

# ISG15 and STAT1 are promising druggable targets for treating Hepatitis C that control activity of FOXO1, NFKB1 and MAX transcription factors on promoters of differentially expressed genes in liver tissue

Demo User

geneXplain GmbH

info@genexplain.com

Data received on 13/08/2019 ; Run on 27/11/2024 ; Report generated on 27/11/2024

Genome Enhancer release 3.5 (TRANSFAC®, TRANSPATH® and HumanPSD™ release 2024.2)

---



## Abstract

In the present study we applied the software package "Genome Enhancer" to a data set that contains *transcriptomics* data obtained from *liver* tissue. The study is done in the context of *Hepatitis C*. The goal of this pipeline is to identify potential drug targets in the molecular network that governs the studied pathological process. In the first step of analysis pipeline discovers transcription factors (TFs) that regulate genes activities in the pathological state. The activities of these TFs are controlled by so-called master regulators, which are identified in the second step of analysis. After a subsequent druggability checkup, the most promising master regulators are chosen as potential drug targets for the analyzed pathology. At the end the pipeline comes up with (a) a list of known drugs and (b) investigational active chemical compounds with the potential to interact with selected drug targets.

From the data set analyzed in this study, we found the following TFs to be potentially involved in the regulation of the differentially expressed genes: FOXO1, NFKB1, TFCEP2, MAX, HMGA1 and ERF. The subsequent network analysis suggested

- isg15(h):UbcH8(h):ISG15 E3 ligases
- isg15:UbcH8:ISG15 E3 ligases
- ISGylated host proteins
- STAT1
- UBP43

as the most promising molecular targets for further research, drug development and drug repurposing initiatives on the basis of identified molecular mechanism of the studied pathology. Having checked the actual druggability potential of the full list of identified targets, both, via information available in medical literature and via cheminformatics analysis of drug compounds, we have identified the following drugs as the most promising treatment candidates for the studied pathology: Sorafenib, seliciclib and Bortezomib.

# 1. Introduction

Recording "-omics" data to measure gene activities, protein expression or metabolic events is becoming a standard approach to characterize the pathological state of an affected organism or tissue. Increasingly, several of these methods are applied in a combined approach leading to large "multiomics" datasets. Still the challenge remains how to reveal the underlying molecular mechanisms that render a given pathological state different from the norm. The disease-causing mechanism can be described by a re-wiring of the cellular regulatory network, for instance as a result of a genetic or epigenetic alterations influencing the activity of relevant genes. Reconstruction of the disease-specific regulatory networks can help identify potential master regulators of the respective pathological process. Knowledge about these master regulators can point to ways how to block a pathological regulatory cascade. Suppression of certain molecular targets as components of these cascades may stop the pathological process and cure the disease.

Conventional approaches of statistical "-omics" data analysis provide only very limited information about the causes of the observed phenomena and therefore contribute little to the understanding of the pathological molecular mechanism. In contrast, the "upstream analysis" method [1-4] applied here has been devised to provide a casual interpretation of the data obtained for a pathology state. This approach comprises two major steps: (1) analysing promoters and enhancers of differentially expressed genes for the transcription factors (TFs) involved in their regulation and, thus, important for the process under study; (2) re-constructing the signaling pathways that activate these TFs and identifying master regulators at the top of such pathways. For the first step, the database TRANSFAC® [6] is employed together with the TF binding site identification algorithms Match [7] and CMA [8]. The second step involves the signal transduction database TRANSPATH® [9] and special graph search algorithms [10-11] implemented in the software "Genome Enhancer".

The "upstream analysis" approach has now been extended by a third step that reveals known drugs suitable to inhibit (or activate) the identified molecular targets in the context of the disease under study. This step is performed by using information from HumanPSD™ database [5]. In addition, some known drugs and investigational active chemical compounds are subsequently predicted as potential ligands for the revealed molecular targets. They are predicted using a pre-computed database of spectra of biological activities of chemical compounds of a library of 2245 known drugs and investigational chemical compounds from HumanPSD™ database. The spectra of biological activities for these compounds are computed using the program PASS on the basis of a (Q)SAR approach [12-14]. These predictions can be used for the research purposes - for further drug development and drug repurposing initiatives.

## 2. Data

For this study the following experimental data was used:

Table 1. Experimental datasets used in the study

File name	Data type
E01_Transcriptomics_LogFC-Table	Transcriptomics

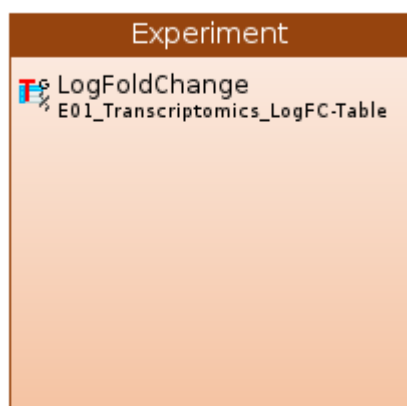


Figure 1. Annotation diagram of experimental data used in this study. With the colored boxes we show those sub-categories of the data that are compared in our analysis.

### 3. Results

We have analyzed the following condition: Experiment.

#### 3.1. Identification of target genes

In the first step of the analysis **target genes** were identified from the uploaded experimental data. Genes were ranked according to the expression value and 300 genes with highest value (see Table 2) and 300 genes with lowest value (see Table 3) were selected for further analysis.

Table 2. Top ten high expressed genes in Experiment.

[See full table](#) →

ID	Gene description	Gene symbol	LogFoldChange
ENSG00000137959	interferon induced protein 44 like	IFI44L	6.19
ENSG00000169245	C-X-C motif chemokine ligand 10	CXCL10	6.02
ENSG00000134321	radical S-adenosyl methionine domain containing 2	RSAD2	5.97
ENSG00000137965	interferon induced protein 44	IFI44	3.78
ENSG00000133106	epithelial stromal interaction 1	EPSTI1	3.77
ENSG00000185745	interferon induced protein with tetratricopeptide repeats 1	IFIT1	3.71
ENSG00000187608	ISG15 ubiquitin like modifier	ISG15	3.63
ENSG00000185201	interferon induced transmembrane protein 2	IFITM2	3.54
ENSG00000185885	interferon induced transmembrane protein 1	IFITM1	3.54
ENSG00000135114	2'-5'-oligoadenylate synthetase like	OASL	3.48

Table 3. Top ten low expressed genes in Experiment.

[See full table](#) →

ID	Gene description	Gene symbol	LogFoldChange
ENSG00000167910	cytochrome P450 family 7 subfamily A member 1	CYP7A1	-1.09
ENSG00000169282	potassium voltage-gated channel subfamily A regulatory beta subunit 1	KCNAB1	-1.04
ENSG00000171560	fibrinogen alpha chain	FGA	-0.98
ENSG00000152133	G-patch domain containing 11	GPATCH11	-0.96
ENSG00000182372	CLN8 transmembrane ER and ERGIC protein	CLN8	-0.91
ENSG00000130649	cytochrome P450 family 2 subfamily E member 1	CYP2E1	-0.88
ENSG00000253327	RAD21 antisense RNA 1	RAD21-AS1	-0.88
ENSG00000170323	fatty acid binding protein 4	FABP4	-0.87
ENSG00000175390	eukaryotic translation initiation factor 3 subunit F	EIF3F	-0.86
ENSG00000261609	gigaxonin	GAN	-0.8

#### 3.2. Functional classification of genes

A functional analysis of differentially expressed genes was done by mapping the top high expressed and top low expressed genes to several known ontologies, such as Gene Ontology (GO), disease ontology (based on HumanPSD™ database) and the ontology of signal transduction and metabolic pathways from the [TRANSPATH®](#) database. Statistical significance was computed using a binomial test.

Figures 2-7 show the most significant categories.



# TRANSPATH® Pathways (2024.2)

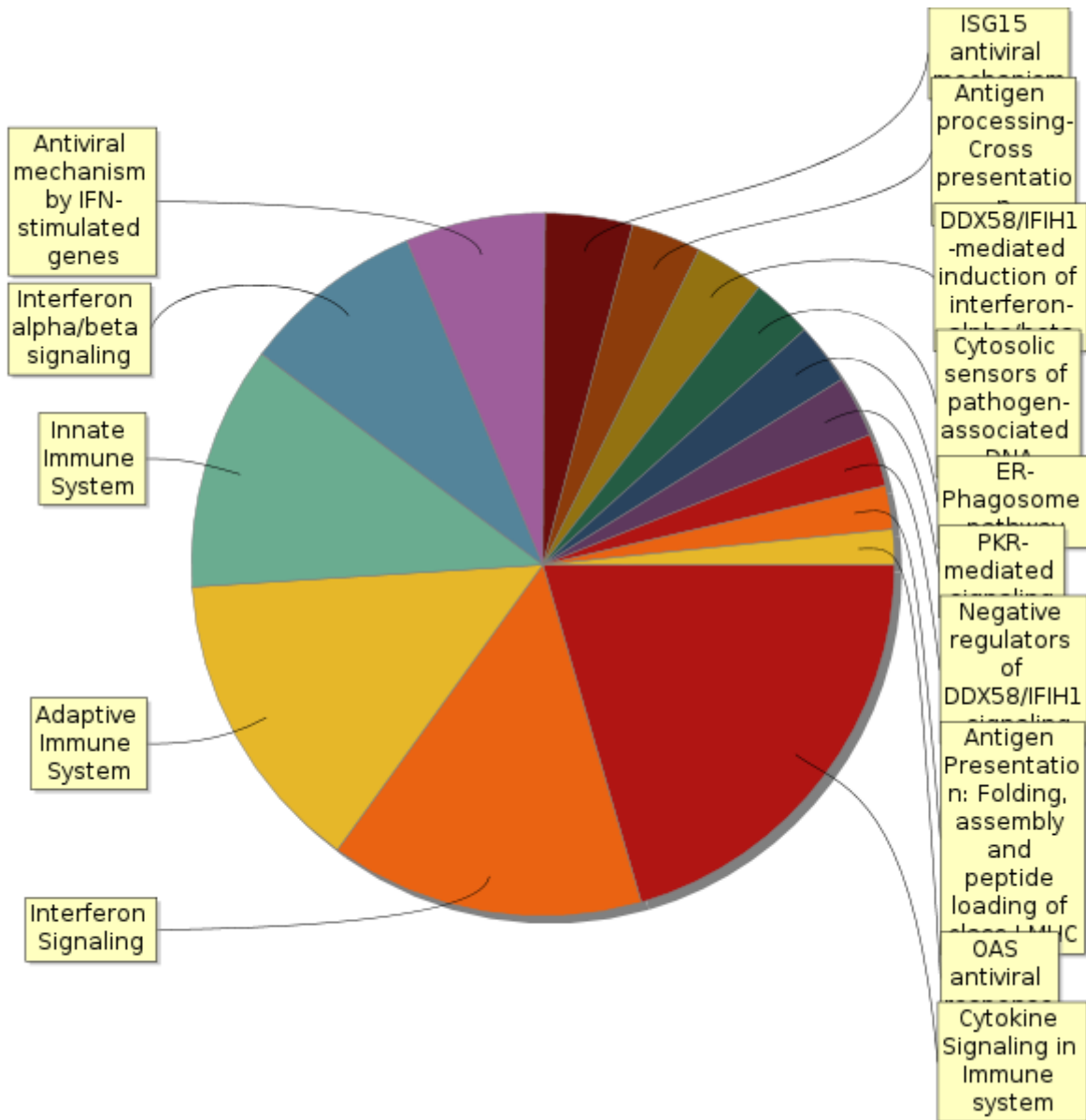


Figure 3. Enriched TRANSPATH® Pathways (2024.2) of high expressed genes in Experiment. [Full classification](#) →

