

International HapMap Project

Home | About the Project | Data | Publications | Conference

FR | English | Français | 日本語 | Yoruba

The International HapMap Project is a partnership of scientists and funding agencies from Canada, China, Japan, Nigeria, the United Kingdom and the United States to develop a public resource that will help researchers find genes associated with human disease and response to pharmaceuticals. See "About the International HapMap Project" for more information.

Project Information

About the Project

HapMap Publications

HapMap Conference

HapMap Meeting List

HapMap Project Participants

HapMap Mirror Site in Japan

News

- 2005-03-01: HapMap public release #16. ATTN: This the so-called Phase I data freeze which marks a major milestone of the project: a genotyped common SNP every 500 kb in all populations under study. Data available for bulk download and graphical browsing. Summary of genotyped SNPs:

Populations	CEU	HCB	JPT	YRI
Genotyped SNPs	1,077,693	1,044,998	1,044,418	1,034,205

- 2005-02-09: HapMap News Volume 1, 2004. This is the first in a series of newsletters to be published by the Coriell Institute for Medical Research to inform communities how their samples are being used. Each issue of the newsletter will be available in the primary language of all the participating communities.
- 2005-02-07: International HapMap Consortium Expands Mapping Effort. The International HapMap Consortium, boosted by an additional 3.3 million in public-private support, announces plans to create an even more powerful map of human genetic variation than originally envisioned. The map will accelerate the discovery of genes related to common diseases, such as asthma, cancer, diabetes and heart disease.
- Old News

www.hapmap.org

Finding SNPs: HapMap Browser

Finding SNPs: HapMap Browser

SNP info: rs1143643 with alleles A/G in dbSNP (dbSNP report | Ensembl SNPView)

Chromsome location: Chr2:113683448..113683448 (-) strand relative to the human reference sequence

Frequency report:

Population	Ref homozygote genotype freq	Heterozygote count	Other homozygote genotype freq	Total genotype freq. count		Allele frequencies		Total allele freq. count	Total count									
				Ref-allele	Other-allele	Ref-allele	Other-allele											
CEU	G/G	0.387	22	A/G	0.483	29	A/A	0.150	9	60	G	0.608	73	A	0.392	47	120	retrieve genotypes
HCB	G/G	0.222	10	A/G	0.533	24	A/A	0.244	11	45	G	0.489	44	A	0.511	46	90	retrieve genotypes
JPT	G/G	0.227	10	A/G	0.523	23	A/A	0.250	11	44	G	0.489	43	A	0.511	45	88	retrieve genotypes
YRI	G/G	0.700	42	A/G	0.300	18	A/A	0.000	0	60	G	0.800	102	A	0.150	18	120	retrieve genotypes

Note: the reference allele is the base observed in the reference genome sequence at this location.

Population descriptions:
 CEU: CEH (Utah residents with ancestry from northern and western Europe)
 HCB: Han Chinese in Beijing, China
 JPT: Japanese in Tokyo, Japan
 YRI: Yoruba in Ibadan, Nigeria

Please see this page for more information about the populations, as well as a general discussion of the populations under study in the project.

Finding SNPs: HapMap Genotypes

```

#chr2:113683448..113683448 (-) strand relative to the human reference sequence
#file: /data/hapmap/genotype/chr2/113683448..113683448 (-) strand relative to the human reference sequence
#format: chr2:113683448..113683448 (-) strand relative to the human reference sequence
#info: rs1143643 with alleles A/G in dbSNP (dbSNP report | Ensembl SNPView)
#populations: CEU HCB JPT YRI
#genotypes:
CEU 113683448..113683448 (-) strand relative to the human reference sequence
HCB 113683448..113683448 (-) strand relative to the human reference sequence
JPT 113683448..113683448 (-) strand relative to the human reference sequence
YRI 113683448..113683448 (-) strand relative to the human reference sequence

```

Bulk data downloads

The following directories contain project data that have been made publicly available. (See [HapMap Data Access Policy](#) for more information). More details about each dataset can be found in READMEs in the respective directories:

- Genotypes: Individual genotype data submitted to the DCC to date.
- LD Data: Linkage disequilibrium properties D', LOD, R² compiled from the genotype data to date.
- Allocated SNPs: dbSNP reference SNP clusters that have been picked and prioritized for genotyping according to several criteria (see info on how SNPs were selected). The file DOREADME contains per-chromosome SNP counts and further details.
- Frequencies: Allele & genotype frequencies compiled from genotyping data submitted to the DCC to date. These have also been submitted to dbSNP and should be available in the next dbSNP build.
- SNP assays: Details about assays submitted to the DCC to date. PCR primers, extension probes etc., specific to each genotyping platform.
- Protocols: Information on assay design, genotyping and other protocols used in the project.
- Samples/Individuals: Information on the samples used in the project and the individuals from which they were drawn. (See About the project: Which Populations are Being Sampled).
- XML docs: Documentation on the XML format used in the project.

Finding SNPs: HapMap Browser

Minimal SNP information for genotyping/characterization

- What is the SNP? Flanking sequence and alleles.
 - FASTA format
 - >snp_name
 - ACCGAGTACCCAG
 - [A/G]
 - ACTGGGATAGAAC
- dbSNP reference SNP # (rs #)
- Where is the SNP mapped? Exon, promoter, UTR, etc.
 - picture of gene with mapped to the gene structure.
- How was it discovered? Method
- What assurances do you have that it is real? Validated how?
- What population African, European, etc?
- What is the allele frequency of each SNP? Common (>10%), rare
- Are other SNPs associated? Genotyping data!

Finding SNPs: HapMap Browser

1. HapMap datasets are useful because individual genotype data can be used to determine optimal genotyping strategies (tagSNPs) or perform population genetic analyses (linkage disequilibrium)
2. Data are specific produced by those projects (not all dbSNP)
 - ✓ HapMap data is available in dbSNP
3. HapMap data (Phase II) can be accessed released prior to dbSNPs
4. Easier visualization of data and direct access to SNP data, individual genotypes, and LD analysis

Finding SNPs: Databases and Extraction

How do I find and download SNP data for analysis/genotyping?

1. Entrez Gene
 - dbSNP
 - Entrez SNP
2. HapMap Genome Browser
3. NIEHS Environmental Genome Project (EGP) Candidate gene website
4. NIEHS web applications and other tools
 - GeneSNPs, PolyDoms, TraFac, PolyPhen, ECR Browser, GVS

Finding SNPs: NIEHS SNPs Candidate Genes

Haplotyping Data PHASE Output Phased Individual Haplotypes Sorted by Frequency

Visual Haplotype

rsid	rsid	rsid	rsid
1	1	1	1
2	2	2	2
3	3	3	3
4	4	4	4
5	5	5	5
6	6	6	6
7	7	7	7
8	8	8	8
9	9	9	9
10	10	10	10
11	11	11	11
12	12	12	12
13	13	13	13
14	14	14	14
15	15	15	15
16	16	16	16
17	17	17	17
18	18	18	18
19	19	19	19
20	20	20	20
21	21	21	21
22	22	22	22
23	23	23	23
24	24	24	24
25	25	25	25
26	26	26	26
27	27	27	27
28	28	28	28
29	29	29	29
30	30	30	30
31	31	31	31
32	32	32	32
33	33	33	33
34	34	34	34
35	35	35	35
36	36	36	36

LD Linkage Data LD Select (Tag SNPs)

African Descent European Descent Hispanic Descent Asian Descent

Bin	total_sites	average_misorder_allele_frequency
Bin 1	13	294
Bin 2	4	394
Bin 3	4	234
Bin 4	3	334
Bin 5	3	84
Bin 6	2	494
Bin 7	2	264

