

PDGFRA and DUSP4 are promising druggable targets for treating Ovarian Neoplasms that control activity of E2F1, HSF2 and ELK1 transcription factors on of differentially expressed genes in ovary tissue

Demo User

geneXplain GmbH

info@genexplain.com

Data received on 14/08/2019 ; Run on 11/06/2020 ; Report generated on 11/06/2020

Genome Enhancer release 2.0 (TRANSFAC®, TRANSPATH® and HumanPSD™ release 2020.2)



Abstract

In the present study we applied the software package "Genome Enhancer" to a multiomics data set that contains *transcriptomics* and *epigenomics* data obtained from *ovary* tissue. The study is done in the context of *Ovarian Neoplasms*. 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: E2F1, HSF2, EGR1, ELK1 and HSF1. The subsequent network analysis suggested

- PP1-gamma1
- PDGFRalpha
- MKP-2
- Chk2

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: Pazopanib, Vitamin E, Paclitaxel and 9-Aminophenanthrene.

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) reconstructing 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] 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 2507 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 [11-13]. 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
GSM385721.CEL	Transcriptomics
GSM385722.CEL	Transcriptomics
GSM385723.CEL	Transcriptomics
GSM385724.CEL	Transcriptomics
GSM385725.CEL	Transcriptomics
GSM385726.CEL	Transcriptomics
GSM385727.CEL	Transcriptomics
GSM385728.CEL	Transcriptomics
GSM385729.CEL	Transcriptomics
GSM385730.CEL	Transcriptomics
GSM385747_CpG_NM.fixed.hg38.top300	Epigenomics

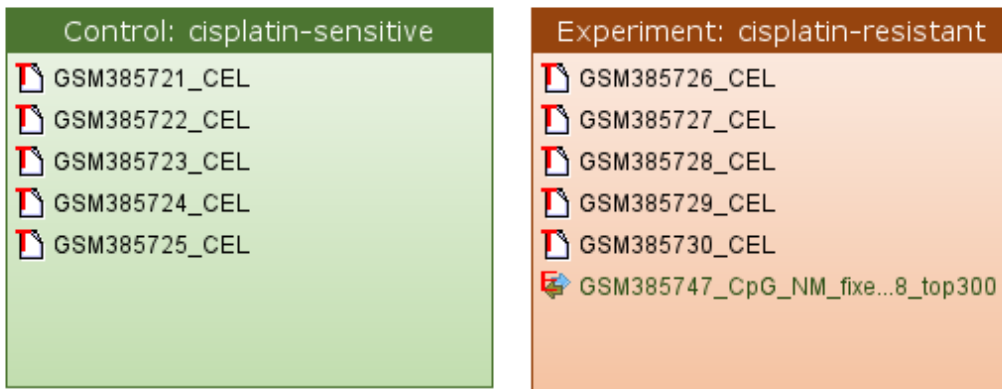


Figure 1. Annotation diagram of experimental data used in this study. With the colored boxes we show those sub-categories of the data that are compared in our analysis.

3. Results

We have compared the following conditions: Experiment: cisplatin-resistant *versus* Control: cisplatin-sensitive.

3.1. Identification of target genes

In the first step of the analysis **target genes** were identified from the uploaded experimental data. We applied the Limma tool (R/Bioconductor package integrated into our pipeline) and compared gene expression in the following sets: "Experiment: cisplatin-resistant" with "Control: cisplatin-sensitive". Limma calculated the LogFC (the logarithm to the base 2 of the fold change between different conditions), the p-value and the adjusted p-value (corrected for multiple testing) of the observed fold change. As a result, we detected 12732 upregulated genes (LogFC>0) out of which 8575 genes were found as significantly upregulated (p-value<0.1) and 12588 downregulated genes (LogFC<0) out of which 8399 genes were significantly downregulated (p-value<0.1). See tables below for the top significantly up- and

downregulated genes. Below we call **target genes** the full list of up- and downregulated genes revealed in our analysis (see tables in [Supplementary section](#)).

Table 2. Top ten significant **up-regulated** genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive.

[See full table](#) →

ID	Gene symbol	Gene description	logFC	P.Value	adj.P.Val
ENSG00000123700	KCNJ2	potassium inwardly rectifying channel subfamily J member 2	5.31	2.04E-15	3.7E-12
ENSG00000064218	DMRT3	doublesex and mab-3 related transcription factor 3	5.17	2.28E-16	9.89E-13
ENSG00000099139	PCSK5	proprotein convertase subtilisin/kexin type 5	4.46	8.64E-13	3.21E-10
ENSG00000196507	TCEAL3	transcription elongation factor A like 3	3.98	4.73E-16	1.33E-12
ENSG00000197705	KLHL14	kelch like family member 14	3.67	4.28E-15	5.95E-12
ENSG00000103449	SALL1	spalt like transcription factor 1	3.4	4.01E-12	9.14E-10
ENSG00000138378	STAT4	signal transducer and activator of transcription 4	3.39	8.94E-12	1.8E-9
ENSG00000164692	COL1A2	collagen type I alpha 2 chain	3.29	7.01E-15	7.4E-12
ENSG00000133083	DCLK1	doublecortin like kinase 1	3.29	6E-15	6.9E-12
ENSG00000126950	TMEM35A	transmembrane protein 35A	3.16	3.42E-15	5.09E-12

Table 4. Top ten significant **down-regulated** genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive.

[See full table](#) →

ID	Gene symbol	Gene description	logFC	P.Value	adj.P.Val
ENSG00000149968	MMP3	matrix metalloproteinase 3	-6.61	1.64E-18	3.42E-14
ENSG00000127324	TSPAN8	tetraspanin 8	-6.08	1.76E-14	1.59E-11
ENSG00000139292	LGR5	leucine rich repeat containing G protein-coupled receptor 5	-5.53	1.28E-16	8.09E-13
ENSG00000153233	PTPRR	protein tyrosine phosphatase receptor type R	-5.29	2.34E-16	9.89E-13
ENSG00000169908	TM4SF1	transmembrane 4 L six family member 1	-4.66	2.7E-18	3.42E-14
ENSG00000106511	MEOX2	mesenchyme homeobox 2	-4.63	9.68E-16	2.45E-12
ENSG00000163359	COL6A3	collagen type VI alpha 3 chain	-4.54	1.66E-17	1.4E-13
ENSG00000060718	COL11A1	collagen type XI alpha 1 chain	-4.53	3.04E-14	2.26E-11
ENSG00000166670	MMP10	matrix metalloproteinase 10	-4.29	1.44E-15	3.11E-12
ENSG00000145431	PDGFC	platelet derived growth factor C	-4.09	3.82E-16	1.21E-12

3.2. Regulatory regions of target genes

We mapped the uploaded Epigenomic peaks on the **target genes** and selected those peaks only that were found located in the body of the gene (in exons or introns of the genes) or in the 5000 nucleotide long flanking regions of the genes. In the tables below we demonstrate localization of such potential regulatory regions in the top up-regulated and down-regulated genes.

Table 3. Top 3 **up-regulated** genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive with epigenomic peaks.

[See full table](#) →











ID	Gene symbol	Gene schematic representation
ENSG00000260774	AC021087.3	
ENSG00000027075	PRKCH	
ENSG000000186684	CYP27C1	

Table 5. Top 7 **down-regulated** genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive with epigenomic peaks.

[See full table](#) →

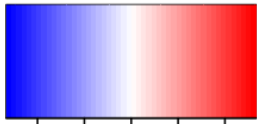
ID	Gene symbol	Gene schematic representation
ENSG000000170558	CDH2	
ENSG000000197921	HES5	
ENSG000000197822	OCLN	
ENSG000000146648	EGFR	
ENSG000000145476	CYP4V2	
ENSG000000237765	FAM200B	
ENSG000000118495	PLAGL1	

3.3. Functional classification of genes

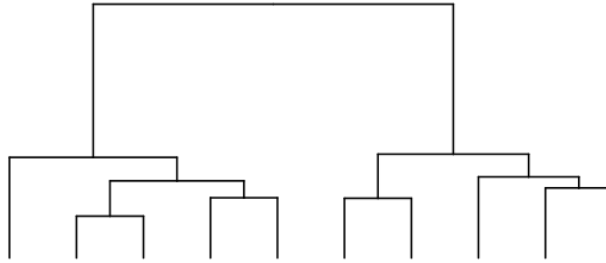
A functional analysis of differentially expressed genes was done by mapping the significant up-regulated and significant down-regulated genes to several known ontologies, such as Gene Ontology (GO), disease ontology (based on HumanPSD™ database) and the ontology of signal transduction and metabolic pathways from the [TRANSPATH®](#) database. Statistical significance was computed using a binomial test. Figures 3-8 show the most significant categories.

