

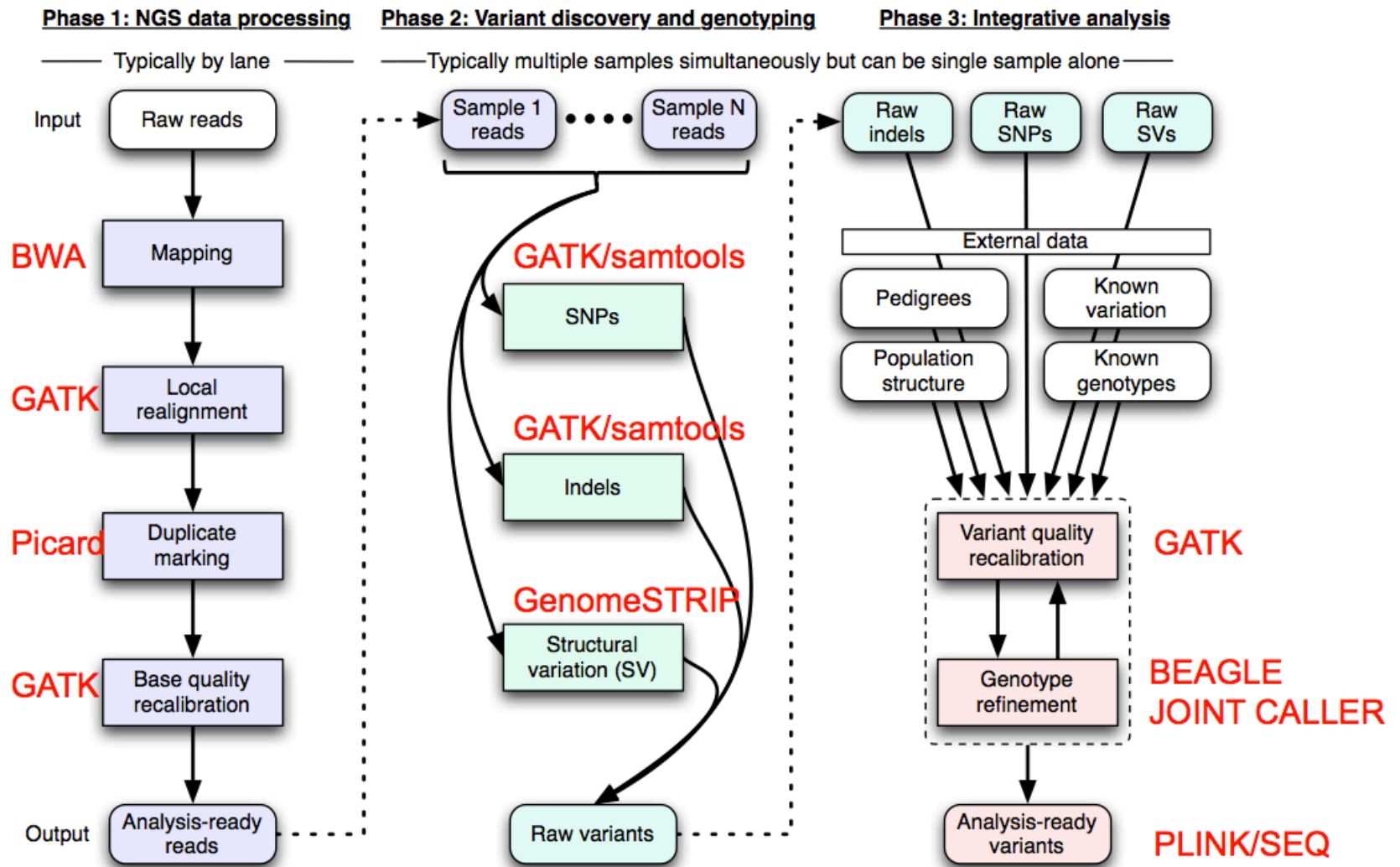
# Analysis of next-generation sequencing data: PSEQ practical

Shaun Purcell

[shaun.purcell@mssm.edu](mailto:shaun.purcell@mssm.edu)

Boulder Course, CO  
March 2013

# The Picard/GATK NGS analysis pipeline



## *example.vcf*

```
##fileformat=VCF4.0
##INFO=<ID=DP,Number=1,Type=Integer,Description="Total read depth">
##INFO=<ID=HM,Number=0,Type=Flag,Description="Seen in HapMap 2 or 3">
##FILTER=<ID=q10,Description="Quality below 10">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype quality">
#CHROM  POS    ID      REF  ALT  QUAL   FILTER  INFO    FORMAT  S001  S002
1       1001  .       A    G    55     PASS    DP=82   GT:GQ   0/0:56 0/1:80
1       8888  rs123   T    A    5       q10     DP=8;HM GT:GQ   ./.:.   0|1:58
```

### **Numeric allele encoding & multi-allelic sites:**

If REF is A, ALT is T      0/0 is A/A      0/1 is A/T

If REF is A, ALT is T,C    0/2 is A/C      2/1 is C/T

Missing genotype      ./.

### **Representing haplotype phase**

0/1      unphased heterozygote

1|0      phased with respect to previous site (implies 01/10 haplotypes )

*(Encoding of haplotypic information will be more explicit in future VCF specs.)*

- What does this genotype mean?

```
GT : AD : DP : GQ : PL  
0/0:366,11:200:99:0,600,5980
```

- GT hard genotype call
- AD and DP read-depth information
- GQ quality score
- PL is (phred-scaled) *genotype-likelihoods* (soft-calls)
  
- VCF likely to be primary out of imputation packages in the future, as it can represent hard-calls (most likely genotype) but also the expected dosage and/or posterior probabilities (and also  $R^2$  in the INFO field, etc)

# Phred-scale

**Phred quality scores are logarithmically linked to error probabilities**

<b>Phred Quality Score</b>	<b>Probability of incorrect base call</b>	<b>Base call accuracy</b>
10	1 in 10	90 %
20	1 in 100	99 %
30	1 in 1000	99.9 %
40	1 in 10000	99.99 %
50	1 in 100000	99.999 %

$$Q = -10 \log_{10} P$$

$$P = 10^{\frac{-Q}{10}}$$

**GT : AD : DP : GQ : PL** 0/1:28,35:62:99:1151,0,889

*Heterozygote genotype; 28 reference reads, 35 alternate; 62 of 63 reads used in calling; genotype quality score of 99; Phred-scaled genotyped likelihoods {1151,0,889}*

Genotype	VCF	PL Phred-scaled genotype likelihood	GL Genotype likelihoods P( read data   true genotype )	P( genotype   read data ) assuming flat priors
Homozygous reference	0/0	1151	0.000	0.000
Heterozygote	0/1	0	1.000	1.000
Homozygous alternate	1/1	889	0.000	0.000

**GT : AD : DP : GQ : PL** 0/0:174,11:1:3:0,3,25

*Homozygous reference genotype; 174 reference reads, 11 alternate; but only 1 read used for calling; genotype quality score of 3 (0-99); Phred-scaled genotyped likelihoods {0,3,25}*

Genotype	VCF	PL Phred-scaled genotype likelihood	GL Genotype likelihoods P( read data   true genotype )	P( genotype   read data ) assuming flat priors
Homozygous reference	0/0	0	1.000	0.66
Heterozygote	0/1	3	0.501	0.33
Homozygous alternate	1/1	25	0.003	0.00

