# Introduction to Genomics England

## MED676 MSc in Genomic Medicine

**Matthew Parker, Ph. D**
Lead Bioinformatician
Sheffield Diagnostic Genetics Service

Sheffield Children's **NHS**
NHS Foundation Trust

- Company owned wholly by the department of health
- Delivery of 100k genome sequences – around 40,000 affected individuals
  - Rare disease families
  - Cancer tumour/normal pairs
- NHS transformation
- Education
- Collaborations with Research/Pharma
- Currently around 17,000 genomes sequenced

Genomics England Limited are tasked with delivering the 100KGP

**Sheffield Children's** **NHS**
NHS Foundation Trust

"

**"Sequence 100,000 genomes from patients with cancer, rare disorders, and infectious disease, and to link the sequence data to a standardised, extensible account of diagnosis, treatment, and outcomes. "**

"

# Aims

- **Patient benefit:** providing clinical diagnosis and in time, new or more effective treatments for NHS patients.

- **New scientific insights and discovery:** with the consent of patients, creating a database of 100,000 whole genome sequences linked to continually updated long term patient health and personal information for analysis by researchers.

- Accelerating the **uptake of genomic medicine in the NHS**: working with NHSE and other partners to deliver a scalable WGS and informatics platform to enable these services to be made widely available for NHS patients & creating a mechanism to both continually improve the accuracy and reliability of information fed back to patients and add to knowledge of the genetic basis of disease.

- Stimulating and enhancing **UK industry and investment**: by providing access to this unique data resource by industry for the purpose of developing new knowledge, methods of analysis, medicines, diagnostics and devices.

- Increasing **public knowledge** and support for genomic medicine: delivering an ethical and transparent program which has public trust and confidence and working with a range of partners to increase knowledge of genomics.

## Why Genomes?

- Why Sequence Whole Genomes?

- Patients without the "normal" genetic aberrations associated with their phenotype – where to look?
- Cost reducing
- Non-coding changes could result in disease
- Extra Value:
  - More reliable
    - Small Insertion Deletions Calling
    - Copy Number Calling
    - Structural Variant Calling
- Better representation of true variant allele frequency (cancer)
- Easy to look at "new" causes of disease

HiSeq Xten WGS Sequecning

**Can anyone think of any reasons why genomes might not give us all the answers?**

Sheffield Children's **NHS**
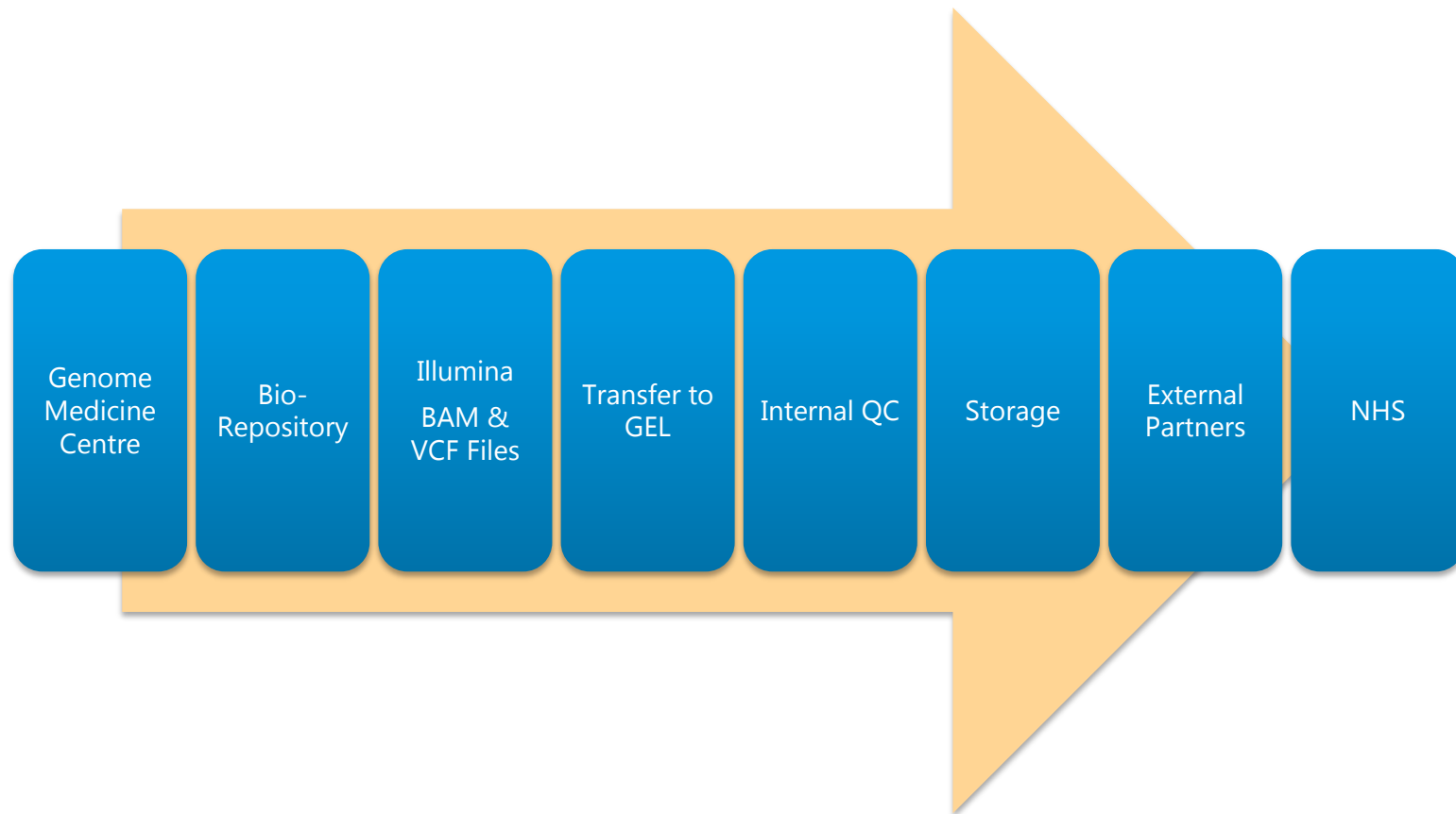NHS Foundation Trust

