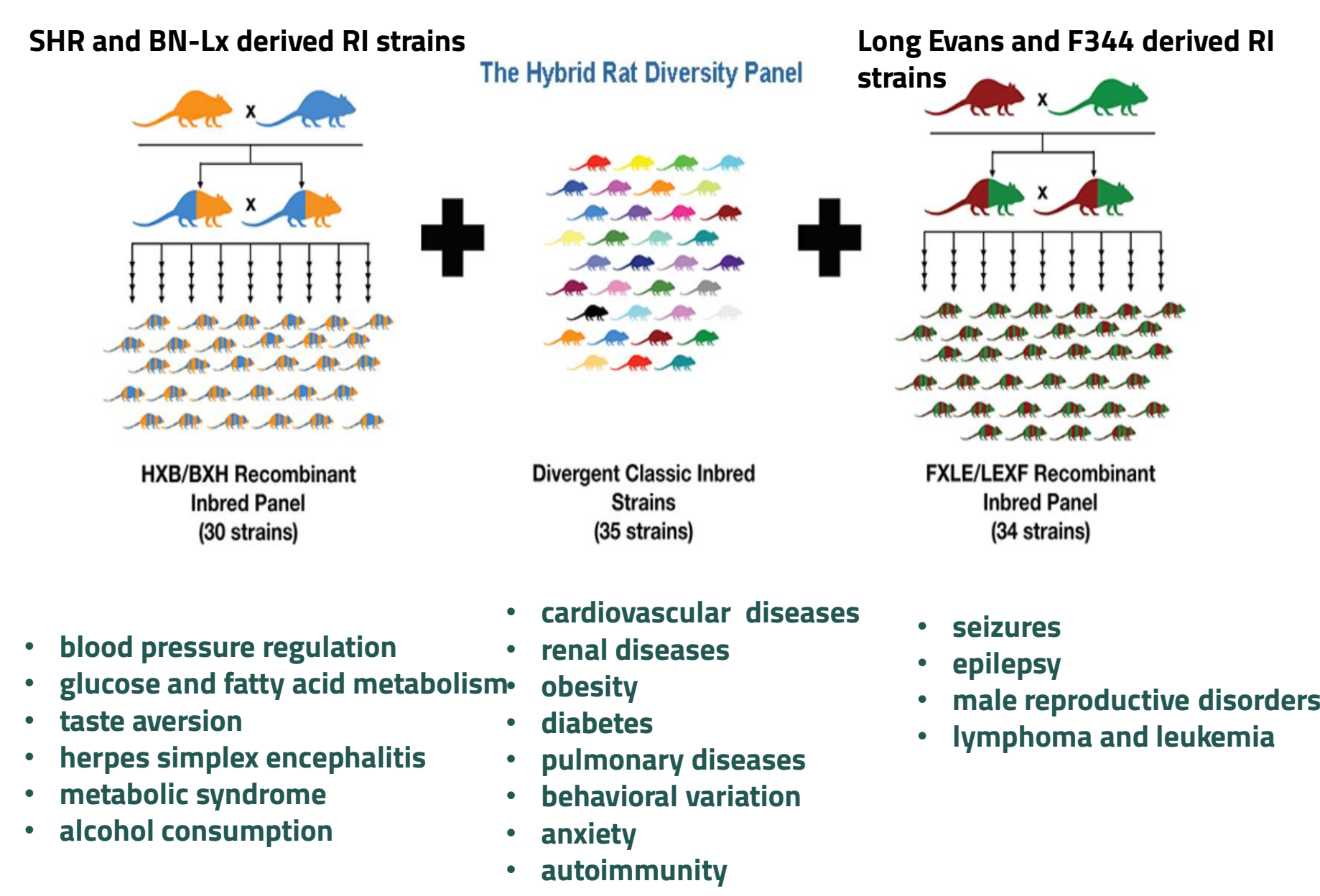


Abstract

The Hybrid Rat Diversity Panel (HRDP) is a group of 96 rat inbred strains selected to study mechanisms of complex traits similar in their pathology to common human diseases. We have analyzed the HRDP whole genomic sequencing data (Illumina short reads) using the high-quality rat reference mRatBN7.2 and variant discovery GATK4 Best Practices recommendations (Broad Institute 2019). The average sequence coverage ranges from 15x to 69x per rat strain sample. We have found that more than 21 mln of germline variants, ~8 mln of short indels and ~13 mln of SNVs, characterize the rat cohort. In addition, we observed a remarkable drop in the number of indels discovered with the new genome reference compared with the old rat assembly from 2014, Rnor6. Our results provide high confidence