# Studying Genetic Variation II: Computational Techniques

Jim Mullikin, PhD
Genome Technology Branch
NHGRI

# The Human Genome Project

The reference sequence describes just one copy of the genome

But everyone has two copies

## What makes us all different?

- Largely due to the differences in our genome sequence.

- In fact, even the two copies of the genome in our cells differ.

- Between any two unrelated genomes, there is about 1 difference every 1000 bases.

## Overview of Topics

- Genome variation origins
- Types of polymorphisms
- Discovery methods
- Access to genetic variation data
- How to find SNPs in a region of interest
- Haplotype Map project

## *Genome variation origins*

- Mutations are fundamentally produced by errors in DNA replication.

- DNA is replicated in the production of the egg and sperm cells.

- Thus, a child does not receive exact copies of information from mother and father.

## *Types of polymorphisms*

- Single Nucleotide Polymorphisms (SNPs) are single base changes and occur at a rate of about 30 - 60 sites per genome per generation.

ACTCCTCT**T**ATCCCTGC
ACTCCTCT**C**ATCCCTGC

ACTCCTCT[C/T]ATCCCTGC

## *Types of polymorphisms*

- Short Tandem Repeats (STRs) are specific repeated segments of sequence.

GGTTTTTGCC------TATATATATAAGTAGGA
GGTTTTTGCC----TATATATATATAAGTAGGA
GGTTTTTGCC--TATATATATATATAAGTAGGA
GGTTTTTGCCTATATATATATATATAAGTAGGA

TTGCC[(TA)5/(TA)6/(TA)7/(TA)8]AGT

## *STRs continued*

- These types of polymorphisms are used by the FBI for forensic testing.
- The 13 STRs that are used are all tetrameric, e.g., D7S280 is "gata" repeated 6-15 times.

## *STRs continued*

- These sites are especially variable in the human genome.

- From the 13 sites used by the FBI for DNA fingerprinting there are more possible combinations than the number of people on earth by a factor of one million.

## *Types of polymorphisms*

- Deletion/Insertion Polymorphisms (DIPs) are deletions or insertions of 1 base to as large as a few kilobases.

```
CATAAAAAAAGAACAAAATC
CATAAAAAAA-AACAAAATC

CATAAAAAAA[G/-]AACAAAATC
```

## Beyond polymorphisms

- When a mutational event is sufficiently large, these events are classified as chromosomal rearrangements.

- There are many examples of these as seen in karyotypes.

- These larger scale rearrangements, duplications or deletions are often associated with various diseases and severe abnormalities.

## Discovery methods

- The primary method for discovering polymorphisms is by sequencing DNA and comparing the sequences.
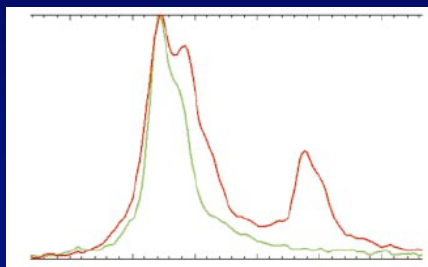
Trace Archive Search query:
trace_name='S213P602260RC9.T0' OR trace_name='50766946'

# *Discovery methods*

- There are other ways to find out if two or more DNAs differ, for example:
  - single strand conformational polymorphism (SSCP)
  - denaturing high performance liquid chromatography (DHPLC)

- These methods do not give specific sequence changes, however, they can be used for rapid mutation screening.

Hecker KH, et al. Anal Biochem. 1999 Aug 1;272(2):156-64.

# BRAF oncogene

- The Sanger Institute's Cancer Genome Project used the heteroduplex variation screening method to discover a missense mutation in the BRAF gene that is present in 66% of malignant melanomas.



Davies H, et al. Nature. 2002 Jun 27;417(6892):949-54

## Mining SNPs from sequence

- EST mining
- Clone overlap
- The SNP Consortium (TSC)
- Targeted resequencing
- Haplotype Map Project (HapMap)
- Other

## Expressed Sequence Tag Mining

- These sequences are primarily associated with coding regions of genes.

- By clustering these sequences, selected differences are identified as SNPs.

- There are over 100,000 SNPs in dbSNP from a variety of species detected from clustered ESTs.

- The following example is from the CGAP SNP project (see refs).