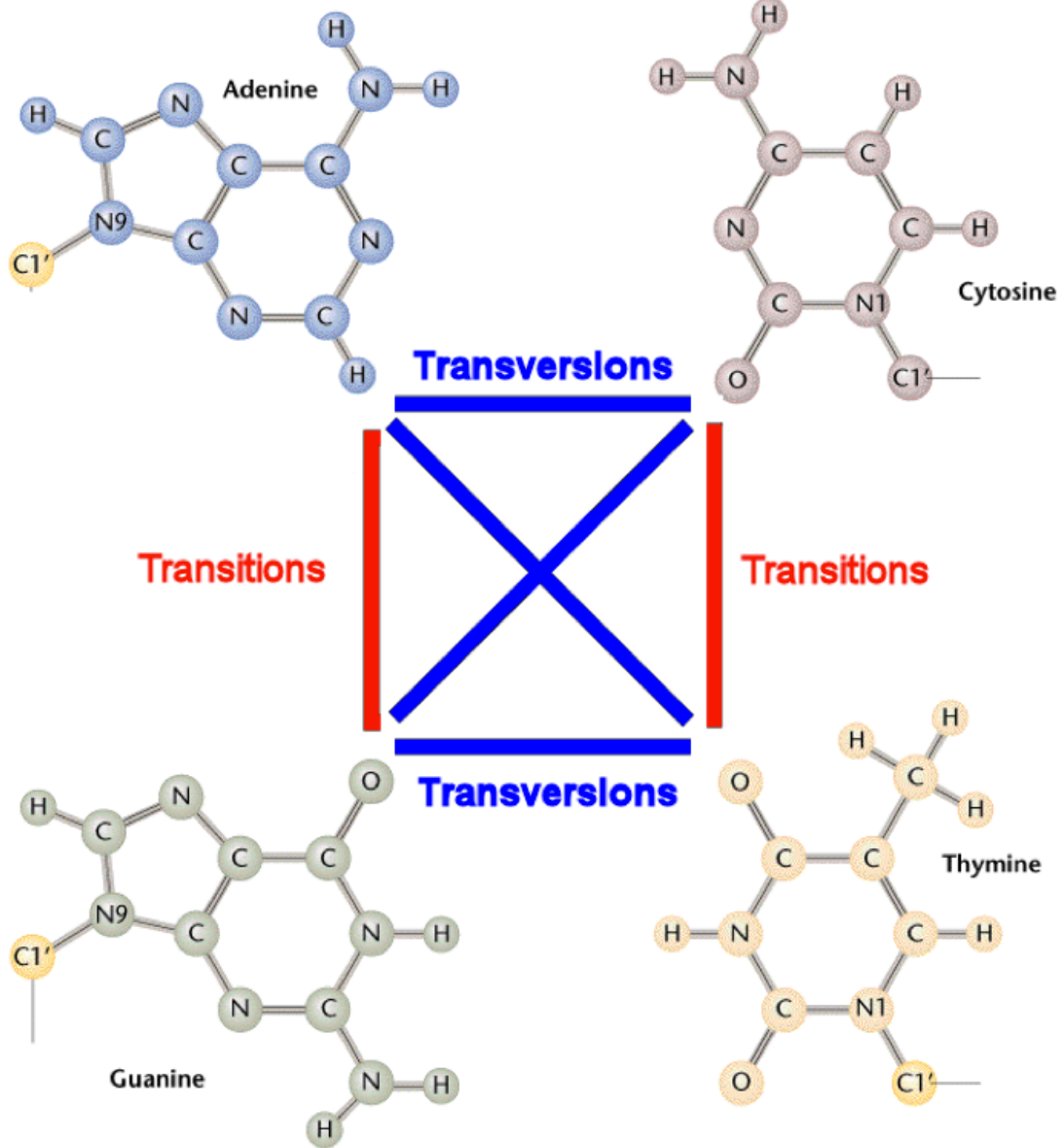
# Variant (and genotype) annotations

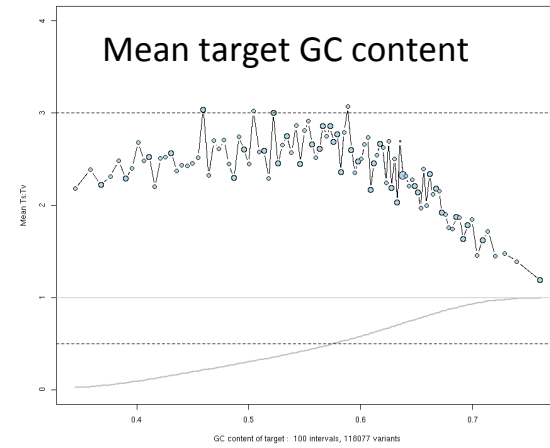
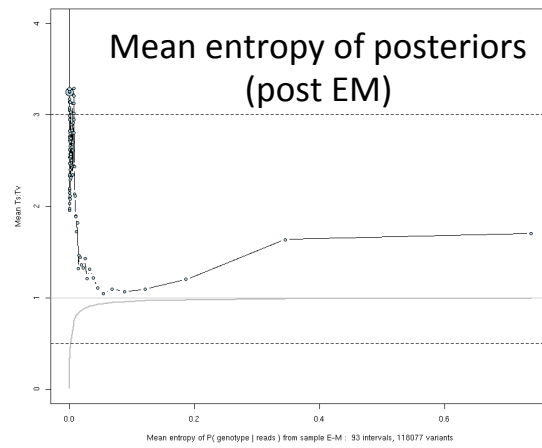
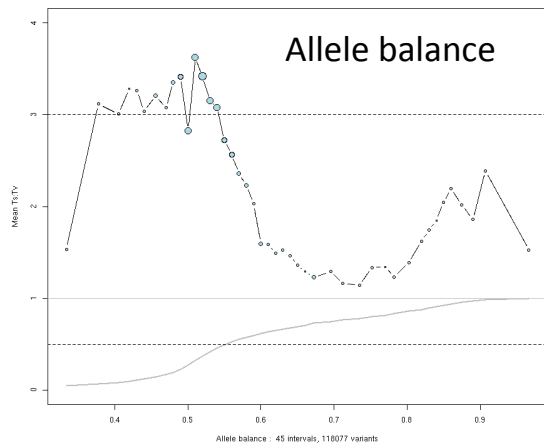
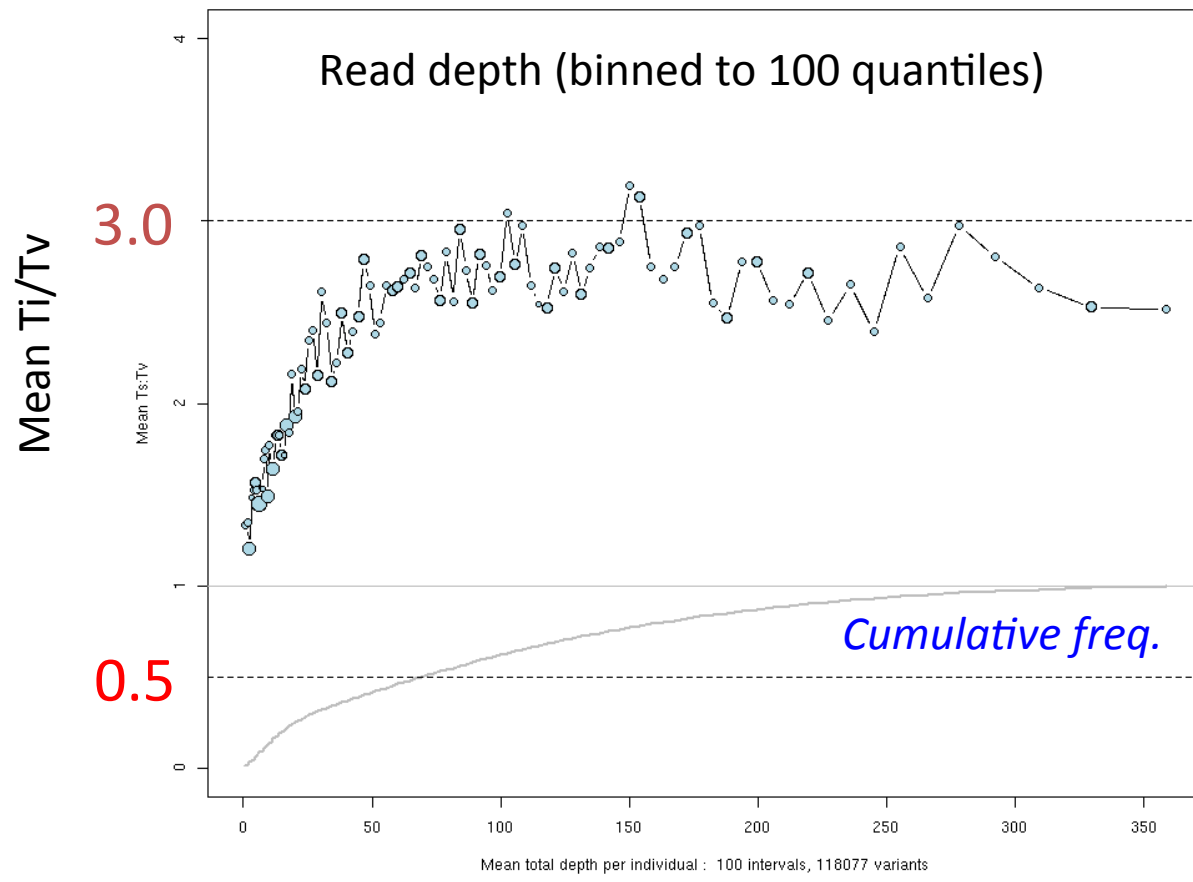
- **Technical**
  - Read depth, allele balance, mean mapping quality, etc
  - QC Filters (PASS, or reasons for exclusion)
  - Hardy-Weinberg disequilibrium
  - (Distribution of) individual genotype likelihoods
- **Population**
  - **Novelty (presence in dbSNP and/or 1000Genomes)**
  - Population frequencies
  - Linkage disequilibrium/phase information on individual genotypes
  - Prior disease/functional associations
- **Genomic**
  - **Gene and coding status**
  - **Transition/transversion**
  - Ancestral/derived allele status