Heatmap of differentially expressed genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive

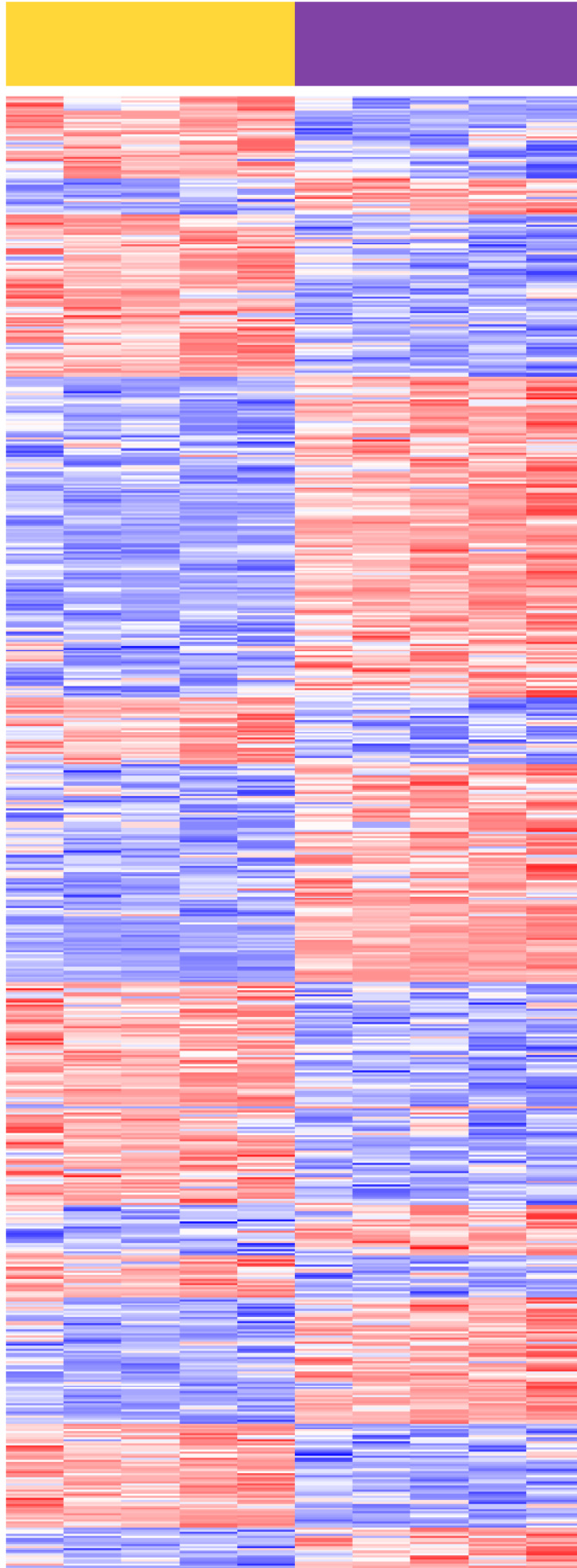
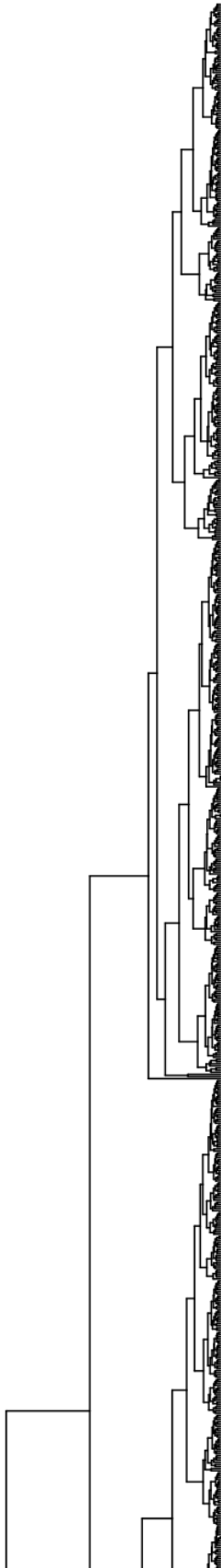
A heatmap of all differentially expressed genes playing a potential regulatory role in the system (enriched in [TRANSPATH®](#) pathways) is presented in Figure 2.



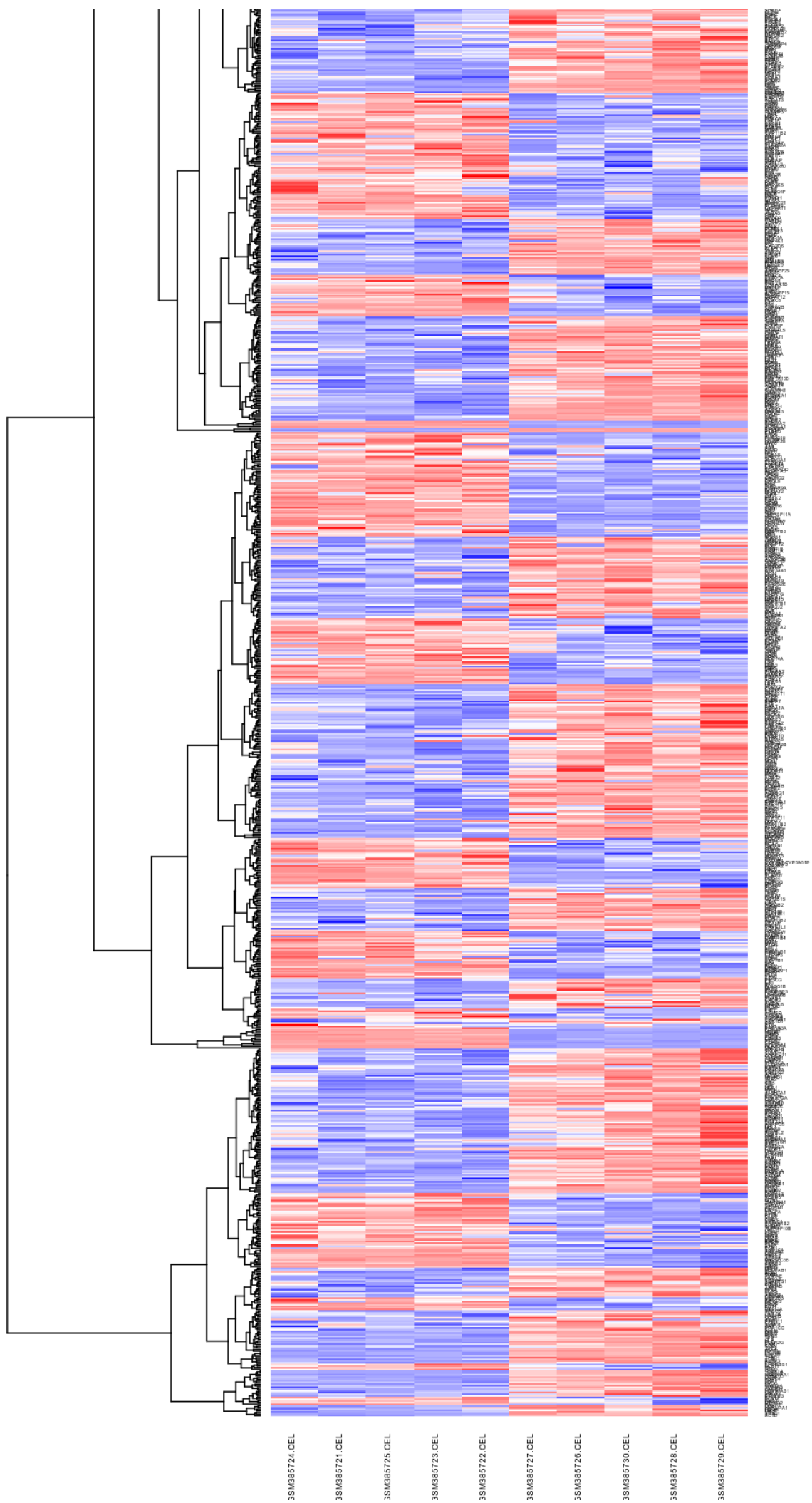
Gene Expression Normalized by rows



Control: cisplatin-sensitive
Experiment: cisplatin-resistant



11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100



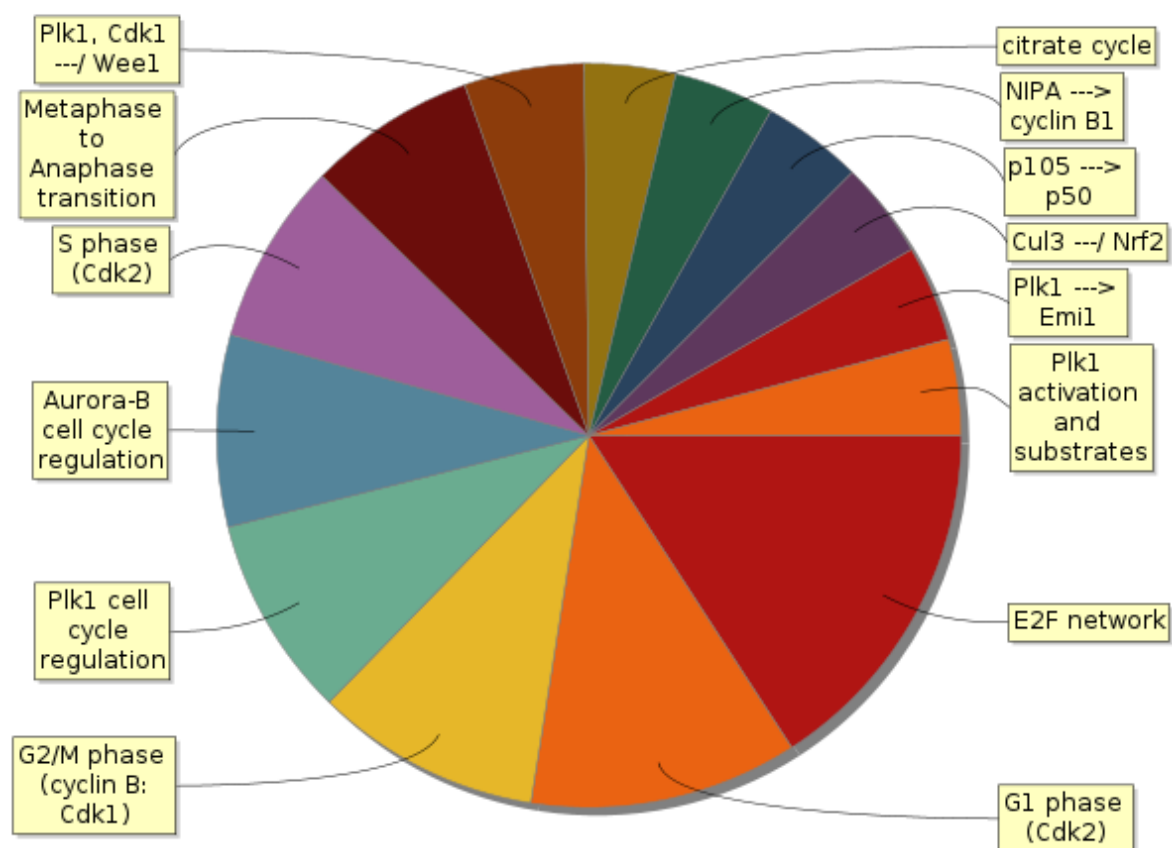


Figure 4. Enriched TRANSPATH® Pathways (2020.2) of up-regulated genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive.

