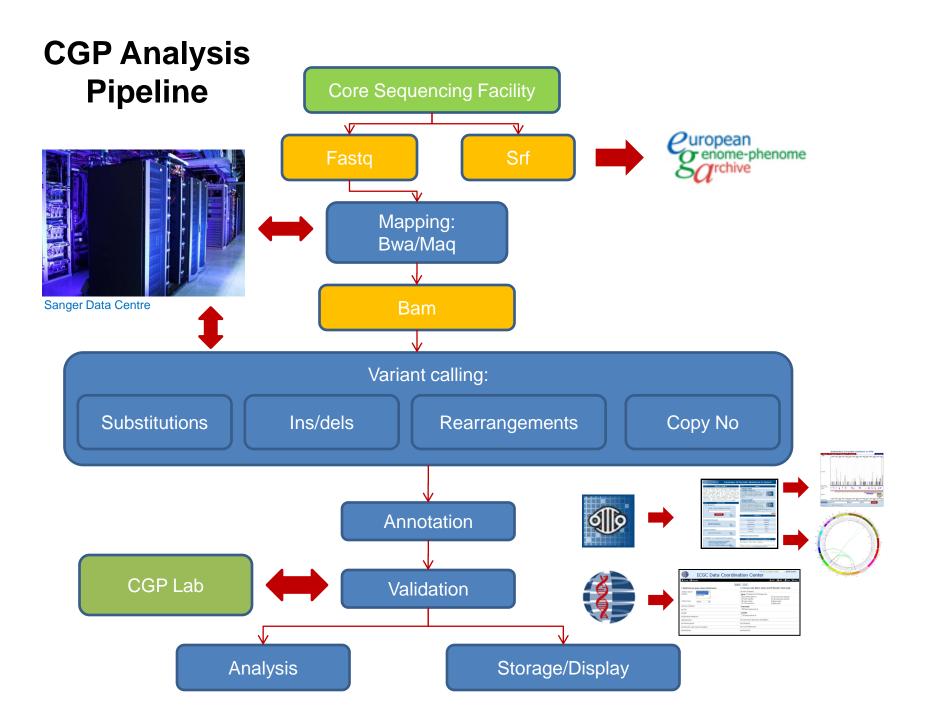# T5B: Developing a substitution calling algorithm to analyse breast cancer exomes by next generation sequencing

David Jones & Andy Menzies

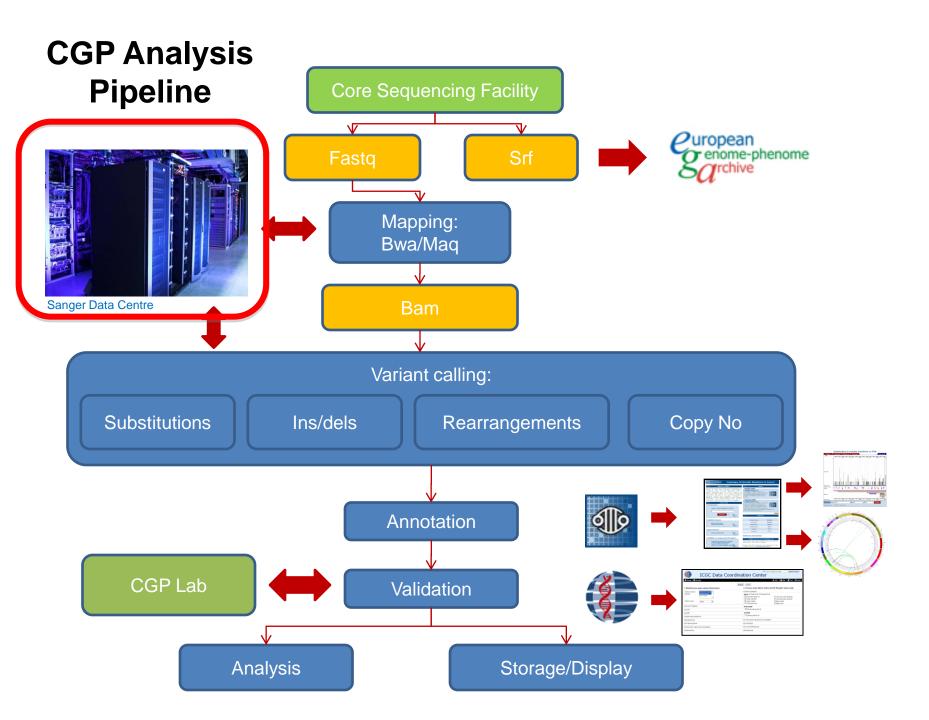Cancer Genome Project

Wellcome Trust Sanger Institute

# Introduction

- Overview of the analysis pipeline
  - Why do we need a pipeline?
  - Initial sequence analysis
  - Data release and presentation

