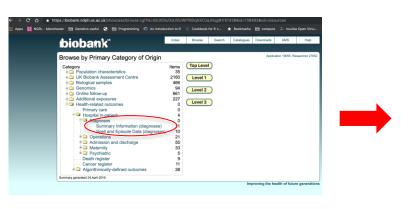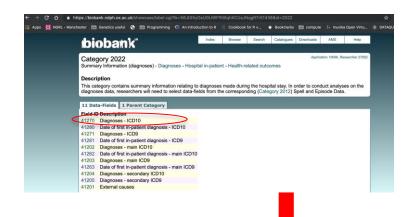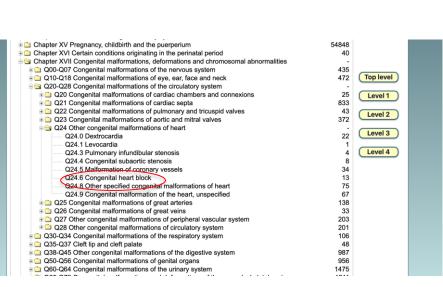# UK Biobank – UoM

## Getting started with hospital episode statistics (HES) data
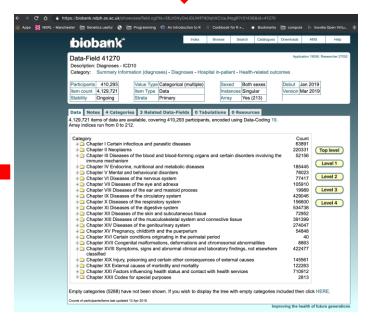
27/06/19

# HES – what diagnoses are present?

# HES – where to access the data once project has been approved

## HES tables

**HESIN**
- main primary hospital diagnoses
- episode/admission dates
- primary ICD9/ICD10/OPCS4 codes

**HESIN_DIAG10**
- secondary ICD10 diagnoses

**HESIN_OPER**
- secondary OPCS4 codes

**HESIN_DIAG9**
- secondary ICD9 codes

Data extracted by SQL queries:

- Query specific code e.g. SELECT eid FROM hesin WHERE diag_icd10 = 'Q256'

- Get everything and export e.g. SELECT * FROM hesin

# Rscript to open HES tables and find samples with matching diagnosis – 'Heart block (Q246)'

```
library(data.table)

################################################################################
# load the HES data tables you've exported from UKB SQL page
################################################################################

#read hesin table
hesin=fread("HESIN.tsv")

#read hesin_diag10 table (secondary diagnosis table)
hesin_diag10=fread("HESIN_DIAG10.tsv")

#read hesin_oper table (secondary operations table)
hesin_oper=fread("HESIN_OPER.tsv")


################################################################################
# extract samples that might be of interest based on ICD10 codes
################################################################################

#read list of ICD codes we're interested in – here jut a file with 'Q246' under the 'ICD10' header.
codes<-read.csv("heart_block_phenotype_codes", header=TRUE)

#subset these samples from HESIN table
hit1<-subset(hesin, hesin$diag_icd10 %in% codes$ICD10)

#then HESIN_DIAG10 table - this table contains the secondary diagnosis (only ICD10)
hit2<-subset(hesin_diag10, hesin_diag10$diag_icd10 %in% codes$ICD10)

#you might also search for related ICD9/operation codes in the same way

#extract the eids
eid<-c(hit1$eid,hit2$eid)

#get unique eids
heart_block_samples<-unique(eid)
```
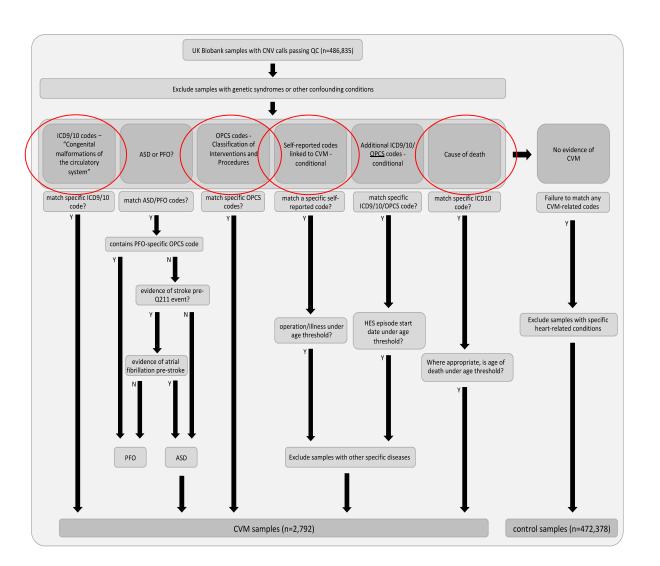
# Getting started genetic data