two-ring **purines**

one-ring **pyrimidines**

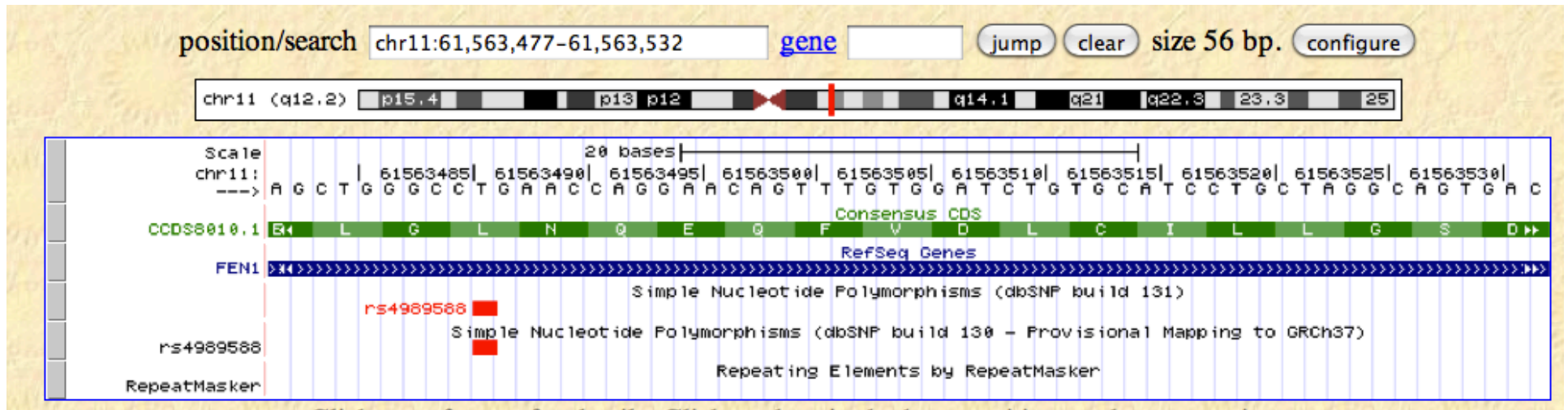




### AMINO ACID TABLE

		Second Position									
		U		C		A		G			
		code	Amino Acid	code	Amino Acid	code	Amino Acid	code	Amino Acid		
First Position	U	UUU	phe	UCU	ser	UAU	tyr	UGU	cys	U	
		UUC		UCC		UAC		UGC		C	
		UUA	leu	UCA		UAA	<b>STOP</b>	UGA	<b>STOP</b>	A	
		UUG		UCG		UAG	<b>STOP</b>	UGG	trp	G	
	C	CUU	leu	CCU	pro	CAU	his	CGU	arg	U	
		CUC		CCC		CAC		CGC		C	
		CUA		CCA		CAA	gln	CGA		A	
		CUG		CCG		CAG		CGG		G	
	A	AUU	ile	ACU	thr	AAU	asn	AGU	ser	U	
		AUC		ACC		AAC		AGC		C	
		AUA		ACA		AAA	lys	AGA	arg	A	
		AUG		ACG		AAG		AGG		G	
	G	GUU	val	GCU	ala	GAU	asp	GGU	gly	U	
		GUC		GCC		GAC		GGC		C	
		GUA		GCA		GAA	glu	GGA		A	
		GUG		GCG		GAG		GGG		G	

