# bcbio validation: build 38, low frequency somatic variants, structural variations

Brad Chapman

Bioinformatics Core, Harvard Chan School
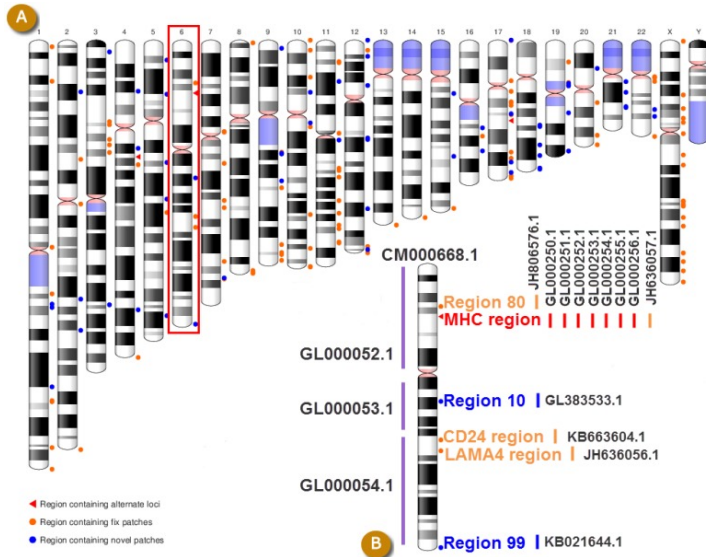
https://bcb.io

http://j.mp/bcbiolinks

5 November 2015

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

Reference assembly influence

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

# Avoiding collapsed repeats

- Build 37 and 38
- Validation sets: Genome in a Bottle, Illumina Platinum Genomes
- Lift-over methods: CrossMap/LiftOver, NCBI Remap
- 38 builds: with/without alternative alleles
- Variant callers: FreeBayes, GATK HaplotypeCaller

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

Genome in a Bottle Consortium

Global Alliance for Genomics & Health

ICGC-TCGA DREAM Mutation Calling challenge

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

hg19/hg38 comparison: NA12878 Platinum Genomes

GRCh37/hg38 comparison: NA12878 Genome in a Bottle

# Small variant results

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

Need conversion approaches for resources not yet available on build 38

- CrossMap:
  http://crossmap.sourceforge.net/

- NCBI remap:
  http://www.ncbi.nlm.nih.gov/genome/tools/remap

- Both performed well

- NCBI remap has additional sensitivity, but needs tuning

# Major histocompatibility complex (MHC) – HLAs



http://www.ebi.ac.uk/ipd/imgt/hla/

http://sciscogenetics.com/technology/human-leukocyte-antigen-complex/

# Alignment: bwa alternative allele support



```
         Read: ATCAGCATC

   ALT ctg 1:      TGAAA---CGAATGCAAATGGTCAATCAGCATCGAACTAGTCACAT
                   ||||| (high div) |||||| (novel ins) ||||||||||||
 Chromosome: GCGTACATGATACGAATCgGCATCATGGTC------------CTAGTCACATCGTAATC
                   |||||||||||||| |||||||||| (novel ins) ||||||||||||
   ALT ctg 2:      TGAATACGAATCgcCATCATGGTCAATCgcCAgCGAACTAGTCACAT

        4 potential hits: ATCAGCATC > ATCgGCATC > ATCgcCATC > ATCgcCAgC
           2 hit groups: {ATCAGCATC,ATCgcCAgC} and {ATCgGCATC,ATCgcCATC}
 Hits considered in mapQ: ATCAGCATC and ATCgGCATC (best from each group)

     In the output SAM: ATCgGCATC as the primary SAM line with mapQ=0
                        ATCAGCATC as a supplementary line with mapQ>0
                        ATCgcCAgC as a supplementary line with mapQ>0
                        ATCgcCATC in an XA tag, not as a separate line
```
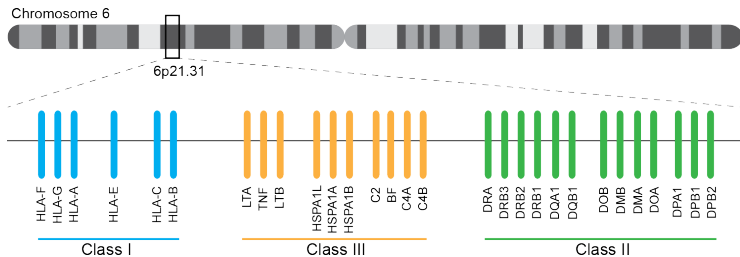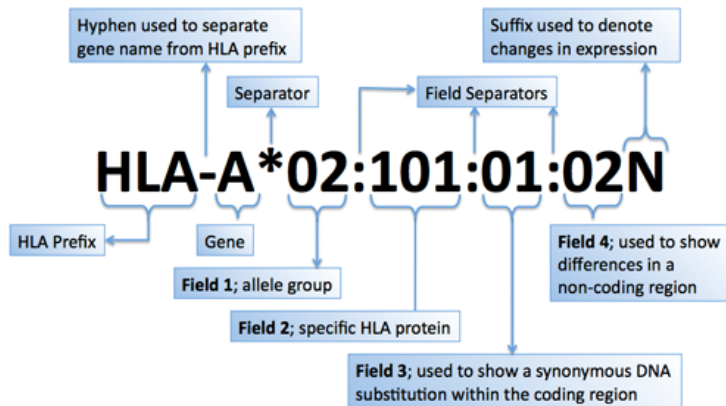
https://github.com/lh3/bwa/blob/master/README-alt.md