[Full classification →](#)

HumanPSD(TM) disease (2020.2)

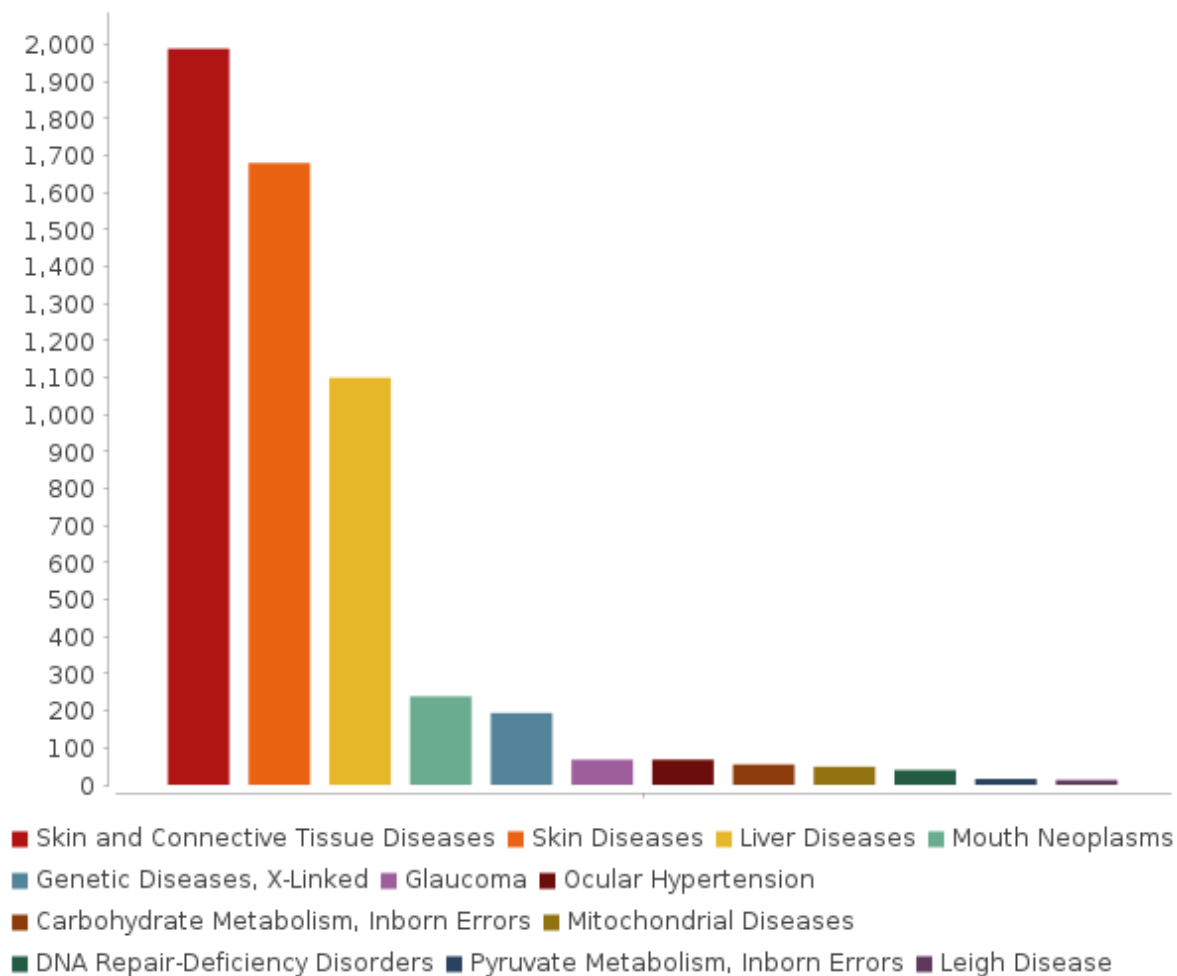


Figure 5. Enriched HumanPSD(TM) disease (2020.2) of up-regulated genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive. The size of the bars correspond to the number of bio-markers of the given disease found among the input set.

[Full classification →](#)

Down-regulated genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive:

8399 significant down-regulated genes were taken for the mapping.

GO (biological process)

[illegible]

Full classification →

TRANSPATH® Pathways (2020.2)

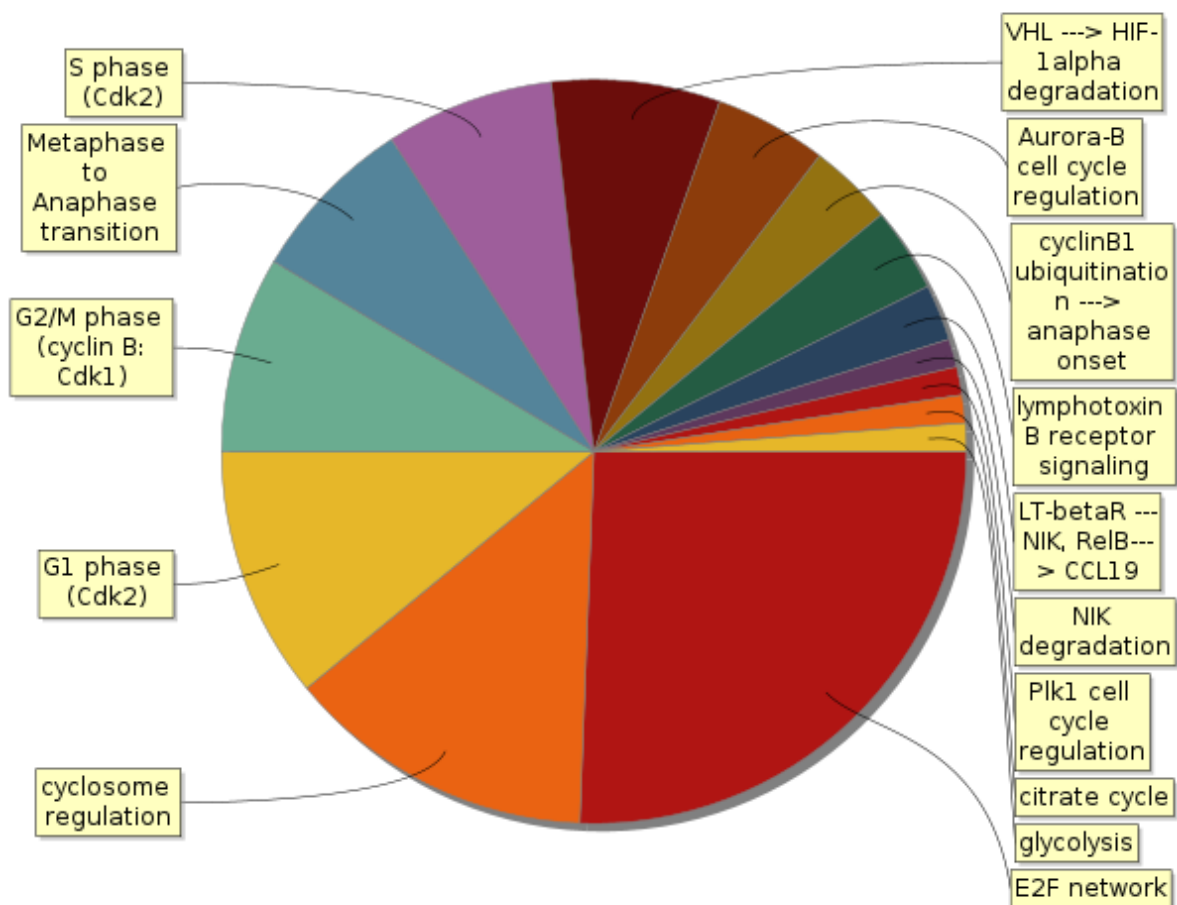


Figure 7. Enriched TRANSPATH® Pathways (2020.2) of down-regulated genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive.

[Full classification →](#)

HumanPSD(TM) disease (2020.2)

