*e!*

# Sequence Variation in Ensembl

Giulietta Spudich
February, 2007

---

*e!*

# Overview

- Genomic Diversity (SNPs)
- Variations in the Ensembl Browser
- Variations in BioMart

*e!*

# Genomic Diversity

**Mutations:**

**base pair substitutions
insertion/deletion (frameshifts)**

**1 in every 300 bp (*human*)**
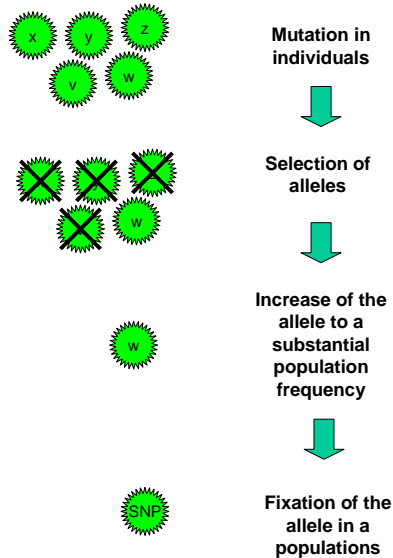
**~3 billion base pairs in mouse
genome!**

---

*e!*

# Single nucleotide polymorphisms (SNPs)

- **Polymorphism: a DNA variation in which each possible sequence is present in at least 1% of the population**

- **Most polymorphisms (~90%) take the form of SNPs: variations that involve just one nucleotide**

# Origin of SNPs



Mutation in individuals

⬇

Selection of alleles

⬇

Increase of the allele to a substantial population frequency

⬇

Fixation of the allele in a populations

Adapted from Bioinformatics for Geneticists, Eds Barnes and Gray

---

# Functional Consequences

| Type | Consequence |
|------|-------------|
| SNPs in coding area that alter aa sequence | Cause of most monogenic disorders, e.g: Cystic fibrosis (CFTR) Hemophilia (F8) |
| SNPs in coding areas that don't alter aa sequence | May affect splicing |
| SNPs in promoter or regulatory regions | May affect the level, location or timing of gene expression |
| SNPs in other regions | No direct known impact on phenotype Useful as markers |

# Studying variation – why?

- Determine disease risk
- Predict responses to environmental changes and drugs (pharmacogenomics)

- Biological markers
- Forensics
- Evolution
- Laboratory: hybridisation studies, marker-assisted breeding

---

# SNPs in Ensembl - Species

- Human
- Mouse
- Rat
- Dog

- Chicken
- Zebrafish
- Mosquito

*e!*

# SNPs in Ensembl

- **Most SNPs imported from dbSNP (rs……):**
  - **Imported data: alleles, frequencies, flanking sequence….**
  - **Calculated data: synonymous status, peptide shift, SNP position….**

- **For mouse also:**
  - **Sanger**

---

*e!*

# dbSNP (NCBI)

- **Main database of SNPs (and short polymorphisms: in-dels)**

- **6,491,554 rs (reference SNPs) in mouse (*11,961,761 in human*).**

- **4,990,170 validated in mouse**
- **(*5,646,244 in human*).**

- **http://www.ncbi.nlm.nih.gov/SNP**

ENTREZ **SNP**
Single Nucleotide Polymorphism

dbSNP - Validation

Method classes organize submissions by a general methodological or experimental approach to assaying for variation in the DNA sequence.

| Method class | Class code and XML |
| --- | --- |
| Denaturing high pressure liquid chromatography (DHPLC) | 1 |
| DNA hybridization | 2 |
| Computational analysis | 3 |
| Single-stranded conformational polymorphism (SSCP) | 5 |
| Other | 6 |
| Unknown | 7 |
| Restriction fragment length polymorphism (RFLP) | 8 |
| Direct DNA sequencing | 9 |



SNPs in Ensembl

SNPs in Ensembl (Views)



TransView

SNPs in sequence

**SNP in different strains**
**Alleles, type (i.e. placement in relation to a gene)**



SARA = <u>S</u>ame <u>A</u>s <u>R</u>eference <u>A</u>ssembly (C57BL/6J)

# *GeneSNPView*

Gene tree info
Gene variation info.
ID history

What SNPs does my gene contain?

**Choose SNP type**

SNPs and variations in region of gene ENSG0000013040...