## HumanPSD(TM) disease (2024.2)

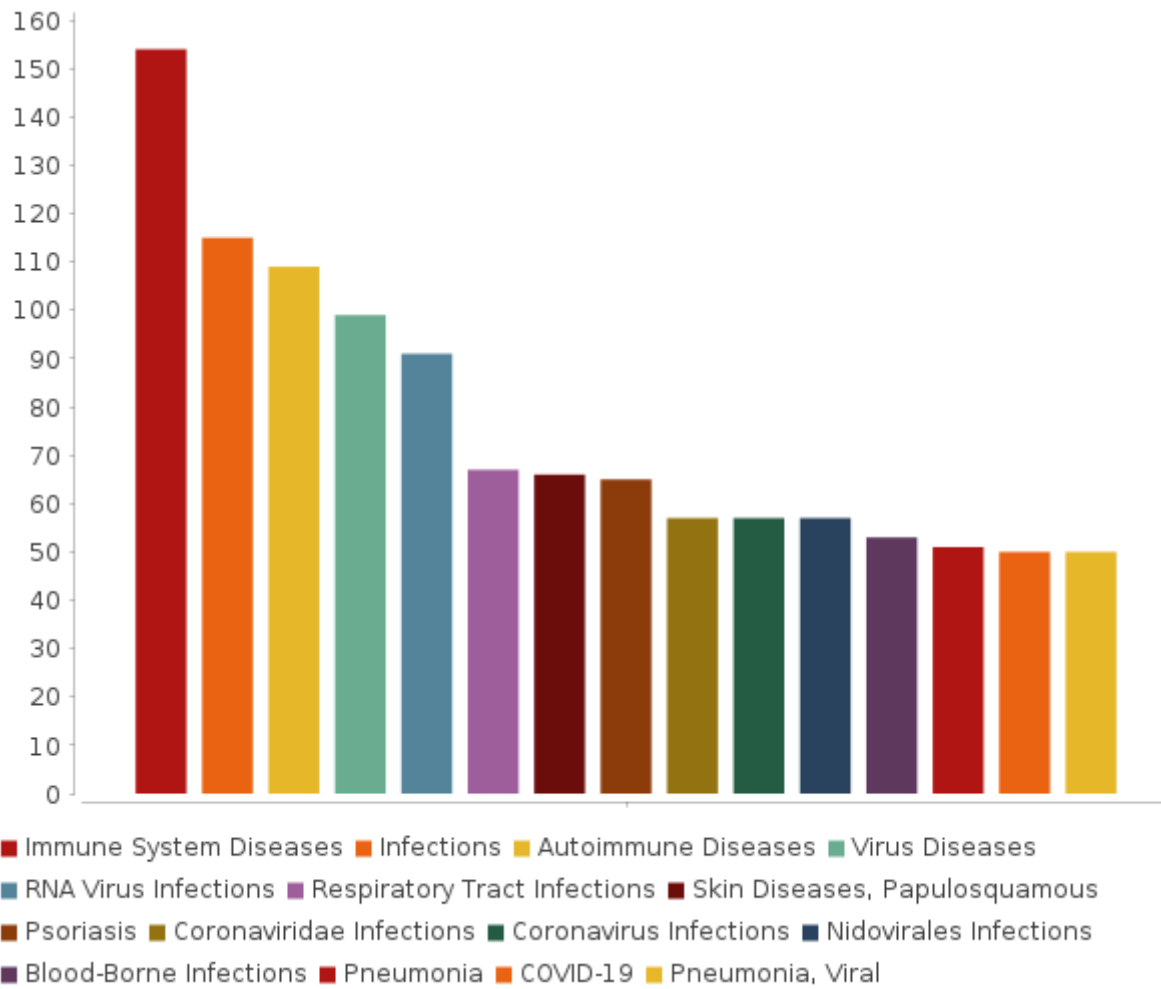


Figure 4. Enriched HumanPSD(TM) disease (2024.2) of high expressed genes in Experiment. The size of the bars correspond to the number of biomarkers of the given disease found among the input set.

[Full classification](#) →

## Low expressed genes in Experiment:

300 top low expressed genes were taken for the mapping.

# GO (biological process)

biological\_process Gene Ontology treemap

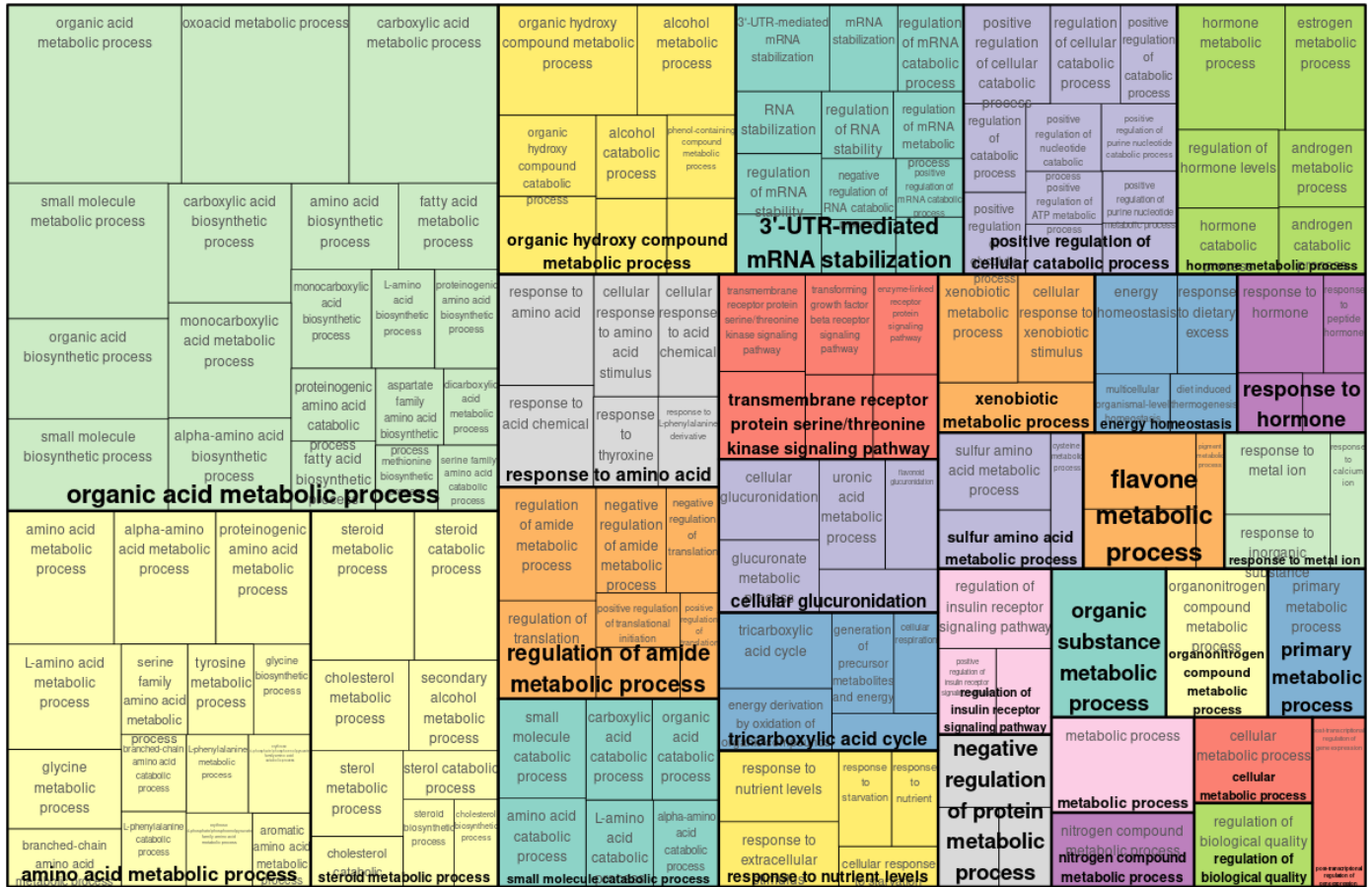


Figure 5. Enriched GO (biological process) of low expressed genes in Experiment.

Full classification →

## TRANSPATH® Pathways (2024.2)

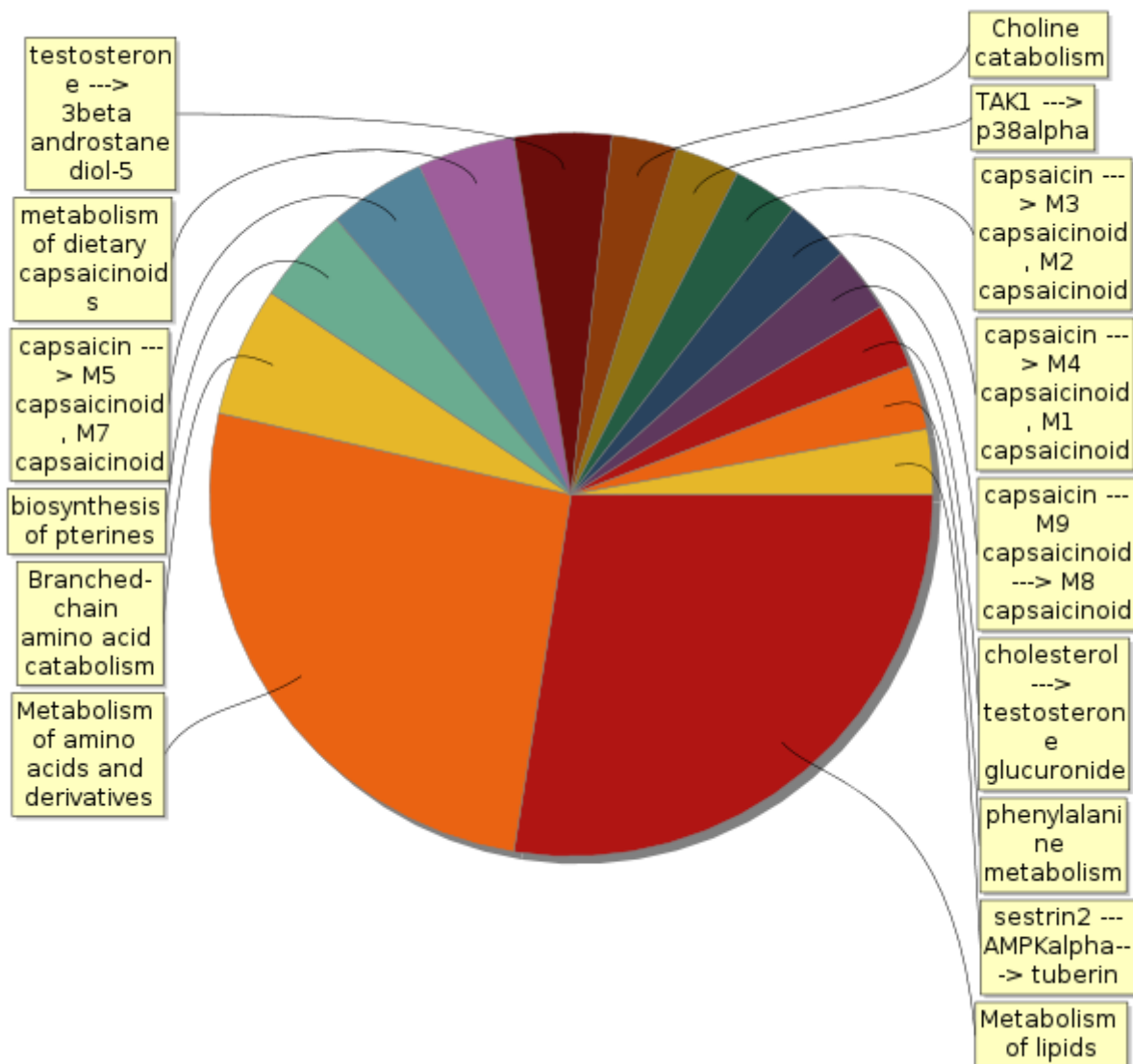


Figure 6. Enriched TRANSPATH® Pathways (2024.2) of low expressed genes in Experiment.