Predictive Analyses Nonsynonymous cSNP Analysis

rsid	AllelePop	a.a. Pos	Residue	Variant	Feat	PPHpredict	SIFTpredict
rsid1	031821	275	R	G	0.01	benign	TOLEPERATED
rsid1	031875	356	O	R	0.01	possibly damaging	INTOLERANT
rsid1	032885	693	D	N	0.05	benign	TOLEPERATED
rsid1	033426	871	P	L	0.41	benign	TOLEPERATED
rsid1	033921	1038	E	G	0.32	benign	INTOLERANT
rsid1	033927	1040	S	N	0.01	benign	TOLEPERATED
rsid1	034226	1140	S	G	0.02	benign	TOLEPERATED
rsid1	034356	1183	R	G	0.33	benign	TOLEPERATED
rsid1	055298	1613	S	G	0.33	benign	TOLEPERATED
rsid1	055319	1620	T	A	0.01	benign	TOLEPERATED

SIFT = Sorting Intolerant From Tolerant
Evolutionary comparison of non-synonymous SNPs

PolyPhen - Polymorphism Phenotyping
Structural protein characteristics and evolutionary comparison

Finding SNPs: NIEHS SNPs Candidate Genes

Download a zip file of all data for this gene

Category	Sub-category	Color FASTA SNP Context	PCR Primers (FASTA) Genbank
Mapping Data	cSNPs cDNA	Color FASTA SNP Context	PCR Primers (FASTA) Genbank
Genotyping Data	Visual Genotype Individual Genotypes	SNP Alleles SNP Allele Frequency	SNP Hardy-Weinberg
Haplotyping Data	PHASE Output Visual Haplotype	Phased individual Haplotypes	Sorted by Frequency
Linkage Data	LD Select (Tag SNPs)	African Descent European Descent Hispanic Descent Asian Descent	
Predictive Analyses	Nonsynonymous cSNP Analysis		

Finding SNPs: NIEHS SNPs Candidate Genes

National Institute of Environmental Health Sciences Environmental Genome Project

NIEHS SNPs

Welcome to the NIEHS SNPs Program

Introduction

The NIEHS Environmental Genome Project is a multi-disciplinary, collaborative effort focused on exploring the relationships between environmental exposures and human disease. The NIEHS SNPs Program is a sub-project of this effort, designed to identify common sequence variants (SNPs) in human genes involved in DNA repair and cell cycle pathways (see the article "Genetics in the repair and maintenance of the cell"). Identifying the specific cell-protein genetic map of human genes that can be applied in evaluating human disease risk with environmental exposures.

Genes/SNPs Database

NIEHS SNPs are available in the Genes/SNPs database, as well as the national database resources, dbSNP, GeneSNP, and the dbSNP. This database provides a list of all available SNP data in human genes and their associated phenotypes in select populations. This information is key in selecting the polymorphic sites needed to examine disease risk in human population studies.

Polymorphism Analysis

Automated DNA sequencing is being used to identify and genotype SNPs in human candidate genes (see PolyPhen). Candidate genes are being presented to identify common sequence variants for functional analysis and population-based studies. Candidate genes were formerly organized across a gene list.

egp.gs.washington.edu

Finding SNPs: NIEHS SNPs Candidate Genes

Data Downloads

All NIEHS SNPs Variation Data

Full Download of All Variation Data Files

WARNING: This is a large file and will take several minutes to download. The file is compressed and "gzipped" and the containing the entire directory of text files. Please use the same text data files which appear in the data pages for each candidate gene in our Project Gene Table. Please use our Usage Policy if this data is to be used in a publication.

Download of Variation Data (Single File)

Global Predefined File

This is a tab-delimited text file in our "bedtools" format that describes all SNPs sites discovered by NIEHS SNPs. The format of this file is:

Line format: chromosome:position:chromosome:HUGO_NAME <+ECG Sample ID> <Allele1 <Allele2>

Example: chromosome:1:100000000:BRCA1 <+ECG Sample ID> C >G

The chromosome position is generated from mapping to the most recent genome assembly available from the USCS Genome Browser.