Features ▼   Source ▼   SNP class ▼   Validation ▼   SNP type ▼   Context ▼   Image size ▼   Export ▼      Help ▼

**Assembly**

**Transcript zoomed**

**SNPs** →

SNP:rs33141009   ✕
SNP properties
bp: 127092040 - 127092040
status:
class: snp
ambiguity code: Y
mapweight: 1
alleles: C/T
source: dbSNP
type: SPLICE_SITE,
SYNONYMOUS_CODING

**Table of Variations**

Variations in ENST00000252723

| ID | Type | Chr: bp | Alleles | Ambiguity | AA | AA change co-ordinate | Class | Source | Validation |
|---|---|---|---|---|---|---|---|---|---|
| rs34144627 | SPRIME_UTR | 7: 100156486-100156485 | -/A | - | - | - | insertion | dbSNP | - |
| rs34907403 | FRAMESHIFT_CODING | 7: 100157590-100157589 | -/T | - | - | 77 (1) | insertion | dbSNP | - |
| rs307592 | INTRONIC | 7: 100157872 | C/T | Y | - | - | snp | HGVbase, dbSNP | cluster, doublehit |
| rs33078705 | INTRONIC | 7: 100158068-100158067 | -/TCAC | - | - | - | insertion | dbSNP | - |
| rs484199 | INTRONIC | 7: 100158087 | G/A | R | - | - | snp | HGVbase, dbSNP | - |
| rs7789679 | INTRONIC | 7: 100158157 | G/A | R | - | - | snp | dbSNP | - |
| rs11976253 | NON_SYNONYMOUS_CODING | 7: 100158317 | C/T | - | P/L | 114 (2) | snp | dbSNP | cluster |
| rs1126887 | NON_SYNONYMOUS_CODING | 7: 100158394 | G/C | S | G/R | 140 (1) | snp | dbSNP | - |
| rs304449 | SPRIME_UTR | 7: 100159074 | T/G | K | - | - | snp | HGVbase, dbSNP, TSC, Affy GeneChip 500K Mapping Array | - |
| rs551238 | DOWNSTREAM | 7: 100159464 | C/A | M | - | - | snp | HGVbase, dbSNP, freq, TSC | cluster, doublehit |
| rs4726606 | DOWNSTREAM | 7: 100159726 | T/C | Y | - | - | snp | HGVbase, dbSNP | cluster, doublehit |
| rs33972223 | DOWNSTREAM | 7: 100159780 | G/- | - | - | - | deletion | dbSNP | - |

© 2006 WTSI / EBI. Ensembl is available to download for public use - please see the code licence for details.

---

# **SNPView**

Info about one specific SNP:

- **SNP Report:**

  **Imported from…** dbSNP, TSC, HGVbase, Affy Chip Data

  **Validation**

  **Genotype and allele frequencies per population**

  **Location in transcripts**

  **Type:coding/noncoding**

SNP Report

Genotype frequencies per population
Allele frequencies per population
SNP rs4149711 is located in the following transcripts
SNP Context - chromosome X 138569741
Individual genotypes for SNP rs4149711

**GeneSeqView**

SNPs in genomic sequence



**MapView**

SNP density on a chromosome

**Example: Mouse chromosome 8**

Obtain a table of genes corresponding to SNPs



Obtain a table of SNPs for Ensembl genes

SNPs in BioMart

SNP FILTER options

SNP Attribute options



Strain-specific SNPs in BioMart

# Summary

- Genomic Diversity (SNPs)
- Variations in the Ensembl Browser
- Variations in BioMart

---

# Ensembl Team

**Leaders**    <u>Ewan Birney</u> **(EBI),** <u>Tim Hubbard</u> **(Sanger Institute)**

<u>Glenn Proctor</u>, Ian Longden, Patrick Meidl, Andreas Kähäri

<u>Arek Kasprzyk</u>, Syed Haider, Richard Holland, Damian Smedley, Benoît Ballester

Eugene Kulesha

<u>Xosé M Fernández</u>, Bert Overduin, Michael Schuster, Giulietta Spudich

<u>James Smith</u>, Fiona Cunningham, Anne Parker, Bethan Pritchard, Stephen Rice, Steve Trevanion, Matt Wood

<u>Abel Ureta-Vidal</u>, Benoit Ballester, Kathryn Beal, Stephen Fitzgerald, Javier Herrero Sánchez, Albert Vilella

<u>Val Curwen</u>, <u>Steve Searle</u>, Bronwen Aken, Julio Banet, Laura Clarke, Sarah Dyer, Kevin Howe, Felix Kokocinski, Jan-Hinnerck Vogel, Simon White

**<u>Paul Flicek</u>, Yuan Chen, Stefan Gräf, Nathan Johnson, Daniel Rios**

<u>Martin Hammond</u>, Dan Lawson, Karyn Megy