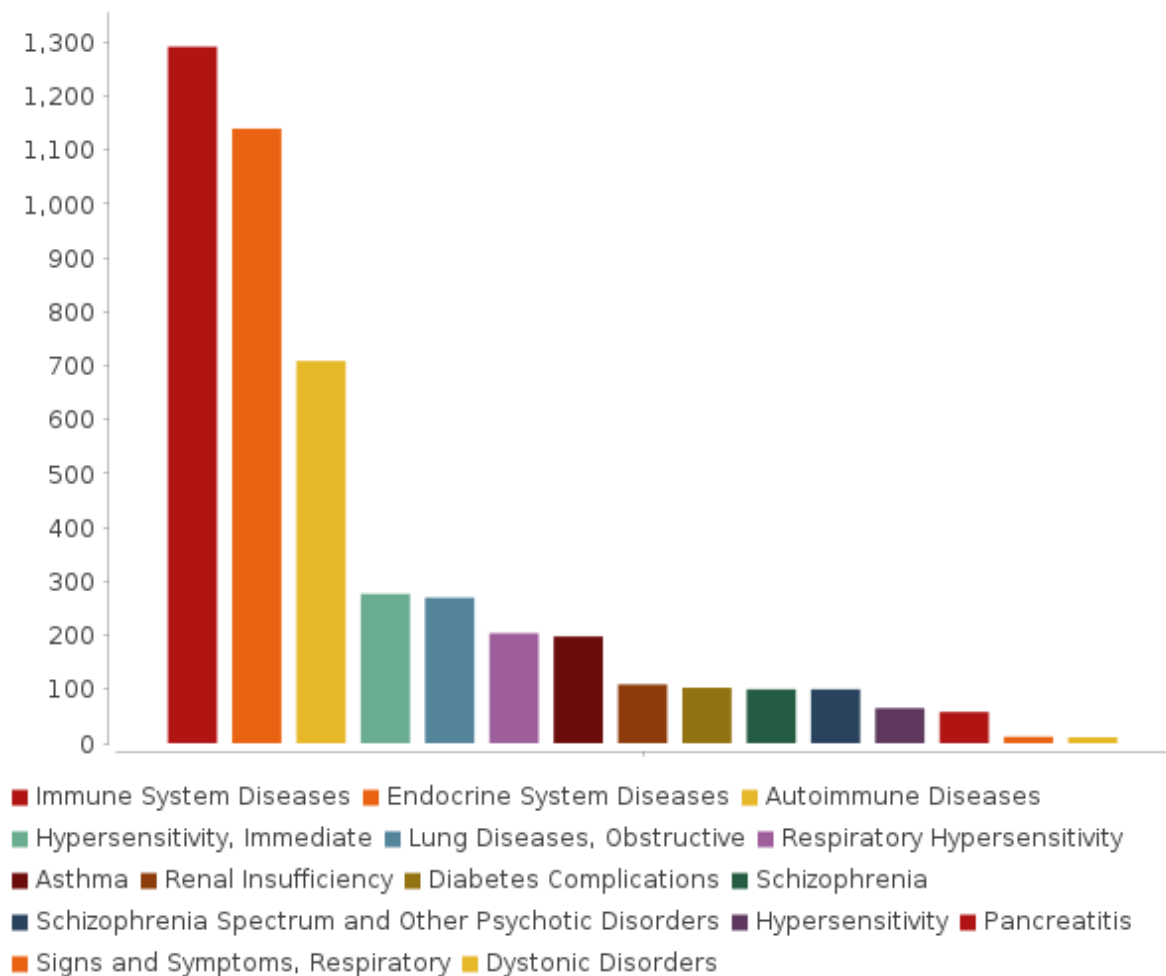
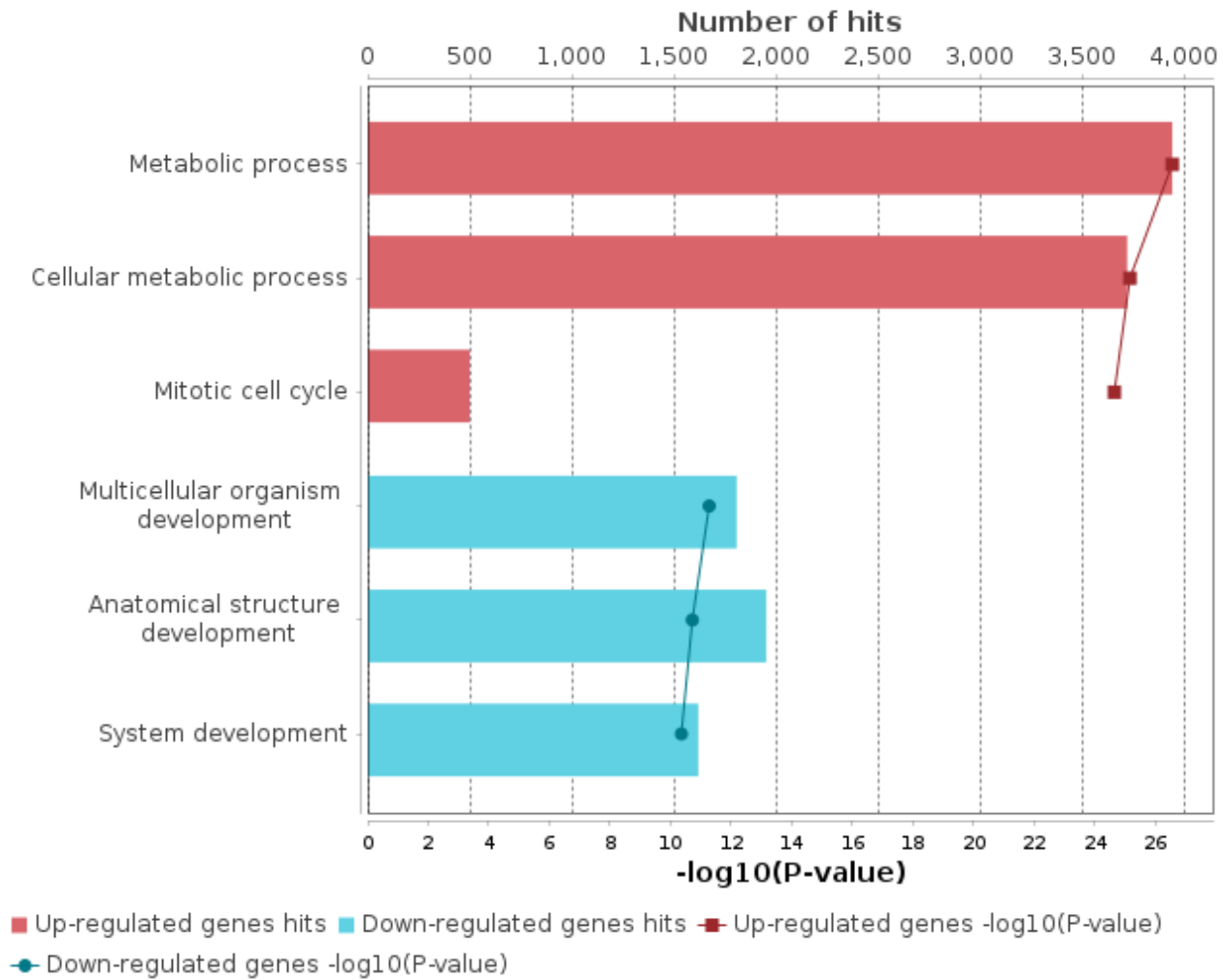


Figure 8. Enriched HumanPSD(TM) disease (2020.2) of down-regulated genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive. The size of the bars correspond to the number of bio-markers of the given disease found among the input set.

[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):



3.4. 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].

In the current work we use the Epigenomics data from the track(s) "GSM385747_CpG_NM.fixed.hg38.top300" to predict positions of potential **enhancers** regulating the differentially expressed genes revealed by comparative transcriptomics analysis. We took genomic regions -550bp upstream and 550bp downstream from the middle point of each interval of the track and check if these regions are located inside the 5kb flanking arias of the differentially expressed genes (or inside the body of the genes). In such cases, these genomic regions are used for the search for potential condition-specific enhancers. In all other cases when the differentially expressed genes did not contain epigenomic peaks in their body or in the 5kb flanking regions we used the upstream regulatory regions of these genes (-1000bp upstream and 100bp downstream of TSS) for the search for condition-specific enhancers.

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 (up-regulated genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive).

To build the most specific composite modules we choose genes as the input of CMA algorithm. The obtained CMA model is then applied to compute CMA score for all up-regulated genes.

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.

Module 1:

V\$NFYA_Q3
0.00; N=3

V\$POU6F1_Q1
0.78; N=3

V\$MTF1_Q5
0.00; N=2

V\$HSF1_Q6_Q1
0.00; N=2

V\$E2F1_Q9
0.97; N=2

V\$E2F4_Q3
0.00; N=2

V\$MZF1_Q5_Q1
0.00; N=3

Module width: 177

Module 2:

V\$NFAT4_Q5
0.00; N=2

V\$EGR1_Q6
0.00; N=3

V\$HSF2_Q1
0.99; N=3

V\$HMG1Y_Q1
0.95; N=2

V\$HLF_Q1
0.84; N=1

V\$BRCA_Q1
0.00; N=2

Module width: 125

Model score ($-\log_{10}(pval)$): 15.65

Wilcoxon p-value (pval): 1.64e-34

Penalty (p): 0.463

Average yes-set score: 7.93

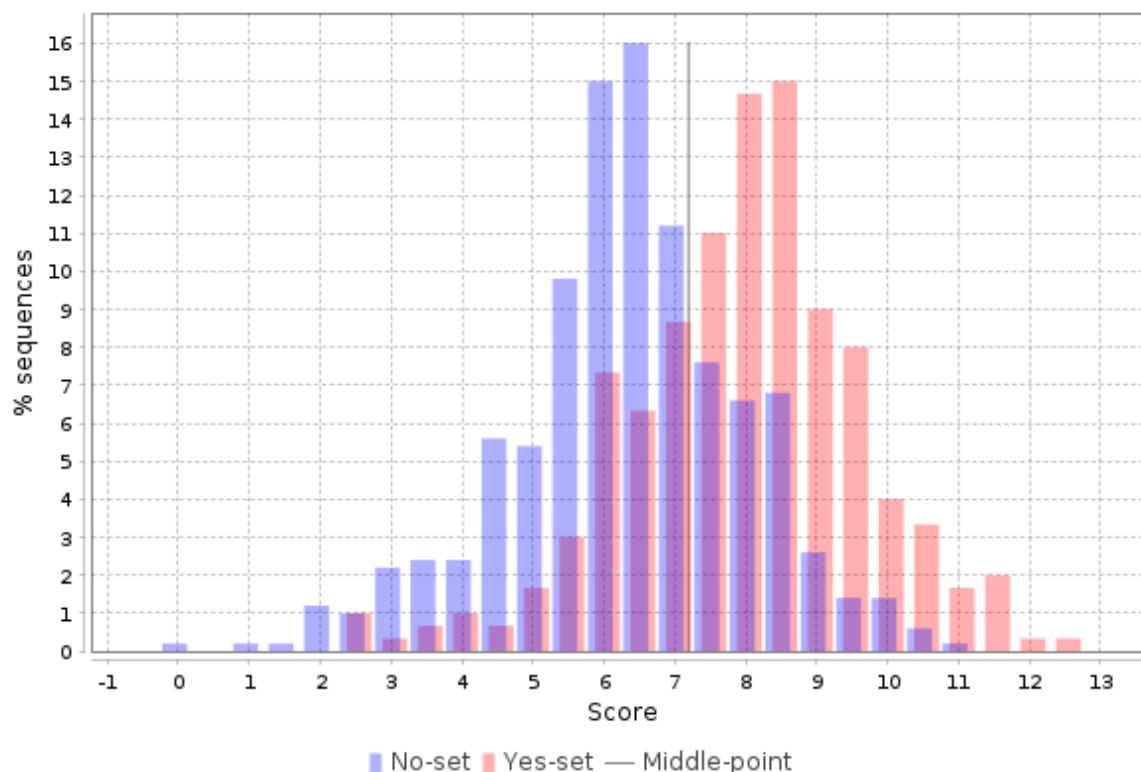
Average no-set score: 6.35

AUC: 0.76

Middle-point: 7.19

False-positive: 28.20%

False-negative: 28.67%



[See model visualization table](#) →

Table 6. List of top ten up-regulated genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive 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
ENSG00000147123	NDUFB11	NADH:ubiquinone oxidoreductase subunit B11	18.41	MTF-1(h), Egr-1(h), NFATc3(h), brca1(h), HSF2(h), HMGIY(h), NF-YA(h)...
ENSG00000161914	ZNF653	zinc finger protein 653	18.39	Egr-1(h), MZF-1(h), brca1(h), MTF-1(h), HSF2(h), HMGIY(h), E2F-4(h)...
ENSG00000225670	CADM3-AS1	CADM3 antisense RNA 1	18.29	MTF-1(h), NF-YA(h), E2F-4(h), Hlf(h), POU6F1(h), brca1(h), HSF1(h)...
ENSG00000110944	IL23A	interleukin 23 subunit alpha	18.21	Hlf(h), Egr-1(h), MZF-1(h), MTF-1(h), E2F-4(h), HSF1(h), HMGIY(h)...
ENSG00000135473	PAN2	poly(A) specific ribonuclease subunit PAN2	18.21	Egr-1(h), MZF-1(h), MTF-1(h), E2F-4(h), HSF1(h), HMGIY(h), NFATc3(h)...
ENSG00000145569	OTULINL	OTU deubiquitinase with linear linkage specificity like	18.18	Hlf(h), POU6F1(h), Egr-1(h), brca1(h), MTF-1(h), MZF-1(h), NFATc3(h)...
ENSG00000272269	AL138724.2	novel transcript, antisense to NUP153	18.03	NF-YA(h), HSF2(h), MTF-1(h), MZF-1(h), NFATc3(h), HSF1(h), HMGIY(h)...
ENSG00000149926	TLCD3B	TLC domain containing 3B	17.95	E2F-4(h), HSF1(h), brca1(h), Egr-1(h), MZF-1(h), E2F-1(h), NF-YA(h)...
ENSG00000102409	BEX4	brain expressed X-linked 4	17.87	NFATc3(h), Egr-1(h), NF-YA(h), E2F-4(h), POU6F1(h), HSF1(h), HMGIY(h)...
ENSG00000129347	KRI1	KRI1 homolog	17.85	E2F-4(h), HMGIY(h), Hlf(h), NFATc3(h), NF-YA(h), brca1(h), POU6F1(h)...

Enhancer model potentially involved in regulation of target genes (down-regulated genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive).

To build the most specific composite modules we choose genes as the input of CMA algorithm. The obtained CMA model is then applied to compute CMA score for all down-regulated genes.

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.

Module 1:

V\$SMAD5_Q5
0.00; N=3

V\$E2F1_Q3_Q1
0.82; N=3

V\$HMGIY_Q1
0.97; N=3

V\$SIRT6_Q1
0.00; N=3

V\$E2F1_Q1
0.97; N=2

V\$NFAT2_Q4
0.00; N=2

V\$ELK1_Q4
0.00; N=2

Module width: 154

Module 2:

V\$SOX5_Q5
0.00; N=3

V\$IRF7_Q1
0.80; N=2

V\$EGR1_Q6
0.85; N=2

V\$TAF1_Q7
0.98; N=2

V\$CDP_Q3
0.00; N=2

V\$HSF1_Q6
0.78; N=3

Module width: 164

Model score ($-\log_{10}(pval)$): 16.66

Wilcoxon p-value (pval): 1.07e-36

Penalty (p): 0.463

Average yes-set score: 11.68

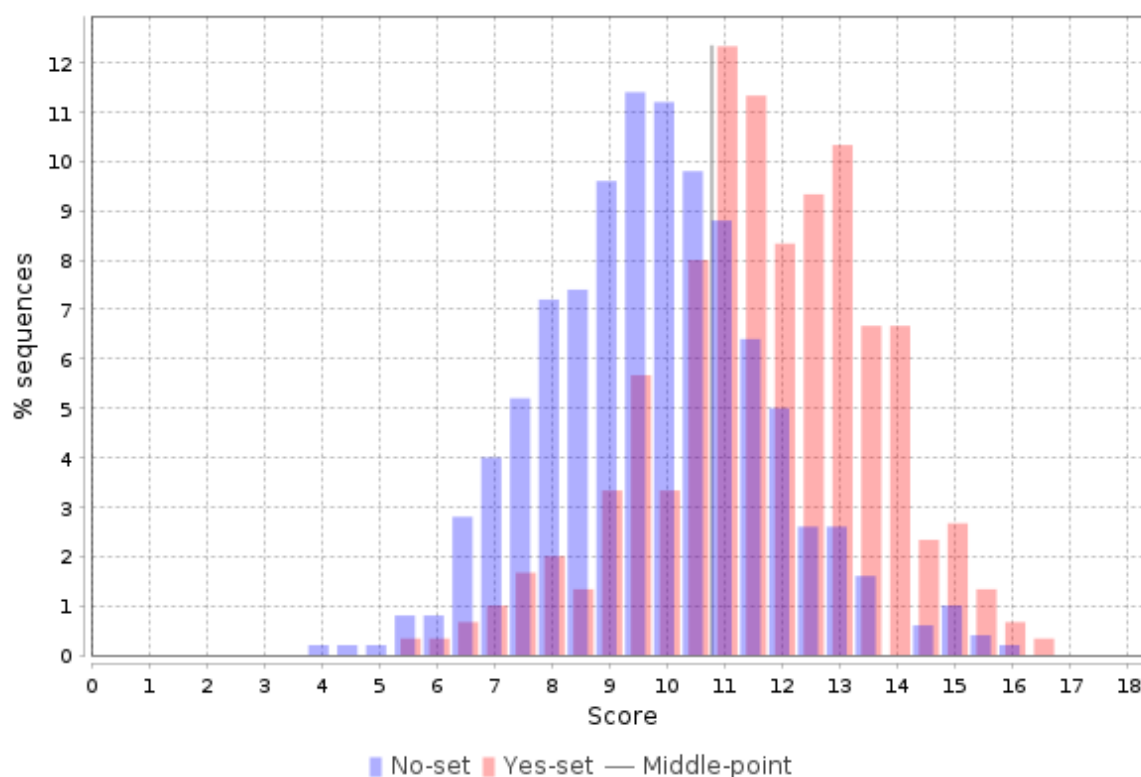
Average no-set score: 9.78

AUC: 0.77

Middle-point: 10.78

False-positive: 28.80%

False-negative: 27.67%



[See model visualization table](#) →

Table 7. List of top ten down-regulated genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive 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
ENSG00000139083	ETV6	ETS variant transcription factor 6	21.41	SIR2L6(h), Sox-5(h), TAFII250(h), Egr-1(h), Smad5(h), E2F-1(h), NFATc1(h)...
ENSG00000225470	JPX	JPX transcript, XIST activator	20.71	Smad5(h), E2F-1(h), Sox-5(h), HSF1(h), CDP(h), HMGIY(h), IRF-7(h)...
ENSG00000155052	CNTNAP5	contactin associated protein like 5	20.55	HSF1(h), Sox-5(h), CDP(h), IRF-7(h), Egr-1(h), Smad5(h), HMGIY(h)...
ENSG00000135083	CCNJL	cyclin J like	20.32	Smad5(h), TAFII250(h), E2F-1(h), Egr-1(h), NFATc1(h), SIR2L6(h), HSF1(h)...
ENSG00000169439	SDC2	syndecan 2	20.09	HMGIY(h), CDP(h), HSF1(h), E2F-1(h), IRF-7(h), NFATc1(h), SIR2L6(h)...
ENSG00000180828	BHLHE22	basic helix-loop-helix family member e22	20.05	Elk-1(h), Sox-5(h), HMGIY(h), Smad5(h), SIR2L6(h), NFATc1(h), E2F-1(h)...
ENSG00000178177	LCORL	ligand dependent nuclear receptor corepressor like	20.01	CDP(h), HSF1(h), IRF-7(h), E2F-1(h), Sox-5(h), SIR2L6(h), Egr-1(h)...
ENSG00000135926	TMBIM1	transmembrane BAX inhibitor motif containing 1	20.01	Egr-1(h), Elk-1(h), Smad5(h), E2F-1(h), TAFII250(h), HMGIY(h), SIR2L6(h)...
ENSG00000108018	SORCS1	sortilin related VPS10 domain containing receptor 1	19.69	E2F-1(h), Egr-1(h), TAFII250(h), IRF-7(h), Sox-5(h), NFATc1(h), Elk-1(h)...
ENSG00000182446	NPLOC4	NPL4 homolog, ubiquitin recognition factor	19.69	TAFII250(h), Egr-1(h), E2F-1(h), Smad5(h), IRF-7(h), NFATc1(h), Elk-1(h)...

