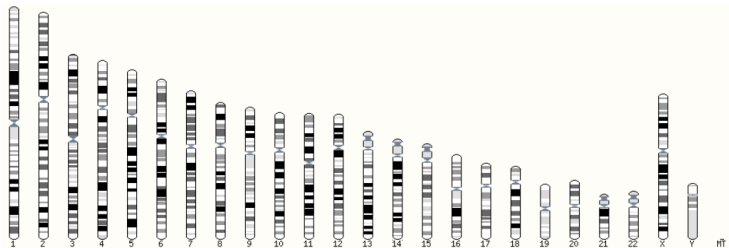https://bcb.io
http://j.mp/bcbiolinks

6 July 2017

- Overview of variant calling tools
- Open source community resources
- bcbio validated variant analysis
- Science
  - Human build 38
  - GATK4 validation
  - Cancer calling of low frequency variants
  - Structural variation
- Practical calling example

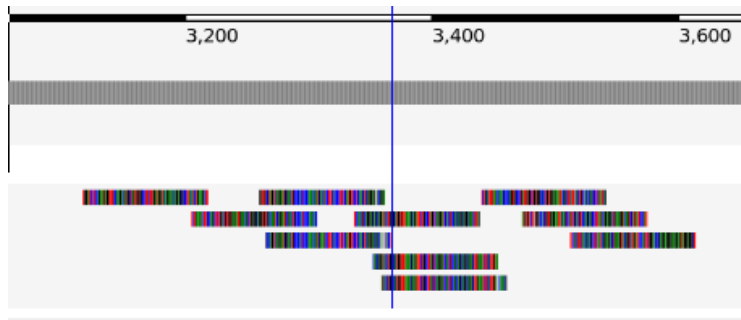# Human whole genome sequencing



http://ensembl.org/Homo_sapiens/Location/Genome

# High throughput sequencing

# Variant calling



Aligned Reads

Reference

http://en.wikipedia.org/wiki/SNV_calling_from_NGS_data

# Scale: exome to whole genome



The haploid human genome sequence
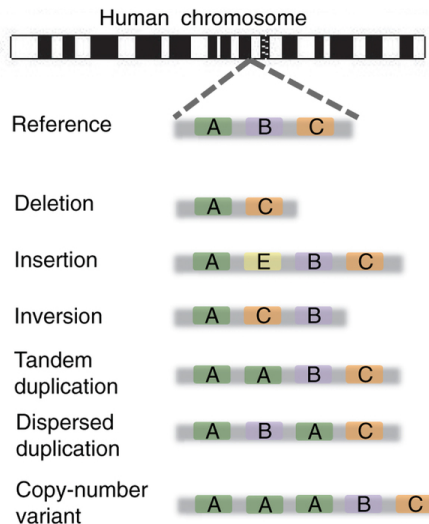
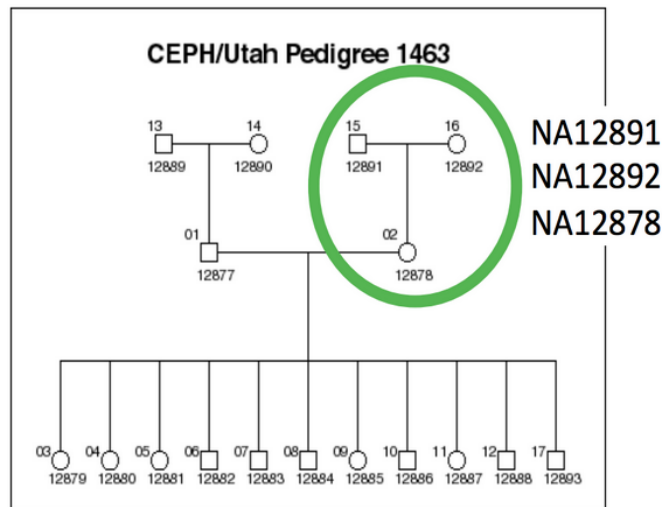https://www.flickr.com/photos/119980645@N06/