- Short reads (150bp) mean Illumina WGS is not great for:
  - Pseudogenes & Gene duplications
  - Repetitive regions
  - Regions of high GC/AT
  - Very complex structural variations
- Incidental findings
- Findings with little evidence of effect
- Slow
- Huge amounts of data that need to be kept for 30yrs



Short reads not always appropriate

## Data Generation

- How do Genomics England generate WGS data for the NHS?

# Overview of Sample → Result

Genome Medicine Centre → Bio-Repository → Illumina BAM & VCF Files → Transfer to GEL → Internal QC → Storage → External Partners → NHS
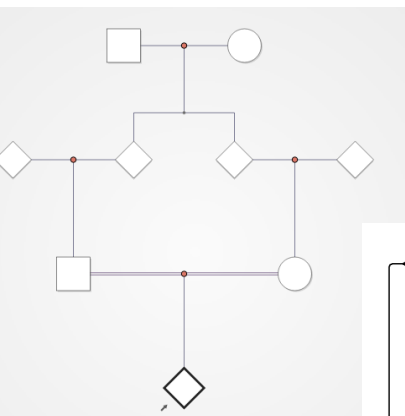
- Multi-Disciplinary teams
- Patient recruitment with complete phenotyping information
- Paid per patient recruited – but only if full information provided
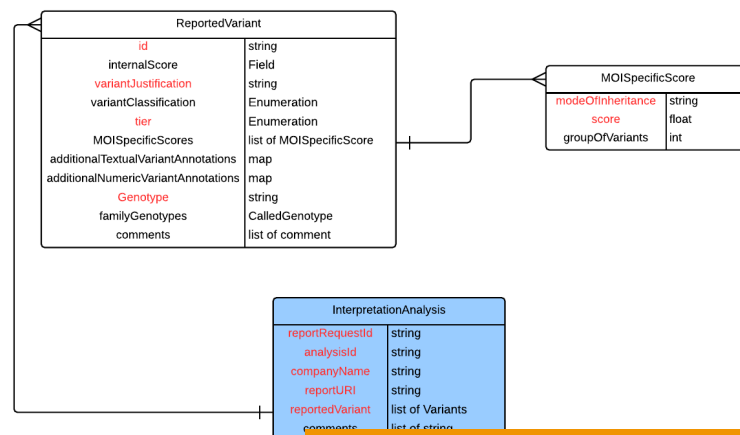- http://compbio.charite.de/hpoweb/showterm?id=HP:0002131
- Reporting to patients
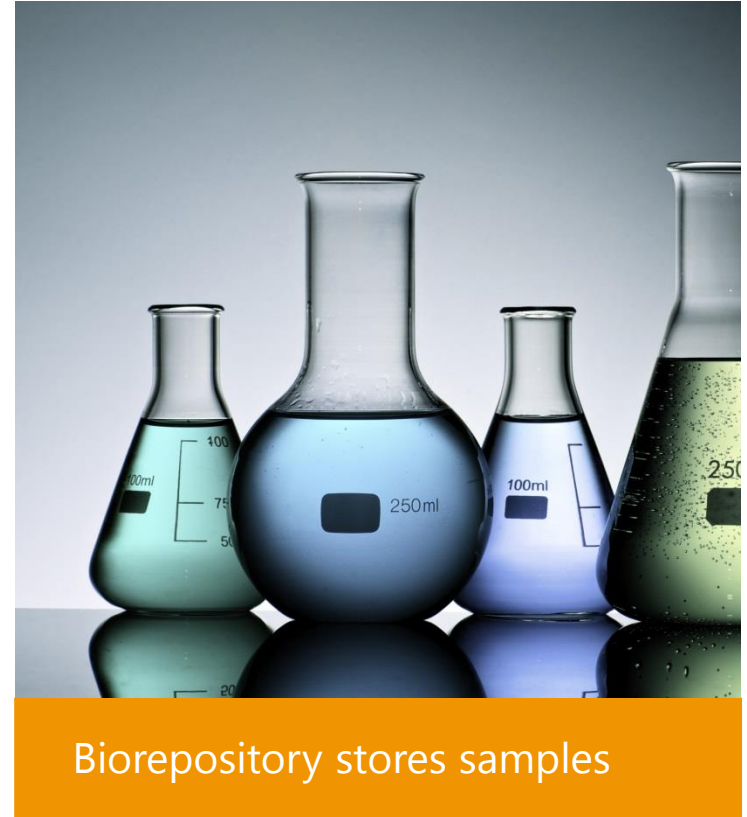
OpenClinica is used to collect patient information

