Jose Blanca COMAV institute bioinf.comav.upv.es





## Genotype matrix

Marker Set Dimension								
		rs10001	rs10002	rs10003	rs10004	rs10005	rs10006	
Sample Set Dimension	Sample1	AA	GG	CC	AC	GG	CC	Cenatyne
	Sample2	AA	GG	СТ	AC	CG	СС	Dimension
	Sample3	AT	GG	Π	AA	CG		
	Sample4	Π	GG	Π	AA	GG	<b>5</b> 3	
	Sample5	AA	GG	Π	AA	GG	Π	
	Sample6	AA	GA	СС	AC	GG	СС	
	Sample7	AT	GG	CC	CC	GG	CC	
	Sample8	AA	GG	СТ	AA	GG	CT	

Genotype matrix: Samples x SNPs

...

A change in a read may due to:

- Sample contamination
- Cloning or PCR artifacts
  - Homopolymer stuttering
  - Hexamers in RNASeq libraries
- Sequencing errors
- Mapping errors
  - Problems related to structural variants
  - Alignment errors
- Reference genome errors
- Real SNVs

Two kinds of errors:

- False positives
- False negatives



## SNP calling errors

In a Hiseq 2500 line there are

if Q= 30:

- 150 millions of sequencing errors



## SNP calling alignment errors

```
Ref
         ...aqqttttataaaac---aattaaqtctacaqaqcaacta...
Sample
         ...aqqttttataaaacAAATaattaaqtctacaqaqcaacta...
         ...aggttttataaaac****aaAtaa
Read1
         ....ggttttataaaac****aaAtaaTt
Read2
Read3
         .....ttataaaacAAATaattaaqtctaca.......
                     CaaaT****aattaagtctacagagcaac.....
read4
                       aaT****aattaagtctacagagcaact....
read5
                         T****aattaaqtctacaqaqcaacta....
read6
```

Strategies to mitigate this problem:

- Fix the problem.
  - GATK, GLIA realignment.
  - It realigns the problematic regions (lots of SNPs or some indels).
  - Computationally slow.
  - It does not fix all problems.
- Avoid using the misaligned positions.
  - Samtools BAQ (calmd).
  - For each position It calculates the probability of being misaligned.
  - Avoid read ends

## **SNP** callers

Simple algorithms based on:

- Allele counts in reads
- Allele rate
- Total coverage

Complex bayesian algorithms add

- Hardy-Weimberg equilibrium
- Linkage disequilibrium

More common SNP callers:

- GATK
- Freebayes

Brad Chapman has a very interesting piece comparing the use of aligners and SNP callers.

- de-duplication with Picard MarkDuplicates, GATK base quality score recalibration and GATK realignment
- Minimal post-processing, with de-duplication using samtools rmdup and no realignment or recalibration.
- FreeBayes (v0.9.9.2-18): A haplotype-based Bayesian caller from the Marth Lab.
- GATK UnifiedGenotyper (2.7-2): GATK's widely used Bayesian caller.
- GATK HaplotypeCaller (2.7-2): GATK's more recently developed haplotype caller which provides local assembly around variant regions



GATK best-practice BAM preparation (recalibration, realignment)

http://bcbio.wordpress.com/



Minimal BAM preparation (samtools de-duplication only)

http://bcbio.wordpress.com/

## **Required coverage**

Depends:

- Confidence required
  - SNP vs Genotypes
  - Clinical diagnosis vs marker search
- Sample complexity
  - Diploid
  - Polyploid
  - Pooled samples
  - Heterogeneous cancer
- Ploidy
- Population structure
  - e.g. F2 population vs random mating population

% of missing calls per SNP

SNP depth or Genotype depth

SNP observed heterozygosity

SNP quality

SNP density

Sample depth

% of missing calls per sample

Sample observed heterozygosity



Number of SNPs called at different missing rates



Number of SNPs with different observed heterozygosities



Observed heterozygosity and called GT rate per sample



Depth distributions per sample for all genotype calls and for the non-missing genotype calls

## **SNP** filtering

Filtering the SNPs after the SNP calling is a critical task.

The objetive is to remove false positives, but keeping the correct SNPs.

The filters and parameteres depend of the data and the application.

Vcftools, vcffilter, Annovar...

Home-made filters



#### SNPs Low quality:

SNP callers usually assign a quality (probability) to the SNPs. We can filter out the SNPs with lower qualities.

#### Genotype Low quality:

It is also possible to filter out not SNPs, but genotypes. In this case the genotype is usually set to not determined

#### **Missing data:**

We could filter the SNPs with large amount of missing genotypes.

#### Number of alleles:

It is possible to remove the monomorphic SNPs or to filter out the SNPs that are not biallelic.

Minor (or Major) Allele Frequency (MAF):

#### **Observed Heterozygosity:**

High heterozygosities in populations are suspicious

#### **High Coverage:**

SNPs with an excessive coverage can be false positives due to duplicated regions in the sequenced sample not found in the reference genome.

#### **Highly Variable Region:**

Having regions with too many SNPs it is also a sign that we are piling up reads from repeated regions.

#### Low Complexity Region:

It has been shown that due to problems with the PCR and the alignment the low complexity regions are particularly prone to false positive SNPs.

#### Linkage Disequilibrium:

If we have genotype a segregant population it could be useful to filter out the SNPs that are not in linkage disequilibrium with their closest SNPs.

#### Kind

We can filter the SNVs according to its type: SNV, indel, complex or structural variation.

#### Aminoacid change:

We can select the SNPs with large impacts in the coded proteins.

#### By genome localization:

We can choose SNPs in a determinate region or kind of features (exon, UTR....)

### **SNP** annotation



### **SNP** annotation

#### RESEARCH ARTICLE

#### Genetic Variation in an Individual Human Exome

Pauline C. Ng, Samuel Levy, Jiaqi Huang, Timothy B. Stockwell, Brian P. Walenz, K...

Synonymous		10,413
Heterozygous	Novel	551
	dbSNP	5,183
Homozygous	Novel	98
	dbSNP	4,581
Nonsynonymous		10,389
Heterozygous*	Novel	557
	dbSNP	5,047
Homozygous	Novel	215
	dbSNP	4,570

<sup>\*</sup>All heterozygous novel nonsynonymous SNPs were manually inspected (see Methods).

doi:10.1371/journal.pgen.1000160.t001









### File format: VCF



### File format: VCF

**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 1 ACG A A--T **Complex events** Alignment VCF representation ACGT POS REF ALT 1 ACG AT A-tT Large structural variants VCF representation POS REF ALT INFO 100 T <DEL> SVTYPE=DEL;END=300

### File format: GFF

Based on a sequence ontology

Annotation of regions in sequences

##gff-version 3
##sequence-region ctg123 1 1497228
ctg123 . gene 1000 9000 . + . ID=gene00001;Name=EDEN
ctg123 . TF\_binding\_site 1000 1012 . + . ID=tfbs00001;Parent=gene00001
ctg123 . mRNA 1050 9000 . + . ID=mRNA00001;Parent=gene00001;Name=EDEN.1

### File format: BED

Annotation of regions in sequences

#### chr22 1000 5000 cloneA 960 + 1000 5000 0 2 567,488, 0,3512 chr22 2000 6000

The BED format uses 0-based coordinates for the starts and 1-based for the ends. So the 1st base on chromosome 1 would be:

chr1 0 1 first\_base