# SNPs and Indels

# Structural variations
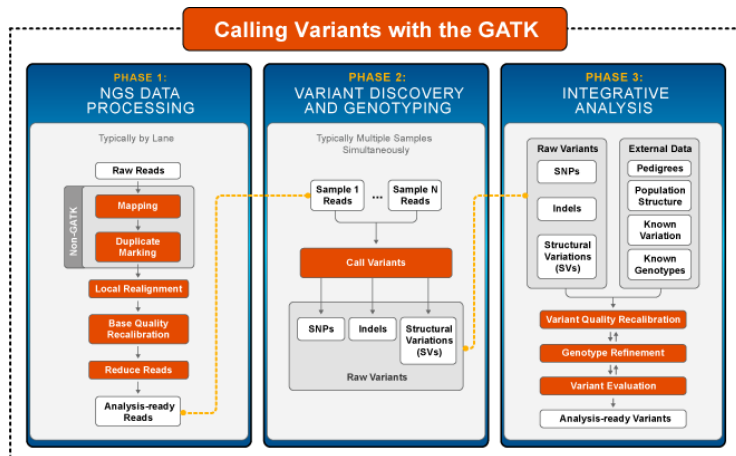
# Germline population calling

# Genome Analysis Toolkit (GATK)

The Genome Analysis Toolkit or GATK is a software package developed at the Broad Institute to analyze high-throughput sequencing data. The toolkit offers a wide variety of tools, with a primary focus on variant discovery and genotyping as well as strong emphasis on data quality assurance. Its robust architecture, powerful processing engine and high-performance computing features make it capable of taking on projects of any size.

https://www.broadinstitute.org/gatk/

# GATK Best Practices



https://www.broadinstitute.org/gatk/guide/best-practices

# HaplotypeCaller

# Joint calling on large populations

# GATK4 now open source for all uses



https://software.broadinstitute.org/gatk/blog?id=9645

# FreeBayes



https://github.com/ekg/freebayes

# Filtering – Variant Quality Score Recalibration



$$\text{VQSLOD(x)} = \text{Log}(\textcolor{blue}{\text{p(x)}}/\textcolor{red}{\text{q(x)}})$$

```
filters = ('((AC[0] / AN) <= 0.5 && DP < 4 && %QUAL < 20) || '
           '(DP < 13 && %QUAL < 10) || '
           '((AC[0] / AN) > 0.5 && DP < 4 && %QUAL < 50)')
```

http://bcb.io/2014/05/12/wgs-trio-variant-evaluation/

# Effects prediction

- snpEff

http://snpeff.sourceforge.net/

- Variant Effect Predictor (VEP) from Ensembl

http://www.ensembl.org/info/docs/tools/vep/index.html

https://github.com/arq5x/gemini

# VCF – overview

## Types of variants

### SNPs

| Alignment | VCF representation | | |
|-----------|-----|-----|-----|
| ACGT | POS | REF | ALT |
| ATGT | 2 | C | T |

### Insertions

| Alignment | VCF representation | | |
|-----------|-----|-----|-----|
| AC-GT | POS | REF | ALT |
| ACTGT | 2 | C | CT |

### Deletions

| Alignment | VCF representation | | |
|-----------|-----|-----|-----|
| ACGT | POS | REF | ALT |
| A--T | 1 | ACG | A |

### Complex events

| Alignment | VCF representation | | |
|-----------|-----|-----|-----|
| ACGT | POS | REF | ALT |
| A-TT | 1 | ACG | AT |

### Large structural variants

| VCF representation | | | |
|-----|-----|-----|-----|
| POS | REF | ALT | INFO |
| 100 | T | <DEL> | SVTYPE=DEL;END=300 |