- Variant analysis techniques
  - Rearrangement detection
  - Insertion/deletion detection
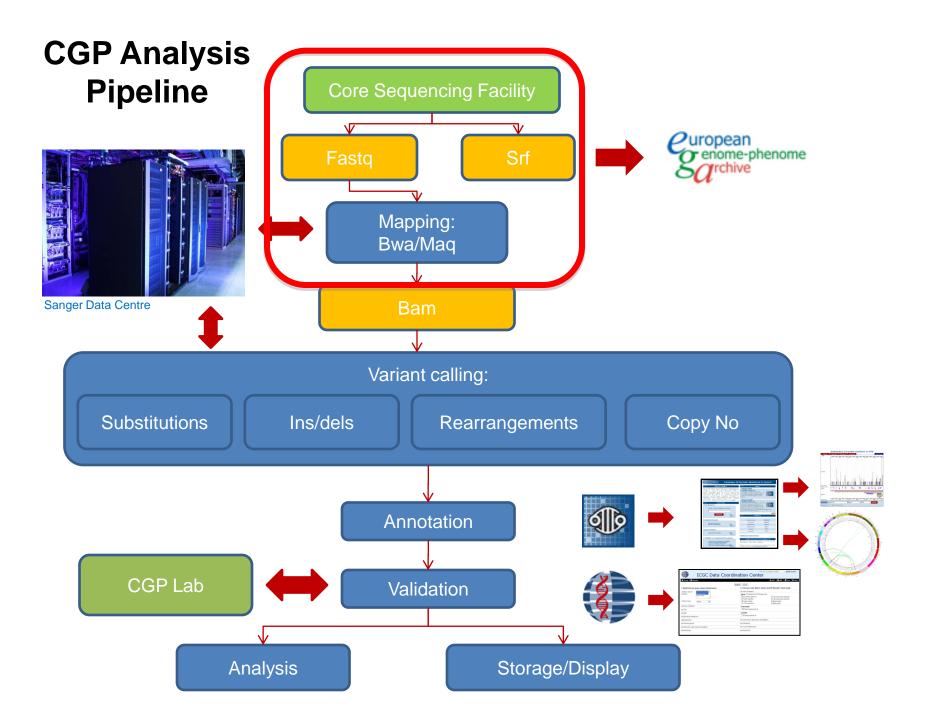  - Substitution detection

# CGP Analysis Pipeline

Core Sequencing Facility

Fastq

Srf

Sanger Data Centre

Mapping: Bwa/Maq

Bam

Variant calling:

Substitutions

Ins/dels

Rearrangements

Copy No

Annotation

CGP Lab

Validation

Analysis

Storage/Display

# Why: Data Volumes

- ICGC – International Cancer Genome Consortium
  - **ICGC Goal:** To obtain a comprehensive description of genomic, transcriptomic and epigenomic changes in 50 different tumor types and/or subtypes which are of clinical and societal importance across the globe.
  - Each ICGC project consists of 500 matched normal/tumour sample pairs

- EU BASIS
  - Breast ER+ve, HER2-ve – 500 norm/tum pairs
- Breast Cancer
  - Triple –ve, lobular & others – 500 norm/tum pairs
- Also assisting with a number of other ICGC projects
- And we have other non-ICGC projects too
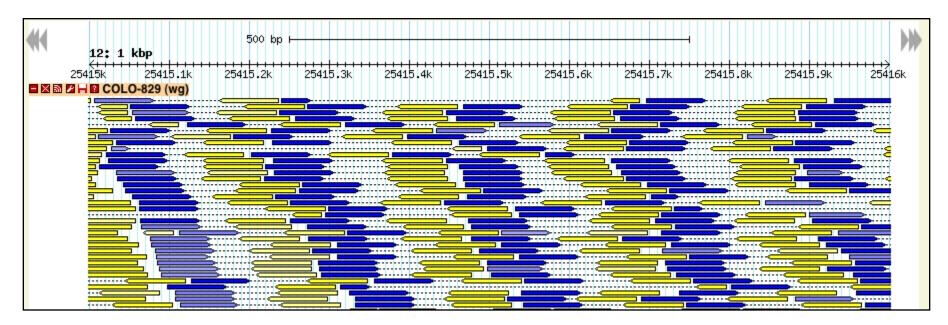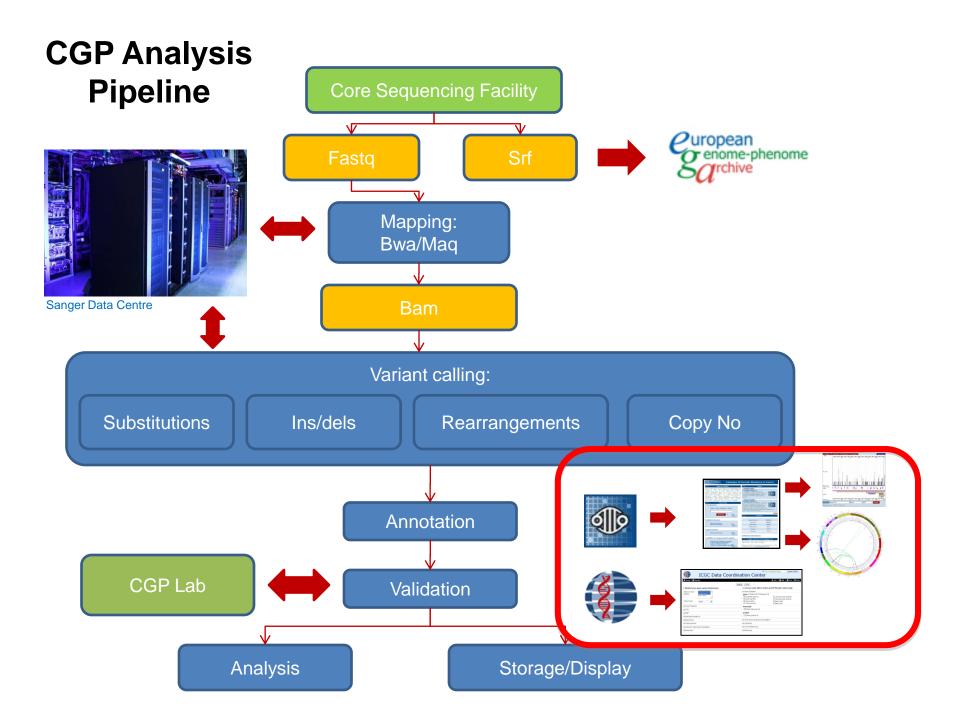
# Hardware: The Farm