# *Clone Overlap*

- The human genome was sequenced from BAC clones (containing about 150kb of sequence each).

- These overlapped to various levels, and within the overlap regions, high quality base differences indicated the position and alleles of SNPs.

# *Clone Overlap*

- About 1.3M SNPs in dbSNP come from mining of clone overlaps.

- Special care was required to insure that the overlapping clones came from different haploids. (see references)

- This can be accomplished by looking at the source DNA for the two clones to see that it originated from different individuals, or if from the same individual, that the variation rate within the overlapping regions indicated that the DNA was from different haploids of one individual.

# *The SNP Consortium*

- A two year effort funded by the Wellcome Trust and 11 pharmaceutical and technological companies to discover 300,000 SNPs randomly distributed across the human genome.

- At its initiation in April 1999, the genome was only 10% finished and 20% in draft form.

- The SNPs were developed from a pool of DNA samples obtained from 24 individuals representing several racial groups.

## *The SNP Consortium*

- With the rapid increase in genome coverage from the public Human Genome Project, the strategies changed to take full advantage of the draft and finished sequence.

- The initial target of 300,000 SNP was passed quickly, and now the sequence generated from that project contributes over 1.3M SNPs to the public archives.

## *More SNPs for HapMap Project*

- This project required many more SNPs than were available when it started in October 2002, which totaled about 2M.

- Additional random shotgun sequencing has brought this to 4.8M SNPs today.

- Plans are to bring this to 6M SNPs by February 2004.

## *Targeted Resequencing*

- Any region of the genome can be targeted for resequencing.  From the finished sequence, PCR primers can be designed to amplify a target followed by sequencing.

- This method generally works from a 1:1 mixture of an individuals two haploids, so the special case of heterozygous base positions must be properly processed.

IMS-JST096911

http://snp.ims.u-tokyo.ac.jp/

Chr 19     PTGER1     gcC/gcT    A/A

# *Targeted Resequencing*

- JSNP database contains 190,562 SNPs detected from resequencing genomic regions containing genes in DNA from 24 Japanese individuals.

- Many groups use this technique for either SNP discovery in their region of interest, or as a way to validate SNPs.

- PolyPhred (see web links) is commonly used for analyzing resequencing traces.

SNP detection by PolyPhred. View of a Consed window with a tag (red=highest ranking SNP tag) marking the consensus position of the SNP in the traces and genotype tags marking each of the samples below (purple=homozygote, pink=heterozygote). On the right trace windows for alternate homozygoes (C/C (top) and G/G (bottom>> and a heterozygoe (C/G) middle).

PolyPhred example from their web site.

# Sequencing Chips



...GCTC**C**GTTT...
...GCTC**T**GTTT...

The Sanger Institute

Perlegen used Affymetrix's chip design process to place 60M probes on a 5x5" chip. From 20 single haploid chromosome 21 chromosomes, they discovered 36k SNPs.



# *Distribution properties*

- EST mining
  - Locates SNPs primarily within coding regions.

- Clone overlap
  - High density of SNPs within overlap regions, absent elsewhere.

- The SNP Consortium (TSC)
  - Randomly distributed across the genome, however, total sequence only covers 50% of the genome

## *Distribution properties*

- Haplotype Map Project (HapMap)
  - Random, like TSC, for first phase that reached 1X coverage
  - Chromosome sorted phase increased coverage from 1X-6X

- Targeted resequencing
  - Focused discovery that has been applied to 100s of individuals

- Chip based resequencing
  - Repetitive elements in the genome are masked

## *Quality of SNPs*

- The SNPs discovered for the TSC and HapMap projects use a method designed to give no more than 5% false positive (FP) SNPs.

- Two recent studies have looked at the quality of SNPs present in dbSNP (see references)
  - One study (Reich, et al., 2003) confirmed these minimum FP rates were achieved.
  - It goes on to show that SNPs with both alleles represented twice in different DNAs can eliminate the FPs.
  - The other study (Carlson, et al. 2003) showed a much lower validation rate, implying either a higher FP rate or that these SNPs were not present in their DNA samples.

# NCBI dbSNP database of genetic variation

- This is the main repository of publicly available polymorphisms.

- You'll also find information on allele frequencies, populations, genotypes assays and much more.

- Most groups submit SNPs to dbSNP and only a few maintain web access to their SNPs.

