

# Know Your Transcriptome

Integrative Bioinformatic Approaches

Anil Jegga

Biomedical Informatics

Contact Information:

Anil Jegga

Biomedical Informatics

Room # 232, S Building 10th Floor

CCHMC

Homepage: <http://anil.cchmc.org>

Tel: 513-636-0261

E-mail: [anil.jegga@cchmc.org](mailto:anil.jegga@cchmc.org)

**Slides and Example data sets available for download at:**

**<http://anil.cchmc.org/dhc.html>**

**Workshop Evaluation**: Please provide your valuable feedback on the evaluation sheet provided along with the hand-outs

This workshop is about the analysis of transcriptome and **does not** cover microarray data analysis

Contact Huan Xu ([huan.xu@cchmc.org](mailto:huan.xu@cchmc.org)) for GeneSpring related questions or microarray data analysis

All the applications/servers/databases used in this workshop are **free** for academic-use. Applications that are not free for use (e.g. Ingenuity Pathway Analysis, MatInspector, etc.) are not covered here. However, we have licensed access to both of these and please contact us if you are interested in using them.

# I have a list of co-expressed mRNAs (Transcriptome)....

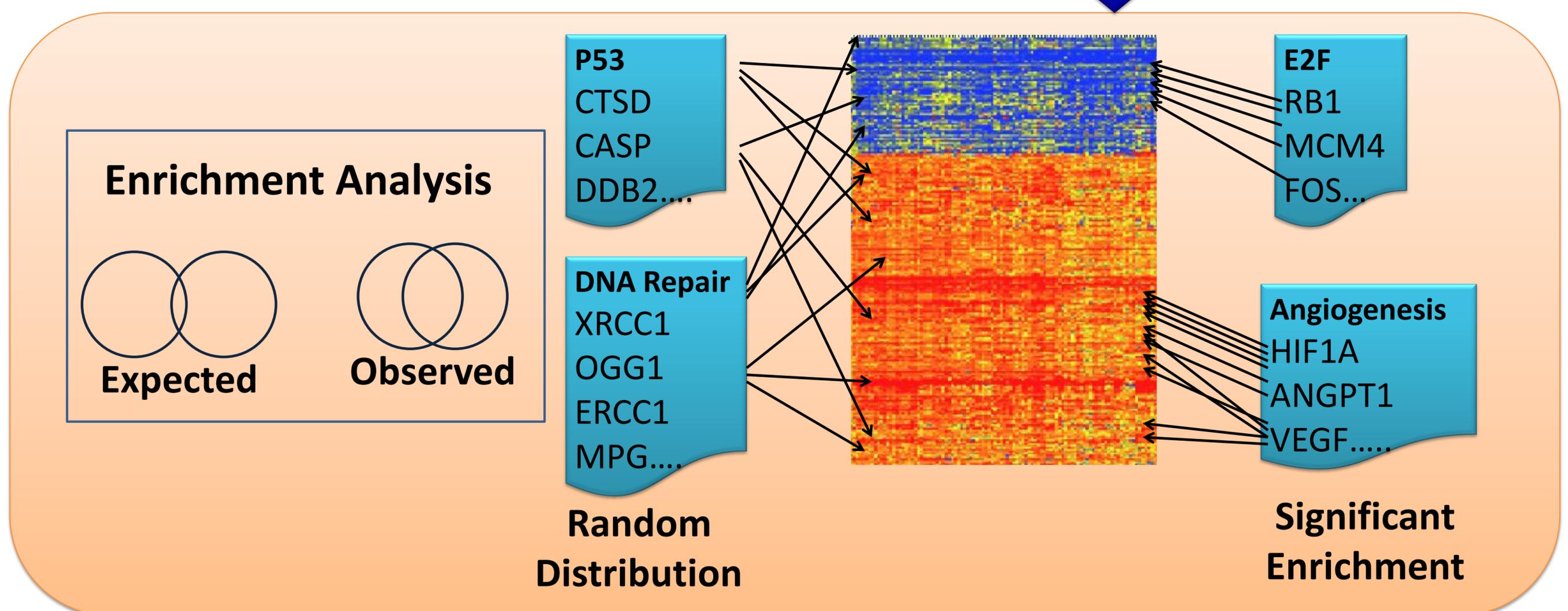
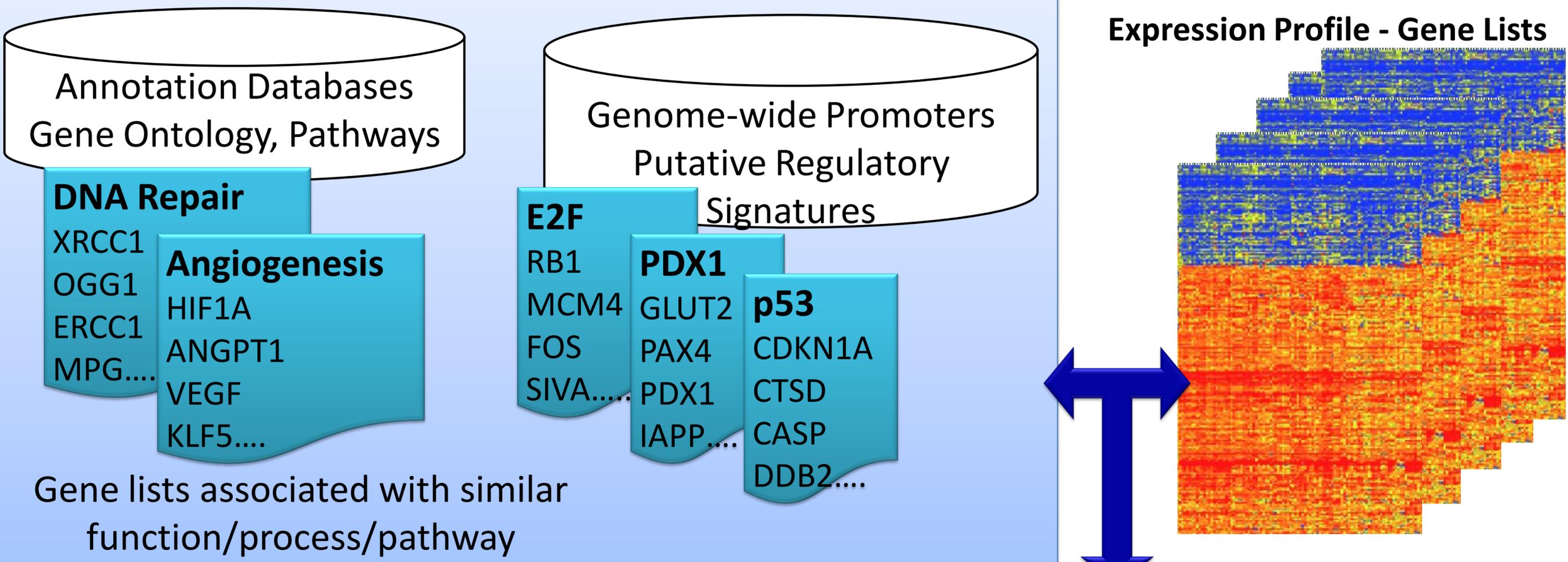
## Now what?

### 1. Identify putative shared regulatory elements

- Known transcription factor binding sites (TFBS)
  - Conserved
  - Non-conserved
- Unknown TFBS or Novel motifs
  - Conserved
  - Non-conserved
- MicroRNAs

### 2. Identify the underlying biological theme

- Gene Ontology
- Pathways
- Phenotype/Disease Association
- Protein Domains
- Protein Interactions
- Expression in other tissues/experiments
- Drug targets
- Literature co-citation...



# I have a list of co-expressed mRNAs (Transcriptome)....

## I want to find the shared cis-elements – Known and Novel

### □ Known transcription factor binding sites (TFBS)

#### ❖ Conserved

- oPOSSUM
- DiRE

#### ❖ Non-conserved

- Pscan
- **MatInspector** (\*Licensed)

### □ Unknown TFBS or Novel motifs

#### ❖ Conserved

- oPOSSUM
- **Weeder-H**

#### ❖ Non-conserved

- **MEME**
- **Weeder**

1. Each of these applications support different forms of input. Very few support probeset IDs.
2. **Red Font**: Input sequence required; Do not support gene symbols, gene IDs, or accession numbers. The advantage is you can use them for scanning sequences from any species.
3. \*Licensed software: We have access to the licensed version.

# I have a list of co-expressed mRNAs (Transcriptome)....

## I want to find the shared cis-elements – Known and Novel

### □ Known transcription factor binding sites (TFBS)

#### ❖ Conserved

- oPOSSUM
- DiRE

#### ❖ Non-conserved

- Pscan
- **MatInspector** (\*Licensed)

### □ Unknown TFBS or Novel motifs

#### ❖ Conserved

- oPOSSUM
- **Weeder-H**

#### ❖ Non-conserved

- **MEME**
- **Weeder**

# oPOSSUM (<http://www.cisreg.ca/oPOSSUM>)



# oPOSSUM

[About](#) | [Contact](#)

## Welcome to oPOSSUM

oPOSSUM is a web-based system for the detection of over-represented transcription factor binding sites in the promoters of sets of genes.

### Human SSA

Enter >>



Human **Single Site Analysis (SSA)** is designed to detect over-represented conserved **single** sites in human and mouse genes.

Reference: Ho Sui, *et al.* (2005). oPOSSUM: Identification of over-represented transcription factor binding sites in co-expressed genes. *NAR*, 33(10):3154-64. PMID: [15933209](#)

### Human CSA (Module analysis)

Enter >>



Human **Combination Site Analysis (CSA)** identifies over-represented **combinations** of conserved transcription factor binding sites in sets of human and mouse genes.

Reference: Huang, S., Fulton, D., *et al.* (2006). Identification of over-represented combinations of transcription factor binding sites in sets of co-expressed genes. *In Advances in Bioinformatics and Computational Biology*, Vol. 3. Imperial College Press, London, UK. 247-56. [PDF](#).

### Worm SSA

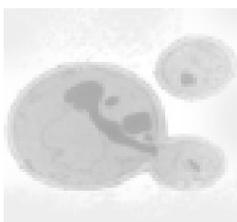
Enter >>



Worm Single Site Analysis (SSA) identifies over-represented conserved transcription factor binding sites in sets of *C. elegans* and *C. briggsae* genes.

### Yeast SSA

Enter >>



Yeast Single Site Analysis (SSA) identifies over-represented transcription factor binding sites in sets of *S. cerevisiae* genes. Phylogenetic footprinting has not been used for yeast.

Supports human and mouse

# oPOSSUM (<http://www.cisreg.ca/oPOSSUM>)

## Select Analysis Parameters

STEP 1: Enter a list of co-expressed genes

Species:

human  mouse

Disadvantage:  
Supports either  
human or mouse  
only

Gene ID type:

Ensembl  HUGO/MGI Symbol/Alias  RefSeq  Entrez Gene

Paste gene IDs:

Use sample genes

Clear

259  
5265  
350  
335  
335  
1558

OR upload a file containing a list of gene identifiers:

Browse...

# oPOSSUM (<http://www.cisreg.ca/oPOSSUM>)

STEP 2: Select transcription factor binding site matrices

## JASPAR CORE Profiles

All profiles with a minimum specificity of  bits (min. 8 bits)

**OR** select by taxonomic supergroup:

plant  vertebrate  insect

**OR** select specific profiles:

ABI4  
Agamous  
AGL3  
Ar  
Arnt  
Arnt-Ahr  
ARR10  
Athb-1

## JASPAR PhyloFACTS Profiles

All profiles with a minimum specificity of  bits (min. 8 bits)

The JASPAR PHYLOFACTS database consists of 174 profiles that were extracted from phylogenetically conserved gene upstream elements. They are a mix of known and as of yet undefined motifs.

### When should it be used?

They are useful when one expects that other factors might determine promoter characteristics and/or tissue specificity.

# oPOSSUM (<http://www.cisreg.ca/oPOSSUM>)

## STEP 3: Select parameters

Level of conservation:

Top 10% of conserved regions (min. conservation 70%) ▼

Matrix match threshold:

80 ▼ %

Amount of upstream / downstream sequence:

2000 / 2000 ▼

Number of results to display:

Top 10 ▼ results

OR only results with **Z-score**  $\geq$  10 ▼ and **Fisher score**  $\leq$  0.01 ▼

Sort results by:

Z-score  Fisher score

Press the **Submit** button to perform the analysis or **Reset** to reset the analysis seconds to a minute or more to perform. Please be patient.

Submit

Reset

The Fisher statistic reflects the proportion of genes that contain the TFBS compared to background.

The Z-score statistic reflects the occurrence of the TFBS in the promoters of the co-expressed set compared to background.

# oPOSSUM (<http://www.cisreg.ca/oPOSSUM>)

## Analysis Results

### Selected Parameters

**Conservation level:** Top 10% of conserved  
**Matrix match score:** 80%  
**Upstream sequence length:** 2000  
**Downstream sequence length:** 2000  
**Number of genes submitted:** 21  
**Number of genes included:** 15  
**Number of genes excluded:** 6

### Target Genes

**Analyzed:** 1356 350 2158 259 383 335 3273 462 1571 5105 229 325 2168 2244 5053  
**Excluded:** 5265 1558 125 3240 3827 5004

### oPOSSUM Analysis

TF	TF Class	TF Supergroup	IC	Background gene hits	Background gene non-hits	Target gene hits	Target gene non-hits	Background TFBS hits	Background TFBS rate	Target TFBS hits	Target TFBS rate	Z-score	Fisher score
<a href="#">HNF1A</a>	HOMEO	vertebrate	15.548	1466	13684	<a href="#">8</a>	7	1860	0.0021	<a href="#">8</a>	0.0136	22.69	2.692e-05
<a href="#">SRY</a>	HMG	vertebrate	9.193	8624	6526	<a href="#">13</a>	2	34149	0.0248	<a href="#">33</a>	0.0361	6.567	1.552e-02
<a href="#">Fos</a>	bZIP	vertebrate	10.670	7001	8149	<a href="#">11</a>	4	16086	0.0104	<a href="#">16</a>	0.0155	4.583	3.175e-02
<a href="#">HLF</a>	bZIP	vertebrate	11.147	3376	11774	<a href="#">7</a>	8	5014	0.0048	<a href="#">9</a>	0.0131	10.72	3.196e-02
<a href="#">Foxq1</a>	FORKHEAD	vertebrate	14.070	3533	11617	<a href="#">7</a>	8	6047	0.0054	<a href="#">9</a>	0.0120	8.207	4.026e-02
<a href="#">NKX3-1</a>	HOMEO	vertebrate	11.127	5391	9759	<a href="#">9</a>	6	12155	0.0069	<a href="#">13</a>	0.0110	4.548	4.696e-02
<a href="#">FOXD1</a>	FORKHEAD	vertebrate	11.926	5516	9634	<a href="#">9</a>	6	11145	0.0072	<a href="#">15</a>	0.0146	7.875	5.417e-02
<a href="#">Pdx1</a>	HOMEO	vertebrate	9.040	9899	5251	<a href="#">13</a>	2	54092	0.0261	<a href="#">47</a>	0.0342	4.571	6.515e-02
<a href="#">Cebpa</a>	bZIP	vertebrate	9.187	5863	9287	<a href="#">9</a>	6	12176	0.0118	<a href="#">14</a>	0.0204	7.21	7.844e-02
<a href="#">Nkx2-5</a>	HOMEO	vertebrate	8.270	10169	4981	<a href="#">13</a>	2	59121	0.0333	<a href="#">52</a>	0.0442	5.46	8.496e-02

[Download as a tab delimited text file](#) (results will be kept on the server for 3 days after analysis)

### Genes Containing Conserved HNF1A Binding Sites:

Gene ID	Ensembl ID	Chr	Strand	TSS	Promoter Start	Promoter End	TFBS Sequence	TFBS Start	TFBS Rel. Start	TFBS End	TFBS Rel. End	TFBS Orientation	TFBS Score
1356	<a href="#">ENSG00000047457</a>	3	-1	150422269	150420270	150424269	GGTTAATGTTTAAAT	150421319	951	150421332	938	1	15.334
350	<a href="#">ENSG00000091583</a>	17	-1	61655974	61653975	61657974	GGTTAATGTTTAAAG	61656032	-58	61656045	-71	-1	13.479
3273	<a href="#">ENSG00000113905</a>	3	1	187866487	187864487	187868486	TGTAATGATTAGT	187866344	-143	187866357	-130	-1	9.708
1571	<a href="#">ENSG00000130649</a>	10	1	135190857	135188857	135192856	GGTTTATTATTAGC	135190745	-112	135190758	-99	-1	14.409
5105	<a href="#">ENSG00000124253</a>	20	1	55569543	55567543	55571542	AGATAATCATTGAA	55569396	-147	55569409	-134	-1	9.903
325	<a href="#">ENSG00000132703</a>	1	1	157824239	157822239	157825284	AGTTATTTATTAGA	157824079	-160	157824092	-147	-1	12.759
2168	<a href="#">ENSG00000163586</a>	2	-1	88208693	88206694	88210693	AGTTAATGTTTGAA	88208792	-99	88208805	-112	-1	12.830
2244	<a href="#">ENSG00000171564</a>	4	1	155703596	155701596	155705595	AGTTAATATTTAAAT	155703524	-72	155703537	-59	-1	14.863

[Download as a tab delimited text file](#)

### Genes Containing Conserved SRY Binding Sites:

Gene ID	Ensembl ID	Chr	Strand	TSS	Promoter Start	Promoter End	TFBS Sequence	TFBS Start	TFBS Rel. Start	TFBS End	TFBS Rel. End	TFBS Orientation	TFBS Score
1356	<a href="#">ENSG00000047457</a>	3	-1	150422269	150420270	150424269	TTAAACATT	150421323	947	150421331	939	-1	6.961
				150422269	150420270	150424269	TGACACAAT	150422361	-92	150422369	-100	1	7.793
				150422269	150420270	150424269	TAAAACAAA	150423255	-986	150423263	-994	-1	9.474
350	<a href="#">ENSG00000091583</a>	17	-1	61655974	61653975	61657974	TAATATAAT	61654150	1825	61654158	1817	1	5.862
				61655974	61653975	61657974	AAAACAAA	61654256	1719	61654264	1711	-1	8.914
2158	<a href="#">ENSG00000101981</a>	X	1	138440561	138438561	138442560	TTGGACAAA	138441494	934	138441502	942	1	6.016
383	<a href="#">ENSG00000118520</a>	6	1	131936059	131934059	131938058	ATGAATAAT	131935824	-235	131935832	-227	1	5.865
3273	<a href="#">ENSG00000113905</a>	3	1	187866487	187864487	187868486	TTAATCAAT	187866435	-52	187866443	-44	1	8.775
462	<a href="#">ENSG00000117601</a>	1	-1	172153139	172151140	172155139	TTAAGCAA	172153193	-54	172153201	-62	1	5.779
				172153139	172151140	172155139	TTAAACAAC	172153216	-77	172153224	-85	-1	7.440
1571	<a href="#">ENSG00000130649</a>	10	1	135190857	135188857	135192856	GAAAATAAT	135188983	-1874	135188991	-1866	-1	8.003
				135190857	135188857	135192856	GCTAATAAT	135190745	-112	135190753	-104	1	6.366
5105	<a href="#">ENSG00000124253</a>	20	1	135200555	135198555	135202554	TAAAACATT	135199099	-1456	135199107	-1448	-1	6.342
				55569543	55567543	55571542	GTACACAAA	55569204	-339	55569212	-331	1	8.214
				55569543	55567543	55571542	ATTAACAAC	55569381	-162	55569389	-154	1	6.352
229	<a href="#">ENSG00000136872</a>	9	-1	103237926	103235927	103239926	TCTCACAAT	103237073	854	103237081	846	1	6.965
				103237926	103235927	103239926	GTAATAAA	103237407	520	103237415	512	1	7.334

### Genes Containing Conserved HNF1A Binding Sites:

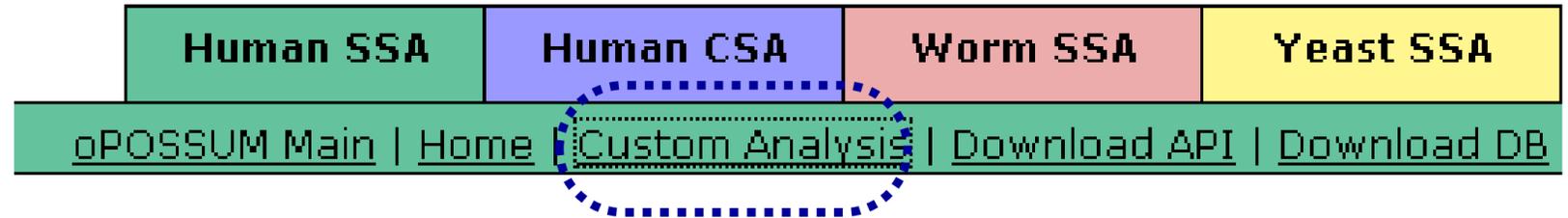
Gene ID	Ensembl ID	Chr	Strand	TSS	Promoter Start	Promoter End	TFBS Sequence	TFBS Start	TFBS Rel. Start	TFBS End	TFBS Rel. End	TFBS Orientation	TFBS Score
1356	<a href="#">ENSG00000047457</a>	3	-1	150422269	150420270	150424269	GGTTAATGTTTAAT	150421319	951	150421332	938	1	15.334
350	<a href="#">ENSG00000091583</a>	17	-1	61655974	61653975	61657974	GGTTAATGTTTAAG	61656032	-58	61656045	-71	-1	13.479
3273	<a href="#">ENSG00000113905</a>	3	1	187866487	187864487	187868486	TGTAAATGATTAGT	187866344	-143	187866357	-130	-1	9.708
1571	<a href="#">ENSG00000130649</a>	10	1	135190857	135188857	135192856	GGTTTATTATTAGC	135190745	-112	135190758	-99	-1	14.409
5105	<a href="#">ENSG00000124253</a>	20	1	55569543	55567543	55571542	AGATAATCATTGAA	55569396	-147	55569409	-134	-1	9.903
325	<a href="#">ENSG00000132703</a>	1	1	157824239	157822239	157825284	AGTTATTTATTAGA	157824079	-160	157824092	-147	-1	12.759
2168	<a href="#">ENSG00000163586</a>	2	-1	88208693	88206694	88210693	AGTTAATGTTTGAA	88208792	-99	88208805	-112	-1	12.830
2244	<a href="#">ENSG00000171564</a>	4	1	155703596	155701596	155705595	AGTTAATATTTAAT	155703524	-72	155703537	-59	-1	14.863

[Download as a tab delimited text file](#)

### Genes Containing Conserved SRY Binding Sites:

Gene ID	Ensembl ID	Chr	Strand	TSS	Promoter Start	Promoter End	TFBS Sequence	TFBS Start	TFBS Rel. Start	TFBS End	TFBS Rel. End	TFBS Orientation	TFBS Score
1356	<a href="#">ENSG00000047457</a>	3	-1	150422269	150420270	150424269	TTAAACATT	150421323	947	150421331	939	-1	6.961
				150422269	150420270	150424269	TGACACAAT	150422361	-92	150422369	-100	1	7.793
				150422269	150420270	150424269	TAAAACAAA	150423255	-986	150423263	-994	-1	9.474
350	<a href="#">ENSG00000091583</a>	17	-1	61655974	61653975	61657974	TAATATAAT	61654150	1825	61654158	1817	1	5.862
				61655974	61653975	61657974	AAAAACAAA	61654256	1719	61654264	1711	-1	8.914
2158	<a href="#">ENSG00000101981</a>	X	1	138440561	138438561	138442560	TTGGACAAA	138441494	934	138441502	942	1	6.016
383	<a href="#">ENSG00000118520</a>	6	1	131936059	131934059	131938058	ATGAATAAT	131935824	-235	131935832	-227	1	5.865
3273	<a href="#">ENSG00000113905</a>	3	1	187866487	187864487	187868486	TTAATCAAT	187866435	-52	187866443	-44	1	8.775
462	<a href="#">ENSG00000117601</a>	1	-1	172153139	172151140	172155139	TTAAGCAA	172153193	-54	172153201	-62	1	5.779
				172153139	172151140	172155139	TTAAACAAC	172153216	-77	172153224	-85	-1	7.440
1571	<a href="#">ENSG00000130649</a>	10	1	135190857	135188857	135192856	GAAAATAAT	135188983	-1874	135188991	-1866	-1	8.003
				135190857	135188857	135192856	GCTAATAAT	135190745	-112	135190753	-104	1	6.366
				135200555	135198555	135202554	TAAAACATT	135199099	-1456	135199107	-1448	-1	6.342
5105	<a href="#">ENSG00000124253</a>	20	1	55569543	55567543	55571542	GTACACAAA	55569204	-339	55569212	-331	1	8.214
				55569543	55567543	55571542	ATTAACAAC	55569381	-162	55569389	-154	1	6.352
229	<a href="#">ENSG00000136872</a>	9	-1	103237926	103235927	103239926	TCTCACAAAT	103237073	854	103237081	846	1	6.965
				103237926	103235927	103239926	GTAATAAAA	103237407	520	103237415	512	1	7.334

# oPOSSUM (<http://www.cisreg.ca/oPOSSUM>)



## Select Custom Analysis Parameters

STEP 1a: Enter a list of co-expressed genes

Species:

human  mouse

Gene ID type:

Ensembl  HUGO/MGI Symbol/Alias  RefSeq  Entrez Gene

Paste gene IDs (max. 1000 genes):

259  
5265  
350  
335  
335  
1558

OR upload a file containing a list of gene identifiers:

STEP 1b: Enter a background list of genes

Use background set of 1000 random genes

OR Paste background gene IDs (max. 1000 genes):

1281  
1281  
1805  
125  
10551

OR upload a file containing a list of gene identifiers:

oPOSSUM Analysis

TF	TF Class	TF Supergroup	IC	Background gene hits	Background gene non-hits	Target gene hits	Target gene non-hits	Background TFBS hits	Background TFBS rate	Target TFBS hits	Target TFBS rate	Z-score	Fisher score
<a href="#">HNF1A</a>	HOMEO	vertebrate	15.548	1	10	8	7	1	0.0025	8	0.0136	20.32	2.426e-02
<a href="#">HLF</a>	bZIP	vertebrate	11.147	1	10	7	8	1	0.0021	9	0.0131	21.68	4.943e-02
<a href="#">NKX3-1</a>	HOMEO	vertebrate	11.127	3	8	9	6	3	0.0037	13	0.0110	10.95	1.042e-01
<a href="#">Bapx1</a>	HOMEO	vertebrate	8.542	4	7	10	5	5	0.0079	20	0.0218	14.26	1.286e-01
<a href="#">Lhx3</a>	HOMEO	vertebrate	12.941	3	8	8	7	5	0.0079	11	0.0120	4.173	1.775e-01
<a href="#">Pdx1</a>	HOMEO	vertebrate	9.040	7	4	13	2	28	0.0295	47	0.0342	2.533	1.826e-01
<a href="#">SRY</a>	HMG	vertebrate	9.193	7	4	13	2	26	0.0410	33	0.0361	-2.303	1.826e-01
<a href="#">Nkx2-5</a>	HOMEO	vertebrate	8.270	7	4	13	2	28	0.0344	52	0.0442	4.865	1.826e-01
<a href="#">FOXI1</a>	FORKHEAD	vertebrate	13.183	4	7	9	6	8	0.0168	17	0.0248	5.555	2.142e-01
<a href="#">RORA_1</a>	NUCLEAR RECEPTOR	vertebrate	13.190	1	10	4	11	1	0.0018	5	0.0061	9.233	2.739e-01

# Exercise 1: Use oPOSSUM to find shared conserved cis-elements in a group of co-expressed genes

1. Download the example dataset (file “Example-Set-1.xls” – rt. click and “save as” from <http://anil.cchmc.org/dhc.html>)
2. Copy 20 or 25 gene IDs from the downloaded file and use them for oPOSSUM analysis

## oPOSSUM Summary:

1. For conserved common cis-elements in a group of genes
2. Supports human or mouse only
3. Uses JASPAR matrices only which are not exhaustive
4. Options to select the regions (max. 10 kb flanking region)
5. Results indicate the TFBSs' positions relative to the TSS and the coordinates are from the current genome assembly
6. Supports selection of background set
7. Does not support upload of your sequences; Input should be standard gene symbols or IDs or accession numbers

# oPOSSUM Summary:

## 4. Options to select the regions (max. 10 kb flanking region)

# DiRE (<http://dire.dcode.org/>)

**DiRE**  
Distant Regulatory Elements of co-regulated genes

Home Details Output example Screenshots Return to submitted job... Citing DiRE Contact us

April 30, 2008. The DiRE tool operates with new ECR Browser alignments now. It is also possible to run DiRE on human (hg18), mouse (mm9), and rat (rn4) genomes.