[Full classification](#) →

## HumanPSD(TM) disease (2024.2)

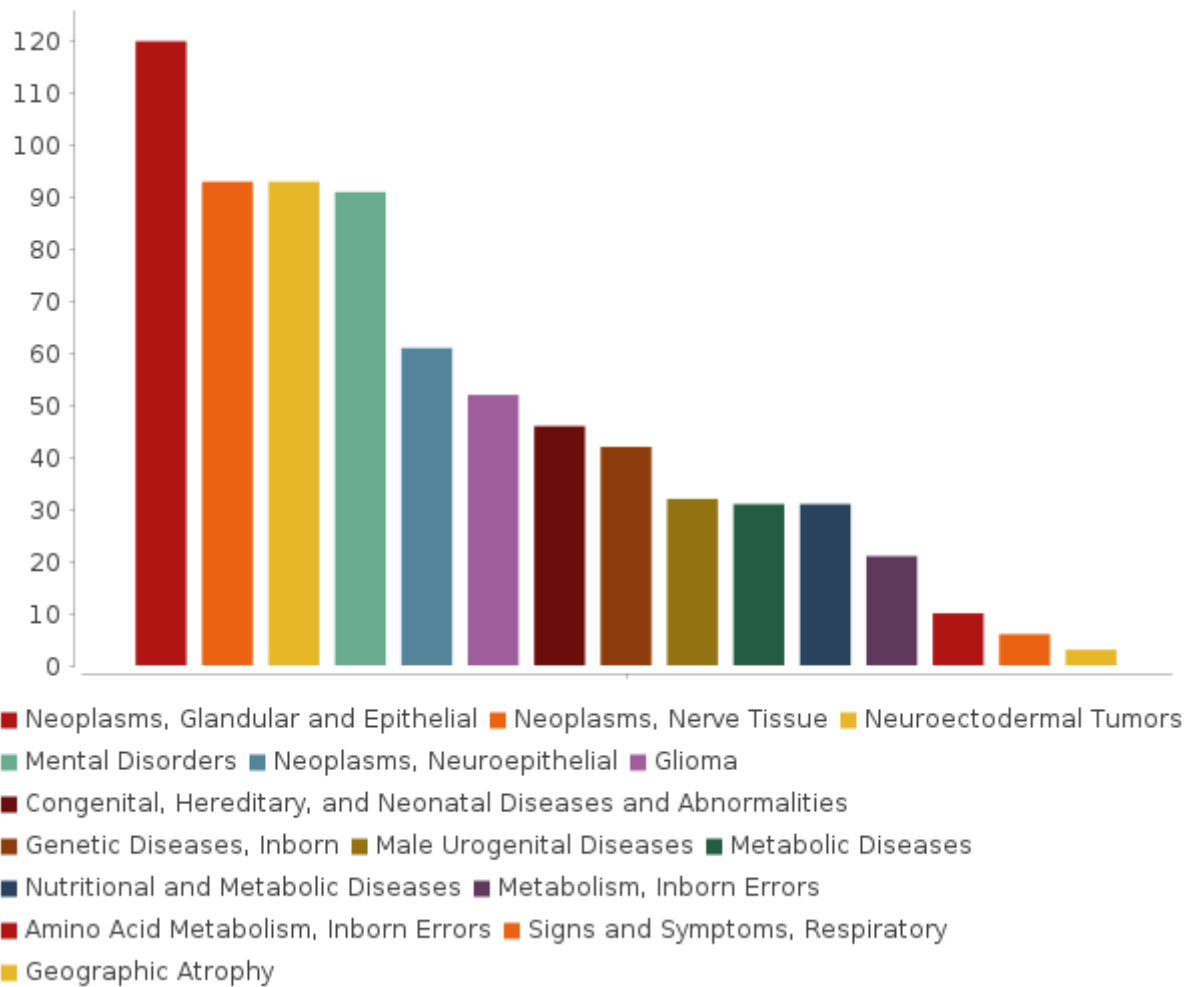
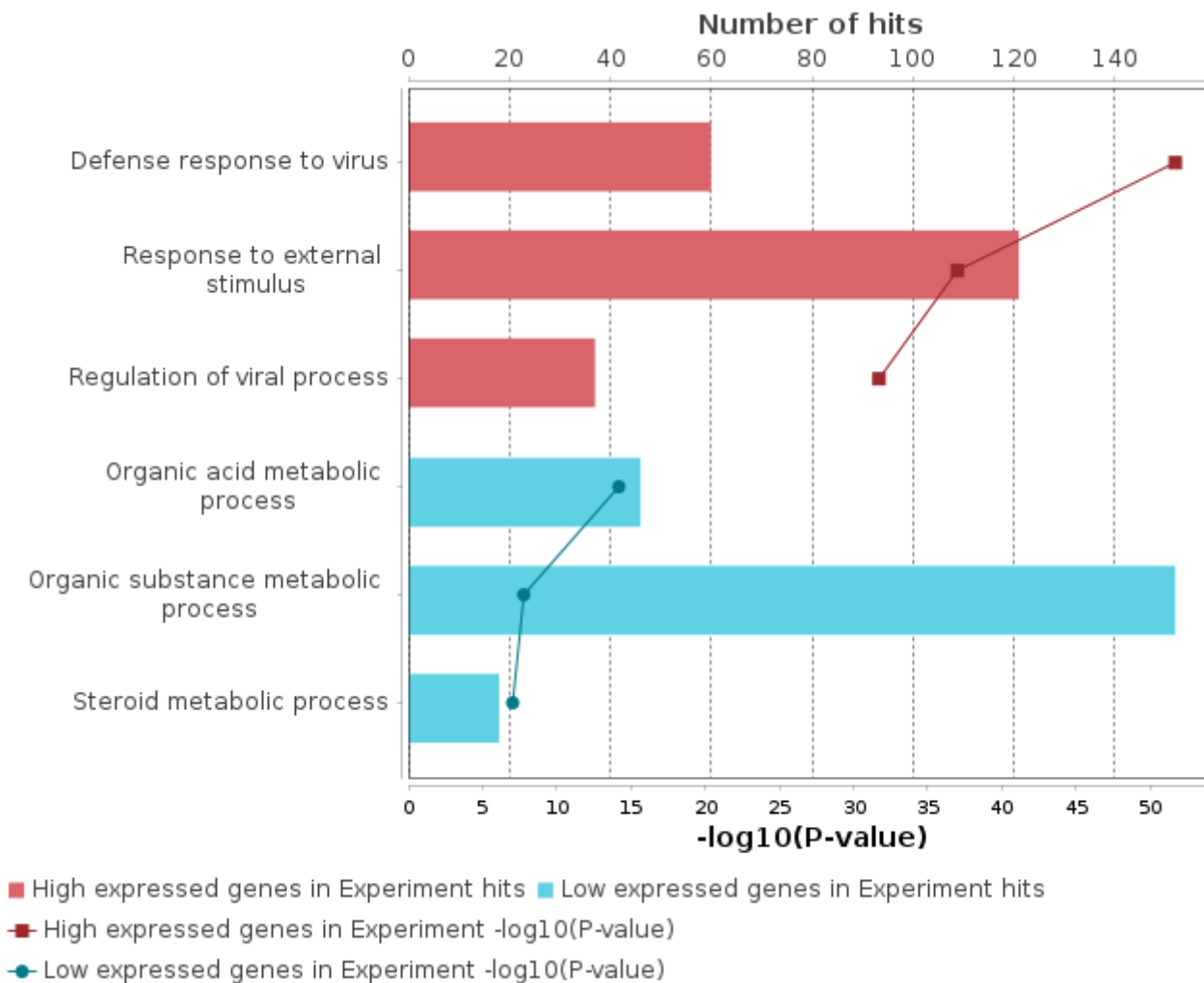


Figure 7. Enriched HumanPSD(TM) disease (2024.2) of low expressed genes in Experiment. The size of the bars correspond to the number of biomarkers of the given disease found among the input set.

[Full classification](#) →

The result of overall Gene Ontology (GO) analysis of the differentially expressed genes of the studied pathology can be summarized by the following diagram, revealing the most significant functional categories overrepresented among the observed (differentially expressed genes):



### **3.3. Analysis of enriched transcription factor binding sites and composite modules**

In the next step a search for transcription factors binding sites (TFBS) was performed in the regulatory regions of the **target genes** by using the TF binding motif library of the [TRANSFAC®](#) database. We searched for so called **composite modules** that act as potential condition-specific **enhancers** of the **target genes** in their upstream regulatory regions (-1000 bp upstream of transcription start site (TSS)) and identify transcription factors regulating activity of the genes through such **enhancers**.

Classically, **enhancers** are defined as regions in the genome that increase transcription of one or several genes when inserted in either orientation at various distances upstream or downstream of the gene [8]. Enhancers typically have a length of several hundreds of nucleotides and are bound by multiple transcription factors in a cooperative manner [9].

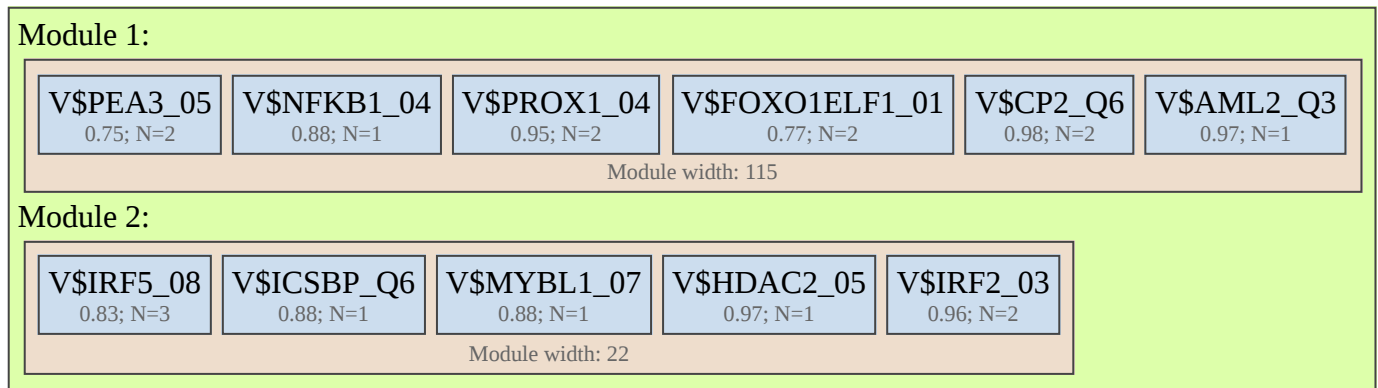
We applied the Composite Module Analyst (CMA) [8] method to detect such potential enhancers, as targets of multiple TFs bound in a cooperative manner to the regulatory regions of the genes of interest. CMA applies a genetic algorithm to construct a generalized model of the enhancers by specifying combinations of TF motifs (from [TRANSFAC®](#)) whose sites are most frequently clustered together in the regulatory regions of the studied genes. CMA identifies the transcription factors that through their cooperation provide a synergistic effect and thus have a great influence on the gene regulation process.

#### **Enhancer model potentially involved in regulation of target genes (high expressed genes in Experiment).**

To build the most specific composite modules we choose top high expressed genes as the input of CMA algorithm.

The model consists of 2 module(s). Below, for each module the following information is shown:

- PWMs producing matches,
- number of individual matches for each PWM,
- score of the best match.



**Model score (-p\*log10(pval)):** 20.94

**Wilcoxon p-value (pval):** 1.01e-43

**Penalty (p):** 0.487

**Average yes-set score:** 4.69

**Average no-set score:** 3.21

**AUC:** 0.79

**Separation point:** 3.45

**False-positive:** 36.00%

**False-negative:** 18.67%

The AUC of the model achieves value significantly higher than expected for a random set of regulatory regions

Z-score = 3.09

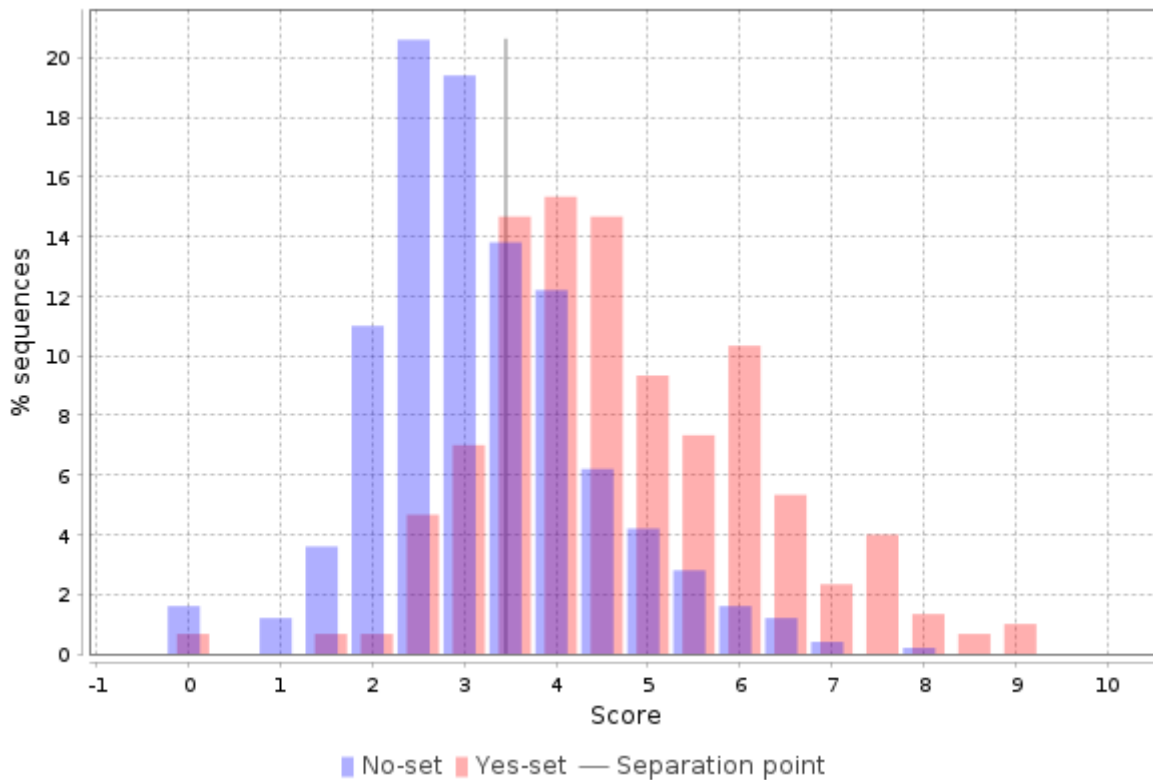


Table 4. List of top ten high expressed genes in Experiment with identified enhancers in their regulatory regions. **CMA score** - the score of the CMA model of the enhancer identified in the regulatory region.