Third Position

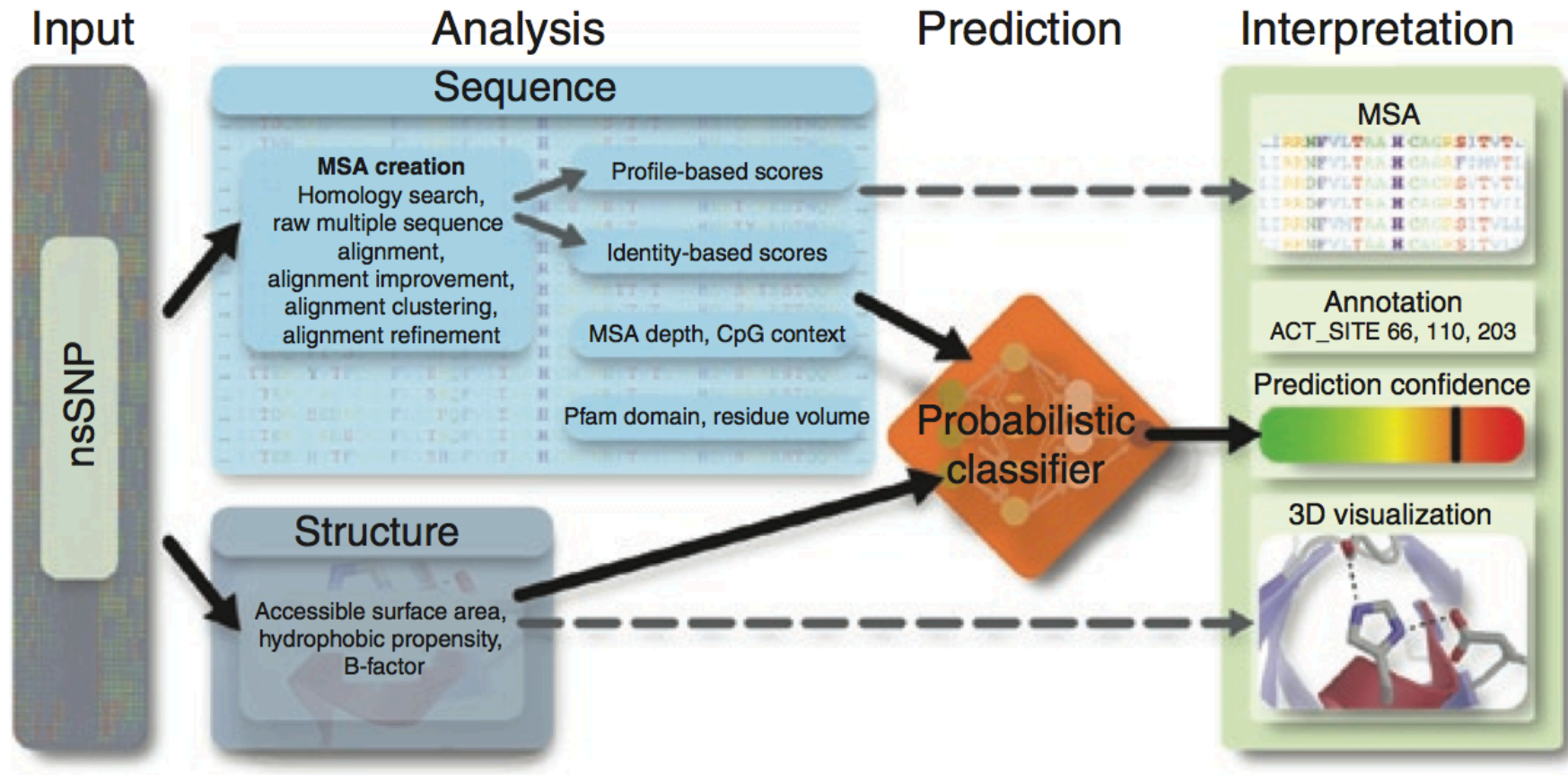


Gene Model(s)							
Function	mRNA				Protein		
	SNP to mRNA	Accession	Position	Allele change	Accession	Position	Residue change
missense	+	<a href="#">NM_004111.4</a>	<a href="#">1025</a>	CTG ⇒ CAG	<a href="#">NP_004102.1</a>	<a href="#">218</a>	L [Leu] ⇒ Q [Gln]

**In addition to assigning missense/nonsense coding status, other complications include:**

- Splice-site, UTR, regions
- Frameshift mutations (indels)
- Multi-nucleotide polymorphisms
- Full haplotype-based per-individual annotation & compound heterozygosity
- Nonsense-mediated decay
- Ranking of missense variants
- Genomic annotation for non-coding variants

# PolyPhen2: predicting the damaging effects of missense mutations



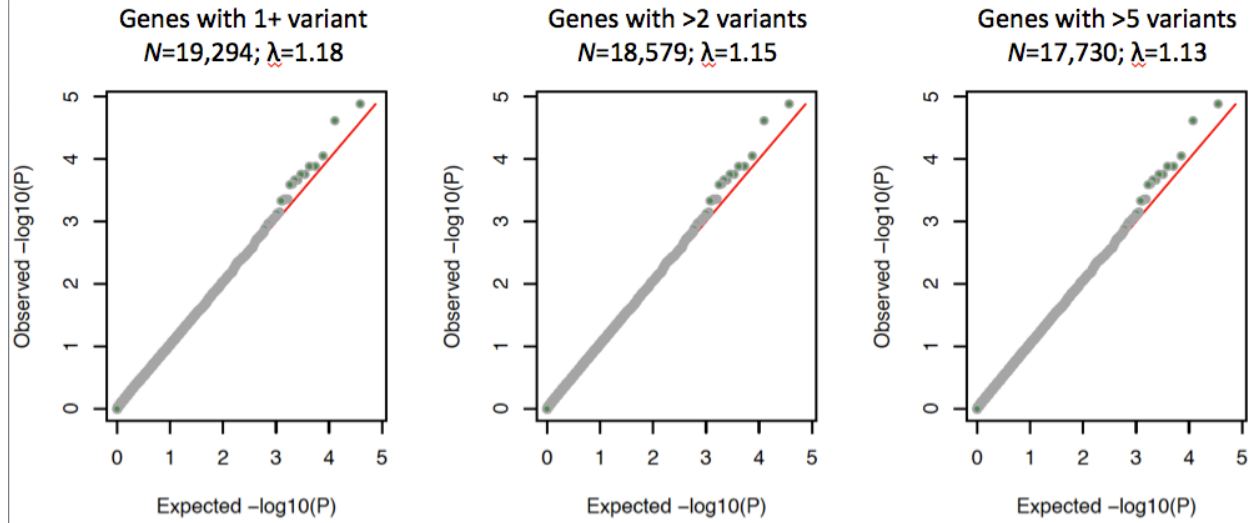
Adzhubei *et al.* Nature Methods (2010).

# Functional annotation of variants

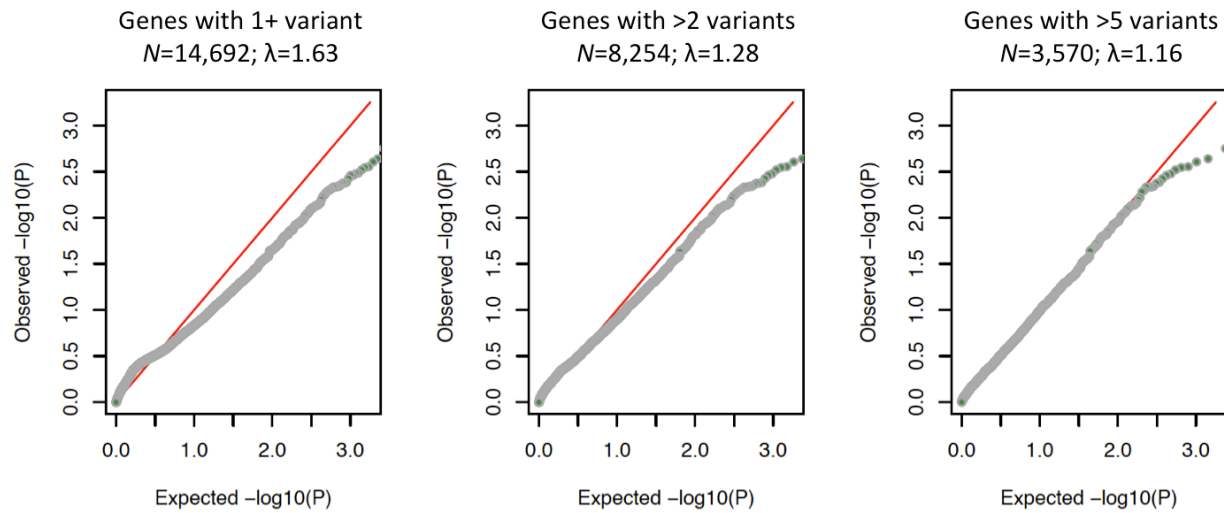
Type	Number	MAF	% singletons
Intronic	101,799	0.066	25
Silent	89,732	0.047	40
Missense	136,847	0.025	51
Benign	60,626	0.038	45
Possibly damaging	22,676	0.020	52
Probably damaging	44,592	0.010	59
Essential splice site	5,020	0.016	54
Nonsense	2,902	0.010	66