- CPUs
    - Mixture of dual and quad core multi processor machines
    - Total of 3920 cores
- Memory
    - Total of 8846Gb memory
    - Average of 2.2Gb per core
- Disk
    - 215TB Luster file system – working space
    - 350TB NFS file system – long term data storage
    - ~1TB to store a Normal/Tumor full genome pair & analysis

**CGP Analysis Pipeline**

Core Sequencing Facility

Fastq

Srf

Mapping: Bwa/Maq

Bam

Sanger Data Centre

Variant calling:

Substitutions

Ins/dels

Rearrangements

Copy No

Annotation

Validation

CGP Lab

Analysis

Storage/Display
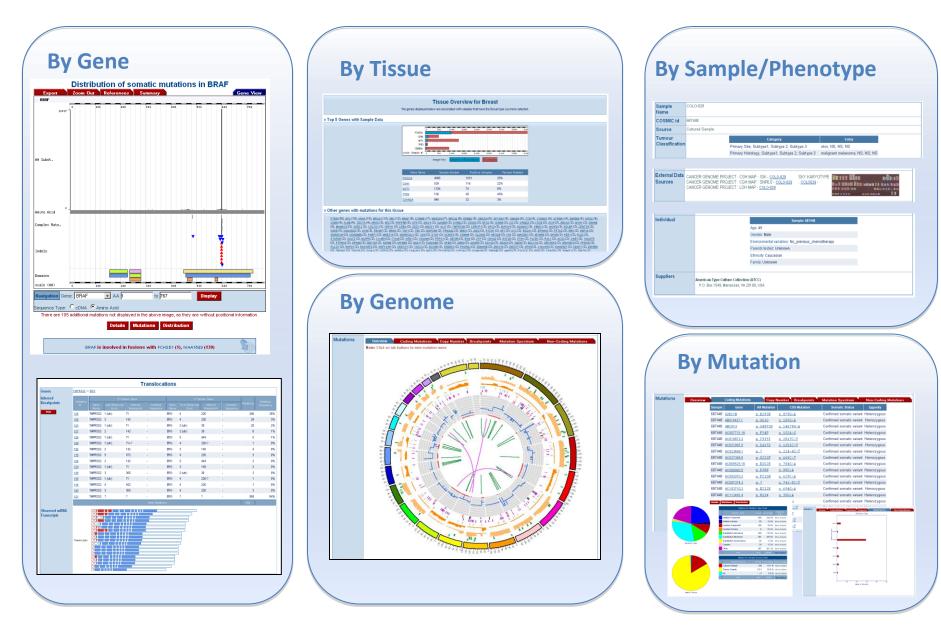
# Initial Data Processing