[See full table](#) →

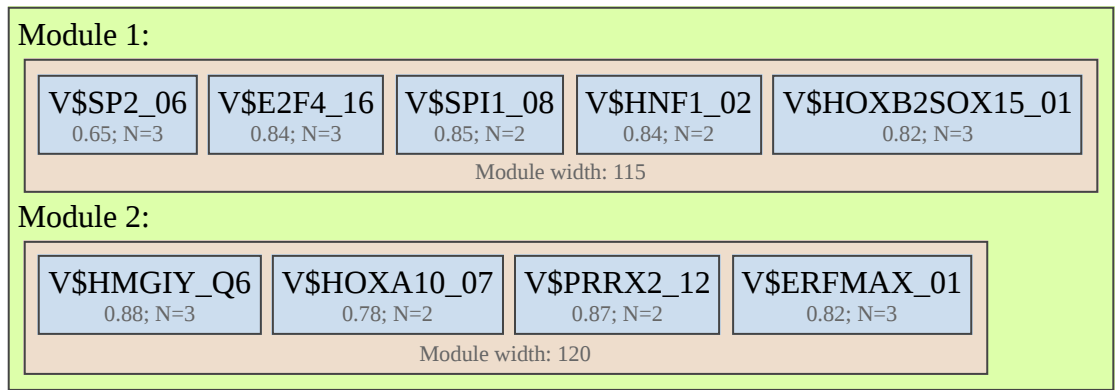
Ensembl IDs	Gene symbol	Gene description	CMA score	Factor names
ENSG00000102081	FMR1	fragile X messenger ribonucleoprotein 1	9.88	PROX-1(h), ELF-1(h), FOXO1(h), Runx3(h), CP2(h), ETV4(h), IRF-2(h), IRF-5(h)...
ENSG00000143093	STRIP1	striatin interacting protein 1	9.3	PROX-1(h), ELF-1(h), FOXO1(h), IRF-5(h), IRF-8(h), IRF-2(h), ETV4(h), NF-kappaB-p105(h)...
ENSG00000136514	RTP4	receptor transporter protein 4	9.22	CP2(h), PROX-1(h), ETV4(h), ELF-1(h), FOXO1(h), IRF-5(h), IRF-2(h), IRF-8(h)...
ENSG00000152778	IFIT5	interferon induced protein with tetratricopeptide repeats 5	9.17	NF-kappaB-p105(h), ETV4(h), ELF-1(h), FOXO1(h), PROX-1(h), IRF-5(h), IRF-8(h), IRF-2(h)
ENSG00000130303	BST2	bone marrow stromal cell antigen 2	8.92	ELF-1(h), FOXO1(h), IRF-5(h), IRF-8(h), IRF-2(h), ETV4(h), NF-kappaB-p105(h), CP2(h)...
ENSG00000166710	B2M	beta-2-microglobulin	8.85	PROX-1(h), ETV4(h), CP2(h), NF-kappaB-p105(h), ELF-1(h), FOXO1(h), IRF-8(h), IRF-2(h)...
ENSG00000188389	PDCD1	programmed cell death 1	8.8	ETV4(h), PROX-1(h), NF-kappaB-p105(h), CP2(h), ELF-1(h), FOXO1(h), IRF-8(h), IRF-2(h)
ENSG00000169871	TRIM56	tripartite motif containing 56	8.79	ETV4(h), CP2(h), IRF-8(h), IRF-2(h), IRF-5(h), ELF-1(h), FOXO1(h), Runx3(h)
ENSG00000228775	WEE2-AS1	WEE2 antisense RNA 1	8.76	A-Myb(h), IRF-8(h), IRF-5(h), ELF-1(h), FOXO1(h), IRF-2(h), CP2(h), ETV4(h)...
ENSG00000089692	LAG3	lymphocyte activating 3	8.64	CP2(h), IRF-8(h), IRF-2(h), ELF-1(h), FOXO1(h), IRF-5(h), A-Myb(h), ETV4(h)...

### Enhancer model potentially involved in regulation of target genes (low expressed genes in Experiment).

To build the most specific composite modules we choose top low expressed genes as the input of CMA algorithm.

The model consists of 2 module(s). Below, for each module the following information is shown:

- PWMs producing matches,
- number of individual matches for each PWM,
- score of the best match.



**Model score (-p\*log10(pval)): 19.99**

**Wilcoxon p-value (pval): 2.28e-39**

**Penalty (p): 0.517**

**Average yes-set score: 9.01**

**Average no-set score: 7.44**

**AUC: 0.78**

**Separation point: 8.78**

**False-positive: 21.40%**

**False-negative: 35.33%**

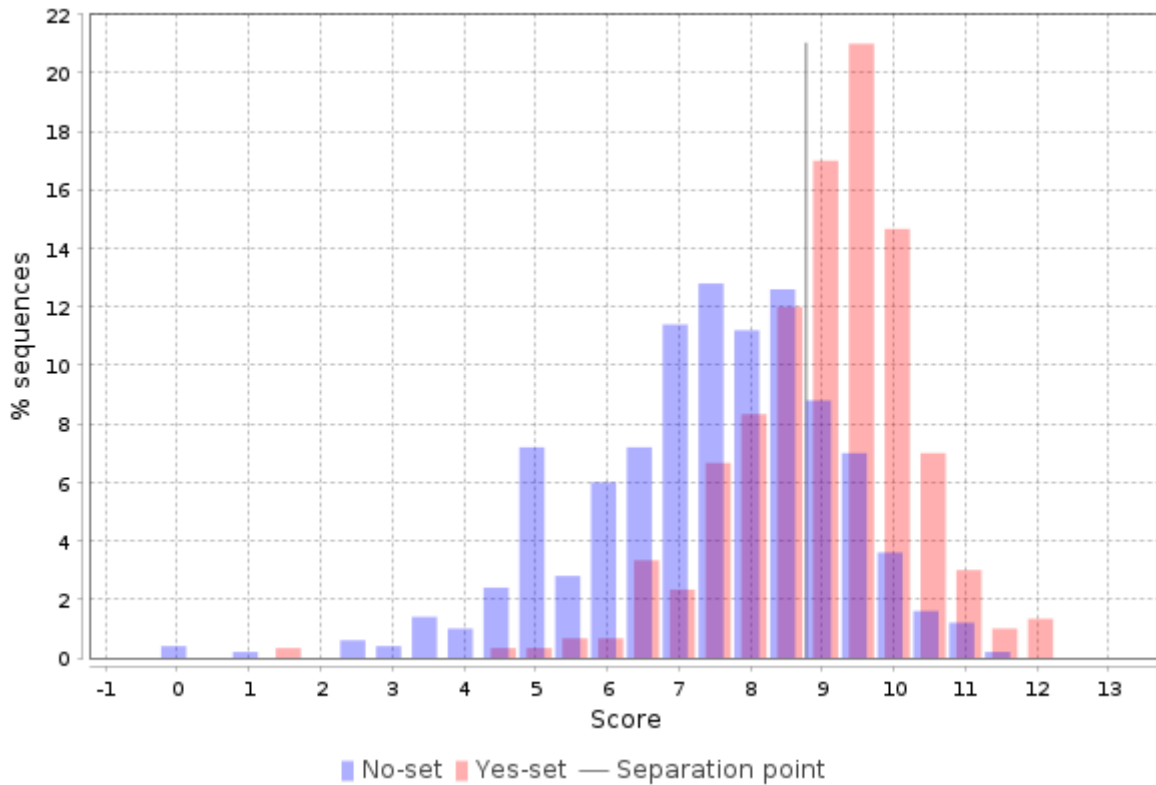


Table 5. List of top ten low expressed genes in Experiment with identified enhancers in their regulatory regions. **CMA score** - the score of the CMA model of the enhancer identified in the regulatory region.

[See full table](#) →

Ensembl IDs	Gene symbol	Gene description	CMA score	Factor names
ENSG00000079150	FKBP7	FKBP prolyl isomerase 7	13.03	Sp2(h), E2F-4(h), HNF-1alpha(h), PRRX-2(h), Hox-A10(h), HMGA1(h),HMGA2(h), Hox-B2(h),SOX-15(h)...
ENSG00000140319	SRP14	signal recognition particle 14	12.81	PU.1(h), Hox-B2(h),SOX-15(h), Sp2(h), E2F-4(h), HMGA1(h),HMGA2(h), HNF-1alpha(h), PRRX-2(h)...
ENSG00000068784	SRBD1	S1 RNA binding domain 1	12.69	E2F-4(h), HNF-1alpha(h), Sp2(h), Hox-A10(h), PU.1(h), HMGA1(h),HMGA2(h), Hox-B2(h),SOX-15(h)...
ENSG00000168646	AXIN2	axin 2	12.62	Sp2(h), E2F-4(h), HNF-1alpha(h), PU.1(h), HMGA1(h),HMGA2(h), PRRX-2(h), Hox-A10(h)
ENSG00000185652	NTF3	neurotrophin 3	12.54	ERF(h),Max(h), Hox-A10(h), HMGA1(h),HMGA2(h), PRRX-2(h), Hox-B2(h),SOX-15(h), E2F-4(h), PU.1(h)...
ENSG00000100483	VCPKMT	valosin containing protein lysine methyltransferase	12.46	Sp2(h), E2F-4(h), HNF-1alpha(h), Hox-A10(h), PRRX-2(h), HMGA1(h),HMGA2(h), Hox-B2(h),SOX-15(h)
ENSG00000198856	OSTC	oligosaccharyltransferase complex non-catalytic subunit	12.45	PU.1(h), E2F-4(h), Sp2(h), Hox-B2(h),SOX-15(h), HMGA1(h),HMGA2(h), Hox-A10(h), PRRX-2(h)
ENSG00000116199	FAM20B	FAM20B glycosaminoglycan xylosylkinase	12.4	HMGA1(h),HMGA2(h), PRRX-2(h), Hox-B2(h),SOX-15(h), PU.1(h), Sp2(h), HNF-1alpha(h), E2F-4(h)
ENSG00000232229	LINC00865	long intergenic non-protein coding RNA 865	12.38	PU.1(h), PRRX-2(h), HMGA1(h),HMGA2(h), Hox-A10(h), Hox-B2(h),SOX-15(h), HNF-1alpha(h), E2F-4(h)...
ENSG00000122779	TRIM24	tripartite motif containing 24	12.29	Hox-B2(h),SOX-15(h), HMGA1(h),HMGA2(h), PRRX-2(h), HNF-1alpha(h), Hox-A10(h), E2F-4(h), PU.1(h)...

On the basis of the enhancer models we identified transcription factors potentially regulating the **target genes** of our interest. We found 12 and 12 transcription factors controlling expression of up- and down-regulated genes respectively (see Tables 6-7).

Table 6. Transcription factors of the predicted enhancer model potentially regulating the differentially expressed genes (high expressed genes in Experiment). **Yes-No ratio** is the ratio between frequencies of the sites in Yes sequences versus No sequences. It describes the level of the enrichment of binding sites for the indicated TF in the regulatory target regions. **Regulatory score** is the measure of involvement of the given TF in the controlling of expression of genes that encode master regulators presented below (through positive feedback loops).

[See full table](#) →

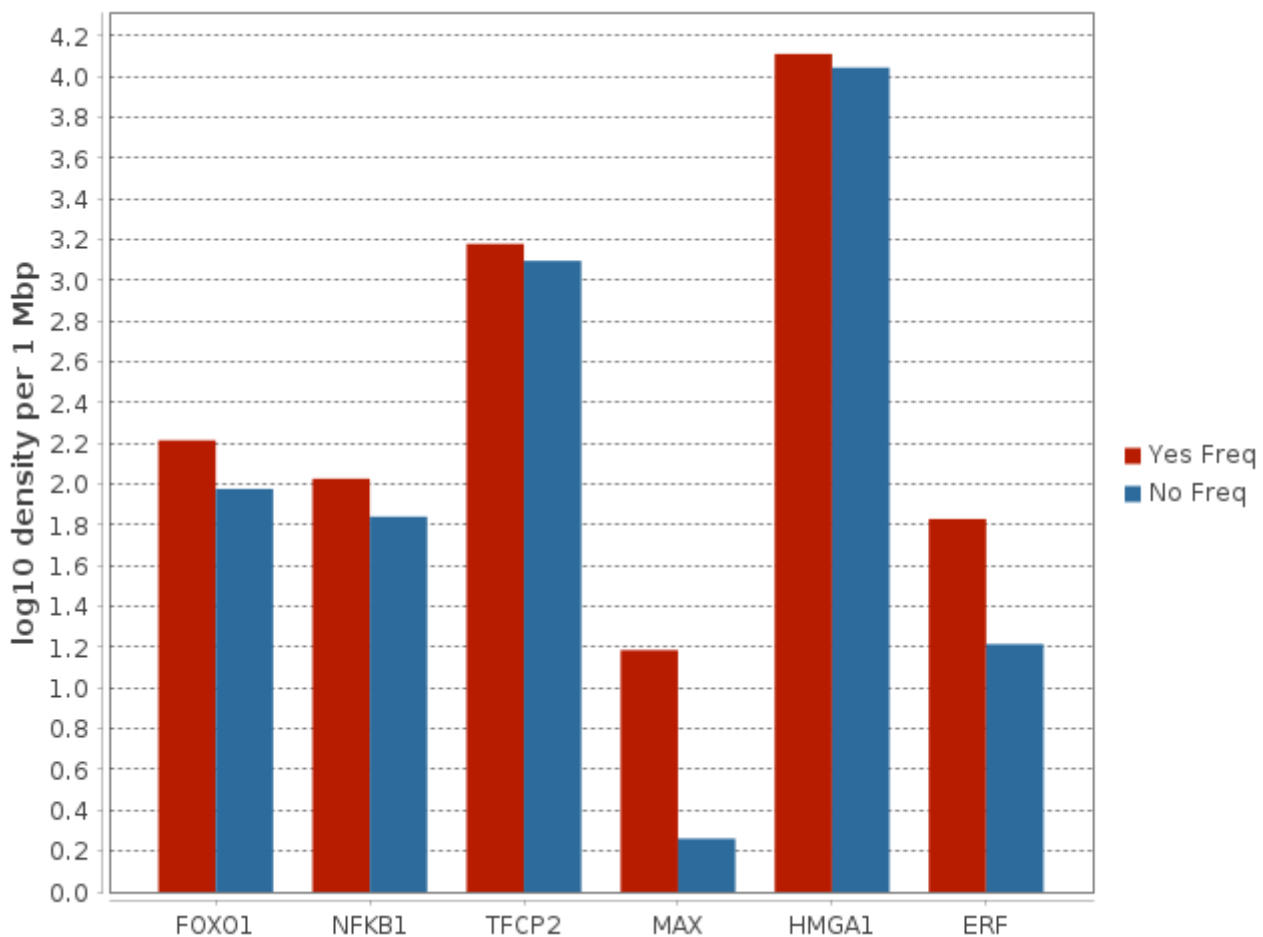
ID	Gene symbol	Gene description	Regulatory score	Yes-No ratio
MO000034454	FOXO1	forkhead box O1	4.42	1.73
MO000019359	NFKB1	nuclear factor kappa B subunit 1	4.02	1.54
MO000117988	TFCP2	transcription factor CP2	3.83	1.21
MO000026238	RUNX3	RUNX family transcription factor 3	3.62	1.96
MO000025410	ELF1	E74 like ETS transcription factor 1	3.27	5.84
MO000046009	ETV4	ETS variant transcription factor 4	3.17	1.95
MO000007691	IRF2	interferon regulatory factor 2	3.11	33.35
MO000046015	MYBL1	MYB proto-oncogene like 1	3.02	2.5
MO000023424	IRF8	interferon regulatory factor 8	2.47	8.34
MO000058923	HDAC2	histone deacetylase 2	2.42	1.52

Table 7. Transcription factors of the predicted enhancer model potentially regulating the differentially expressed genes (low expressed genes in Experiment). **Yes-No ratio** is the ratio between frequencies of the sites in Yes sequences versus No sequences. It describes the level of the enrichment of binding sites for the indicated TF in the regulatory target regions. **Regulatory score** is the measure of involvement of the given TF in the controlling of expression of genes that encode master regulators presented below (through positive feedback loops).

