

# **Introduction to Cancer genomics**



### Darwinian Evolution Neutral Evolution

Raheleh Rahbari

raheleh.rahbari@sanger.ac.uk





# Learning objectives

- Cancer Vs. Normal genomes & sources of genome variation between the two
- Bioinformatic approaches to detect cancer genome variation
- How cancer genome analysis may be used to guide cancer patient management



# Human genome annotated with data related to genes implicated in disease





Circos, Martin Krzywinski et al., Genome Res. 2009

# Ways that the genomes can be modified

Point mutations

A

В

D

- Copy number alteration
- C Structural rearrangements
  - Genes & regulatory elements
- E Pathways
- F Homologous and epigenetic modifications



# The timing of the somatic mutations



Relapse after chemotherapy can be associated with resistance mutations that often predate the initiation of treatment.



Michael Stratton, Peter J. Campbell, and P. Andrew Futreal, Nature 2009

# Somatic mutation frequencies observed in exomes from 3,083 tumor-normal





Lawrence et al., nature 2013

### Hallmarks of cancer

Tumor calls acquire abnormal abilities by co-opting normal cell behavior





Hanahan and Weinberg, Cell 2011



# Applications of next-generation sequencing to cancer



# General DNA sequencing workflow









# **Background mutation detection**

### Drivers

- Implicated in oncogenesis
- Reside in cancer genes
- Inform biology
- Targets for therapy

Known cancer genes Novel cancer genes Germline

### Passengers

- biologically inert
- reflect DNA exposures
- life History of Cancer

### Signature

- spectrum of variants
- sequence context





# Background mutation detection







## **Class of mutations**

129.1 Mb



Chromosome 8

127.6 Mb

Insertions & deletions



#### Rearrangements



# Background: Genetics & Somatic Substitutions

- Base Pair Substitution single base pair is changed
- Most common mutation type



Image from http://rosalind.info/glossary/point-mutation/

- When occur in a gene
  - Silent
  - Missense
  - Nonsense



## Which Somatic Mutation Substitution Caller should I use?

### WGS/WES/Targeted Analysis Caveman with

### **Targeted Shearwater**

- Subclonal variants in deeply sequenced samples, can be <u>unmatched</u>
- Must ensure normal panel is optimised for your specific analysis



ICGC/TCGA Pilot-63, 250,000 of 6.4 Million mutations validated by deep sequencing on Illumina HiSeq





- User can specify a single global copy number estimate. Default is 2/5.
- Alternatively can use copy number inferred by ASCAT

At each position the algorithm then considers all possible genotypes that can be made from a single variant allele. The ploidy and copy number determine possible genotypes

### Space of Genotypes (Normal Ploidy= $N_p$ , Tumour Copy Number= $T_p$ )

Let A represent the reference allele and B the variant allele (here  $N_p = 2$  and  $T_p = C$ )

Somatic Genotypes:  $G_N = AA$  and  $G_T = \underbrace{A \dots A}_{C-k} \underbrace{B \dots B}_{k \ copies}$  where  $1 < k \le C$ SNP Genotypes:  $G_N = AB$  or BB and  $G_T = \underbrace{A \dots A}_{C-k} \underbrace{B \dots B}_{k \ copies}$  where  $0 \le k \le C$ 

Where  $G_N$  and  $G_T$  are the genotypes of the Normal and Tumour respectively.

Doesn't handle:

- Somatic mutations at SNP sites
- Multi-allelic variants



### Caveman directly estimates genotype probabilities

 $P(G_N, G_T | D) \propto P(D | G_N, G_T) P(G_N, G_T)$ 

Where D is the set of overlapping reads at position p. Here we focus on the first term

 $P(D|G_N, G_T) = \prod_{R \in Normal Reads} P^{(n)}(R_p|G_N, G_T, \theta) \prod_{R \in Tumor Reads} P^{(t)}(R_p|G_N, G_T, \theta)$ 

Incorporates prior probability of germline/somatic mutation.

