# PSMA7 and PSMC5 are promising druggable targets for treating Colorectal Neoplasms that control activity of ZEB1, YY1 and NEUROD1 transcription factor on promoters of genes carrying sequence variations

Demo User
geneXplain GmbH
info@genexplain.com
Data received on 28/10/2020 ; Run on 29/01/2021 ; Report generated on 29/01/2021

Genome Enhancer release 2.3 (TRANSFAC®, TRANSPATH® and HumanPSD™ release 2021.1)

## Abstract

In the present study we applied the software package "Genome Enhancer" to a data set that contains *genomics* data. The study is done in the context of *Colorectal 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 genes carrying sequence variations: ZEB1, YY1 and NEUROD1. The subsequent network analysis suggested

- E1
- RPTPepsilon
- ERK1
- 26S proteasome

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: Regorafenib, Tofacitinib, Camptothecin and 2,6-Dihydroanthra/1,9-Cd/Pyrazol-6-One.

## 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 deviced 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 genes carrying sequence variations 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] 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 [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*

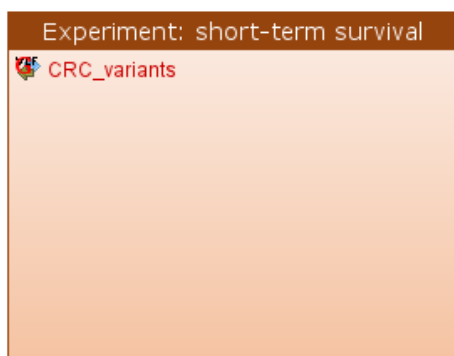| File name | Data type |
|-----------|-----------|
| CRC_variants | Genomics |



*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: short-term survival.

## 3.1. Identification of target genes

In the first step of the analysis **target genes** were identified from the uploaded experimental data. The most frequently mutated genes were used as target genes.

Table 2. Top ten the most frequently mutated genes in Experiment: short-term survival.
**See full table →**

| ID | Gene description | Gene symbol | Gene schematic representation | Number of variations | Gene weight | Weighted score |
|---|---|---|---|---|---|---|
| ENSG00000132570 | pterin-4 alpha-carbinolamine dehydratase 2 | PCBD2 |  | 172 | 171.7 | 257.55 |
| ENSG00000234745 | major histocompatibility complex, class I, B | HLA-B |  | 122 | 109.4 | 218.8 |
| ENSG00000228716 | dihydrofolate reductase | DHFR |  | 56 | 48.2 | 144.6 |
| ENSG00000176890 | thymidylate synthetase | TYMS |  | 44 | 43.7 | 131.1 |
| ENSG00000067057 | phosphofructokinase, platelet | PFKP |  | 92 | 86 | 129 |
| ENSG00000248923 | MT-ND5 pseudogene 11 | MTND5P11 |  | 126 | 121.8 | 121.8 |
| ENSG00000242086 | MUC20 overlapping transcript | MUC20-OT1 |  | 147 | 118.2 | 118.2 |
| ENSG00000169894 | mucin 3A, cell surface associated | MUC3A |  | 68 | 57.2 | 114.4 |
| ENSG00000259755 | novel transcript, antisense to LRRK1 | AC090907.2 |  | 111 | 111 | 111 |
| ENSG00000204525 | major histocompatibility complex, class I, C | HLA-C |  | 71 | 55.4 | 110.8 |

## *3.2. Functional classification of genes*

A functional analysis of genes carrying sequence variations was done by mapping the 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 2-4 show the most significant categories.

## The most frequently mutated genes in Experiment: short-term survival:

300 top mutated genes were taken for the mapping.