- Illumina GA2 – paired end sequencing
- Use BWA and MAQ
    - H Li et al Bioinformatics 2010
    - H Li et al Genome Research 2008
- Need to use fast aligners because of the vast number of individual reads generated by the sequencers
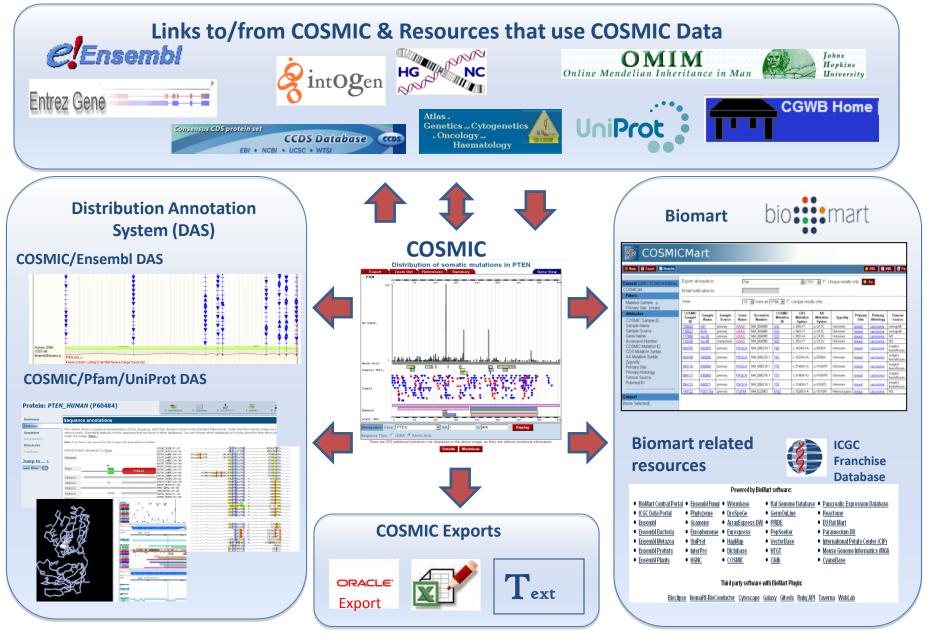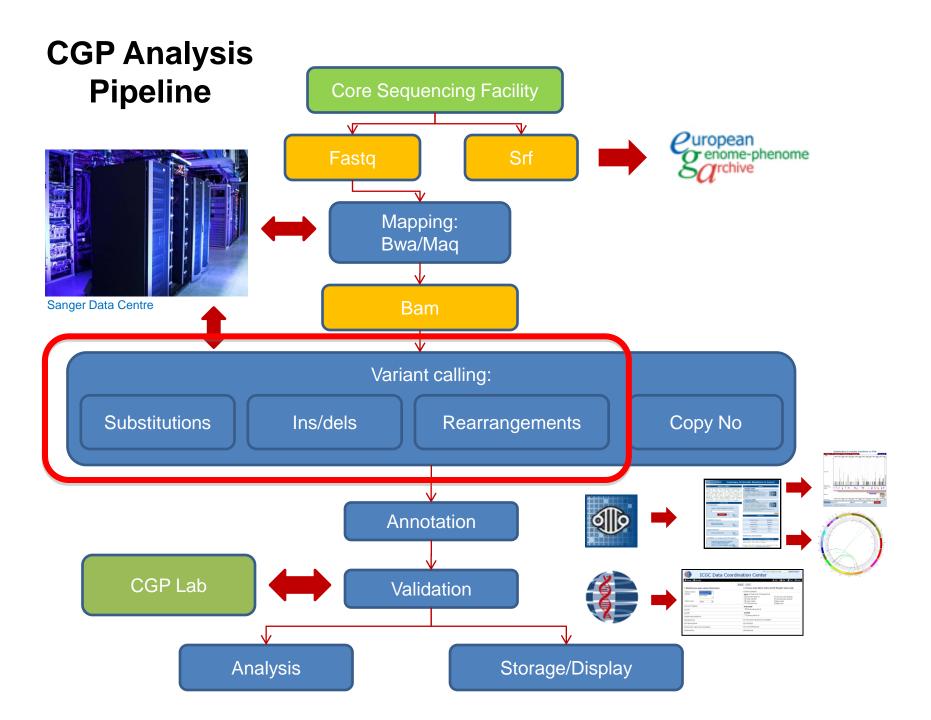
# CGP Analysis Pipeline

Core Sequencing Facility

Fastq

Srf

*european genome-phenome archive*

Sanger Data Centre

Mapping: Bwa/Maq

Bam

Variant calling:

Substitutions

Ins/dels

Rearrangements

Copy No

Annotation

CGP Lab

Validation

Analysis

Storage/Display

# COSMIC Interfaces and Web-based Tools

# COSMIC Data Interoperability

# CGP Analysis Pipeline

Core Sequencing Facility

Fastq

Srf

european genome-phenome archive

Sanger Data Centre

Mapping: Bwa/Maq

Bam

Variant calling:

Substitutions

Ins/dels

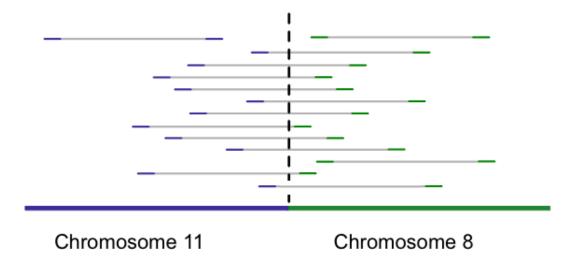Rearrangements

Copy No

Annotation

CGP Lab

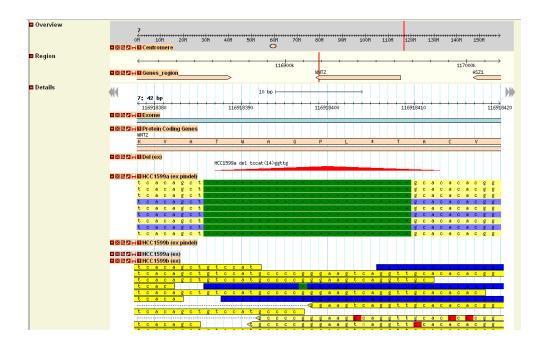Validation

Analysis

Storage/Display

# Variant Calling: Rearrangements

- See poster 'Identifying Structural Rearrangements via Local Assembly of Next-Generation Sequence Data' by John Marshall
- Identify informative read-pairs
  - Individual reads map accurately
  - Read-pair form an unexpected insert size
- Group read-pairs spanning the same putative break point
- Use the Velvet *de novo* assembler to reconstruct the sequence across the break



Chromosome 11          Chromosome 8

# Variant Calling: Insertions/Deletions

- Pindel – K Ye *et al* Bioinformatics 2009
- Uses the read-pairs with:
  - One uniquely mapped read
  - One read is either unmapped or has gapped alignment
- Pindel then realigns the mis-mapped reads expecting greater divergence from the reference than BWA/MAQ

# Variant Calling: Substitution

- Introduction to CaVEMan (**Ca**ncer **V**ariants through **E**xpectation **Ma**ximisatio**n**)

- Usage / Requirements

- Post processing

- Results

# What is CaVEMan?

- Single base substitution calling algorithm

- Java

- Indexed bam as input

- Picard[1]

- 3 way comparison

  - reference

  - normal + tumour

- Expectation Maximisation algorithm[2]

1. http://picard.sourceforge.net/index.shtml

2. C.B. Do & S. Batzoglou (2008). 'What is the expectation maximization algorithm?'. Nat Biotech 26(8):897-899

# Expectation Maximisation Algorithm

- Two step iterative algorithm

- M(aximisation) step
  - Build profile of sequencing errors
  - reference base, called base, base quality, read position, mapping quality, lane as covariates…

- E(xpectation) step
  - Use profile to call substitutions
  - Naïve Bayesian classifier

# What is CaVEMan?

- Modular - can update E-step (sub. calling) parameters without rerunning M-step.
- Flexible
    - More new sequencing technologies coming.
- Can be used on SOLiD (not yet tuned)
- Many optional parameters:

Normal contamination

Include Smith-Waterman reads

SNP probability cut-off

Expected mutation frequency

Reference bias

Mutation probability cut-off

Minimum base quality to include

Expected SNP frequency

# Output of Results

- Result for each genomic position
- Probability assigned for every possible genotype given the reference base
- If above defined cut off - called somatic/SNP, written to file

```
1   50084056   C   1.0e+00 2.5e-06 CC/TT   1.0e+00 CT/TT   2.5e-06   0   26  0   0   0   0   0   36
```

- 'Normal' bases - to file, info for every position.

```
1   14562   A   7.9e-07 8.4e-05 AA/AA   1.0e+00 AC/AA   2.2e-05
```

# Usage / Requirements

# Usage

- Designed for a compute farm environment.
- Initial split step - farm sized chunks.

|            | Exome  | Genome    |
|------------|--------|-----------|
| Jobs       | ~200   | 2000-3000 |
| CPU sec/job | 3372  | 6850      |
| Mem (Mb)/job | 3500 | 3500      |

# Usage

- Somatic Calls (before post processing)

|        | avg.    | min.  | max.   |
|--------|---------|-------|--------|
| Exome  | 8513.6  | 557   | 64824  |
| Genome | 84551.4 | 40178 | 163155 |

# Post processing

# Capillary confirmation before post processing

- Better to overcall than miss potentially interesting substitutions
    - Sensitivity over specificity
    - Computationally less intensive to filter after calling
- Example from exomic data:
    - 935 Putative somatic substitutions
    - 131 confirmed as somatic (14%)

# Minor causes of false positives in exome data

## Slippage at mononucleotide tracts



## Slippage at poly (n) nucleotide tracts

# Minor causes of false positives in exome data

Germline INDELS  (with & without SNPs)

