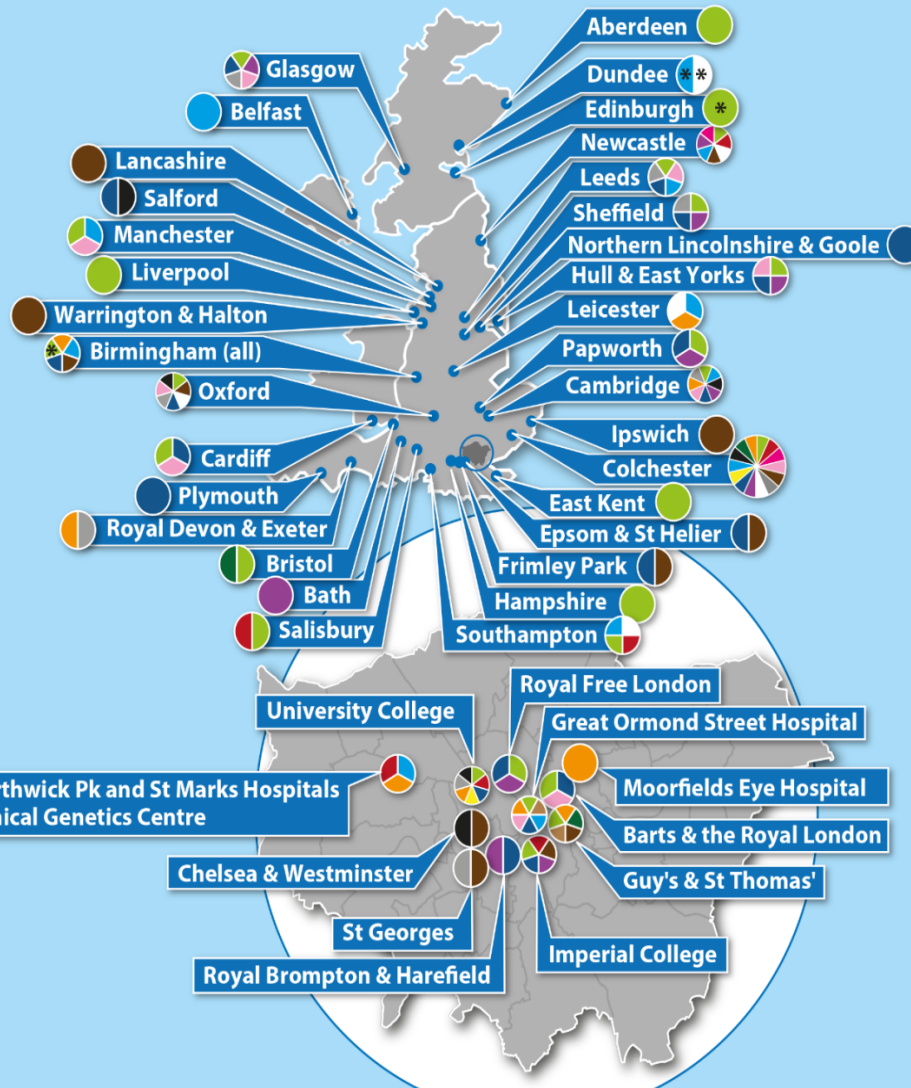# 10K genomes from rare disease cohorts

## Chris Penkett
## University of Cambridge

# The Rare Diseases Pilot
*started in 2013*



- **50** NHS Hospitals, **300** Clinical Care Teams

- PCR-free Whole Genome Sequencing in a **Clinically Accredited** Laboratory

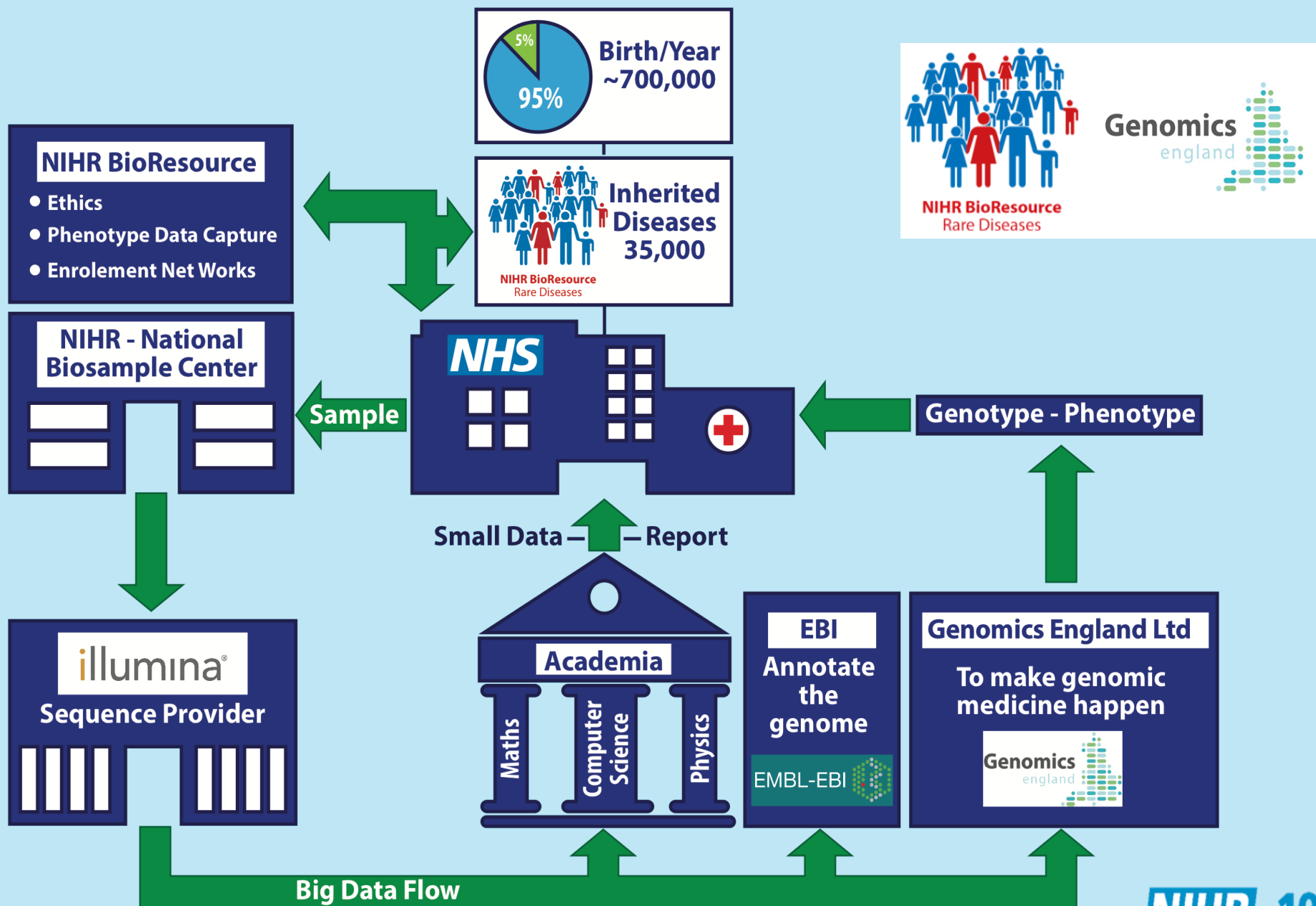- Clinical **Feedback** + Research

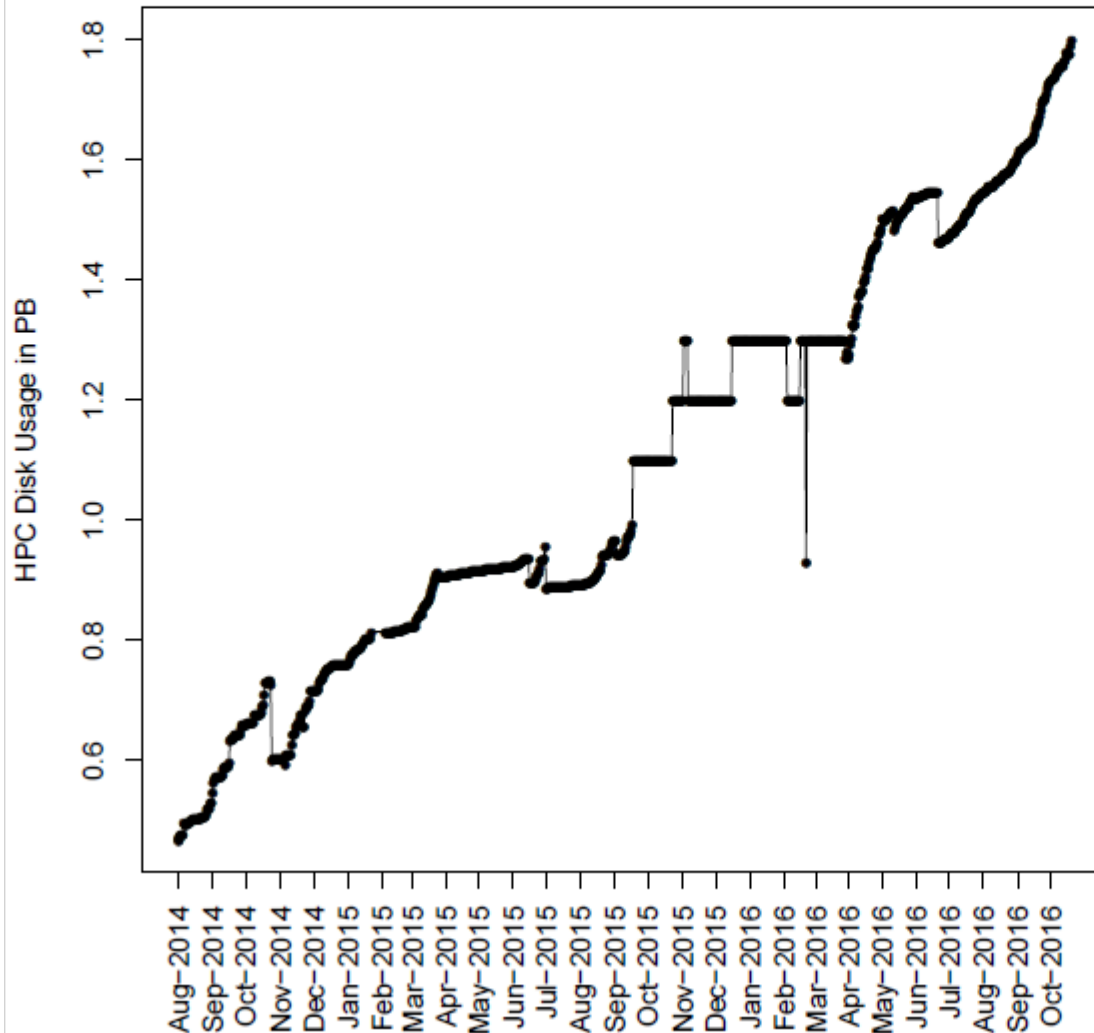# Main Rare Disease Cohorts

- 15 projects in total covering a variety of **Rare Diseases**

- **>10,000** Patients & Relatives

- Each project has a PI and (usually) a Chief Analyst to identify causal variants