# Submitting SNPs to dbSNP

- From their main web page, they have extensive information on how to submit SNPs, genotypes, validation experiments, population frequencies, etc., for any species.

- SNPs that you submit are called Submitter SNPs and get ssIDs.

- If there is a reference sequence available for the species submitted, they will map SNPs to this reference using the flank information you provide.

- SNPs that cluster at the same locus, are merged into Reference SNPs which have unique rsIDs.

**Reference SNP Cluster Report**

| | |
|---|---|
| NCBI SNP CLUSTER ID: | rs3137 |
| Organism: | human (*Homo sapiens*) |
| Variation Class: | SNP: single nucleotide polymorphism |
| Molecule Type: | Genomic |
| dbSNP build of first appearance: | 36 |
| dbSNP build of most recent change to cluster: | 116 |

SNP Details are categorized in the following sections:

Submission | Fasta | Resource | Locus | Map | Variation | Validation

**Submitter records for this RefSNP Cluster**

The submission **ss10401625** has the longest flanking sequence of all cluster members and was used to instantiate sequence for **rs3137** durin current build.

| NCBI Assay ID | Handle|Submitter ID | Validation Status | Entry Date | Update Date | Build Added | Molecule Type | Sequence Orientation | Observed Alleles |
|---|---|---|---|---|---|---|---|---|
| ss3168 | WIAF|WIAF-1477 | ⚔ | 01/23/99 | 10/10/03 | 36 | cDNA | forward | C/T |
| ss8206 | CGAP-GAI|47647 | ⚔ | 08/23/99 | 10/10/03 | 92 | cDNA | forward | C/T |
| ss1531001 | LEE|546510 | | 09/13/00 | 10/10/03 | 92 | cDNA | reverse | G/A |
| ss4395874 | LEE|ge546510 | | 04/25/02 | 10/10/03 | 106 | cDNA | reverse | G/A |
| ss4420318 | LEE|e546510 | | 04/26/02 | 10/10/03 | 106 | cDNA | reverse | G/A |
| ss10401625 | BCM_SSAHASNP|chr7.NT_007933.12_24364459 | ⚔ | 06/29/03 | 10/10/03 | 116 | Genomic | reverse | G/A |

**Fasta sequence  (Legend)**

```
>gnl|dbSNP|rs3137|allelePos=376|totalLen=576|taxid=9606|snpclass=1|alleles='C/T'|mol=Genomic|build=116

GGAAGTGACT CCTGGGTGag gtgagatggc tctcatctcc tgagggcaat tctccagttt
ctggagaaca agagctgtga gcacttgcag cccaaactca gagcagctgg ggatgggggt
ctcagcttgg tagaggggac tggacagggc accaCAGTGT ACAACACATA TGGtcaacaa
atatttattg ggcatttatt gtaagccagg caAGTCAGCA GAAACGGCCT GAGCAGTGCC
CAAGAGCACT CACTCACTCT CCCTAGCAAA CAGGCTCAGA ACTCTCTCAC ACATGTCATC
CTCTTTCCCA CTCAAAACTC CCACCCCAAC CTTCCTGGAA GGCAGGGCTA ACAGGACCTC
CTGCCTGCCT GCTCA
Y
GACTGATTAC TTTCAATCCC AGCTGCAATG CAAACTGAAA CTCATTCTGT ATATCACCAC
TCTACAGGAG AGGTCTATTT CTGGGGCACC CAGAAGTCAG CACACATACT GCTGGGACCA
GGACTCGTAA TTCGCCTTGG TCCAACTCCT TCTATGGGTT TAGCTGCCCT CATTCCTGTG
GGTAATACAA GATCAAACAG
```

**NCBI Resource Links**

**Submitter-Referenced Accessions:**
dbSTS:
GenBank: NT_007933 Hs.110839

**dbSNP Blast Analysis:**
NCBI RefSeq NM (mRNA): NM_014569.2 NM_145102.1
GenBank HTGS Finished: AACC01000011.1 AC005020.5
GenBank STS: G15373.1
GenBank mRNA: AB023232.1 AF170025.1 BX648490.1

**UniGene transcribed sequence cluster:**
UniGene Cluster ID: 110839

19

### LocusLink Analysis

**LocusLink via analysis of contig annotation:** ZFP95 zinc finger protein 95 homolog (mouse)
Click to see [all] [cSNP] [has frequency] [double hit] [haplotye tagged] variations associated with this gene.

