



gemini: a scalable, agile framework for mining genome variation

Uma Devi Paila and Aaron Quinlan

Department of Public Health Sciences

Center for Public Health Genomics, University of Virginia.



cphg.virginia.edu/quinlan/
github.com/arq5x/gemini



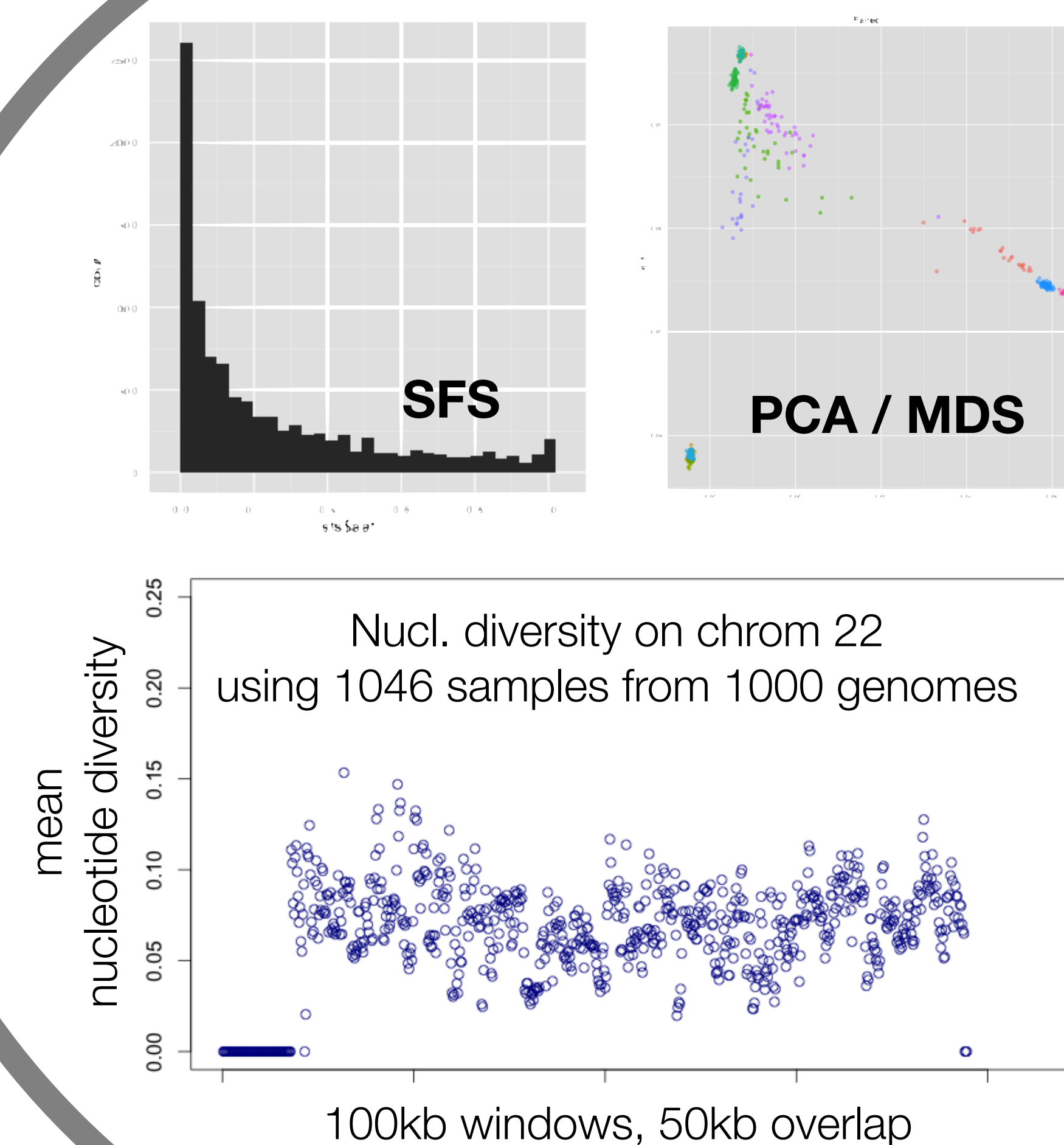
Enhanced SQL engine allows **selection** and **filtering** based on individual genotypes.