[See full table](#) →

ID	Gene symbol	Gene description	Regulatory score	Yes-No ratio
MO000022234	MAX	MYC associated factor X	2.88	8.41
MO000026358	HMGA1	high mobility group AT-hook 1	2.88	1.17
MO000028675	ERF	ETS2 repressor factor	2.8	4.11
MO000023603	E2F4	E2F transcription factor 4	2.74	1.44
MO000085616	SPI1	Spi-1 proto-oncogene	2.72	1.31
MO000082618	HNF1A	HNF1 homeobox A	2.62	2.16
MO000046082	SP2	Sp2 transcription factor	2.5	1.72
MO000255539	HMGA2	high mobility group AT-hook 2	2.48	1.17
MO000089495	HOXA10	homeobox A10	2.43	1.65
MO000219104	PRRX2	paired related homeobox 2	2.32	1.7

The following diagram represents the key transcription factors, which were predicted to be potentially regulating differentially expressed genes in the analyzed pathology: FOXO1, NFKB1, TFCP2, MAX, HMGA1 and ERF.



### **3.4. Finding master regulators in networks**

In the second step of the upstream analysis common regulators of the revealed TFs were identified. These master regulators appear to be the key candidates for therapeutic targets as they have a master effect on regulation of intracellular pathways that activate the pathological process of our study. The identified master regulators are shown in Tables 8-9.

Table 8. Master regulators that may govern the regulation of high expressed genes in Experiment. **Total rank** is the sum of the ranks of the master molecules sorted by keynode score, CMA score, transcriptomics data.

[See full table](#) →

ID	Master molecule name	Gene symbol	Gene description	Total rank	LogFoldChange
<a href="#">MO000143731</a>	UBP43(h)	USP18	ubiquitin specific peptidase 18	54	2.79
<a href="#">MO001084877</a>	isg15(h):UbcH8(h):ISG15 E3 ligases(h)	ISG15, UBE2L6	ISG15 ubiquitin like modifier, ubiquitin conjugating enzyme E2 L6	74	3.63
<a href="#">MO001084878</a>	isg15:UbcH8:ISG15 E3 ligases	ARIH1, HERC5, ISG15, TRIM25, UBE2L6	HECT and RLD domain containing E3 ubiquitin protein ligase 5, ISG15 ubiquitin like modifier, ariadne...	94	3.63
<a href="#">MO000019506</a>	STAT1alpha(h)	STAT1	signal transducer and activator of transcription 1	102	2.51
<a href="#">MO000019521</a>	STAT1(h)	STAT1	signal transducer and activator of transcription 1	112	2.51
<a href="#">MO000143730</a>	UBP43-isoform1(h)	USP18	ubiquitin specific peptidase 18	144	2.79
<a href="#">MO000335346</a>	UBP43-isoform2(h)	USP18	ubiquitin specific peptidase 18	144	2.79
<a href="#">MO001091834</a>	ISGylated host proteins	BECN1, EIF2AK2, IFIT1, ISG15, JAK1, MAPK3, MX1, MX2, PLCG1, RIGI, STAT1, UBE2E1	ISG15 ubiquitin like modifier, Janus kinase 1, MX dynamin like GTPase 1, MX dynamin like GTPase 2, R...	145	3.71
<a href="#">MO001091833</a>	ISGylated host proteins(h)	ISG15	ISG15 ubiquitin like modifier	150	3.63
<a href="#">MO001076186</a>	(PKR(h){pT88}{pT89}{pT90}{pY101}{pY162}{pS242}{pT255}{pT258}{pY293}{pT446}{pT451})2	EIF2AK2	eukaryotic translation initiation factor 2 alpha kinase 2	262	1.05

Table 9. Master regulators that may govern the regulation of low expressed genes in Experiment. **Total rank** is the sum of the ranks of the master molecules sorted by keynode score, CMA score, transcriptomics data.

[See full table](#) →

ID	Master molecule name	Gene symbol	Gene description	Total rank	LogFoldChange
<a href="#">MO000038235</a>	itch(h)	ITCH	itchy E3 ubiquitin protein ligase	231	-0.74
<a href="#">MO000007566</a>	InsR(h)	INSR	insulin receptor	341	-0.47
<a href="#">MO001096056</a>	(MLKL(h){pT357}{pS358})4: (RIP(h){pS166}:RIP3(h){pS199} {pS227})n:itch(h):UbcH7(h) {ubC86}:Ub(h)	ITCH, MLKL, RIPK1, RIPK3, UBE2L3	itchy E3 ubiquitin protein ligase, mixed lineage kinase domain like pseudokinase, receptor interacti...	480	-0.74
<a href="#">MO000082690</a>	Itch-isoform2(h)	ITCH	itchy E3 ubiquitin protein ligase	484	-0.74
<a href="#">MO000021242</a>	TAK1(h)	MAP3K7	mitogen-activated protein kinase kinase kinase 7	513	-0.5
<a href="#">MO000114255</a>	AMPKalpha-2(h)	PRKAA2	protein kinase AMP-activated catalytic subunit alpha 2	577	-0.53
<a href="#">MO000032766</a>	AKT-2(h)	AKT2	AKT serine/threonine kinase 2	604	-0.35
<a href="#">MO000329204</a>	Cdk6(h):cyclinD3-isoform1(h)	CCND3, CDK6	cyclin D3, cyclin dependent kinase 6	627	-0.34
<a href="#">MO000030927</a>	DNA-PKcs(h)	PRKDC	protein kinase, DNA-activated, catalytic subunit	628	-0.52
<a href="#">MO000020073</a>	Ubc5A(h)	UBE2D1	ubiquitin conjugating enzyme E2 D1	638	-0.41

The intracellular regulatory pathways controlled by the above-mentioned master regulators are depicted in Figures 8 and 9. These diagrams display the connections between identified transcription factors, which play important roles in the regulation of differentially expressed genes, and selected master regulators, which are responsible for the regulation of these TFs.

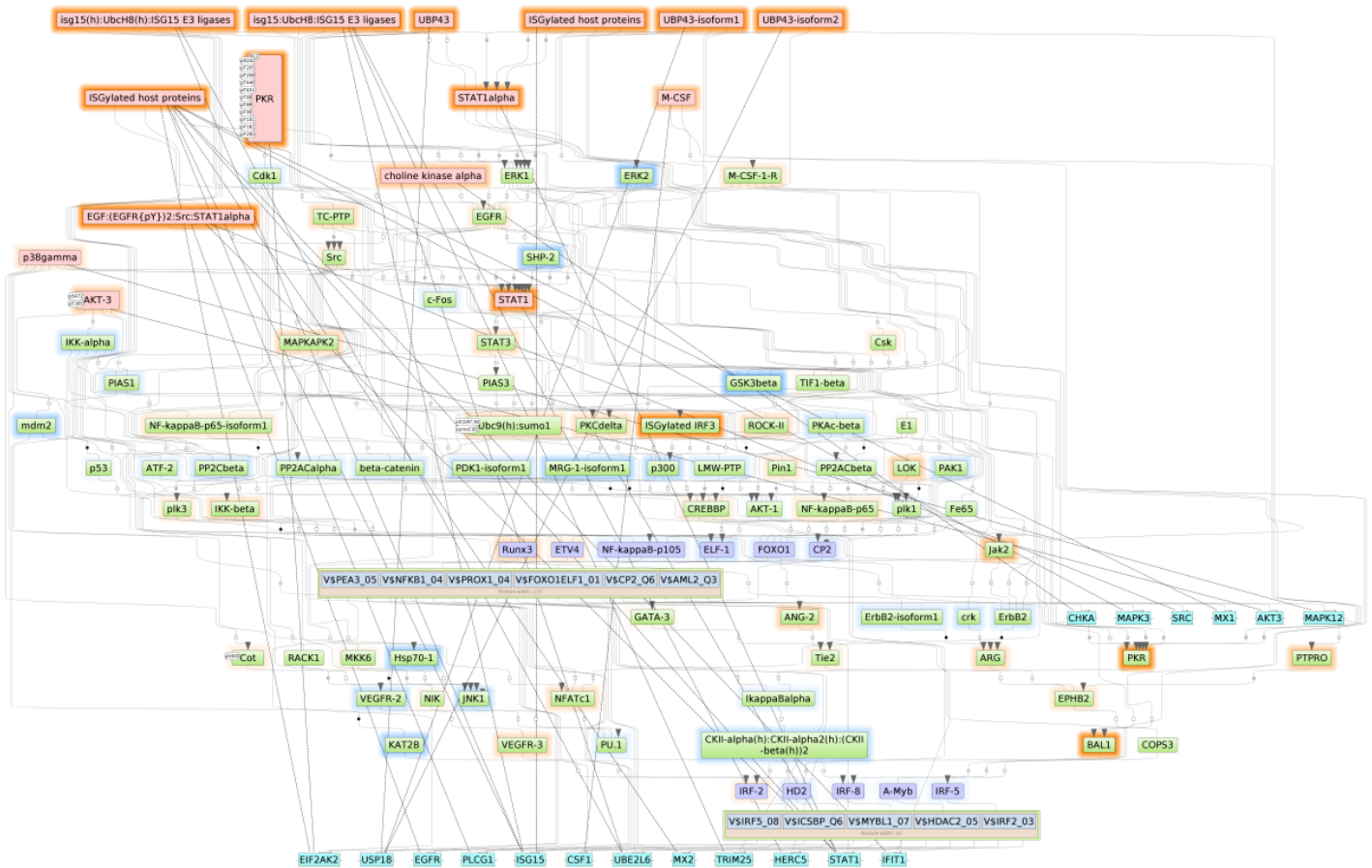


Figure 8. Diagram of intracellular regulatory signal transduction pathways of high expressed genes in Experiment. Master regulators are indicated by red rectangles, transcription factors are blue rectangles, and green rectangles are intermediate molecules, which have been added to the network during the search for master regulators from selected TFs. Orange and blue frames highlight molecules that are encoded by up- and downregulated genes, resp.

[See full diagram](#) →

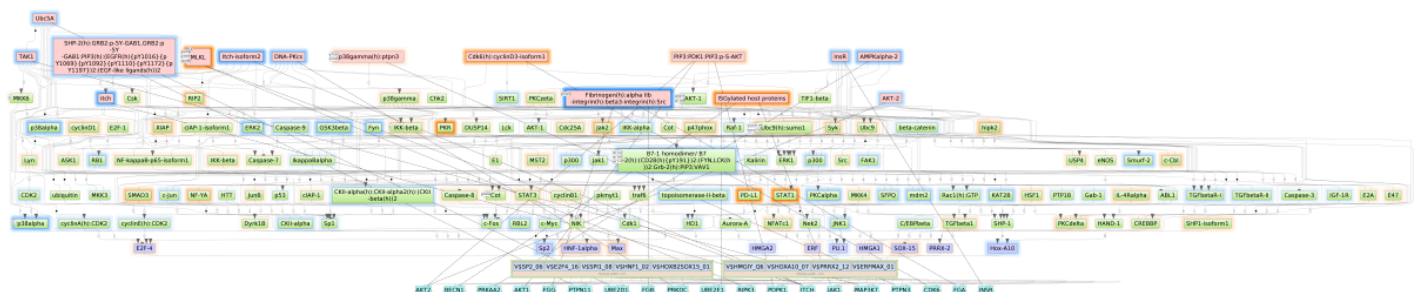


Figure 9. Diagram of intracellular regulatory signal transduction pathways of low expressed genes in Experiment. Master regulators are indicated by red rectangles, transcription factors are blue rectangles, and green rectangles are intermediate molecules, which have been added to the network during the search for master regulators from selected TFs. Orange and blue frames highlight molecules that are encoded by up- and downregulated genes, resp.

[See full diagram](#) →


## 4. Finding prospective drug targets

The identified master regulators that may govern pathology associated genes were checked for druggability potential using HumanPSD™ [5] database of gene-disease-drug assignments and PASS [12-14] software for prediction of biological activities of chemical compounds on the basis of a (Q)SAR approach. Respectively, for each master regulator protein we have computed two Druggability scores: HumanPSD Druggability score and PASS Druggability score. Where Druggability score represents the number of drugs that are potentially suitable for inhibition (or activation) of the corresponding target either according to the information extracted from medical literature (from HumanPSD™ database) or according to cheminformatics predictions of compounds activity against the examined target (from PASS software).