http://vcftools.sourceforge.net/VCF-poster.pdf

- Step by step guide from Broad

https://www.broadinstitute.org/gatk/guide/article?id=1268

- Specification

http://samtools.github.io/hts-specs/

# You want to build a variant calling pipeline



**Best Practices for Germline SNPs and Indels in Whole Genomes and Exomes - June 2016**

https://software.broadinstitute.org/gatk/best-practices/

- Changing tools
- Feature support burden
- Validation

# Build open source communities





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

http://www.open-bio.org/wiki/BOSC_2017

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

# Supported analysis types

# 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

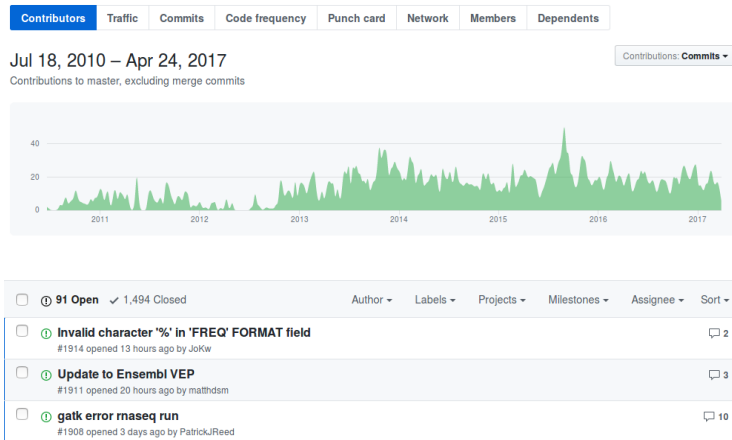# Complex, rapidly changing baseline functionality

# Feature support burden

Table 1: Comparison of Nextflow with other workflow management systems

| Workflow | Nextflow | Galaxy | Toil | Snakemake | Bpipe |
|---|---|---|---|---|---|
| Platform[b] | GroovyJVM | Python | Python | Python | GroovyJVM |
| Native task support[c] | Yes (any) | No | No | Yes (BASH only) | Yes (BASH only) |
| Common workflow language[d] | No | Yes | Yes | No | No |
| Streaming processing[e] | Yes | No | No | No | No |
| Dynamic branch evaluation | Yes | ? | Yes | Yes | Undocumented |
| Code sharing integration[f] | Yes | No | No | No | No |
| Workflow modules[g] | No | Yes | Yes | No | Yes |
| Workflow versioning[h] | Yes | Yes | No | No | No |
| Automatic error failover[i] | Yes | No | Yes | No | No |
| Graphical user interface[j] | No | Yes | No | No | No |
| DAG rendering[k] | Yes | Yes | Yes | Yes | Yes |
| **Container management** | | | | | |
| Docker support[l] | Yes | Yes | Yes | No | No |
| Singularity support[m] | Yes | No | No | No | No |
| Multi-scale containers[n] | Yes | No | No | No | No |
| **Built-in batch schedulers[o]** | | | | | |
| Univa Grid Engine | Yes | Yes | Yes | Partial | Yes |
| PBS/Torque | Yes | Yes | No | Partial | Yes |
| LSF | Yes | Yes | Yes | Partial | Yes |
| SLURM | Yes | Yes | Yes | Partial | No |
| HTCondor | Yes | Yes | No | Partial | No |
| **Built-in distributed cluster[p]** | | | | | |
| Apache Ignite | Yes | No | No | No | No |
| Apache Spark | No | No | Yes | No | No |
| Kubernetes | Yes | No | No | No | No |
| Apache Mesos | No | No | Yes | No | No |
| **Built-in cloud[q]** | | | | | |
| AWS (Amazon Web Services) | Yes | Yes | Yes | No | No |

# Community: sustainability and support



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

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

http://www.genomeinabottle.org/
http://ga4gh.org/#/benchmarking-team
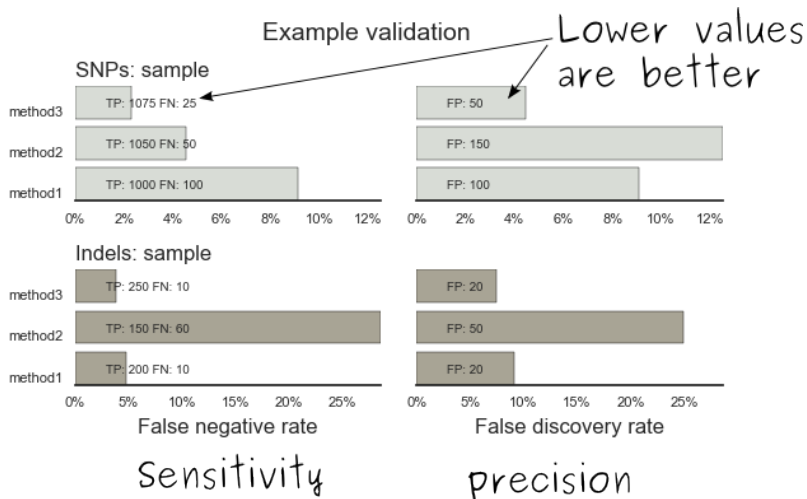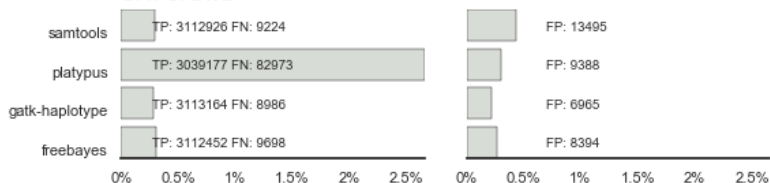https://www.synapse.org/#!Synapse:syn312572