Download of EGP Variation Data by Chromosome

These are tab-delimited text files in our "bedtools" format that describe all SNPs sites discovered by NIEHS SNPs but repeated into files based on chromosome. The format of this file is:

NIEHS SNPs SIFT/PolyPhen Data

Functional changes in a candidate gene's protein function were assessed by taking the nonsynonymous coding SNPs (cSNPs) for each gene and using SIFT and PolyPhen. Generally, each nonsynonymous amino acid change is included in the context of other available similar proteins to determine the likelihood of the polymorphic nonsynonymous change and then statistically classified. These programs classify each coding SNP as being either tolerated (SIFT) or as being possibly damaging or probably damaging (PolyPhen).

Continued SIFT/PolyPhen Data for NIEHS SNPs Nonsynonymous SNPs (Tolerant or Possibly Damaging)

Continued SIFT/PolyPhen Data for NIEHS SNPs Nonsynonymous SNPs (Probably Damaging)

SIFT Data for NIEHS SNPs Nonsynonymous SNPs (Probably Damaging)

PolyPhen Data for NIEHS SNPs Nonsynonymous SNPs (Probably Damaging)

Finding SNPs: Databases and Extraction

How do I find and download SNP data for analysis/genotyping?

1. Entrez Gene
 - dbSNP
 - Entrez SNP
2. HapMap Genome Browser
3. NIEHS Environmental Genome Project (EGP)
Candidate gene website
4. NIEHS web applications and other tools
GeneSNPs, PolyDoms, TraFac, PolyPhen,
ECR Browser, GVS

GeneSNPs

Graphic view of SNPs in context of gene elements
All NIEHS genes presented

- organized by pathway/function
SNPs from dbSNP

- organized by submitter handle
Sequenc context of SNPs presented in Color Fasta
format

Link-outs to EntrezSNP pages
Summary Genome SNPs internal SNP viewer for
one-stop SNP shopping

<http://www.genome.utah.edu/genesnps/>

GeneSNPs: One stop shopping

Gene	SNP ID	Position	Alleles	Frequency	Submitter
ADH1B	rs1044396	123456	A/G	0.8	dbSNP
ADH1B	rs1044397	123457	C/T	0.7	dbSNP
ADH1B	rs1044398	123458	G/A	0.6	dbSNP
ADH1B	rs1044399	123459	T/C	0.5	dbSNP
ADH1B	rs1044400	123460	A/G	0.4	dbSNP
ADH1B	rs1044401	123461	C/T	0.3	dbSNP
ADH1B	rs1044402	123462	G/A	0.2	dbSNP
ADH1B	rs1044403	123463	T/C	0.1	dbSNP
ADH1B	rs1044404	123464	A/G	0.05	dbSNP
ADH1B	rs1044405	123465	C/T	0.02	dbSNP

GeneSNPs: One stop shopping

Gene	SNP ID	Position	Alleles	Frequency	Submitter
ADH1B	rs1044396	123456	A/G	0.8	dbSNP

GeneSNPs: One stop shopping

Gene	SNP ID	Position	Alleles	Frequency	Submitter
ADH1B	rs1044396	123456	A/G	0.8	dbSNP

Polydoms

A web-based application that maps synonymous and non-synonymous SNPs onto known functional protein domains

- SNPs are from dbSNP and GeneSNPs
- Domain structures from NCBI's Conserved Domain Database
- Functional predictions based on SIFT and PolyPhen
- 3 dimensional mapping of SNPs on protein structure using Chime viewer

<http://polydoms.cchmc.org/polydoms>

Polydoms

Polydoms

Mapping of nsSNPs onto protein structure

TraFac: Transcription Factor Binding Site Comparison

A tool for validating cis regulatory elements conserved between human and mouse

- Aligns human and mouse sequences using BLAST%
- Consensus transcription factor binding sequences from Transfac database