Gene Model (contig mRNA transcript) NT_079595->NM_014569:

| Contig accession | Contig position | mRNA accession | mRNA orientation | Protein accession | Function | dbSNP allele | Protein residue | Codon position | Amino acid position |
|---|---|---|---|---|---|---|---|---|---|
| NT_079595 | 24393474 | NM_014569 | forward | | untranslated region | | | | |

### Integrated Maps:

**NCBI MapViewer:** rs3137 maps exactly once on NCBI human chromosome 7

| Chromosome | Contig accession | Contig position | Chromosome position | Hit orientation | Group term | Group label | Contig label |
|---|---|---|---|---|---|---|---|
| 7 | NT_079595.1 | 24393474 | 98120626 | minus strand | alt_assembly_1 | Toronto | Toronto |
| 7 | NT_007933.13 | 24364459 | 98742272 | minus strand | ref_haplotype | reference | reference |

**NCBI Sequence Viewer:** See rs3137 in Sequence Viewer.

**Project Ensembl:** Query rs3137 in Ensembl.

**UC Santa Cruz Genome Assembly:** Query rs3137 on the Santa Cruz Assembly.

### Variation Summary:

Assay sample size (number of chromosomes): 38
Population data sample size (number of chromosomes): 308
Total number of populations with frequency data: 2
Total number of individuals with genotype data: 5  Genotype Detail NEW
Hardy-weinberg Probability: 0.883
Average estimated heterozygosity: 0.491
Average Allele Frequency:
T     0.566
C     0.434

### Validation Summary:

**Validation status:** DoubleHit found by: BCM_SSAHASNP, NCBI
Marker displays Mendelian segregation: UNKNOWN
PCR results confirmed in multiple reactions: UNKNOWN
Homozygotes detected in individual genotype data: UNKNOWN

## Validation summary

| | |
|---|---|
| 🔵 | validated by multiple, independent submissions to the refSNP cluster |
| ✂ | validated by frequency or genotype data: minor alleles observed in at least two chromosomes. |
| 🧪 | validated by submitter confirmation |
| 📊 | all alleles have been observed in at least two chromosomes apiece |

## Viewing SNPs in Browsers

NCBI          Ensembl          UCSC

# How to find SNPs in a region of interest

- Gene based example

- A 2 Mbp region



http://www.ncbi.nlm.nih.gov/SNP/index.html

http://www.ncbi.nlm.nih.gov/entrez/query/Snp/EntrezSNPlegend.html



http://innateimmunity.net/IIPGA2/PGAs/InnateImmunity/MEFV/

23

http://www.ensembl.org/Homo_sapiens

Many submissions, however, possibly all from same source sequences.

| NCBI Assay ID | Handle\|Submitter ID | Validation Status | Entry Date | Update Date |
|---|---|---|---|---|
| ss290959 | KWOK\|OVLP-000621-270987 | | 06/30/00 | 10/10/03 |
| ss508456 | SC_JCM\|AJ003147.1_213692 | | 07/12/00 | 10/10/03 |
| ss1011433 | KWOK\|OVLP-000804-197113 | | 09/02/00 | 10/10/03 |
| ss1780721 | KWOK\|OVLP-000925-363908 | | 10/05/00 | 10/10/03 |
| ss1829272 | KWOK\|OVLP-000925-377600 | | 10/05/00 | 10/10/03 |
| ss2421405 | HGBASE\|SNP000002845 | | 11/07/00 | 10/10/03 |

IMS-JST095225

**Submitter records for this RefSNP Cluster**

The submission **ss4929937** has the longest flanking sequence of all cluster BLAST analysis for the current build.

| NCBI Assay ID | Handle\|Submitter ID | Validation Status | Entry Date | Update Date |
|---|---|---|---|---|
| ss4929937 | YUSUKE\|IMS-JST095225 | | 08/01/02 | 10/10/03 |

**Analysis of the three most common MEFV mutations in 412 patients with familial Mediterranean fever.**

Zaks N, Shinar Y, Padeh S, Lidar M, Mor A, Tokov I, Pras M, Langevitz P, Pras E, Livneh A.

Heller Institute of Medical Research, Sheba Medical Center, Tel Hashomer, Israel.