# HLA typing

- bwakit implementation
- 1000 genomes: build 38 + IMGT/HLA-3.18.0
- bwa extracts HLA reads
- fermi de novo assembly
- Remap assemblies back to HLA choices
- Call HLA types

https://github.com/lh3/bwa/blob/master/README-alt.md#hla-typing

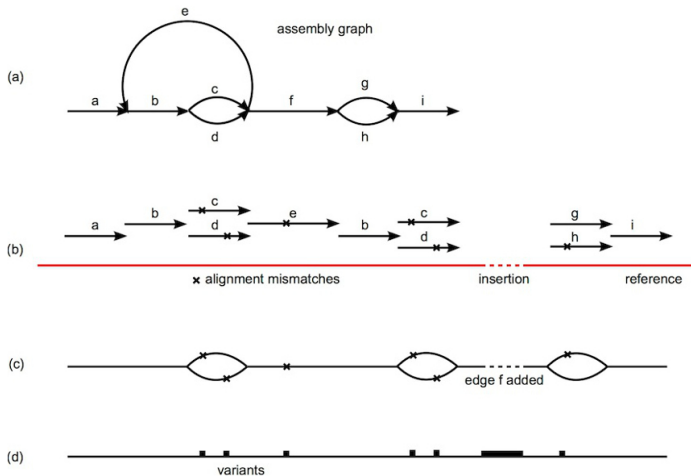# HLA nomenclature



https://www.ebi.ac.uk/ipd/imgt/hla/
http://hla.alleles.org/alleles/p_groups.html

- Omixon example data
- bwakit calls on exome and deep targeted data
- P-group resolution
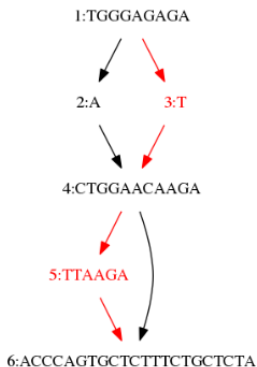- Good results for exome
- Assembly problems with deep targeted

http://www.omixon.com/hla-typing-example-data/

https://gist.github.com/chapmanb/8e2a18c7bbbee3167395

# Genome graphs and variation

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

# Cancer somatic calling

# Cancer heterogeneity



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

# VarDict

- AstraZeneca
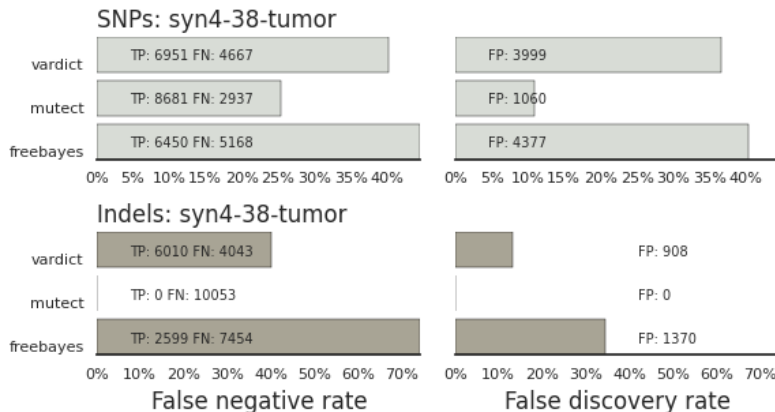- SNP + Insertion/Deletions
- Works on very deep targeted data