# What is available?

**Arrays**

- 488,766 individuals
- 820,967 SNP and indel markers included
- 2 arrays
  - Affymetrix Axiom UK BiLEVE array (~50,000)
  - Affymetrix Axiom UK Biobank array (~450,000)

- Description of the files available:
http://www.ukbiobank.ac.uk/wp-content/uploads/2017/07/ukb_genetic_file_description.txt

**Versions**
- V1 – 1st release ~150,000 samples
- V2 – full release ~500,000 samples
- V3 – for imputation files only – due to an error in the initial V2 imputation release these files were re-processed and re-released

**CSF3 central dataset**

To access you must be added to a group with correct permissions (dataset-ukbiobank-full group)

Email the research IT team with confirmation of approved UKB data access

module load tools/env/ukbiobank-full-release-2018

Sets a number of environment variables:
UKBB_FILELIST
/mnt/data-sets/ukbiobank/full-release/filelist.2018.txt – list paths to files

More info:
http://ri.itservices.manchester.ac.uk/csf-apps/software/applications/ukbiobank/

# Calls

- The genotype calls are in binary PLINK format (.bed, .bim, .fam) - see https://www.cog-genomics.org/plink/1.9/formats for details of the formats.

- The **BIM** file determines the order of markers in the calls and all of the other genotype data sets. The SNP_id is the rsid where it is available or the Affymetrix_SNP_id otherwise.

- The positions are **GRCh37** coordinates.

- The **FAM** file determines the order of samples in the calls and all of the other genotype data sets. The FAM file includes 'Batch' in the Phenotype field (6th column).  - this file is project specific – the eids are different between projects

# Imputed data

- The imputed genotype calls are in BGEN v1.2 format (.bgen, .sample, .bgi)

# Files

- Calls BED — ukb_cal_chrN_v2.bed
- Calls BIM — ukb_snp_chrN_v2.bim
- Calls FAM — ukbA_cal_v2_sP.fam → Project-specific – download this yourself "ukbgene evc -c1 –m" from linux command line
- Marker-QC — ukb_snp_qc.txt
- Sample-QC — ukb_sqc_v2.txt
- Relatedness — ukbA_rel_sP.txt
- Imputation BGEN — ukb_imp_chrN_v3.bgen
- Imputation BGI — ukb_bgi_chrN_v3.bgi
- Imputation MAF+info — ukb_mfi_chrN_v3.txt
- Imputation sample — ukbA_imp_autosome_v3_sP.sample
- Haplotypes BGEN — ukb_hap_chrN_v3.bgen
- Haplotypes BGI — ukb_hbg_chrN_v3.bgi
- HLA Imputation — ukb_hla_v2.txt
- Intensity — ukb_int_chrN_v2.bin
- Confidences — ukb_con_chrN_v2.txt
- CNV log2r — ukb_l2r_chrN_v2.txt
- CNV baf — ukb_baf_chrN_v2.txt
- SNP-posterior — ukb_snp_posterior_chrN.bin
- Batch — ukb_snp_posterior.batch

# Whole Exome Sequencing

- 1st 50,000 released
- 39Mbp exome
- 75bp paired end reads
- Illumina NovaSeq 6000
- Mapped to GRCh38 reference

- Variant called through two pipelines:
  - 'FE' – 'Functional Equivalent' pipeline (GATK)
  - 'SPB' - Regeneron's Seal Point Balinese pipeline

- PLINK format release of all samples together
- Individual gVCFs can also be downloaded
- Download using 'ukbgene' utility

## Category 170
Exome sequences - Genomics

**Description**

The first tranche of UKBiobank whole exome sequencing (WES) is now available for ~50,000 UK Biobank participants.