|  | Projects | Targeted organ/state |
|---|---|---|
| **BPD** | Bleeding and platelets | Blood, coagulation etc. |
| **PAH** | Pulmonary Arterial Hypertension | Lung blood vessel |
| **PID** | Primary Immune Disorders | Immunity |
| **SPEED** | Retinal dystrophy + neurological disorders | Retinal, neurological |

# The NHS 100,000 Genome Project

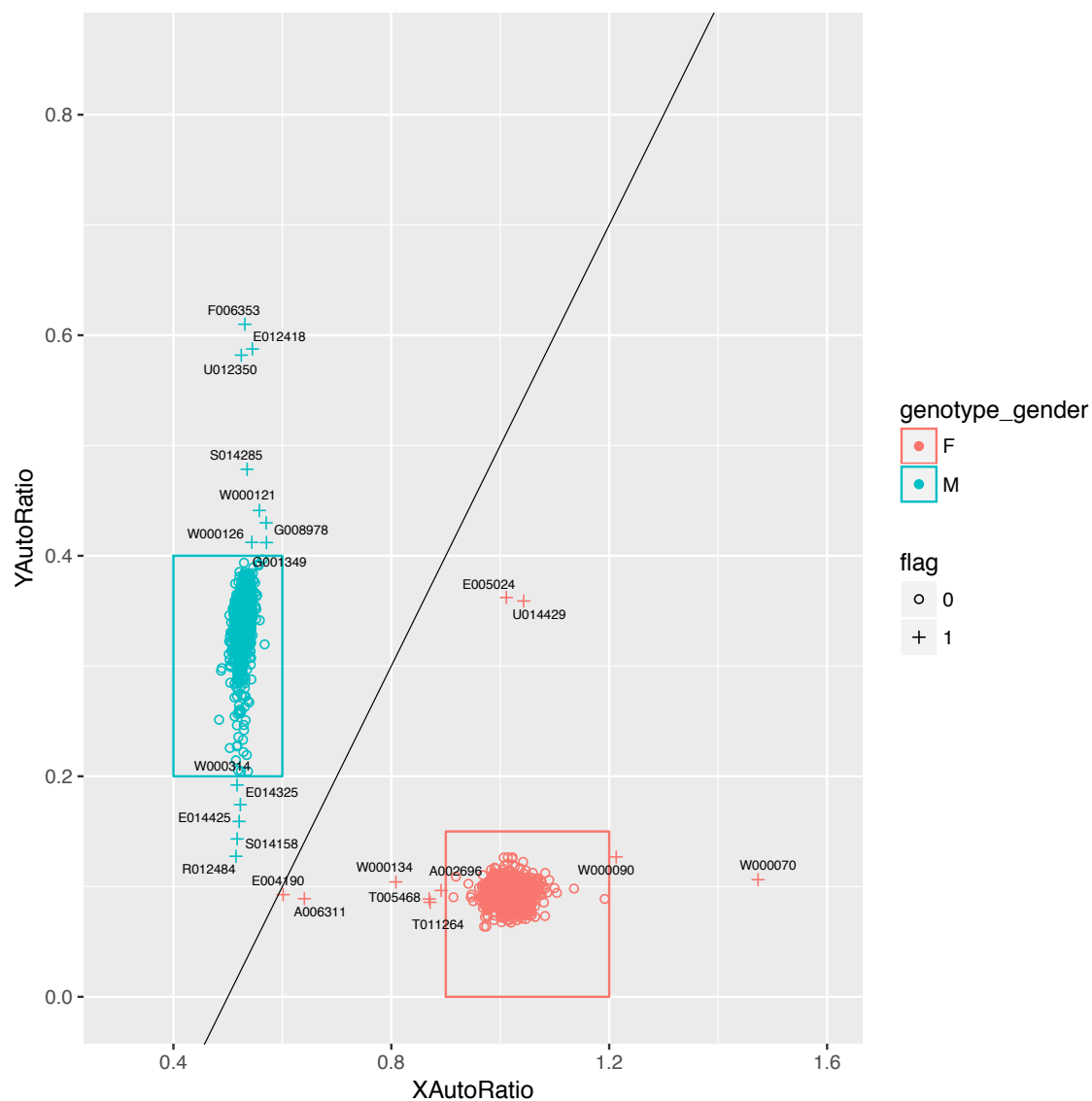# NIHR BioResource – The Data



**Big Data**

Now over 2 PB in > 3 years

~70GB raw data per genome

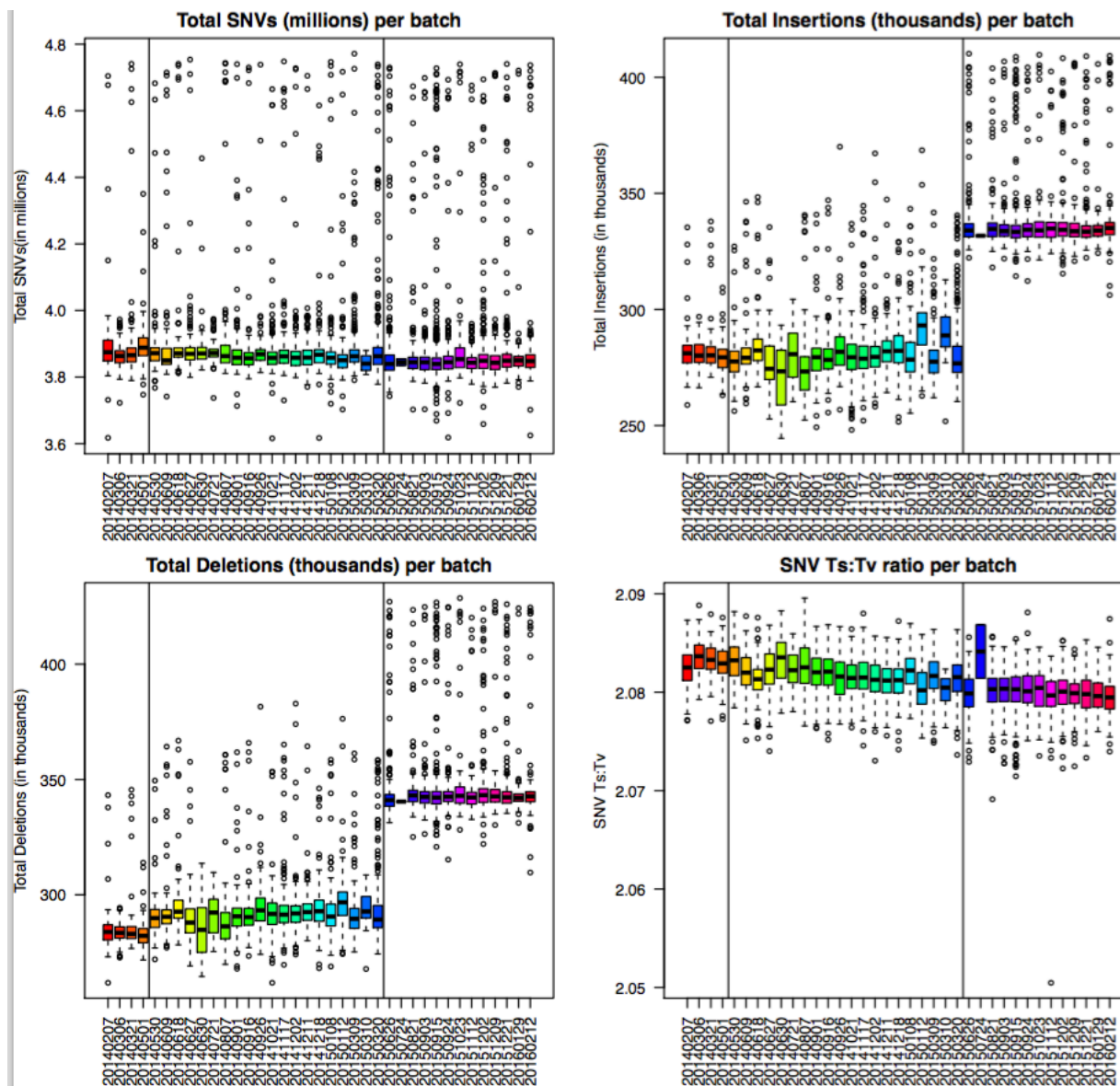~0.5GB "just the variants" per genome

# WGS Data QC

Data is checked/annotated for:

- Errors introduced during data transfer (md5 checksum)
- Create smaller cram files to archive at EGA (EBI) without raw data loss
- WGS quality metrics
- Genetic gender vs manifest gender
- Relatedness
    Check family structure as reported
    Identity unknown family recruitments
    Detect duplicates and/or identical twins
    Identify subset of unrelated individuals for unbiased
        allele frequency calculation
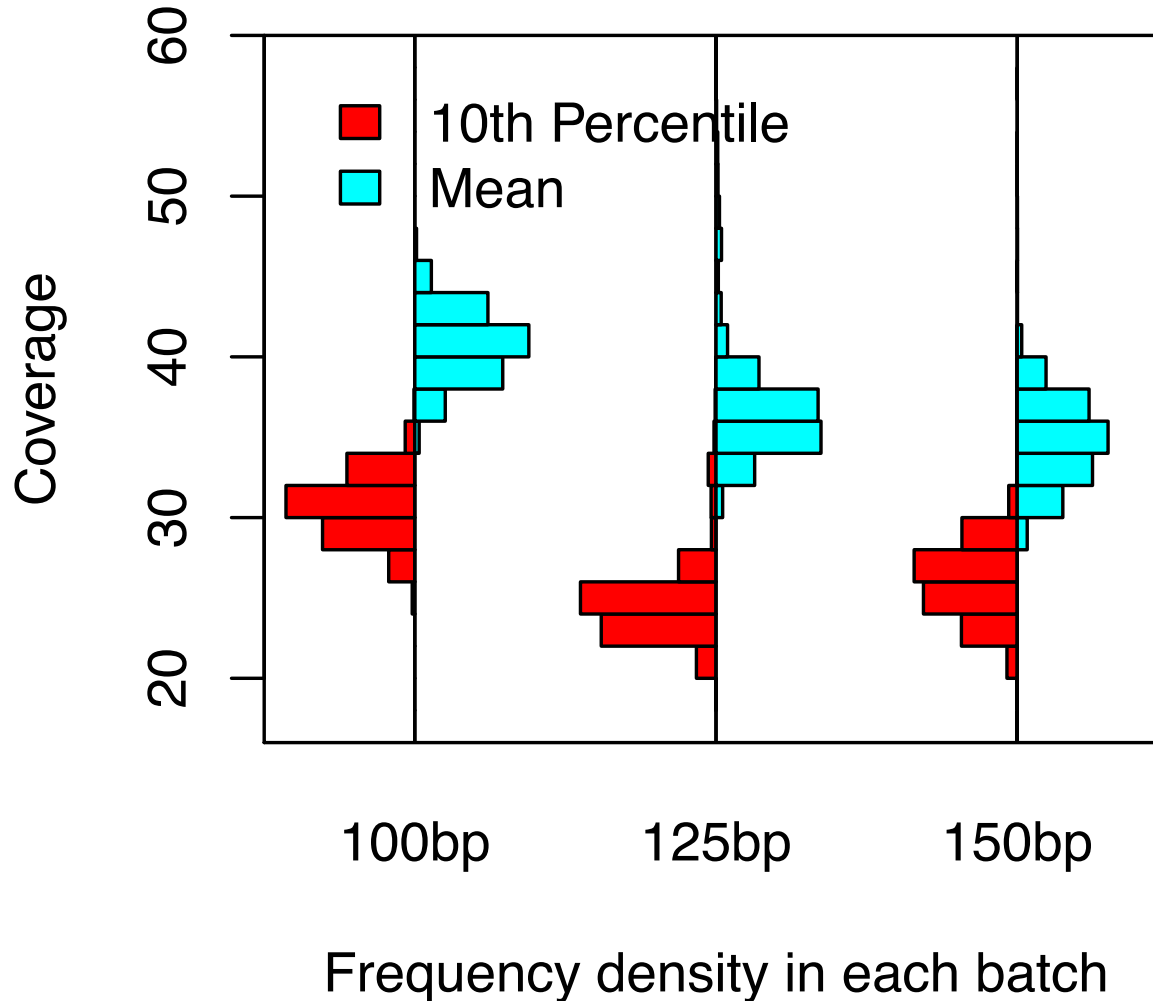- Ethnicity