variants that represent the sequence cohort heterogeneity but also variants that require additional quality testing. Low coverage variants, a fraction of private variants of Brown Norway reference rat strain and genomic regions with accumulation of heterozygous variants will require further analysis. Currently we are incorporating the data into Rat Genome Database (<https://rgd.mcw.edu>) to provide functional annotations to the rat community, as well as assist in prioritization of potential disease-causing mutations. Thus, researchers can access information about strain specific variants on gene report pages, in Variant Visualizer tool and in genomic browser. The data in vcf format are available for retrieval in the 'Download' section of the RGD site.

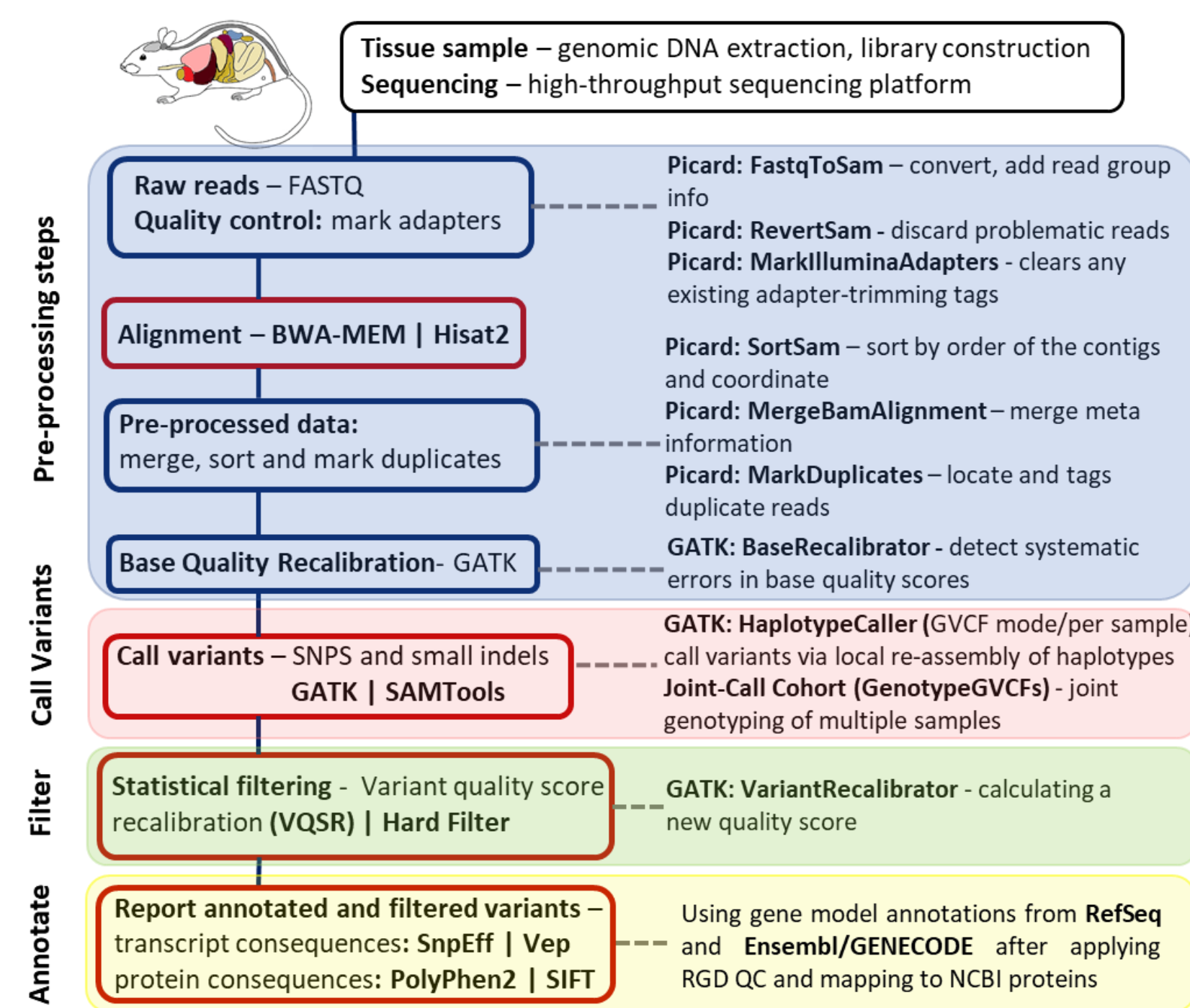
What is the Hybrid Rat Diversity Panel ?