<http://trafac.cchmc.org/trafac>

TraFac: Transcription Factor Binding Site Comparison

All TFBS in common

TFBS in parallel

ECR Browser: Evolutionary Conserved Regions

Aligns sequences to Mouse, Rat, Dog, Opposum, Chicken, Fugu and Drosophila

Gene annotations from UCSC Genome Browser

Easy retrieval of ECR sequences and alignments

Pre-computed transcription factor binding sites

<http://ecrbrowser.dcode.org>

ECR Browser: Evolutionary Conserved Regions

ECR Browser on Human (hg17) Settings

Display alignments with:

- Mouse (mm8)
- Rat (rn3)
- Opposum (morfDom1)
- Dog (canfam1)
- Chicken (gal2)
- Frog (ata)
- Fugu (fugu)
- Zebrafish 4.0 (zfa4)
- Tetraodon (toc)

Graph type: Smooth Pip-plot

Display gene features:

- mRNA
- Ensemble
- Known Genes from UCSC
- RefSeq

Number of layers: 100

Layer height: 100

Detect ECRs (Evolutionary Conserved Regions):

- min length: 100 bp
- min identity: 70 %

Coordinates: Relative Absolute

Show conserved TF binding sites (Optimized for function; TFBS search; TFBS visualization is limited to 20kb regions):

- select to show sites in alignments with:
 - Frog [ata]
 - Chicken [gal2]
 - Mouse [mm8]
 - Mouse [mm5]
 - Frog [ata]

GVS: Genome Variation Server beta Sponsored by SeattleSNPs

Display Results

Gene Name: **NDR2A**
 Allele Frequency Output (v2) 3
 Population: **POP_10000-PAN01, Submitter: EGP_3KBP3**

SNP ID	Ref Allele	Alt Allele	Frequency (%)	MAF	rsID	Gene	Function
2118499	TTCTTTG	G	29	0.46	3355	NDR2A	missense
2118813	TTGGATG	T	28	0.41	690	NDR2A	missense
2118635	TTTCTTT	A	13	0.23	623	NDR2A	missense
2118687	TTGGATG	A	13	0.23	632	NDR2A	missense
2118755	TTGGATG	T	13	0.23	622	NDR2A	missense
2118755	TTGGATG	T	13	0.23	622	NDR2A	missense
2111440	TTTCTTT	T	36	0.46	636	NDR2A	missense
2111572	TTTCTTT	G	41	0.49	116	NDR2A	missense
2111294	TTTCTTT	A	28	0.39	616	NDR2A	missense
2112359	TTTCTTT	G	12	0.23	623	NDR2A	missense
2112362	TTTCTTT	A	25	0.36	615	NDR2A	missense
2112362	TTTCTTT	A	17	0.29	647	NDR2A	missense
2112331	TTTCTTT	A	13	0.23	622	NDR2A	missense
2112331	TTTCTTT	A	17	0.29	647	NDR2A	missense
2112354	TTTCTTT	A	25	0.36	133	NDR2A	missense
2112355	TTTCTTT	A	38	0.47	426	NDR2A	missense
2112379	TTTCTTT	A	13	0.23	623	NDR2A	missense
2112392	TTTCTTT	G	42	0.44	146	NDR2A	missense
2112401	TTTCTTT	G	38	0.47	671	NDR2A	missense
2112401	TTTCTTT	G	24	0.41	630	NDR2A	missense
2112437	TTTCTTT	G	38	0.47	671	NDR2A	missense
2112470	TTTCTTT	G	28	0.41	630	NDR2A	missense
2112504	TTTCTTT	T	1	0.16	639	NDR2A	missense
2112470	TTTCTTT	G	13	0.23	623	NDR2A	missense
2112470	TTTCTTT	G	13	0.23	623	NDR2A	missense

Finding SNPs: Databases and Extraction

- One stop shopping
 - NIEHS SNPs and GeneSNPs
- Prediction of functional variations
 - Polydoms and PolyPhen
- Identification of transcription factor binding sites in Evolutionary Conserved Regions
 - TraFac and the ECR browser
- Visualization and analysis of LD and TagSNPs
 - GVS