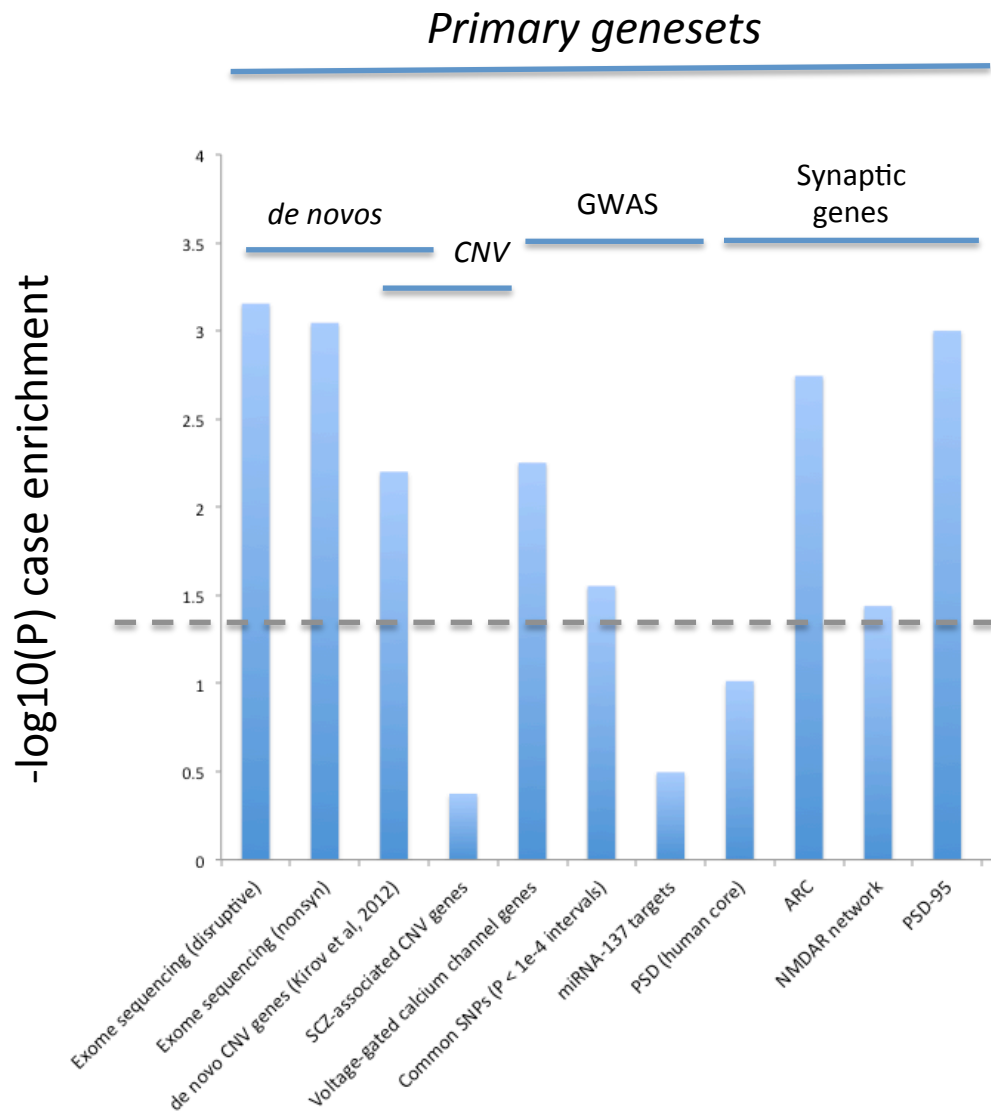
*(based on RefSeq transcripts and hg19; missense ranking w/ PolyPhen2)*