Hybrid Rat Diversity Panel was selected to:

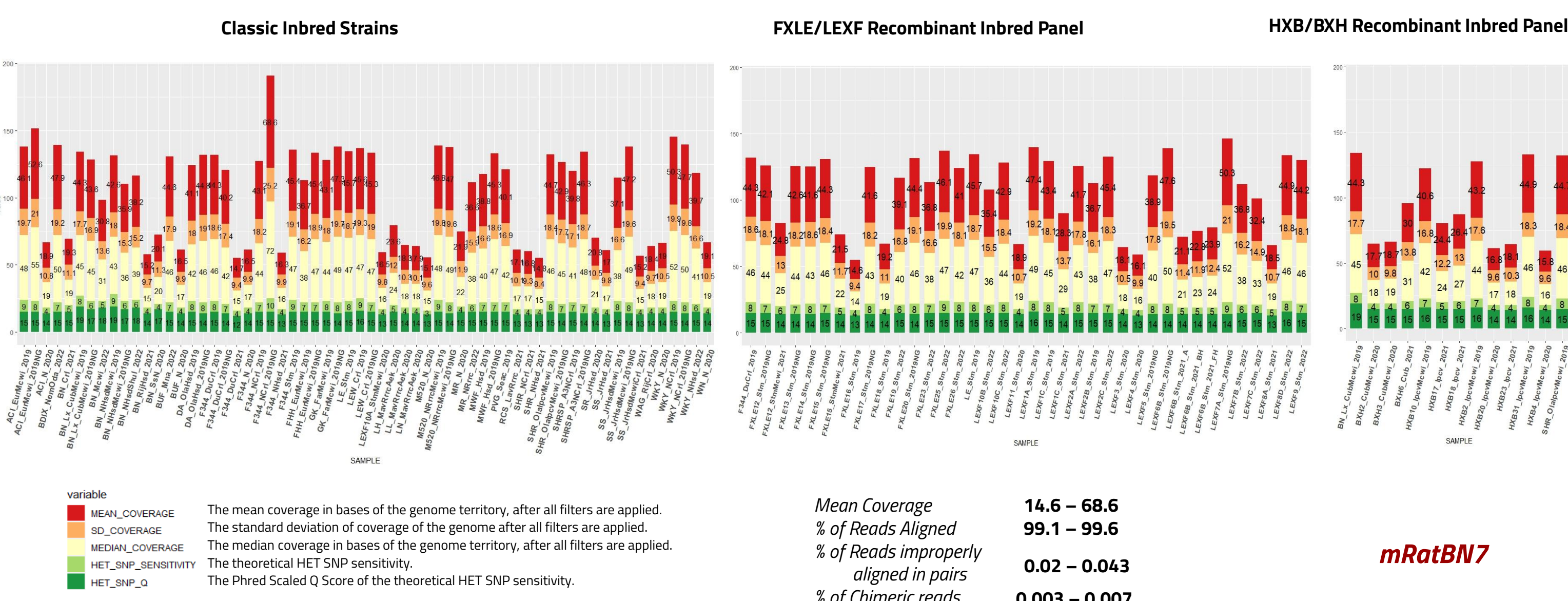
1. Provide stable genetic and phenotypic strains to allow researchers to conduct reproducible experiments
2. Maximize the genetic diversity among strains and to maximize power to detect specific genetic loci associated with a complex trait (QTL mapping resolution)
3. Extend the whole genome sequencing to all HRDP inbred rat strains with susceptibility to different complex diseases
4. Facilitate the translation of disease-related genetics and genomics research to pre-clinical and clinical studies