Let genotype  $\psi(G)=B(G)/(A(G)+B(G))$  represent the proportion<sup>\*</sup> of non-ref base (B) in genotype G. Let c be the proportion of normal contamination of tumour. The covariates for read R at position p are denoted by  $\theta$  and let the called base be denoted C



Implementation also includes reference bias correction

# Background: Genetics & somatic ins/dels

- Small Insertion & Deletions
- Less than 200base pairs
- · Are known to drive cancer development
- 3 Classes:



Can be disruptive if occur in genes (switch genes on/off)

- Inframe
- Frameshift



#### Background: Genetics & somatic ins/dels GAT GAT Target GAT GA1 Reference g.12:130454983 130454985del Control GAT GAT GAT GAT GAT GAT GAT

#### **Artefact Mutations**

- Sequence errors (bubbles on flow cell)
- D PCR errors
- Sample Prep errors (FFPE, 8oxoG)

#### **Analysis Challenges**

- Clonal composition (Aberrant cell fraction)
- □ Changes in Copy Number
- Data Analysis Mapping
  - $\hfill\square$  Accuracy of mapping
  - Alignment left or right justified

#### Sensitivity:

□ Call mutations at ~10% VAF



# Which Somatic Mutation Ins/dels Caller should I use?

#### 1.00 broad\_snowman sanger dkfz wustl 0.75 broad\_mutect WGS/WES/Targeted Cancer smufin Analysis precision Pindel 0.50 novobreak crg\_clindel 0.25 0.00 0.25 1.00 0.50 0.75 0.00 sensitivity



### Pindel



- Pattern Growth Approach (sequential pattern mining algorithm)
- Adapted to handle complex structural variants

Pindel optimised to call somatic small ins/dels in cancer samples Requires Tumour and Matched-Normal Bam files Detect Deletions > 10kb Detect Insertions (factor of read length) ~50bp Can also call tandem duplications, inversions





### Pindel: How it works for an insertion **Read does not map** Sample Reference Anchor Read Anchor Insertion in sample reads Read Read is split into 3 with velicome inserted sequence removed and now maps

# Background: Copy Number Algorithms

Algorithms use the read depth to detect Copy Number changes

- Increase number of reads -> Amplification
- Decreased number of reads -> Deletion One Copy (Loss of Heterozygosity, LOH)
- No reads -> Deletion Both Copies (Homozygous Deletion)



In Addition, by using Germline Information (SNPs) can also work out which allele is affected (paternal or maternal)

# **Types of Structural Variants**



Scherer SW et al. Challenges and standards in integrating surveys of structural variation. Nat Genet. 2007 Jul;39(7 Suppl):S7-15.

# DNA Breakpoint: Signposts for Structural Variants

**DNA Breakpoint**: All Structural Variants are characterised by two unexpected fragments of DNA being next to one another



Breakpoint – Deletions, Tandem Duplications
Breakpoints – Insertions, Inversions, Translocation
Many Breakpoints – Amplicon, Chromothripsis, Chromoplexy



# 1. Paired End Analysis (aberrant read mapping)

Expect 2 reads to map with set gap (insert size), if not they could indicate structural rearrangement



### 2. Read depth





# Allele-Specific Copy number Analysis of Tumours

- Developed by Peter Van Loo et al
- Can detect Copy Number Alterations (CNAs) in Tumours
- Allelic-specific paternal or maternal changes (Often cannot determine this so algorithm identifies major/minor copy number changes)
- Estimates the Absolute Copy Number
- Estimate normal & tumour composition of a sample



# **ASCAT results - Plots**

### **Sunrise Plot**

### **Copy Number Profile**





**Copy Number Profile** – This is the final result with integer (clonal) copy number states. Red indicates the major and green the minor allele copy number

**Sunrise Plot** – This plot shows the goodness of fit against sample ploidy and purity. The green cross indicates the area of best fit (dark blue)

# **Brass Overall Workflow**







# **Reading list**

- Somatic mutant clones colonize the human esophagus with age: http://science.sciencemag.org/content/362/6417/911/tab-article-info
- Tumor evolution. High burden and pervasive positive selection of somatic mutations in normal human skin.

https://www.ncbi.nlm.nih.gov/pubmed/25999502

• Tissue-specific mutation accumulation in human adult stem cells during life <a href="https://www.ncbi.nlm.nih.gov/pubmed/27698416">https://www.ncbi.nlm.nih.gov/pubmed/27698416</a>

raheleh.rahbari@sanger.ac.uk

