

H_Boger_CMV_WGHum-SeqWholeExome_180608_1

General Methods Summary:

Sample QC:

2% gels, OD260 readings and volume checks are done upon sample receipt at CIDR to confirm adequate quantity and quality of genomic DNA. In addition, samples were processed with an **Illumina InfiniumQCArray-24v1-0** array to confirm gender, identify unexpected duplicates and relatedness, confirm study duplicates and relatedness, provide sample performance information and sample identity confirmation against the sequencing data. Problems are noted in the problems report (Sample_Info directory).

Exome Capture:

The Agilent SureSelectXT HumanAllExon V6 S07604514 was used. Details below.

Library Preparation, Enrichment:

Low_Input:

A low input library prep protocol developed at CIDR (Marosy et al) was performed. Libraries are prepared from 50ng of genomic DNA, sheared for 80s using the Covaris LE220plus instrument (Covaris). The Kapa Hyper prep kit is used to process the sheared DNA into amplified indexed adapter ligated fragments. All processing was done in 96 well plate formats using robotics (Beckman FXp, Perkin Elmer Janus, Agilent Bravo, Beckman NX). 'With Bead' clean ups were used following shearing and adapter ligation. 750ng of amplified library was used in an enrichment reaction following Agilent protocols (24 hour hybridization). Post-capture washing was done using the Agilent protocol except the 'off-bead' catch process from Fisher et al., was incorporated (samples are not eluted off the DynaBeads (Invitrogen) directly adding post-capture PCR master mix and indexes to the beads). Post-capture PCR was done according to the Agilent protocol, with the adjustment of water volume and PCR cycles where needed.

Sequencing:

Libraries were sequenced on the HiSeq2500 platform with template generation on the cBot. 72 samples per flowcell, 125 bp paired end runs and sequencing chemistry kits HiSeq PE Cluster Kit v4 and HiSeq SBS kit v4.

Primary Analysis:

Intensity analysis and base calling were performed through the Illumina Real Time Analysis (RTA) software (version 1.18.66.4). Basecall files were demultiplexed from a binary format (BCL) to single sample fastq files using a demultiplexer written at CIDR as part of CIDRSeqSuite version 7.5.0 (unpublished).

Secondary Analysis:

Fastq files were aligned with BWA mem (Li H. 2013) version 0.7.15 to the 1000 genomes phase 2 (GRCh37) human genome reference. Duplicate molecules were flagged with Picard version 2.17.0. Base call quality score recalibration and binning (2,10,20,30) were performed using the Genome Analysis Toolkit (GATK) (McKenna et al., 2010) version v4.0.1.1. Cram files were generated using SAMTools version 1.5. GATK's reference confidence model workflow was used to perform joint sample genotyping using GATK version 3.7. Briefly this workflow entails; 1) Producing a gVCF (genomic VCF) for each sample individually using HaplotypeCaller (--emitRefConfidence GVCF) and --max_alternate_alleles was set to 3 for all bait intervals to generate likelihoods that the sites are homozygote reference or not 2) Joint genotyping the single sample gVCFs together with GenotypeGVCFs to produce a multi-sample VCF file.

Variant filtering was done using the Variant Quality Score Recalibration (VQSR) method (DePristo et al., 2011). For SNVs, the annotations of MQRankSum, QD, FS, ReadPosRankSum, MQ and SOR were used in the adaptive error model. HapMap3.3, Omni2.5 and 1000G phase high confidence snp calls were used as training sites with HapMap3.3 and Omni2.5 used as the truth set. SNVs were filtered to obtain all variants up to the 99.5th percentile of truth sites (0.5% false negative rate). For indels, the annotations of FS, ReadPosRankSum, MQRankSum, QD and SOR were used in the adaptive error model (4 max Gaussians allowed). A set of curated indels obtained from the GATK resource bundle (Mills_and_1000G_gold_standard.indels.b37.vcf) were used as training and truth sites. Indels were filtered to obtain all variants up to the 99th percentile of truth sites (1% false negative rate).

An additional/optional VCF file was created where genotypes for biallelic SNPs were further refined using CalculateGenotypePosteriors using allele frequency information from 1000 genomes phase 3 data

(ALL.wgs.phase3_shapeit2_mvncall_integrated_v5.20130502.sites.vcf) as well as Exome Aggregation Consortium data (ExAC.r0.3.sites.vcf).

Variant Annotation and Summary Statistics:

All variants in the final multi-sample VCF file were annotated using Annovar (version 2013_02_21) against a variety of data sources including gene annotation, function prediction and frequency information (see dictionary file).

Summary Statistics:

Summary statistics (for SNVs and INDELS) on the multi-sample .vcf file were calculated for each variant (both PASS and FAIL) including counts and frequencies of alleles and genotypes, missing rates, overall quality scores, and mean depth.

References:

Fisher et al; A scalable, fully automated process for construction of sequence-ready human exome targeted capture libraries. *Genome Biology* 2011, 12:R1

Li H. (2013) Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. arXiv:1303.3997v1 [q-bio.GN].

Church D., Deanna M., Schneider V. *et al.* Modernizing reference genome assemblies. *PLOS Biol.* 9, e1001091 (2011)

International Human Genome Sequencing Consortium. Finishing the euchromatic sequence of the human genome. *Nature* 43, 931–945 (2004)

McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D, Gabriel S, Daly M, DePristo MA (2010). The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* 20:1297-303

DePristo M, Banks E, Poplin R, Garimella K, Maguire J, Hartl C, Philippakis A, del Angel G, Rivas MA, Hanna M, McKenna A, Fennell T, Kernytsky A, Sivachenko A, Cibulskis K, Gabriel S, Altshuler D and Daly, M (2011). A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nature Genetics.* 43:491-498

From FastQ Data to High-Confidence Variant Calls: The Genome Analysis Toolkit Best Practices Pipeline
Van der Auwera GA, Carneiro M, Hartl C, Poplin R, del Angel G, Levy-Moonshine A, Jordan T, Shakir K, Roazen D, Thibault J, Banks E, Garimella K, Altshuler D, Gabriel S, DePristo M, 2013 CURRENT PROTOCOLS IN BIOINFORMATICS 43:11.10.1-11.10.33

Wang K, Li M, Hakonarson H. ANNOVAR: Functional annotation of genetic variants from next-generation sequencing data, *Nucleic Acids Research*, 38:e164, 2010

An integrated map of genetic variation from 1,092 human genomes, McVean et Al, *Nature* 491, 56–65 (01 November 2012) doi:10.1038/nature11632

Exome Aggregation Consortium (ExAC), Cambridge, MA (URL: <http://exac.broadinstitute.org>)