On the basis of the enhancer models we identified transcription factors potentially regulating the **target genes** of our interest. We found 13 and 12 transcription factors controlling expression of up- and down-regulated genes respectively (see Tables 8-9).

Table 8. Transcription factors of the predicted enhancer model potentially regulating the differentially expressed genes (up-regulated genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive). **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](#) →

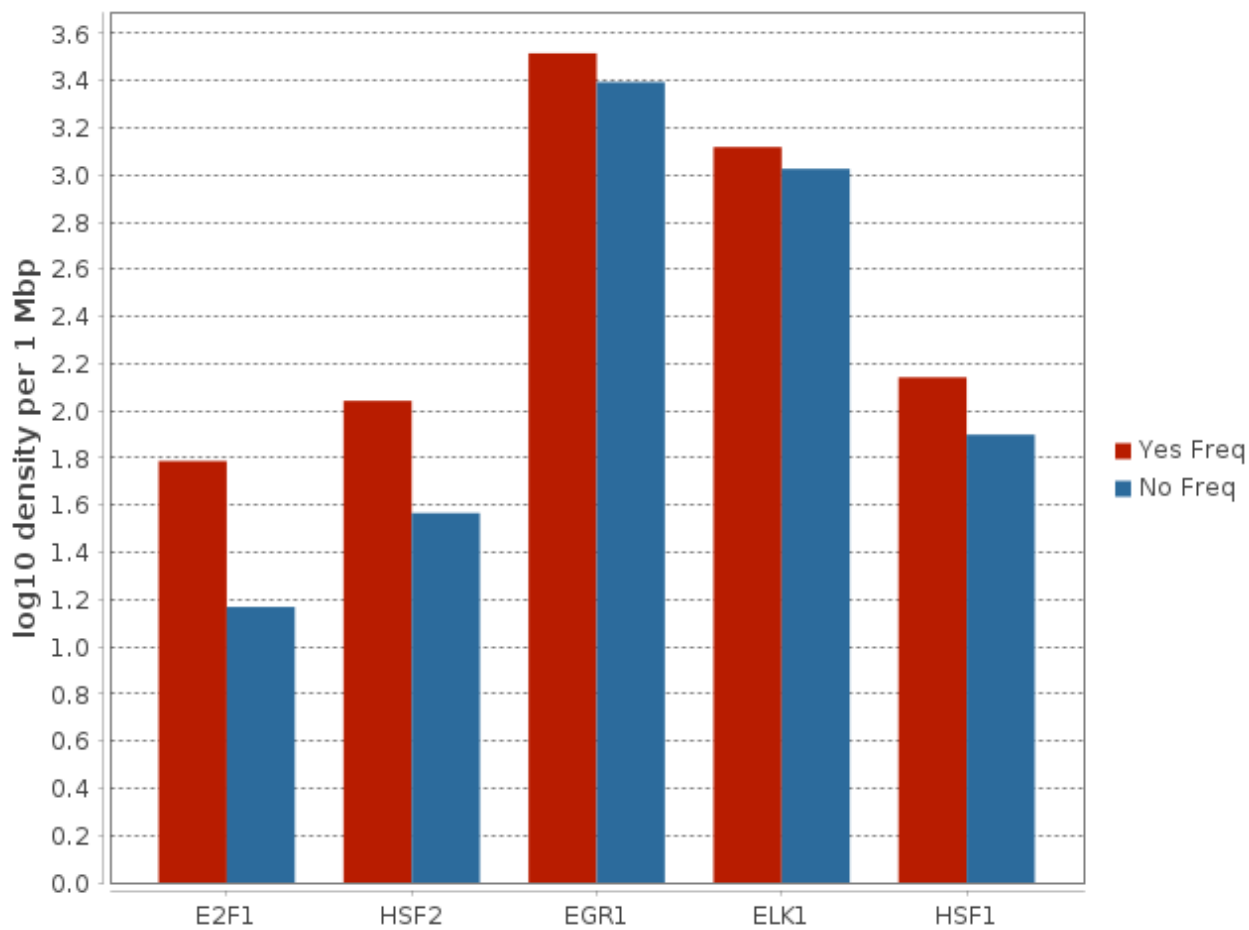
ID	Gene symbol	Gene description	Regulatory score	Yes-No ratio
MO000004274	E2F1	E2F transcription factor 1	5.94	4.15
MO000046011	HSF2	heat shock transcription factor 2	5.44	2.99
MO000017914	EGR1	early growth response 1	5.16	1.32
MO000033378	HSF1	heat shock transcription factor 1	5.09	4.06
MO000021981	BRCA1	BRCA1 DNA repair associated	4.98	1.54
MO000020739	NFATC3	nuclear factor of activated T cells 3	4.66	1.27
MO000025939	NFYA	nuclear transcription factor Y subunit alpha	4.63	1.43
MO000026358	HMGA1	high mobility group AT-hook 1	4.42	1.22
MO000023603	E2F4	E2F transcription factor 4	4.22	1.62
MO000028320	null	null	3.66	1.75

Table 9. Transcription factors of the predicted enhancer model potentially regulating the differentially expressed genes (down-regulated genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive). **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
MO000019544	ELK1	ETS transcription factor ELK1	4.76	1.24
MO000017914	EGR1	early growth response 1	4.71	1.27
MO000033378	HSF1	heat shock transcription factor 1	4.68	1.75
MO000020760	NFATC1	nuclear factor of activated T cells 1	4.58	1.3
MO000004274	E2F1	E2F transcription factor 1	4.41	2.75
MO000007703	IRF7	interferon regulatory factor 7	3.97	1.29
MO000026358	HMGA1	high mobility group AT-hook 1	3.94	1.2
MO000020635	SMAD5	SMAD family member 5	3.88	2.07
MO000024708	CUX1	cut like homeobox 1	3.62	1.28
MO000142283	SIRT6	sirtuin 6	3.53	1.24

The following diagram represents the key transcription factors, which were predicted to be potentially regulating differentially expressed genes in the analyzed pathology: E2F1, HSF2, EGR1, ELK1 and HSF1.



3.5. 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 10-11.

Table 10. Master regulators that may govern the regulation of **up-regulated** genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive. **Total rank** is the sum of the ranks of the master molecules sorted by keynode score, CMA score, transcriptomics and epigenomics data.

[See full table](#) →

ID	Master molecule name	Gene symbol	Gene description	logFC	Total rank
MO000030895	Chk2(h)	CHEK2	checkpoint kinase 2	0.88	81
MO000081890	Chk2-isoform1(h)	CHEK2	checkpoint kinase 2	0.88	115
MO000032652	MKP-2(h)	DUSP4	dual specificity phosphatase 4	1.17	121
MO000021981	brca1(h)	BRCA1	BRCA1 DNA repair associated	0.88	122
MO000081925	Chk2-xbb12(h)	CHEK2	checkpoint kinase 2	0.88	141
MO000019376	Cot(h)	MAP3K8	mitogen-activated protein kinase kinase kinase 8	1.87	152
MO000031112	Chk2(h){pT68}	CHEK2	checkpoint kinase 2	0.88	155
MO000030927	DNA-PKcs(h)	PRKDC	protein kinase, DNA-activated, catalytic subunit	0.58	246
MO000162677	PHLPP(h)	PHLPP1	PH domain and leucine rich repeat protein phosphatase 1	0.78	261
MO000010977	PDGFRalpha(h)	PDGFRA	platelet derived growth factor receptor alpha	2.93	300

Table 11. Master regulators that may govern the regulation of **down-regulated** genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive. **Total rank** is the sum of the ranks of the master molecules sorted by keynode score, CMA score, transcriptomics and epigenomics data.

[See full table](#) →

ID	Master molecule name	Gene symbol	Gene description	logFC	Total rank
MO000022222	MKP-1(h)	DUSP1	dual specificity phosphatase 1	-1.22	47
MO000083769	MKP-1(h)	DUSP1	dual specificity phosphatase 1	-1.22	148
MO000042839	ptpn21(h)	PTPN21	protein tyrosine phosphatase non-receptor type 21	-1.32	271
MO000039099	IL-1beta-p17:IL-1RI:IL-1RAcP:MyD88:tollip:IRAK-1{pS376}{pT387}:IRAK-4:IRAK-2	AC093012.1, IL1B, IL1R1, IL1RAP, IRAK1, IRAK2, MYD88, TOLLIP	MYD88 innate immune signal transduction adaptor, interleukin 1 beta, interleukin 1 receptor accessor...	-0.97	286
MO000031202	Cdc14A(h)	CDC14A	cell division cycle 14A	-0.49	357
MO000081777	Pellino2(h)	PELI2	pellino E3 ubiquitin protein ligase family member 2	-0.65	372
MO000101468	LRRK2(h)	LRRK2	leucine rich repeat kinase 2	-1.02	375
MO000122463	mTOR(h):ricor(h):mLST8(h):SIN1(h)	MAPKAP1, MLST8, MTOR, RICTOR	MAPK associated protein 1, MTOR associated protein, LST8 homolog, RPTOR independent companion of MTO...	-0.53	401
MO000019375	IL-1beta(h)	IL1B	interleukin 1 beta	-0.64	402
MO000017291	integrins	ITGA1, ITGA2B, ITGA3, ITGA4, ITGA5, ITGA6, ITGA8, ITGA9, ITGAL, ITGAV, ITGB1, ITGB2, ITGB3, ITGB4, I...	integrin subunit alpha 1, integrin subunit alpha 2b, integrin subunit alpha 3, integrin subunit alph...	-1.04	403

The intracellular regulatory pathways controlled by the above-mentioned master regulators are depicted in Figures 9 and 10. 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.