# Genomic Gender from
# X/Y vs Autosome Coverage

# Summary of Data Parameters for 3 Read Lengths

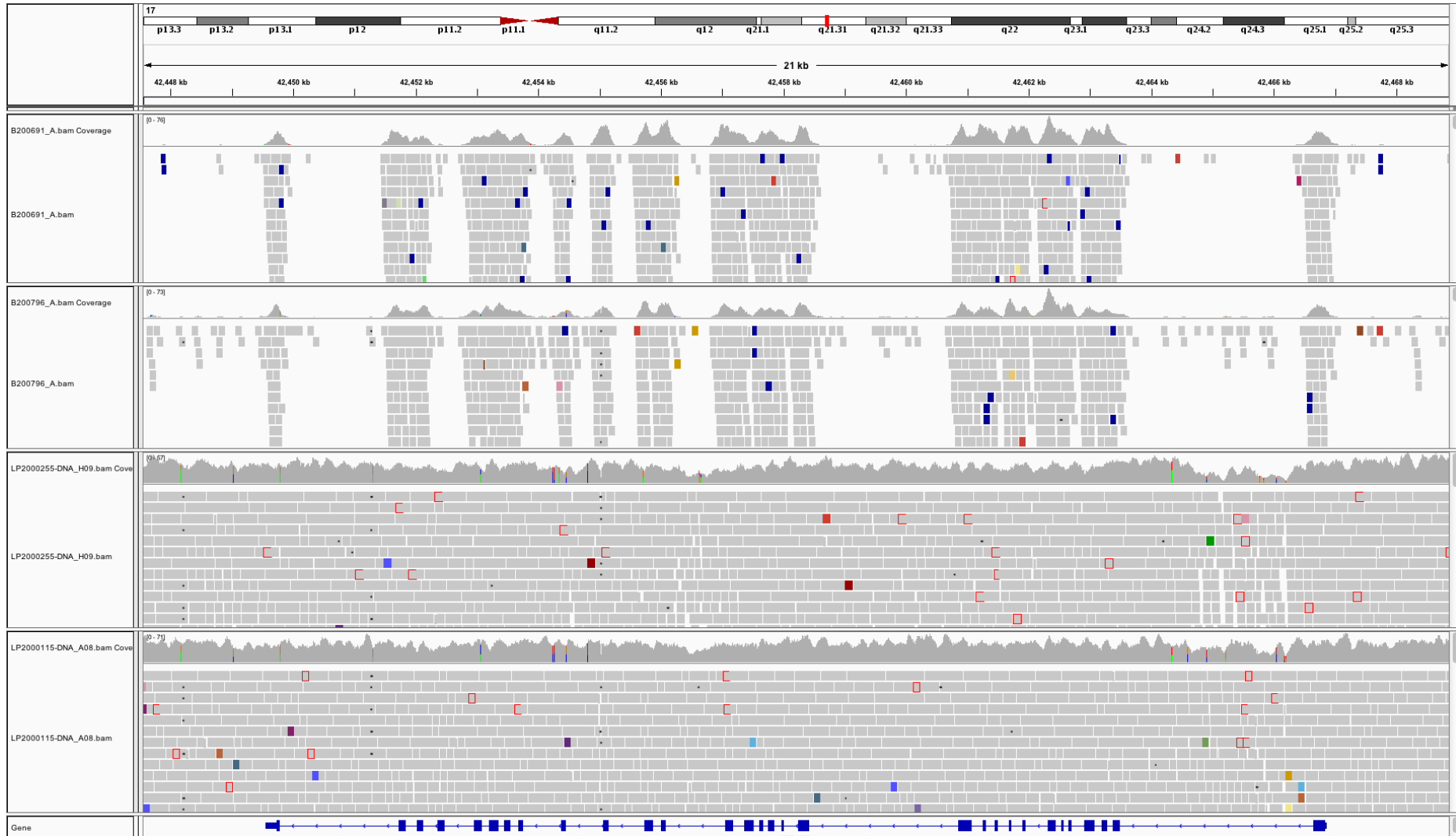# Overall Genome Coverage



Frequency density in each batch

**Coverage** is the average number of reads representing a given nucleotide in the reconstructed sequence
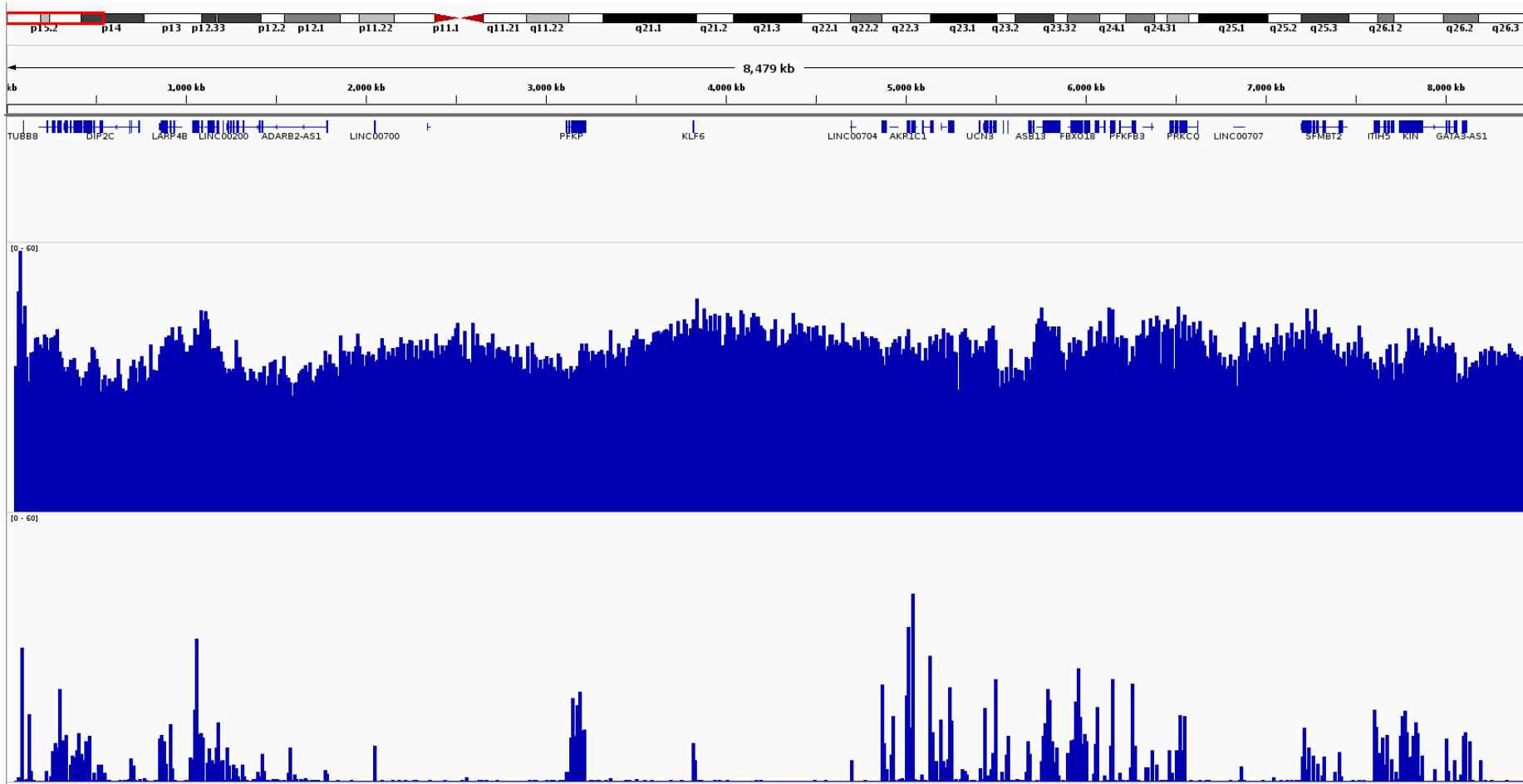
For detection of SNPs and rearrangements, publications recommend from 10× to 30× depth of coverage
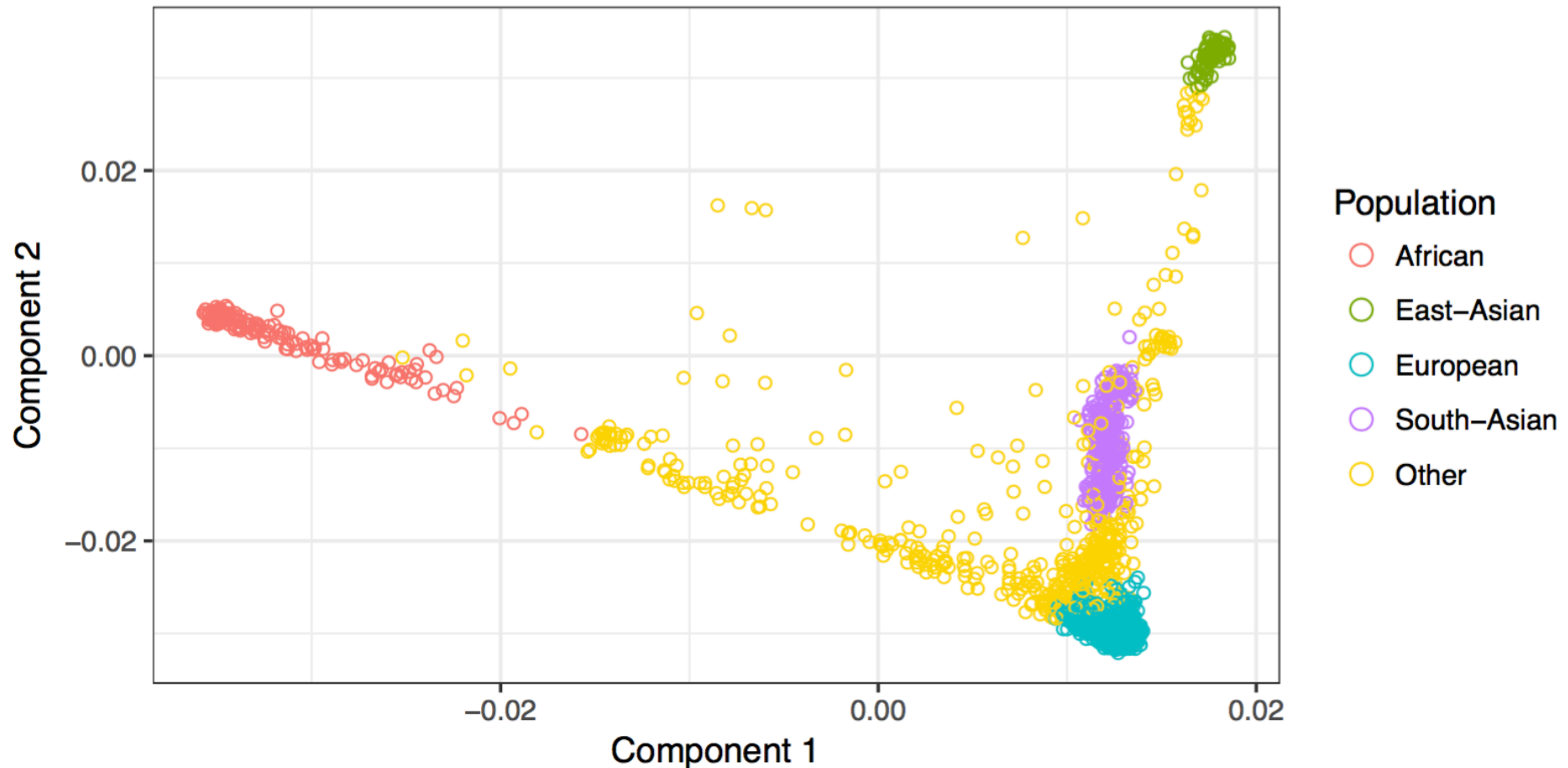
More than 90% of the genome has 25x coverage