CURATED CLINICAL PHENOTYPES & HEALTH DATA

Clinical Data

Family Pedigree

Interpreted Genome

# Bio-Repository

- Tasked with sample collection and QC for sequencing
- Act as a barrier for poor quality DNA
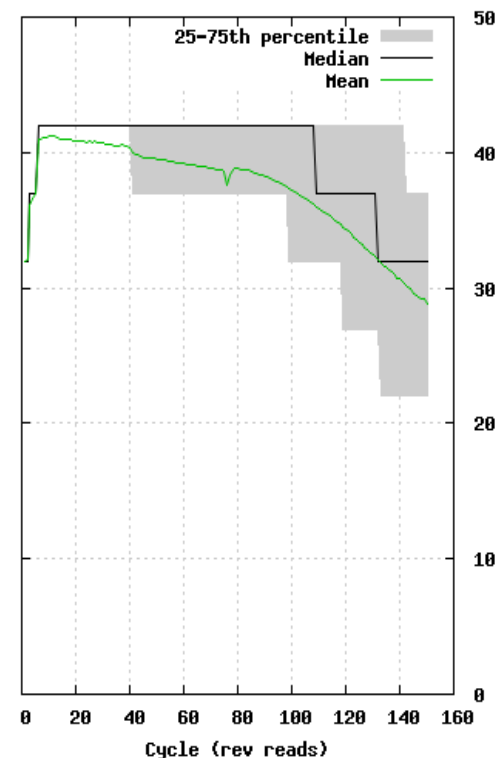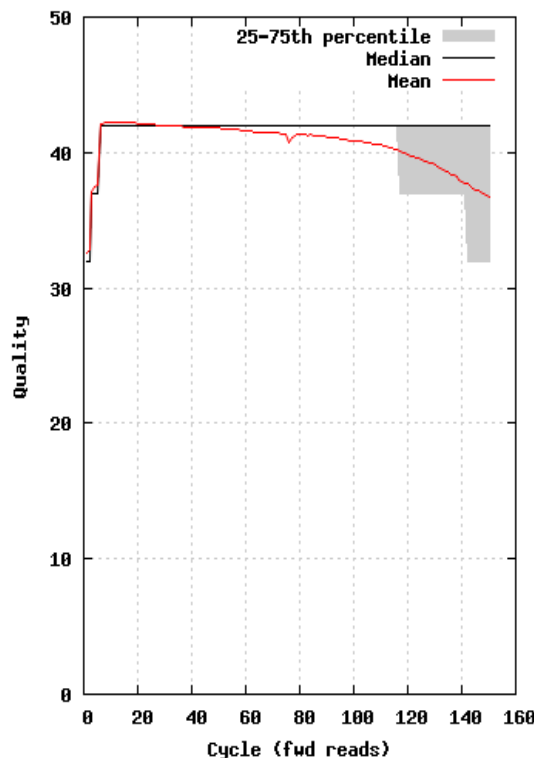- Standardised tubes, labels etc ready for receipt by illumina