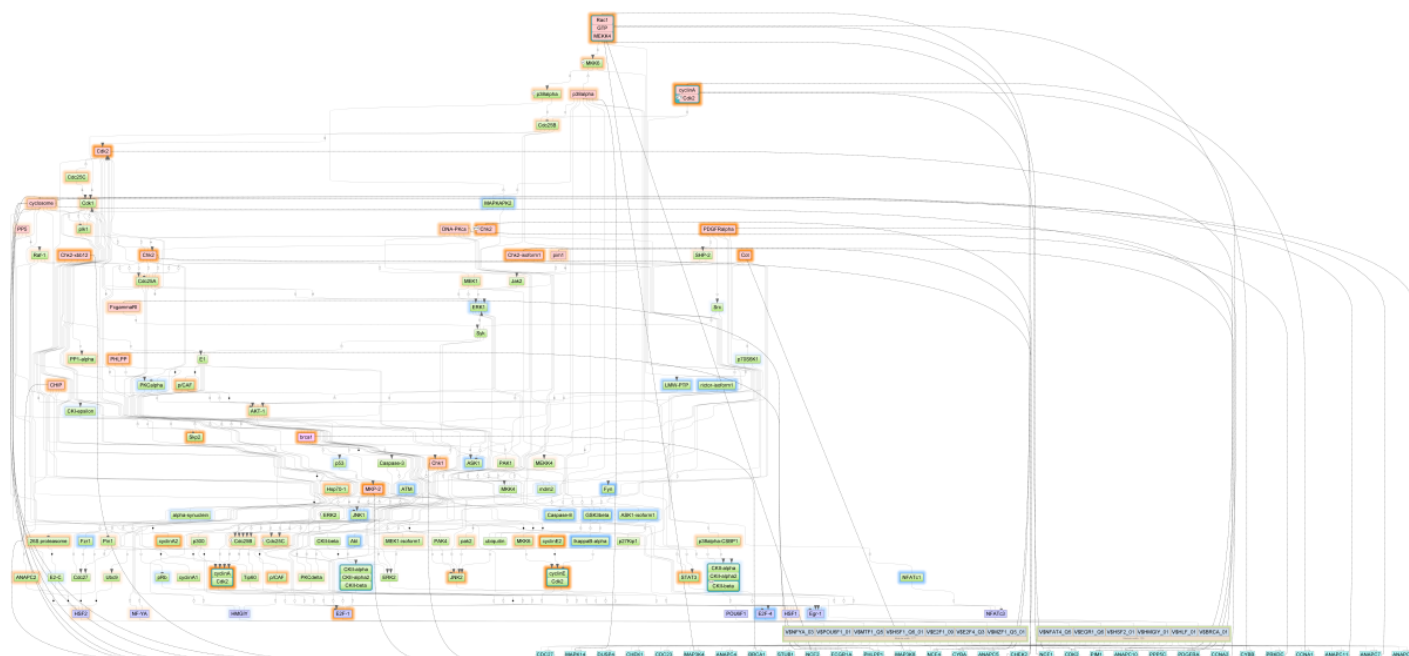


Figure 9. Diagram of intracellular regulatory signal transduction pathways of up-regulated genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive. Master regulators are indicated by red rectangles, transcription factors are blue rectangles, and green rectangles are intermediate molecules, 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 →](#)

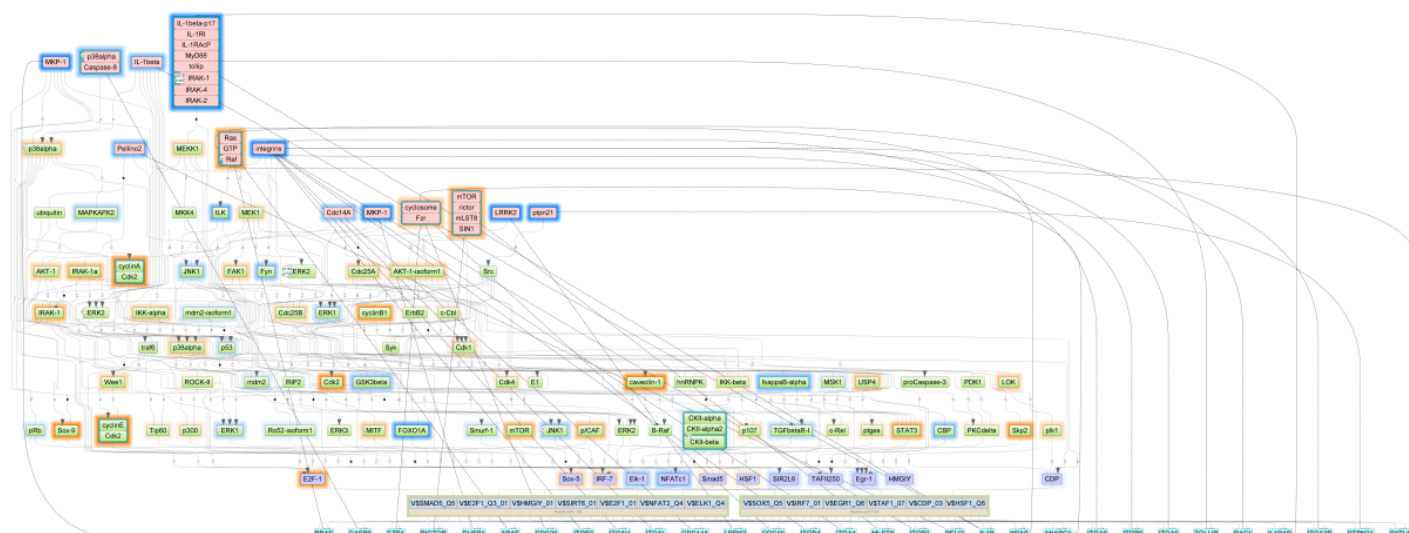


Figure 10. Diagram of intracellular regulatory signal transduction pathways of down-regulated genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive. Master regulators are indicated by red rectangles, transcription factors are blue rectangles, and green rectangles are intermediate 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 [11-13] 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 [11-13] on the basis of a (Q)SAR approach.

If both druggability scores were below defined thresholds (see Method 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 12. 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	logFC	Total rank
PDGFRA	platelet derived growth factor receptor alpha	8	2.93	300
PPP1CC	protein phosphatase 1 catalytic subunit gamma	4	0.41	702
PSMA7	proteasome 20S subunit alpha 7	3	0.44	732
KAT2B	lysine acetyltransferase 2B	3	0.62	811
AURKB	aurora kinase B	3	0.83	816
CSNK1G2	casein kinase 1 gamma 2	3	0.56	877



Table 13. 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.

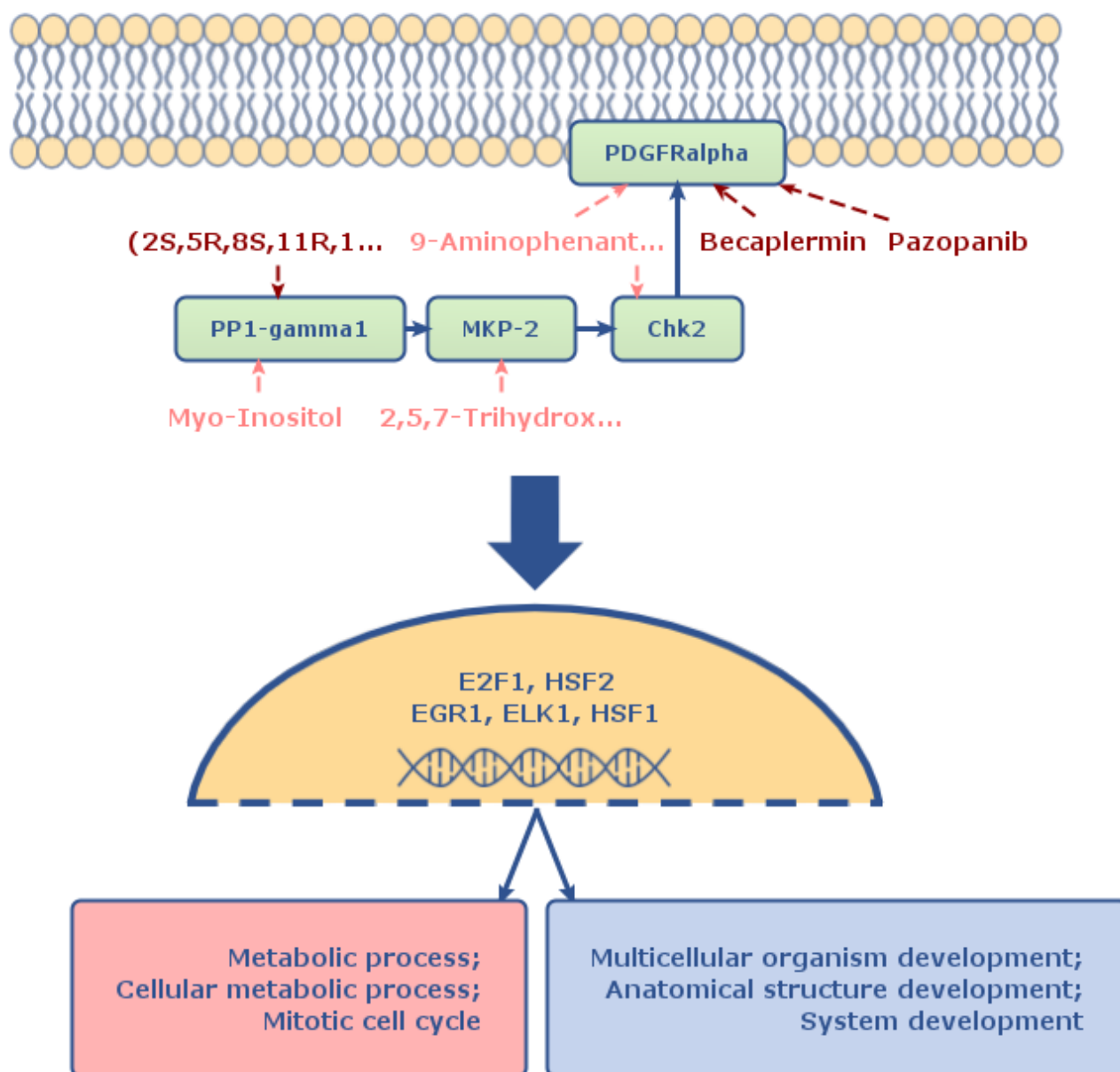
[See full table](#) →

Gene symbol	Gene Description	Druggability score	logFC	Total rank
DUSP4	dual specificity phosphatase 4	45.98	1.17	121
PDGFRA	platelet derived growth factor receptor alpha	81.19	2.93	300
MELK	maternal embryonic leucine zipper kinase	55.87	0.69	422
CYLD	CYLD lysine 63 deubiquitinase	0	1.05	527
HSPA1A	heat shock protein family A (Hsp70) member 1A	121.23	0.58	543
PRKDC	protein kinase, DNA-activated, catalytic subunit	55.2	0.58	611