**Co-regulated genes:**  
Copy and paste gene names (or accession numbers) (example)  
AMBP  
SERPINA1  
APOH  
APOA1  
APOA1  
CYP2C8  
CYP2E1  
ALDOB  
SERPINC1  
ADH1B  
HP  
PCK1  
SERPINA1  
HRG  
FGB  
F9

**Background (control) genes:**  
Select the source of background genes:  
 random set of 100 genes  
 copy and paste background genes

**Target elements:**  
 top 3 ECRs + promoter ECRs [default]  
 UTR ECRs + promoter ECRs  
 promoter ECRs only

Gene annotation: Gene symbols (GATA2, PAX6, etc)  
Genome: human (hg18)

**Submit**

job ID: 0926091038152309  
22 signal genes; 100 background genes  
Identifying distant regulatory elements...  
parameter optimization: 82%

Automatic status update every 3 seconds.  
Click here to manually refresh the page.

Request ID... 0926091006373607  
perm link: <http://dire.dcode.org/?id=0926091006373607>

**68 Potential Regulatory Elements...**  
intergenic 40 (59%)  
promoter 10 (15%)  
utr 12 (18%)  
intron 6 (9%)

Detailed description of regulatory elements (in tabulated textual format)  
Chromosomal distribution

**Candidate Transcription Factors...**  
10 top TFs  
occurrence importance

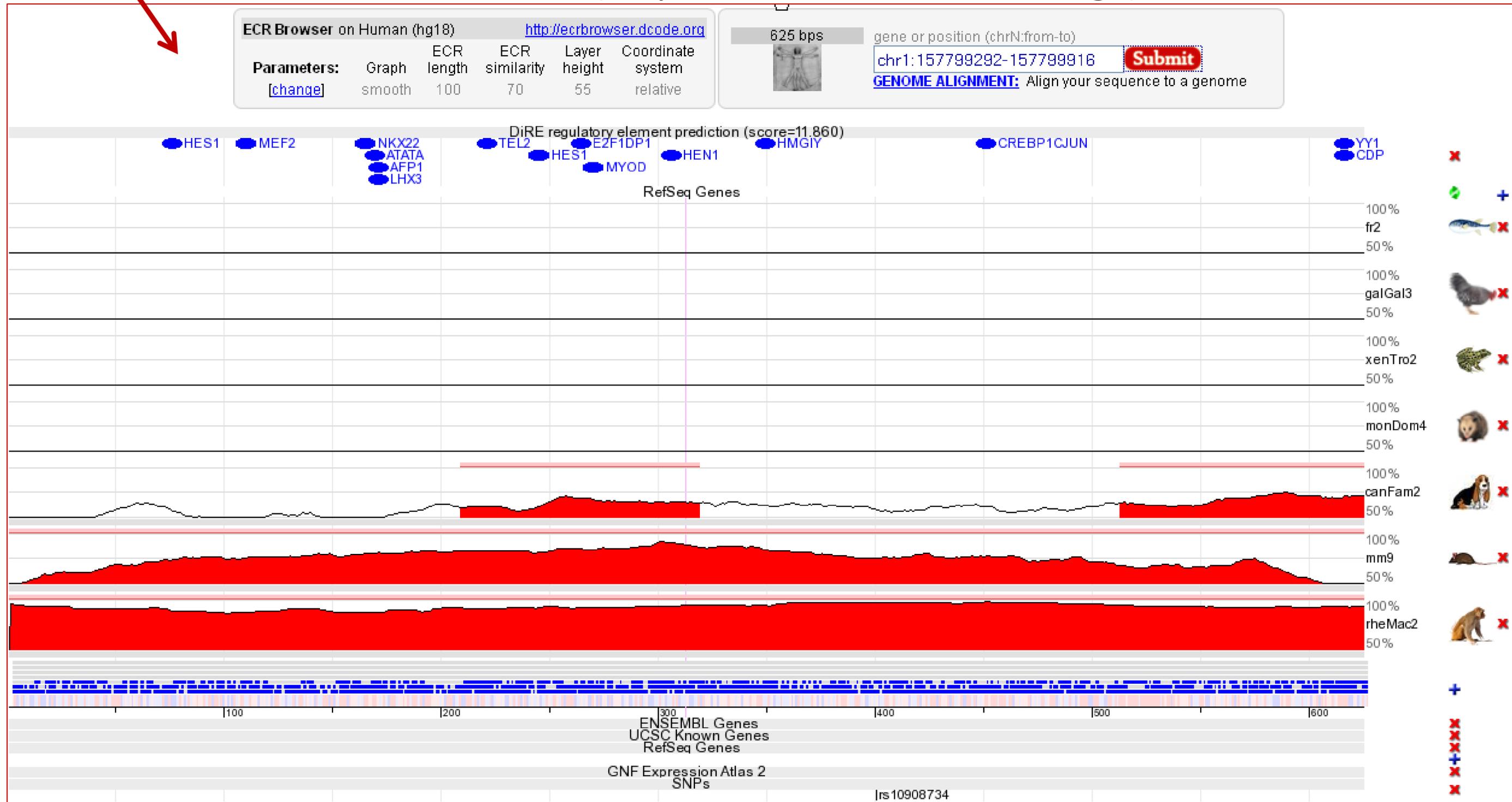
#	Transcription Factor	Occurrence	Importance
1	NRSF	14.29%	0.36786
2	VMAF	10.00%	0.34781
3	HSF2	2.86%	0.28571
4	YY1	17.14%	0.27750
5	ALX4	7.14%	0.27679
6	IRF7	5.71%	0.19643
7	RUSH1A	5.71%	0.19643
8	ISRE	10.00%	0.19375
9	R	5.71%	0.17661
10	HMG1Y	11.43%	0.16571
11	CDP	8.57%	0.16286
12	LXR	5.71%	0.16161
13	POU6F1	5.71%	0.15732
14	PPARG	15.71%	0.15714
15	CEBP	21.43%	0.15000
16	STAT	11.43%	0.14857

#	Regulatory element	Type	Score	Locus	Gene	Candidate transcription factor binding sites (relative positions)
1	chr1:157799292-157799916	intergenic	11.860	chr1:157772437-157948651	APCS	15 :: HES1(71) MEF2(105) NKX22(160) ATATA(164) AFP1(166) LHX3(166) TEL2(216) HES1(240) E2F1DP1(259) MYOD(265) HEN1(301) HMG1Y(344) CREBP1CJUN(446) YY1(611) CDP(611)
2	chr1:157822774-157823321	promoter	1.327	chr1:157772437-157948651	APCS	4 :: AREB6(60) TBX5(62) CHX10(365) TBX5(520)
3	chr1:157823477-157823620	promoter	0.843	chr1:157772437-157948651	APCS	3 :: ARP1(12) LEF1(20) CMAF(53)
4	chr1:157823967-157824176	promoter	6.780	chr1:157772437-157948651	APCS	8 :: STAT3(62) RORA2(76) PPARA(78) LXR(82) ER(83) PPARA(83) LXR_DR4(83) HNF1(117) 39 :: PAX2(56) STAT3(65) OSF2(106) PEBP(107) AML1(107) RFX1(340) NFMUE1(396) TAL1BETAITE2(398) CHOP(403) YY1(507) GATA1(598) CDP(611) CLOX(611) NRSF(649)

# DiRE (<http://dire.dcode.org/>)

#	Regulatory element	Type	Score	Locus	Gene	Candidate transcription factor binding sites (relative positions)
1	<a href="#">chr1:157799292-157799916</a>	intergenic	11.860	chr1:157772437-157948651	APCS	15 :: HES1(71) MEF2(105) NKX22(160) ATATA(164) AFP1(166) LHX3(166) TEL2(216) HES1(240) E2F1DP1(259) MYOD(265) HEN1(301) HMG1Y(344) CREBP1CJUN(446) YY1(611) CDP(611)
2	<a href="#">chr1:157822774-157823321</a>	promoter	1.327	chr1:157772437-157948651	APCS	4 :: AREB6(60) TBX5(62) CHX10(365) TBX5(520)
3	<a href="#">chr1:157823477-157823620</a>	promoter	0.843	chr1:157772437-157948651	APCS	3 :: ARP1(12) LEF1(20) CMAF(53)
4	<a href="#">chr1:157823967-157824176</a>	promoter	6.780	chr1:157772437-157948651	APCS	8 :: STAT3(62) RORA2(76) PPARA(78) LXR(82) ER(83) PPARA(83) LXR_DR4(83) HNF1(117) 39 :: PAX2(56) STAT3(65) OSF2(106) PEBP(107) AML1(107) RFX1(340) NFMUE1(396) TAL1BETAITE2(398) CHOP(403) YY1(507) GATA4(598) CDP(611) CLOX(611) NRSE(649)

## ECR-Browser (<http://ecrbrowser.dcode.org/>)



# Exercise 2: Use DiRE to find shared conserved cis-elements in a group of co-expressed genes

Use the same example dataset (downloaded file “Example-Set-1.xls”) and identify putative distant regulatory regions using DiRE

## DiRE Summary:

1. DiRE's unique feature is the detection of conserved REs outside of proximal promoter regions, as it takes advantage of the full gene locus to conduct the search.
2. Supports human, mouse, and rat
3. Uses TRANSFAC matrices which are more exhaustive than JASPAR matrices
4. Limited options to select the regions for scanning
5. Results indicate the context (promoter, intronic, or UTR, etc.) and the coordinates are from the current genome assembly
6. Supports selection of background set
7. Does not support upload of your sequences; Input should be standard gene symbols or IDs or accession numbers
8. Connects to genome browser

# Pscan (<http://159.149.109.9/pscan>)

Insert Gene/Sequence ID list: ([help](#)) **PSCAN**

Select Organism:

Select Region:

Select Descriptors:

- Jaspar
- Jaspar\_Fam
- Transfac
- User Defined

Messages:



## Pscan Web Interface

Use the input form on the left to set up your query. The results will be displayed in this window.

[If you need HELP please click here.](#)

**Source:**  
[Download Pscan source code](#)

**Reference:**  
F.Zambelli, G.Pesole, G.Pavesi  
[Pscan: Finding Over-represented Transcription Factor Binding Site Motifs in Sequences from Co-Regulated or Co-Expressed Genes.](#)  
*Nucleic Acids Research* 2009 37(Web Server issue):W247-W252.

**Contacts:**  
[giulio.pavesi@unimi.it](mailto:giulio.pavesi@unimi.it)  
[federico.zambelli@unimi.it](mailto:federico.zambelli@unimi.it)

### Sample data

List of MYC target genes. MYCxx indicates that xx percent of the genes in the list are MYC targets, while the others are random genes added to the set to assess the performance of the algorithm.

[MYC100](#) [MYC90](#) [MYC80](#) [MYC75](#) [MYC65](#)  
[MYC55](#)

---

List of NFkB target genes, collected from literature. NFkBxx should be read as in the MYC dataset.

[NFkB100](#) [NFkB90](#) [NFkB80](#) [NFkB70](#) [NFkB60](#)  
[NFkB50](#) [NFkB40](#)

---

List of NRF1 target genes. NRFxx should be read as in the MYC dataset. Use the NRF1 matrix with the link provided below to test these datasets (save the matrix as a text file).

[NRF1\\_100](#) [NRF1\\_90](#) [NRF1\\_80](#) [NRF1\\_70](#)  
[NRF1\\_60](#) [NRF1\\_50](#) [NRF1\\_40](#)

[NRF1 Matrix](#)

# Pscan (<http://159.149.109.9/pscan>)

Insert Gene/Sequence ID list: ([help](#)) **PSCAN**

NM\_006408  
NM\_006418  
NM\_006439  
NM\_006475  
NM\_001285  
NM\_000668  
NM\_000667  
NM\_000669  
NM\_000668

Select Organism:

Select Region:

Select Descriptors:  
 Jaspar  
 Jaspar\_Fam  
 Transfac  
 User Defined

Messages:

6 (out of 84) gene ID(s) not found:  
NM\_138298  
NM\_138299  
NM\_024416  
XM\_936565  
XM\_941953  
XM\_930062

Working on 78 gene promoter(s).

Select Organism:

Select Region:

Select Descriptors:

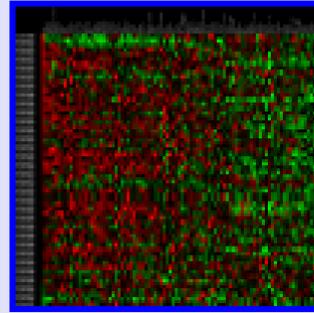
Jaspar  
 Jaspar\_Fam  
 Transfac  
 User Defined

# Pscan (<http://159.149.109.9/pscan>)

[View Text Results](#)

97 TF profiles used

Matrix Name	P-value
<a href="#">TBP</a>	1.59074e-08
<a href="#">Foxa2</a>	0.000274079
<a href="#">FOXL1</a>	0.000657034
<a href="#">MEF2A</a>	0.000657227
<a href="#">Hand1-Tcfe2a</a>	0.000697277
<a href="#">Nobox</a>	0.000790445
<a href="#">FOXJ1</a>	0.000804377
<a href="#">PBX1</a>	0.00124224
<a href="#">SRF</a>	0.00124647
<a href="#">Evi1</a>	0.00128699
<a href="#">TEAD1</a>	0.00212538
<a href="#">Lhx3</a>	0.00303459
<a href="#">Foxq1</a>	0.00355502
<a href="#">Prx2</a>	0.00486451
<a href="#">Lhx3</a>	0.00527407
<a href="#">NKX3-1</a>	0.00590862
<a href="#">NFIL3</a>	0.00642618
<a href="#">REL</a>	0.00685234
<a href="#">Pax6</a>	0.00765503
<a href="#">Foxd3</a>	0.00776631
<a href="#">HNF1A</a>	0.00783389
<a href="#">Cebpa</a>	0.00920516
<a href="#">Nkx2-5</a>	0.00940039



## Matrix Info

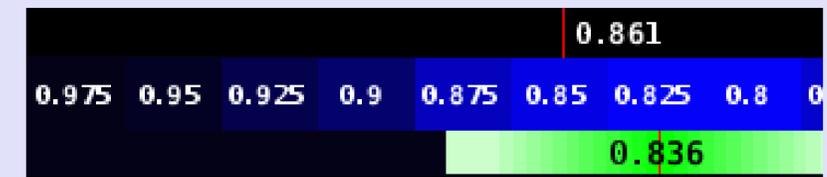
ID	<a href="#">MA0047</a>
Name	Foxa2
Class	FORKHEAD
Species	Rattus norvegicus
Inf. Content	12.43
SuperGroup	vertebrate
Protein Acc.	<a href="#">P32182</a>
Type	COMPILED
PMID	<a href="#">8139574</a>
Report Occurrences	<input type="button" value="Go!"/>

## MA0047

	1	2	3	4	5	6	7	8	9	10	11	12
<b>A</b>	6	11	8	0	11	0	0	0	9	0	1	0
<b>C</b>	7	3	3	1	0	1	1	0	0	9	1	5
<b>G</b>	3	1	1	0	6	0	0	6	8	0	1	0
<b>T</b>	1	2	5	16	0	16	16	11	0	8	14	12



## Sample Mean Score



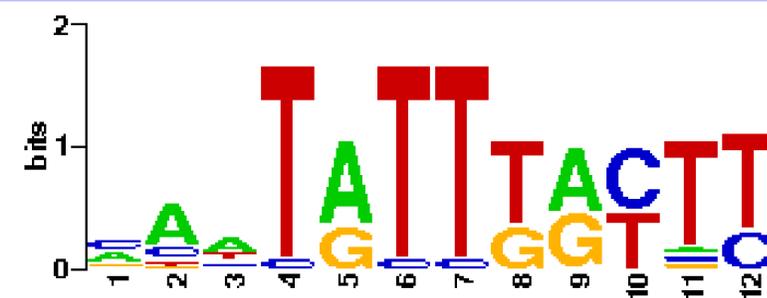
## Background Score Distribution

## Sample Statistics

p-value	0.000274079
Bonferroni p-value	0.026585663
Mean	0.861172
Std Dev	0.0563177
Size	60

## Compare with... (using Welch's t-test) [help](#)

Mean	<input type="text"/>	<input type="button" value="Go!"/>
Std Dev	<input type="text"/>	
Size	<input type="text"/>	

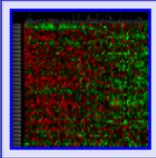


# Pscan (<http://159.149.109.9/pscan>)

[View Text Results](#)

97 TF profiles used

Matrix Name	P-value
TBP	1.59074e-08
Foxa2	0.000274079
FOXL1	0.000657034
MEF2A	0.000657227
Hand1-Tcfe2a	0.000697277
Nobox	0.000790445
FOX1	0.000804377
PBX1	0.00124224
SRF	0.00124647
Evi1	0.00128699
TEAD1	0.0012538
Lhx3	0.00300459
Foxq1	0.0035502
Prrx2	0.00486451
Lhx3	0.00527407
NKX3-1	0.00590862
NFIL3	0.00642618
REL	0.00685234
Pax6	0.00765503
Foxd3	0.00776631
HNF1A	0.00783389
Cebpa	0.00920516
Nkx2-5	0.00940039



[View Text Results](#)

Name	Score	Position	Sequence	Strand
<a href="#">hg18 refGene NM 002345</a>	0.98343	-197	CAATATTGATTT	-
<a href="#">hg18 refGene NM 000668</a>	0.982619	-145	AAATATTGACTT	-
<a href="#">hg18 refGene NM 006408</a>	0.951013	-258	CTTTATTTACTT	-
<a href="#">hg18 refGene NM 000667</a>	0.944804	-34	ATTTATTTATTT	-
<a href="#">hg18 refGene NM 000609</a>	0.939791	-428	ACTTGTTTGCTT	+
<a href="#">hg18 refGene NM 001033888</a>	0.939791	-428	ACTTGTTTGCTT	+
<a href="#">hg18 refGene NM 199168</a>	0.939791	-428	ACTTGTTTGCTT	+
<a href="#">hg18 refGene NM 021010</a>	0.935506	-261	GAGTATTTACTT	-
<a href="#">hg18 refGene NM 133477</a>	0.932417	-356	AAACATTTATTT	+
<a href="#">hg18 refGene NM 194435</a>	0.931745	-216	CTTTGTTTGTTT	+
<a href="#">hg18 refGene NM 003381</a>	0.931745	-216	CTTTGTTTGTTT	+
<a href="#">hg18 refGene NM 005603</a>	0.922637	-143	GAATATTTACAT	+
<a href="#">hg18 refGene NM 004616</a>	0.9222	-56	ATCTGTTTACTT	+
<a href="#">hg18 refGene NM 004295</a>	0.919060	220	ATTTATTTACTT	+

### Matrix Info

ID	<a href="#">MA0047</a>
Name	Foxa2
Class	FORKHEAD
Species	Rattus norvegicus
Inf. Content	12.43
SuperGroup	vertebrate
Protein Acc.	<a href="#">P32182</a>
Type	COMPILED
PMID	<a href="#">8139574</a>
Report Occurrences	<input type="button" value="Go!"/>

**MA0047**

	1	2	3	4	5	6	7	8	9	10	11	12
A	6	11	8	0	11	0	0	0	9	0	1	0
C	7	3	3	1	0	1	1	0	0	9	1	5
G	3	1	1	0	0	0	0	6	8	0	1	0
T	1	2	5	10	0	16	16	11	0	8	14	12

Sample Mean Score: 0.861

Background Score Distribution: 0.975 0.95 0.925 0.9 0.875 0.85 0.825 0.8 0.775 0.75 0.725 0.7 0.675 0.65 0.625 0.6 0.575 0.55 0.525 0.5 0.475 0.45 0.425 0.4 0.375 0.35 0.325 0.3 0.275 0.25 0.225 0.2 0.175 0.15 0.125 0.1 0.075 0.05 0.025 0 0

Background Score Distribution: 0.836

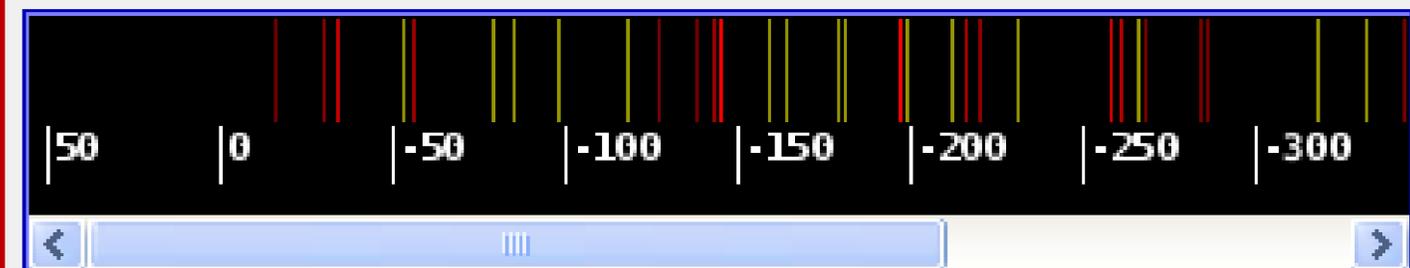
**Sample Statistics**

p-value	0.000274079
Bonferroni p-value	0.026585663
Mean	0.861172
Std Dev	0.0563177
Size	60

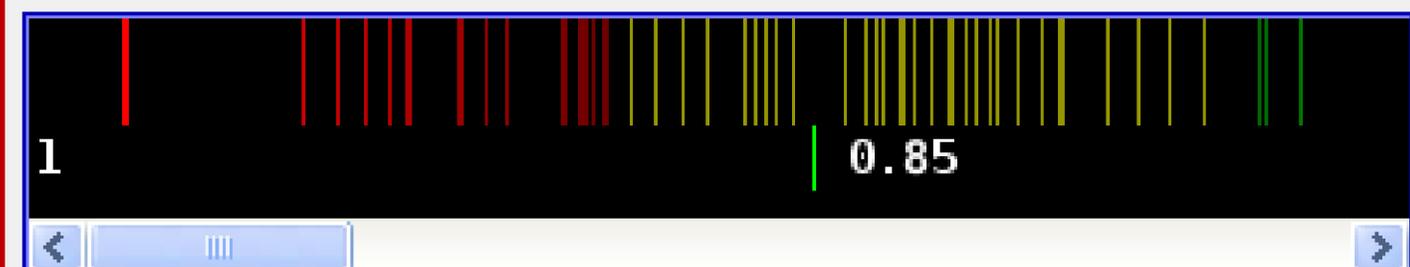
**Compare with... (using Welch's t-test) [help](#)**

Mean	<input type="text"/>	<input type="button" value="Go!"/>
Std Dev	<input type="text"/>	
Size	<input type="text"/>	

Occurrences Position Distribution (score  $\geq 0.836$ )



Occurrences Score Distribution

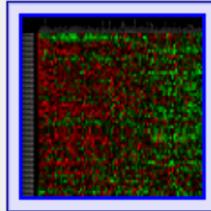


# Pscan (<http://159.149.109.9/pscan>)

[View Text Results](#)

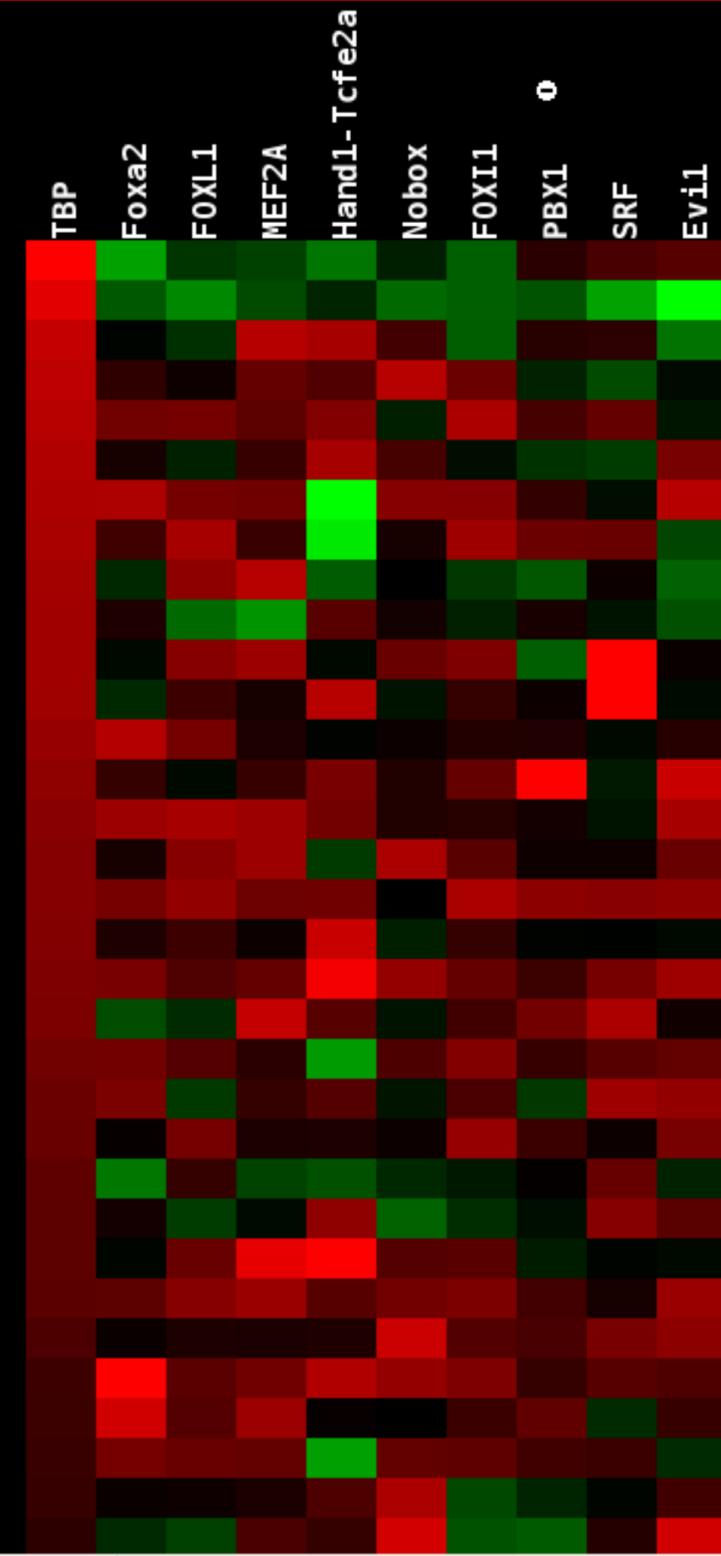
97 TF profiles used

Matrix Name	P-value
<a href="#">TBP</a>	1.59074e-08
<a href="#">Foxa2</a>	0.000274079
<a href="#">FOXL1</a>	0.000657034
<a href="#">MEF2A</a>	0.000657227
<a href="#">Hand1-Tcfe2a</a>	0.000697277
<a href="#">Nobox</a>	0.000790445
<a href="#">FOXI1</a>	0.000804377
<a href="#">PBX1</a>	0.00124224
<a href="#">SRF</a>	0.00124647
<a href="#">Evi1</a>	0.00128699
<a href="#">TEAD1</a>	0.00212538
<a href="#">Lhx3</a>	0.00303459
<a href="#">Foxq1</a>	0.00355502
<a href="#">Prrx2</a>	0.00486451
<a href="#">Lhx3</a>	0.00527407
<a href="#">NKX3-1</a>	0.00590862
<a href="#">NFIL3</a>	0.00642618
<a href="#">REL</a>	0.00685234
<a href="#">Pax6</a>	0.00765503
<a href="#">Foxd3</a>	0.00776631
<a href="#">HNF1A</a>	0.00783389
<a href="#">Cebpa</a>	0.00920516
<a href="#">Nkx2-5</a>	0.00940039



```

>hg18_refGene_NM_000088
>hg18_refGene_NM_058175
>hg18_refGene_NM_022844
>hg18_refGene_NM_002354
>hg18_refGene_NM_001937
>hg18_refGene_NM_207373
>hg18_refGene_NM_194435
>hg18_refGene_NM_001443
>hg18_refGene_NM_032413
>hg18_refGene_NM_013372
>hg18_refGene_NM_001613
>hg18_refGene_NM_001615
>hg18_refGene_NM_021010
>hg18_refGene_NM_000900
>hg18_refGene_NM_005603
>hg18_refGene_NM_199512
>hg18_refGene_NM_000090
>hg18_refGene_NM_003013
>hg18_refGene_NM_000587
>hg18_refGene_NM_006274
>hg18_refGene_NM_002667
>hg18_refGene_NM_006418
>hg18_refGene_NM_006439
  