Analysis workflow

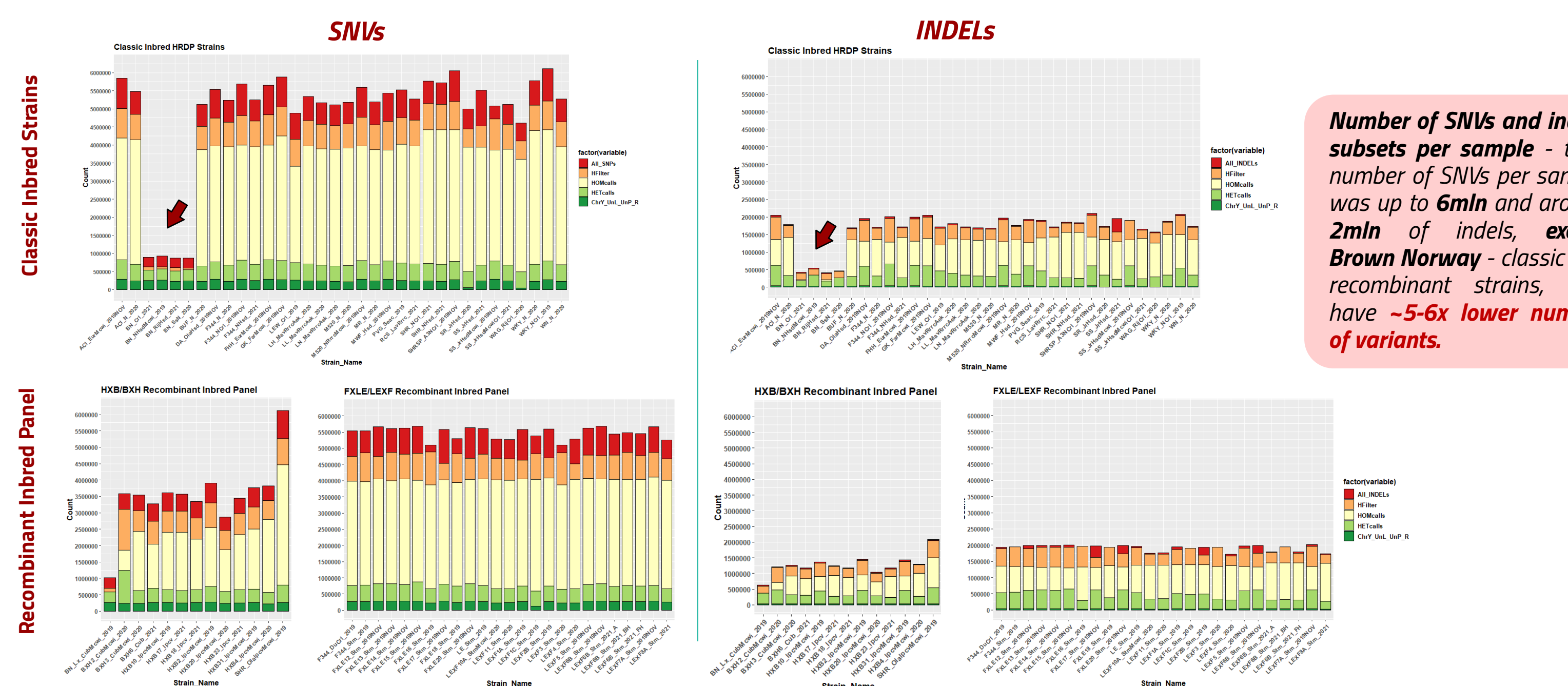


1. Pre-processing steps involve marking adapter sequences, alignment to the new rat genome reference **mRatBN7** and marking duplicates.
2. Base Quality Scores Recalibration corrects biases introduced by sequencing platforms and assigns scores empirically determined from the read data using validated variants.
3. Variant calling is accomplished by running the GATK HaplotypeCaller that simultaneously detects SNVs and indels via local de-novo assembly of haplotypes (method to increase accuracy of the variant call comparing with position-based algorithm).
4. In the filtering process we remove less reliable variant calls: variants with low coverage, low quality, strand biased, located in SNV clusters, and supported by low-confidence read alignment.