**GO (biological process)**

*Figure 2. Enriched GO (biological process) of the most frequently mutated genes in Experiment: short-term survival.*
**Full classification →**
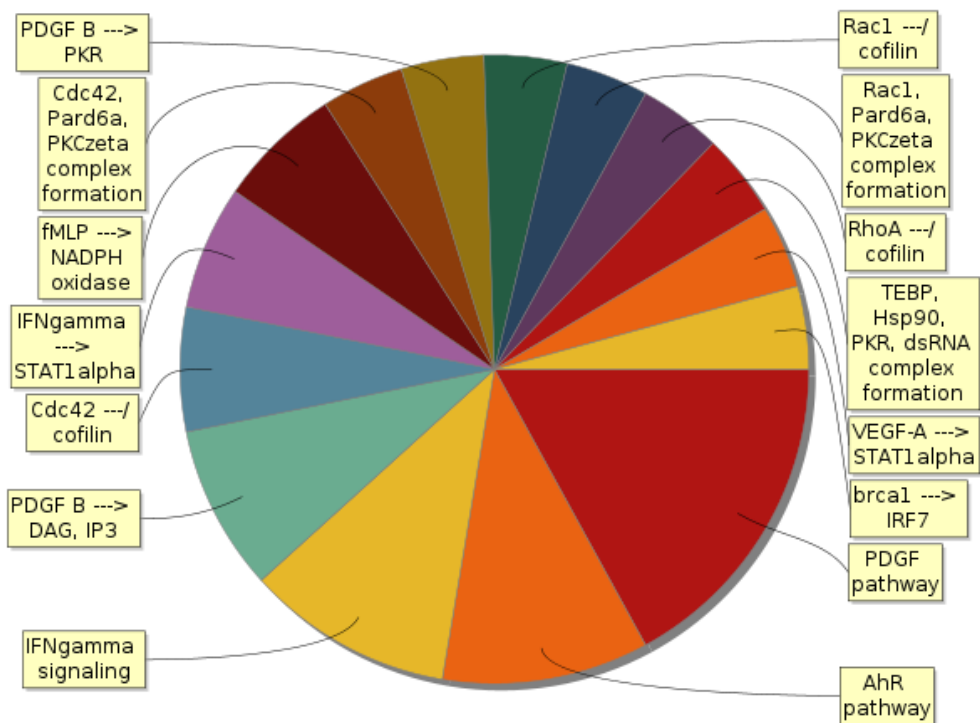
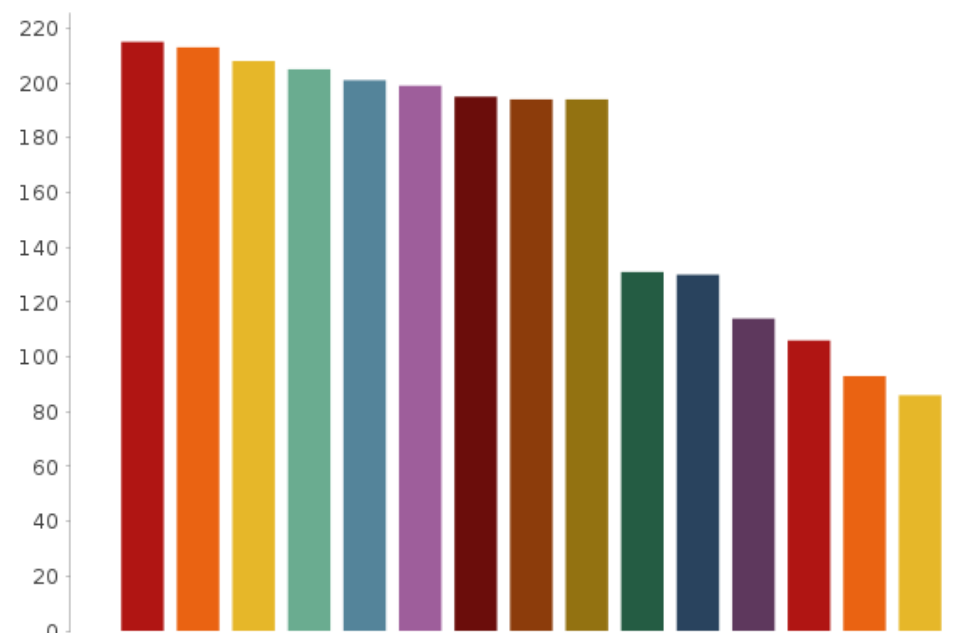**TRANSPATH® Pathways (2021.1)**