# Genomes versus Exomes
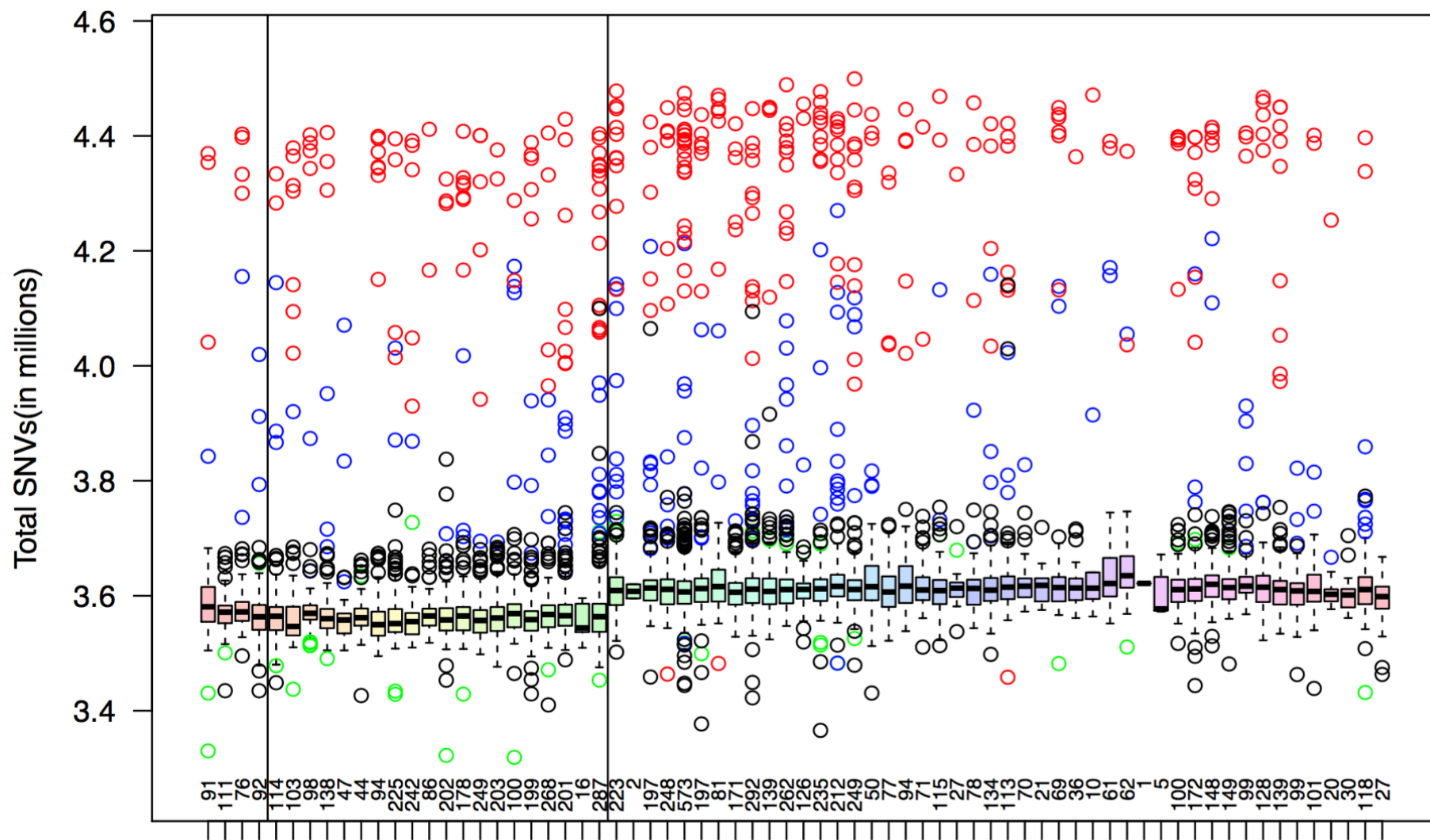
# Example WGS10K Coverage vs WES

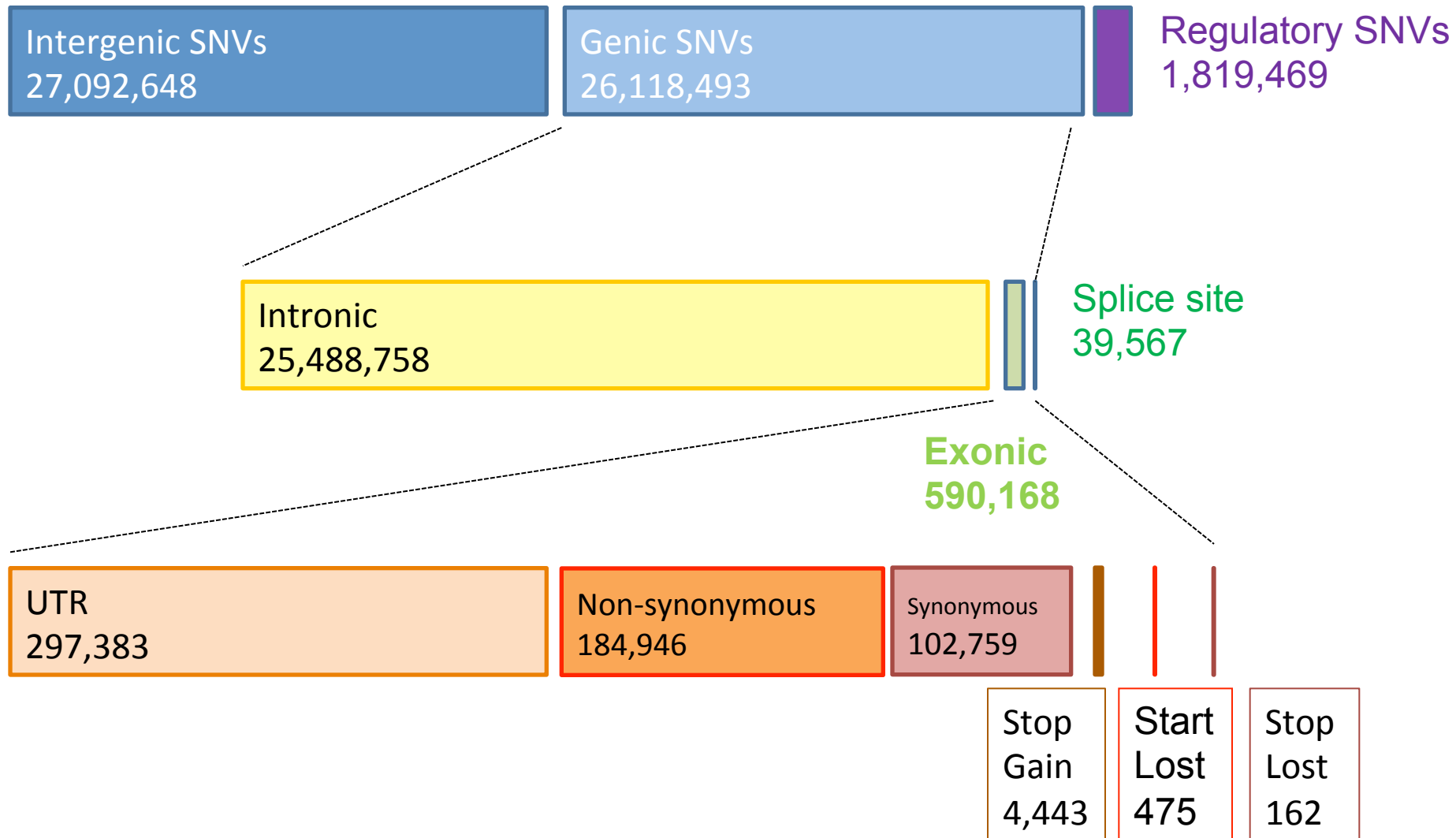# Genomic Ethnicity Determination from Principal Components Analysis

# Batch Data, SNVs per Person and Ethnicity



Average 3.9 Million SNVs per person
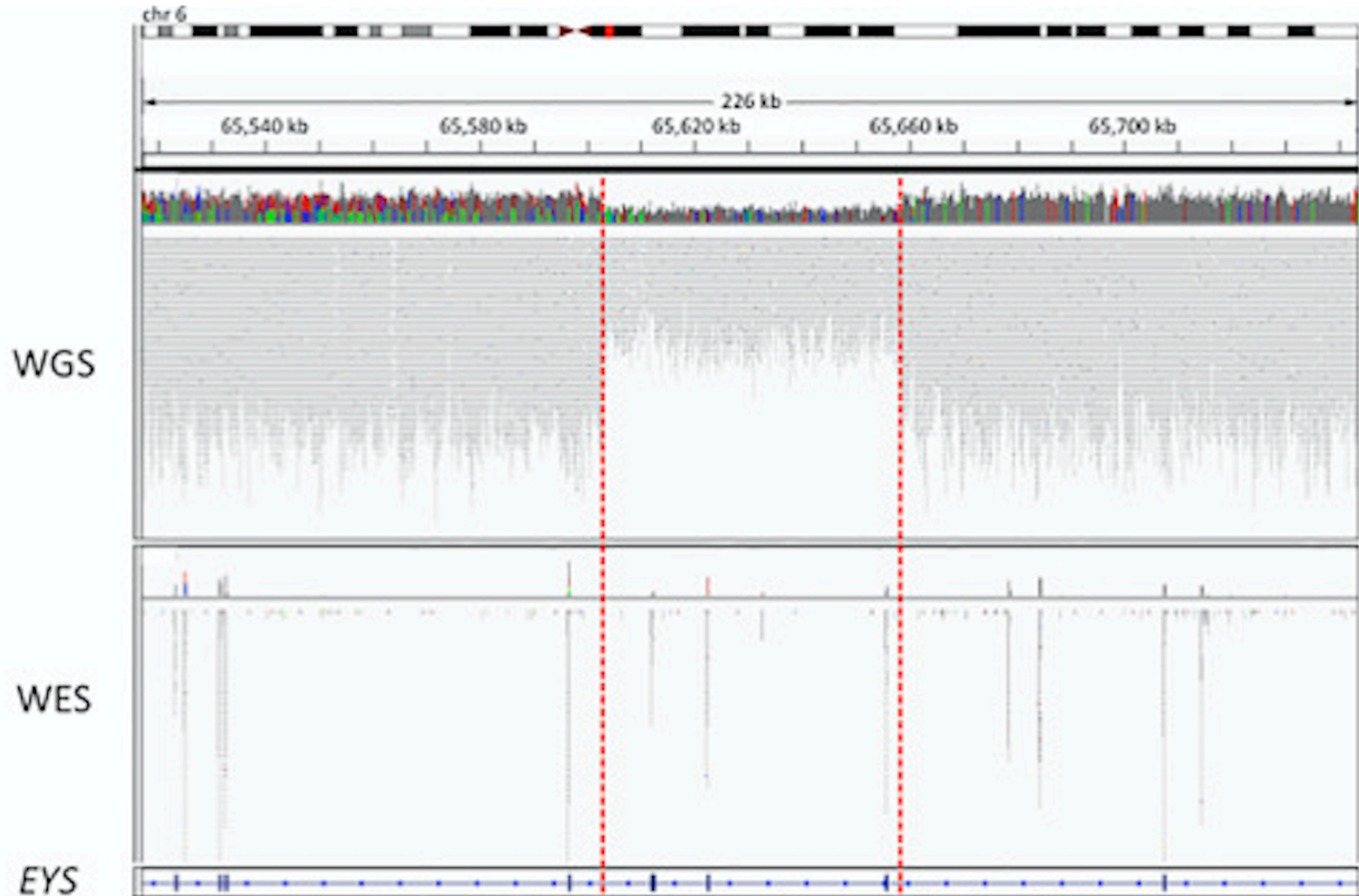Average 3.6 Million post QC filtering