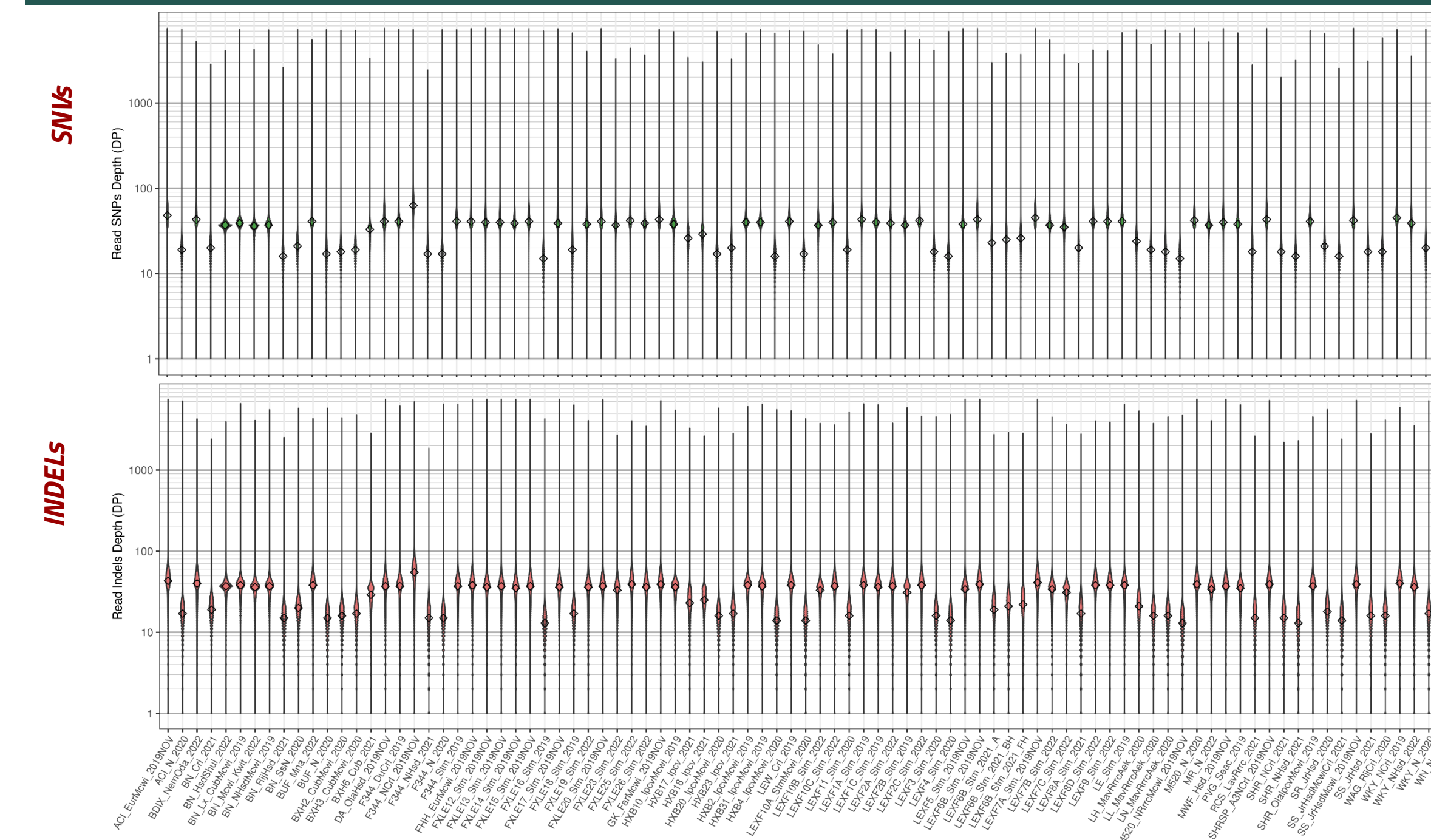
Sequencing Coverage and Variants Discovery



Mapping analysis shows the coverage and performance of the whole genome sequencing. **Mean base coverage** ranges from 15x to 69x and the estimated **sensitivity to detect heterozygous sites** (as a function of coverage and base quality distribution ranges) drops quickly when mean coverage is below 20X.