*Figure 3. Enriched TRANSPATH® Pathways (2021.1) of the most frequently mutated genes in Experiment: short-term survival.*
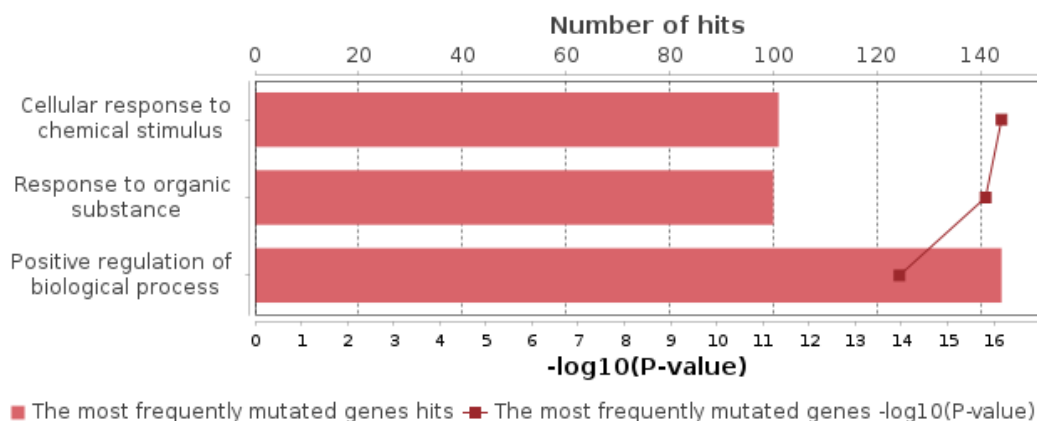
**Full classification →**

**HumanPSD(TM) disease (2021.1)**

Figure 4. Enriched HumanPSD(TM) disease (2021.1) of the most frequently mutated genes in Experiment: short-term survival. 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 genes carrying sequence variations of the studied pathology can be summarized by the following diagram, revealing the most significant functional categories overrepresented among the observed (genes carrying sequence variations):



## 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].

In the current work, we use the Genomics data from the "Yes VCF track" track to predict positions of potential **enhancers** where the observed sequence variations may influence the gene expression in the pathology under

study. We scan 5kb flanking regions and the body of all genes caring the variations, with a sliding window of 1100bp size and find the position of the window with the maximal sum of the mutation weights, where we then perform the search for potential condition-specific enhancers (CMA model search).

We analyzed mutations that were revealed in the potential enhancers located upstream, downstream or inside the **target genes** (see Table 3). We identified 25516 mutations potentially affecting gene regulation. Table 4 shows the following lists of PWMs whose sites were lost or gained due to these mutations. These PWMs were put in focus of the CMA algorithm that constructs the model of the enhancers by specifying combinations of TF motifs (see more details of the algorithm in the Method section).

*Table 3. Mutations revealed in the most frequently mutated genes*
**See full table →**

| ID | Gene symbol | Gene schematic representation | Number of variations |
|---|---|---|---|
| ENSG00000132570 | PCBD2 | | 660 |
| ENSG00000248923 | MTND5P11 | | 459 |
| ENSG00000230021 | AL669831.3 | | 404 |
| ENSG00000247627 | MTND4P12 | | 374 |
| ENSG00000249119 | MTND6P4 | | 279 |
| ENSG00000242086 | MUC20-OT1 | | 252 |
| ENSG00000234745 | HLA-B | | 246 |
| ENSG00000198868 | MTND4LP30 | | 245 |
| ENSG00000263963 | AC008670.1 | | 245 |
| ENSG00000154237 | LRRK1 | | 230 |