# Example of CNV Detection for Genome and Exome Data
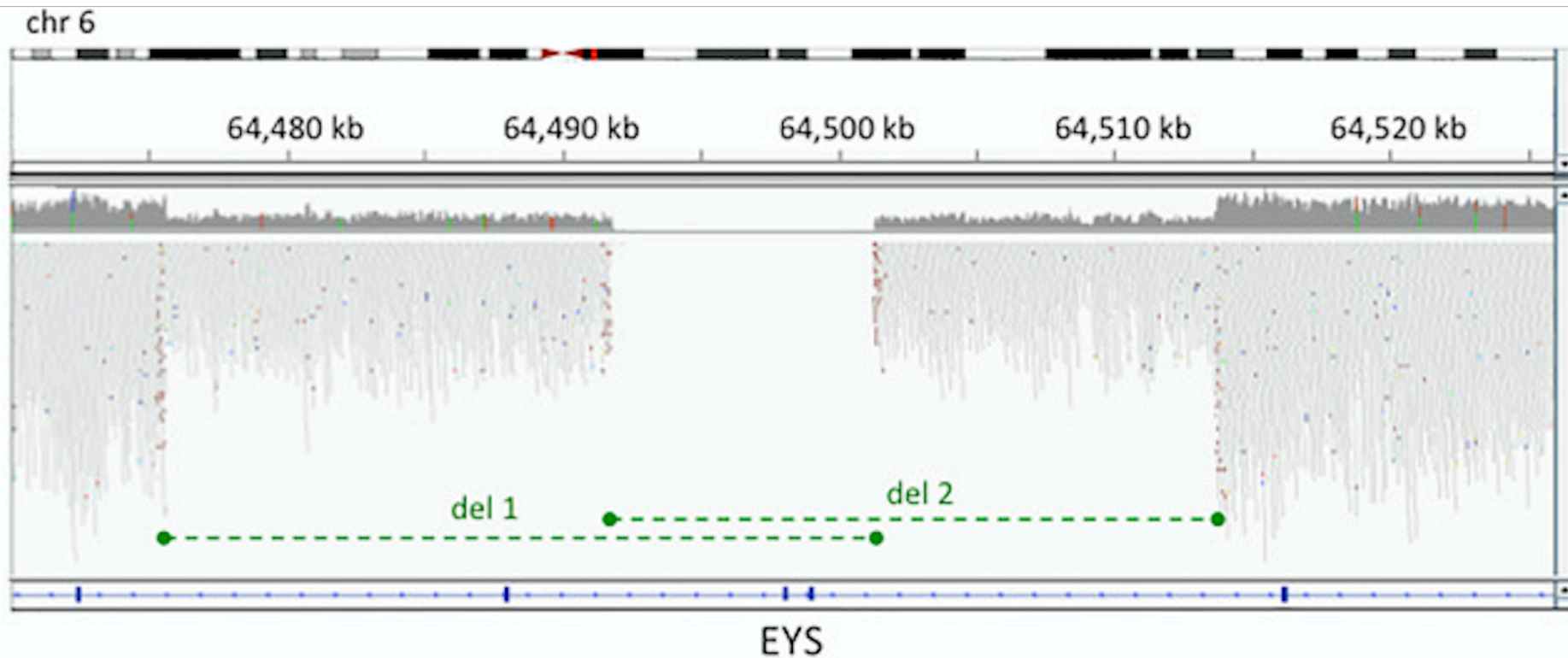
# Example of Two Large Overlapping Deletions in WGS Data

# MDTs

- Multi-Disciplinary Team meetings
  - 2 clinician/clinical geneticist, 1 chairing
  - Project coordinator
  - Pertinent Finding team member
  - Analyst

# Reporting Across the Cohorts

| Projects | Pre-screening | Cases reviewed at MDT | Diagnostic yield % |
|---|---|---|---|
| BPD – Bleeding and Platelet Disorders | Extensive | 1162 | 10 % |
| PAH – Pulmonary Arterial Hypertension | Moderate | 781 | 18 % |
| PID – Primary Immunodeficiency | Extensive | 720 | 9 % |
| SPEED RD – Retinal Dystrophy | Moderate | 284 | 63 % |
| SPEED neuro – Paediatric Neurodevelopmental | Limited | 243 | 32 % |

**Novel variants: ~50%**

# The Teams

Sri Deevi
Salih Tuna
Olga Shamardina
Fengyuan Hu
Kathy Stirrups
Stuart Meacham
Tony Attwood
Stefan Gräf
Matthias Haimel
Marta Bleda
Ernest Turro
Daniel Greene
Keren Carss
Alba Sanchis-Juan
Hana Lango Allen
Karyn Megy
Louise Daugherty
Tim Young
Roger James
Catherine Titterton
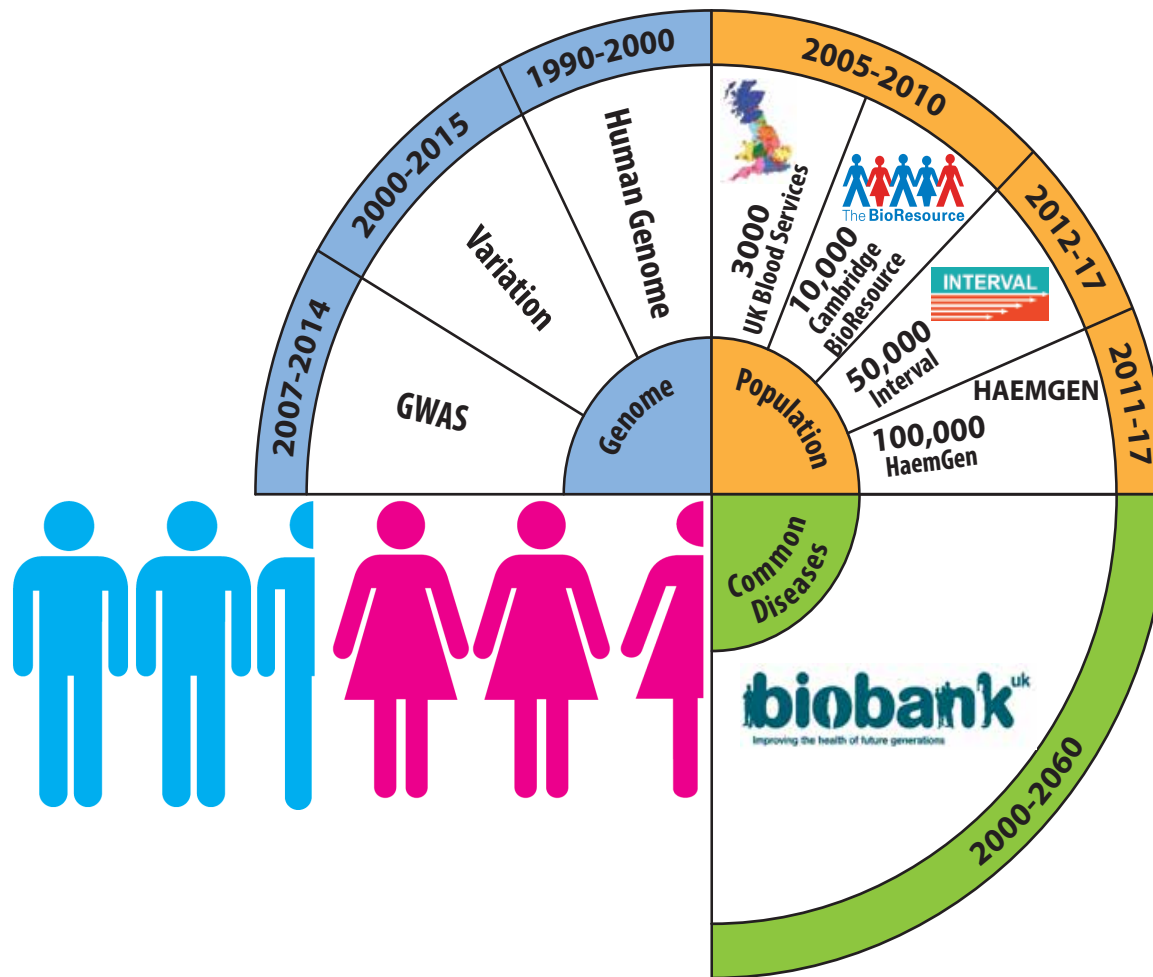Lucy Raymond
Willem Ouwehand

- CATGO
- NIHR-RD Enrolment team
    - Sofie Ashford
    - Sofia Papadia
- HPC
    - Stuart Rankin
    - Wojciech Turek
    - Paul Calleja
    - Nacho (Ignacio) Medina (GEL/EBI)
    - Jacobo Coll (GEL)
- BLUEPRINT team
- Admin team
- Illumina sequencing/bioinformatics teams
    - Russell Grocock
    - John Peden
    - Christian Bourne
    - Sean Humphray
    - Terry Gerighty
- GEL bioinformatics team
    - Augusto Rendon
    - Katherine Smith
- EGA team at EBI
    - Jeff Almeida-King