Scales to 1000s of samples.

```
SELECT chrom, start,
gene, impact, hwe,
in_dbsnp, in_omim,
gts.NA12878
FROM variants
WHERE is_lof = 1
AND aaf <= 0.01
AND call_rate \
>= 0.95
```

Ad hoc queries and filters using custom SQL engine

```
> gemini query -q [QUERY] my.db
```



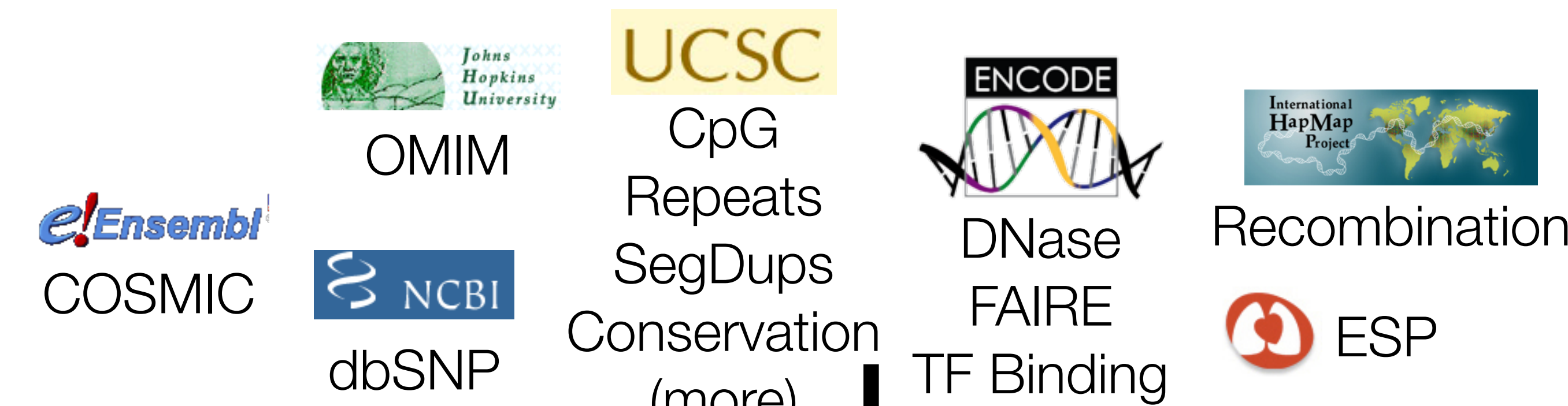
“Shortcuts” for common analyses, QC, visualization

```
> gemini stats --sfs my.db
> gemini stats --mds my.db
> gemini stats --tstv my.db
> gemini stats --vars-by-sample my.db
> gemini stats --gts-by-sample my.db
```

Motivation

- ★ Tools for population-scale studies of genome-wide variation require **efficient** analytical **frameworks**.
- ★ Existing approaches are ill-suited to the analysis of variation in the context of **genome annotations**.
- ★ Researchers need a simple, efficient, and flexible framework for large-scale data exploration.