```



TF_NAME	MATRIX_ID	Z_SCORE	P_VALUE	SAMPLE_AVERAGE	BACKGROUND_AVERAGE	SAMPLE_DEVSTD	SAMPLE_SIZE
TBP	MA0108	5.52625	1.59074e-08	0.859494	0.816811	0.0493175	60
Foxa2	MA0047	3.45141	0.000274079	0.861172	0.836541	0.0563177	60
FOXL1	MA0033	3.21019	0.000657034	0.918922	0.891892	0.0523769	60
MEF2A	MA0052	3.20941	0.000657227	0.810217	0.777574	0.0710495	60
Hand1-Tcfe2a	MA0092	3.19081	0.000697277	0.895392	0.879315	0.0428231	60
Nobox	MA0125	3.15645	0.000790445	0.880691	0.854103	0.0553916	60
FOXI1	MA0042	3.15094	0.000804377	0.85506	0.825768	0.0653533	60
PBX1	MA0070	3.02199	0.00124224	0.802727	0.781876	0.048527	60
SRF	MA0083	3.02072	0.00124647	0.759912	0.740766	0.0462187	60
Evi1	MA0029	3.01105	0.00128699	0.769179	0.746618	0.0546341	60
TEAD1	MA0090	2.85352	0.00212538	0.827634	0.810446	0.0565813	60
Lhx3	MA0134	2.7417	0.00303459	0.834887	0.804852	0.0683236	60
Foxq1	MA0040	2.68955	0.00355502	0.825721	0.802203	0.0522346	60
Prrx2	MA0075	2.58254	0.00486451	0.92434	0.89143	0.0897038	60
Lhx3	MA0135	2.55439	0.00527407	0.798641	0.773954	0.0704547	60
NKX3-1	MA0124	2.51437	0.00590862	0.853125	0.829214	0.0728452	60
NFIL3	MA0025	2.48463	0.00642618	0.80459	0.783263	0.0663844	60
REL	MA0101	2.46099	0.00685234	0.8843	0.869993	0.0500624	60
Pax6	MA0069	2.42112	0.00765503	0.778947	0.766632	0.0436411	60
Foxd3	MA0041	2.4171	0.00776631	0.857889	0.837233	0.0604468	60
HNF1A	MA0046	2.41362	0.00783389	0.790719	0.770942	0.0615964	60

# Pscan (<http://159.149.109.9/pscan>)

## Comparing different input gene sets:

1. In the detailed output for a given matrix, you can compare the results obtained with the matrix on the gene set just submitted with the results the matrix had produced on another gene set. The latter could be a "negative" gene set (or vice versa).
2. To perform the comparison, you have to fill in the "Compare with..." box fields with mean, standard deviation and sample size values of the other analysis - for the current one you can find them in the "Sample Data Statistics" box or in the overall text output that can be downloaded from the main output page.
3. **Warning:** Make sure that the values you input are correct, and especially that they were obtained by using the same matrix. Once you have clicked the "Go!" button, an output window will pop up and report if either of the two means is significantly higher than the other, together with a confidence p-value computed with a Welch t-test.

Matrix Info	
ID	MA0047
Name	Foxa2
Class	FORKHEAD
Species	Rattus norvegicus
Inf. Content	12.43
SuperGroup	vertebrate
Protein Acc.	P32182
Type	COMPILED
PMID	8139574
Report Occurrences	<input type="button" value="Go!"/>

MA0047												
	1	2	3	4	5	6	7	8	9	10	11	12
A	6	11	8	0	11	0	0	0	9	0	1	0
C	7	3	3	1	0	1	1	0	0	9	1	5
G	3	1	1	0	6	0	0	6	8	0	1	0
T	1	2	5	16	0	16	16	11	0	8	14	12

Sample Mean Score	
	0.861
	0.975 0.95 0.925 0.9 0.875 0.85 0.825 0.8 0
	0.836

Sample Statistics	
p-value	0.000274079
Bonferroni p-value	0.026585663
Mean	0.861172
Std Dev	0.0563177
Size	60

Compare with... (using Welch's t-test) <a href="#">help</a>	
Mean	<input type="text"/>
Std Dev	<input type="text"/>
Size	<input type="text"/>
<input type="button" value="Go!"/>	

Matrix Info	
ID	MA0047
Name	Foxa2
Class	FORKHEAD
Species	Rattus norvegicus
Inf. Content	12.43
SuperGroup	vertebrate
Protein Acc.	P32182
Type	COMPILED
PMID	8139574
Report Occurrences	<input type="button" value="Go!"/>

MA0047												
	1	2	3	4	5	6	7	8	9	10	11	12
A	6	11	8	0	11	0	0	0	9	0	1	0
C	7	3	3	1	0	1	1	0	0	9	1	5
G	3	1	1	0	6	0	0	6	8	0	1	0
T	1	2	5	16	0	16	16	11	0	8	14	12

Sample Mean Score	
	0.861
	0.975 0.95 0.925 0.9 0.875 0.85 0.825 0.8 0
	0.836

Sample Statistics	
p-value	0.000274079
Bonferroni p-value	0.026585663
Mean	0.861172
Std Dev	0.0563177
Size	60

Compare with... (using Welch's t-test) <a href="#">help</a>	
Mean	<input type="text"/>
Std Dev	<input type="text"/>
Size	<input type="text"/>
<input type="button" value="Go!"/>	

Sample Statistics	
Mean	0.815468
Std Dev	0.0541239
Size	334

Compare with... (using Welch's t-test) <a href="#">help</a>	
Mean	0.83
Std Dev	0.07
Size	250
<input type="button" value="Go!"/>	

### Welch's t-test

Sample 1 Mean = 0.815468  
Sample 2 Mean = 0.83

t	-2.72829142951
Degrees of Freedom	453.764151261
Sample 2 is more enriched than Sample 1 (this one) with a p-value confidence of	<b>0.003307</b>

# Exercise 3: Use Pscan to find shared cis-elements (Transfac) in a group of co-expressed genes

1. Use the same example data set (downloaded file “Example-Set-1.xls”) and find the enriched JASPAR and Transfac TFBS. How do the outputs differ?

## PScan Summary:

1. Pscan supports a variety of TFBS matrices (e.g. JASPAR, Transfac) including user input matrix.
2. Supports human, mouse, drosophila, and yeast
3. Limited options to select the regions for scanning
4. Cannot select the background set although comparisons can be computed
5. Does not support upload of your sequences; Input options are very limited
6. Variety of user-friendly output formats including heat map view

# I have a list of co-expressed mRNAs (Transcriptome)....

## I want to find the shared cis-elements – Known and Novel

### □ Known transcription factor binding sites (TFBS)

#### ❖ Conserved

- oPOSSUM
- DiRE

#### ❖ Non-conserved

- Pscan
- **MatInspector** (\*Licensed)

### □ Unknown TFBS or Novel motifs

#### ❖ Conserved

- oPOSSUM
- **Weeder-H**

#### ❖ Non-conserved

- **MEME**
- **Weeder**

# oPOSSUM (<http://www.cisreg.ca/oPOSSUM>)

## Select Analysis Parameters

STEP 1: Enter a list of co-expressed genes

**Species:**

human  mouse

**Gene ID type:**

Ensembl  HUGO/MGI Symbol/Alias  RefSeq  Entrez Gene

Paste gene IDs:

Use sample genes

Clear

259

5265

350

335

335

1558

**OR** upload a file containing a list of gene identifiers:

Browse...

# oPOSSUM (<http://www.cisreg.ca/oPOSSUM>)

## STEP 2: Select transcription factor binding site matrices

### JASPAR CORE Profiles

All profiles with a minimum specificity of  bits (min. 8 bits)

**OR** select by taxonomic supergroup:

plant  vertebrate  insect

**OR** select specific profiles:



The JASPAR PHYLOFACTS database consists of 174 profiles that were extracted from phylogenetically conserved gene upstream elements. They are a mix of known and as of yet undefined motifs.

### When should it be used?

They are useful when one expects that other factors might determine promoter characteristics and/or tissue specificity.

### JASPAR PhyloFACTS Profiles

All profiles with a minimum specificity of  bits (min. 8 bits)

# oPOSSUM (<http://www.cisreg.ca/oPOSSUM>)

## STEP 3: Select parameters

Level of conservation:

Top 10% of conserved regions (min. conservation 70%) ▼

Matrix match threshold:

80 ▼ %

Amount of upstream / downstream sequence:

2000 / 2000 ▼

Number of results to display:

Top 10 ▼ results

OR only results with **Z-score**  $\geq$  10 ▼ and **Fisher score**  $\leq$  0.01 ▼

Sort results by:

Z-score  Fisher score

Press the **Submit** button to perform the analysis or **Reset** to reset the analysis. It may take several seconds to a minute or more to perform. Please be patient.

Submit

Reset

# oPOSSUM (<http://www.cisreg.ca/oPOSSUM>)

## oPOSSUM Analysis

TF	TF Class	TF Supergroup	IC	Background gene hits	Background gene non-hits	Target gene hits	Target gene non-hits	Background TFBS hits	Background TFBS rate	Target TFBS hits	Target TFBS rate	Z-score	Fisher score
<a href="#">RYTAAWNNNTGAY</a>	Unknown	mammals	16.655	2636	12514	<a href="#">9</a>	6	3909	0.0041	<a href="#">15</a>	0.0237	27.75	2.676e-04
<a href="#">TAATTA</a>	Unknown	mammals	12.000	7400	7750	<a href="#">13</a>	2	27227	0.0132	<a href="#">31</a>	0.0226	7.458	2.848e-03
<a href="#">TATAAA</a>	Unknown	mammals	12.000	9219	5931	<a href="#">14</a>	1	47951	0.0232	<a href="#">47</a>	0.0342	6.638	6.209e-03
<a href="#">RGTTAMWNATT</a>	Unknown	mammals	17.072	2277	12873	<a href="#">6</a>	9	3189	0.0028	<a href="#">8</a>	0.0107	13.33	1.705e-02
<a href="#">RTAAACA</a>	Unknown	mammals	13.000	7918	7232	<a href="#">12</a>	3	25209	0.0142	<a href="#">29</a>	0.0246	7.953	2.670e-02
<a href="#">YATTNATC</a>	Unknown	mammals	13.061	6858	8292	<a href="#">11</a>	4	18528	0.0119	<a href="#">19</a>	0.0185	5.394	2.682e-02
<a href="#">CTTTGA</a>	Unknown	mammals	12.000	10591	4559	<a href="#">14</a>	1	54148	0.0262	<a href="#">47</a>	0.0342	4.553	3.478e-02
<a href="#">YCATTAA</a>	Unknown	mammals	13.004	7484	7666	<a href="#">11</a>	4	22958	0.0129	<a href="#">22</a>	0.0187	4.57	5.404e-02
<a href="#">AACWWCAANK</a>	Unknown	mammals	15.858	3060	12090	<a href="#">6</a>	9	4631	0.0037	<a href="#">8</a>	0.0097	8.817	6.381e-02
<a href="#">TGGAAA</a>	Unknown	mammals	12.000	11182	3968	<a href="#">14</a>	1	67892	0.0328	<a href="#">43</a>	0.0313	-0.7882	6.656e-02

## Genes Containing Conserved RYTAAWNNNTGAY Binding Sites:

Gene ID	Ensembl ID	Chr	Strand	TSS	Promoter Start	Promoter End	TFBS Sequence	TFBS Start	TFBS Rel. Start	TFBS End	TFBS Rel. End	TFBS Orientation	TFBS Score
1356	<a href="#">ENSG00000047457</a>	3	-1	150422269	150420270	150424269	ACTAAATTGTGTC	150422362	-93	150422375	-106	-1	8.802
383	<a href="#">ENSG00000118520</a>	6	1	131936059	131934059	131938058	GTACAAGTTTGAC	131937518	1460	131937531	1473	1	6.292
3273	<a href="#">ENSG00000113905</a>	3	1	187866487	187864487	187868486	ACTAATCATTAC	187866344	-143	187866357	-130	1	8.969
462	<a href="#">ENSG00000117601</a>	1	-1	172146654	172144655	172148654	CTTAATATCTGTC	172144991	1664	172145004	1651	1	6.144
				172153139	172151140	172155139	GTCAAAGGCTGAT	172153165	-26	172153178	-39	1	11.059
1571	<a href="#">ENSG00000130649</a>	10	1	135190857	135188857	135192856	TTCAAAGGCTGAT	135190716	-141	135190729	-128	-1	6.038
5105	<a href="#">ENSG00000124253</a>	20	1	55569543	55567543	55571542	ACTAAACCTTGAC	55569306	-237	55569319	-224	-1	13.603
				55569543	55567543	55571542	GTTAATGAATGCT	55569374	-169	55569387	-156	-1	8.174
				55569543	55567543	55571542	GATAATCATTGAA	55569396	-147	55569409	-134	-1	6.601
325	<a href="#">ENSG00000132703</a>	1	1	157824239	157822239	157825284	ATTAAATACAGAC	157822921	-1318	157822934	-1305	-1	10.324
2168	<a href="#">ENSG00000163586</a>	2	-1	88208693	88206694	88210693	GTTAATGTTTGAA	88208792	-99	88208805	-112	-1	12.880
				88208693	88206694	88210693	CTTTATCATTGAC	88208819	-126	88208832	-139	-1	6.066
				88208693	88206694	88210693	ATTAATGTTTGCT	88208867	-174	88208880	-187	-1	11.146
2244	<a href="#">ENSG00000171564</a>	4	1	155703596	155701596	155705595	GTTAATATTTAAT	155703524	-72	155703537	-59	-1	11.267
				155703596	155701596	155705595	GCTAATGTAAGAT	155703971	376	155703984	389	1	7.064

# I have a list of co-expressed mRNAs (Transcriptome)....

## I want to find the shared cis-elements – Known and Novel

### □ Known transcription factor binding sites (TFBS)

#### ❖ Conserved

- oPOSSUM
- DiRE

#### ❖ Non-conserved

- Pscan
- **MatInspector** (\*Licensed)

### □ Unknown TFBS or Novel motifs

#### ❖ Conserved

- oPOSSUM
- **Weeder-H**

#### ❖ Non-conserved

- **MEME**
- **Weeder**

1. Each of these applications support different forms of input. Very few support probeset IDs.
2. **Red Font:** Input sequence required; Do not support gene symbols, gene IDs, or accession numbers. The advantage is you can use them for scanning sequences from any species.
3. \*Licensed software: We have access to the licensed version.

How to fetch promoter/upstream sequence – single/multiple?

# Genome Browser (<http://genome.ucsc.edu>)

## UCSC Genome Bioinformatics

[Genomes](#) - [Blat](#) - [Tables](#) - [Gene Sorter](#) - [PCR](#) - [VisiGene](#) - [Proteome](#) - [Session](#) - [FAQ](#) - [Help](#)

[Genome Browser](#)

[ENCODE](#)

[Blat](#)

[Table Browser](#)

[Gene Sorter](#)

[In Silico PCR](#)

[Genome Graphs](#)

[Galaxy](#)

[VisiGene](#)

[Proteome Browser](#)

[Utilities](#)

[Downloads](#)

[Release Log](#)

[Custom Tracks](#)

[Archaeal Genomes](#)

[Mirrors](#)

[Archives](#)

### About the UCSC Genome Bioinformatics Site

Welcome to the UCSC Genome Browser website. This site contains the reference sequence and working draft assemblies for a large collection of genomes. It also provides a portal to the ENCODE project.

We encourage you to explore these sequences with our tools. The [Genome Browser](#) zooms and scrolls over chromosomes, showing the work of annotators worldwide. The [Gene Sorter](#) shows expression, homology and other information on groups of genes that can be related in many ways. [Blat](#) quickly maps your sequence to the genome. The [Table Browser](#) provides convenient access to the underlying database. [VisiGene](#) lets you browse through a large collection of *in situ* mouse and frog images to examine expression patterns. [Genome Graphs](#) allows you to upload and display genome-wide data sets.

The UCSC Genome Browser is developed and maintained by the Genome Bioinformatics Group, a cross-departmental team within the Center for Biomolecular Science and Engineering ([CBSE](#)) at the University of California Santa Cruz ([UCSC](#)). If you have feedback or questions concerning the tools or data on this website, feel free to contact us on our [public mailing list](#).

### News

[News Archives](#) ►

To receive announcements of new genome assembly releases, new software features, updates and training seminars by email, subscribe to the [genome-announce](#) mailing list.

#### 9 September 2009 - Changes to the bigBed/bigWig data formats

If you have been taking advantage of the new bigBed format (for very large data sets), you'll be happy to hear that we have considerably slimmed down the memory footprint of the program that converts BED files into bigBed files: bedToBigBed. Because it now uses a multi-pass approach, it now takes only 1/4 the amount of RAM as the size of the uncompressed BED input file (instead of the 5x RAM it needed previously!). Read more [here](#). Pick up the new bedToBigBed executable [here](#).

In conjunction with this change, there is also a change to the way you must specify your bigBed or bigWig Custom Track. When you specify the location of your local bigBed/bigWig file (on your web-accessible http, https, or ftp server), use this designation: bigDataUrl (instead of the old designation: dataUrl).

```
e.g. track type=bigBed name="My Big Bed" description="Some Data from My Lab" bigDataUrl=http://myorg.edu/mylab/myBigBed.bb
```

Additionally, we would like to announce a companion program to the previously-announced wigToBigWig program: bedGraphToBigWig. This program converts bedGraph files into bigWig files. The bedGraph format allows display of sparse or varying-size data. Read more [here](#). You can download the new bedGraphToBigWig utility [here](#).

The main advantage of the bigBed and bigWig formats is that only the portions of the files needed to display a particular region are transferred to UCSC, so for large data sets, displaying bigBed/bigWig data is considerably faster than regular BED/wig data. The bigBed/bigWig file remains on your web accessible server (http, https, or ftp), not on the UCSC server. Consequently, creating your Custom Track is very fast. Only the portion that is needed for the chromosomal position you are currently viewing is locally cached at UCSC as a "sparse file".

# Genome Browser (<http://genome.ucsc.edu>)

**Human (*Homo sapiens*) Genome Browser Gateway**

The UCSC Genome Browser was created by the Genome Bioinformatics Group of UC Santa Cruz. Software copyright (c) The Regents of the University of California. All rights reserved.

1 clade: Vertebrate (dropdown menu open)

2 genome: Human (dropdown menu open)

3 assembly: Mar. 2006 (dropdown menu open)

4 position or search term: put symbol, keyword, ID here (input field)

5 image width: 620 (input field)

6 configure tracks and display (button)

**Configure Image**

image width: 620 text size: small (dropdown menu open)

Display chromosome ideograms in main graphic.

Show light blue vertical guide lines between tracks.

Display labels to the left of tracks.

Display track description above each track.

Submit

## Genome Browser Gateway choices:

1. Select Clade
2. Select genome/species: You can search only one species at a time
3. Assembly: the official backbone DNA sequence
4. Position: location in the genome to examine or search term (gene symbol, accession number, etc.)
5. Image width: how many pixels in display window; 5000 max
6. Configure: make fonts bigger + other options

# Genome Browser (<http://genome.ucsc.edu>)

UCSC Genome Bioinformatics

**Genomes** - Blat - Tables - Gene Sorter - PCR - VisiGene - Proteome - Session - FAQ - Help

Genome Browser  
ENCODE  
Blat  
Table Browser

### About the UCSC Genome Bioinformatics Site

Welcome to the UCSC Genome Browser website. This site contains the reference sequence and working data for the human genome. We encourage you to explore these sequences with our tools. The [Genome Browser](#) zooms and scrolls over the genome in various ways. [Blat](#) quickly maps your sequence to the genome. The [Table Browser](#) provides convenient access to upload and display genome-wide data sets.

clade genome assembly position or search term image width

Vertebrate Human Mar. 2006 chrX:151,073,054-151,383,976 620 submit

[Click here to reset](#) the browser user interface settings to their defaults.

add custom tracks configure tracks and display clear position

clade genome assembly position or search term image width

Vertebrate Human Mar. 2006 PDX1 620 submit

[Click here to reset](#) the browser user interface settings to their defaults.

add custom tracks configure tracks and display clear position

# Genome Browser (http://genome.ucsc.edu)

UCSC Genome Browser on Human Mar. 2006 Assembly (hg18)

move <<< << < > >> >>> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x

position/search chr9:103,218,857-103,241,688 jump clear size 22,832 bp. configure

chr9 (q31.1) 24.1 9p23 p21.3 p21.1 12 11.2 9q12 q21.13 q32 q33.1 33.2

Scale chr9: 10 kb | 103225000 | 103230000 | 103235000 | 103240000 |

RepeatMasker RefSeq Genes Repeating Elements by RepeatMasker

move start < 2.0 > Click on a feature for details. Click or drag in the base position track to zoom in. Click gray/blue bars on left for track options and descriptions. move end < 2.0 >

default tracks hide all add custom tracks configure reverse refresh

collapse all Use drop-down controls below and press refresh to alter tracks displayed. expand all Tracks with lots of items will automatically be displayed in more compact modes.

+	Mapping and Sequencing Tracks	refresh
+	Phenotype and Disease Associations	refresh
+	Genes and Gene Prediction Tracks	refresh
+	mRNA and EST Tracks	refresh
+	Expression	refresh
+	Regulation	refresh
+	Comparative Genomics	refresh
+	Variation and Repeats	refresh
+	Pilot ENCODE Regions and Genes	refresh
+	Pilot ENCODE Transcription	refresh
+	Pilot ENCODE Chromatin Immunoprecipitation	refresh
+	Pilot ENCODE Chromatin Structure	refresh
+	Pilot ENCODE Comparative Genomics and Variation	refresh

refresh

**Explore the tracks**

## UCSC Genes

[PDX1 \(uc001urt.1\)](#) at chr13:27392157-27397394 -  
[PDX1 \(uc001mvt.1\)](#) at chr11:34894741-34974094 -  
[SPOP \(uc002ipg.1\)](#) at chr17:45031245-45110524 -  
[SPOP \(uc002ipf.1\)](#) at chr17:45031245-45110524 -  
[SPOP \(uc002ipe.1\)](#) at chr17:45031245-45110524 -  
[SPOP \(uc002ipd.1\)](#) at chr17:45031245-45110524 -  
[SPOP \(uc002ipc.1\)](#) at chr17:45031245-45110524 -  
[SPOP \(uc002ipb.1\)](#) at chr17:45031245-45110524 -

## RefSeq Genes

[PDX1 at chr13:27392157-27397394](#) - (NM\_000209) pa

## Non-Human RefSeq Genes

[Pdx1 at chr13:27392267-27397054](#) - (NM\_022852) pa  
[PDX1 at chr13:27392276-27396838](#) - (NM\_001081478)  
[pdx1 at chr13:27392612-27396613](#) - (NM\_131443) pa  
[Pdx1 at chr13:27392153-27398338](#) - (NM\_008814) pa

## Alias of STS Marker

[PDX1 at chr7:34489430-34689640](#) - (AFM067XA9)

## Non-Human Aligned mRNA Search Results

[U73854](#) - Mesocricetus auratus, homeodomain prote  
[BC103572](#) - Mus musculus, pancreatic and duodenal  
[BC103581](#) - Mus musculus, pancreatic and duodenal homeobox 1  
[BC103582](#) - Mus musculus, pancreatic and duodenal homeobox 1  
[BC105642](#) - Mus musculus, pancreatic and duodenal homeobox 1  
[BC078192](#) >hg18\_dna range=chr13:27391157-27397394 5'pad=1000 box 1

```
GCCAAGCACAGATGTTATCATGGAAAATGCAGCGTTTTTATTCTTTTTT  
TAAATATGTAACCTCTTCCCTCCACTTCCCCTCTCCTGCTTGCCTTATTC  
AATTGCAAGCAGAAGAGAGTGTGTTCTCTGCGGCAAACTCCGCCAGG  
GTCCCGGCCCGTAGAGAGTCGTCAAGGGTCTGGAACCCCGTGCCAACAC  
CTGCCCTGCTTCGCAGCCCCAAGAGGAAGGCCGCGTCTTCCCCTCGC  
TGTATTGGGAAGCTACGTTCCGGGCTGGCCAAATGGGCCCAATTTTCCA  
AAACCCAAATTTGTAATACCCTTCaatttttttaaaaaaagaatttaaaa  
aaGTCTCTGTGAATGCTTCAGAAGTTACCGTTTACACCCCGAAGTACTT  
GCAGCACATCCACAAGTAAAAACACACAACGAATGCCAGAGTTTCGTGTG  
TTTTTTAACCGACATCTTTGTGGCTGTGAACAAACTTCATAAATAAATA  
GAATCAAATGCTTCTGACCTAGAGAGCTGGGTCTGCAAACTTTTTTTTA  
TCGTATTCCGCAACAGTTAAATAAAAAATTA AAAACTCAACATGTCTCCT  
TGTAACACTACATCAATTAACAAACACACTATGTCCATTATCAAATATAAT  
AGAAAAAATATAGGAAAAATAGAAAAATAGAAAAATATAGGAAAAATAGAAAC  
TTTTAAGCCACGGTGAAAAATGTTTCTATAAATGAGTGGTTCTAATGTTT  
CGTGAGCGCCCATTTTGGGGAGCACCGCCAGCTGCCCGTTCAGGAGTGTG  
CAGCAAACTCAGCTGAGAGAGAAAAATGGAACAAAAGCAGGTGCTCGCGG  
GTACCTGGGCCTAGCCTCTTAGTGCGGCCAGCCAGGCCAATCACGGCCCC  
CGGCTGAACCACGTGGGGCCCCGCGGAGCCTATGGTGCGGCGGCCGCCCC  
GCCGCTCCGCGCTGGCTGTGGGTTCCCTCTGAGATCAGTGCGGAGCTGTC
```

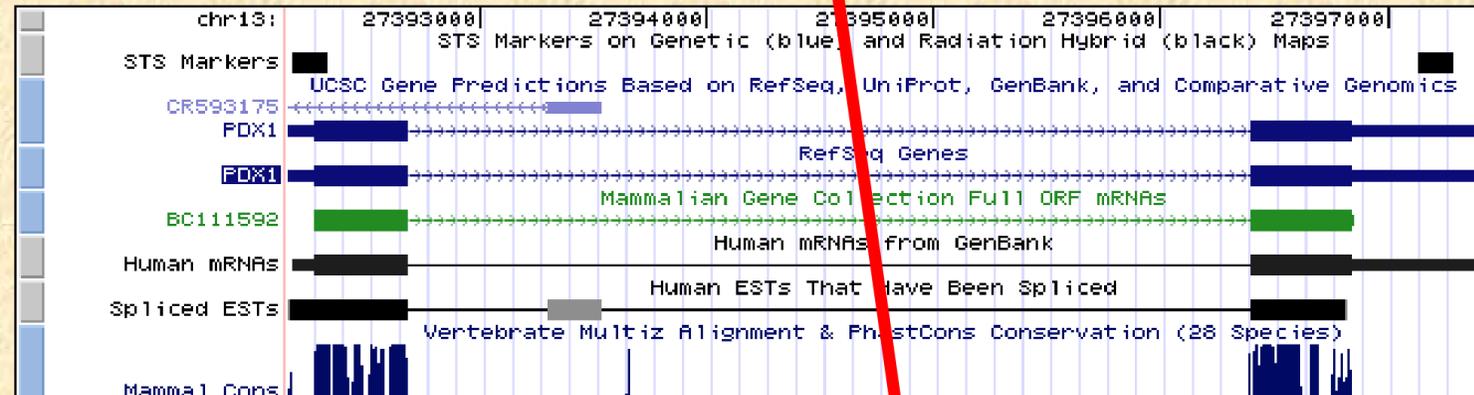
Home Genomes Blat Tables Gene Sorter PCR **DNA** Convert Ensembl NCBI PDF/PS Session Help

## UCSC Genome Browser on Human Mar. 2006 Assembly

move <<< << < > >> >>> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x

position/search    size 5,238 bp.

chr13 (q12.2) 13 12 q31.1 q34



Home Genomes Genome Browser Blat Tables Gene Sorter PCR Session FAQ Help

### Get DNA in Window

### Get DNA for

Position

Note: if you would prefer to get DNA for features of a particular track or table, try the [Table Browser](#) using the output format sequence.