The cheminformatics druggability check is done using a pre-computed database of spectra of biological activities of chemical compounds from a library of all small molecular drugs from HumanPSD™ database, 2507 pharmaceutically active known chemical compounds in total. The spectra of biological activities has been computed using the program PASS [12-14] on the basis of a (Q)SAR approach.


If both Druggability scores were below defined thresholds (see Methods section for the details) such master regulator proteins were not used in further analysis of drug prediction.

As a result we created the following two tables of prospective drug targets (top targets are shown here):

 *Table 10. Prospective drug targets selected from full list of identified master regulators filtered by Druggability score from HumanPSD™ database. **Druggability score** contains the number of drugs that are potentially suitable for inhibition (or activation) of the target. The drug targets are sorted according to the **Total rank** which is the sum of three ranks computed on the basis of the three scores: keynode score, CMA score and expression change score (logFC, if present). See Methods section for details.*

[See full table →](#)

Gene symbol	Gene Description	Druggability score	Total rank	LogFoldChange
ISG15	ISG15 ubiquitin like modifier	1	74	3.63
TRIM25	tripartite motif containing 25	1	94	3.63
STAT1	signal transducer and activator of transcription 1	4	102	2.51
EIF2AK2	eukaryotic translation initiation factor 2 alpha kinase 2	28	145	3.71
MAPK3	mitogen-activated protein kinase 3	137	145	3.71
PLCG1	phospholipase C gamma 1	3	145	3.71

 *Table 11. Prospective drug targets selected from full list of identified master regulators filtered by Druggability score predicted by PASS software. Here, the **Druggability score** for master regulator proteins is computed as a sum of PASS calculated probabilities to be active as a target for various small molecular compounds. The drug targets are sorted according to the **Total rank** which is the sum of three ranks computed on the basis of the three scores: keynode score, CMA score and expression change score (logFC, if present). See Methods section for details.*

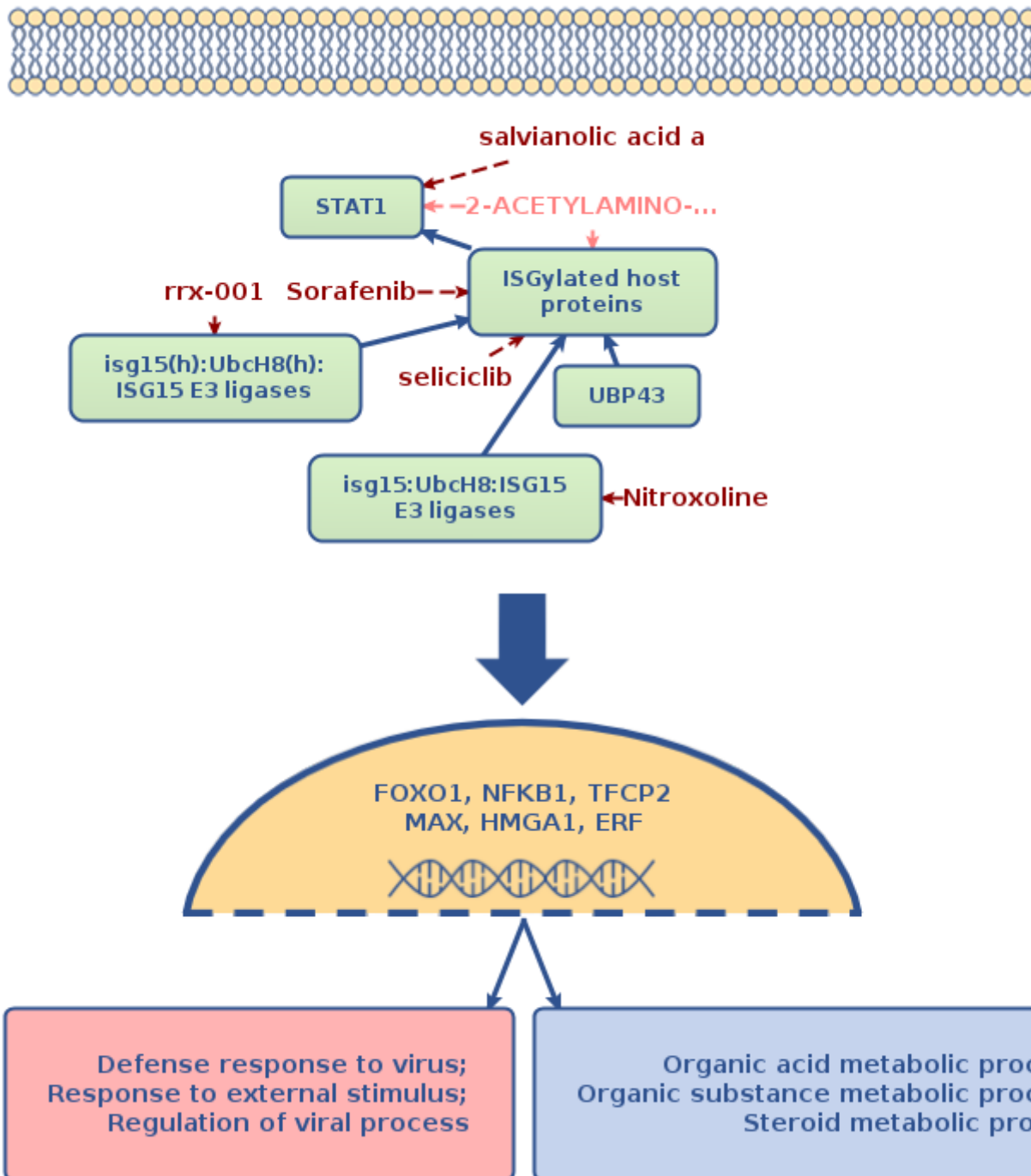
[See full table →](#)

Gene symbol	Gene Description	Druggability score	Total rank	LogFoldChange
STAT1	signal transducer and activator of transcription 1	3.01	102	2.51
EIF2AK2	eukaryotic translation initiation factor 2 alpha kinase 2	1.07	145	3.71
MAPK3	mitogen-activated protein kinase 3	11.73	145	3.71
PLCG1	phospholipase C gamma 1	30.95	145	3.71
STAT2	signal transducer and activator of transcription 2	3.01	325	2.51
JAK2	Janus kinase 2	1.07	342	0.54

Below we represent schematically the main mechanism of the studied pathology. In the schema we considered the top two drug targets of each of the two categories computed above. In addition we have added two top identified master regulators for which no drugs may be identified yet, but that are playing the crucial role in the molecular mechanism of the studied pathology. Thus the molecular mechanism of the studied pathology was predicted to be mainly based on the following key master regulators:

- isg15(h):UbcH8(h):ISG15 E3 ligases
- isg15:UbcH8:ISG15 E3 ligases
- ISGylated host proteins
- STAT1
- UBP43

This result allows us to suggest the following schema of affecting the molecular mechanism of the studied pathology:



Drugs which are shown on this schema: rrx-001, Sorafenib, seliciclib, salvianolic acid a, 2-ACETYLAMINO-4-METHYL-PENTANOIC ACID [1-(1-FORMYL-PENTYL-CARBAMOYL)-3-METHYL-BUTYL]-AMIDE and Nitroxoline, should be considered as a prospective research initiative for further drug repurposing and drug development. These drugs were selected as top matching treatments to the most prospective drug targets of the studied pathology, however, these results should be considered with special caution and are to be used for research purposes only, as there is not enough clinical information for adapting these results towards immediate treatment of patients.

The drugs given in dark red color on the schema are FDA approved drugs or drugs which have gone through various phases of clinical trials as active treatments against the selected targets.

The drugs given in pink color on the schema are drugs, which were cheminformatically predicted to be active against the selected targets.

## 5. Identification of potential drugs

In the last step of the analysis we strived to identify known activities as well as drugs with cheminformatically predicted activities that are potentially suitable for inhibition (or activation) of the identified molecular targets in the context of specified human diseases(s).

Proposed drugs are top ranked drug candidates, that were found to be active on the identified targets and were selected from 4 categories:

1. FDA approved drugs or used in clinical trials drugs for the studied pathology;
2. Repurposing drugs used in clinical trials for other pathologies;
3. Drugs, predicted by PASS to be active against identified drug targets and against the studied pathology;
4. Drugs, predicted by PASS to be active against identified drug targets but for other pathologies.

Proposed drugs were selected on the basis of Drug rank which was computed from the ranks sum based on the individual ranks of the following scores:

- Target activity score (depends on ranks of all targets that were found for the selected drug);
- Disease activity score (weighted sum of number of clinical trials on disease(s) under study where the selected drug is known to be applied or PASS Disease activity score - cheminformatically predicted property of the compound to be active against the studied disease(s));
- Clinical validity score (applicable only for drugs predicted on the basis of literature curation in HumanPSD™ database (Tables 12 and 13), reflects the number of the highest clinical trials phase on which the drug was tested for any pathology).

You can refer to the Methods section for more details on drug ranking procedure.

Based on the Drug rank, a numerical value of Drug score was calculated, which reflects the potential activity of the respective drug on the overall molecular mechanism of the studied pathology. Drug score values belong to the range from 1 to 100 and are calculated as a quotient of maximum drug rank and the drug rank of the given drug multiplied by 100.

Top drugs of each category are given in the tables below:

## Drugs approved in clinical trials



Table 12. FDA approved drugs or drugs used in clinical trials for the studied pathology (most promising treatment candidates selected for the identified drug targets on the basis of literature curation in [HumanPSD™](#) database)

[See full table](#) →

Name	Target names	Drug score	Disease activity score	Disease trial phase
<a href="#">Sorafenib</a>	TEC, IKBKE, JAK3, PRKACA, MAP3K11, SYK, EPHA1, NEK6, EIF2AK2, LIMK1, MAPK3, MAP2K6, AXL, ACVR2B, SRC, CAMKK2, PRKD3, PIP4K2B, CSNK1G2, MAP4K1, DMPK, STK11, PRKCZ, CSNK1G1, MAPK12, PLK3, WEE1, MAPK7, PRKCD, LYN, CHEK2, MUSK, STK3, MAPK10, NTRK1, CSNK1E, PDGFRA, AURKB, CDK7, CSNK1D, HCK, CHEK1, PTK6, NUA2, HIPK2, RIPK2, CDK9, PAK4, TYRO3, ITK, MAP3K4, FGR, NLK, RET, ABL2, CSF1R, STK10, ROCK2, MARK3, BLK, EPHA2, CAMK2B, CSK, TEK, PKN1, TYK2, CAMK4, TNIK, MAP3K5, CSNK2A2, MAPK4, RIPK1, STK4, CAMK2A, MAPKAPK2, IKBKB, CAMK2G, MET, NTRK3, PRKCQ, EPHA4, LATS1, MAP2K4, RPS6KA1, FLT4, EGFR, ACVR2A, JAK2, PKMYT1, RPS6KB1, ALK, EPHB2, MERTK, EPHA3, EPHB4, PRKCE, BRAF, FER, MAPK13, AKT3, ZAP70, PRKD2	98	1	small molecule, approved, investigational
<a href="#">Sirolimus</a>	IKKBK, MAPK10, ROCK2, MARK3, PRKACA, ITGAL, IL10, AURKB, RPS6KA1, CSNK1D, TGM2, NFE2L2, FKBP1A, NEK6, CHEK1, CSK, MAPK3, RPS6KB1, HIPK2, CAMKK2, PAK4, MAPK13, PRKCZ, MAPK12, MAPKAPK2, CHEK2, STK3	94	1	small molecule, approved, investigational
<a href="#">Everolimus</a>	ZEB1, AKT3, BCL2, CASP8, RICTOR, RPS6KB1, RPTOR	83	1	small molecule, approved
<a href="#">Pirfenidone</a>	MAPK12, TGFB1, TNF, MAPK13, FURIN	83	1	small molecule, investigational, approved
<a href="#">Rifaximin</a>	MAPK10, TGFB1, TNF, IL6, CXCL8, NR1I2	81	3	small molecule, approved, investigational

The **Disease trial phase** column reflects the maximum clinical trials phase in which the drug was studied for the analyzed pathology.

