# Data Mining in Ensembl with BioMart

# BioMart- Data mining

- BioMart is a search engine that can find multiple terms and put them into a table format.

- Such as: mouse gene (IDs), chromosome and base pair position

- No programming required!

# General or Specific Data-Tables

- All the genes for one species

- Or… only genes on one specific region of a chromosome

- Or… genes on one region of a chromosome associated with an InterPro domain

# The First Step: Choose the Dataset

# The Second Step: Filters



**Filters define which genes we are looking at.**

# Attributes attach information



Determine output columns with Attributes.

# Results



**Tables or sequences**

# Query:

- For all mouse genes on chromosome 10 that are protein coding, I would like to know the **IDs** in both **Ensembl** and **MGI**.

  Are there **Illumina probes** and **GO IDs** for these genes?

- In the query:

  Filters: what we know

  Attributes: what we want to know.

# Query:

- For all **mouse genes** on **chromosome 10** that are **protein coding**, I would like to know the IDs in both Ensembl and MGI.

  Are there Illumina probes and GO IDs for these genes?

- In the query:

  **Filters: what we know**

  Attributes: what we want to know.

# Query:

- For all mouse genes on chromosome 10 that are protein coding, I would like to know the **IDs** in both **Ensembl** and **MGI**.

  Are there **Illumina probes** and **GO IDs** for these genes?

- In the query:

  Filters: what we know

  **Attributes: what we want to know.**

# A Brief Example

# Select the genes with Filters

# Filters (selecting the genes)

# Filters (selecting the genes)

# Attributes (Output Options)



We would like GO terms and IDs in MGI (the Mouse Genome Informatics site).

# Attributes (Output)



Click 'Results'

Scroll down to add 'Illumina v1' probes that map to these genes.

# The Results Table - Preview



'Results' shows Gene IDs, GO terms, and Illumina probes for all protein coding mouse genes on chromosome 10.

# Full Result Table

**Ensembl Gene and Transcript IDs**

**GO terms**

**MGI symbol**

**Illumina probes**

| Ensembl Gene ID | Ensembl Transcript ID | GO ID | GO description | MGI symbol | Illumina v1 |
|---|---|---|---|---|---|
| ENSMUSG00000015202 | ENSMUST00000015346 | GO:0005515 | protein binding | Cnksr3 | scl38236.13.428_30-S |
| ENSMUSG00000015202 | ENSMUST00000015346 | GO:0005737 | cytoplasm | Cnksr3 | scl38236.13.428_30-S |
| ENSMUSG00000015202 | ENSMUST00000015346 | GO:0009966 | regulation of signal transduction | Cnksr3 | scl38236.13.428_30-S |
| ENSMUSG00000015202 | ENSMUST00000015346 | GO:0016020 | membrane | Cnksr3 | scl38236.13.428_30-S |
| ENSMUSG00000015202 | ENSMUST00000105621 | GO:0005515 | protein binding | | scl38236.13.428_30-S |
| ENSMUSG00000015202 | ENSMUST00000105621 | GO:0005737 | cytoplasm | | scl38236.13.428_30-S |
| ENSMUSG00000015202 | ENSMUST00000105621 | GO:0009966 | regulation of signal transduction | | scl38236.13.428_30-S |
| ENSMUSG00000015202 | ENSMUST00000105621 | GO:0016020 | membrane | | scl38236.13.428_30-S |
| ENSMUSG00000064065 | ENSMUST00000086896 | | | A130090K04Rik | |
| ENSMUSG00000064065 | ENSMUST00000058132 | | | A130090K04Rik | ri|B930094H08|PX00167G09|AK |
| ENSMUSG00000064065 | ENSMUST00000058132 | | | A130090K04Rik | scl0018390.2_215-S |
| ENSMUSG00000064065 | ENSMUST00000105617 | | | A130090K04Rik | scl38237.12_618-S |
| ENSMUSG00000064065 | ENSMUST00000105617 | | | A130090K04Rik | scl0018390.2_215-S |
| ENSMUSG00000064065 | ENSMUST00000105617 | | | A130090K04Rik | ri|A130090K04|PX00125I14|AK |
| ENSMUSG00000064065 | ENSMUST00000105618 | | | A130090K04Rik | scl0018390.2_215-S |
| ENSMUSG00000064065 | ENSMUST00000105618 | | | A130090K04Rik | ri|A130090K04|PX00125I14|AK |
| ENSMUSG00000064065 | ENSMUST00000078070 | | | A130090K04Rik | scl0018390.2_215-S |
| ENSMUSG00000000766 | ENSMUST00000105615 | GO:0001584 | rhodopsin-like receptor activity | Oprm1 | scl39174.18_263-S |
| ENSMUSG00000000766 | ENSMUST00000105615 | GO:0004872 | receptor activity | Oprm1 | scl39174.18_263-S |
| ENSMUSG00000000766 | ENSMUST00000105615 | GO:0004966 | galanin receptor activity | Oprm1 | scl39174.18_263-S |
| ENSMUSG00000000766 | ENSMUST00000105615 | GO:0004982 | N-formyl peptide receptor activity | Oprm1 | scl39174.18_263-S |
| ENSMUSG00000000766 | ENSMUST00000105615 | GO:0004983 | neuropeptide Y receptor activity | Oprm1 | scl39174.18_263-S |
| ENSMUSG00000000766 | ENSMUST00000105615 | GO:0004985 | opioid receptor activity | Oprm1 | scl39174.18_263-S |

# Original Query:

- For all mouse genes on chromosome 10 that are protein coding, I would like to know the **IDs** in both **Ensembl** and **MGI**.

  Are there **Illumina probes** and **GO IDs** for these genes?

- In the query:

  Filters: what we know

  Attributes: columns in the **Result Table**

# Other Export Options (Attributes)

- ❖ Sequences: UTRs, flanking sequences, cDNA and peptides, etc

- ❖ Gene IDs from Ensembl and external sources (MGI, Entrez, etc)

- ❖ Microarray data

- ❖ Protein Functions/descriptions (Interpro, GO)

- ❖ Orthologous gene sets

- ❖ SNP/ Variation Data

# BioMart Data Sets

- Ensembl genes
  - Vega genes
    - Variations

# BioMart around the world…



BioMart started at
Ensembl…
To where has it travelled?

# Central Portal

**Powered by BioMart software:**

- BioMart Central Portal
- Ensembl
- HapMap
- HTGT

- Dictybase
- Wormbase
- Gramene
- Europhenome

- Rat Genome Database
- DroSpeGe
- ArrayExpress DW
- Eurexpress

- GermOnLine
- PRIDE
- PepSeeker
- VectorBase

- Pancreatic Expression Database
- Reactome
- EU Rat Mart
- Paramecium DB

**Third party software with BioMart Plugin:**

Bioclipse  biomaRt-BioConductor  Cytoscape  Galaxy  Taverna  WebLab

*www.biomart.org*

# WormBase

# HapMap

# DictyBase

# GRAMENE



*www.gramene.org*

# How to Get There

http://www.biomart.org/biomart/martview

http://www.ensembl.org/biomart/martview

- Or click on 'BioMart' from Ensembl