BACKGROUND: Familial Mediterranean fever is an autosomal recessive disease characterized by recurrent attacks of fever and serositis. The disease is caused by mutations in the MEFV gene, presumed to act as a down-regulator of inflammation within the polymorphonuclear cells. OBJECTIVES: To present the results of 412 FMF patients genotyped for three MEFV mutations, M694V, V726A and E148Q. RESULTS: The most frequent mutation, M694V, was detected in 47% of the carrier chromosomes. This mutation, especially common among North African Jewish FMF patients, was not found in any of the Ashkenazi (East European origin) patients. Overall, one of the three mutations was detected in 70% of the carrier chromosomes. M694V/M694V was the most common genotype (27%), followed by M694V/V726A (16%). The full genotype could be assessed in 57% of the patients, and one disease-causing mutation in an additional 26%. Only one patient with the E148Q/E148Q genotype was detected despite a high carrier rate for this mutation in the Jewish population, a finding consistent with a low penetrance of this genotype. The M694V/M694V genotype was observed in 15 patients with amyloidosis compared to 4 amyloidosis patients with other genotypes (P < 0.0001). CONCLUSIONS: Because of low penetrance and as yet other undetermined reasons, mutation analysis of the most common MEFV mutations supports a clinical diagnosis in only about 60% of patients with definite FMF.

Publication Types:
- Comment

Isr Med Assoc J. 2003 Aug;5(8):585-8.

PMID: 12929299 [PubMed - indexed for MEDLINE]

27

Double hit SNP minor allele frequency characteristics

Credit: Dr. Paul Hardenbol, Parallele Bioscience



http://www.ensembl.org/Multi/martview?species=Homo_sapiens

| Chromosome Name | Start Position (bp) | Strand | Reference ID | Allele | Mapweight | Heterozygosity | Ensembl gene name |
|---|---|---|---|---|---|---|---|
| 2 | 37848035 | -1 | 2231503 | C/G | 1 | 0 | ENSG00000163171.1 |
| 2 | 38018879 | 1 | 4670779 | C/T | 1 | 0 | ENSG00000177956.1 |
| 2 | 38019365 | 1 | 4670218 | C/G | 1 | 0 | ENSG00000177956.1 |
| 2 | 38153669 | 1 | 4670800 | A/G | 1 | 0 | ENSG00000115841.3 |
| 2 | 38272674 | -1 | 1800440 | A/G | 1 | 0.22283 | ENSG00000138061.1 |
| 2 | 38272685 | -1 | 1056837 | A/C/T | 1 | 0.412616 | ENSG00000138061.1 |
| 2 | 38272704 | -1 | 4986888 | C/G | 1 | 0.035188 | ENSG00000138061.1 |
| 2 | 38272711 | -1 | 4986887 | C/G | 1 | 0.0117367 | ENSG00000138061.1 |
| 2 | 38272738 | -1 | 1056836 | C/G | 1 | 0.417813 | ENSG00000138061.1 |
| 2 | 38272918 | 1 | 4398252 | C/T | 1 | 0 | ENSG00000138061.1 |
| 2 | 38276712 | -1 | 1056827 | G/T | 1 | 0 | ENSG00000138061.1 |
| 2 | 38276925 | -1 | 10012 | C/G | 1 | 0.44473 | ENSG00000138061.1 |
| 2 | 38382501 | 1 | 68352 | C/T | 1 | 0.5 | ENSG00000177744.1 |
| 2 | 38500195 | 1 | 7582826 | C/G | 1 | 0 | ENSG00000119787.2 |
| 2 | 38500195 | 1 | 7582826 | C/G | 1 | 0 | ENSG00000119787.2 |
| 2 | 38578886 | -1 | 3731847 | C/T | 1 | 0 | ENSG00000119787.2 |
| 2 | 38578886 | -1 | 3731847 | C/T | 1 | 0 | ENSG00000119787.2 |
| 2 | 38683723 | 1 | 7559613 | C/T | 1 | 0 | ENSG00000175340.1 |
| 2 | 38891505 | 1 | 6741892 | A/T | 1 | 0 | ENSG00000143891.2 |
| 2 | 38891505 | 1 | 6741892 | A/T | 1 | 0 | ENSG00000143891.2 |
| 2 | 38983484 | 1 | 1056104 | A/G | 1 | 0.44145 | ENSG00000152147.1 |
| 2 | 39056879 | 1 | 7598922 | C/T | 1 | 0 | ENSG00000183254.2 |
| 2 | 39056879 | 1 | 7598922 | C/T | 1 | 0 | ENSG00000163214.3 |
| 2 | 39198647 | 1 | 8192671 | C/T | 1 | 0 | ENSG00000115904.1 |
| 2 | 39489827 | -1 | 1061687 | A/G | 1 | 0 | ENSG00000011566.1 |