*Table 4. PWMs whose sites were lost or gained due to mutations in the most frequently mutated genes*
**See full table →**

| ID | P-value (gains) | P-value (losses) | yesCount (gains) | yesCount (losses) |
|---|---|---|---|---|
| V$MAFG_01 | 2.44E-2 | 7.94E-17 | 29 | 5331 |
| V$BRACH_01 | 2.03E-2 | 6.47E-7 | 276 | 416 |
| V$E2F1EOMES_02 | 4.74E-3 | 1.94E-6 | 28 | 843 |
| V$MAFG_02 | 2.92E-3 | 9.22E-12 | 2223 | 3902 |
| V$ZNF76_03 | 1.9E-3 | 3.22E-10 | 35 | 4969 |
| V$NFIB_02 | 2.99E-4 | 1.74E-27 | 78 | 4834 |
| V$FKLF_02 | 2.5E-9 | 2.34E-9 | 73 | 199 |
| V$WT1_03 | 7.05E-12 | 2.7E-5 | 138 | 80 |
| V$E2F3HES7_01 | 1.17E-20 | 2.26E-5 | 247 | 82 |
| V$SMAD6_02 | 4.94E-33 | | 246 | null |
| V$WT1_Q6_02 | 1.18E-39 | | 413 | null |
| V$TBX2_06 | 1.53E-42 | | 2268 | null |
| V$CENPB_01 | 2.39E-43 | | 791 | null |
| V$GCM1_06 | 4.79E-48 | | 1877 | null |
| V$GCM1_08 | 2.79E-48 | | 1912 | null |
| V$E2F3_09 | 2.69E-53 | | 4545 | null |
| V$E2F2_06 | 1.03E-53 | | 983 | null |
| V$MECP2_02 | 2.36E-56 | | 1741 | null |
| V$TFDP1_02 | 3.87E-58 | | 1163 | null |
| V$MYOGNF1_01 | | 3.78E-6 | null | 7518 |

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 (the most frequently mutated genes in Experiment: short-term survival).

To build the most specific composite modules we choose top mutated genes as the input of CMA algorithm. The obtained CMA model is then applied to compute CMA score for all the most frequently mutated 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$TBX2_06 | V$RUNX2_04 | V$TBX3_05 | V$NKX25_13 | V$NEUROD_02 |
|---|---|---|---|---|
| 0.97; N=3 | 0.92; N=1 | 0.94; N=3 | 0.84; N=3 | 0.97; N=3 |

Module width: 176

**Module 2:**

| V$ZNF462_01 | V$JUND_10 | V$AREB6_01 | V$ATF4CEBPD_01 | V$YY1_06 |
|---|---|---|---|---|
| 0.92; N=2 | 0.79; N=3 | 0.91; N=2 | 0.69; N=3 | 0.89; N=3 |