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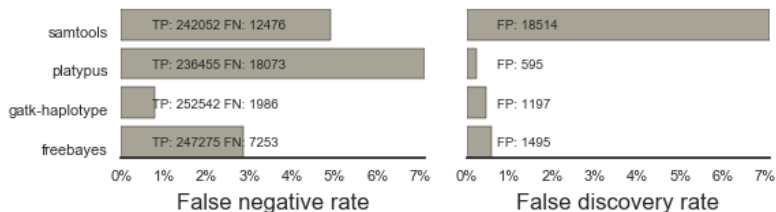  - NA24631 – Chinese

# Validation graphs

NA12878: Genome in a Bottle whole genome validation

# NA24385



NA24385: Genome in a Bottle whole genome validation

SNPs: bwa

| | |
|---|---|
| samtools | TP: 3020996 FN: 3122 | FP: 15760 |
| platypus | TP: 2976681 FN: 47437 | FP: 9516 |
| gatk-haplotype | TP: 3022152 FN: 1966 | FP: 5594 |
| freebayes | TP: 3020514 FN: 3604 | FP: 8510 |

0%  0.2% 0.4% 0.6% 0.8%  1%  1.2% 1.4%        0%  0.2% 0.4% 0.6% 0.8%  1%  1.2% 1.4%

Indels: bwa

| | |
|---|---|
| samtools | TP: 236173 FN: 9974 | FP: 27733 |
| platypus | TP: 233249 FN: 12898 | FP: 810 |
| gatk-haplotype | TP: 244901 FN: 1246 | FP: 795 |
| freebayes | TP: 242538 FN: 3609 | FP: 1085 |

0%   2%   4%   6%   8%   10%        0%   2%   4%   6%   8%   10%

False negative rate                    False discovery rate

# Validation results

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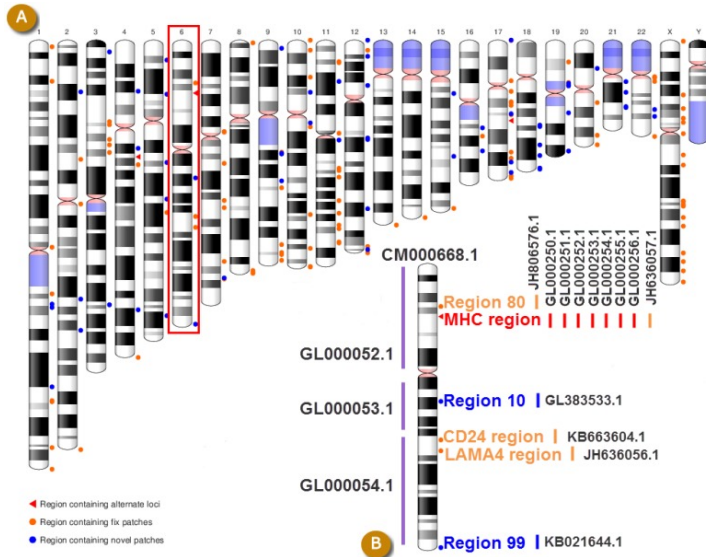  - samtools Indel false positive rates
  - Platypus SNP sensitivity
- PrecisionFDA: https://precision.fda.gov/

- **Human build 38**
- GATK4 validation
- Low frequency somatic calling
- Structural variation

# GRCh37/hg19

# GRCh38 – graph based, many more alternative loci

Reference assembly influence

Sample

Ref
Assembly

http://www.slideshare.net/GenomeRef/transitioning-to-grch38

# Avoiding collapsed repeats

- Build 37 and 38
- Validation sets: Genome in a Bottle, Illumina Platinum Genomes
- 38 builds: with/without alternative alleles
- Variant callers: FreeBayes, GATK HaplotypeCaller

http://bcb.io/2015/09/17/hg38-validation/

hg19/hg38 comparison: NA12878 Platinum Genomes

- SNPs: build 38 more sensitive
- SNPs: build 38 reduces false positives
- Indels: build 38 detected more
- Indels: work on sensitivity and precision