High Confidence Variants

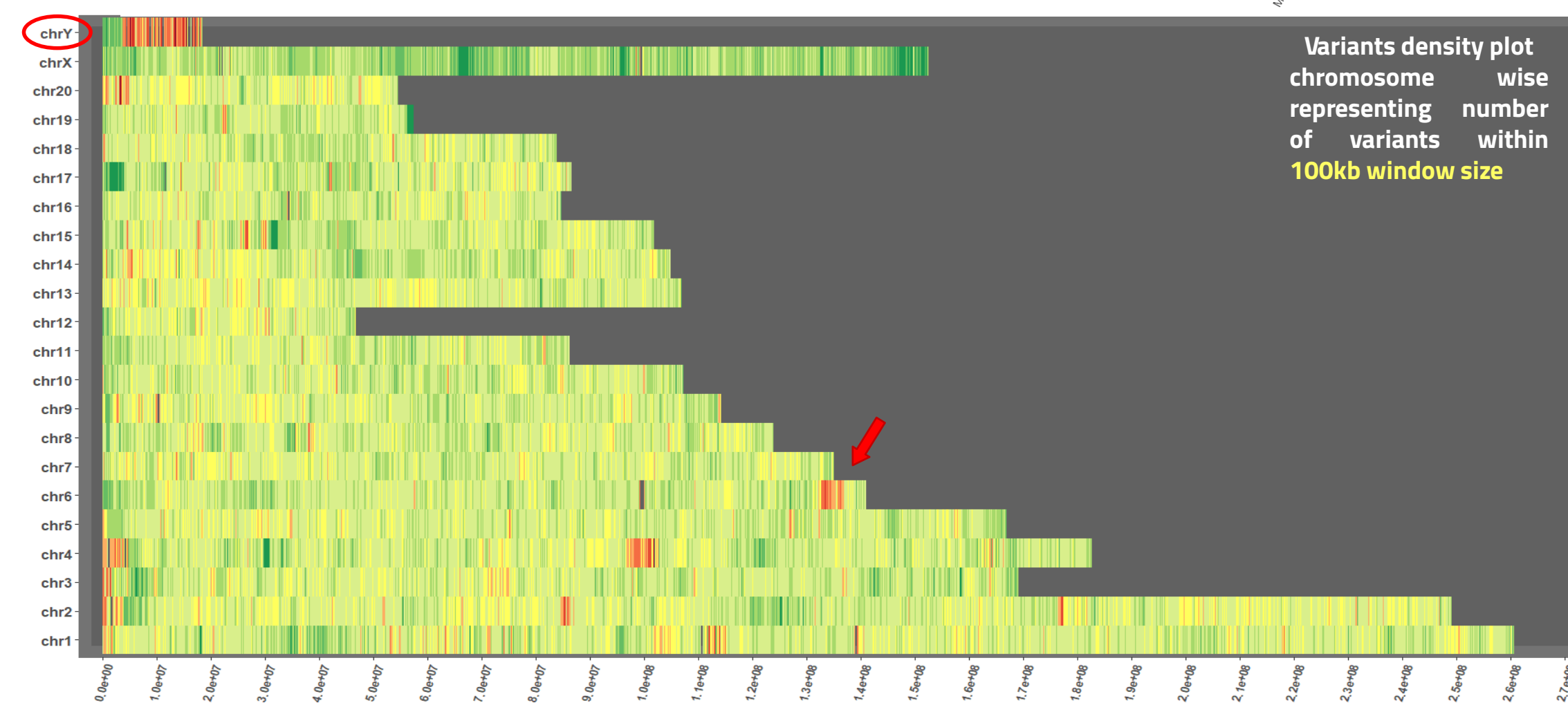


Hard Filter - 6 parameters PASS

1. QualByDepth (QD)
2. FisherStrand (FS)
3. StrandOddsRatio (SOR)
4. RMSMappingQuality (MQ)
5. MappingQualityRankSumTest (MQRankSum)
6. ReadPosRankSumTest (ReadPosRankSum)

All SNVs	17,550,340
HF SNVs	13,205,288
All INDELS	9,428,709
HF INDELS	6,373,427

High confidence variants - Number of variants decrease after applying hard filtering parameters and only variants supported by more than 10 reads were selected as high confidence set.