Biorepository stores samples

# Illumina

- Chosen as WGS sequencing partner for Genomics England and subsequently the NHS
- Provide Sequencing
- Provide Primary Analysis Pipeline: Codename NORTHSTAR
  - Mapping
  - Variant Calling
- Data delivered to GEL data centre through fast connections



**Bridget Ogilvie Building**

- As a minimum for a '30X' germline genome:

    - 85x$10^9$ good bases
    - **>95% of autosomal genome covered at ≥15X with "good" bases**

- Current average stats:
    - 97x$10^9$ bases
    - **>97.3% of autosomal genome at ≥15X**

- Tumor sequenced at '75X'

## "We're Gonna Need a Bigger Truck"

- 200 Genomes a DAY!!
- Logistically this is a massive challenge

Illumina Sends Data to GEL


BERTHA Manages Pipeline

- Illumina sends genomes and associated files like:
  - Variant calls
  - QC – coverage etc
- Checked for integrity and to ensure illumina are upholding the agreement of the contract – certain number of bases at a certain quality
- Once happy GEL "accepts" the delivery

GEL Storage



Data Processed in 1 day = 20 Petabytes

- Rare Disease:
  - Each Genome: 100Gb
  - Trio is preferred so 300Gb per participant
- Cancer
  - Germline: 100Gb
  - Tumour: 200Gb
  - 300Gb per patient
- 10,000,000Gb = 10 Petabytes

# Databases

- Along with raw data files from illumina – variants are stored in a database called OpenCGA
- Allows annotation (adding information to variants like presence in public datasets)
- Allows fast internal frequency information
- Better for querying and producing reports



OpenCB – a suite of tools for genomic data

# Catalog

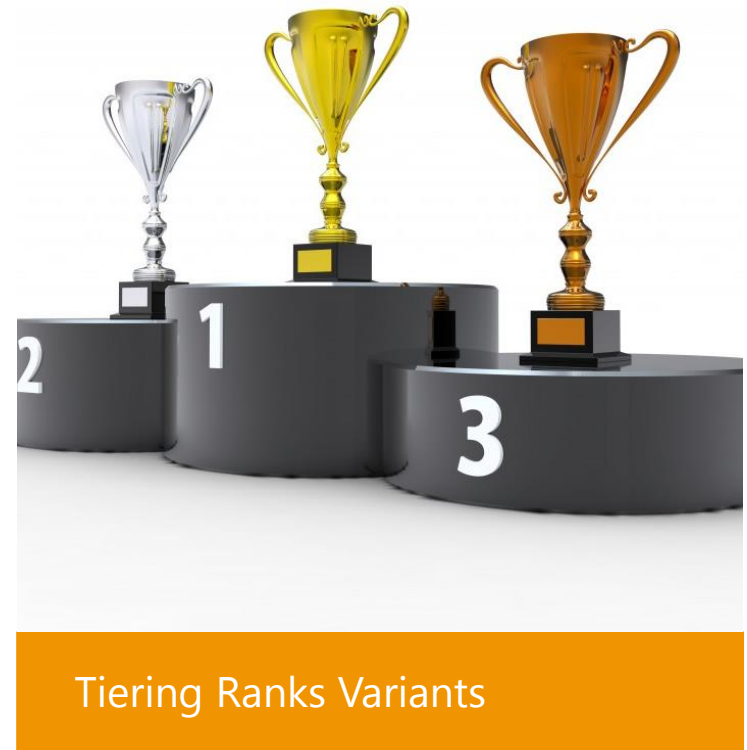| | |
|---|---|
| ▼ ▣ (2) 57156 | { 22 fields } |
| ⊞ _id | 57156 |
| ⊞ id | 57156 |
| " " name | LP2000755-DNA_B04.bam |
| " " type | FILE |
| " " format | BAM |
| " " bioformat | ALIGNMENT |
| null uri | null |
| " " path | by_date/2015-06-25/0000526305/CancerLP |
| " " ownerId | arendon |
| " " creationDate | 20150726125651 |
| " " modificationDate | 20150625111017 |
| " " description | |
| " " status | READY |
| ⊞ diskUsage | 204013908708 |
| ⊞ experimentId | -1 |
| ▼ ▣ sampleIds | Array [1] |
| ⊞ 0 | 56904 |
| ⊞ jobId | -1 |
| ▶ ▣ acl | Array [0] |
| null index | null |
| ▼ ▣ stats | { 6 fields } |
| ▶ ▣ BAM_HEADER_MACHINE | Array [3] |
| ▶ ▣ BAM_HEADER_OTHER | { 2 fields } |
| ▼ ▣ SAMTOOLS_STATS_ALL | { 31 fields } |
| ⊞ SAMTOOLS_READS_UNMAPPED | 105430116 |
| ⊞ SAMTOOLS_FILTERED_SEQUENCES | 0 |
| " " SAMTOOLS_ERROR_RATE | 1.045282E_02 |
| ⊞ SAMTOOLS_1ST_FRAGMENTS | 1330261830 |
| ⊞ SAMTOOLS_MISMATCHES | 3796815961 |
| ⊞ SAMTOOLS_BASES_DUPLICATED | 52353539700 |
| ⊞ SAMTOOLS_AVERAGE_LENGTH | 150 |
| ⊞ SAMTOOLS_SEQUENCES | 2661143644 |