G/A SNP

C/T artefact

# Major causes of false positives in exome data

# i) Poor quality data at ends of reads
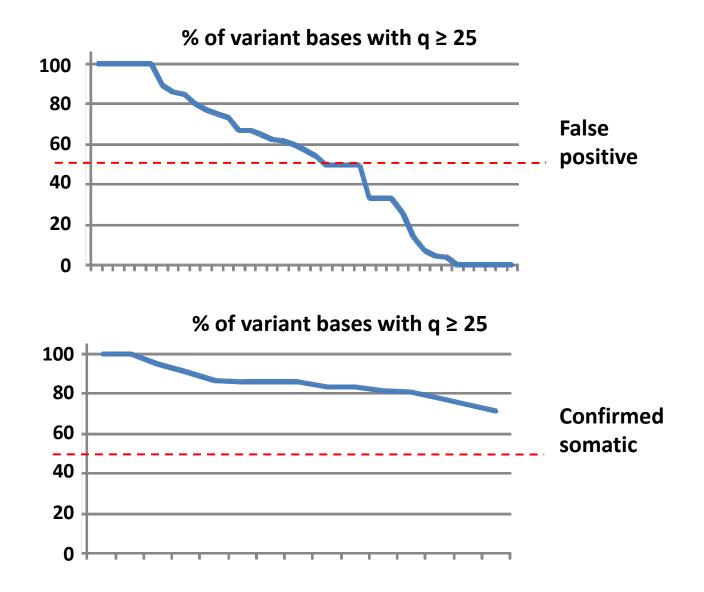
# Post processing substitutions

If coverage (in tumour) ≥10,

≥1 base call reporting a variant in the 2nd third of the read.

OR

If coverage (in tumour) ≤10,
≥ 1 base call reporting a variant in the 1st or 2nd third of the read.

# ii) Low quality bases called as mutations

### % of variant bases with q ≥ 25



False positive

### % of variant bases with q ≥ 25



Confirmed somatic

**Post processing substitutions**

If coverage (in tumour) ≥10,
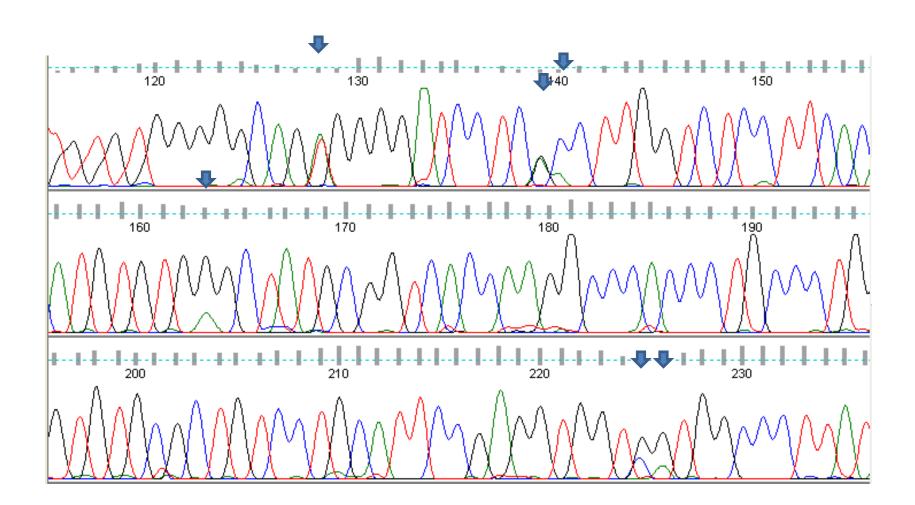
≥1 base call reporting a variant in the 2nd third of the read.

OR

If coverage (in tumour) ≤10,
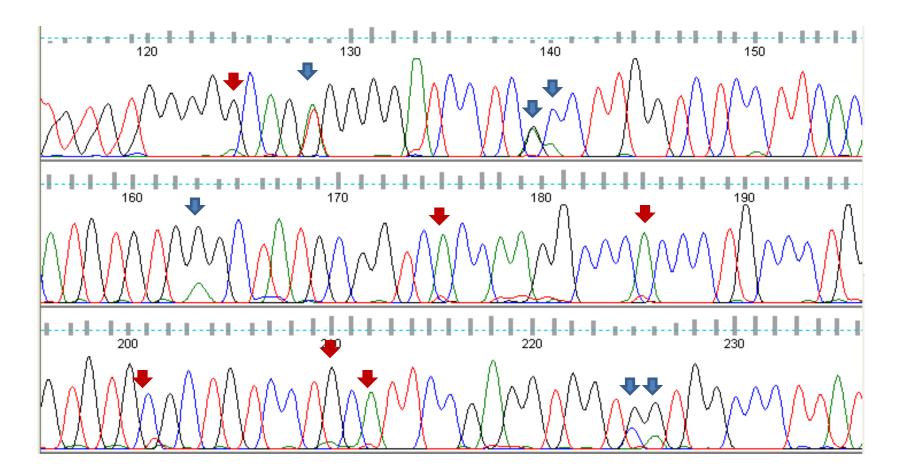≥ 1 base call reporting a variant in the 1st or 2nd third of the read.