- Human build 38
- **GATK4 validation**
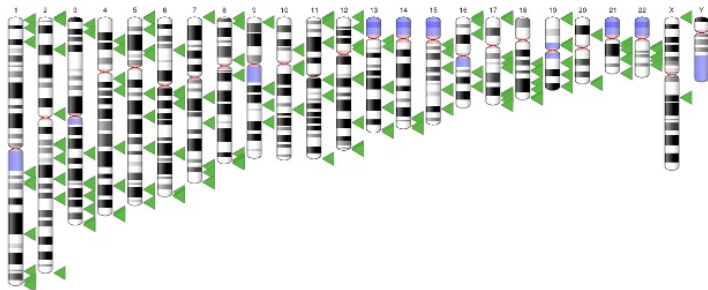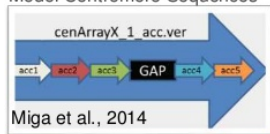- Low frequency somatic calling
- Structural variation

NA12878 hg38: GATK4

SNPs: NA12878

| | |
|---|---|
| vardict | TP: 2976078 FN: 13079 |
| haplotyper | TP: 2981646 FN: 7511 |
| gatk4-haplotype-sor | TP: 2975527 FN: 13630 |
| gatk4-haplotype | TP: 2983914 FN: 5243 |
| gatk3-haplotype | TP: 2992792 FN: 8138 |
| freebayes | TP: 2994258 FN: 6672 |

| | |
|---|---|
| FP: 14004 |
| FP: 3964 |
| FP: 2661 |
| FP: 4780 |
| FP: 4340 |
| FP: 6813 |

Indels: NA12878

| | |
|---|---|
| vardict | TP: 214028 FN: 90254 |
| haplotyper | TP: 277097 FN: 27185 |
| gatk4-haplotype-sor | TP: 257647 FN: 46635 |
| gatk4-haplotype | TP: 257653 FN: 46629 |
| gatk3-haplotype | TP: 284086 FN: 38588 |
| freebayes | TP: 284349 FN: 38325 |

| | |
|---|---|
| FP: 18890 |
| FP: 10660 |
| FP: 2103 |
| FP: 2096 |
| FP: 13415 |
| FP: 10847 |

False negative rate

False discovery rate

- Comparable sensitivity and specificity to GATK3
- Removed a recommended filter
  - Strand Odds Ratio (SOR) – strand bias
  - Improves sensitivity
  - ~6000 TPs vs ~2000 FPs
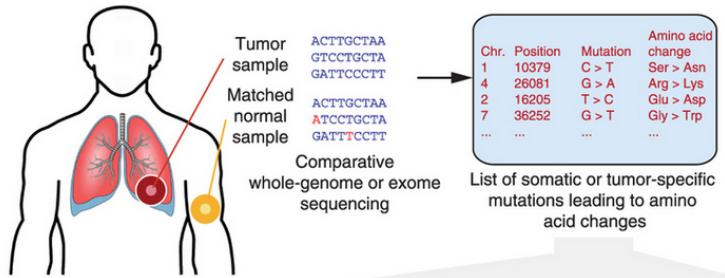- Indels in GATK need additional tuning
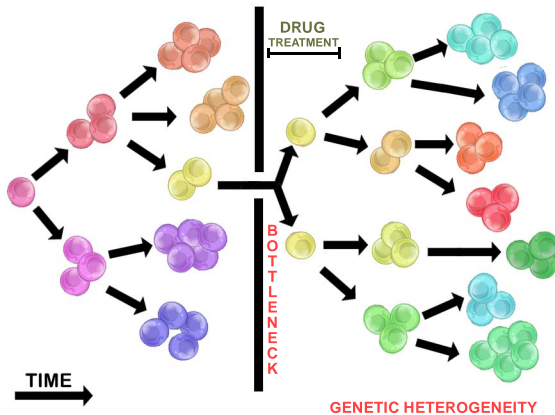  - Sensitivity/specificity tradeoff
  - ~26k TPs vs ~11k FPs