Genome annotations place variants in context



snpEff
VEP
VAT*

VCF file

PED+ file

(sample relationships)

samples	sample_genotypes
pk	pk
sample_id	sample_id
name	name
family_id	family_id
maternal_id	maternal_id
sex	sex
phenotype	phenotype
ethnicity	ethnicity

variants	variant_impacts
pk	pk
variant_id	variant_id
chrom	chrom
start	start
end	end
ref	ref
alt	alt
qual	qual
filter	filter
type	type
gt_bases	gt_bases
gt_types	gt_types
gt_phases	gt_phases
call_rate	call_rate
in_dbsnp	in_dbsnp
in_omim	in_omim
in_ensembl	in_ensembl
in_ucsc	in_ucsc
in_cosmic	in_cosmic
in_ncbi	in_ncbi
in_repeats	in_repeats
in_segdup	in_segdup
in_conservation	in_conservation
in_encode	in_encode
in_dnase	in_dnase
in_faice	in_faice
in_tfbinding	in_tfbinding
in_recombination	in_recombination
in_esp	in_esp
is_lof	is_lof
is_lof_pred	is_lof_pred
is_lof_score	is_lof_score
is_lof_pval	is_lof_pval
is_lof_qval	is_lof_qval
is_lof_risk	is_lof_risk
is_lof_hwe	is_lof_hwe
is_lof_aaf	is_lof_aaf
is_lof_call_rate	is_lof_call_rate
is_lof_in_dbsnp	is_lof_in_dbsnp
is_lof_in_omim	is_lof_in_omim
is_lof_in_ensembl	is_lof_in_ensembl
is_lof_in_ucsc	is_lof_in_ucsc
is_lof_in_cosmic	is_lof_in_cosmic
is_lof_in_ncbi	is_lof_in_ncbi
is_lof_in_repeats	is_lof_in_repeats
is_lof_in_segdup	is_lof_in_segdup
is_lof_in_conservation	is_lof_in_conservation
is_lof_in_encode	is_lof_in_encode
is_lof_in_dnase	is_lof_in_dnase
is_lof_in_faice	is_lof_in_faice
is_lof_in_tfbinding	is_lof_in_tfbinding
is_lof_in_recombination	is_lof_in_recombination
is_lof_in_esp	is_lof_in_esp
is_lof_in_1000g	is_lof_in_1000g
is_lof_in_1000g_pred	is_lof_in_1000g_pred
is_lof_in_1000g_score	is_lof_in_1000g_score
is_lof_in_1000g_pval	is_lof_in_1000g_pval
is_lof_in_1000g_qval	is_lof_in_1000g_qval
is_lof_in_1000g_risk	is_lof_in_1000g_risk
is_lof_in_1000g_hwe	is_lof_in_1000g_hwe
is_lof_in_1000g_aaf	is_lof_in_1000g_aaf
is_lof_in_1000g_call_rate	is_lof_in_1000g_call_rate
is_lof_in_1000g_in_dbsnp	is_lof_in_1000g_in_dbsnp
is_lof_in_1000g_in_omim	is_lof_in_1000g_in_omim
is_lof_in_1000g_in_ensembl	is_lof_in_1000g_in_ensembl
is_lof_in_1000g_in_ucsc	is_lof_in_1000g_in_ucsc
is_lof_in_1000g_in_cosmic	is_lof_in_1000g_in_cosmic
is_lof_in_1000g_in_ncbi	is_lof_in_1000g_in_ncbi
is_lof_in_1000g_in_repeats	is_lof_in_1000g_in_repeats
is_lof_in_1000g_in_segdup	is_lof_in_1000g_in_segdup
is_lof_in_1000g_in_conservation	is_lof_in_1000g_in_conservation
is_lof_in_1000g_in_encode	is_lof_in_1000g_in_encode
is_lof_in_1000g_in_dnase	is_lof_in_1000g_in_dnase
is_lof_in_1000g_in_faice	is_lof_in_1000g_in_faice
is_lof_in_1000g_in_tfbinding	is_lof_in_1000g_in_tfbinding
is_lof_in_1000g_in_recombination	is_lof_in_1000g_in_recombination
is_lof_in_1000g_in_esp	is_lof_in_1000g_in_esp
is_lof_in_1000g_in_1000g	is_lof_in_1000g_in_1000g
is_lof_in_1000g_in_1000g_pred	is_lof_in_1000g_in_1000g_pred
is_lof_in_1000g_in_1000g_score	is_lof_in_1000g_in_1000g_score
is_lof_in_1000g_in_1000g_pval	is_lof_in_1000g_in_1000g_pval
is_lof_in_1000g_in_1000g_qval	is_lof_in_1000g_in_1000g_qval
is_lof_in_1000g_in_1000g_risk	is_lof_in_1000g_in_1000g_risk
is_lof_in_1000g_in_1000g_hwe	is_lof_in_1000g_in_1000g_hwe
is_lof_in_1000g_in_1000g_aaf	is_lof_in_1000g_in_1000g_aaf
is_lof_in_1000g_in_1000g_call_rate	is_lof_in_1000g_in_1000g_call_rate
is_lof_in_1000g_in_1000g_in_dbsnp	is_lof_in_1000g_in_1000g_in_dbsnp
is_lof_in_1000g_in_1000g_in_omim	is_lof_in_1000g_in_1000g_in_omim
is_lof_in_1000g_in_1000g_in_ensembl	is_lof_in_1000g_in_1000g_in_ensembl
is_lof_in_1000g_in_1000g_in_ucsc	is_lof_in_1000g_in_1000g_in_ucsc
is_lof_in_1000g_in_1000g_in_cosmic	is_lof_in_1000g_in_1000g_in_cosmic
is_lof_in_1000g_in_1000g_in_ncbi	is_lof_in_1000g_in_1000g_in_ncbi
is_lof_in_1000g_in_1000g_in_repeats	is_lof_in_1000g_in_1000g_in_repeats
is_lof_in_1000g_in_1000g_in_segdup	is_lof_in_1000g_in_1000g_in_segdup
is_lof_in_1000g_in_1000g_in_conservation	is_lof_in_1000g_in_1000g_in_conservation
is_lof_in_1000g_in_1000g_in_encode	is_lof_in_1000g_in_1000g_in_encode
is_lof_in_1000g_in_1000g_in_dnase	is_lof_in_1000g_in_1000g_in_dnase
is_lof_in_1000g_in_1000g_in_faice	is_lof_in_1000g_in_1000g_in_faice
is_lof_in_1000g_in_1000g_in_tfbinding	is_lof_in_1000g_in_1000g_in_tfbinding
is_lof_in_1000g_in_1000g_in_recombination	is_lof_in_1000g_in_1000g_in_recombination
is_lof_in_1000g_in_1000g_in_esp	is_lof_in_1000g_in_1000g_in_esp
is_lof_in_1000g_in_1000g_in_1000g	is_lof_in_1000g_in_1000g_in_1000g
is_lof_in_1000g_in_1000g_in_1000g_pred	is_lof_in_1000g_in_1000g_in_1000g_pred
is_lof_in_1000g_in_1000g_in_1000g_score	is_lof_in_1000g_in_1000g_in_1000g_score
is_lof_in_1000g_in_1000g_in_1000g_pval	is_lof_in_1000g_in_1000g_in_1000g_pval
