# SNP Comparison
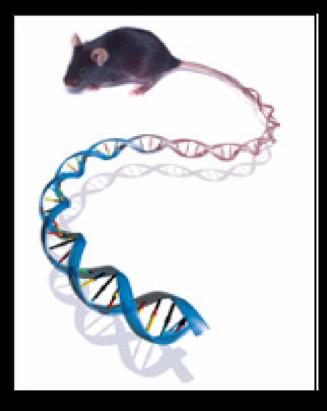
Group Members
Amira
Jhelum
Rahul
Shweta

# Chimpanzee Or Mouse

- The earlier goal of the project was to compare SNP distribution over certain genes in human and chimpanzee
- But due to unavailability of sufficient amount of data on chimpanzee, we had to change our focus to mouse

# Reasons for choosing Mouse

- – Mouse Genome and SNP data is readily available on the NCBI website
- – Mouse is the most important animal model and is widely used in the study of human diseases.
- – Mouse carries virtually the same set of genes as the human and more then 90% of the mouse genome can be lined up with a region on the human genome.

# Accomplishments Till Date

- Homology searching for Human and Mouse Genes on the NCBI website.

- List of potential 484 genes

- Analyzed these homologous genes for SNIPS
  - Total number of SNIPS
  - SNIPS in the coding region

- Tabulated data
  - Unfiltered data
  - Filtered data

# NCBI Website

# List of Homologenes

# Searching for SNIPS

# Brief Description of IRAK4 Gene

# SNIPS for IRAK4 Gene

# Total SNIPS for IRAK4 Gene

# Tabulation Of the Data

| ORGANISM | GENE NAME | GENE ID | SNIPS IN CODING REGION | TOTAL NO. OF SNIPS | ALLELES FOR SNIPS IN CODING REGION |
|---|---|---|---|---|---|
| HUMAN | PAXIP1L | 22976 | 6 | 140 | [G/A] [G/A] [C/T] [A/C] [A/G] [C/A] |
| MOUSE | Paxip1 | 55982 | 0 | 2 | N/A |
| | | | | | |
| HUMAN | POLI | 11201 | 11 | 177 | [A/G] [G/A] [C/T] [G/A] [A/G] [A/T] [G/A] [G/A] [C/T] [C/T] [A/G] |
| MOUSE | Poli | 26447 | 0 | 0 | N/A |
| | | | | | |
| HUMAN | TFCP2 | 7024 | 0 | 292 | N/A |
| MOUSE | Tcfcp2 | 21422 | 0 | 2 | N/A |
| | | | | | |
| HUMAN | ADH5 | 128 | 3 | 72 | [G/T] [A/T] [T/G] |
| MOUSE | Adh5 | 11532 | 0 | 0 | N/A |
| | | | | | |
| HUMAN | CYP4V2 | 285440 | 8 | 110 | [C/G] [G/C] [T/A] [A/C] [G/T] [C/T] [T/C] [C/G] |
| MOUSE | Cyp4v3 | 102294 | 11 | 44 | [T/C] [C/T] [A/C] [T/C] [G/A] [G/C] [G/A] [C,T,A] [C/T] [C/T] [G/A] |
| | | | | | |
| | NOL6 | 65083 | 2 | 41 | [G/T] [T/C] |
| MOUSE | Nol6 | 230082 | 1 | 7 | [C/T] |
| | | | | | |
| HUMAN | DBI | 1622 | 5 | 31 | [G/A] [A/G] [G/A] [C/T] [A/G] |
| MOUSE | Dbi | 13167 | 2 | 7 | [T/A] [G/A] |
| | | | | | |
| HUMAN | DPEP3 | 64180 | 2 | 7 | [A/G] [A/G] |
| MOUSE | Dpep3 | 71854 | N/A | N/A | N/A |
| | | | | | |
| HUMAN | LANCL2 | 55915 | 5 | 230 | [C/G] [C/A] [G/A] [G/A] [T/C] |
| MOUSE | Lancl2 | 71835 | 0 | 14 | N/A |

# Data Analysis

- Concentrated our study on genes which have SNIPS in their coding region for both mouse and human.

- After first round of analysis, we reduced our data from 484 genes to 82 genes based on number of SNIPS in the coding region.

# Filtered Data

| Organism | Gene Name | Gene ID | SNPs in Coding region | Total SNPs | Alleles for SNPS in coding region |
|---|---|---|---|---|---|
| | | | | | |
| | | _ | | | |
| | | | | | |
| HUMAN | FKBP9 | 11328 | 7 | 313 | [A/ G/] [ G/ A] [G/C] [C/T] [A/G] [C/t] [T/C] |
| MOUSE | Fkbp9 | 27055 | 3 | 33 | [T/ C] [ G/A]  [C/G] |
| | | | | | |
| HUMAN | TPI1 | 7167 | 3 | 19 | [T/ G] [C/G] [G/T] |
| MOUSE |  Tpi | 21991 | 2 | 3 | [G/A] [A/T] |
| | | | | | |
| HUMAN | ITCH | 83737 | 6 | 485 | [C/T] [G/T] [T/C] [T/A] [A/G] [T/A] |
| MOUSE | Itch | 16369 | 1 | 27 | [T/G] |
| | | | | | |
| HUMAN | KRTHB1 | 3887 | 8 | 39 | [A/A] [G G] [T/C] [G/A] [C/T] [T/C] [G/T] [C G] |
| MOUSE | Krt2-19 | 64818 | 8 | 10 | [T/C] [T/C] [A/G] [G/A] [T/C] [A/C] [T/C] [C/T] |
| | | | | | |
| HUMAN | NRF1 | 4899 | 2 | 398 | [G/T] [T/C] |
| MOUSE | Nrf1 | 18181 | 1 | 2 | [A/G] |
| | | | | | |
| HUMAN | EIF4A2 | 1974 | 6 | 20 | [G/T] [T/C] [C/A] [T/C] [G/A] [A/T] |
| MOUSE | Eif4a2 | 13682 | 1 | 23 | [NA/T] |
| | | | | | |
| HUMAN | CLOCK | 9575 | 4 | 411 | [C/T] [G/A] [A/C] [A/G] |
| MOUSE | Clock | 12753 | 1 | 18 | [A/G] |
| | | | | | |
| HUMAN | ROR2 | 4920 | 5 | 770 | [A/G] [T/C] [T/C] [G/A] [C/T] |
| MOUSE | Ror2 | 26564 | 1 | 37 | [C/G] |

# Interim Results

- Consolidated data from Human and Mouse
  - Based on Homologous genes between the two species.
  - Distribution of SNIPS on the homologenes.
- Analyzed the data to select the genes with SNIPS in the coding region for both the species.

# Next Steps

- Analyze the selected 82 genes and draw statistical conclusions of biological significance from the above data.

- Further filter the data
  - To study the distribution of SNIPS on potential genes for both the species using a parser.

- The complete data for all 484 genes and the selected 82 genes is available on our website
  - http://www.angelfire.com/sk3/compbio601/