https://github.com/AstraZeneca-NGS/VarDictJava
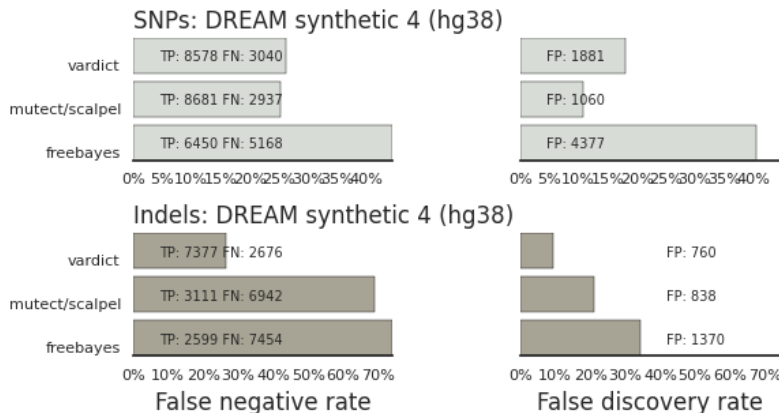
# DREAM synthetic dataset 4

| 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



SNPs: syn4-38-tumor

vardict — TP: 6951 FN: 4667 — FP: 3999

mutect — TP: 8681 FN: 2937 — FP: 1060

freebayes — TP: 6450 FN: 5168 — FP: 4377

0% 5% 10% 15% 20% 25% 30% 35% 40%     0% 5% 10% 15% 20% 25% 30% 35% 40%

Indels: syn4-38-tumor

vardict — TP: 6010 FN: 4043 — FP: 908

mutect — TP: 0 FN: 10053 — FP: 0

freebayes — TP: 2599 FN: 7454 — FP: 1370

0% 10% 20% 30% 40% 50% 60% 70%     0% 10% 20% 30% 40% 50% 60% 70%

False negative rate     False discovery rate

# VarDict sensitivity/precision after



SNPs: DREAM synthetic 4 (hg38)

| | | |
|---|---|---|
| vardict | TP: 8578 FN: 3040 | FP: 1881 |
| mutect/scalpel | TP: 8681 FN: 2937 | FP: 1060 |
| freebayes | TP: 6450 FN: 5168 | FP: 4377 |

0% 5%10%15%20%25%30%35%40%    0% 5%10%15%20%25%30%35%40%

Indels: DREAM synthetic 4 (hg38)

| | | |
|---|---|---|
| vardict | TP: 7377 FN: 2676 | FP: 760 |
| mutect/scalpel | TP: 3111 FN: 6942 | FP: 838 |
| freebayes | TP: 2599 FN: 7454 | FP: 1370 |