### Sequence Retrieval Region Options:

Add  extra bases upstream (5') and  extra downstream (3')

Note: if a feature is close to the beginning or end of a chromosome and upstream/downstream bases are added, they may be truncated in order to avoid extending past the edge of the chromosome.

### Sequence Formatting Options:

- All upper case.
- All lower case.
- Mask repeats:  to lower case  to N
- Reverse complement (get '-' strand sequence)

Note: The "Mask repeats" option applies only to "get DNA", not to "extended case/color options".

# Genome Browser (<http://genome.ucsc.edu>)

Genome Browser

ENCODE

Blat

Table Browser

Gene Sorter

In Silico PCR

Genome Graphs

Galaxy

VisiGene

Proteome Browser

Utilities

Downloads

Release Log

Custom Tracks

Archaeal Genomes

Mirrors

Archives

### About the UCSC Genome Bioinformatics Site

Welcome to the UCSC Genome Browser website. This site contains the reference sequence and working draft assemblies for a large collection of genomes. It also provides a portal to the ENCODE project.

We encourage you to explore these sequences with our tools. The [Genome Browser](#) zooms and scrolls over chromosomes, showing the work of annotators worldwide. The [Gene Sorter](#) shows expression, homology and other information on groups of genes that can be related in many ways. [Blat](#) quickly maps your sequence to the genome. The [Table Browser](#) provides convenient access to the underlying database. [VisiGene](#) lets you browse through a large collection of *in situ* mouse and frog images to examine expression patterns. [Genome Graphs](#) allows you to upload and display genome-wide data sets.

The UCSC Genome Browser is developed and maintained by the Genome Bioinformatics Group, a cross-departmental team within the Center for Biomolecular Science and Engineering ([CBSE](#)) at the University of California Santa Cruz ([UCSC](#)). If you have feedback or questions concerning the tools or data on this website, feel free to contact us on our [public mailing list](#).

### News

[News Archives](#) ►

To receive announcements of new genome assembly releases, new software features, updates and training seminars by email, subscribe to the [genome-announce](#) mailing list.

#### 9 September 2009 - Changes to the bigBed/bigWig data formats

If you have been taking advantage of the new bigBed format (for very large data sets), you'll be happy to hear that we have considerably slimmed down the memory footprint of the program that converts BED files into bigBed files: bedToBigBed. Because it now uses a multi-pass approach, it now takes only 1/4 the amount of RAM as the size of the uncompressed BED input file (instead of the 5x RAM it needed previously!). Read more [here](#). Pick up the new bedToBigBed executable [here](#).

In conjunction with this change, there is also a change to the way you must specify your bigBed or bigWig Custom Track. When you specify the location of your local bigBed/bigWig file (on your web-accessible http, https, or ftp server), use this designation: bigDataUrl (instead of the old designation: dataUrl).

```
e.g. track type=bigBed name="My Big Bed" description="Some Data from My Lab" bigDataUrl=http://myorg.edu/mylab/myBigBed.bb
```

Additionally, we would like to announce a companion program to the previously-announced wigToBigWig program: bedGraphToBigWig. This program converts bedGraph files into bigWig files. The bedGraph format allows display of sparse or varying-size data. Read more [here](#). You can download the new bedGraphToBigWig utility [here](#).

The main advantage of the bigBed and bigWig formats is that only the portions of the files needed to display a particular region are transferred to UCSC, so for large data sets, displaying bigBed/bigWig data is considerably faster than regular BED/wig data. The bigBed/bigWig file remains on your web accessible server (http, https, or ftp), not on the UCSC server. Consequently, creating your Custom Track is very fast. Only the portion that is needed for the chromosomal position you are currently viewing is locally cached at UCSC as a "sparse file".

# Genome Browser (<http://genome.ucsc.edu>)

1

2

3

6

5

4

### Table Browser

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to see [Using the Table Browser](#) for a description of the controls in this form, the [User's Guide](#) for general information and sample presentation of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL](#) restrictions associated with these data.

clade: Mammal genome: Human assembly: Mar. 2006

**group:** Genes and Gene Prediction Tracks **track:** RefSeq Genes

**table:** refFlat

**region:**  genome  ENCODE  position chrX:151073054-151383976

**identifiers (names/accessions):**

**filter:**

**intersection:**

**output format:** all fields from selected table  Send output to [Galaxy](#)

**output file:** (leave blank to keep output in browser)

**file type returned:**  plain text  gzip compressed

To reset all user cart settings (including custom tracks), [click here](#).

### Table Browser

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to see [Using the Table Browser](#) for a description of the controls in this form, the [User's Guide](#) for general information and sample presentation of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL](#) restrictions associated with these data.

clade: Mammal genome: Human assembly: Mar. 2006

**group:** Genes and Gene Prediction Tracks **track:** RefSeq Genes

**table:** refFlat

**region:** position chrX:151073054-151383976

**identifiers (names/accessions):**

**filter:**

**intersection:**

**output format:** table  Send output to [Galaxy](#)

**output file:** (leave blank to keep output in browser)

**file type returned:**  plain text  gzip compressed

To reset all user cart settings (including custom tracks), [click here](#).

### Using the Table Browser

### Paste In Identifiers for refFlat

Please paste in the identifiers you want to include. The items must be values of the "geneName" field in the "refFlat" table schema. Some example values:

LOC133874  
PLXDC1  
SNORD115-17

AMBP  
SERPINC1  
APOH  
APOA1  
CYP2C8  
CYP2E1  
ALDOB  
SERPINC1

### Sequence Retrieval Region Options:

Promoter/Upstream by 500 bases

5' UTR Exons

CDS Exons

3' UTR Exons

Introns

Downstream by 1000 bases

One FASTA record per gene.

One FASTA record per region (exon, intron, etc.) with 0 extra bases upstream (5') and 0 extra downstream (3')

Split UTR and CDS parts of an exon into separate FASTA records

Note: if a feature is close to the beginning or end of a chromosome and upstream/downstream bases are added, they may be truncated.

### Sequence Formatting Options:

Exons in upper case, everything else in lower case.

CDS in upper case, UTR in lower case.

All upper case.

All lower case.

Mask repeats:  to lower case  to N

```
>hg18_refFlat_SERPINC1 range=chr1:172153140-172153639 5'pad=0 3'pad=0
gaaacggccttggcagcagccgagccctgctctctccctgtccc
accactcagggtgctgggaatgggtctctgtggccacaggtgta
accattgtgtttccctgtctgtgccaggacacctggcactcagatgc
ctgaaggtagcagcttgcctcttgccttctctaatagatattctc
tctctctccctctctccataaagaaaactatgagagaggggtggatg
aaccaagttgtttccctgggtagtttcttaaccaagtttgagggtatga
acatactctctcttccctttctataaagctgaggagaagagtgaggag
tgtgggaagagaggtggctcaggctttccctggccctgatgaaactta
aaactctctactaataaacacactgggctctacactttgcttaacc
tgggaactggctcagcctttgacctcagttccctcctgaccagctc
>hg18_refFlat_ALDOB range=chr9:103237884-103238383 5'pad=0 3'pad=0
ttcttgaatgcctactgcgtacagacactatacaactcacttaatgctcc
tgtgaaatgagatcccactttacagagaagaaaataacggctcagacaa
gtaaaaaacaaaacaaaacccctcaccagttcctatgtgtttgtgaa
gttacagctacactaacattctctagatcatcttaaatggacctagag
ctccatcattgactagctgtatgactaaaggaacctctctctctct
gccttagtttccctcactgtgaaatggagggtctggattagaaactccac
ggctaaagacattccttgcagctttaaactctatgagcataaggagta
tatggcagatgatttaaggactggttgttatgagcaatcagaggtgtt
gaataaacacctccctactaggtcaaggtagaagggaggggcaaatat
ggaaaaaaaacacatgatgagaagctataaaaaatgtgtgctacaaa
>hg18_refFlat_AMBP range=chr9:115880574-115881073 5'pad=0 3'pad=0
tctgggcaagtcacttccctcttggagactcaactcctctgctgtaaa
ggggcaaatagcactacttgcagctgtgtgtgggatttactggac
acctgtggcagctgtaagtgctgtgacacagccagcactagcaggtg
ctcagtgaaagagctcgtaacagctcatctgttaagtctctctat
ttgcgaaggttaggcaagcctggggctggagaaggaactctggagaca
ggatacctgggtctggtcttggctctgcccctagactccctgaggtct
ctggcaagtcacttccctccttgggctcagtttttctgtgagaaat
gggcccactcaggcttgacctgctggtctcacagagctgtgaggtccagat
ggcatgactcagcagaagggcaggtcgtccccatcctcgcagcctc
tgtgggggacttttggggggggagcccaatcctggtgctccagggcct
>hg18_refFlat_CYP2C8 range=chr10:96819245-96819744 5'pad=0 3'pad=0
tatacattttagctaaatattgatattgtaattcaacatgt
atgagttatattcactatttcatgtttaggcagctgtatttaagtgaac
tatacctaaatattgaaaggttttggatcaagggtcaagctcctatt
ttttgatagcattacaatgtacatttttatacacaataatagaata
cactgatttccctcaaggtcataaatccccactggctcattaatctgaga
atattgaaatttgagatatttcaacatagaatcatttactcaggtttt
ctccatcatcacagcacttggacaaccagggtctttaaataaaaaa
cctgggctccaatccaatacaataaacacagaatctcctagattggcact
ggaaaagaggttaggacaaaagaacattttatctatccatgggcaa
agtcactcagaaaaaagataaaatggactcaggtgatgtttacttt
```

### Table Browser (Input Identifiers)

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to see [Using the Table Browser](#) for a description of the controls in this form, the [User's Guide](#) for general information and sample presentation of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL](#) restrictions associated with these data.

clade: Mammal genome: Human assembly: Mar. 2006

**group:** Genes and Gene Prediction Tracks **track:** RefSeq Genes

**table:** refFlat

**region:**  genome  ENCODE  position chrX:151073054-151383976

**identifiers (names/accessions):**

**filter:**

**intersection:**

**output format:** sequence

**output file:** all fields from selected table  
selected fields from primary and related tables

**file type returned:** sequence

GTF - gene transfer format  
BED - browser extensible data  
custom track  
hyperlinks to Genome Browser

To reset all user cart settings (including custom tracks), [click here](#).



Home Genomes Genome Browser Blat Tables Gene Sorter PCR Session FAQ Help

### Table Browser

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to retrieve DNA sequence covered by a track. For help in using this application see [Using the Table Browser](#) for a description of the controls in this form, the [User's Guide](#) for general information and sample queries, and the OpenHelix Table Browser [tutorial](#) for a narrated presentation of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL server](#). Refer to the [Credits](#) page for the list of contributors and usage restrictions associated with these data.

clade: Vertebrate genome: Human assembly: Mar. 2006

group: All Tables database: hg18

table: refFlat describe table schema

region:  genome  position chrX:151073054-151383976 lookup define regions

identifiers (names/accessions): paste list upload list

filter: create

intersection with ct\_UserTrack: edit clear

output format: all fields from selected table  Send

output file: (leave blank to keep output in browser)

file type returned:  plain text  gzip compressed

Note: Intersection doesn't work with all fields or selected fields output.

get output summary/statistics

To reset all user cart settings (including custom tracks), [click here](#).

1

1. Select "refFlat" under "table"
2. Ensure that "region" is "genome"
3. Click on "paste list"

Home Genomes Genome Browser Blat Tables Gene Sorter PCR Session FAQ Help

### Paste In Identifiers for refFlat

Please paste in the identifiers you want to include. The items must be values of the **geneName** field of the currently selected table, **refFlat**. (The "describe table schema" button shows more information about the table fields.) Some example values:

AADACL3  
GPX4  
ACF

1. Paste the gene symbols

2. Remember it is case-sensitive:

- Human: all upper case (e.g. XRCC1)
- Mouse: lower case (first letter upper case. E.g. Xrcc1)

ABCC3  
ABHD2  
ACY1  
ADH1C  
APCS  
APOC4  
AQP9  
ASL  
AZGP1

submit clear cancel

2

Home Genomes Genome Browser Blat Tables Gene Sorter PCR Session FAQ Help

### Output refFlat as Custom Track

Custom track header:

name= UpSt\_1kb

description= table browser query on refFlat

visibility= pack

url=

Describe your track

4

Create one BED record per:

Whole Gene

Upstream by 1000 bases

Exons plus 0 bases at each end

Introns plus 0 bases at each end

5' UTR Exons

Coding Exons

3' UTR Exons

Downstream by 200 bases

Note: if a feature is close to the beginning or end of a chromosome and upstream/downstream bases are added, they may be truncated in order to avoid extending past the edge of the chromosome.

get custom track in table browser get custom track in file

get custom track in genome browser cancel

Enter number of bp you want to analyze/download

Home Genomes Genome Browser Blat Tables Gene Sorter PCR Session FAQ Help

### Table Browser

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to retrieve DNA sequence covered by a track. For help in using this application see [Using the Table Browser](#) for a description of the controls in this form, the [User's Guide](#) for general information and sample queries, and the OpenHelix Table Browser [tutorial](#) for a narrated presentation of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL server](#). Refer to the [Credits](#) page for the list of contributors and usage restrictions associated with these data.

clade: Vertebrate genome: Human assembly: Mar. 2006

group: All Tables database: hg18

table: refFlat describe table schema

region:  genome  position chrX:151073054-151383976 lookup define regions

identifiers (names/accessions): paste list upload list clear list

filter: create

intersection: create

output format: custom track  Send output to Galaxy

output file: (leave blank to keep output in browser)

file type returned:  plain text  gzip compressed

get output summary/statistics

To reset all user cart settings (including custom tracks), [click here](#).

3

- Select the output format as "custom track"

Home Genomes Genome Browser Blat Tables Gene Sorter PCR Session FAQ Help

### Table Browser

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to retrieve DNA sequence covered by a track. For help in using this application see [Using the Table Browser](#) for a description of the controls in this form, the [User's Guide](#) for general information and sample queries, and the OpenHelix Table Browser [tutorial](#) for a narrated presentation of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL server](#). Refer to the [Credits](#) page for the list of contributors and usage restrictions associated with these data.

clade: Vertebrate genome: Human assembly: Mar. 2006

**group:** Variation and Repeats **track:** SNPs (126)

**table:** snp126 describe table schema

region:  genome  position chrX:151073054-151383976 lookup define regions

identifiers (names/accessions): paste list upload list

filter: create

**intersection:** create

correlation: create

output format: custom track

output file: (leave blank to keep out)

file type returned:  plain text  gzip compressed

get output summary/statistics

To reset all user cart settings (including custom tracks), [click here](#).

5

1. Select "Variation and Repeats" under "Group"
2. Click on "create" under "intersection"

Home Genomes Genome Browser Blat Tables Gene Sorter PCR Session FAQ Help

### Intersect with SNPs (126)

Select a group, track and table to intersect with:

**group:** Genes and Gene Prediction Tracks **track:** UCSC Genes

**table:** UCSC Genes (knownGene)

6

These combinations will maintain the gene/alignment structure (if any) of SNPs (126):

- All SNPs (126) records that have any overlap with UCSC Genes
- All SNPs (126) records that have no overlap with UCSC Genes
- All SNPs (126) records that have at least 80 % overlap with UCSC Genes
- All SNPs (126) records that have at most 80 % overlap with UCSC Genes

These combinations will discard the gene/alignment structure (if any) of SNPs (126) and produce a simple list of position ranges.

- Base-pair-wise intersection (AND) of SNPs (126) and UCSC Genes
- Base-pair-wise union (OR) of SNPs (126) and UCSC Genes

Check the following boxes to complement one or both tables. To complement a table means to include a row in the intersection if it is *not* included in the table.

Complement SNPs (126) before intersection/union

Complement UCSC Genes before intersection/union

submit cancel

Change the "group" to "Custom Tracks" and select the appropriate "track" and "table"

Home Genomes Genome Browser Blat Tables Gene Sorter PCR Session FAQ Help

### Table Browser

Use this program to retrieve the data associated with a track in text format, to calculate intersections between tracks, and to retrieve DNA sequence covered by a track. For help in using this application see [Using the Table Browser](#) for a description of the controls in this form, the [User's Guide](#) for general information and sample queries, and the OpenHelix Table Browser [tutorial](#) for a narrated presentation of the software features and usage. For more complex queries, you may want to use [Galaxy](#) or our [public MySQL server](#). Refer to the [Credits](#) page for the list of contributors and usage restrictions associated with these data.

clade: Vertebrate genome: Human assembly: Mar. 2006

**group:** Variation and Repeats **track:** SNPs (126)

**table:** snp126 describe table schema

region:  genome  position chrX:151073054-151383976 lookup define regions

identifiers (names/accessions): paste list upload list

filter: create

intersection with ct\_UserTrack: edit clear

correlation: create

**output format:** BED - browser extensible data  Send output to [Galaxy](#)

output file: (leave blank to keep output in browser)

file type returned:  plain text  gzip compressed

Note: Intersection doesn't work with all fields or selected fields output.

get output summary/statistics

To reset all user cart settings (including custom tracks), [click here](#).

8

Try GTF output too

Home Genomes Genome Browser Blat Tables Gene Sorter PCR Session FAQ Help

### Intersect with SNPs (126)

Select a group, track and table to intersect with:

**group:** Custom Tracks **track:** User Track

**table:** User Track (ct\_UserTrack)

7

These combinations will maintain the gene/alignment structure (if any) of SNPs (126):

- All SNPs (126) records that have any overlap with User Track
- All SNPs (126) records that have no overlap with User Track
- All SNPs (126) records that have at least 80 % overlap with User Track
- All SNPs (126) records that have at most 80 % overlap with User Track

These combinations will discard the gene/alignment structure (if any) of SNPs (126) and produce a simple list of position ranges.

- Base-pair-wise intersection (AND) of SNPs (126) and User Track
- Base-pair-wise union (OR) of SNPs (126) and User Track

Check the following boxes to complement one or both tables. To complement a table means to include a row in the intersection if it is *not* included in the table.

Complement SNPs (126) before intersection/union

Complement User Track before intersection/union

submit cancel

Home Genomes Genome Browser Blat Tables Gene Sorter PCR Session FAQ Help

Output snp126 as BED

Include custom track header:

name=

description=

visibility=

url=

9

Create one BED record per:

Whole Gene

Upstream by  bases

Downstream by  bases

Note: if a feature is close to the beginning or end of a chromosome and upstream/downstream bases are added, they may be truncated in order to avoid extending past the edge of the chromosome.

chr1	157823462	157824462
chr1	157823893	157824893
chr1	157822934	157823934
chr1	157822978	157823978
chr1	157823993	157824993
chr1	157823138	157824138
chr1	157823144	157824144
chr4	100493180	100494180
chr4	100492308	100493308
chr4	100493487	100494487
chr4	100492593	100493593
chr4	100493772	100494772
chr4	100492918	100493918
chr7	65176569	65177569
chr7	99411067	99412067
chr15	56215720	56216720
chr15	56216839	56217839
chr15	56217082	56218082
chr15	56216231	56217231
chr15	56216582	56217582

rs3753869_up_1000_chr1_157823463_r	0	-
rs3753868_up_1000_chr1_157823894_r	0	-
rs2592882_up_1000_chr1_157822935_f	0	+
rs2808655_up_1000_chr1_157822979_f	0	+
rs3753867_up_1000_chr1_157823994_r	0	-
rs28383571_up_1000_chr1_157823139_f	0	+
rs6689429_up_1000_chr1_157823145_f	0	+
rs4147541_up_1000_chr4_100493181_r	0	-
rs1789924_up_1000_chr4_100492309_f	0	+
rs34774688_up_1000_chr4_100493488_r	0	-
rs17586163_up_1000_chr4_100492594_f	0	+
rs11499823_up_1000_chr4_100493773_r	0	-
rs1629838_up_1000_chr4_100492919_f	0	+
rs10247708_up_1000_chr7_65176570_f	0	+
rs4727442_up_1000_chr7_99411068_f	0	+
rs11629801_up_1000_chr15_56215721_f	0	+
rs1554203_up_1000_chr15_56216840_r	0	-
rs3840843_up_1000_chr15_56217083_r	0	-
rs8027250_up_1000_chr15_56216232_f	0	+
rs9920853_up_1000_chr15_56216583_f	0	+
rs3742955_up_1000_chr15_56216602_f	0	+
rs293370_up_1000_chr15_87431481_r	0	-
rs293369_up_1000_chr15_87432137_r	0	-
rs4148404_up_1000_chr17_46065400_f	0	+
rs35467079_up_1000_chr17_46065407_f	0	+
rs9895420_up_1000_chr17_46066037_f	0	+
rs4793665_up_1000_chr17_46066086_f	0	+
rs12721101_up_1000_chr19_50135470_f	0	+
rs35366732_up_1000_chr19_50135470_f	0	+
rs4803773_up_1000_chr19_50135582_f	0	+
rs12721105_up_1000_chr19_50136135_f	0	+

10

Home Genomes Blat Tables Gene Sorter PCR DNA Convert Ensembl NCBI PDF/PS Session Help

## UCSC Genome Browser on Human Mar. 2006 Assembly

move <<< << < > >> >>> zoom in 1.5x 3x 10x base zoom out 1.5x 3x 10x

position/search    size 5,010 bp

chr1 (q23.2)

Genome Browser view that lists all the SNPs lying within the upstream 1 kb (the region we queried) region of one of the genes analyzed.

One drawback with this output is it doesn't tell you which SNPs are in the upstream region of which gene. However, since the positions of SNPs are included, you can compare them with the gene coordinates and figure it out.

**Exercise 4: Download upstream 500 bp sequence for a list of genes (use the same list as before).**

**Exercise 5: Download all SNPs overlapping with these genes.**

**Exercise 6: Download the orthologous promoter sequences (human, mouse, and rat) for the gene SLC7A1.**

**Exercise 7: Are there any putative microRNA regulators for SLC7A1? If yes, download all of them using table browser.**

# I have a list of co-expressed mRNAs (Transcriptome)....

## I want to find the shared cis-elements – Known and Novel

### □ Known transcription factor binding sites (TFBS)

#### ❖ Conserved

- oPOSSUM
- DiRE

#### ❖ Non-conserved

- Pscan
- **MatInspector** (\*Licensed)

### □ Unknown TFBS or Novel motifs

#### ❖ Conserved

- oPOSSUM
- **Weeder-H**

#### ❖ Non-conserved

- **MEME**
- **Weeder**

1. Each of these applications support different forms of input. Very few support probeset IDs.
2. **Red Font:** Input sequence required; Do not support gene symbols, gene IDs, or accession numbers. The advantage is you can use them for scanning sequences from any species.
3. \*Licensed software: We have access to the licensed version.

Use the fetched promoter/upstream sequences for the following analyses

# WeederH (<http://159.149.109.9/pscan>)

## WeederH

Motif discovery in sequences from **homologous** genes  
Version **beta** running.

[Click here to switch to Weeder](#)

Please, avoid submitting a large number of jobs (> 5) simultaneously. For large-scale analyses, you're welcome to download the standalone version.

**NEW** If you are looking for over-represented motifs in promoter sequences, perhaps you can also find our brand new tool, [Pscan](#) useful.

Enter your e.mail address

Input **exactly one** sequence in each box

Reference sequence (FASTA)	<input type="text"/>	from Homo sapiens
Homologous sequence n. 1 (FASTA)	<input type="text"/>	from Homo sapiens
Homologous sequence n. 2 (FASTA)	<input type="text"/>	from Homo sapiens
Homologous sequence n. 3 (FASTA)	<input type="text"/>	from Homo sapiens Mus musculus Rattus norvegicus Canis familiaris Yeast (any) Drosophila (any) Caenorhabditis (any) Anopheles gambiae Arabidopsis thaliana Ciona intestinalis Danio rerio Fugu rubripes Gallus gallus Xenopus tropicalis P. falciparum Magnaporthe grisea

For technical reasons, it's better that you directly paste sequence in the text boxes, rather than uploading a file.

Your sequences are

Name of this job:

1. Supports large number of species.
2. Does not support multiple sequences (multifasta) input. You have to enter each sequence separately.
3. Good for small number of sequences where you expect a potential novel (or not included in the TFBS libraries) conserved motif.

# Weeder (<http://159.149.109.9:8080/weederweb2006>)

The screenshot shows the Weeder web application interface. At the top left, there is a logo with the text 'MOTIF TOOLZ' and a 'Home' link. Below this, there are sections for 'On-line tools' (Weeder & WeederH, RNAProfile) and 'Additional tools' (Motif locator, Motif p-value calculator). The main content area is titled 'Tools for MOTIF Discovery in nucleotide sequences' and features a 'Weeder' input form. The form includes a text area for 'Enter your e.mail address', a section for 'Input at least two sequences' with a 'Browse...' button and an 'Upload' button, and a checkbox for 'Check here if you want to look for motifs in both strands of the input sequences'. The browser window title is 'Weeder Input Form - Windows Internet Explorer' and the address bar shows 'http://159.149.109.16:8080/weederweb2006/input.faces'.

## Weeder

Motif discovery in sequences from **co-regulated** genes  
Version 1.3.1 running.

[Click here to switch to WeederH](#)

Please note: submitting simultaneously a large (> 10) number of jobs has the effect of slowing down the server, for your jobs, as well as the jobs submitted by other users. If you plan to use Weeder extensively, you can download the stand-alone version. Client IPs and e-mail addresses generating high workloads on the server (as defined in the previous sentence) might have their jobs terminated before completion without notice.

Enter your e.mail address

Input at least two sequences

**Input sequences (FASTA)**

```
gattgtaatgactaatctgtgtccatgaggcacagagccaaggaagaga
tgctgctgtagccagaaaggccgctgtgatcatgcacagtacactgga
```

from:

To upload a file, first locate it by using the browse button, then click on Upload.

Check here if you want to look for motifs in **both strands** of the input sequences

Check here if you want motifs to appear in **all the sequences** (default is in **some**)  
Hint: don't try this option even if you're pretty much sure that all your sequences share a motif.

Check here if you think that the motif might appear **more than once** in a single sequence (without, you expect zero or one occurrence per sequence)

And, finally, you'd like:

a quick scan (short motifs, no longer than 8 nts) of your sequences

a normal scan of your sequences

a complete and thorough scan

**Important:** input larger than 20K will be limited to quick analysis. For larger jobs, you can download the source code by following the link in the home page.

**Quick scan:** results will be ready in a few minutes **Normal scan:** results will be ready in one-two hours **Thorough scan:** results will be ready in a few hours **However:** try the normal scan first. If nothing interesting comes out, try the thorough one.