≥1/3 base calls reporting a variant (in tumour) high quality (≥ 25)

## iii) Mapping errors called as mutations

# iii) Mapping errors called as mutations



| Match | Percent | Chr | Start | Finish |
|---|---|---|---|---|
| 19-334 | 98.8% | **8** | 11226692 | 11227007 |
| 19-334 | 98.5% | **18** | 11634476 | 11634791 |
| 19-334 | 97.8% | **18** | 11600261 | 11600576 |
| 19-334 | 97.5% | **17** | 30544458 | 30544773 |
| 19-334 | 97.2% | **17** | 7326694 | 7327009 |

## Post processing substitutions

If coverage (in tumour) ≥10,

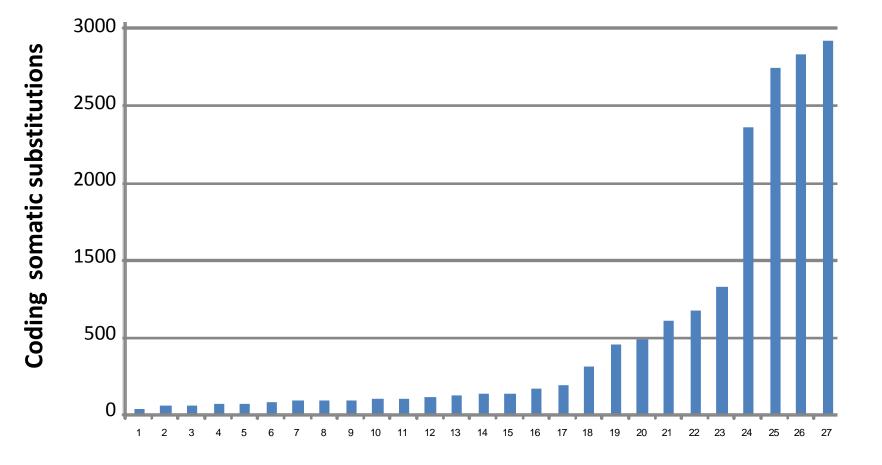≥1 base call reporting a variant in the 2nd third of the read.

OR

If coverage (in tumour) ≤10,
≥ 1 base call reporting a variant in the 1st or 2nd third of the read.
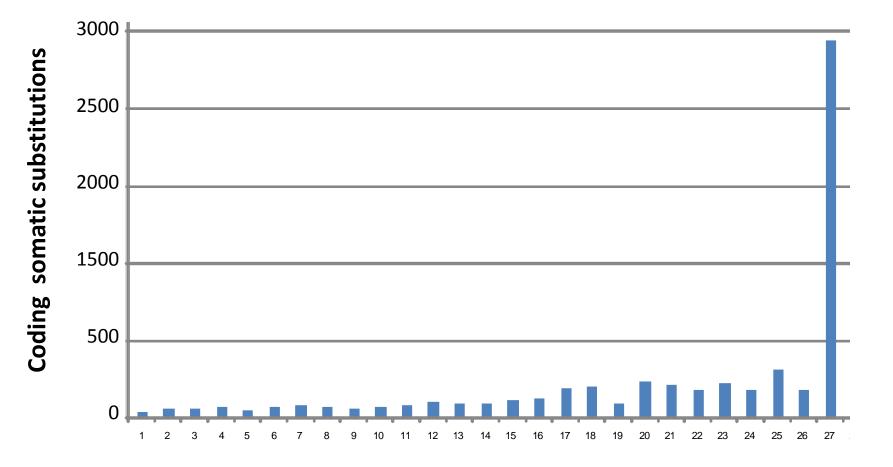
≥1/3 base calls reporting a variant (in tumour) high quality (≥ 25)