### SKAT | all variants | MDS covariates



### SKAT | nonsyn. (strict) variants | MDS covariates

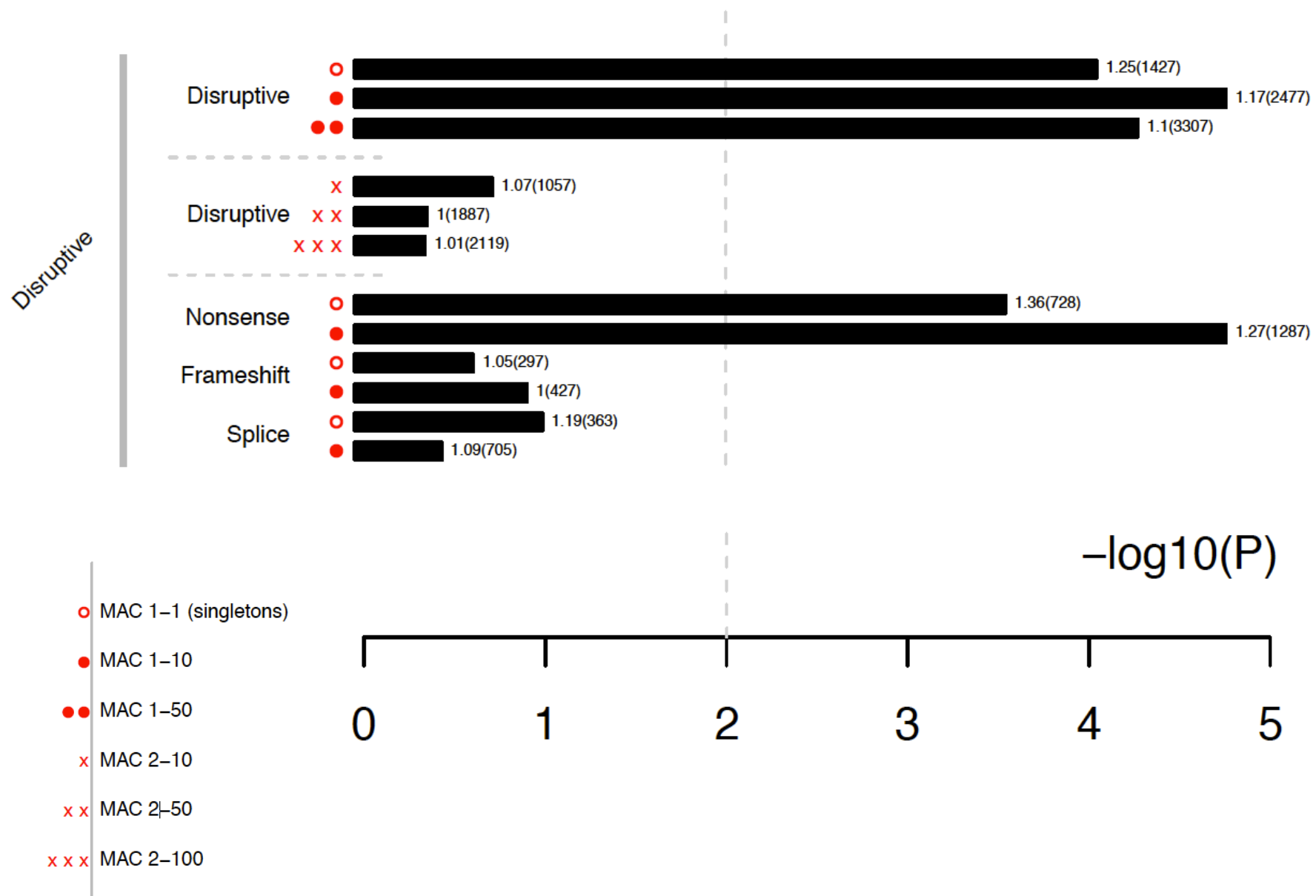




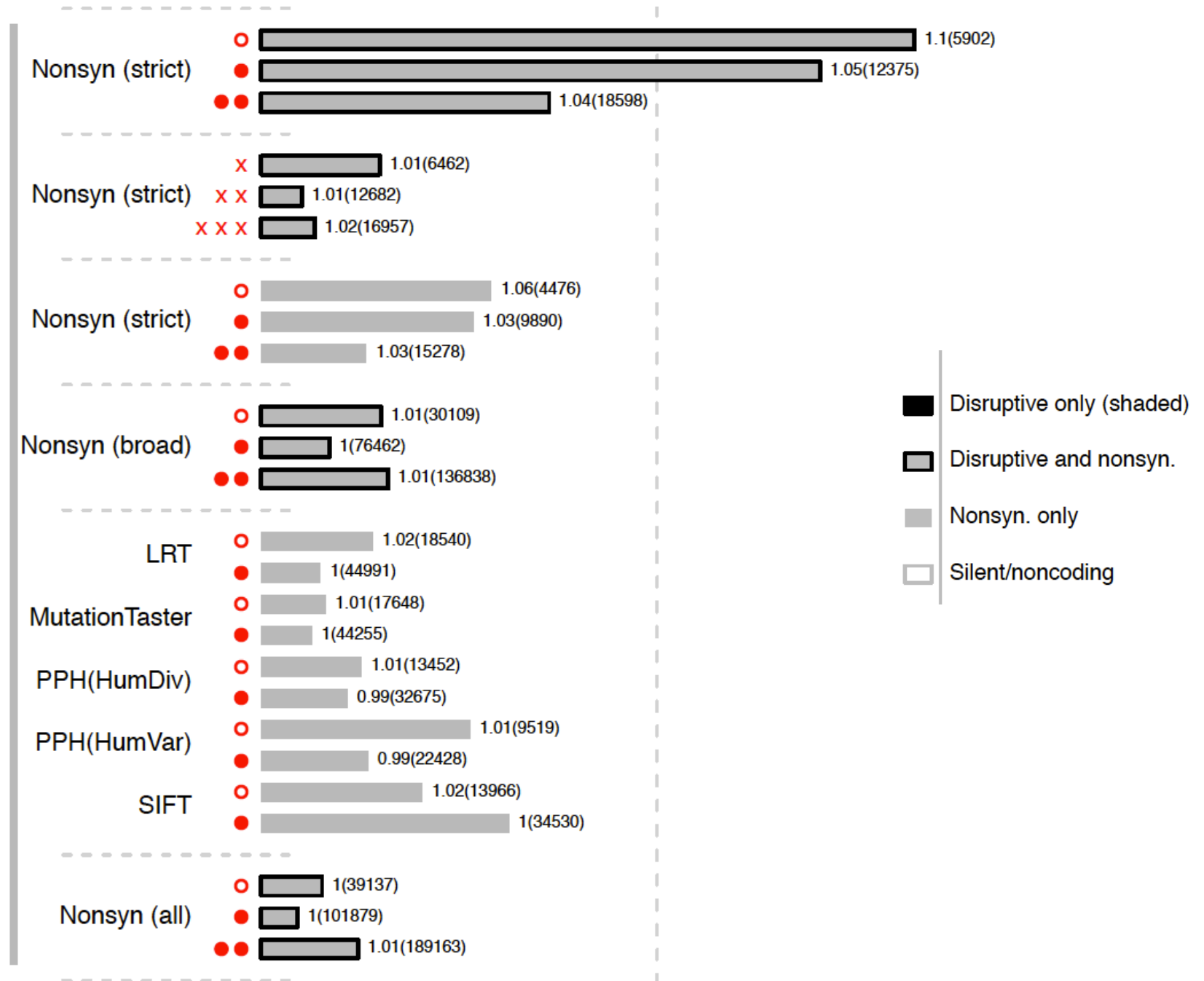
Summary of case burden geneset enrichment results for disruptive mutations (MAF<0.1%)