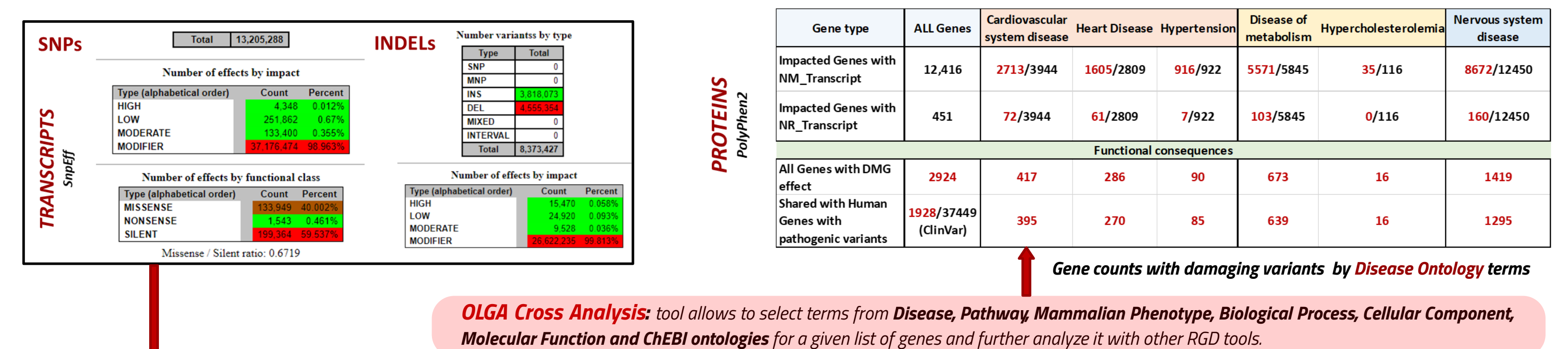
Variants density plot shows genome location with high number of variants accumulation that require further analysis.

HRDP variants in RGD

	Atanur et al. 2013	Hermesen et al. 2015	RGD 2018	HRDP 2022	
Number of analyzed rat strains	27	40	25	47	89
Number of identified high confidence SNVs	9,665,340	9,183,702	8,953,897	11,585,641	13,205,288
Reference Genome	RGSC 3.4 - chromosomes	RGSC 5.0 - whole	RGSC 6.0 - whole	RGSC 6.0 - whole	mRatBN7v2 - whole
Alignment Software	BWA mem	BWA mem	BWA mem	BWA mem	BWA mem
Genomic variants call	GATK v. 1.0.6001	GATK HaplotypeCaller v2.8-1	GATK HaplotypeCaller v3.6	GATK HaplotypeCaller v4.1.3.0	
Variant quality recalibration (VQSR) - true training set	Top 30% of high quality SNVs	Not defined	273,568 selected SNVs	In preparation	
dbSNPs	dbSNP125 - 41,658 (1,291) (35,186)	dbSNP138 - 5,076,239 (5,043,831)	dbSNP149 - 5,075,461 (5,042,280) 4,721,043	dbSNP149 - 5,075,461 (5,042,280) 4,721,043	dbSNP_EVAv3 9653928
Tranche sensitivity threshold	99.0	99.5	95.0	In preparation	

Comparison of previously reported analysis of single nucleotide variants identified in different rat populations.