0% 10% 20%30% 40%50%60% 70%    0% 10% 20%30% 40%50% 60% 70%

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)))
```
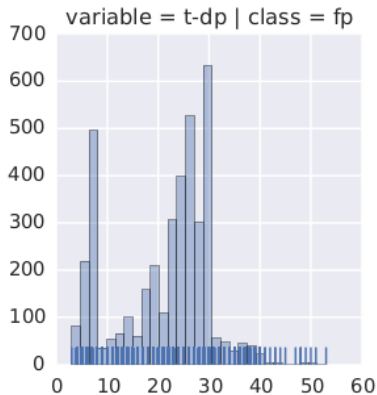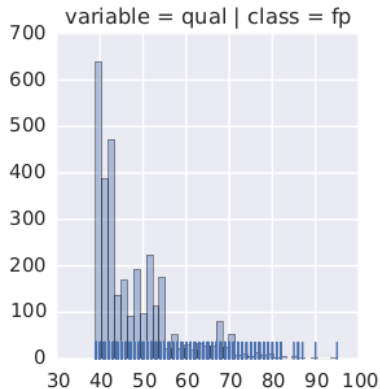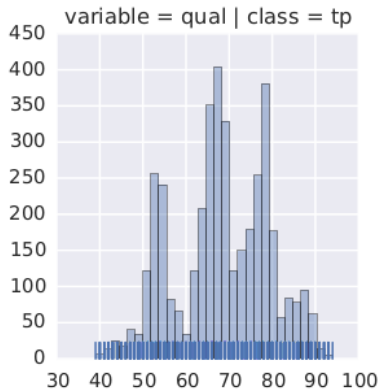
# Filter: mapping quality and number of mismatches

# Filter: low depth

# Filter: low quality

- Incorporate machine learning methods
- Generalize with additional datasets
- AML31: http://aml31.genome.wustl.edu/

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

# Structural variants critical in cancer

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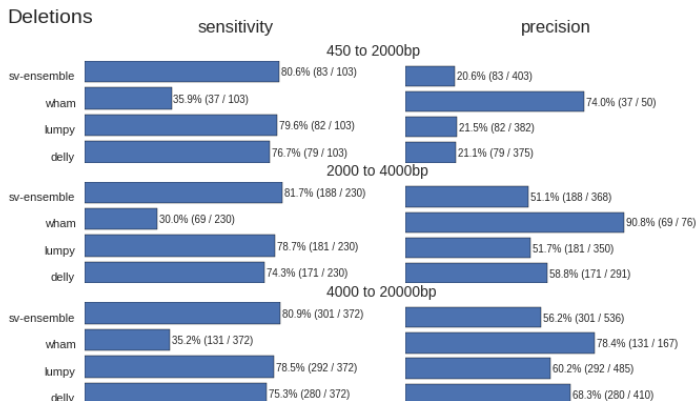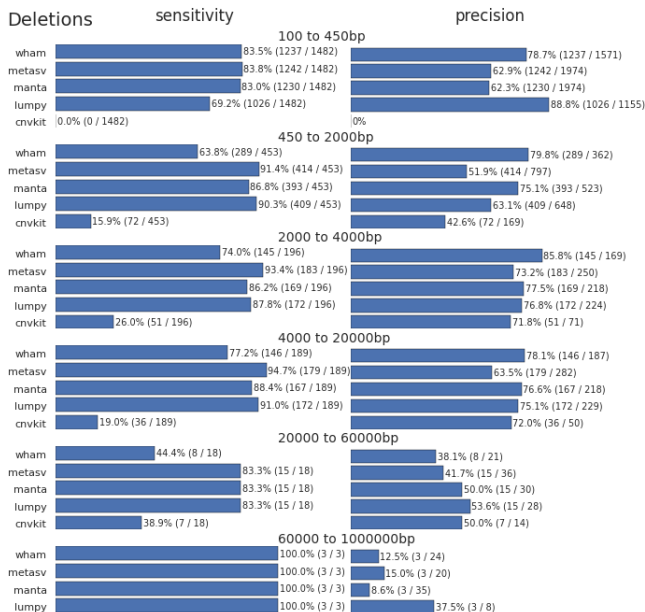- Manta: https://github.com/Illumina/manta
- CNVkit: https://github.com/etal/cnvkit
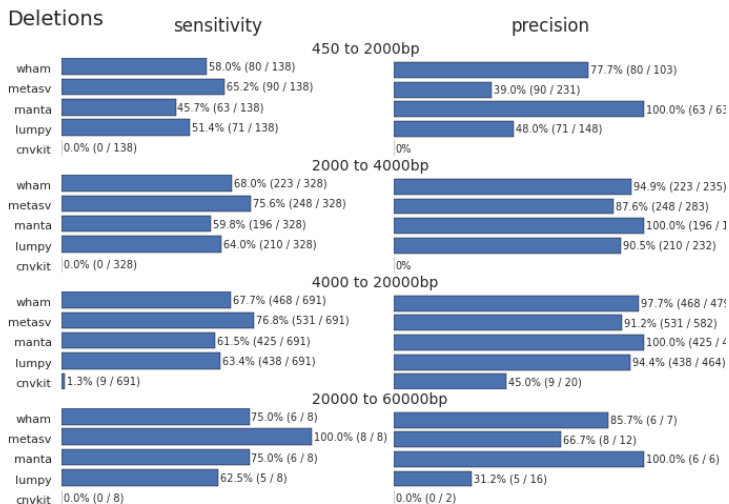- WHAM: https://github.com/zeeev/wham
- MetaSV: https://github.com/bioinform/metasv

# Results: Germline deletions

# Results: Somatic deletions

# Results: Somatic insertions

# Public cancer variant databases

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

# Summary

- Demonstrate current validation work in bcbio
- Human build 38
- HLA typing
- Low frequency cancer calling
- Structural variations + prioritization

http://bcb.io