**To ensure equality of access the individual level data is currently embargoed to allow all researchers an opportunity to download the PLINK formatted data. The VCF files will be released by early-April followed by the CRAM files. Researchers who already have access to UK Biobank genetic data do NOT have to submit new baskets to request exome data - this will be done for them automatically by the Access team.**

This sample set prioritizes individuals with whole body MRI imaging data, enhanced baseline measurements, hospital episode statistics (HES), and/or linked primary care records. Additionally, one disease area was selected for enrichment: individuals with admission to hospital with a primary diagnosis of asthma (ICD10 J45 or J46). The sequenced set includes 194 parent-offspring pairs, 613 full-sibling pairs, including 26 trios, 1 monozygotic twin pair and 195 second degree genetically determined relationships.

Exomes were captured with the IDT xGen Exome Research Panel v1.0 including supplemental probes. The basic design targets 39Mbp of the human genome (19,396 genes). Multiplexed samples were sequenced with dual-indexed 75x75bp paired-end reads on the Illumina NovaSeq 6000 platform using S2 flow cells. In each sample and among targeted bases, coverage exceeds 20X at 94.6% of sites on average. Complete sequencing protocols are described in detail by the summary manuscript (add link when available). This manuscript also fully describes the "SPB" primary and secondary analysis pipeline that converts raw sequencing data to a quality-controlled set of population variation. The SPB pipeline first converted all raw sequencing data to FASTQs according to Illumina NovaSeq best practices and aligned those reads to the GRCh38 reference genome with BWA-mem to generate a CRAM file for each sample. After read-duplicate marking, SNVs and indels were called for with WeCall (GenomicsPLC), generating a gVCF per sample. These gVCFs were joint genotyped using GLnexus (https://www.biorxiv.org/content/10.1101/572347v1) to create a single, unfiltered project-level VCF (pVCF). Genotype depth filters (SNV DP>7, indel DP>10) were applied prior to variant site filters requiring at least one variant genotype passing an allele balance filter (heterozygous SNV AB>0.15, heterozygous indel<0.20), resulting in a second 'filtered' pVCF. A total of 4,735,722 variants are identified within targeted regions, with 9,693,536 variants identified across all covered bases including 100bp regions flanking the capture targets.

To maximize data utility and ease of use, an additional "Functionally Equivalent" (FE) pVCF was generated from FASTQs, following the primary analysis protocol described in the 2018 manuscript (PMID: 30279509) and then subject to GATK 3.0 variant calling and hard filtering of variants with inbreeding coefficient<-0.03 or without at least one variant genotype of DP≥10, GQ≥20 and, if heterozygous, AB≥0.20.

| 10 Data-Fields | 1 Parent Category | 1 Resource |
|---|---|---|

**Field ID Description**

| | |
|---|---|
| 23170 | Population-level SPB variants, PLINK format ‡ |
| 23160 | Population-level FE variants, PLINK format ‡ |
| 23171 | Exome SPB variant call files (VCFs) ‡ |
| 23172 | Exome SPB variant calls indices ‡ |
| 23173 | Exome SPB CRAM files ‡ |
| 23174 | Exome SPB CRAM indices ‡ |
| 23161 | Exome FE variant call files (VCFs) ‡ |
| 23162 | Exome FE variant calls indices ‡ |
| 23163 | Exome FE CRAM files ‡ |
| 23164 | Exome FE CRAM indices ‡ |

- Fields marked ‡ are blob/bulk.

Improving the health of future generations

Variant call data for 50,000 exomes (sample level **VCF files**) ~10TB

CRAM files also available ~150TB

Capacity for downloading centrally…..

Initial comparison of arrays and exomes…

Mike Weedon (twitter) - analysis of 3000 QC'd exome SNPS reveals most SNPs on array with MAF <0.005% are FPs