is_lof_in_1000g_in_1000g_in_1000g_qval	is_lof_in_1000g_in_1000g_in_1000g_qval
is_lof_in_1000g_in_1000g_in_1000g_risk	is_lof_in_1000g_in_1000g_in_1000g_risk
is_lof_in_1000g_in_1000g_in_1000g_hwe	is_lof_in_1000g_in_1000g_in_1000g_hwe
is_lof_in_1000g_in_1000g_in_1000g_aaf	is_lof_in_1000g_in_1000g_in_1000g_aaf
is_lof_in_1000g_in_1000g_in_1000g_call_rate	is_lof_in_1000g_in_1000g_in_1000g_call_rate
is_lof_in_1000g_in_1000g_in_1000g_in_dbsnp	is_lof_in_1000g_in_1000g_in_1000g_in_dbsnp
is_lof_in_1000g_in_1000g_in_1000g_in_omim	is_lof_in_1000g_in_1000g_in_1000g_in_omim
is_lof_in_1000g_in_1000g_in_1000g_in_ensembl	is_lof_in_1000g_in_1000g_in_1000g_in_ensembl
is_lof_in_1000g_in_1000g_in_1000g_in_ucsc	is_lof_in_1000g_in_1000g_in_1000g_in_ucsc
is_lof_in_1000g_in_1000g_in_1000g_in_cosmic	is_lof_in_1000g_in_1000g_in_1000g_in_cosmic
is_lof_in_1000g_in_1000g_in_1000g_in_ncbi	is_lof_in_1000g_in_1000g_in_1000g_in_ncbi
is_lof_in_1000g_in_1000g_in_1000g_in_repeats	is_lof_in_1000g_in_1000g_in_1000g_in_repeats
is_lof_in_1000g_in_1000g_in_1000g_in_segdup	is_lof_in_1000g_in_1000g_in_1000g_in_segdup
is_lof_in_1000g_in_1000g_in_1000g_in_conservation	is_lof_in_1000g_in_1000g_in_1000g_in_conservation
is_lof_in_1000g_in_1000g_in_1000g_in_encode	is_lof_in_1000g_in_1000g_in_1000g_in_encode
is_lof_in_1000g_in_1000g_in_1000g_in_dnase	is_lof_in_1000g_in_1000g_in_1000g_in_dnase
is_lof_in_1000g_in_1000g_in_1000g_in_faice	is_lof_in_1000g_in_1000g_in_1000g_in_faice
is_lof_in_1000g_in_1000g_in_1000g_in_tfbinding	is_lof_in_1000g_in_1000g_in_1000g_in_tfbinding
is_lof_in_1000g_in_1000g_in_1000g_in_recombination	is_lof_in_1000g_in_1000g_in_1000g_in_recombination
is_lof_in_1000g_in_1000g_in_1000g_in_esp	is_lof_in_1000g_in_1000g_in_1000g_in_esp
is_lof_in_1000g_in_1000g_in_1000g_in_1000g	is_lof_in_1000g_in_1000g_in_1000g_in_1000g
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_pred	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_pred
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_score	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_score
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_pval	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_pval
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_qval	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_qval
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_risk	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_risk
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_hwe	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_hwe
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_aaf	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_aaf
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_call_rate	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_call_rate
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_dbsnp	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_dbsnp
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_omim	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_omim
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_ensembl	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_ensembl
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_ucsc	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_ucsc
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_cosmic	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_cosmic
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_ncbi	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_ncbi
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_repeats	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_repeats
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_segdup	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_segdup
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_conservation	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_conservation
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_encode	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_encode
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_dnase	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_dnase
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_faice	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_faice
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_tfbinding	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_tfbinding
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_recombination	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_recombination
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_esp	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_esp
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_pred	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_pred
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_score	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_score
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_pval	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_pval
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_qval	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_qval
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_risk	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_risk
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_hwe	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_hwe
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_aaf	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_aaf
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_call_rate	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_call_rate
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_dbsnp	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_dbsnp
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_omim	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_omim
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_ensembl	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_ensembl
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_ucsc	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_ucsc
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_cosmic	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_cosmic
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_ncbi	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_ncbi
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_repeats	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_repeats
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_segdup	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_segdup
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_conservation	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_conservation
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_encode	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_encode
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_dnase	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_dnase
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_faice	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_faice
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_tfbinding	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_tfbinding
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_recombination	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_recombination
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_esp	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_esp
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_pred	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_pred
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_score	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_score
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_pval	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_pval
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_qval	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_qval
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_risk	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_risk
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_hwe	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_hwe
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_aaf	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_aaf
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_call_rate	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_call_rate
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_dbsnp	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_dbsnp
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_omim	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_omim
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_ensembl	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_ensembl
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_ucsc	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_ucsc
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_cosmic	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_cosmic
is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_ncbi	is_lof_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_1000g_in_ncbi