# Using HaploReg and RegulomeDB to mine ENCODE data:

(Updated 17 May 2013, Mike Pazin)

HaploReg and RegulomeDB are ENCODE-funded tools described in recent publications that retrieve ENCODE annotations at SNPs of interest, as well as annotations from work by other researchers and projects.

## HaploReg v2:

Go to the HaploReg site, and enter the name of the SNP of interest (Arrow 1).
(Using the "Set Options" tab, the user can configure values such as the LD threshold and the population used from 1000 Genomes data used to calculate LD.)
Click on the submit button (Arrow 2)



HaploReg retrieves the ENCODE and Roadmap Epigenomics annotations for the selected SNP, as well as other SNPs in LD (arrow 3).

**RegulomeDB**:
Go to the RegulomeDB site and enter the name of the SNP of interest (Arrow 1).



Click on the submit button (Arrow 2).
RegulomeDB calculates a score for the regulatory potential of this region.



Clicking on the score (arrow 3) retrieves the ENCODE (and other) annotation for the region, including transcription factor binding, chromatin structure (DNase, FAIRE, and histone modifications), transcription factor motifs and eQTL.

**RegulomeDB Disease Association Database**, a database of predicted functional SNPs, organized by disease/trait and by SNP, is available at:
http://regulome.stanford.edu/GWAS

There is a list of over 4700 SNPs associated with human traits and disease (arrow 1), as well as a list of over 470 human traits and diseases (arrow 2).



Clicking on a trait/disease returns a list of SNPs that have been associated with that trait or disease:



Clicking on a SNP (red arrow) returns the evidence for the association:



As well as the annotation for the lead SNP, and other SNPs in LD that, based on functional annotation, are candidates for the functional variant:

**Lead SNP**
rs3024505
**Position:** chr1:206,939,904 (Open in UCSC Genome Browser)
**Distance to nearest TSS:** 18,466 bp
**GENCODE v7 location:** Intergenic region
**RegulomeDB Score:** 2b - ChIP-seq peak + any motif + matched DNase Footprint + DNaseI-seq peak (Open in RegulomeDB)

**Linkage disequilibrium region**
Linkage disequilibrium threshold:
- In all HapMap 2 populations: $r^2 \geq 0.8$ $r^2 \geq 0.9$ $r^2 = 1.0$
- In the HapMap 2 CEU population $r^2 \geq 0.8$ $r^2 \geq 0.9$ $r^2 = 1.0$
**SNPs in the linkage disequilibrium region sorted by decreasing amount of evidence supporting a functional role for the SNP:**
rs3024493
**Position:** chr1:206,943,968 (Open in UCSC Genome Browser)
**Distance to lead SNP:** 4,064 bp
**Distance to nearest TSS:** 22,530 bp
**GENCODE v7 location:** Intron
**RegulomeDB Score:** 2b - ChIP-seq peak + any motif + matched DNase Footprint + DNaseI-seq peak (Open in RegulomeDB)
**Linkage disequilibrium with Lead SNP (HapMap 2):** CEU: D'=1.0, r2=1.0 / CHB: D'=1.0, r2=1.0 / JPT: D'=1.0, r2=1.0 / YRI: D'=1.0, r2=1.0
rs3024495
**Position:** chr1:206,942,413 (Open in UCSC Genome Browser)
**Distance to lead SNP:** 2,509 bp
**Distance to nearest TSS:** 20,975 bp
**GENCODE v7 location:** Intron
**RegulomeDB Score:** 5a - ChIP-seq peak (Open in RegulomeDB)
**Linkage disequilibrium with Lead SNP (HapMap 2):** CEU: D'=1.0, r2=1.0 / CHB: D'=1.0, r2=1.0 / JPT: D'=1.0, r2=1.0 / YRI: D'=1.0, r2=1.0

One can follow the links to view the genomic annotation of these SNPs in the genome browser.