Module width: 82

**Model score (-p*log10(pval)):** 21.31
**Wilcoxon p-value (pval):** 3.02e-43
**Penalty (p):** 0.501
**Average yes-set score:** 11.18
**Average no-set score:** 9.31
**AUC:** 0.75
**Middle-point:** 10.62
**False-positive:** 25.16%
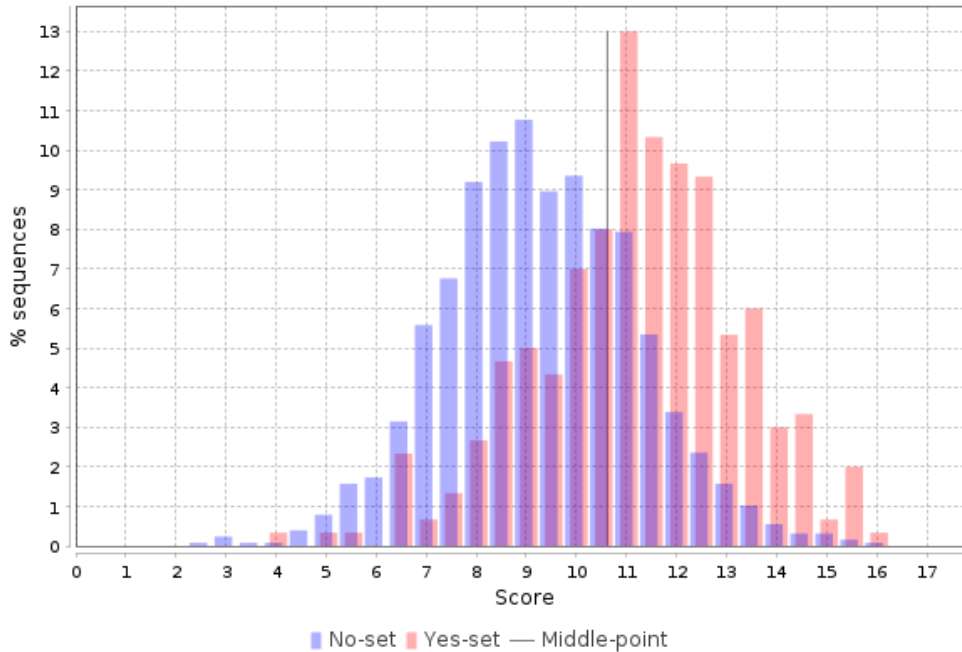**False-negative:** 34.00%

*Table 5. List of top ten the most frequently mutated genes in Experiment: short-term survival 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 →**

