

Data processing and output of the Qatar Genome Project Pilot phase to be made available to the PPM program

For 3000 individuals, raw whole genome read data is generated by Illumina HiSeq X Ten¹ sequencers and converted from the native BCL format to paired end FASTQ² format using bcl2fastq³[v2.16]. The quality of the raw data is then assessed using fastqc⁴. Data passing quality control is then aligned to the reference genome sequence (build GRCh37 (hs37d5)⁵) using the bwa-kit⁶ aligner[v7.12]. Variant calling is performed using GATK⁷ haplotype caller[v3.3] and annotation of the resulting VCF⁸ is performed using snpeff⁹[v4.1b] and the following databases(dbsnp¹⁰ v138 and dbNSFP¹¹ v2.9) . In the future, we will also provide mapping and variant calling using the new GRCh38 reference genome¹². Annotation using VEP¹³ tool will be made available as well.

The data types which will be provided to the PPM investigators are BAM and VCF.

¹<http://www.illumina.com/systems/hiseq-x-sequencing-system/system.html>

²https://en.wikipedia.org/wiki/FASTQ_format

³<https://support.illumina.com/downloads/bcl2fastq-conversion-software-v216.html>

⁴<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>

⁵ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/reference/phase2_reference_assembly_sequence/README_human_reference_20110707

⁶<https://github.com/lh3/bwa/tree/master/bwakit>

⁷<https://www.broadinstitute.org/gatk/>

⁸<https://github.com/samtools/hts-specs>

⁹<http://snpeff.sourceforge.net/>

¹⁰http://www.ncbi.nlm.nih.gov/SNP/snp_summary.cgi?view+summary=view+summary&build_id=138

¹¹

¹¹<https://sites.google.com/site/jpopgen/dbNSFP>

¹²<http://www.ncbi.nlm.nih.gov/projects/genome/assembly/grc/human/>

¹³<http://www.ensembl.org/info/docs/tools/vep/index.html>