- Human build 38
- GATK4 validation
- **Low frequency somatic calling**
- Structural variation

# Cancer somatic calling

# Cancer heterogeneity



http://en.wikipedia.org/wiki/Tumour_heterogeneity

# VarDict

- AstraZeneca
- Germline + Cancer calling
- SNP + Insertion/Deletions
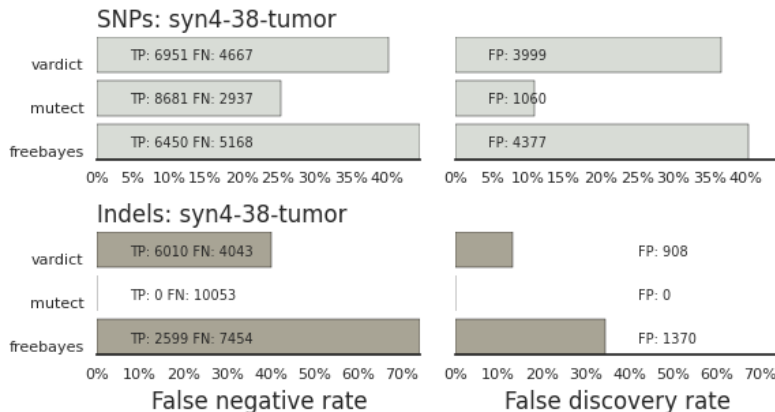- Whole genome + exome
- Also works on deep targeted data

https://github.com/AstraZeneca-NGS/VarDictJava
http://nar.oxfordjournals.org/content/early/2016/04/07/
nar.gkw227.full

# DREAM synthetic dataset 4

| in silico 3 | in silico 4 |
|---|---|
| BWA Backtrack | BWA MEM |
| SNV, SV (deletions, duplications, insertions, inversions) & INDEL | SNV, SV (deletions, duplications, inversions) & INDEL |
| 100% | 80% |
| 50%, 33%, 20% | 50%, 35% (effectively 30% and 15% due to cellularity) |
| Female | Male |
| HCC1143 BL from TCGA Benchmark 4 | CPCG0102R (Provided by ICGC) |

https://www.synapse.org/#!Synapse:syn312572/wiki/62018

# VarDict sensivitity/precision before

# VarDict sensivitity/precision after



SNPs: DREAM synthetic 4 (hg38)

Indels: DREAM synthetic 4 (hg38)

False negative rate

False discovery rate
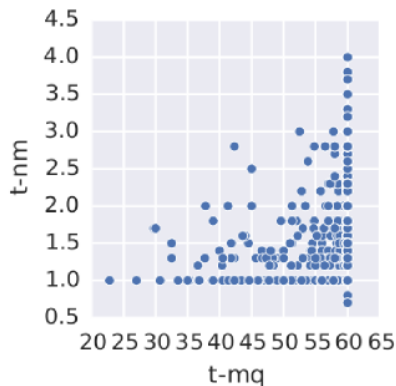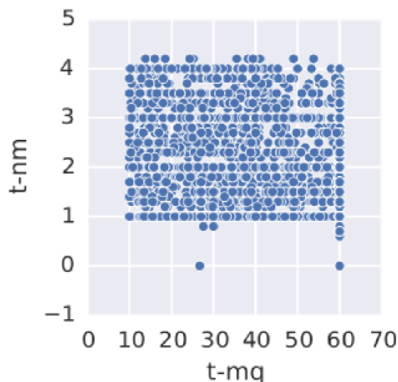
```
((AF * DP < 6) &&
 ((MQ < 55.0 && NM > 1.0) ||
  (MQ < 60.0 && NM > 2.0) ||
  (DP < 10) ||
  (QUAL < 45)))
```