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:

- PP1-gamma1
- PDGFRalpha
- MKP-2
- Chk2

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: 2,5,7-Trihydroxynaphthoquinone, Myo-Inositol, 9-Aminophenanthrene, Becaplermin, (2S,5R,8S,11R,12S,15S,18S,19S,E)-8-ISOBUTYL-18-((5S,6S)-6-METHOXY-3,5-DIMETHYL-7-PHENYLHEPTYL)-1,2,5,12,15,19-HEXAMETHYL-3,6,9,13,16,20,25-HEPTAOXO-1,4,7,10,14,17,21-HEPTAAZACYCLOPENTACOS-21-ENE-11,22-DICARBOXYLIC ACID and Pazopanib, 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 two 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)).

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

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

Drugs approved in clinical trials



Table 14. 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 rank	Disease activity score	Phase 4	Status (provided by Drugbank)
Pazopanib	ITK, KDR, FLT1, PDGFRB, PDGFRA	7	7	Carcinoma, Renal Cell, Neoplasms, Noma	small molecule, approved
Sunitinib	KDR, PDGFRB, FLT1, PDGFRA	36	2	Carcinoma, Renal Cell, Gastrointestinal Neoplasms, Gastrointestinal Stromal Tumors, Intestinal Neoplasms, Lung Neoplasms, Neoplasms, Neuroendocrine Tumors...	small molecule, approved, investigational
Regorafenib	KDR, FLT1, PDGFRB, PDGFRA, RAF1	39	2	Colorectal Neoplasms, Gastrointestinal Stromal Tumors, Neoplasms, Rectal Neoplasms	small molecule, approved
Imatinib	PDGFRB, PDGFRA	59	3	Breast Neoplasms, Gastrointestinal Stromal Tumors, Leukemia, Leukemia, Lymphoid, Leukemia, Myelogenous, Chronic, BCR-ABL Positive, Leukemia, Myeloid, Mastocytosis...	small molecule, approved
Sorafenib	KDR, PDGFRB, FLT1, RAF1	74	4	Carcinoma, Hepatocellular, Carcinoma, Renal Cell, Liver Neoplasms, Neoplasms, Noma, Thrombosis	small molecule, approved, investigational

Repurposing drugs



Table 15. 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 rank	Phase 4	Status (provided by Drugbank)
Vitamin E	PPP2CB, PPP2CA	116	Angina Pectoris, Variant, Asphyxia, Cicatrix, Cicatrix, Hypertrophic, Diabetes Mellitus, Dyslipidemias, Epilepsy...	small molecule, approved, nutraceutical
Panobinostat	HDAC8, HDAC6, HDAC9, HDAC3	136	Brain Abscess, Multiple Myeloma	small molecule, approved, investigational
Minocycline	CASP3, CASP1, CYCS	170	Acne Vulgaris, Affect, Alopecia, Autistic Disorder, Bacterial Infections, Bipolar Disorder, Chronic Periodontitis...	small molecule, approved, investigational
Mesalazine	IKBKB, CHUK	183	Colitis, Colitis, Ulcerative, Diarrhea, Diverticulum, Irritable Bowel Syndrome, Ulcer	small molecule, approved
Plerixafor	CXCR4	185	Hodgkin Disease, Lymphoma, Lymphoma, Non-Hodgkin	small molecule, approved



Table 16. Prospective drugs, predicted by [PASS](#) software to be active against the identified drug targets with predicted activity against the studied disease(s) (drug candidates predicted with the cheminformatics tool PASS)

[See full table](#) →

Name	Target names	Drug rank	Target activity score
Paclitaxel	GH1, PRKD3, KDR, PRKCE, PRKACA, PRKDC, IL10...	428	0.28
Docetaxel	GH1, PRKD3, KDR, PRKCE, PRKACA, PRKDC, IL10...	439	0.21
Cyclophosphamide	PIK3CG, MTOR, PIK3CA, BCL2, PIK3R5, KDR, FLT1	445	9.32E-2



Table 17. 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 rank	Target activity score
9-Aminophenanthrene	BMPR1A, CSF2RA, PAK2, IL6ST, CDC27, ITGA1, UBE2N...	59	6.41
6-AMINO-BENZO[DE]ISOQUINOLINE-1,3-DI...	BMPR1A, CSF2RA, PAK2, IL6ST, CDK4, CDC27, ITGA1...	62	6.12
2,5,7-Trihydroxynaphthoquinone	CSF2RA, CDC27, TRIM32, PPM1B, DUSP4, TP53BP2, NEK6...	64	6.16
5,8-Di-Amino-1,4-Dihydroxy-Anthraqui...	CSF2RA, CDK4, CDC27, UBE2N, TRIM32, PPM1B, DUSP4...	73	5.16
Aminoanthracene	GH1, BMPR1A, CSF2RA, PAK2, IL6ST, CDC27, ITGA1...	74	6.2

As the result of drug search we propose the following drugs as most promising candidates for treating the pathology under study: Pazopanib, Vitamin E, Paclitaxel and 9-Aminophenanthrene. These drugs were selected for acting on the following targets: PDGFRA,

PPP2CB and PRKDC, 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 multi-omics data set that contains *transcriptomics* and *epigenomics* data obtained from *ovary* tissue. The study is done in the context of *Ovarian Neoplasms*. 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:



Pazopanib, Vitamin E, Paclitaxel and 9-Aminophenanthrene

These drugs were selected for acting on the following targets: PDGFRA, PPP2CB and PRKDC, 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:



PP1-gamma1, PDGFRalpha, MKP-2 and Chk2

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: 2,5,7-Trihydroxynaphthoquinone, Myo-Inositol, 9-Aminophenanthrene, Becaplermin, (2S,5R,8S,11R,12S,15S,18S,19S,E)-8-ISOBUTYL-18-((5S,6S)-6-METHOXY-3,5-DIMETHYL-7-PHENYLHEPTYL)-1,2,5,12,15,19-HEXAMETHYL-3,6,9,13,16,20,25-HEPTAOXO-1,4,7,10,14,17,21-HEPTAAZACYCLOPENTACOS-21-ENE-11,22-DICARBOXYLIC ACID and Pazopanib. 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.

- PP1-gamma1
- PDGFRalpha
- MKP-2
- Chk2

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®](https://genexplain.com/transfac) library, release 2020.2 (geneXplain GmbH, Wolfenbüttel, Germany) (<https://genexplain.com/transfac>).

The master regulator search uses the [TRANSPATH®](https://genexplain.com/transpath) database (BIOBASE), release 2020.2 (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®](https://genexplain.com/transpath).

The information about drugs corresponding to identified drug targets and clinical trials references were extracted from [HumanPSD™](https://genexplain.com/humanpsd) database, release 2020.2 (<https://genexplain.com/humanpsd>).

The Ensembl database release Human99.38 (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.

Genomic data processing

When analyzing a list of genomic variations (from vcf file or computed by Genome Enhancer from fastq files), first of all, we compute a specific mutation weight (w) for each variation depending on its location in gene body and gene flanking regions (-1000 upstream and +1000 downstream of the gene body).

$w = 0.7$ for variations in exon area

$w = 1.3$ for variations in promoter region (-1000bp upstream and 100bp downstream of TSS),

$w = 1.0$ for variations in other locations.

Total Gene mutation weight is the sum of the weights w of all variations located inside the gene body and in the gene flanking regions.

Next, a weighted score is calculated for all genes with the following formula:

Weighted score = $\text{In_disease} * \text{In_transpath} * \text{Gene mutation weight}$, where

$\text{In_disease} = 1.5$ for genes assigned to selected diseases,

$\text{In_transpath} = 2.0$ for genes mapped to Transpath pathways,

and $\text{In_disease} = \text{In_transpath} = 1.0$ in all other cases.

At the next step, 300 genes with highest weighted score are selected for further CMA model search.

The mutation weights (w) are also used to find the regulatory regions of the genes most affected by the variations. A sliding window of 1100 bp is used to scan through the intronic, 5' and 3' regions of the genes and a region is selected with the highest sum of the mutation weights.

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 sum of two other ranks:

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

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

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

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, $rank(t)$ is rank of given target, $maxRank(T)$ equals $max(rank(t))$ for all targets t in T .

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

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

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

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:

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

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:

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

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
0. 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
1. 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.
2. Filimonov DA, Poroikov VV. Prognosis of specters of biological activity of organic molecules. *Russian chemical journal*. **2006**;50(2):66-75 (russ)
3. Filimonov D, Poroikov V, Borodina Y, Gloriovova 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

Supplementary material

1. [Supplementary table 1 - Up-regulated genes](#)
2. [Supplementary table 2 - Down-regulated genes](#)
3. [Supplementary table 3 - Detailed report. Composite modules and master regulators \(up-regulated genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive\).](#)
4. [Supplementary table 4 - Detailed report. Composite modules and master regulators \(down-regulated genes in Experiment: cisplatin-resistant vs. Control: cisplatin-sensitive\).](#)

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.