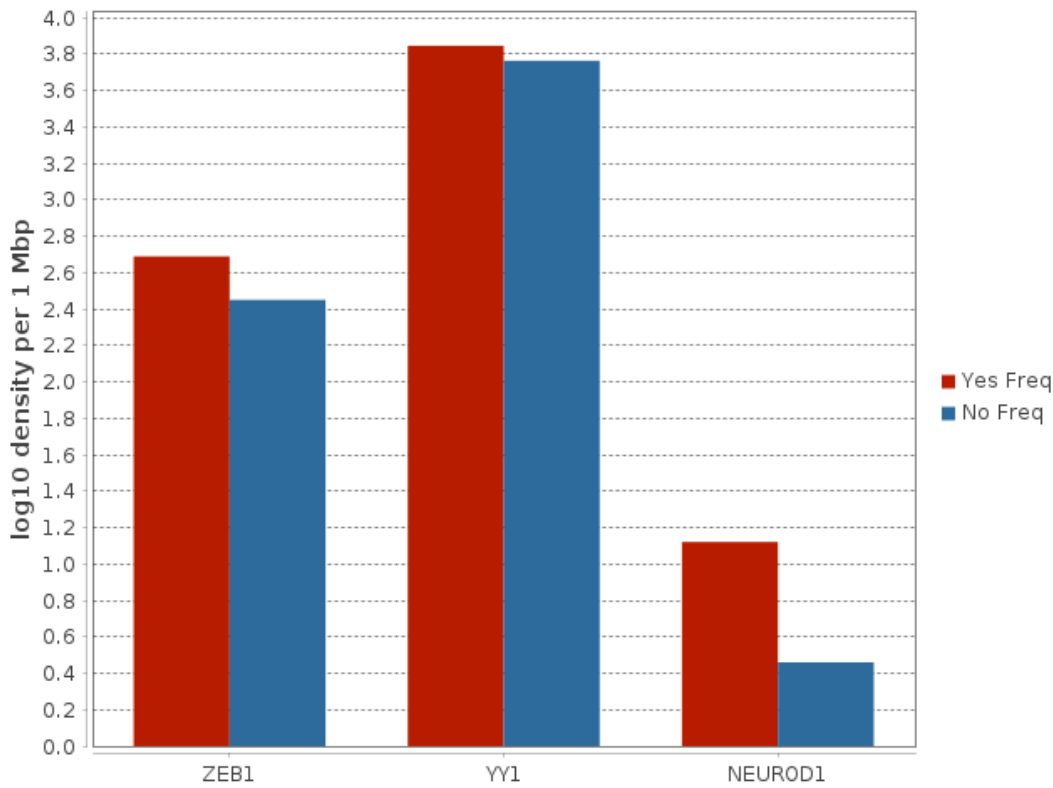| Ensembl IDs | Gene symbol | Gene description | CMA score | Factor names |
|---|---|---|---|---|
| ENSG00000198894 | CIPC | CLOCK interacting pacemaker | 18.32 | CSX(h), AML3(h), ZEB1(h), NeuroD(h), ATF-4(h),C/EBPdelta(h), YY1(h), JunD(h)... |
| ENSG00000259164 | AC007375.2 | novel protein | 18.32 | CSX(h), AML3(h), ZEB1(h), NeuroD(h), ATF-4(h),C/EBPdelta(h), YY1(h), JunD(h)... |
| ENSG00000158555 | GDPD5 | glycerophosphodiester phosphodiesterase domain containing 5 | 18.05 | ATF-4(h),C/EBPdelta(h), ZEB1(h), ZNF462(h), CSX(h), YY1(h), JunD(h), NeuroD(h)... |
| ENSG00000099783 | HNRNPM | heterogeneous nuclear ribonucleoprotein M | 17.79 | CSX(h), TBX2(h), ZEB1(h), Tbx3(h), NeuroD(h), AML3(h), JunD(h)... |
| ENSG00000003756 | RBM5 | RNA binding motif protein 5 | 17.61 | CSX(h), Tbx3(h), ZNF462(h), NeuroD(h), TBX2(h), AML3(h), ATF-4(h),C/EBPdelta(h)... |
| ENSG00000156639 | ZFAND3 | zinc finger AN1-type containing 3 | 17.57 | NeuroD(h), CSX(h), TBX2(h), Tbx3(h), AML3(h), ZNF462(h), JunD(h)... |
| ENSG00000198937 | CCDC167 | coiled-coil domain containing 167 | 17.55 | Tbx3(h), YY1(h), ATF-4(h),C/EBPdelta(h), JunD(h), NeuroD(h), ZNF462(h), TBX2(h)... |
| ENSG00000137177 | KIF13A | kinesin family member 13A | 17.1 | JunD(h), YY1(h), AML3(h), TBX2(h), CSX(h), Tbx3(h), ZNF462(h)... |
| ENSG00000278214 | LINC02139 | long intergenic non-protein coding RNA 2139 | 16.96 | NeuroD(h), TBX2(h), CSX(h), Tbx3(h), AML3(h), ATF-4(h),C/EBPdelta(h), JunD(h)... |
| ENSG00000100379 | KCTD17 | potassium channel tetramerization domain containing 17 | 16.9 | ATF-4(h),C/EBPdelta(h), NeuroD(h), ZEB1(h), ZNF462(h), YY1(h), JunD(h), AML3(h)... |

On the basis of the enhancer models we identified transcription factors potentially regulating the ***target genes*** of our interest. We found 11 transcription factors controlling expression of the genes associated with genomic variations (see Table 6).

*Table 6. Transcription factors of the predicted enhancer model potentially regulating the genes carrying sequence variations (the most frequently mutated genes in Experiment: short-term survival).* **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 |
|---|---|---|---|---|
| MO000139677 | ZEB1 | zinc finger E-box binding homeobox 1 | 3.13 | 1.73 |
| MO000078913 | YY1 | YY1 transcription factor | 3.01 | 1.21 |
| MO000028384 | NEUROD1 | neuronal differentiation 1 | 2.86 | 4.58 |
| MO000007834 | JUND | JunD proto-oncogene, AP-1 transcription factor subunit | 2.86 | 1.55 |
| MO000019140 | ATF4 | activating transcription factor 4 | 2.84 | 2.03 |
| MO000026285 | RUNX2 | RUNX family transcription factor 2 | 2.69 | 2.44 |
| MO000002641 | CEBPD | CCAAT enhancer binding protein delta | 2.65 | 2.21 |
| MO000092587 | ZNF462 | zinc finger protein 462 | 2.65 | 1.17 |
| MO000028181 | NKX2-5 | NK2 homeobox 5 | 2.5 | 1.79 |
| MO000028209 | TBX2 | T-box transcription factor 2 | 2.07 | 18.31 |

The following diagram represents the key transcription factors, which were predicted to be potentially regulating genes carrying sequence variations in the analyzed pathology: ZEB1, YY1 and NEUROD1.