Name of this job:

Click submit once to start the computation. Click reset to clear all the fields.

Do not use Groupwise mail when submitting large number of sequences because the results are sent "in the mail" and not as an attachment. And Groupwise mail truncates messages if they are very long. Use Gmail instead. A link to the results page used to be sent earlier.

## Weeder

Thank you!  
You submitted 33 sequences from Homo sapiens

You asked to process both strands of the input sequences  
You asked for a normal scan

A confirmation e-mail and the final results will be sent to the following e.mail address:  
**anil.jegga@gmail.com**

# Weeder (<http://159.149.109.9:8080/weederweb2006>)

## \*\*\* Your Weeder Web Results \*\*\*

The name of this job was Fetal\_Liver\_33\_27

Input sequences from H. sapiens

You asked to include both strands of the input sequences  
 You asked for a normal scan of your sequences

Confused about this output? Click [here](#)

Searching for motifs of length 6 with 1 mutations.....

- 1) CAATTA 0.81
- 2) TAAACG 0.70
- 3) ATTGAT 0.67
- 4) TATGAT 0.63
- 5) GATTTA 0.61
- 6) ATGGTA 0.60
- 7) TCATTG 0.59
- 8) TGGTAT 0.59
- 9) TGATTA 0.59
- 10) TGATAT 0.58

Searching for motifs of length 8 with 2 mutations.....

- 1) CGTTIAGA 0.93
- 2) ACTAAACG 0.88
- 3) GATAAACT 0.87
- 4) TATGGTAT 0.87
- 5) CTAAACGT 0.87
- 6) AGTATTTC 0.84
- 7) ACATTGAT 0.82
- 8) GTAATACT 0.80
- 9) CTAGCAAT 0.79
- 10) ATAGTTCG 0.78

\*\*\* Interesting motifs (highest-ranking) seem to be :

**GATAAACT**  
**AGTTTATC**

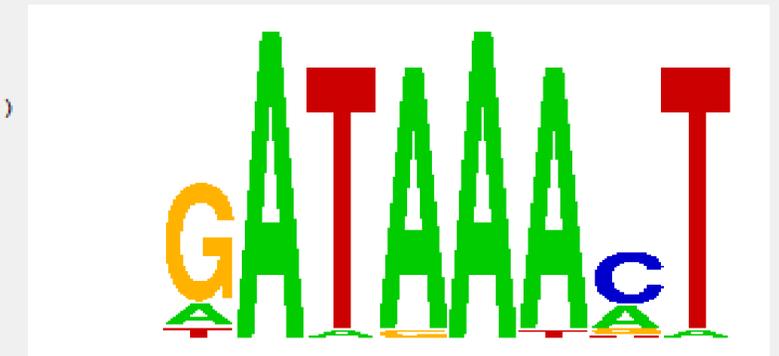
0 redundant motifs found:

Best occurrences (match percentage):

```
Seq St oligo pos match
1 + .GAAAAACT. 205 (92.84)
1 + [AATAAATT] 676 (85.29)
1 + [GATTAACT] 922 (88.60)
1 - .TATAAACT. 786 (92.79)
1 - .AATAAACT. 697 (92.36)
1 - [GATAATAT] 169 (85.17)
2 + [TAAAAACT] 508 (85.63)
2 + [TATAAATT] 944 (85.73)
2 - [GAAAAAGT] 956 (85.28)
2 - [AATAAATT] 776 (85.29)
2 - .GATGAACT. 652 (90.33)
4 + [AATAAAAT] 546 (87.13)
4 + [GAGAAAAT] 786 (85.24)
5 + [AATAAATT] 393 (85.29)
5 - [GAGAAAAT] 260 (85.24)
6 + [TATAAAAT] 733 (87.56)
7 - .GATAAAAT. 430 (94.77)
8 + [AATAAAAT] 307 (87.13)
8 + [AATAAAAT] 791 (87.13)
8 - [AAAAAACT] 808 (85.19)
8 - [AATAAATT] 484 (85.29)
8 - [TATAAAGT] 285 (85.24)
8 - [TATAAAAT] 13 (87.56)
9 + .GATAAACT. 603 (100.00)
9 + [GAGAAAAT] 615 (85.24)
9 - .GATAAAAT. 438 (94.77)
10 + .GATAAACT. 603 (100.00)
10 + [GAGAAAAT] 615 (85.24)
10 - .GATAAAAT. 438 (94.77)
11 + [GATGAAAT] 148 (85.10)
11 + [GAAAAATT] 205 (85.77)
12 + .GATAAATT. 143 (92.93)
12 + .TATAAACT. 271 (92.79)
12 + [AATAAAAT] 286 (87.13)
12 + .GATAAACA. 523 (90.60)
12 + .GATAAAAT. 896 (94.77)
12 - [AATAAATT] 347 (85.29)
13 + .AATAAACT. 549 (92.36)
13 - [GATAAGCT] 832 (88.34)
13 - [GATAAGCT] 577 (88.34)
14 + .AATAAACT. 549 (92.36)
14 - [GATAAGCT] 832 (88.34)
14 - [GATAAGCT] 577 (88.34)
16 + .GATAAAAT. 161 (94.77)
16 + [GACAAACT] 316 (89.17)
16 + [AATAAATT] 814 (85.29)
16 - .AATAAACT. 967 (92.36)
16 - .GATAATCT. 943 (90.40)
18 - [GTTAAACT] 637 (89.17)
18 - .GATAAAAT. 555 (94.77)
```

### Frequency Matrix

	All Occs				Best Occs			
	A	C	G	T	A	C	G	T
1	28	16	167	31	4	0	20	2
2	201	8	17	16	26	0	0	0
3	33	14	19	176	1	0	0	25
4	201	6	21	14	25	0	1	0
5	208	6	9	19	26	0	0	0
6	198	10	13	21	25	0	0	1
7	43	146	25	28	7	16	2	1
8	22	17	5	198	1	0	0	25



# MEME (<http://meme.sdsc.edu>)

MEME takes as input a group of DNA or protein sequences and outputs as many motifs as requested. MEME uses statistical modeling techniques to automatically choose the best width, number of occurrences, and description for each motif.

Your MEME results consist of:

- your MEME results in HTML format
- your MEME results in XML format
- your MEME results in TEXT format
- and the MAST results of searching your input sequences for the motifs found by MEME using MAST.



Version 4.1.1

Use this form to submit DNA or protein sequences to MEME. MEME will analyze your sequences for similarities among them and produce a description (**motif**) for each pattern it discovers.

## Data Submission Form

### Required

Your **e-mail address**:

Re-enter **e-mail address**:

Please enter the **sequences** which you believe share one or more motifs. The sequences may contain no more than **60000 characters** total in any of a large number of **formats**.

Enter the **name of a file** containing the sequences here:

or

the **actual sequences** here (**Sample Protein Input Sequences**):

```
>SERPINC1 range=chr1:172153140-172154139
ggtgacagatgagcctctgggcattctggggcccttgggacagttctcg
gctactgttcacactgcacatggagggtcgaaggatccagggtctgaa
tcagccaagaacttagacacagcttcagtcaggagagatgtcctctc
atgaacgaagagctctcgaatgacctatgactgcaacaactacaag
gagtcctgatcacacagcaggaggcacatgcacctgtgaactgttg
```

How do you think the occurrences of a single motif are **distributed** among the sequences?

- One per sequence  
 Zero or one per sequence  
 Any number of repetitions

**Note:** The maximum number of occurrences of a motif is limited to 300.

MEME will find the optimum **width** of each motif within the limits you specify here:

**Minimum width** ( $\geq 2$ )  
 **Maximum width** ( $\leq 300$ )

**Maximum number of motifs** to find

### Optional

**Description** of your sequences:

MEME will find the optimum **number of sites** for each motif within the limits you specify here:

**Minimum sites** ( $\geq 2$ )  
 **Maximum sites** ( $\leq 300$ )

**Shuffle** sequence letters

Enter the name of a file containing a **background Markov model**:

**DNA-ONLY OPTIONS**  
(Ignored for protein searches)

- Search given **strand** only  
 Look for **palindromes** only

Your job id is: **app1254080196482**

You can view your job results at: [http://meme.nbc.net/meme4\\_1\\_1/cgi-bin/querystatus.cgi?jobid=app1254080196482&service=MEME](http://meme.nbc.net/meme4_1_1/cgi-bin/querystatus.cgi?jobid=app1254080196482&service=MEME)

You can view server activity [here](#).

- Sequence file: **pasted\_sequences**
- Distribution of motif occurrences: **Zero or one per sequence**
- Number of different motifs: **20**
- Minimum motif width: **5**
- Maximum motif width: **20**
- Statistics on your dataset:

type of sequence	<b>dna</b>
number of sequences	<b>20</b>
shortest sequence (residues)	<b>1000</b>
longest sequence (residues)	<b>1000</b>
average sequence length (residues)	<b>1000.0</b>
total dataset size (residues)	<b>20000</b>

You will also receive a confirming message at your email address: **anil.jegga@cchmc.org**





**Exercise 8: Use the downloaded SLC7A1 ortholog promoter sequences to find out common motifs using WeederH**

**Exercise 9: Use the downloaded promoter sequences (from Exercise 4) to find out common motifs using Weeder and MEME**

**Exercise 10: Does any of the motifs found by Meme match known TFBS?**

I have found a miRNA enriched in my gene list or I am interested in a specific gene and I want to identify putative regulatory regions for miRNA/gene

GenomeTrafac: <http://genometrafac.cchmc.org>

## GenomeTraFaC

A comparative genomics-based resource for initial characterization of gene models and the identification of putative cis-regulatory regions of RefSeq Gene Orthologs

- Cis-element clusters within BlastZ Alignments  
Find conserved *cis*-element clusters within BlastZ-identified conserved sequence alignment blocks.
- Cis-elements shared between any gene pair  
Find shared *cis*-elements between user-selected gene segment pairs.
- Conserved Cis-Element Scanner  
Genome-wide ortholog conserved Cis-element module search

**Note:** If you publish results obtained using GenomeTrafac, please cite [Jegga et al., Nucleic Acids Res. 2006 Dec 18; \[Epub ahead of print\]](#)

OR

[Jegga et al., Genome Research 12: 1408-1417, September 2002](#)

# GenomeTrafac: <http://genometrafac.cchmc.org>



## Basic Search

Description

mir-122a

## Search by disease, gene ontology, pathway, gene family, or custom groups

Select

Query

- Disease (*Always use [Disease Selector](#)*)
- Pathway (*Always use [Pathway Selector](#)*)
- Gene ontology (*Always use [Ontology Selector](#)*)
- Mammalian phenotype (*Always use [Phenotype Selector](#)*)
- Select gene family from the list
- Select custom group from the list

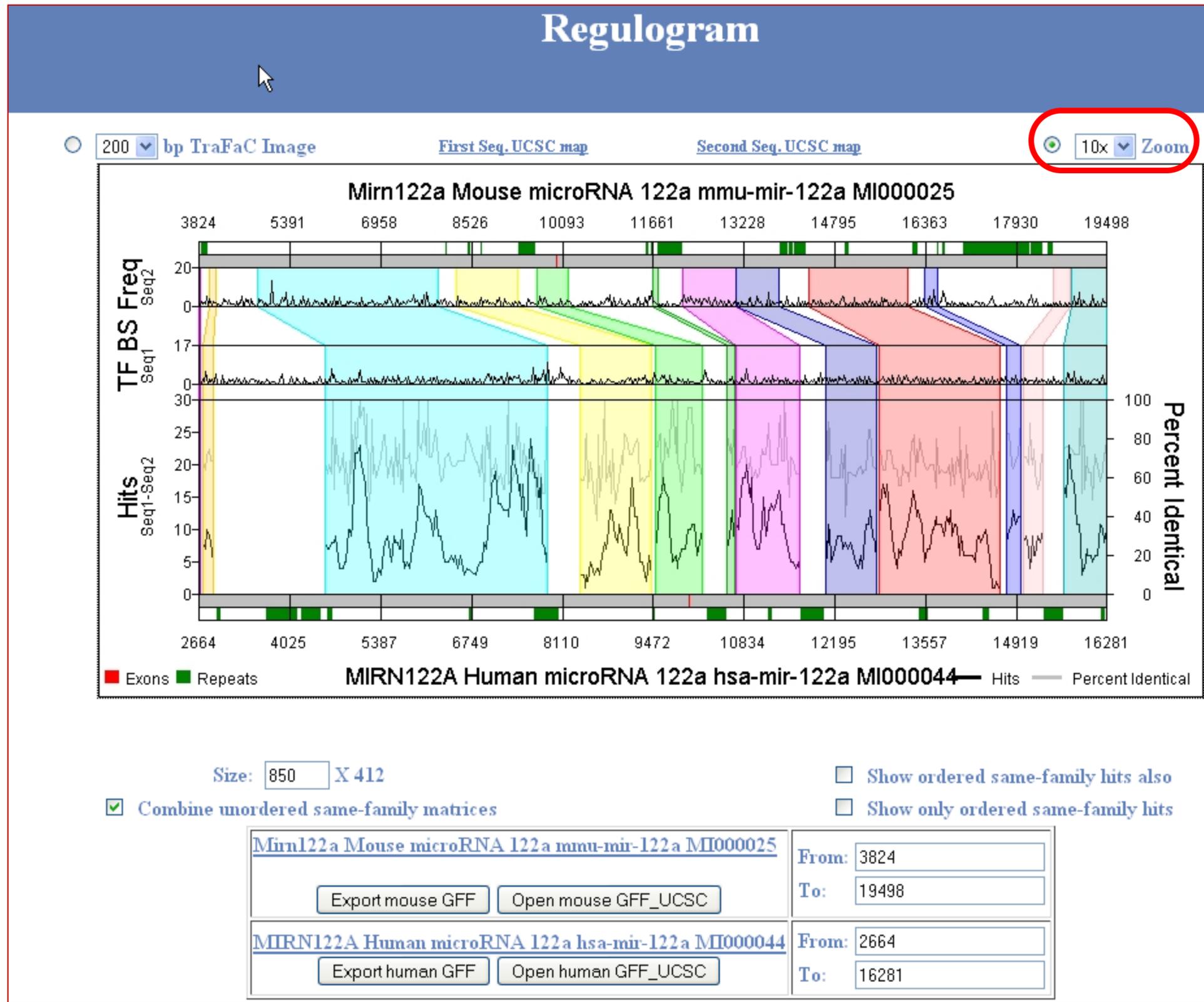
Query took 1.514 s

(2 genes meet the search criteria)

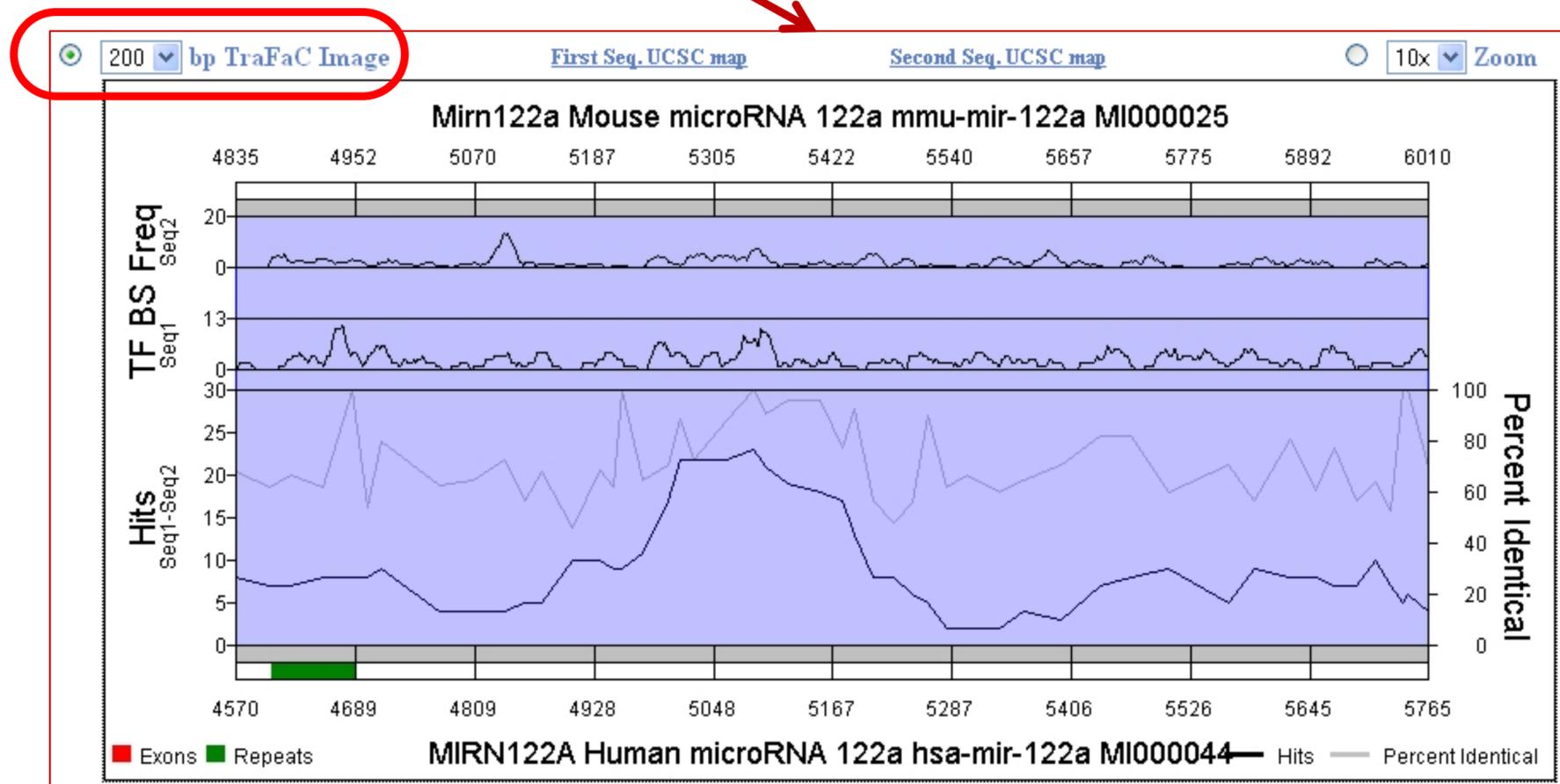
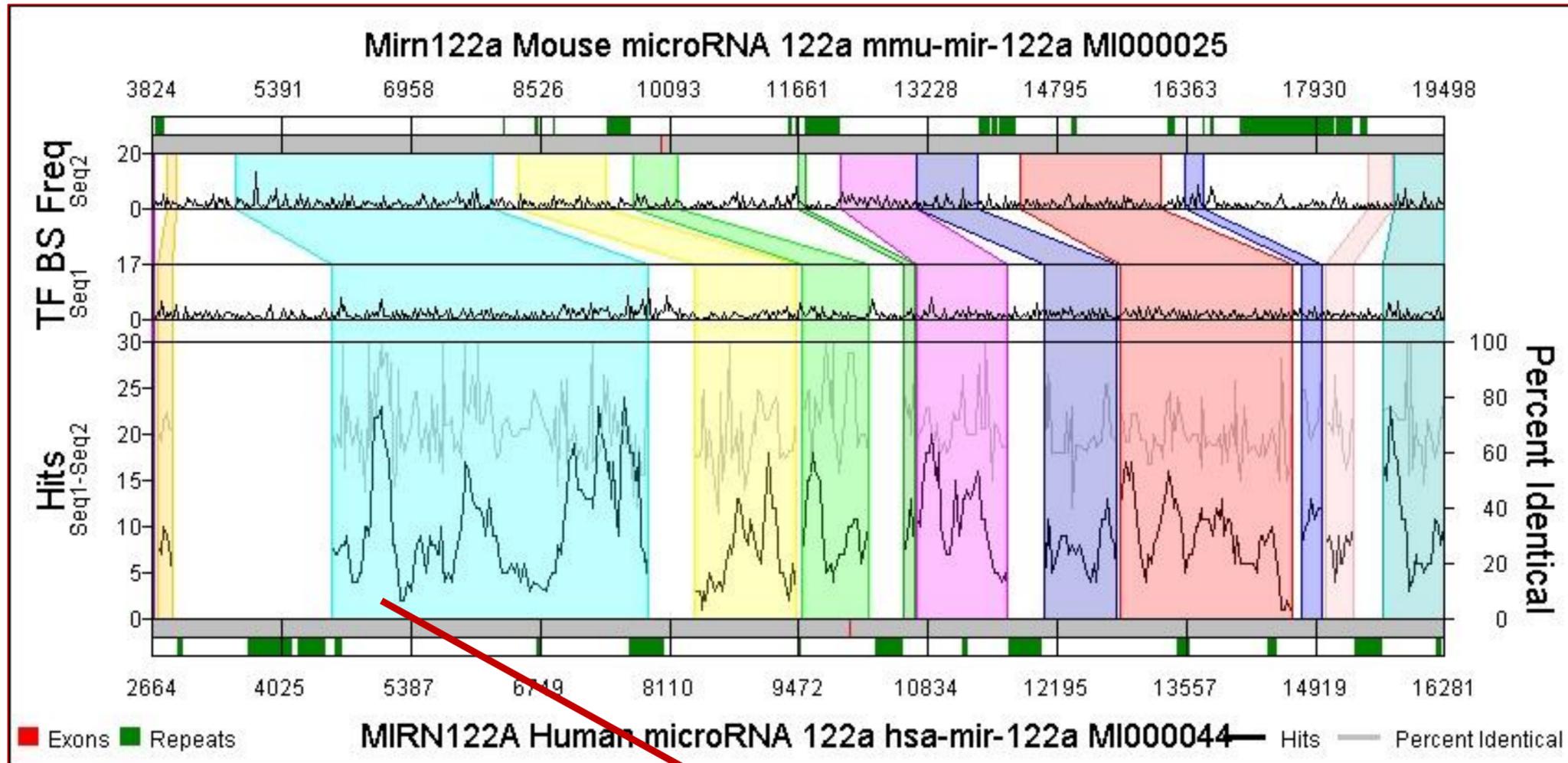
<input checked="" type="checkbox"/>	Query Term	Accession Number	Name
<input checked="" type="checkbox"/>	MIR-122A	hgMIRN122A	MIRN122A Human microRNA 122a hsa-mir-122a MI000044
<input checked="" type="checkbox"/>	MIR-122A	mgMirn122a	Mirn122a Mouse microRNA 122a mmu-mir-122a MI000025

# GenomeTrafac: <http://genometrafac.cchmc.org>

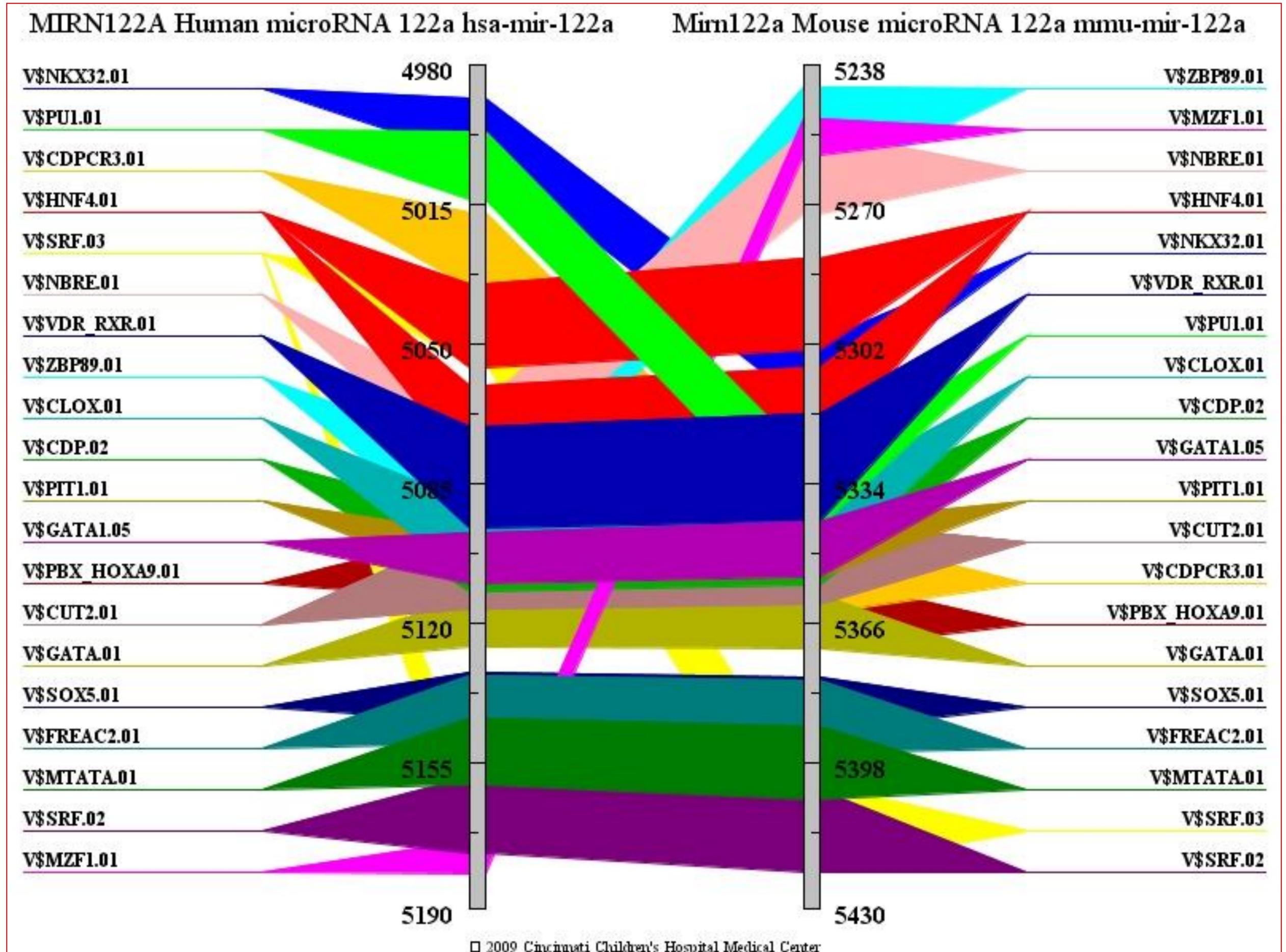
First Sequence	Second Sequence	Timestamp	Action
hgMIRN122A, MIRN122A Human microRNA 122a hsa-mir-122a MI000044	mgMirn122a, Mirn122a Mouse microRNA 122a mmu-mir-122a MI000025	07/25/2006 12:00	<a href="#">View</a> <a href="#">Regulogram</a>



# GenomeTrafac: <http://genometrafac.cchmc.org>



# GenomeTrafac: <http://genometrafac.cchmc.org>



# GenomeTrafac: <http://genometrafac.cchmc.org>

## Shared *Cis*-elements



(Genomatix Matrix Family Library Version 5.0 (January 2005))

(For details and annotations of TFBS-PWMs, please register at [Genomatix](#))

Family/Matrix	Description	<a href="#">hgMIRN122A</a>				<a href="#">mgMirn122a</a>			
		Begin	End	Sequence		Begin	End	Sequence	
<a href="#">V\$NKXH/V\$NKX32.01</a>	Homeodomain protein NKX3.2 (BAPX1, NKX3B, Bagpipe homolog)	4993	5007	CCCCACTCAGCAGA	-	5301	5315	CTGACTTAGTGGACT	+
<a href="#">V\$ETSF/V\$PU1.01</a>	Pu. 1 (Pu120) Ets-like transcription factor identified in lymphoid B-cells	5001	5017	CAGCAGAGGAATGGACT	+	5326	5342	CCTCTCTCCCCCACA	-
<a href="#">V\$CLOX/V\$CDPCR3.01</a>	Cut-like homeodomain protein	5020	5038	CCAATCTTGCTGAGTGTGT	-	5343	5361	TCGATAATTTAATGTGACT	-
<a href="#">V\$HNF4/V\$HNF4.01</a>	Hepatic nuclear factor 4	5037	5057	GTTTGACCAAAGGTGGTGCTG	+	5283	5303	GTTTGACCAAAGGTGACTCTG	+
<a href="#">V\$SRFF/V\$SRF.03</a>	Serum responsive factor	5038	5056	TTTGACCAAAGGTGGTGCT	-	5399	5417	GGATCCCATAAAGGGAGAG	-
<a href="#">V\$HNF4/V\$HNF4.01</a>	Hepatic nuclear factor 4	5061	5081	TAGTGGCCTAAGGTCGTGCC	+	5307	5327	TAGTGGACTAAGGTCATGCC	+
<a href="#">V\$RORA/V\$NBRE.01</a>	Monomers of the nur subfamily of nuclear receptors (nur77, nurr1, nor-1)	5065	5083	GGCCTAAGGTCGTGCCCTC	+	5255	5273	GGGAGCTGGACCTTCGGTT	-
<a href="#">V\$RXRF/V\$VDR_RXR.01</a>	VDR/RXR Vitamin D receptor RXR heterodimer site	5071	5095	AGGTCGTGCCCTCCCTCCCCACTG	-	5317	5341	AGGTCATGCCCTCTCTCCCCACA	-
<a href="#">V\$ZBPF/V\$ZBP89.01</a>	Zinc finger transcription factor ZBP-89	5077	5099	TGCCCTCCCTCCCCACTGAATC	+	5245	5267	GGGGCATGGGGGAGCTGGACCT	-
<a href="#">V\$CLOX/V\$CLOX.01</a>	Clox	5089	5107	CCCCTGAATCGATAAATA	+	5334	5352	CCCCACAATCGATAATTT	+

# I have a list of co-expressed mRNAs (Transcriptome)....

## Now what?

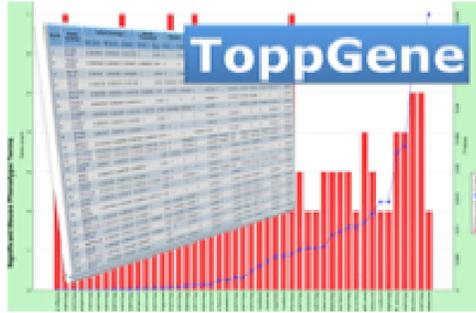
### 1. Identify putative shared regulatory elements

- Known transcription factor binding sites (TFBS)
  - Conserved
  - Non-conserved
- Unknown TFBS or Novel motifs
  - Conserved
  - Non-conserved
- MicroRNAs

### 2. Identify the underlying biological theme

- Gene Ontology
- Pathways
- Phenotype/Disease Association
- Protein Domains
- Protein Interactions
- Expression in other tissues/experiments
- Drug targets
- Literature co-citation...