# NIHR BioResource - Rare Diseases

| Rare Disease/Condition | Acronym/ Approved in month/year | Lead Investigator | Capacity allocated | Samples sent for sequencing | WGS10K Samples received at the HPC |
|---|---|---|---|---|---|
| Bleeding and Platelet Disorders | BPD 12/12 | Prof Willem Ouwehand | 1250 | 839 (127 WES) | 759 |
| Cerebral small vessel | CSVD 04/14 | Prof Hugh Markus | 250 | 125 | 110 |
| Ehler-Danlos Syndrome | EDS 07/13 | Prof Tim Aitman | 400 | 90* WES | 0 |
| Genomics England pilot | GEL Pilot 11/13 | Prof Mark Caulfield | 2000 (+ 3000) | 4092 (10 dups) | 2000 |
| Hypertrophic Cardiomyopathy | HCM 05/12 | Prof Hugh Watkins | 300 | 146 | 134 |
| Intrahepatic Cholestasis of Pregnancy | ICP 01/14 | Prof Catherine Williams | 270 | 90 | 76 |
| Multiple Primary Malignant Tumours | MPMT 11/13 | Prof Eamonn Maher | 700 | 297 | 263 |
| Neuropathic Pain Disorders | NPD 10/14 | Prof Geoff Woods | 250 | 41 | 39 |
| Primary Immune Disorders | PID 12/12 | Prof Ken Smith | 1250 | 1119 (26 WES) | 1080 |
| Primary Membranoproliferative Glomerulonephritis | PMG 05/13 | Dr Danny Gale | 213 | 97 | 97 |
| Pulmonary Arterial Hypertension | PAH 12/12 | Prof Nick Morrell | 1250 | 789 | 751 |
| Specialist Pathology | SPEED 12/12 | Prof Lucy Raymond | 1250 | 1065 (188 WES) | 1043 |
| Stem Cell and Myeloid Disorders | SMD 10/14 | Prof Irene Roberts | 600 | 94 | 63 |
| Steroid Resistant Nephrotic Syndrome | SRNS 03/13 | Dr Ania Koziell | 250 | 93 | 84 |
| Leber Resistant Nephrotic Syndrome | LHON | Prof Patrick Chinnery | 70 | 0 | 0 |

# UK Biobank – 0.5 Million Genotyped Volunteers
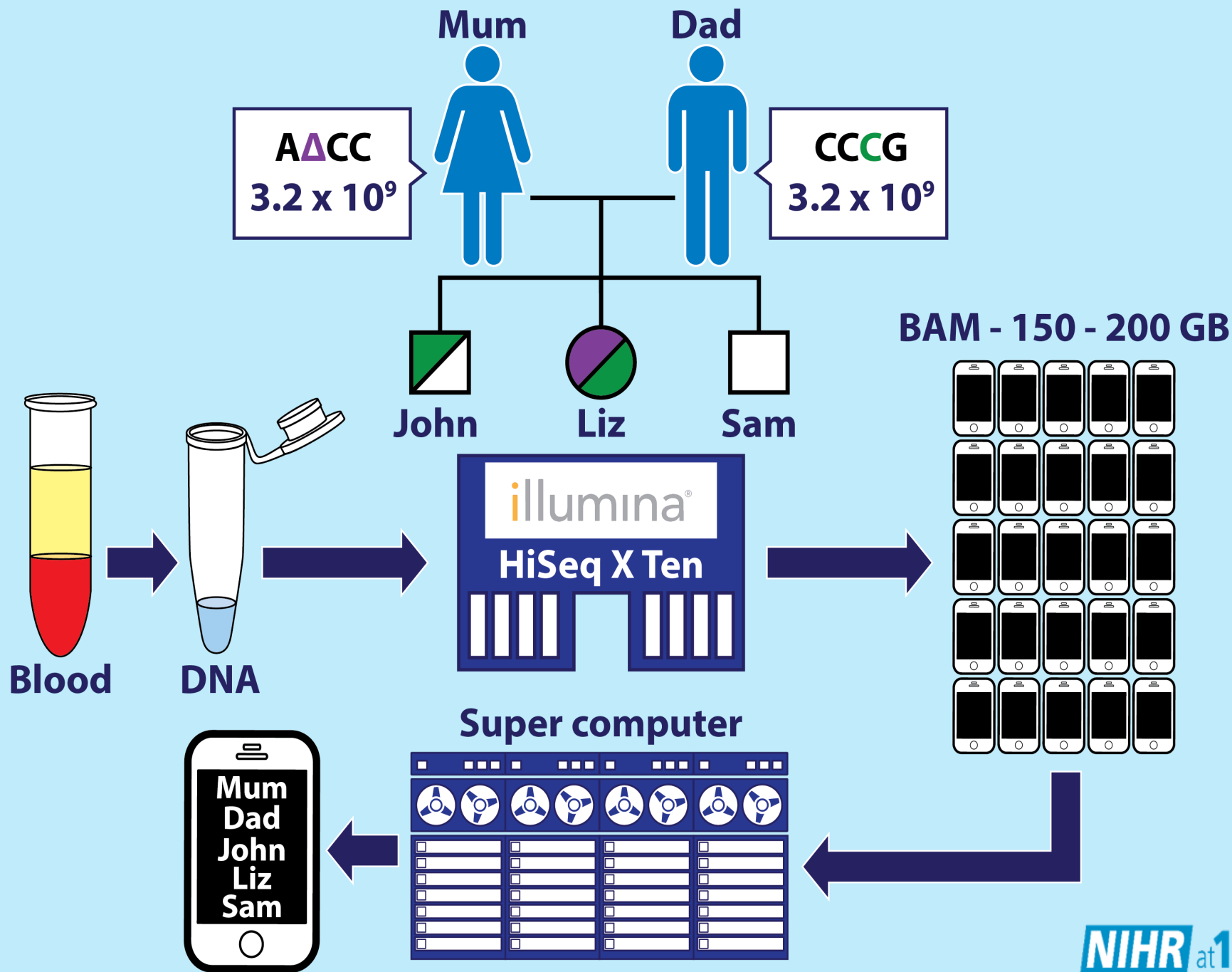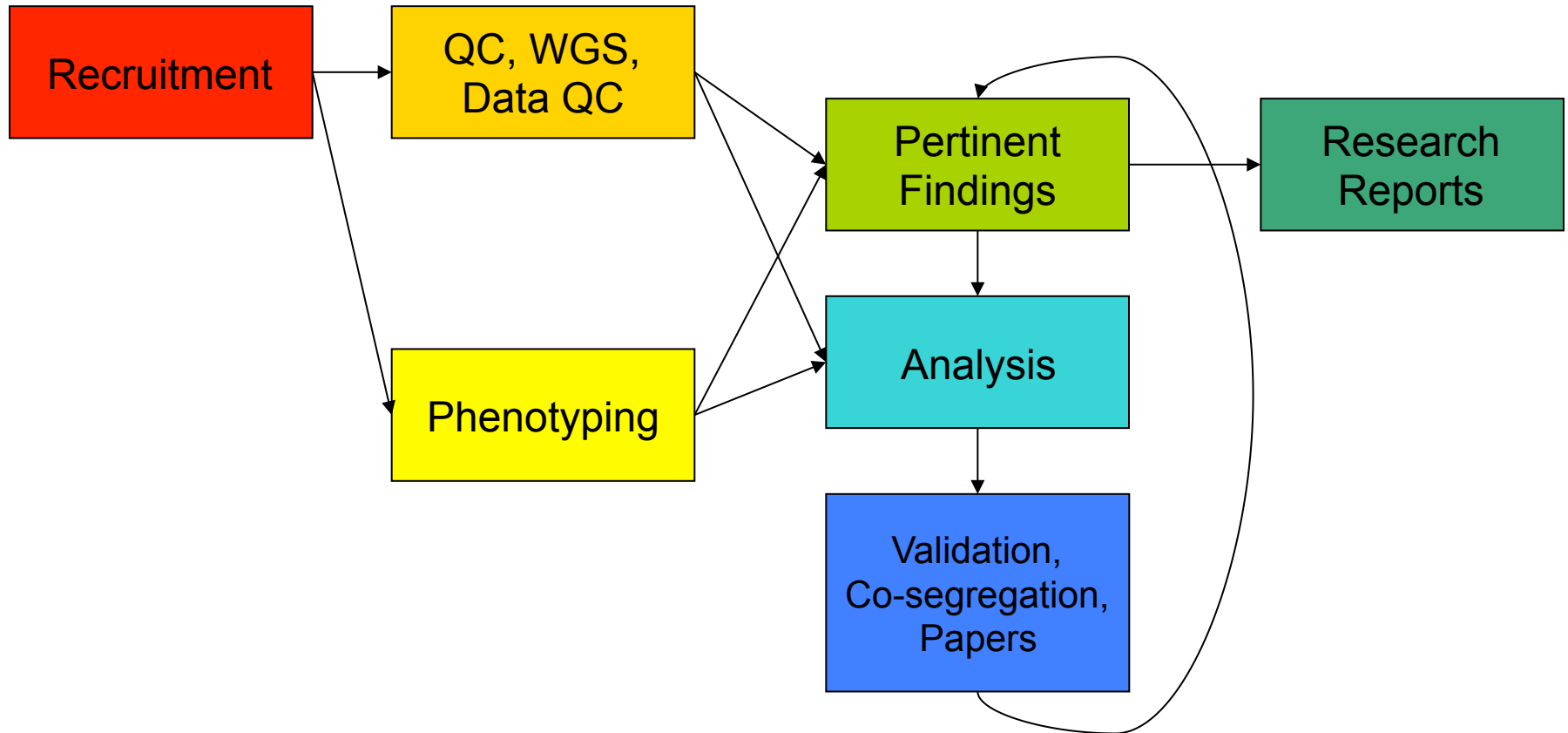