Information on quality is stored in the database along with where to find the files

- Key to a good clinical genome and high fidelity result: **Quality Control**
  - Delivery integrity
  - BAM & VCF validity
  - Coverage
  - Number of bases
  - Sex
  - Mendelian Errors
  - Inbreeding estimates
  - IBD estimation
  - Ancestry

Clinical Quality Genomes

- NHS were worried about overloading clinical labs with too much validation work
- Prioritisation allows rules to be passed to GMCs on what they must validate, report etc
- Tiering based on mode of inheritance, phenotype, frequency in populations, frequency in disease groups etc



Tiering Ranks Variants

- A number of external providers chosen after "bake-off"
- Provide interfaces so that GMCs can easily access and interpret the variants from the WGS data
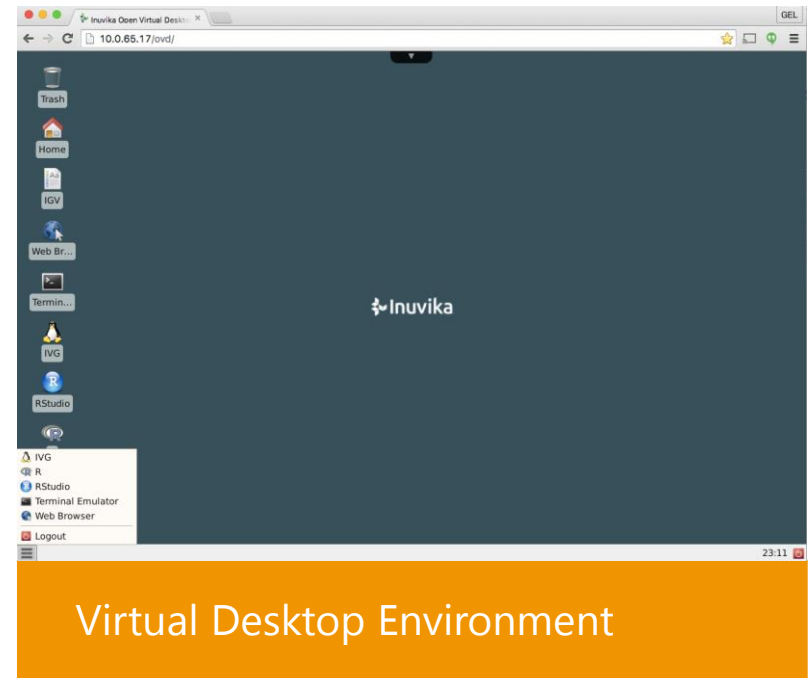


Congenica Interpretation Interface

- Genomics England & Illumina not ISO accredited to deliver clinical results
- Findings must be validated by an orthogonal method
- Reports must be written by clinical scientist