## Repurposing drugs



Table 13. Repurposed drugs used in clinical trials for other pathologies (prospective drugs against the identified drug targets on the basis of literature curation in [HumanPSD™](#) database)

[See full table](#) →

Name	Target names	Drug score	Maximum trial phase
<a href="#">seliciclib</a>	TEC, IKBKE, JAK3, PRKACA, MAP3K11, CDK4, SYK, EPHA1, NEK6, EIF2AK2, LIMK1, MAPK3, MAP2K6, AXL, ACVR2B, SRC, CAMKK2, PRKD3, PIP4K2B, CSNK1G2, MAP4K1, DMPK, STK11, PRKCZ, CSNK1G1, MAPK12, PLK3, WEE1, MAPK7, PRKCD, LYN, CHEK2, MUSK, STK3, MAPK10, NTRK1, CSNK1E, PDGFRA, AURKB, CDK7, CSNK1D, HCK, CHEK1, PTK6, NUAKE2, HIPK2, RIPK2, CDK9, PAK4, TYRO3, ITK, MAP3K4, FGR, NLK, RET, ABL2, CSF1R, STK10, ROCK2, MARK3, BLK, EPHA2, CAMK2B, CSK, TEK, PKN1, TYK2, CAMK4, TNIK, MAP3K5, CSNK2A2, MAPK4, RIPK1, STK4, CAMK2A, MAPKAPK2, IKBKB, CAMK2G, MET, NTRK3, PRKCQ, EPHA4, LATS1, MAP2K4, RPS6KA1, FLT4, EGFR, ACVR2A, JAK2, PKMYT1, RPS6KB1, ALK, EPHB2, MERTK, EPHA3, EPHB4, PRKCE, BRAF, FER, MAPK13, AKT3, ZAP70, PRKD2	96	PHASE2: Cystic Fibrosis, Cysts, Fibrosis
<a href="#">ruboxistaurin</a>	TEC, IKBKE, JAK3, PRKACA, MAP3K11, SYK, EPHA1, NEK6, EIF2AK2, LIMK1, MAPK3, MAP2K6, AXL, ACVR2B, PRKCG, SRC, CAMKK2, PRKD3, PIP4K2B, CSNK1G2, MAP4K1, DMPK, STK11, PRKCZ, CSNK1G1, MAPK12, PLK3, WEE1, MAPK7, PRKCD, LYN, CHEK2, MUSK, STK3, MAPK10, NTRK1, CSNK1E, PDGFRA, AURKB, CDK7, CSNK1D, HCK, CHEK1, PTK6, NUAKE2, HIPK2, RIPK2, CDK9, PAK4, TYRO3, ITK, MAP3K4, FGR, NLK, RET, ABL2, CSF1R, STK10, ROCK2, MARK3, BLK, EPHA2, CAMK2B, CSK, TEK, PKN1, TYK2, CAMK4, TNIK, MAP3K5, CSNK2A2, MAPK4, RIPK1, STK4, CAMK2A, MAPKAPK2, IKBKB, CAMK2G, MET, NTRK3, PRKCQ, EPHA4, LATS1, MAP2K4, RPS6KA1, FLT4, EGFR, ACVR2A, JAK2, PKMYT1, RPS6KB1, ALK, EPHB2, MERTK, EPHA3, EPHB4, PRKCE, BRAF, FER, MAPK13, AKT3, ZAP70, PRKD2	95	PHASE1: Diabetes Mellitus, Diabetes Mellitus, Type 2, Heart Failure
<a href="#">1-(5-Tert-Butyl-2-P-Tolyl-2h-Pyrazol-3-Yl)-3-[4-(2-Morpholin-4-Yl-Ethoxy)-Naphthalen-1-Yl]-Urea</a>	TEC, IKBKE, JAK3, PRKACA, MAP3K11, SYK, EPHA1, NEK6, EIF2AK2, LIMK1, MAPK3, MAP2K6, AXL, ACVR2B, SRC, CAMKK2, PRKD3, PIP4K2B, CSNK1G2, MAP4K1, DMPK, STK11, PRKCZ, CSNK1G1, MAPK12, PLK3, WEE1, MAPK7, PRKCD, LYN, CHEK2, MUSK, STK3, MAPK10, NTRK1, CSNK1E, PDGFRA, AURKB, CDK7, CSNK1D, HCK, CHEK1, PTK6, NUAKE2, HIPK2, RIPK2, CDK9, PAK4, TYRO3, ITK, MAP3K4, FGR, NLK, RET, ABL2, CSF1R, STK10, ROCK2, MARK3, BLK, EPHA2, CAMK2B, CSK, TEK, PKN1, TYK2, CAMK4, TNIK, MAP3K5, CSNK2A2, MAPK4, RIPK1, STK4, CAMK2A, MAPKAPK2, IKBKB, CAMK2G, MET, NTRK3, PRKCQ, EPHA4, LATS1, MAP2K4, RPS6KA1, FLT4, EGFR, ACVR2A, JAK2, PKMYT1, RPS6KB1, ALK, EPHB2, MERTK, EPHA3, EPHB4, PRKCE, BRAF, FER, MAPK13, AKT3, ZAP70, PRKD2	95	PHASE2: Arthritis, Arthritis, Rheumatoid, Psoriasis
<a href="#">pi-103</a>	TEC, IKBKE, JAK3, PRKACA, MAP3K11, SYK, EPHA1, NEK6, EIF2AK2, LIMK1, MAPK3, MAP2K6, AXL, ACVR2B, SRC, CAMKK2, PRKD3, PIP4K2B, CSNK1G2, MAP4K1, DMPK, STK11, PRKCZ, CSNK1G1, MAPK12, PLK3, WEE1, MAPK7, PRKCD, LYN, CHEK2, MUSK, STK3, MAPK10, NTRK1, CSNK1E, PDGFRA, AURKB, CDK7, CSNK1D, HCK, CHEK1, PTK6, NUAKE2, HIPK2, RIPK2, CDK9, PAK4, TYRO3, ITK, MAP3K4, FGR, NLK, RET, ABL2, CSF1R, STK10, ROCK2, MARK3, BLK, EPHA2, CAMK2B, CSK, TEK, PKN1, TYK2, CAMK4, TNIK, MAP3K5, CSNK2A2, MAPK4, RIPK1, STK4, CAMK2A, MAPKAPK2, IKBKB, CAMK2G, MET, NTRK3, PRKCQ, EPHA4, LATS1, MAP2K4, RPS6KA1, FLT4, EGFR, ACVR2A, JAK2, PKMYT1, RPS6KB1, ALK, EPHB2, MERTK, EPHA3, EPHB4, PRKCE, BRAF, FER, MAPK13, AKT3, ZAP70, PRKD2	95	N/A

Erlotinib	TEC, IKBKE, JAK3, PRKACA, MAP3K11, SYK, EPHA1, NEK6, EIF2AK2, LIMK1, MAPK3, MAP2K6, AXL, ACVR2B, SRC, CAMKK2, PRKD3, PIP4K2B, CSNK1G2, MAP4K1, DMPK, STK11, CSNK1G1, MAPK12, PLK3, WEE1, MAPK7, PRKCD, LYN, MUSK, STK3, MAPK10, NTRK1, CSNK1E, PDGFRA, AURKB, CDK7, CSNK1D, HCK, CHEK1, PTK6, BIRC5, NUA2, ERBB3, RIPK2, CDK9, PAK4, TYRO3, ITK, MAP3K4, FGR, NLK, RET, ABL2, CSF1R, STK10, MARK3, BLK, EPHA2, CAMK2B, CSK, TEK, PKN1, TYK2, CAMK4, TNIK, MAP3K5, CSNK2A2, MAPK4, RIPK1, STK4, CAMK2A, MAPKAPK2, CAMK2G, MET, NTRK3, PRKCQ, EPHA4, LATS1, MAP2K4, RPS6KA1, FLT4, ILK, EGFR, ACVR2A, JAK2, PKMYT1, ALK, EPHB2, MERTK, EPHA3, EPHB4, PRKCE, ERBB4, BRAF, FER, AKT3, ZAP70, PRKD2, NR1I2	95	NA: Carcinoma, Non-Small-Cell Lung, Carcinoma, Squamous Cell, Head and Neck Neoplasms, Lung Neoplasms, Neoplasms
-----------	---	----	--

The **Maximum trial phase** column reflects the maximum clinical trials phase in which the drug was studied for any pathology.



No prospective drugs were found, which would be predicted by PASS software to be active against the identified drug targets and would be predicted to have biological activity against the studied disease(s).



Table 14. Prospective drugs, predicted by PASS software to be active against the identified drug targets, though without cheminformatically predicted activity against the studied disease(s) (drug candidates predicted with the cheminformatics tool PASS)

[See full table](#) →

Name	Target names	Drug score	Target activity score
Bortezomib	NFKB2, PSMC5, PSMA7, PRSS1, F2, PSMC3, PSMD4, ITGB3, ITGA2B, RELA	91	0.55
2-ACETYLAMINO-4-METHYL-PENTANOIC ACID [1-(1-FORMYL-PENTYL-CARBAMOYL)-3-METHYL-BUTYL]-AMIDE	PSMC5, PSMA7, STAT5A, PSMC3, TNFRSF10A, STAT2, STAT1, STAT5B, NGF, ITGA2B, PADI2, STAT3, PRSS1, TNFSF10, IFNAR2, PSMD4, TNF, ITGB3, STAT6	89	0.91
1-ETHOXYCARBONYL-D-PHE-PRO-2(4-AMINO-BUTYL)HYDRAZINE	STAT5A, STAT3, STAT2, STAT1, ITGB3, STAT5B, ITGA2B, STAT6	88	2.21
TI-3-093	PSMC5, PSMA7, STAT5A, PSMC3, STAT2, STAT1, STAT5B, CASP1, ITGA2B, STAT3, PRSS1, PSMD4, ITGB3, STAT6	86	0.71
Edotecarin	CAMK4, PRKCG, PNPT1, CAMK2G, CAMK2A, PRKD3, PRKACA, CAMK2B, ART1, MAP2K6, PRKCZ	84	1.01

As the result of drug search we propose the following drugs as most promising candidates for treating the pathology under study: Sorafenib, seliciclib and Bortezomib. These drugs were selected for acting on the following targets: EIF2AK2 and PSMC5, which were predicted to be active in the molecular mechanism of the studied pathology.

The selected drugs are top ranked drug candidates from each of the four categories of drugs: (1) FDA approved drugs or used in clinical trials drugs for the studied pathology; (2) repurposing drugs used in clinical trials for other pathologies; (3) drugs, predicted by PASS software to be active against the studied pathology; (4) drugs, predicted by PASS software to be repurposed from other pathologies.

## 6. Conclusion

We applied the software package "Genome Enhancer" to a data set that contains *transcriptomics* data obtained from *liver* tissue. The study is done in the context of *Hepatitis C*. The data were pre-processed, statistically analyzed and differentially expressed genes were identified. Also checked was the enrichment of GO or disease categories among the studied gene sets.

We propose the following drugs as most promising candidates for treating the pathology under study:



**Sorafenib, seliciclib and Bortezomib**

These drugs were selected for acting on the following targets: EIF2AK2 and PSMC5, which were predicted to be involved in the molecular mechanism of the pathology under study.

The identified molecular mechanism of the studied pathology was predicted to be mainly based on the following key drug targets:



**isg15(h):UbcH8(h):ISG15 E3 ligases, isg15:UbcH8:ISG15 E3 ligases, ISGylated host proteins, STAT1 and UBP43**

These potential drug targets should be considered as a prospective research initiative for further drug repurposing and drug development purposes. The following drugs were predicted as, matching those drug targets: rrx-001, Sorafenib, seliciclib, salvianolic acid a, 2-ACETYLAMINO-4-METHYL-PENTANOIC ACID [1-(1-FORMYL-PENTYL-CARBAMOYL)-3-METHYL-BUTYL]-AMIDE and Nitroxoline. These drugs should be considered with special caution for research purposes only.

In this study, we came up with a detailed signal transduction network regulating differentially expressed genes in the studied pathology. In this network we have revealed the following top master regulators (signaling proteins and their complexes) that play a crucial role in the molecular mechanism of the studied pathology, which can be proposed as the most promising molecular targets for further drug repurposing and drug development initiatives.

- isg15(h):UbcH8(h):ISG15 E3 ligases
- isg15:UbcH8:ISG15 E3 ligases
- ISGylated host proteins
- STAT1
- UBP43

Potential drug compounds which can be affecting these targets can be found in the "Finding prospective drug targets" section.

## 7. Methods

### Databases used in the study

Transcription factor binding sites in promoters and enhancers of differentially expressed genes were analyzed using known DNA-binding motifs described in the **TRANSFAC®** library, release 2024.2 (geneXplain GmbH, Wolfenbüttel, Germany) (<https://genexplain.com/transfac>).

The master regulator search uses the **TRANSPATH®** database (BIOBASE), release 2024.2 (geneXplain GmbH, Wolfenbüttel, Germany) (<https://genexplain.com/transpath>). A comprehensive signal transduction network of human cells is built by the software on the basis of reactions annotated in **TRANSPATH®**.

The information about drugs corresponding to identified drug targets and clinical trials references were extracted from **HumanPSD™** database, release 2024.2 (<https://genexplain.com/humanpsd>).

The Ensembl database release Human112.38 (hg38) (<http://www.ensembl.org>) was used for gene IDs representation and Gene Ontology (GO) (<http://geneontology.org>) was used for functional classification of the studied gene set.

## Methods for the analysis of enriched transcription factor binding sites and composite modules

Transcription factor binding sites in promoters and enhancers of differentially expressed genes were analyzed using known DNA-binding motifs. The motifs are specified using position weight matrices (PWMs) that give weights to each nucleotide in each position of the DNA binding motif for a transcription factor or a group of them.

We search for transcription factor binding sites (TFBS) that are enriched in the promoters and enhancers under study as compared to a background sequence set such as promoters of genes that were not differentially regulated under the condition of the experiment. We denote study and background sets briefly as Yes and No sets. In the current work we used a workflow considering promoter sequences of a standard length of 1100 bp (-1000 to +100). The error rate in this part of the pipeline is controlled by estimating the adjusted p-value (using the Benjamini-Hochberg procedure) in comparison to the TFBS frequency found in randomly selected regions of the human genome (adj.p-value < 0.01).

We have applied the CMA algorithm (Composite Module Analyst) for searching composite modules [7] in the promoters and enhancers of the Yes and No sets. We searched for a composite module consisting of a cluster of 10 TFs in a sliding window of 200-300 bp that statistically significantly separates sequences in the Yes and No sets (minimizing Wilcoxon p-value).