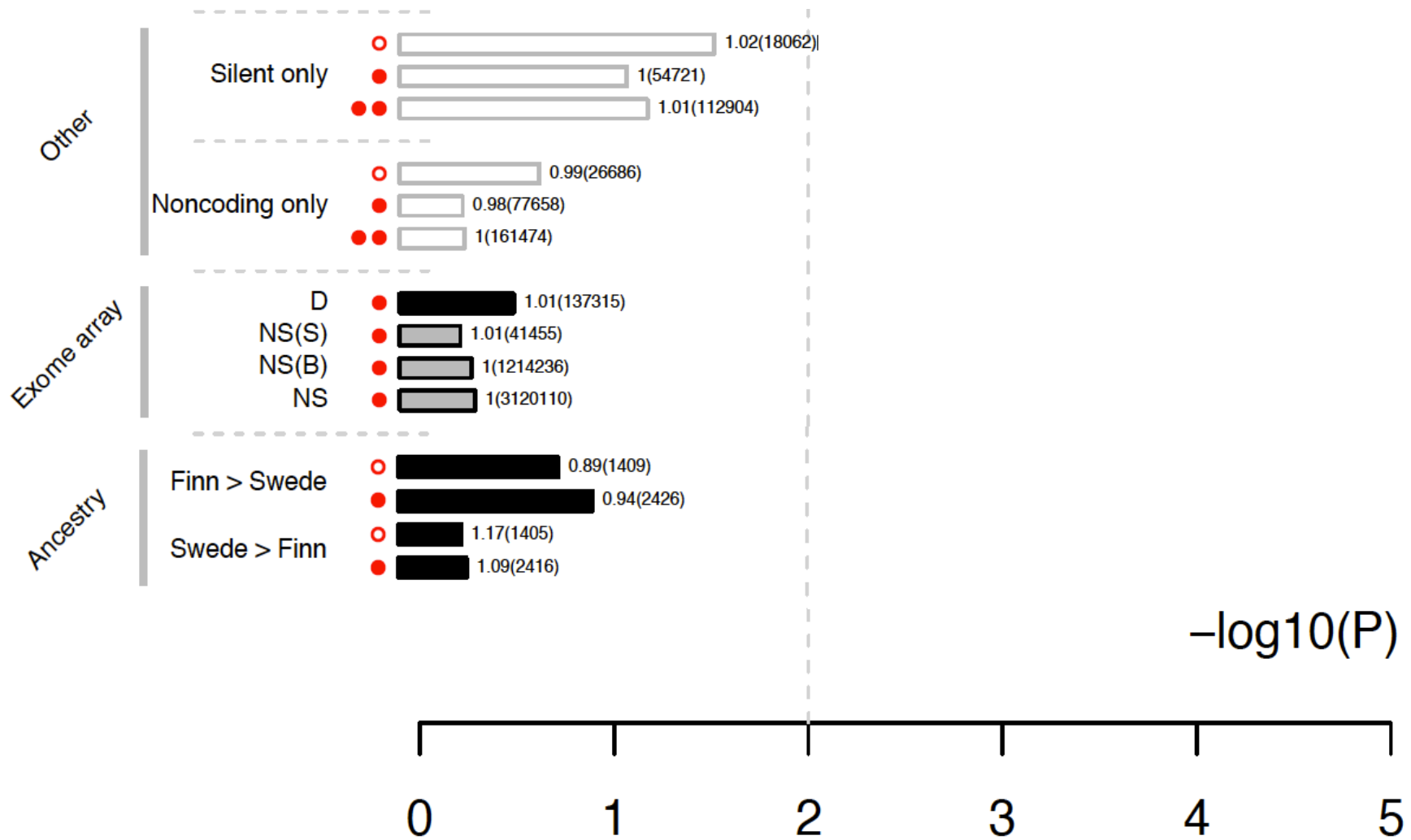


## Characterizing this increased burden: a “core” gene-set of 1858 genes



Nonsyn.



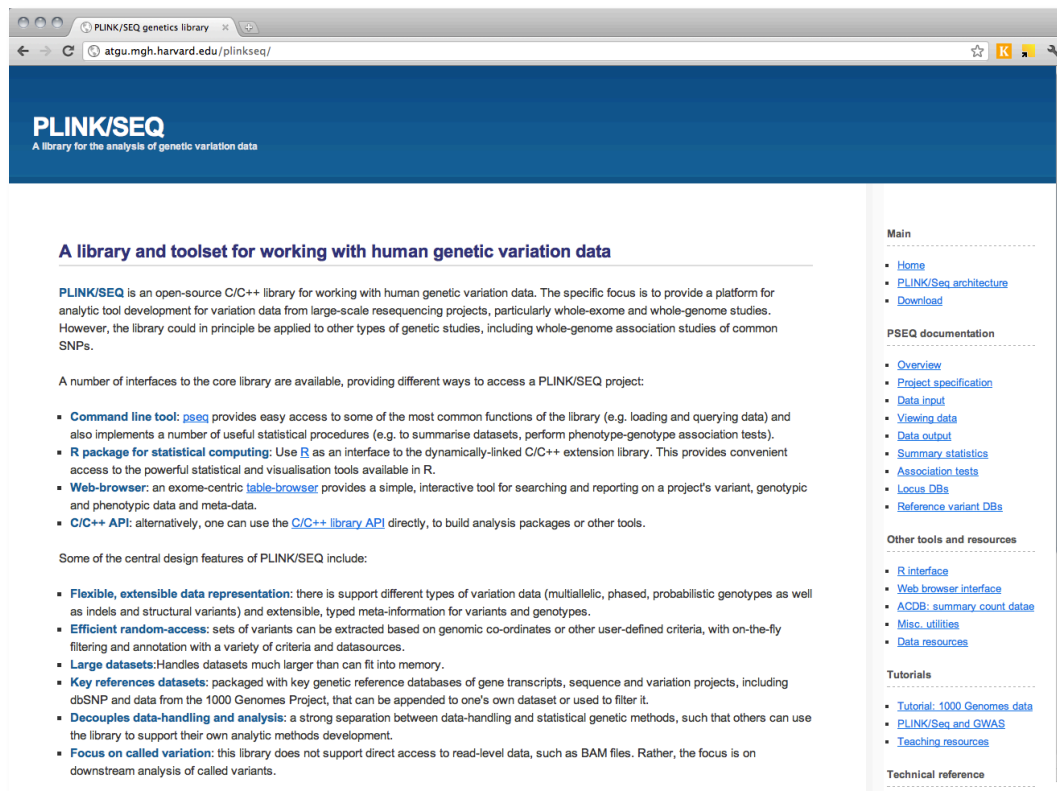


# Working with NGS data / VCFs

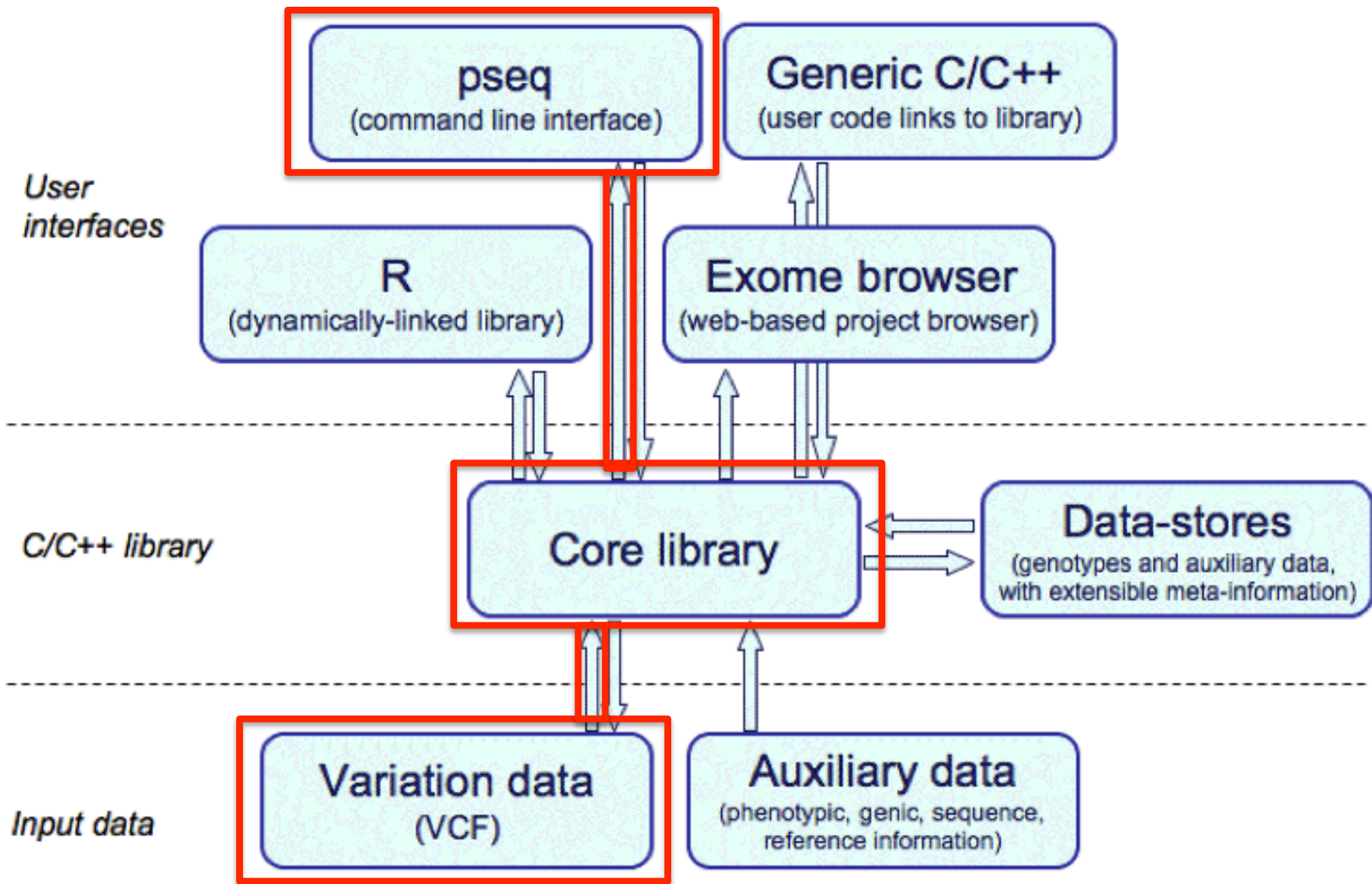
- “How do I...”
  - ... find variants in a specific gene in my dataset?
  - ... lookup a list of variants?
  - ... annotate a list of sites?
  - ... get summary QC metrics for regions, samples?