GEL/Illumina Not ISO accredited

- How can we access the data?
- Designed to protect patients
- Names & other identifiers not stored with WGS data
- All activity monitored & recorded
- **Diagnostics**
  - GMCs returning results to patients have full access - identifieable
  - Data kept in "reading library"
  - Electronic data access agreement
- **Researchers**
  - Non-identifiable data within "reading library"
  - Access Review Committee
  - GeCIPs
  - Electronic data access agreement
- **GENE – Pharma Collaborators**
  - Non-identifiable data within "reading library"
  - Electronic data access agreement
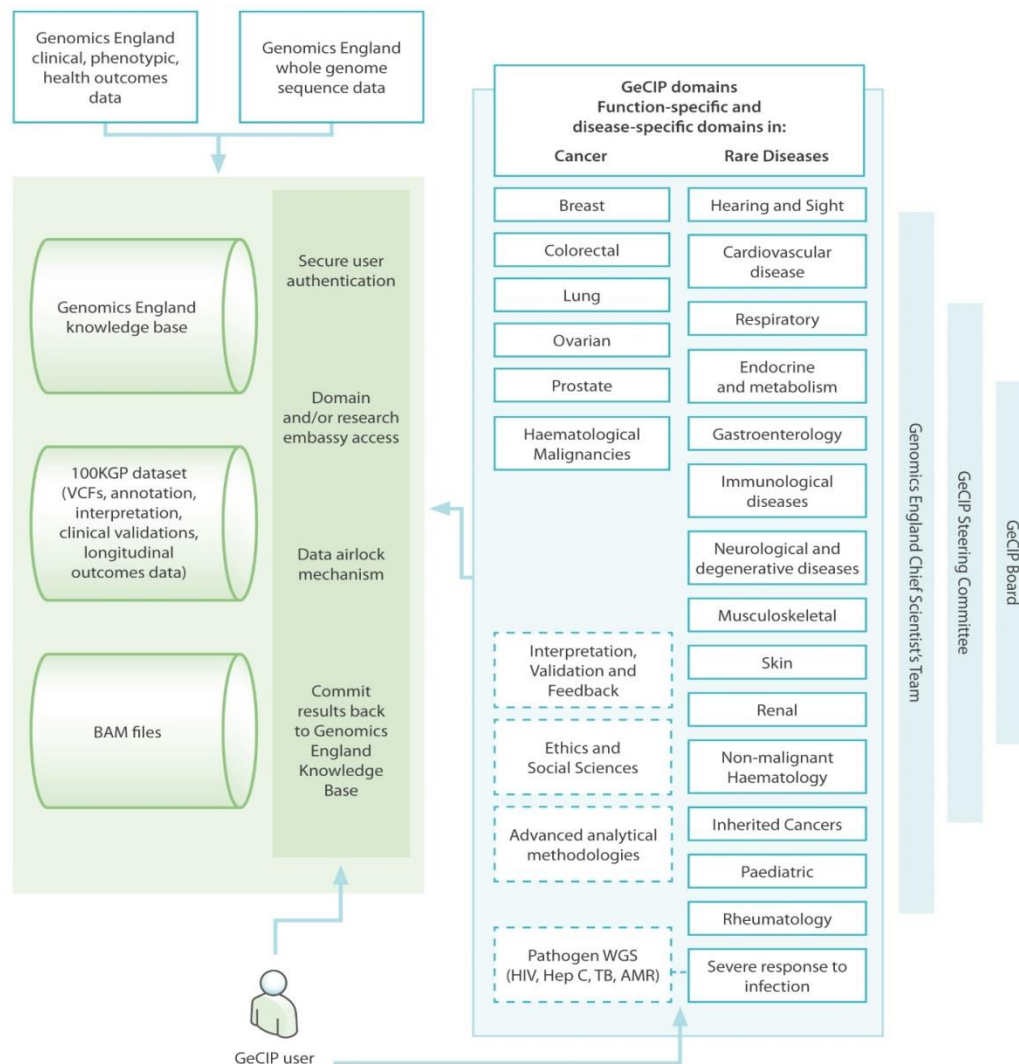  - Selected Pharma partners
  - To promote cures for rare diseases



Virtual Desktop Environment

# GeCIPs

Figure 1: The structure of GeCIP

- Ethical/Governance Issues with WGS?
- Logistical Issues?

Sheffield Children's **NHS**
NHS Foundation Trust

- Data is generated as a result of a long co-ordinated pipeline combining disparate organisations
- WGS + Health & Phenotype Data a valuable resource
- Data is stored at GEL and is only accessible through a virtual desktop
- But…. Important to remember WGS is not a magic bullet