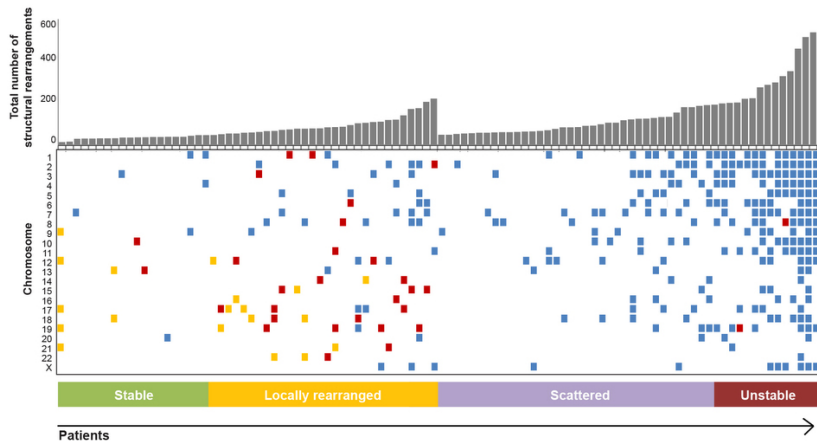# Example filter: mapping quality and number of mismatches

- Human build 38
- GATK4 validation
- Low frequency somatic calling
- **Structural variation**

# Structural variants critical in cancer

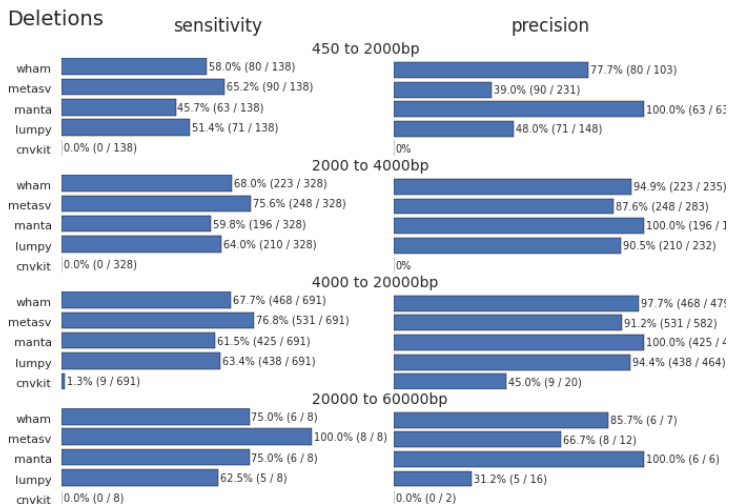- Manta: https://github.com/Illumina/manta
- CNVkit: https://github.com/etal/cnvkit
- Lumpy: https://github.com/arq5x/lumpy-sv
- WHAM: https://github.com/zeeev/wham
- MetaSV: https://github.com/bioinform/metasv

# Results: Somatic deletions

# Public cancer variant databases

- CIViC: https://civic.genome.wustl.edu
- IntOGen: http://www.intogen.org



http://www.amazon.com/The-Biology-Cancer-Robert-Weinberg/dp/0815340761

- Small dataset – single chromosome, exome
- Cancer sample from DREAM synthetic dataset 3
- Call against build 38

https://www.synapse.org/#!Synapse:syn312572/wiki/58893

- Somatic tumor/normal samples
- SNP and indel calling at lower frequency
- Structural variant detection
- Prioritization with CIViC
- HLA typing

# bcbio configuration file

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

https://bcbio-nextgen.readthedocs.org/en/latest/contents/
configuration.html

# bcbio template file – CSV

```
samplename,description,batch,phenotype,sex,variant_regions
sample1,ERR256785,batch1,normal,female,/path/to/regions.bed
sample2,ERR256786,batch1,tumor,,/path/to/regions.bed
```

https://bcbio-nextgen.readthedocs.org/en/latest/contents/configuration.html#automated-sample-configuration

```
bcbio_nextgen.py -w template \
    tumor-paired.yaml project1.csv \
    sample1.bam sample2_1.fq sample2_2.fq
```

https://bcbio-nextgen.readthedocs.org/en/latest/contents/configuration.
html#automated-sample-configuration

`bcbio_nextgen.py project1.yaml -n 8`

https://bcbio-nextgen.readthedocs.org/en/latest/contents/testing.html

https://bcbio-nextgen.readthedocs.org/en/latest/
contents/teaching.html

- Pre-downloaded and analysis run
- AMI (ami-5e84fe34)

# Summary

- Overview of variant calling tools
- Open source community resources
- bcbio validated variant analysis
- Science
  - Human build 38
  - GATK4 validation
  - Cancer calling of low frequency variants
  - Structural variation
- Practical calling example

http://bcb.io