

## 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 Ten<sup>1</sup> sequencers and converted from the native BCL format to paired end FASTQ<sup>2</sup> format using bcl2fastq<sup>3</sup>[v2.16]. The quality of the raw data is then assessed using fastqc<sup>4</sup>. Data passing quality control is then aligned to the reference genome sequence (build GRCh37 (hs37d5)<sup>5</sup>) using the bwa-kit<sup>6</sup> aligner[v7.12]. Variant calling is performed using GATK<sup>7</sup> haplotype caller[v3.3] and annotation of the resulting VCF<sup>8</sup> is performed using snpeff<sup>9</sup>[v4.1b] and the following databases(dbsnp<sup>10</sup> v138 and dbNSFP<sup>11</sup> v2.9) . In the future, we will also provide mapping and variant calling using the new GRCh38 reference genome<sup>12</sup>. Annotation using VEP<sup>13</sup> tool will be made available as well.

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

---

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

<sup>2</sup>[https://en.wikipedia.org/wiki/FASTQ\\_format](https://en.wikipedia.org/wiki/FASTQ_format)

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

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

<sup>5</sup>[ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/reference/phase2\\_reference\\_assembly\\_sequence/README\\_human\\_reference\\_20110707](ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/reference/phase2_reference_assembly_sequence/README_human_reference_20110707)

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

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

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

<sup>9</sup><http://snpeff.sourceforge.net/>

<sup>10</sup>[http://www.ncbi.nlm.nih.gov/SNP/snp\\_summary.cgi?view+summary&view+summary&build\\_id=138](http://www.ncbi.nlm.nih.gov/SNP/snp_summary.cgi?view+summary&view+summary&build_id=138)

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

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

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