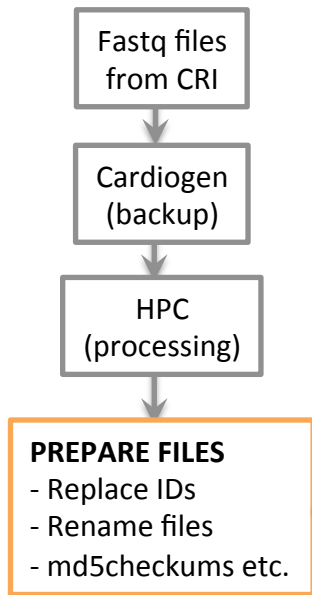= 100,000 males       = 100,000 females

Enrolment across the UK from 2006-10
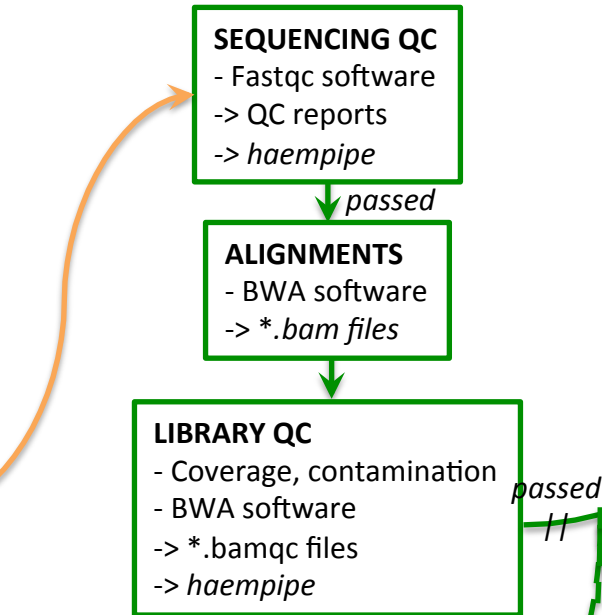Age at enrolment between 40-69 yrs
Linkage to NHS GP and Hospital records
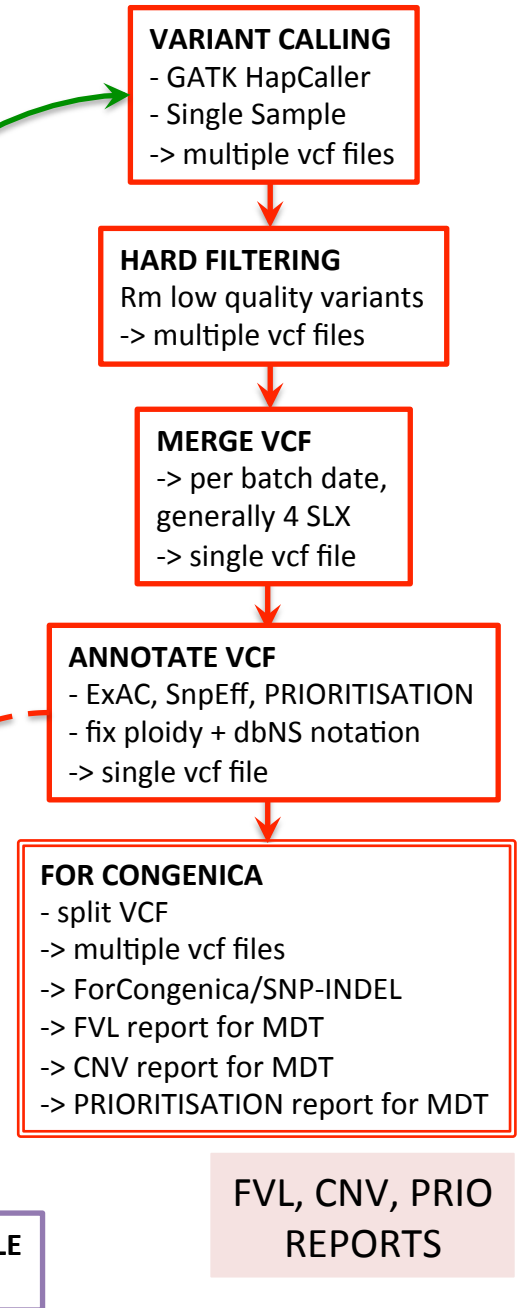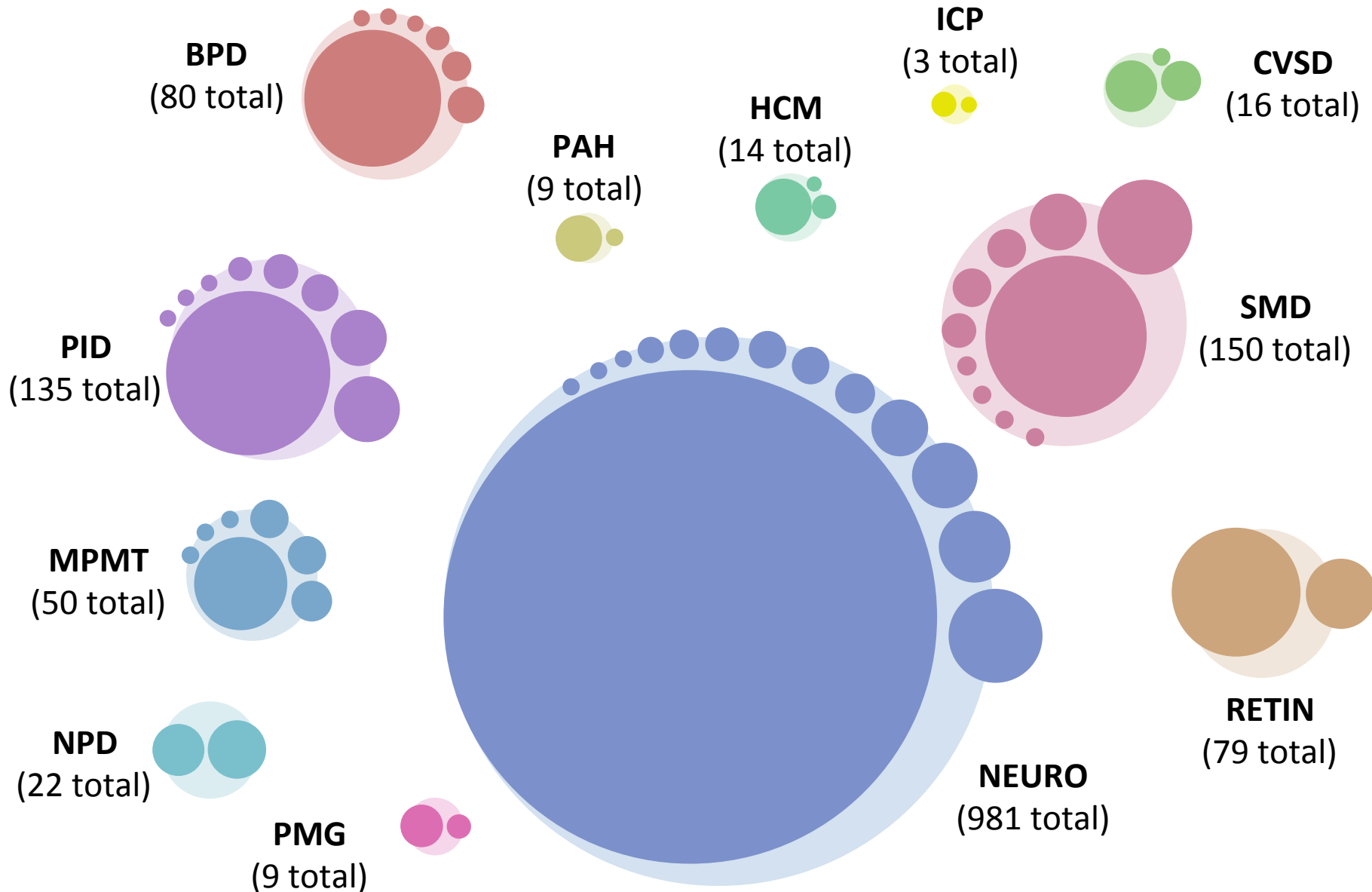
# NIHR BioResource – The Process

**Pre-pipeline**

Fastq files from CRI

Cardiogen (backup)

HPC (processing)

**PREPARE FILES**
- Replace IDs
- Rename files
- md5checkums etc.

**PHASE I**

**SEQUENCING QC**
- Fastqc software
-> QC reports
-> *haempipe*

*passed*

**ALIGNMENTS**
- BWA software
-> *.bam files*

**LIBRARY QC**
- Coverage, contamination
- BWA software
-> *.bamqc files
-> *haempipe*

*passed*
//

**COV. PLOT REPORT**

**Coverage plots**

**CNV REPORT**

**CNV**
- ExomeDepth
-> *.vcf files*

**ETHNICITY GENDER**

**SAMPLE SUMMARY FILE**
(relatedness)

**PHASE II**

**VARIANT CALLING**
- GATK HapCaller
- Single Sample
-> multiple vcf files

**HARD FILTERING**
Rm low quality variants
-> multiple vcf files

**MERGE VCF**
-> per batch date, generally 4 SLX
-> single vcf file

**ANNOTATE VCF**
- ExAC, SnpEff, PRIORITISATION
- fix ploidy + dbNS notation
-> single vcf file

**FOR CONGENICA**
- split VCF
-> multiple vcf files
-> ForCongenica/SNP-INDEL
-> FVL report for MDT
-> CNV report for MDT
-> PRIORITISATION report for MDT

**FVL, CNV, PRIO REPORTS**

**Bioinformatics Pipeline**
03-2017 – v1

# 1,400 genes for clinical reporting



BPD
(80 total)

ICP
(3 total)

CVSD
(16 total)

HCM
(14 total)

PAH
(9 total)

SMD
(150 total)

PID
(135 total)

MPMT
(50 total)

NEURO
(981 total)

RETIN
(79 total)

NPD
(22 total)

PMG
(9 total)

# Research report