# ToppGene Suite (<http://toppgene.cchmc.org>)



## ToppGene Suite

A one-stop portal for gene list enrichment analysis and candidate gene prioritization based on functional annotations and protein interactions network

- Home
- Links
- Database details
- Supplementary
- Help
- Publications
- Terms of Use
- Contacts

Supported by:

Computational  
Medicine  
Center



- **ToppFun:** Transcriptome, ontology, phenotype, proteome, and pharmacome annotations based gene list functional enrichment analysis  
Detect functional enrichment of your gene list based on Transcriptome, Proteome, Regulome (TFBS and miRNA), Ontologies (GO, Pathway), Phenotype (human disease and mouse phenotype), Pharmacome (Drug-Genes associations), literature co-citation, and other features.
- **ToppGene:** Candidate gene prioritization  
Prioritize or rank candidate genes based on functional similarity to training gene list.
- **ToppNet:** Relative importance of candidate genes in networks  
Prioritize or rank candidate genes based on topological features in protein-protein interaction network.
- **ToppGenet:** Prioritization of neighboring genes in protein-protein interaction network  
Identify and prioritize the neighboring genes of the seeds in protein-protein interaction network based on functional similarity to the "seed" list (ToppGene) or topological features in protein-protein interaction network (ToppNet).

# ToppGene Suite (<http://toppgene.cchmc.org>)

ToppFun: Transcriptome, ontology, phenotype, proteome, and pharmacome annotations based gene list functional enrichment analysis

Select your gene identifier type, paste your sets below or select example set, then submit.

Entry Type:

Example gene sets: [HGNC Symbol](#) [Entrez ID](#)  
(click on "HGNC Symbol" or "Entrez ID" to use the example training and test set of genes)

Training Gene Set:

259
5265
350
335
335
1558
1571
229
462
125
3240
5105
5265
3273
2244
2158
5053
125
1356
3827
383

Clear

Submit Query

## Input Gene List (81 / 97)

Entered	Human Symbol	Gene ID
259	AMBP	259
5265	SERPINA1	5265
350	APOH	350
335	APOA1	335
1558	CYP2C8	1558
1571	CYP2E1	1571
229	ALDOB	229
462	SERPINC1	462
125	ADH1B	125
3240	HP	3240
5105	PCK1	5105
3273	HRG	3273
2244	FGB	2244
2158	F9	2158
5053	PAH	5053
1356	CP	1356
3827	KNG1	3827
383	ARG1	383
5004	ORM1	5004
2168	FABP1	2168
325	APCS	325

## Genes Not Found

Entered	Status
335	Duplicated
5265	Duplicated
125	Duplicated
1571	Duplicated
1571	Duplicated
1373	Duplicated
1356	Duplicated
462	Duplicated
2243	Duplicated
3827	Duplicated
125	Duplicated
229	Duplicated
6822	Duplicated
2328	Duplicated

1. Supports variety of inputs
2. Supports symbol correction
3. Eliminates any duplicates
4. Drawback: Supports human and mouse genes only

# ToppGene Suite (<http://toppgene.cchmc.org>)

## Calculations

Feature	Correction	p-Value cutoff	Gene Limits
<input checked="" type="checkbox"/> All	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> GO: Molecular Function	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> GO: Biological Process	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> GO: Cellular Component	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Human Phenotype	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Mouse Phenotype	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Domain	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Pathway	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Pubmed	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Interaction	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Cytoband	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> TFBS	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Gene Family	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Coexpression	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Computational	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> MicroRNA	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Drug	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Disease	Bonferroni	0.05	1 ≤ n ≤ 1500

Home

Modify Query

Submit

1. Gene list analyzed for as many as 17 features!
2. Single-stop enrichment analysis server for both regulatory elements (TFBSs and miRNA) and biological themes
3. Back-end has an exhaustive, normalized data resources compiled and integrated
4. Bonferroni correction is “too stringent”; FDR with 0.05 is preferable.
5. TFBS are based on conserved cis-elements and motifs within  $\pm 2\text{kb}$  region of TSS in human, mouse, rat, and dog.
6. miRNA-targets are based on TargetScan and PicTar

# ToppGene Suite (<http://toppgene.cchmc.org>)

<b>GO Biological Process</b>		<b>Human Phenotype</b>		<b>Mouse Phenotype</b>	
Annotations:	16,372	Annotations:	9,551	Annotations:	6,203
Genes:	15,079	Genes:	2,531	Genes:	5,590
	Updated Aug 26, 2009		Updated Sep 10, 2009		Updated Aug 25, 2009
<b>GO Cellular Component</b>					
Annotations:	2,335				
Genes:	16,728				
	Updated Aug 26, 2009				
<b>GO Molecular Function</b>					
Annotations:	8,583				
Genes:	15,948				
	Updated Aug 26, 2009				
<b>Pathways</b>		<b>Domains</b>		<b>Pubmed</b>	
Annotations:	1,672	Annotations:	10,223	Annotations:	221,282
	BioCyc 164		Gene3D 285	Genes:	22,176
Aug 25, 2009	CGAP BioCarta 314		InterPro 4,859		Updated Aug 25, 2009
	GenMAPP 67		PROSITE 1,351		
Jun 15, 2009	KEGG pathway 202		Pfam 2,774		
May 10, 2009	MSigDB 431		ProDom 385		
	PantherDB 150		SMART 569		
Aug 25, 2009	Pathway Ontology 306	Genes:	12,430		
	Reactome 25				
	SigmaAldrich 2				
	Signalling Transduction KE 11				
Genes:	6,697				
<b>Interactions</b>		<b>Cytoband</b>		<b>TFBS</b>	
Annotations:	18,047	Annotations:	382	Annotations:	615
	BIND 4,370	Genes:	29,821	Genes:	9,770
	BioGRID 7,602				
	HPRD 6,075				
Genes:	5,541				
<b>miRNA</b>		<b>Gene Families</b>		<b>Coexpression</b>	
Annotations:	740	Annotations:	151	Annotations:	1,203
	MSigDB 313	Genes:	6,098		Body Map 23
	PicTar 178				mSigDB 1,180
	TargetScan 249			Genes:	12,694
Genes:	11,618				
<b>Computational Gene Set</b>		<b>Drugs</b>		<b>Disease</b>	
Annotations:	427	Annotations:	13,141	Annotations:	3,789
Genes:	4,712	Aug 28, 2009	CTD 4,977	Aug 28, 2009	CTD 1,006
		Aug 25, 2009	Drug Bank 2,009		GWAS 291
		Aug 25, 2009	Stitch 6,155		OMIM 2,492
		Genes:	14,836	Genes:	4,385
<b>Master Gene Info File</b>					
For All Annotations	35,449				
	Updated Aug 28, 2009				

1. Database updated regularly
2. Exhaustive collection of annotations

# ToppGene Suite (<http://toppgene.cchmc.org>)

Results

[Go To Start Page](#)

**Input Parameters** [\[Show Detail\]](#)

**Training Results** [\[Show All\]](#) [\[Download All\]](#) [\[Sparse Matrix\]](#)

- 1: GO: Molecular Function [\[Display Chart\]](#) [\[Show Detail\]](#)
- 2: GO: Biological Process [\[Display Chart\]](#) [\[Show Detail\]](#)
- 3: GO: Cellular Component [\[Display Chart\]](#) [\[Show Detail\]](#)
- 4: Human Phenotype [\[Display Chart\]](#) [\[Show Detail\]](#)
- 5: Mouse Phenotype [\[Display Chart\]](#) [\[Show Detail\]](#)
- 6: Domain [\[Display Chart\]](#) [\[Show Detail\]](#)
- 7: Pathway [\[Display Chart\]](#) [\[Show Detail\]](#)
- 8: Pubmed [\[Display Chart\]](#) [\[Show Detail\]](#)
- 9: Interaction [\[Display Chart\]](#) [\[Show Detail\]](#)
- 10: Cytoband [\[Display Chart\]](#) [\[Show Detail\]](#)
- 11: TFBS [\[Display Chart\]](#) [\[Show Detail\]](#)
- 12: Gene Family [\[Display Chart\]](#) [\[Show Detail\]](#)
- 13: Coexpression [\[Display Chart\]](#) [\[Show Detail\]](#)
- 14: Computational [\[Display Chart\]](#) [\[Show Detail\]](#)
- 15: MicroRNA [\[Display Chart\]](#) [\[Show Detail\]](#)
- 16: Drug [\[Display Chart\]](#) [\[Show Detail\]](#)
- 17: Disease [\[Display Chart\]](#) [\[Show Detail\]](#)

**Input Parameters** [\[Hide Detail\]](#)

Number of genes in training set:	81																																																																																										
Number of genes in test set:	0																																																																																										
Correction and Cutoff:	<table border="1" style="width: 100%; border-collapse: collapse;"> <thead> <tr> <th>category</th> <th>Correction</th> <th>Cutoff</th> <th>Min</th> <th>Max</th> </tr> </thead> <tbody> <tr><td>GO: Molecular Function</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> <tr><td>GO: Biological Process</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> <tr><td>GO: Cellular Component</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> <tr><td>Human Phenotype</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> <tr><td>Mouse Phenotype</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> <tr><td>Domain</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> <tr><td>Pathway</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> <tr><td>Pubmed</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> <tr><td>Interaction</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> <tr><td>Cytoband</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> <tr><td>TFBS</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> <tr><td>Gene Family</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> <tr><td>Coexpression</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> <tr><td>Computational</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> <tr><td>MicroRNA</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> <tr><td>Drug</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> <tr><td>Disease</td><td>Bonferroni</td><td>0.05</td><td>1</td><td>1500</td></tr> </tbody> </table>	category	Correction	Cutoff	Min	Max	GO: Molecular Function	Bonferroni	0.05	1	1500	GO: Biological Process	Bonferroni	0.05	1	1500	GO: Cellular Component	Bonferroni	0.05	1	1500	Human Phenotype	Bonferroni	0.05	1	1500	Mouse Phenotype	Bonferroni	0.05	1	1500	Domain	Bonferroni	0.05	1	1500	Pathway	Bonferroni	0.05	1	1500	Pubmed	Bonferroni	0.05	1	1500	Interaction	Bonferroni	0.05	1	1500	Cytoband	Bonferroni	0.05	1	1500	TFBS	Bonferroni	0.05	1	1500	Gene Family	Bonferroni	0.05	1	1500	Coexpression	Bonferroni	0.05	1	1500	Computational	Bonferroni	0.05	1	1500	MicroRNA	Bonferroni	0.05	1	1500	Drug	Bonferroni	0.05	1	1500	Disease	Bonferroni	0.05	1	1500
category	Correction	Cutoff	Min	Max																																																																																							
GO: Molecular Function	Bonferroni	0.05	1	1500																																																																																							
GO: Biological Process	Bonferroni	0.05	1	1500																																																																																							
GO: Cellular Component	Bonferroni	0.05	1	1500																																																																																							
Human Phenotype	Bonferroni	0.05	1	1500																																																																																							
Mouse Phenotype	Bonferroni	0.05	1	1500																																																																																							
Domain	Bonferroni	0.05	1	1500																																																																																							
Pathway	Bonferroni	0.05	1	1500																																																																																							
Pubmed	Bonferroni	0.05	1	1500																																																																																							
Interaction	Bonferroni	0.05	1	1500																																																																																							
Cytoband	Bonferroni	0.05	1	1500																																																																																							
TFBS	Bonferroni	0.05	1	1500																																																																																							
Gene Family	Bonferroni	0.05	1	1500																																																																																							
Coexpression	Bonferroni	0.05	1	1500																																																																																							
Computational	Bonferroni	0.05	1	1500																																																																																							
MicroRNA	Bonferroni	0.05	1	1500																																																																																							
Drug	Bonferroni	0.05	1	1500																																																																																							
Disease	Bonferroni	0.05	1	1500																																																																																							
Random sampling size in analysis:	0																																																																																										
Minimum feature count in test set:	2																																																																																										
Analysis took:	2 seconds																																																																																										
Analysis finished at:	Sun Sep 27 16:45:06 EDT 2009																																																																																										

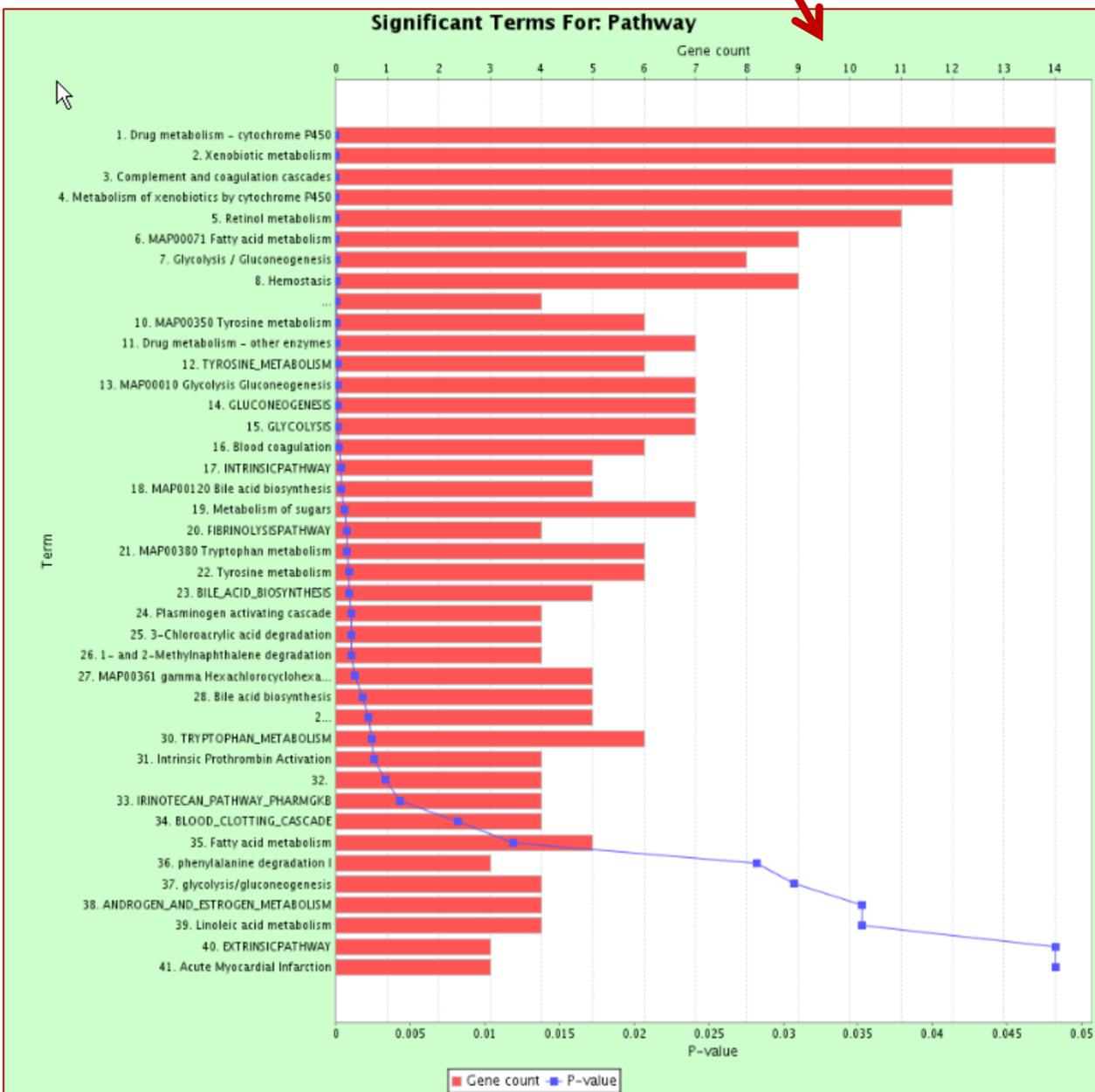
**2: GO: Biological Process** [\[Display Chart\]](#) [\[Hide Detail\]](#)

	ID	Name	Source	P-value	Term in Query	Term in Genome
1	GO:0009605	response to external stimulus		0	27	893
2	GO:0007513	blood coagulation		0	12	115
3	GO:0006629	lipid metabolic process		0	25	874
4	GO:0044255	cellular lipid metabolic process		0	23	720
5	GO:0050817	coagulation		0	12	119
6	GO:0007599	hemostasis		0	12	120
7	GO:0009611	response to wounding		0	20	542
8	GO:0042060	wound healing		0	13	185
9	GO:0050878	regulation of body fluid levels		0	12	151
10	GO:0055114	oxidation reduction		0	19	624
11	GO:0019752	carboxylic acid metabolic process		0	18	570

# ToppGene Suite (<http://toppgene.cchmc.org>)

2: GO: Biological Process [Display Chart] [Hide Detail]

ID	Name	Source	P-value	Term in Query	Term in Genome
1	GO:0009605 response to external stimulus		0	27	893
2	GO:0007514 blood coagulation		0	12	115
3	GO:0006629 lipid metabolic process		0	25	874
4	GO:0044255 cellular lipid metabolic process		0	23	720
5	GO:0050817 coagulation		0	12	119
6	GO:0007599 hemostasis		0	12	120
7	GO:0009611 response to wounding		0	20	542
8	GO:0042060 wound healing		0	13	185
9	GO:0050878 regulation of body fluid levels		0	12	151
10	GO:0055114 oxidation reduction		0	19	624
11	GO:0019752 carboxylic acid metabolic process		0	18	570



response to external stimulus; GO:0009605

	Entrez Gene ID	Gene Symbol	Gene Name	Original Symbol
1	126	ADH1C	alcohol dehydrogenase 1C (class I), gamma polypeptide	126
2	335	APOA1	apolipoprotein A-I	335
3	350	APOH	apolipoprotein H (beta-2-glycoprotein I)	350
4	2158	F9	coagulation factor IX	2158
5	5950	RBP4	retinol binding protein 4, plasma	5950
6	197	AHSG	alpha-2-HS-glycoprotein	197
7	2243	FGA	fibrinogen alpha chain	2243
8	213	ALB	albumin	213
9	2244	FGB	fibrinogen beta chain	2244
10	629	CFB	complement factor B	629
11	3158	HMGCS2	3-hydroxy-3-methylglutaryl-Coenzyme A synthase 2 (mitochondrial)	3158
12	5444	PON1	paraoxonase 1	5444
13	1361	CPB2	carboxypeptidase B2 (plasma)	1361
14	3078	CFHR1	complement factor H-related 1	3078
15	5265	SERPINA1	serpin peptidase inhibitor, clade A (alpha-1 antiproteinase, antitrypsin), member 1	5265
16	3827	KNG1	kininogen 1	3827
17	325	APCS	amyloid P component, serum	325
18	2538	G6PC	glucose-6-phosphatase, catalytic subunit	2538
19	4153	MBL2	mannose-binding lectin (protein C) 2, soluble (opsonic defect)	4153
20	735	C9	complement component 9	735
21	462	SERPINC1	serpin peptidase inhibitor, clade C (antithrombin), member 1	462
22	3273	HRG	histidine-rich glycoprotein	3273
23	5340	PLG	plasminogen	5340
24	5004	ORM1	orosomucoid 1	5004
25	316	AOX1	aldehyde oxidase 1	316
26	3053	SERPIND1	serpin peptidase inhibitor, clade D (heparin cofactor), member 1	3053
27	1356	CP	ceruloplasmin (ferroxidase)	1356

# ToppGene Suite (<http://toppgene.cchmc.org>)

**ToppGene Result Page**

Number of genes in training set: 81  
Number of genes in test set: 0

Correction and Cutoff:

category	Correction	Cutoff	Min	Max
GO: Molecular Function	Bonferroni	0.05	1	1500
GO: Biological Process	Bonferroni	0.05	1	1500
GO: Cellular Component	Bonferroni	0.05	1	1500
Human Phenotype	Bonferroni	0.05	1	1500
Mouse Phenotype	Bonferroni	0.05	1	1500
Domain	Bonferroni	0.05	1	1500
Pathway	Bonferroni	0.05	1	1500
Pubmed	Bonferroni	0.05	1	1500
Interaction	Bonferroni	0.05	1	1500
Cytoband	Bonferroni	0.05	1	1500
TFBS	Bonferroni	0.05	1	1500
Gene Family	Bonferroni	0.05	1	1500
Coexpression	Bonferroni	0.05	1	1500
Computational	Bonferroni	0.05	1	1500
MicroRNA	Bonferroni	0.05	1	1500
Drug	Bonferroni	0.05	1	1500
Disease	Bonferroni	0.05	1	1500

Random sampling size in analysis: 0  
Minimum feature count in test set: 2  
Analysis took: 2 seconds  
Analysis finished at: Sun Sep 27 16:45:06 EDT 2009

Enter name of file to save to...

Save in: Desktop

My Recent Documents  
Desktop  
My Documents  
My Computer  
My Network

Document2  
TOB1\_p53  
Refrig-Bulb  
miRNAPromoters\_192  
Anil  
My Computer  
Imp\_Dates  
Bioinfo\_Workshop-2009  
Unused Desktop Shortcuts  
Disease CVs  
p53mhw  
songs  
New Folder  
Photos  
Misc

MyWebSite  
My Network Places  
My Computer  
My Documents

File name: LiverGenes\_ToppFun.txt  
Save as type: Text Document

Save Cancel

## Training Results [\[Show All\]](#) [\[Download All\]](#) [\[Sparse Matrix\]](#)

- 1: GO: Molecular Function [\[Display Chart\]](#) [\[Show Detail\]](#)
- 2: GO: Biological Process [\[Display Chart\]](#) [\[Show Detail\]](#)
- 3: GO: Cellular Component [\[Display Chart\]](#) [\[Show Detail\]](#)
- 4: Human Phenotype [\[Display Chart\]](#) [\[Show Detail\]](#)
- 5: Mouse Phenotype [\[Display Chart\]](#) [\[Show Detail\]](#)
- 6: Domain [\[Display Chart\]](#) [\[Show Detail\]](#)
- 7: Pathway [\[Display Chart\]](#) [\[Show Detail\]](#)
- 8: Pubmed [\[Display Chart\]](#) [\[Show Detail\]](#)

# ToppGene Suite (<http://toppgene.cchmc.org>)

I have a list of 200 over-expressed genes and I want to prioritize them for experimental validation (apart from using the fold change as a parameter).....

## ToppGene Suite

A one-stop portal for gene list enrichment analysis and candidate gene prioritization based on functional annotations and protein interactions network

- **ToppFun:** Transcriptome, ontology, phenotype, proteome, and pharmacome annotations based gene list functional enrichment analysis

Detect functional enrichment of your gene list based on Transcriptome, Proteome, Regulome (TFBS and miRNA), Ontologies (GO, Pathway), Phenotype (human disease and mouse phenotype), Pharmacome (Drug-Genes associations), literature co-citation, and other features.

- **ToppGene:** Candidate gene prioritization

Prioritize or rank candidate genes based on functional similarity to training gene list.

- **ToppNet:** Relative importance of candidate genes in networks

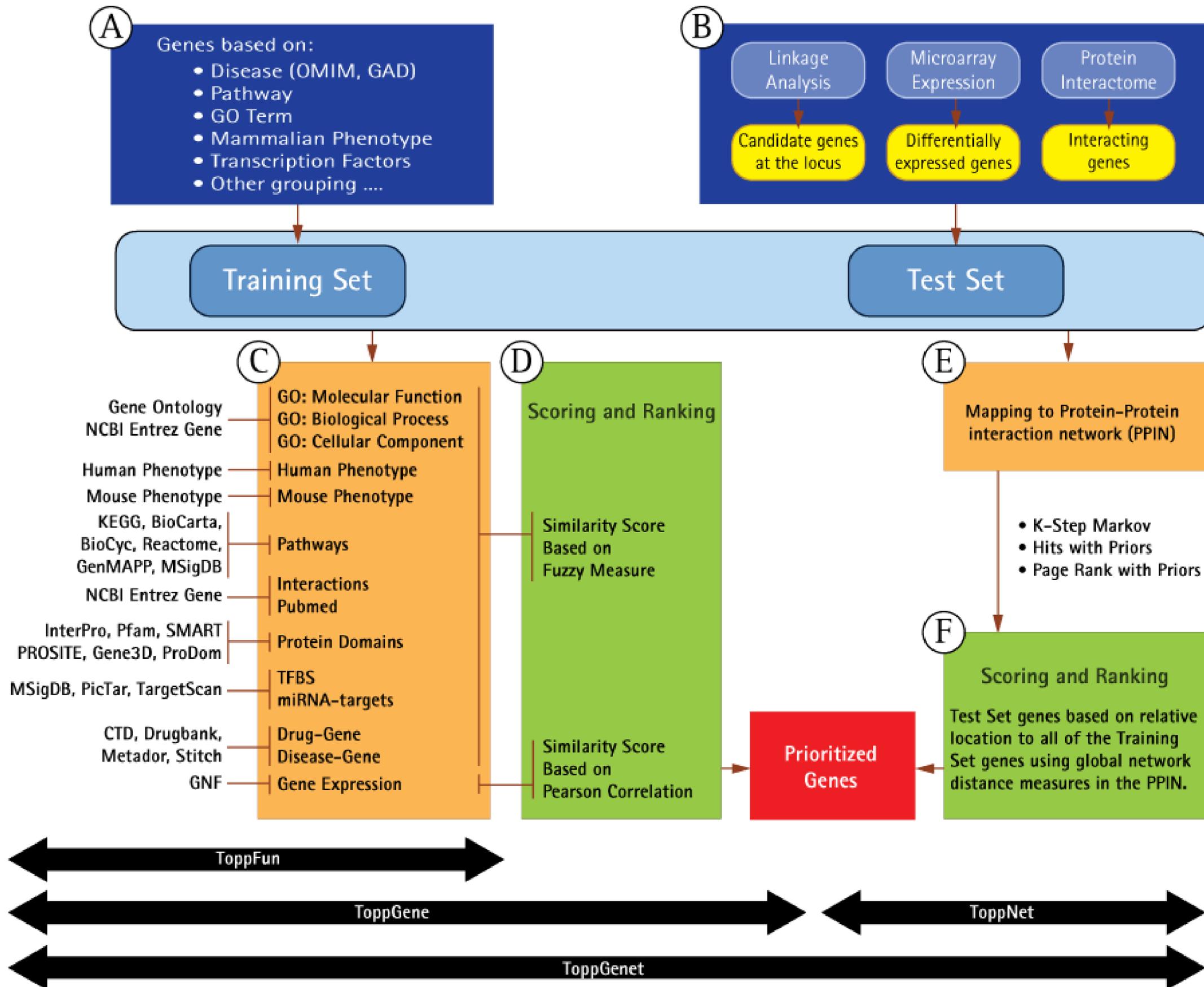
Prioritize or rank candidate genes based on topological features in protein-protein interaction network.