http://www.ensembl.org/Homo_sapiens/transview?transcript=ENST00000281950&db=core

31

# *Haplotype Map project*

- What is a Haplotype?

- What is Linkage Disequilibrium (LD)?

- What is the Haplotype Map Project?

# *What is a Haplotype?*

- A set of closely linked genetic markers present on one chromosome which tend to be inherited together (not easily separable by recombination).

- Recombination occurs between homologous chromosomes when cells divide.

- It is believed that recombination is not equally likely across the genome, but that it is punctuated by hot-spots.

From: Goldstein DB. Islands of linkage disequilibrium. Nat Genet. 2001 Oct;29(2):109-11.

# *What is Linkage Disequilibrium?*

- When the observed frequencies of genetic markers in a population does not agree with haplotype frequencies predicted by multiplying together the frequency of individual genetic markers in each haplotype.

| | |
|---|---|
| 139 | 0.352 |
| 140 | 0.5 |
| 141 | 0.499 |
| 142 | 0.5 |
| 143 | 0.499 |
| 144 | 0.453 |
| 145 | 0.499 |
| 146 | 0.497 |

CAACTCAT .217    $0.352*0.5^7=0.00275$
TGGTCTGC .365    $0.648*0.5^7=0.00534$
TGGTCCGC .127    $0.648*0.5^7=0.00534$
TAACTCAT .266    $0.648*0.5^7=0.00534$

0.975

## *Haplotype Map project*

- The goal of the International HapMap Project is to develop a haplotype map of the human genome, the HapMap, which will describe the common patterns of human DNA sequence variation.

- The HapMap is expected to be a key resource for researchers to use to find genes affecting health, disease, and responses to drugs and environmental factors.

- The information produced by the Project will be made freely available.

http://www.hapmap.org/abouthapmap.html

## *HapMap Strategy*

- To develop the HapMap, samples from 270 individuals will be genotyped for at least 1 million SNPs across the human genome.

- DNA samples come from:
  - Nigeria (30 both-parent-and-adult-child trios)
  - Japanese in Tokyo (45 unrelated individuals)
  - Han Chinese in Beijing (45 unrelated individuals)
  - CEPH (30 trios, Northern and Western Europe ancestry)

## *Data Analysis*

- Genotyped SNPs are analyzed for association using standard measures, such as *D'* and $r^2$.
  - Deviation from equilibrium between two markers is denoted by *D*.
  - When normalized it is called *D'* and has a range from -1 to +1.
  - $r^2$ uses a different normalization method and ranges between 0 and 1.
  - See URL below for a good description of these measures.

  http://www.ucl.ac.uk/~ucbhdjm/courses/b242/2+Gene/2+Gene.html

## *Current status of HapMap*

- November 1st, 2003: First major public data release!
  - Over 13 million genotypes from 145,554 SNPs
  - Associated allele frequency and assay data have been released for public download
- Here's an example of generating haplotype information from the current data release…

# *HaploView*

- Developed and maintained by Jeffrey Barrett in Mark Daly's lab at The Broad Institute.
- Haploview currently allows users to:
  - examine block structures
  - generate haplotypes in these blocks
  - run association tests
  - and save the data in a number of formats.

# *Perlegen's haplotype map*

- Used chip based resequencing to discover SNPs across the genome.

- Applied this technology to single haploid copies of each chromosome from a number of different individuals.

- From this information, haplotypes can be deduced directly from the data.

- They have released data from chromosome 21 for public use.

# *Concluding remarks*

- Along with the emergence of the human genome, we also have a growing database of variations that are critical to the overall value of the human genome sequence.

- These variations are what make us all (phenotypically) different, and impart different levels of resistance and susceptibility to disease.

- The collection of human sequence variation information will continue to evolve rapidly.