# PLINK/Seq : a toolset for NGS variation datasets

<http://atgu.mgh.harvard.edu/plinkseq/>

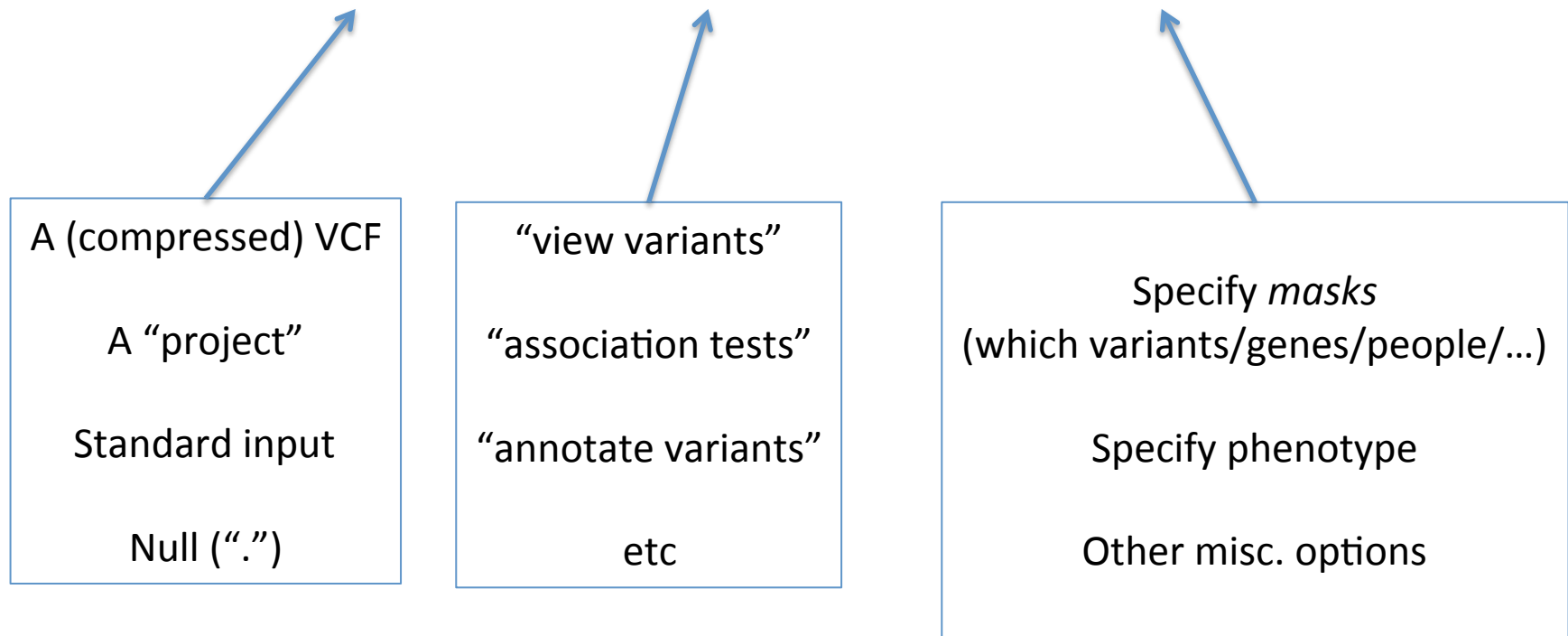


- VCF as primary input
- Focus on analysis of rare variation
- Extensible meta-information on locus, genotypes, individuals
- Bundled with key reference databases that can be directly intersected with one's own project
- Command-line and R library; web-based GUI under-development



## Basic structure of *PSEQ* commands

`pseq`   `data-source`   `command`   `{ --options }`



Note: other types of genotypic data can be incorporated into existing projects:

PLINK files (BED/BIM/FAM format)

"dosage" files, post imputation

BCF2 files

**Three main modes of operation :  
given various inclusion/exclusion masks, to iterate over all ...**

**1) Variants**

- e.g. viewing/filtering a VCF
- calculating various summary statistics

**2) Groups of variants (e.g. genes, all nonsense variants)**

- primarily gene-based (or set-based) association tests

**3) Individuals**

- per-individual statistics, burden of rare variants, etc



## Core project databases

VARDB	Variant information
INDDB	Individual phenotype/covariate information
LOCDB	Gene/transcript information (e.g. RefSeq, CCDS)
SEQDB	Human genome sequence (e.g. hg19)
REFDB	(Annotated) known variants (e.g. dbSNP, HGMD)
NETDB	Network information (e.g. PPI)
PROTDB	Protein domain/motif information (e.g. Pfam)
IBDDB	Pairwise identity-by-descent (IBD) information
SEGDB	Known variants ( <u>ex</u> clude novel)

## Examples of common *masks* (following `--mask`)

<code>reg=chr2</code>	Variants on chromosome 2
<code>reg=chr2:123456</code>	A specific variant(s) by position
<code>reg=chr2:1000000..20000000</code>	Variants by range
<code>reg=chr2,chr7</code>	List of ranges/positions
<code>reg=@pos.txt</code>	List of regions in file <code>pos.txt</code>
<code>id=rs123456</code>	Specific variant by ID (from VCF)
<code>novel</code>	Variants with no known ID
<code>novel.ex</code>	Known variants ( <u>ex</u> clude novel)
<code>var=coding</code>	Variants by <i>variant set</i> “coding” from VARDB
<code>loc=refseq</code>	Variants in <i>locus group</i> “refseq” from LOCDB

## Exome data used in this practical

Deep-coverage, whole-exome sequence data from the 1000 Genomes Project, Phase 1

GBR Great British, N=63

TSI Tuscans in Italy, N=60

LWK Luhya in Webuye, Kenya, N=27

## Overview of practical

Look at 1000 Genomes VCF files directly

Filter and output a new VCF

Look at resource files

Create a PLINK/SEQ “project”

Variant level summary statistics

Individual level summary statistics

Extracting meta-information

Experimenting with filters/masks

Gene-based summaries

Annotation

Single-site metrics

Comparing two VCFs

Ancestry inference and relatedness

Single-site association

Gene-based association

Using pbrowse and Rplinkseq interfaces

Geneset enrichment analyses

- Main script to follow:
- `~pshaun/2013/pseq/extra/commands.txt`