- **ToppGenet:** Prioritization of neighboring genes in protein-protein interaction network

Identify and prioritize the neighboring genes of the seeds in protein-protein interaction network based on functional similarity to the "seed" list (ToppGene) or topological features in protein-protein interaction network (ToppNet).

# ToppGene Suite (<http://toppgene.cchmc.org>)

I have a list of 200 over-expressed genes and I want to prioritize them for experimental validation (apart from using the fold change as a parameter).....



# ToppGene Suite (<http://toppgene.cchmc.org>)

## ToppGene: Candidate gene prioritization

Select your gene identifier type, paste your training and test gene sets below or select example sets, then submit.

Example gene sets: [HGNC Symbol](#) [Entrez ID](#)  
(click on "HGNC Symbol" or "Entrez ID" to use the example training and test set of genes)

Symbol Types

Training Gene Set:

NKX2-5  
MEF2A  
GATA4  
HAND1  
HAND2  
TBX5  
SRF

Test gene set:

ACVR1  
ACVR2B  
ADAM19  
ADM  
ADRA1A  
ADRA1B  
ADRBK1  
ALDH1A2  
ALPK3  
ATP6VOA1  
BMP10  
BMP2  
BMP4  
BMPR1A  
CALCRL  
CASP3  
CASP7  
CASP8  
CASQ2  
CENTA2  
CHD7

Clear

Submit Query

# ToppGene Suite (<http://toppgene.cchmc.org>)

**Training set (7 / 7)**

Entered	Human Symbol	Gene ID
NKX2-5	NKX2-5	1482
MEF2A	MEF2A	4205
GATA4	GATA4	2626
HAND1	HAND1	9421
HAND2	HAND2	9464
TBX5	TBX5	6910
SRF	SRF	6722

**Test set (146 / 158)**

Entered	Human Symbol	Gene ID
ACVR1	ACVR1	90
ACVR2B	ACVR2B	93
ADAM19	ADAM19	8728
ADM	ADM	133
ADRA1A	ADRA1A	148
ADRA1B	ADRA1B	147
ADRBK1	ADRBK1	156
ALDH1A2	ALDH1A2	8854
ALPK3	ALPK3	57538
ATP6V0A1	ATP6V0A1	535
BMP10	BMP10	27302
BMP2	BMP2	650
BMP4	BMP4	652
BMPR1A	BMPR1A	657
CALCRL	CALCRL	10203
CASP3	CASP3	836
CASP7	CASP7	840
CASP8	CASP8	841
CASQ2	CASQ2	845
CHD7	CHD7	55636
CITED2	CITED2	10370

Entered	Suggestions
CENTA2	<input checked="" type="checkbox"/> ADAP2 - ArfGAP with dual PH domains 2 <a href="#">Human Synonym</a>
CMYA1	<input checked="" type="checkbox"/> XIRP1 - xin actin-binding repeat containing 1 <a href="#">Human Synonym</a>
GJA7	<input checked="" type="checkbox"/> GJC1 - gap junction protein, gamma 1, 45kDa <a href="#">Human Synonym</a>
HOP	<input checked="" type="checkbox"/> HOPX - HOP homeobox <a href="#">Human Synonym</a> <input type="checkbox"/> ST13 - suppression of tumorigenicity 13 (colon carcinoma) (Hsp70 interacting protein) <a href="#">Human Synonym</a> <input type="checkbox"/> STIP1 - stress-induced-phosphoprotein 1 <a href="#">Human Synonym</a>
PPARBP	<input checked="" type="checkbox"/> MED1 - mediator complex subunit 1 <a href="#">Human Synonym</a>
RBPSUH	RBPJ Duplicated

Check All

**Ignored**

Entered	Status
CENTA2	Not Found
CMYA1	Not Found
GATA4	In Training Set
GJA7	Not Found
HAND1	In Training Set
HAND2	In Training Set
HOP	Not Found
NKX2-5	In Training Set
PPARBP	Not Found
RBPSUH	Not Found
SRF	In Training Set
TBX5	In Training Set

[Find alternatives for missing symbols](#)

Entered Suggestions

CENTA2  ADAP2 - ArfGAP with dual PH domains 2 [Human Synonym](#)

CMYA1  XIRP1 - xin actin-binding repeat containing 1 [Human Synonym](#)

GJA7  GJC1 - gap junction protein, gamma 1, 45kDa [Human Synonym](#)

HOP  HOPX - HOP homeobox [Human Synonym](#)

ST13 - suppression of tumorigenicity 13 (colon carcinoma) (Hsp70 interacting protein) [Human Synonym](#)

STIP1 - stress-induced-phosphoprotein 1 [Human Synonym](#)

PPARBP  MED1 - mediator complex subunit 1 [Human Synonym](#)

RBPSUH RBPJ Duplicated

Check All

Ignored

Entered	Status
CENTA2	Not Found
CMYA1	Not Found
GATA4	In Training Set
GJA7	Not Found
HAND1	In Training Set
HAND2	In Training Set
HOP	Not Found
NKX2-5	In Training Set
PPARBP	Not Found
RBPSUH	Not Found
SRF	In Training Set
TBX5	In Training Set

[Find alternatives for missing symbols](#)

# ToppGene Suite (<http://toppgene.cchmc.org>)

Training parameters

Feature	Correction	p-Value cutoff	Gene Limits
<input checked="" type="checkbox"/> All	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> GO: Molecular Function	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> GO: Biological Process	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> GO: Cellular Component	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Human Phenotype	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Mouse Phenotype	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Domain	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Pathway	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Pubmed	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Interaction	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Cytoband	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> TFBS	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Gene Family	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Coexpression	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Computational	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> MicroRNA	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Drug	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Disease	Bonferroni <input type="button" value="v"/>	0.05 <input type="button" value="v"/>	1 ≤ n ≤ 1500

Test parameter

Random sampling size: 1500 (6% of genome)

Min. feature count: 2

Home

Modify Query

Start prioritization

ToppGene is processing your query



Estimating p-Values

To see the training results before the test set is complete, [click here](#).

# ToppGene Suite (<http://toppgene.cchmc.org>)

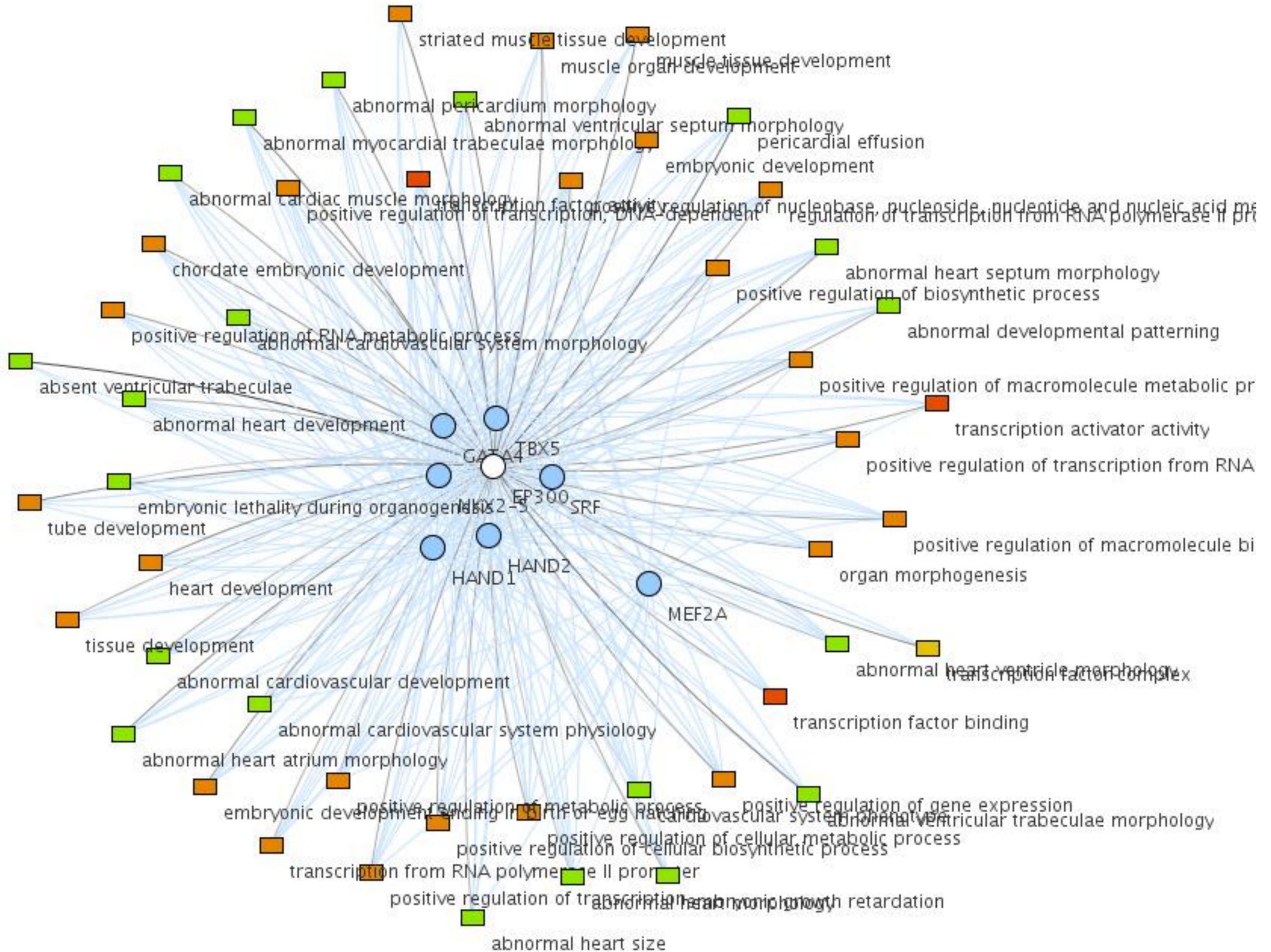
Rank	Gene Symbol	Gene ID	GO: Molecular Function		GO: Biological Process		GO: Cellular Component		Human Phenotype		Mouse Phenotype		Domain		Pathway		Pubmed		Interaction		Cytoband	
			Score	pValue	Score	pValue	Score	pValue	Score	pValue	Score	pValue	Score	pValue	Score	pValue	Score	pValue	Score	pValue	Score	pValue
1	EP300	2033	0.7136999	0.0235641	0.9999726	0.0029455	0.4305814	0.001	0	0.5	0.9999881	0.005			0	0.5049834	0.8410056	0.001	0.7991753	0.001	0	0.5004907
2	TEAD1	7003	0.5804123	0.0309278	0.997771	0.0103093	0.4305814	0.001	0	0.5	0.9989337	0.035	0	0.5			0.7993714	0.001	0.7160863	0.001	0	0.5004907
3	HIF1A	3091	0.9391513	0.0029455	1	0.001	0.4305814	0.001			0.9997067	0.02	0.8885697	0.001	0	0.5049834	0	0.4995093	0	0.505	0	0.5004907
4	CTNNB1	1499	0.7136999	0.0235641	1	0.001	0.4305814	0.001	0	0.5	0.9529218	0.06			0.6371253	0.001	0	0.4995093	0	0.505	0	0.5004907
5	TBX20	57057	0.5804123	0.0309278	0.9999964	0.0014728	0	0.5022091			0.9999902	0.005	0	0.5			0	0.4995093			0	0.5004907
6	ZFPM2	23414	0.6308709	0.0250368	0.9999978	0.001	0	0.5022091	0	0.5	1	0.001	0	0.5			0	0.4995093			0	0.5004907
7	BMP4	652	0	0.5493373	1	0.001	0	0.5022091	0	0.5	0.9999435	0.01	0	0.5	0.6478057	0.001	0.7993714	0.001	0	0.505	0	0.5004907
8	TBX1	6899	0.9660807	0.001	0.999997	0.001	0	0.5022091	0	0.5	0.9996966	0.02	0	0.5			0	0.4995093			0	0.5004907
9	TBX2	6909	0.5418852	0.0397644	0.9943991	0.0162003	0.4305814	0.001			0.9993508	0.035	0	0.5	0	0.5049834	0	0.4995093			0	0.5004907
10	TGFB2	7042	0.8603852	0.005891	1	0.001	0	0.5022091			0.9999998	0.005	0	0.5	0.3937178	0.0033223	0	0.4995093	0	0.505	0	0.5004907

Rank	Gene Symbol
1	EP300
2	TEAD1
3	HIF1A
4	CTNNB1
5	TBX20
6	ZFPM2
7	BMP4
8	TBX1
9	TBX2
10	TGFB2

Average score	Overall P-value
0.3417445	0.0000003
0.3015437	0.0000058
0.3041435	0.0000062
0.2489552	0.0000788
0.3207447	0.0000893
0.2749466	0.000112
0.229808	0.0001787
0.2395215	0.0002528
0.2566618	0.0002615
0.2503156	0.0002619
0.3307561	0.0002975

# ToppGene Suite (<http://toppgene.cchmc.org>)

## Why is a test set gene ranked higher?



# ToppGene Suite (<http://toppgene.cchmc.org>)

I have a list of 200 over-expressed genes and I want to prioritize them for experimental validation (apart from using the fold change as a parameter).....

## ToppGene Suite

A one-stop portal for gene list enrichment analysis and candidate gene prioritization based on functional annotations and protein interactions network

- **ToppFun:** Transcriptome, ontology, phenotype, proteome, and pharmacome annotations based gene list functional enrichment analysis

Detect functional enrichment of your gene list based on Transcriptome, Proteome, Regulome (TFBS and miRNA), Ontologies (GO, Pathway), Phenotype (human disease and mouse phenotype), Pharmacome (Drug-Gene associations), literature co-citation, and other features.

- **ToppGene:** Candidate gene prioritization

Prioritize or rank candidate genes based on functional similarity to training gene list.

- **ToppNet:** Relative importance of candidate genes in networks

Prioritize or rank candidate genes based on topological features in protein-protein interaction network.

- **ToppGenet:** Prioritization of neighboring genes in protein-protein interaction network

Identify and prioritize the neighboring genes of the seeds in protein-protein interaction network based on functional similarity to the "seed" list (ToppGene) or topological features in protein-protein interaction network (ToppNet).

# ToppGene Suite (<http://toppgene.cchmc.org>)

Graph prioritization parameters

Prioritization method:

Step Size(normally 4-8):

Training gene neighborhood subnetwork visualization parameters

Neighborhood distance:

When training set is big, the training gene neighborhood subnetwork can be huge.

[Home](#)

[Modify Query](#)

[Start prioritization](#)

## Test Genes [\[Hide All\]](#)

Rank	ID	Name	Interactant count	Score
1	2033	EP300	129	0.008192
2	23414	ZFPM2	4	0.004724
3	4776	NFATC4	4	0.004615
4	7003	TEAD1	9	0.003739
5	9734	HDAC9	20	0.002319
6	10014	HDAC5	33	0.002317
7	23054	NCOA6	49	0.001991
8	93649	MYOCD	2	0.0016
9	57496	MKL2	2	0.0016
10	1499	CTNNB1	138	0.001546



**Exercise 11: Use the gene list from the downloaded file (“Example-Set-2”) and find out:**

- a. How many of these genes are transcription factors?**
- b. What are the enriched TFBSs and miRNAs?**
- c. What gene families are enriched in this list?**
- d. Are there are salivary gland development associated genes present in this list?**
- e. How many and which genes from this list are associated with non-insulin dependent diabetes mellitus (NIDDM)?**

**Exercise 12: Prioritize the 721 genes (“Example-Set-2”) using “stomach genes” from the “Example-Set-1”.**

- a. What are the top 10 ranked genes using ToppGene and ToppNet?**
- b. Why is TFF3 ranked among the top 5 in ToppGene prioritization? What is its rank in ToppNet?**

# What if I want to compare several gene lists at a time?

## ToppCluster (<http://toppcluster.cchmc.org>)



- Navigation**
  - Main
  - Alternative Entry Methods
  - Cluster Dataset
- Information**
  - Disclaimer
  - ToppGene
- Support**
  - Documentation

Paste input list    Load Sample Data    Genes

Symbols are

Cluster Name

# ToppCluster (<http://toppcluster.cchmc.org>)

**Paste input list** [Load Sample Data](#)

Symbols are

Cluster Name

Genes
259
5265
350
335
335
1558
1571
229
462
125
3240
5105
5265
3273
2244
2158
5053
125
1356
3827
383
5004
2168
1571
325
5950
127
1373
213
735
64065
3855
100129410
8842
3514
55504
26298
64065
245973

Symbols are

Cluster Name

Options

Feature	Correction	p-Value cutoff	Gene Limits
<input checked="" type="checkbox"/> All	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> GO: Molecular Function	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> GO: Biological Process	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> GO: Cellular Component	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Human Phenotype	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Mouse Phenotype	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Domain	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Pathway	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Pubmed	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Interaction	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Cytoband	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> TFBS	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Coexpression	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Computational	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> MicroRNA	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Drug	Bonferroni	0.05	1 ≤ n ≤ 1500
<input checked="" type="checkbox"/> Disease	Bonferroni	0.05	1 ≤ n ≤ 1500

Annotations must have at least  gene(s)

Chose Toppcluster output format:

Gene Sets

Liver		
44 known - 0 unknown		
Original	Human Symbol	Entrez ID
1576	CYP3A4	1576
1577	CYP3A5	1577
7036	TFR2	7036
229	ALDOB	229
341	APOC1	341
126	ADH1C	126
125	ADH1B	125
7448	VTN	7448
5053	PAH	5053
3240	HP	3240
197	AHSG	197
3078	CFHR1	3078
383	ARG1	383

Salivary_Glands		
46 known - 0 unknown		
Original	Human Symbol	Entrez ID
2591	GALNT3	2591
9073	CLDN8	9073
486	FXD2	486
54959	ODAM	54959
5349	FXD3	5349
155006	TMEM213	155006
100129410	LOC100129410	100129410
54097	FAM3B	54097
352999	C6orf58	352999
57535	KIAA1324	57535
26298	EHF	26298
999	CDH1	999
360	AQP3	360

# ToppCluster (http://toppcluster.cchmc.org)

ToppCluster

Processing Salivary\_Glands

Navigate

Jump To ...

Links

Back to Start  
Shareable Link

Cytoscape

Include Orphaned Genes  
 Include Super Category  
 XGMML  
 Build

Re-Enrich

- Navigate
- Jump To ...
- Jump To ...
  - GO: Molecular Function
  - GO: Biological Process
  - GO: Cellular Component
  - Human Phenotype
  - Mouse Phenotype
  - Domain
  - Pathway
  - Pubmed
  - Interaction
  - Cytoband
  - TFBS
  - Coexpression
  - Computational
  - Drug
  - Disease
  - MicroRNA

Category	ID	Title (or Source)		<input type="checkbox"/>	liver_logP	salivary gland_logP	stomach cardiac_logP
<b>GO: Molecular Function</b>							
					<b>pValues</b>		
GO: Molecular Function	GO:0016491	oxidoreductase activity		<input type="checkbox"/>	<input type="checkbox"/>		
					10.0000		
GO: Molecular Function	GO:0005201	extracellular matrix structural constituent		<input type="checkbox"/>			<input type="checkbox"/> 5.3780
GO: Molecular Function	GO:0005506	iron ion binding		<input type="checkbox"/>	<input type="checkbox"/> 5.3161		
GO: Molecular Function	GO:0004497	monooxygenase activity		<input type="checkbox"/>	<input type="checkbox"/> 5.2535		
GO: Molecular Function	GO:0004022	alcohol dehydrogenase activity		<input type="checkbox"/>	<input type="checkbox"/> 4.8034		
GO: Molecular Function	GO:0046906	tetrapyrrole binding		<input type="checkbox"/>	<input type="checkbox"/> 4.5687		
GO: Molecular Function	GO:0020037	heme binding		<input type="checkbox"/>	<input type="checkbox"/> 4.5687		
GO: Molecular Function	GO:0016705	oxidoreductase activity, acting on paired donors, with incorporation or reduction of molecular oxygen		<input type="checkbox"/>	<input type="checkbox"/> 4.4053		
GO: Molecular Function	GO:0019825	oxygen binding		<input type="checkbox"/>	<input type="checkbox"/> 3.9166		
GO: Molecular Function	GO:0004866	endopeptidase inhibitor activity		<input type="checkbox"/>	<input type="checkbox"/> 3.9162		
GO: Molecular Function	GO:0030414	peptidase inhibitor activity		<input type="checkbox"/>	<input type="checkbox"/> 3.8272		
GO: Molecular Function	GO:0004024	alcohol dehydrogenase activity, zinc-dependent		<input type="checkbox"/>	<input type="checkbox"/> 3.5535		
GO: Molecular Function	GO:0043499	eukaryotic cell surface binding		<input type="checkbox"/>	<input type="checkbox"/> 3.4320		
GO: Molecular Function	GO:0030246	carbohydrate binding		<input type="checkbox"/>	<input type="checkbox"/> 2.0565		<input type="checkbox"/> 1.3457
GO: Molecular Function	GO:0008289	lipid binding		<input type="checkbox"/>	<input type="checkbox"/> 3.1619		
GO: Molecular Function	GO:0004857	enzyme inhibitor activity		<input type="checkbox"/>	<input type="checkbox"/> 2.8756		
GO: Molecular Function	GO:0004867	serine-type endopeptidase inhibitor activity		<input type="checkbox"/>	<input type="checkbox"/> 2.8632		
GO: Molecular Function	GO:0048407	platelet-derived growth factor binding		<input type="checkbox"/>			<input type="checkbox"/> 2.7407
GO: Molecular Function	GO:0008201	heparin binding		<input type="checkbox"/>	<input type="checkbox"/> 2.6553		

# ToppCluster (<http://toppcluster.cchmc.org>)

EHF  
COL15A1  
LOC100130100  
IGHA1  
LTF  
IGKC  
IGL@  
FAM129A  
ATP8B1  
IGLC2

Network View – Shared and specific genes and annotations between different gene lists  
Cytoscape  
(<http://cytoscape.org>)  
installation required

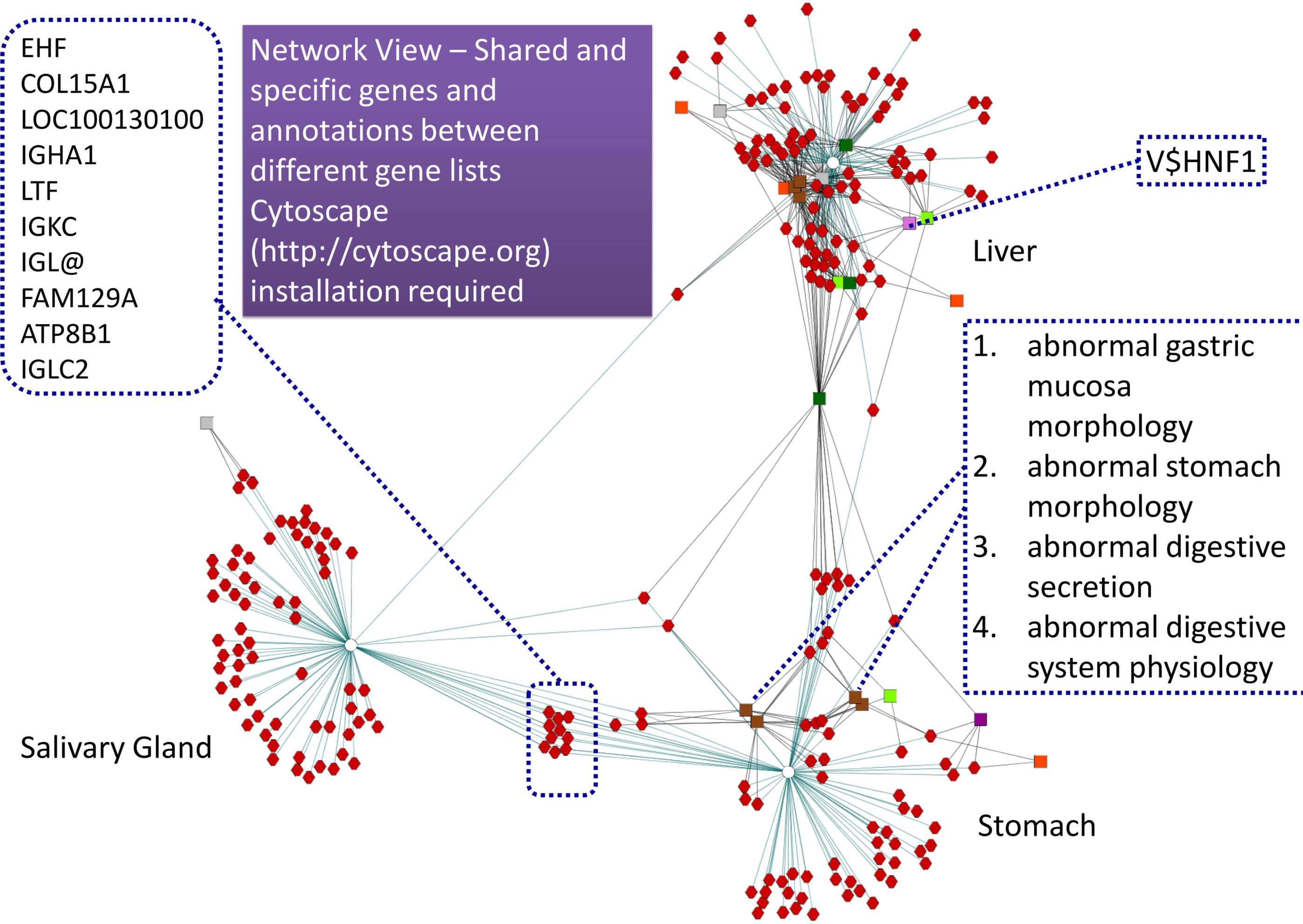
Salivary Gland

Liver

1. abnormal gastric mucosa morphology
2. abnormal stomach morphology
3. abnormal digestive secretion
4. abnormal digestive system physiology