## Methods for finding master regulators in networks

We searched for master regulator molecules in signal transduction pathways upstream of the identified transcription factors. The master regulator search uses a comprehensive signal transduction network of human cells. The main algorithm of the master regulator search has been described earlier [3,4]. The goal of the algorithm is to find nodes in the global signal transduction network that may potentially regulate the activity of a set of transcription factors found at the previous step of the analysis. Such nodes are considered as most promising drug targets, since any influence on such a node may switch the transcriptional programs of hundreds of genes that are regulated by the respective TFs. In our analysis, we have run the algorithm with a maximum radius of 12 steps upstream of each TF in the input set. The error rate of this algorithm is controlled by applying it 10000 times to randomly generated sets of input transcription factors of the same set-size. Z-score and FDR value of ranks are calculated then for each potential master regulator node on the basis of such random runs (see detailed description in [9]). We control the error rate by the FDR threshold 0.05.

## Methods for analysis of pharmaceutical compounds

We seek for the optimal combination of molecular targets (key elements of the regulatory network of the cell) that potentially interact with pharmaceutical compounds from a library of known drugs and biologically active chemical compounds, using information about known drugs from HumanPSD™ and predicting potential drugs using PASS program.

### *Method for analysis of known pharmaceutical compounds*

We selected compounds from HumanPSD™ database that have at least one target. Next, we sort compounds using "Drug rank" that is the sum of the following ranks:

1. ranking by "Target activity score" ( $T\text{-score}_{PSD}$ ),
2. ranking by "Disease activity score" ( $D\text{-score}_{PSD}$ ),
3. ranking by "Clinical validity score".

"Target activity score" ( $T\text{-score}_{PSD}$ ) is calculated as follows:

$$T\text{-score}_{PSD} = -\frac{|T|}{|T| + w(|AT| - |T|)} \sum_{t \in T} \log_{10} \left( \frac{\text{rank}(t)}{1 + \text{maxRank}(T)} \right),$$

where  $T$  is set of all targets related to the compound intersected with input list,  $|T|$  is number of elements in  $T$ ,  $AT$  and  $|AT|$  are set set of all targets related to the compound and number of elements in it,  $w$  is weight multiplier,  $\text{rank}(t)$  is rank of given target,  $\text{maxRank}(T)$  equals  $\text{max}(\text{rank}(t))$  for all targets  $t$  in  $T$ .

We use following formula to calculate "Disease activity score" ( $D\text{-score}_{PSD}$ ):

$$D\text{-score}_{PSD} = \begin{cases} \sum_{d \in D} \sum_{p \in P} \text{phase}(d, p) \\ 0, D = \emptyset \end{cases},$$

where  $D$  is the set of selected diseases, and if  $D$  is empty set,  $D\text{-score}_{PSD}=0$ .  $P$  is a set of all known phases for each disease,  $\text{phase}(p,d)$  equals to the phase number if there are known clinical trials for the selected disease on this phase and zero otherwise.

The clinical validity score reflects the number of the highest clinical trials phase (from 1 to 4) on which the drug was ever tested for any pathology.

### *Method for prediction of pharmaceutical compounds*

In this study, the focus was put on compounds with high pharmacological efficiency and low toxicity. For this purpose, comprehensive library of chemical compounds and drugs was subjected to a SAR/QSAR analysis. This library contains 13040 compounds along with their pre-calculated potential pharmacological activities of those substances, their possible side and toxic effects, as well as the possible mechanisms of action. All biological activities are expressed as probability values for a substance to exert this activity ( $Pa$ ).

We selected compounds that satisfied the following conditions:

1. Toxicity below a chosen toxicity threshold (defines as  $Pa$ , probability to be active as toxic substance).
2. For all predicted pharmacological effects that correspond to a set of user selected disease(s)  $Pa$  is greater than a chosen effect threshold.
3. There are at least 2 targets (corresponding to the predicted activity-mechanisms) with predicted  $Pa$  greater than a chosen target threshold.

The maximum  $Pa$  value for all toxicities corresponding to the given compound is selected as the "Toxicity score". The maximum  $Pa$  value for all activities corresponding to the selected diseases for the given compound is used as the "Disease activity score". "Target activity score" (T-score) is calculated as follows:

$$T\text{-score}(s) = \frac{|T|}{|T| + w(|AT| - |T|)} \sum_{m \in M(s)} \left( pa(m) \sum_{g \in G(m)} IAP(g) optWeight(g) \right),$$

where  $M(s)$  is the set of activity-mechanisms for the given structure (which passed the chosen threshold for activity-mechanisms  $Pa$ );  $G(m)$  is the set of targets (converted to genes) that corresponds to the given activity-mechanism ( $m$ ) for the given compound;  $pa(m)$  is the probability to be active of the activity-mechanism ( $m$ ),  $IAP(g)$  is the invariant accuracy of prediction for gene from  $G(m)$ ;  $optWeight(g)$  is the additional weight multiplier for gene.  $T$  is set of all targets related to the compound intersected with input list,  $|T|$  is number of elements in  $T$ ,  $AT$  and  $|AT|$  are set set of all targets related to the compound and number of elements in it,  $w$  is weight multiplier.

"Druggability score" (D-score) is calculated as follows:

$$D\text{-score}(g) = IAP(g) \sum_{s \in S(g)} \sum_{m \in M(s,g)} pa(m),$$

where  $S(g)$  is the set of structures for which target list contains given target,  $M(s,g)$  is the set of activity-mechanisms (for the given structure) that corresponds to the given gene,  $pa(m)$  is the probability to be active of the activity-mechanism ( $m$ ),  $IAP(g)$  is the invariant accuracy of prediction for the given gene.

## 8. References

1. Kel A, Voss N, Jauregui R, Kel-Margoulis O, Wingender E. Beyond microarrays: Finding key transcription factors controlling signal transduction pathways. *BMC Bioinformatics*. **2006**;7(S2), S13. doi:10.1186/1471-2105-7-s2-s13
2. Stegmaier P, Voss N, Meier T, Kel A, Wingender E, Borlak J. Advanced Computational Biology Methods Identify Molecular Switches for Malignancy in an EGF Mouse Model of Liver Cancer. *PLoS ONE*. **2011**;6(3):e17738. doi:10.1371/journal.pone.0017738
3. Koschmann J, Bhar A, Stegmaier P, Kel A, Wingender E. "Upstream Analysis": An Integrated Promoter-Pathway Analysis Approach to Causal Interpretation of Microarray Data. *Microarrays*. **2015**;4(2):270-286. doi:10.3390/microarrays4020270.

4. Kel A, Stegmaier P, Valeev T, Koschmann J, Poroikov V, Kel-Margoulis OV, and Wingender E. Multi-omics “upstream analysis” of regulatory genomic regions helps identifying targets against methotrexate resistance of colon cancer. *EuPA Open Proteom.* **2016**;13:1-13. doi:10.1016/j.euprot.2016.09.002
5. Michael H, Hogan J, Kel A et al. Building a knowledge base for systems pathology. *Brief Bioinformatics.* **2008**;9(6):518-531. doi:10.1093/bib/bbn038
6. Matys V, Kel-Margoulis OV, Fricke E, Liebich I, Land S, Barre-Dirrie A, Reuter I, Chekmenev D, Krull M, Hornischer K, Voss N, Stegmaier P, Lewicki-Potapov B, Saxel H, Kel AE, Wingender E. TRANSFAC and its module TRANSCompel: transcriptional gene regulation in eukaryotes. *Nucleic Acids Res.* **2006**;34(90001):D108-D110. doi:10.1093/nar/gkj143
7. Kel AE, Gössling E, Reuter I, Cheremushkin E, Kel-Margoulis OV, Wingender E. MATCH: A tool for searching transcription factor binding sites in DNA sequences. *Nucleic Acids Res.* **2003**;31(13):3576-3579. doi:10.1093/nar/gkg585
8. Waleev T, Shtokalo D, Konovalova T, Voss N, Cheremushkin E, Stegmaier P, Kel-Margoulis O, Wingender E, Kel A. Composite Module Analyst: identification of transcription factor binding site combinations using genetic algorithm. *Nucleic Acids Res.* **2006**;34(Web Server issue):W541-5.
9. Krull M, Pistor S, Voss N, Kel A, Reuter I, Kronenberg D, Michael H, Schwarzer K, Potapov A, Choi C, Kel-Margoulis O, Wingender E. TRANSPATH: an information resource for storing and visualizing signaling pathways and their pathological aberrations. *Nucleic Acids Res.* **2006**;34(90001):D546-D551. doi:10.1093/nar/gkj107
10. Boyarskikh U, Pintus S, Mandrik N, Stelmashenko D, Kiselev I, Evshin I, Sharipov R, Stegmaier P, Kolpakov F, Filipenko M, Kel A. Computational master-regulator search reveals mTOR and PI3K pathways responsible for low sensitivity of NCI-H292 and A427 lung cancer cell lines to cytotoxic action of p53 activator Nutlin-3. *BMC Med Genomics.* **2018**;11(1):12. doi:10.1186/1471-2105-7-s2-s13
11. Kel, A., Boyarskikh, U., Stegmaier, P., Leskov, L.S., Sokolov, A.V., Yevshin, I., Mandrik, N., Stelmashenko, D., Koschmann, J., Kel-Margoulis, O. and Krull, M. Walking pathways with positive feedback loops reveal DNA methylation biomarkers of colorectal cancer. *BMC bioinformatics.* Cambridge (UK): RSC Publishing. **2019**;20(Suppl 4):119:1-20. doi:10.1186/s12859-019-2687-7
12. Filimonov D, Poroikov V. Probabilistic Approaches in Activity Prediction. Varnek A, Tropsha A. *Cheminformatics Approaches to Virtual Screening.* Cambridge (UK): RSC Publishing. **2008**;:182-216.
13. Filimonov DA, Poroikov VV. Prognosis of specters of biological activity of organic molecules. *Russian chemical journal.* **2006**;50(2):66-75 (russ)
14. Filimonov D, Poroikov V, Borodina Y, Glorizova T. Chemical Similarity Assessment Through Multilevel Neighborhoods of Atoms: Definition and Comparison with the Other Descriptors. *ChemInform.* **1999**;39(4):666-670. doi:10.1002/chin.199940210

## Thank you for using the Genome Enhancer!

In case of any questions please contact us at [support@genexplain.com](mailto:support@genexplain.com)

## Supplementary material

1. [Supplementary table 1 - Detailed report. Composite modules and master regulators \(high expressed genes in Experiment\).](#)
2. [Supplementary table 2 - Detailed report. Composite modules and master regulators \(low expressed genes in Experiment\).](#)
3. [Supplementary table 3 - Detailed report. Pharmaceutical compounds and drug targets.](#)

## Disclaimer

Decisions regarding care and treatment of patients should be fully made by attending doctors. The predicted chemical compounds listed in the report are given only for doctor’s consideration and they cannot be treated as prescribed medication. It is the physician’s responsibility to independently decide whether any, none or all of the predicted compounds can be used solely or in combination for patient treatment purposes, taking into account all applicable information regarding FDA prescribing recommendations for any therapeutic and the patient’s condition, including, but not limited to, the patient’s and family’s medical history, physical examinations, information from various diagnostic

tests, and patient preferences in accordance with the current standard of care. Whether or not a particular patient will benefit from a selected therapy is based on many factors and can vary significantly.

The compounds predicted to be active against the identified drug targets in the report are not guaranteed to be active against any particular patient's condition. GeneXplain GmbH does not give any assurances or guarantees regarding the treatment information and conclusions given in the report. There is no guarantee that any third party will provide a refund for any of the treatment decisions made based on these results. None of the listed compounds was checked by Genome Enhancer for adverse side-effects or even toxic effects.

The analysis report contains information about chemical drug compounds, clinical trials and disease biomarkers retrieved from the HumanPSD™ database of gene-disease assignments maintained and exclusively distributed worldwide by geneXplain GmbH. The information contained in this database is collected from scientific literature and public clinical trials resources. It is updated to the best of geneXplain's knowledge however we do not guarantee completeness and reliability of this information leaving the final checkup and consideration of the predicted therapies to the medical doctor.

The scientific analysis underlying the Genome Enhancer report employs a complex analysis pipeline which uses geneXplain's proprietary Upstream Analysis approach, integrated with TRANSFAC® and TRANSPATH® databases maintained and exclusively distributed worldwide by geneXplain GmbH. The pipeline and the databases are updated to the best of geneXplain's knowledge and belief, however, geneXplain GmbH shall not give a warranty as to the characteristics or to the content and any of the results produced by Genome Enhancer. Moreover, any warranty concerning the completeness, up-to-dateness, correctness and usability of Genome Enhancer information and results produced by it, shall be excluded.

The results produced by Genome Enhancer, including the analysis report, severely depend on the quality of input data used for the analysis. It is the responsibility of Genome Enhancer users to check the input data quality and parameters used for running the Genome Enhancer pipeline.

Note that the text given in the report is not unique and can be fully or partially repeated in other Genome Enhancer analysis reports, including reports of other users. This should be considered when publishing any results or excerpts from the report. This restriction refers only to the general description of analysis methods used for generating the report. All data and graphics referring to the concrete set of input data, including lists of mutated genes, differentially expressed genes/proteins/metabolites, functional classifications, identified transcription factors and master regulators, constructed molecular networks, lists of chemical compounds and reconstructed model of molecular mechanisms of the studied pathology are unique in respect to the used input data set and Genome Enhancer pipeline parameters used for the current run.