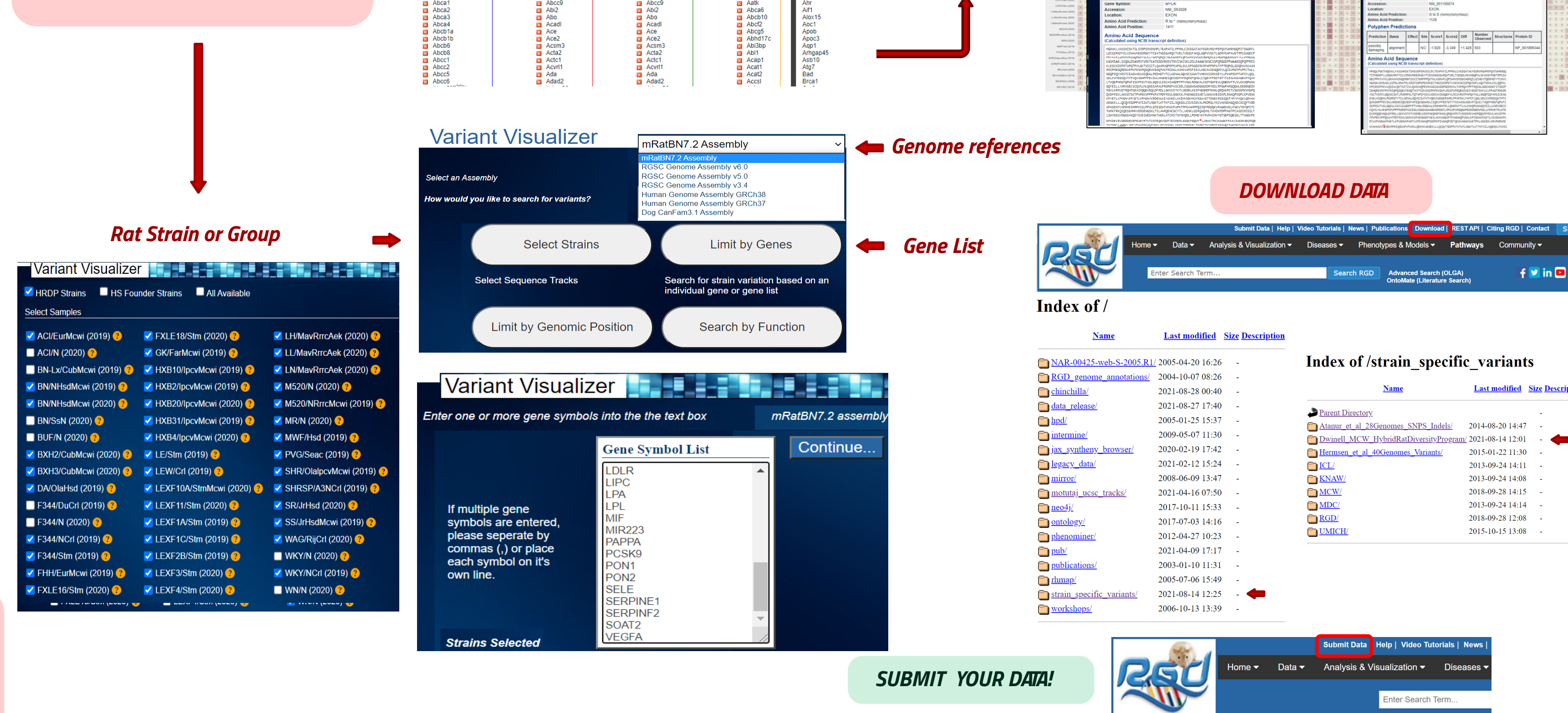
Variants Functional Consequences



OLGA Cross Analysis: tool allows to select terms from Disease, Pathway, Mammalian Phenotype, Biological Process, Cellular Component, Molecular Function and ChEBI ontologies for a given list of genes and further analyze it with other RGD tools.

Variant Visualizer

Variant Visualizer shows genome-wide variants distribution. You can select organism/ breeds/ strains of interest, or input gene list, you can also define genomic position, set parameters/filters for the type(s) of desired variants and the tool will return all of the SNVs and/or indels which match the input criteria, including information on read depth, zygosity, conservation score and more.



Grant support: R24OD022617, R01HL064541

Summary

- Whole genome sequencing of rat diversity panel provides high confidence variants that represent the sequence cohort heterogeneity but also variants that require additional quality testing.
- Genetic variation analysis in rat strains selected for HRDP helps to generate high resolution association mapping and complete systems genetics on complex traits
- We are building the HRDP Portal within RGD that will provide data mining and visualization functions for genomic and phenotypic data