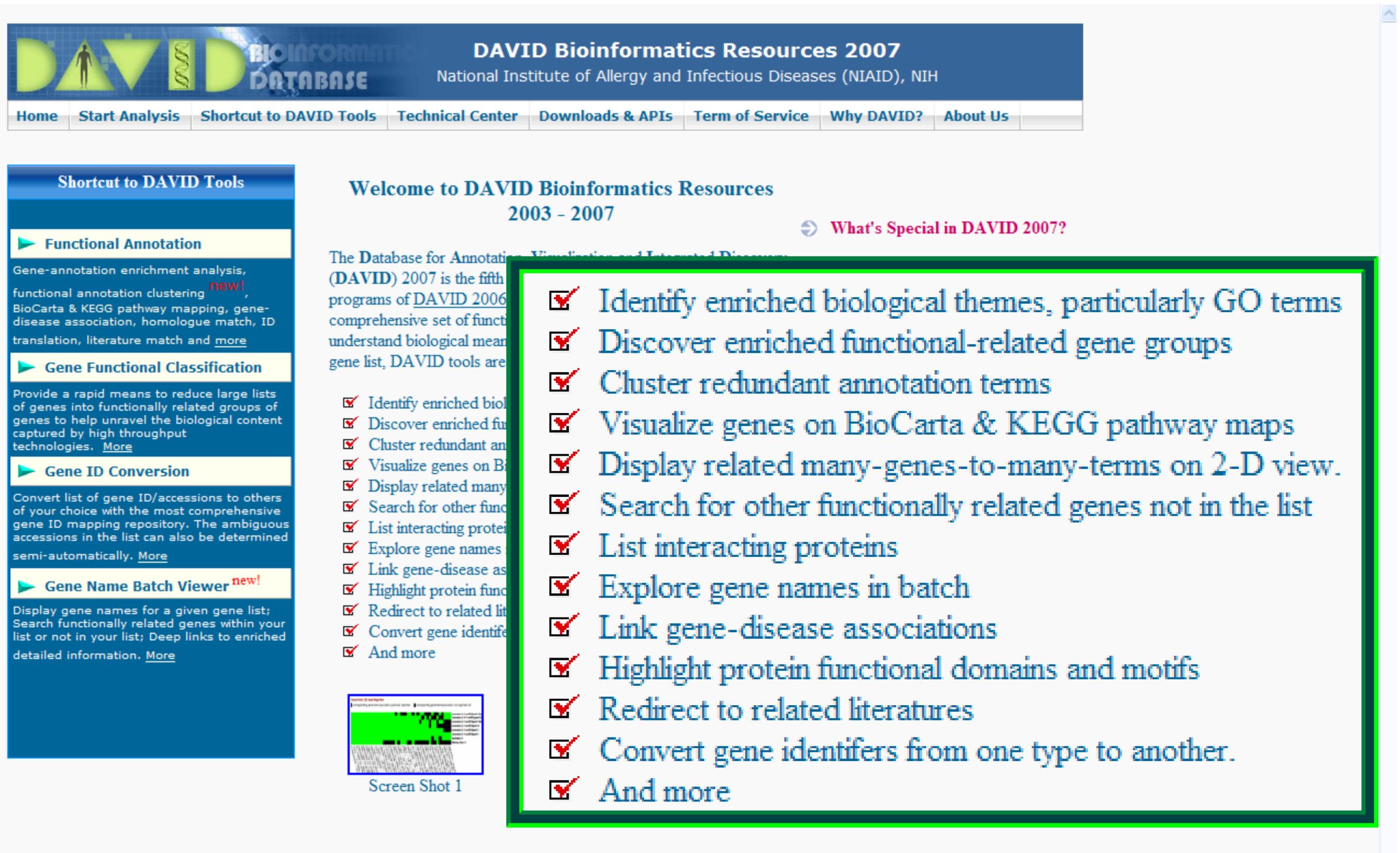
Stomach

V\$HNF1



# DAVID (<http://david.abcc.ncifcrf.gov>)

## Database for Annotation, Visualization and Integrated Discovery



**DAVID Bioinformatics Resources 2007**  
National Institute of Allergy and Infectious Diseases (NIAID), NIH

Home | Start Analysis | Shortcut to DAVID Tools | Technical Center | Downloads & APIs | Term of Service | Why DAVID? | About Us

### Shortcut to DAVID Tools

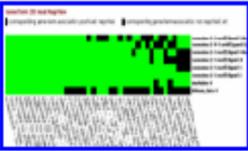
- Functional Annotation**  
Gene-annotation enrichment analysis, functional annotation clustering **new!**, BioCarta & KEGG pathway mapping, gene-disease association, homologue match, ID translation, literature match and [more](#)
- Gene Functional Classification**  
Provide a rapid means to reduce large lists of genes into functionally related groups of genes to help unravel the biological content captured by high throughput technologies. [More](#)
- Gene ID Conversion**  
Convert list of gene ID/accessions to others of your choice with the most comprehensive gene ID mapping repository. The ambiguous accessions in the list can also be determined semi-automatically. [More](#)
- Gene Name Batch Viewer **new!****  
Display gene names for a given gene list; Search functionally related genes within your list or not in your list; Deep links to enriched detailed information. [More](#)

### Welcome to DAVID Bioinformatics Resources 2003 - 2007

[What's Special in DAVID 2007?](#)

The Database for Annotation, Visualization and Integrated Discovery (DAVID) 2007 is the fifth comprehensive set of functional annotation tools to help understand biological meaning of a gene list, DAVID tools are:

- Identify enriched biological themes, particularly GO terms
- Discover enriched functional-related gene groups
- Cluster redundant annotation terms
- Visualize genes on BioCarta & KEGG pathway maps
- Display related many-genes-to-many-terms on 2-D view.
- Search for other functionally related genes not in the list
- List interacting proteins
- Explore gene names in batch
- Link gene-disease associations
- Highlight protein functional domains and motifs
- Redirect to related literatures
- Convert gene identifiers from one type to another.
- And more



Screen Shot 1



# DAVID (<http://david.abcc.ncifcrf.gov>)

 **DAVID Bioinformatics Resources 2007**  
National Institute of Allergy and Infectious Diseases (NIAID), NIH

[Home](#) | [Start Analysis](#) | [Shortcut to DAVID Tools](#) | [Technical Center](#) | [Downloads & APIs](#) | [Term of Service](#) | [Why DAVID?](#) | [About Us](#)

### Shortcut to DAVID Tools

- Functional Annotation**  
Gene-annotation enrichment analysis, functional annotation clustering **new!**, BioCarta & KEGG pathway mapping, gene-disease association, homologue match, ID translation, literature match and [more](#)
- Gene Functional Classification**  
Provide a rapid means to reduce large lists of genes into functionally related groups of genes to help unravel the biological content captured by high throughput technologies. [More](#)
- Gene ID Conversion**  
Convert list of gene ID/accessions to others of your choice with the most comprehensive gene ID mapping repository. The ambiguous accessions in the list can also be determined semi-automatically. [More](#)
- Gene Name Batch Viewer **new!****  
Display gene names for a given gene list; Search functionally related genes within your list or not in your list; Deep links to enriched detailed information. [More](#)

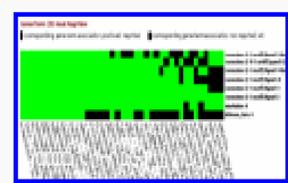
## Welcome to DAVID Bioinformatics Resources 2003 - 2007

[What's Special in DAVID 2007?](#)

The Database for Annotation, Visualization and Integrated Discovery (DAVID) 2007 is the fifth programs of [DAVID 2006](#) comprehensive set of functions to understand biological meaning of a gene list, DAVID tools are

- ✓ Identify enriched biological themes, particularly GO terms
- ✓ Discover enriched functional-related gene groups
- ✓ Cluster redundant annotation terms
- ✓ Visualize genes on BioCarta & KEGG pathway maps
- ✓ Display related many-genes-to-many-terms on 2-D view.
- ✓ Search for other functionally related genes not in the list
- ✓ List interacting proteins
- ✓ Explore gene names in batch
- ✓ Link gene-disease associations
- ✓ Highlight protein functional domains and motifs
- ✓ Redirect to related literatures
- ✓ Convert gene identifiers from one type to another.
- ✓ And more

- ✓ Identify enriched biological themes, particularly GO terms
- ✓ Discover enriched functional-related gene groups
- ✓ Cluster redundant annotation terms
- ✓ Visualize genes on BioCarta & KEGG pathway maps
- ✓ Display related many-genes-to-many-terms on 2-D view.
- ✓ Search for other functionally related genes not in the list
- ✓ List interacting proteins
- ✓ Explore gene names in batch
- ✓ Link gene-disease associations
- ✓ Highlight protein functional domains and motifs
- ✓ Redirect to related literatures
- ✓ Convert gene identifiers from one type to another.
- ✓ And more



Screen Shot 1

# DAVID (<http://david.abcc.ncifcrf.gov>)

## Convert NCBI Entrez Gene IDs to RefSeq Accession Numbers

**Gene ID Conversion Tool**

Submit your gene list to start conversion!

**The Cross-Conversion of Gene ID Types:**

- Entrez Gene ID
- Affy ID
- GenBank Accession
- Genpept Accession
- NCBI GI
- PIR Accession
- PIR ID
- PIR NREF ID
- RefSeq Genomic Accession
- RefSeq mRNA Accession
- RefSeq Protein Accession
- RefSeq RNA Accession
- Unigene
- UNIPROT Accession
- UNIPROT ID
- UNIREF100 ID
- Official Gene Symbol *new!*
- Not Sure *new!*

**Upload** | **List** | Background

Upload Gene List

[Demolist 1](#) [Demolist 2](#)

[Upload Help](#)

Step 1: Enter Gene List

A: Paste a list

Clear

Or

B: Choose From a File

Browse...

Step 2: Select Identifier

AFFY\_ID

Step 3: List Type

Gene List

Background

Step 4: Submit List

Submit List

**Upload** | **List** | Background

Upload Gene List

[Demolist 1](#) [Demolist 2](#)

[Upload Help](#)

Step 1: Enter Gene List

A: Paste a list

3493

3512

10562

3535

5284

Clear

Or

B: Choose From a File

Browse...

Step 2: Select Identifier

ENTREZ\_GENE\_ID

Step 3: List Type

Gene List

Background

Step 4: Submit List

Submit List

[Tell us how you like the tool](#)

[Technical notes of the tool](#)

[Contact us for questions](#)

# DAVID (<http://david.abcc.ncifcrf.gov>)

**Upload** **List** **Background**

Gene List Manager

Select to limit annotations by one or more species [Help](#)

- Use All Species -  
 HOMO SAPIENS(68)  
 UNKNOWN(4)  
 UNIDENTIFIED(1)

Select

List Manager [Help](#)

Uploaded List\_1

Select List to:

Use Rename  
 Remove Combine

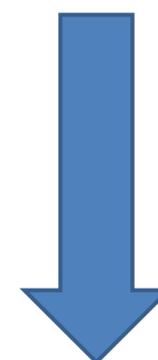
Show Gene List <sup>new!</sup>

[View Unmapped Ids](#)



Convert the gene list being selected in left panel to

Submit



## Gene ID Conversion Tool Result

[Right-click to Download the result](#) [Help](#)

[Submit Converted List to DAVID as a Gene List](#) [Submit Converted List to DAVID as a Background](#)

Gene Accession Conversion Statistics [Help](#)

Conversion Summary		
ID Count	In DAVID DB	Conversion
<a href="#">60</a>	Yes	Successful
<a href="#">8</a> IDs	Yes	None
<a href="#">0</a> IDs	No	None
<a href="#">0</a> IDs	Ambiguous	Pending
<b>Total Unique User IDs: 68</b>		
Summary of Ambiguous Gene IDs		
ID Count	Possible Source	Convert All
All Possible Sources For Ambiguous IDs		
Ambiguous ID	Possibility	Convert

From	To	Species	David Gene Name
202	<a href="#">NM_001624</a>	HOMO SAPIENS	ABSENT IN MELANOMA 1
72	<a href="#">NM_001613</a>	HOMO SAPIENS	ACTIN, ALPHA 2, SMOOTH MUSCLE, AORTA
72	<a href="#">NM_001615</a>	HOMO SAPIENS	ACTIN, ALPHA 2, SMOOTH MUSCLE, AORTA
27299	<a href="#">NM_014479</a>	HOMO SAPIENS	ADAM-LIKE, DECYSIN 1
125	<a href="#">NM_000667</a>	HOMO SAPIENS	ALCOHOL DEHYDROGENASE 1A (CLASS I), ALPHA POLYPEPTIDE
125	<a href="#">NM_000668</a>	HOMO SAPIENS	ALCOHOL DEHYDROGENASE 1A (CLASS I), ALPHA POLYPEPTIDE
126	<a href="#">NM_000669</a>	HOMO SAPIENS	ALCOHOL DEHYDROGENASE 1A (CLASS I), ALPHA POLYPEPTIDE
126	<a href="#">NM_000668</a>	HOMO SAPIENS	ALCOHOL DEHYDROGENASE 1A (CLASS I), ALPHA POLYPEPTIDE
125	<a href="#">NM_000669</a>	HOMO SAPIENS	ALCOHOL DEHYDROGENASE 1A (CLASS I), ALPHA POLYPEPTIDE

**Exercise 13: Convert affymetrix probeset IDs to gene symbols**

**Exercise 14: What are the enriched pathways and diseases for this gene set?**

*From the same example data set (“Example-Set-1.xls”), use the probe set IDs (2<sup>nd</sup> column) and extract their RefSeq accession numbers*

# PANTHER (<http://www.pantherdb.org/>)

## Protein Analysis Through Evolutionary Relationships



- Quick links**
- [Browse PANTHER](#)
  - [Search PANTHER](#)
  - [Batch search](#)
  - [Browse pathways](#)
  - [Community Curation](#)
  - [My Workspace](#)
  - [Gene expression tools](#)
  - [HMM scoring](#)
  - [cSNP analysis](#)
  - [Downloads](#)
  - [Site map](#)

Find PANTHER-classified genes, transcripts, and proteins by uploading a list of IDs

### Batch ID Search

Enter IDs:

separate IDs by a space or comma - [supported IDs](#)

Upload IDs:

- [file format](#)

Select upload ID type:

Select File Type:  ID List  Previously exported text search results

Result page:  Genes  Transcripts/Proteins

Select datasets:

Celera:	<input type="checkbox"/> H. sapiens	<input type="checkbox"/> M. musculus	<input type="checkbox"/> R. norvegicus
NCBI:	<input checked="" type="checkbox"/> H. sapiens	<input type="checkbox"/> M. musculus	<input type="checkbox"/> R. norvegicus
FlyBase:	<input type="checkbox"/> D. melanogaster		

### GENE EXPRESSION DATA ANALYSIS

Our expression analysis tools can be used for microarray data interpretation. Multiple gene lists can be mapped to PANTHER molecular function and biological process categories, as well as to biological pathways. Our pathway visualization tool will display your experimental results on detailed diagrams of the relationships between genes/proteins in known pathways.

➤ [Compare gene lists](#)   
Upload lists of genes or gene products and statistically compare them to a reference list to look for under- and over-represented functional categories.

➤ [Analyze a list of genes with expression values](#)   
Upload a list of genes and their corresponding fold-change values from a differential expression experiment.

**You can compare *multiple* lists!**

## Compare Classifications of Lists ?

Map lists of genes to a PANTHER ontology. For pathways, you can then view the gene expression values overlaid on top of a pathway diagram, where genes will be colored differently for different clusters of genes.

Use the binomial statistics tool to compare classifications of multiple clusters of lists to a reference list to statistically determine over- or under- representation of PANTHER classification categories. Each list is compared to the reference list using the binomial test (Cho & Campbell, TIGs 2000) for each molecular function, biological process, or pathway term in PANTHER.

### Steps:

1. Select list(s) to analyze
2. Select reference list

### 1. Select Lists to Compare to a Reference List

For example, each selected list may be a cluster of co-expressed genes under a particular set of conditions.

Select list(s)

selected: FetalLiverSpecific.txt  
FetalBrainSpecific.txt  
AdultHeartSpecific.txt

### 2. Select Reference List

For example, the reference list may be the set of all genes in the experiment, or the set of all genes in the genome being analyzed.

Select reference list

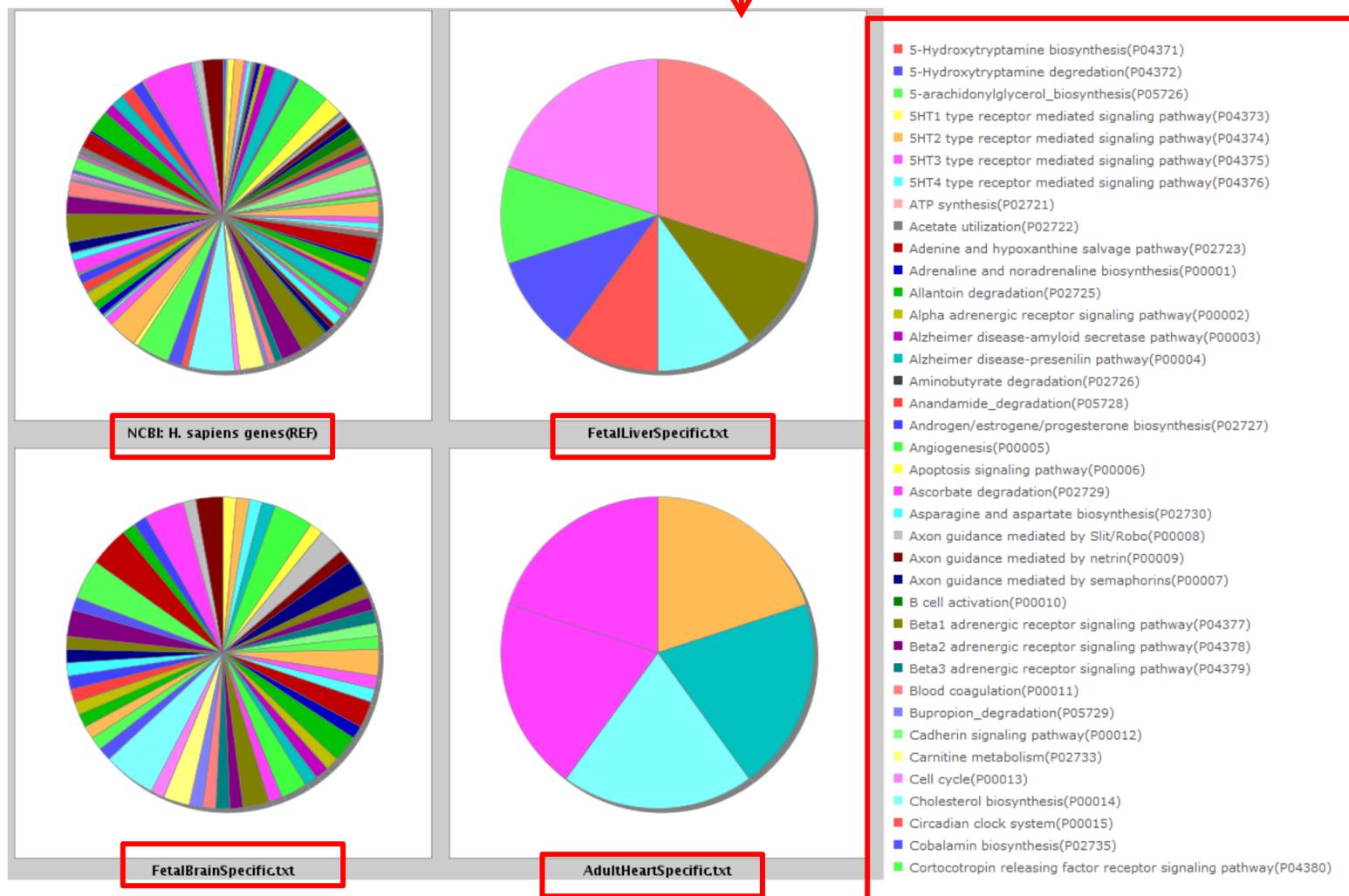
default: NCBI: H. sapiens genes

### Search options

PANTHER Ontology:

- Pathways
- Biological Process
- Molecular Function

Use the Bonferroni correction for multiple testing ?



# PANTHER (<http://www.pantherdb.org/>)

## Protein Analysis Through Evolutionary Relationships

### Results ?

Colors for viewing genes in pathway diagrams:

Example-Set-3.txt:  ▼

*gray: components only in the reference list*

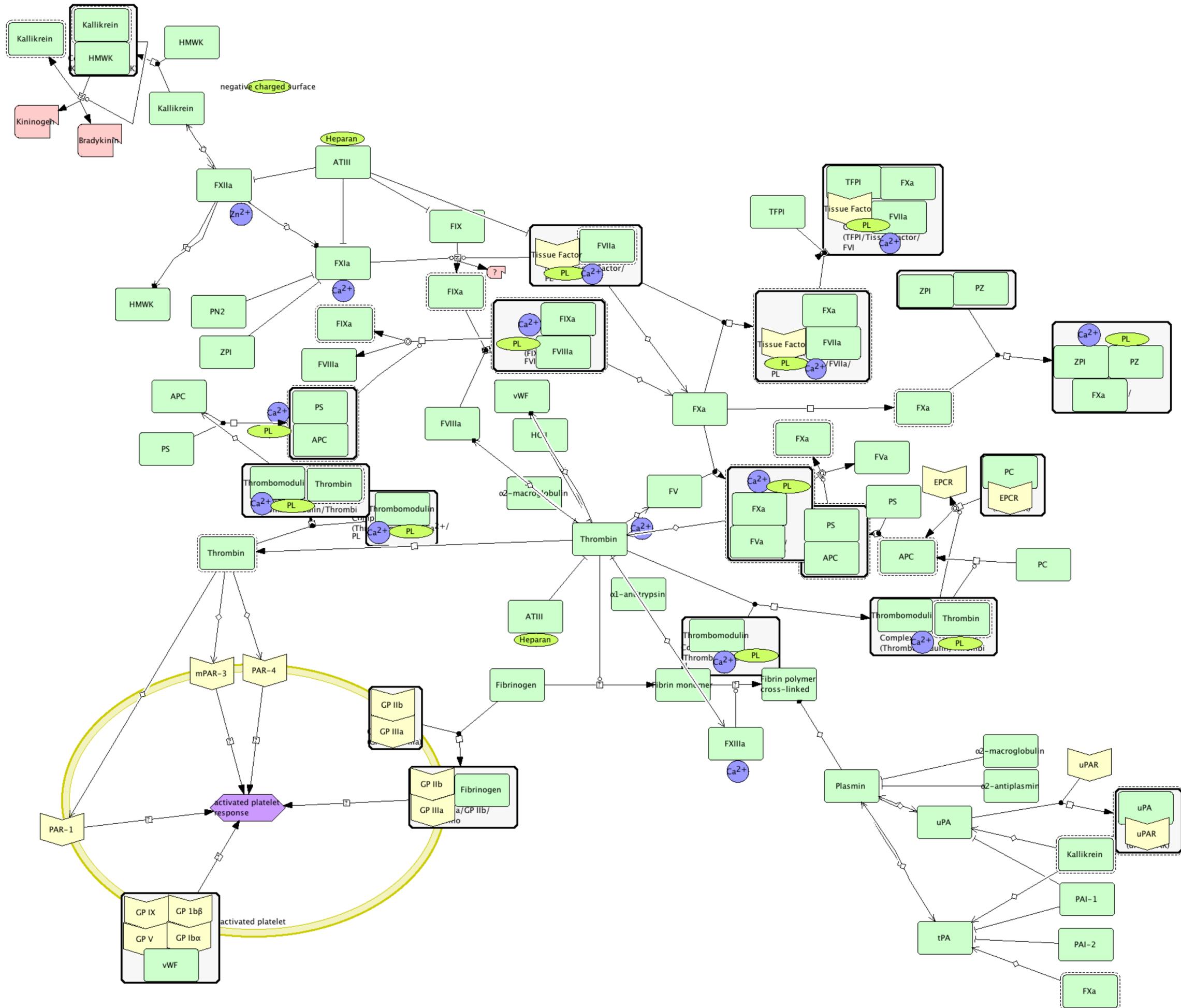
	Reference list	Example-Set-3.txt
Mapped IDs:	<a href="#">25431</a>	<a href="#">135</a>
Unmapped IDs:	<a href="#">0</a>	<a href="#">15</a>

Click on pathway name to see genes highlighted on pathway diagram

**Export results** View:  ▼

<a href="#">Pathways</a>	NCBI: H. sapiens genes (REF)	Example-Set-3.txt			
	#	#	expected	+/-	▲ P value
<a href="#">Blood coagulation</a>	<a href="#">55</a>	<a href="#">8</a>	.29	+	1.38E-07
<a href="#">Plasminogen activating cascade</a>	<a href="#">21</a>	<a href="#">4</a>	.11	+	9.37E-04

# PANTHER (<http://www.pantherdb.org/>)



# Summary

## Cis-Element Finding Matrix

	CONSERVED	NON-CONSERVED
KNOWN TFBS	oPOSSUM DiRE	Pscan MatInspector*
NOVEL/UNKNOWN TFBS OR MOTIFS	oPOSSUM WEEDER-H	MEME WEEDER

# RESOURCES - URLs: Summary

Application/Resource	URL
oPOSSUM	<a href="http://burgundy.cmmmt.ubc.ca/oPOSSUM/">http://burgundy.cmmmt.ubc.ca/oPOSSUM/</a>
DiRE	<a href="http://dire.dcode.org/">http://dire.dcode.org/</a>
Weeder-H	<a href="http://159.149.109.9/modtools/">http://159.149.109.9/modtools/</a>
Weeder	<a href="http://159.149.109.9/modtools/">http://159.149.109.9/modtools/</a>
Pscan	<a href="http://159.149.109.9/modtools/">http://159.149.109.9/modtools/</a>
MEME	<a href="http://meme.sdsc.edu/">http://meme.sdsc.edu/</a>
MatInspector	<a href="http://www.genomatix.de/">http://www.genomatix.de/</a>
GenomeTrafac	<a href="http://genometrafac.cchmc.org">http://genometrafac.cchmc.org</a>
ToppGene	<a href="http://toppgene.cchmc.org">http://toppgene.cchmc.org</a>
ToppCluster	<a href="http://toppcluster.cchmc.org">http://toppcluster.cchmc.org</a>
DAVID	<a href="http://david.abcc.ncifcrf.gov">http://david.abcc.ncifcrf.gov</a>
PANTHER	<a href="http://www.pantherdb.org">http://www.pantherdb.org</a>
Genome Browser	<a href="http://genome.ucsc.edu">http://genome.ucsc.edu</a>
ECR Browser	<a href="http://ecrbrowser.dcode.org">http://ecrbrowser.dcode.org</a>
Slides/Exercises	<a href="http://anil.cchmc.org/dhc.html">http://anil.cchmc.org/dhc.html</a>

# Exercises - Summary

- Exercise 1:** Use oPOSSUM to find shared conserved cis-elements in a group of co-expressed genes
- Exercise 2:** Use DiRE to find shared conserved cis-elements in a group of co-expressed genes
- Exercise 3:** Use Pscan to find shared cis-elements (Transfac) in a group of co-expressed genes
- Exercise 4:** Download upstream 500 bp sequence for a list of genes
- Exercise 5:** Download all SNPs overlapping with these genes
- Exercise 6:** Download the orthologous promoter sequences (human, mouse, and rat) for the gene SLC7A1
- Exercise 7:** Are there any putative microRNA regulators for SLC7A1? If yes, download all of them using table browser
- Exercise 8:** Use the downloaded SLC7A1 ortholog promoter sequences to find out common motifs using WeederH
- Exercise 9:** Use the downloaded promoter sequences to find out common motifs using Weeder and MEME
- Exercise 10:** Does any of the motifs found by Meme match known TFBS?
- Exercise 11:** Use the gene list from the downloaded file (“Example-Set-2”) and find out:
  - How many of these genes are transcription factors?
  - What are the enriched TFBSs and miRNAs?
  - What gene families are enriched in this list?
  - Are there any salivary gland development associated genes present in this list?
  - How many and which genes from this list are associated with non-insulin dependent diabetes mellitus (NIDDM)?
- Exercise 12:** Prioritize the 721 genes (“Example-Set-2”) using “stomach genes” from the “Example-Set-1”.
  - What are the top 10 ranked genes using ToppGene and ToppNet?
  - Why is TFF3 ranked among the top 5 in ToppGene prioritization? What is its rank in ToppNet?
- Exercise 13:** Convert Affymetrix probeset IDs to gene symbols
- Exercise 14:** What are the enriched pathways and diseases for this gene set?

**For additional exercises, see <http://anil.cchmc.org/dhc.html>**