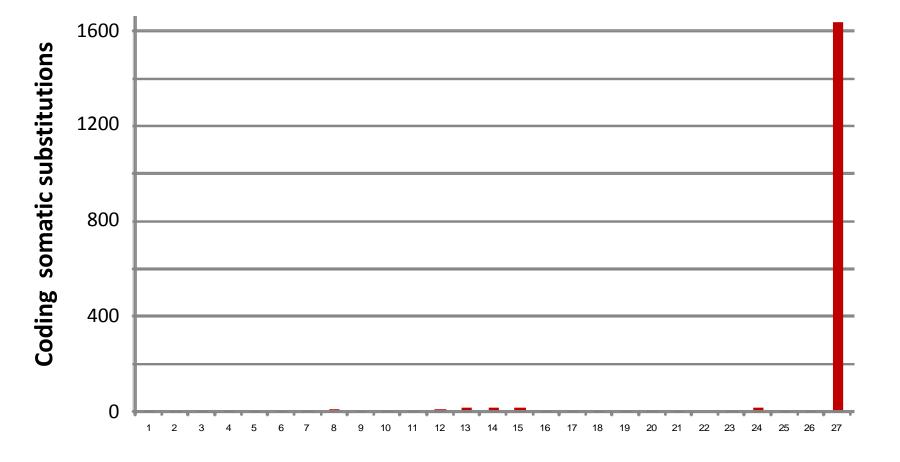≤1 base calls reporting a variant in (ANY of 28 normals) high quality (≥20)
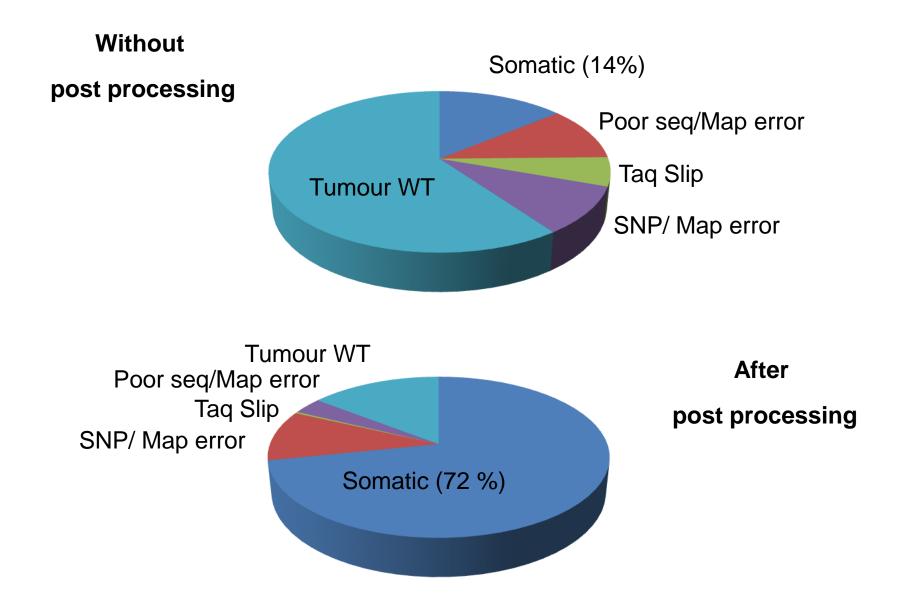
**Variants called:- without post processing**

# Variants called:- with post processing

# Variants called:- with known SNP I.D

**Variants called:-**

**Without**

**post processing**

Somatic (14%)

Poor seq/Map error

Taq Slip

SNP/ Map error

Tumour WT

Tumour WT
Poor seq/Map error
Taq Slip
SNP/ Map error

Somatic (72 %)

**After**

**post processing**

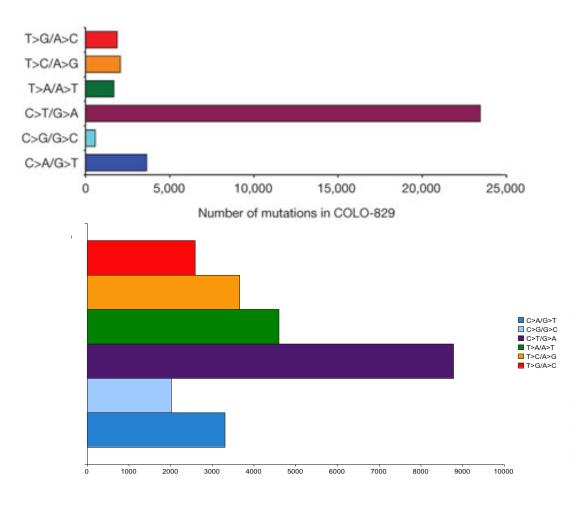# Results

# COLO-829

- E.D. Pleasance *et al.*

- 522 validated substitutions - NCBI36[1]

- CaVEMan - GRCh37 (liftover)
  - Missed               3
  - Called               519
  - Novel (unvalidated)        24965

1. E.D. Pleasance *et al* (2009). A comprehensive catalogue of somatic mutations from a human cancer genome. Nature 463(7278):191-196

# COLO-829 - Novel Calls

- 24965 novel (unvalidated)    subs
- Does the mutation spectrum match?

27 primary breast tumour and matching normal exomes

# 27 Breast Exomes

• Known cancer genes at frequencies in concordance with literature

    • PIK3CA, TP53, AKT1, NF1, MAP2K4, GATA3, PTEN and CDH1

• ~500 subs in  over 400 genes, many in >1 sample

    • Currently validating and evaluating with more breast exomes

# Summary

- Pipeline
- CaVEMan
  - Substitution calling Expectation Maximisation algorithm attempting to deal with sequencing errors.
  - Performance / Useage
  - False positives + Post processing
  - Results - COLO-829 and 27 Breast exomes.

# Thanks to.....

Peter Campbell

Phil Stephens

Keiran Raine

Serena Nik-Zainal

Ignacio Varela

Adam Butler

Jon Teague

Mike Stratton

Andy Futreal

Rest of the CGP

Wellcome Trust