# *References*

### EST SNPs

Hu G, Modrek B, Riise Stensland HM, Saarela J, Pajukanta P, Kustanovich V, Peltonen L, Nelson SF, Lee C., Efficient discovery of single-nucleotide polymorphisms in coding regions of human genes. Pharmacogenomics J. 2002;2(4):236-42.

Clifford R, Edmonson M, Hu Y, Nguyen C, Scherpbier T, Buetow KH., Expression-based genetic/physical maps of single-nucleotide polymorphisms identified by the cancer genome anatomy project. Genome Res. 2000 Aug;10(8):1259-65.

Irizarry K, Kustanovich V, Li C, Brown N, Nelson S, Wong W, Lee CJ., Genome-wide analysis of single-nucleotide polymorphisms in human expressed sequences. Nat Genet. 2000 Oct;26(2):233-6.

### Clone Overlaps/TSC

The International SNP Map Working Group, A map of human genome sequence variation containing 1.4 million SNPs. Nature 15 February 2001, v409, 928 - 933

Ning Z, Cox AJ, Mullikin JC, SSAHA: a fast search method for large DNA databases. Genome Res. 2001 Oct;11(10):1725-9.

Marth G, Schuler G, Yeh R, Davenport R, Agarwala R, Church D, Wheelan S, Baker J, Ward M, Kholodov M, Phan L, Czabarka E, Murvai J, Cutler D, Wooding S, Rogers A, Chakravarti A, Harpending HC, Kwok PY, Sherry ST. Sequence variations in the public human genome data reflect a bottlenecked population history. Proc Natl Acad Sci U S A. 2003 Jan 7;100(1):376-81.

### Heteroduplex analysis

Hecker KH, Taylor PD, Gjerde DT. Mutation detection by denaturing DNA chromatography using fluorescently labeled polymerase chain reaction products. Anal Biochem. 1999 Aug 1;272(2):156-64.

Davies H, Bignell GR, Cox C, Stephens P, Edkins S, et al. Mutations of the BRAF gene in human cancer. Nature. 2002 Jun 27;417(6892):949-54.

# *References*

### Targeted Resequencing

Haga H, Yamada R, Ohnishi Y, Nakamura Y, Tanaka T. Gene-based SNP discovery as part of the Japanese Millennium Genome Project: identification of 190,562 genetic variations in the human genome. Single-nucleotide polymorphism. J Hum Genet. 2002;47(11):605-10.

### Chip based SNP discovery

Patil N, Berno AJ, Hinds DA, Barrett WA, Doshi JM, Hacker CR, Kautzer CR, Lee DH, Marjoribanks C, McDonough DP, Nguyen BT, Norris MC, Sheehan JB, Shen N, Stern D, Stokowski RP, Thomas DJ, Trulson MO, Vyas KR, Frazer KA, Fodor SP, Cox DR. Blocks of limited haplotype diversity revealed by high-resolution scanning of human chromosome 21. Science. 2001 Nov 23;294(5547):1719-23.

### SNP quality

Reich DE, Gabriel SB, Altshuler D. Quality and completeness of SNP databases. Nat Genet. 2003 Apr;33(4):457-8.

Carlson CS, Eberle MA, Rieder MJ, Smith JD, Kruglyak L, Nickerson DA. Additional SNPs and linkage-disequilibrium analyses are necessary for whole-genome association studies in humans. Nat Genet. 2003 Apr;33(4):518-21.

### Haplotype Map Project

Dennis C. Special section on human genetics: the rough guide to the genome. Nature. 2003 Oct 23;425(6960):758-9.

Goldstein DB. Islands of linkage disequilibrium. Nat Genet. 2001 Oct;29(2):109-11.

# *WEB pages*

snp.cshl.org : The SNP Consortium web pages

http://droog.mbt.washington.edu/PolyPhred.html

http://www.ncbi.nlm.nih.gov/SNP/index.html : dbSNP home page

http://www.ensembl.org : Ensembl home page

http://www.ucl.ac.uk/~ucbhdjm/courses/b242/2+Gene/2+Gene.html

http://www.hapmap.org/: Haplotype Map Project home page

http://www.hapmap.org/cgi-perl/gbrowse/gbrowse/hapmap

http://www.broad.mit.edu/personal/jcbarret/haploview/

http://www.perlegen.com/haplotype/ : Perlegen's chr21 HapMap