## 3.4. Finding master regulators in networks

In the second step of the upstream analysis common regulators of the revealed TFs were identified. We identified 169 signaling proteins whose structure and function is highly damaged by the mutations (see Table 7).

*Table 7. Signaling proteins whose structure and function is damaged by the mutations in the most frequently mutated genes*
**See full table →**

| ID | Title | Mutation count | Consequence | Codons |
|---|---|---|---|---|
| MO000138949 | Drp1(h) | 13 | NMD_transcript_variant,stop_gained | Gaa/Taa |
| MO000019673 | p85alpha(h) | 9 | stop_gained | Cga/Tga |
| MO000113258 | MYPT1(h) | 8 | NMD_transcript_variant,frameshift_variant | aga/aAga |
| MO000127741 | SMC4L1(h) | 8 | stop_gained | Cga/Tga |
| MO000214698 | MS4A6A(h) | 8 | NMD_transcript_variant,frameshift_variant | -/T,tta/ttTa |
| MO000035319 | kinectin(h) | 7 | NMD_transcript_variant,frameshift_variant | -/A |
| MO000144675 | NULP1(h) | 7 | NMD_transcript_variant,frameshift_variant | -/A |
| MO000145695 | Anamorsin(h) | 7 | NMD_transcript_variant,frameshift_variant | -/A |
| MO000206935 | C11orf74(h) | 7 | stop_gained | Gaa/Taa |
| MO000068933 | HLA-G(h) | 6 | NMD_transcript_variant,splice_region_variant,stop_lost | Tga/Aga |

Top 100 mutated proteins for the most frequently mutated genes were used in the algorithm of master regulator search as a list of nodes of the signal transduction network that are removed from the network during the search of master regulators (see more details about the algorithm in the Method section). 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 Table 8.

Table 8. Master regulators that may govern the regulation of the most frequently mutated genes in Experiment: short-term survival. **Total rank** is the sum of the ranks of the master molecules sorted by keynode score, CMA score, genomics data.
**See full table →**

| ID | Master molecule name | Gene symbol | Gene description | Total rank |
|---|---|---|---|---|
| MO000019948 | E1(h) | UBA1 | ubiquitin like modifier activating enzyme 1 | 23 |
| MO000089081 | E1-isoform1(h) | UBA1 | ubiquitin like modifier activating enzyme 1 | 74 |
| MO000004672 | ERK1(h) | MAPK3 | mitogen-activated protein kinase 3 | 126 |
| MO000019259 | c-Cbl(h) | CBL | Cbl proto-oncogene | 138 |
| MO000018962 | ErbB2(h) | ERBB2 | erb-b2 receptor tyrosine kinase 2 | 143 |
| MO000078913 | YY1(h) | YY1 | YY1 transcription factor | 145 |
| MO000031189 | PKCdelta(h) | PRKCD | protein kinase C delta | 153 |
| MO000021344 | Jak3(h) | JAK3 | Janus kinase 3 | 186 |
| MO000031003 | ERK1(h){p} | MAPK3 | mitogen-activated protein kinase 3 | 198 |
| MO000059577 | PKCdelta(h) | PRKCD | protein kinase C delta | 198 |

The intracellular regulatory pathways controlled by the above-mentioned master regulators are depicted in Figure 5. This diagram displays the connections between identified transcription factors, which play important roles in the regulation of genes carrying sequence variations, and selected master regulators, which are responsible for the regulation of these TFs.
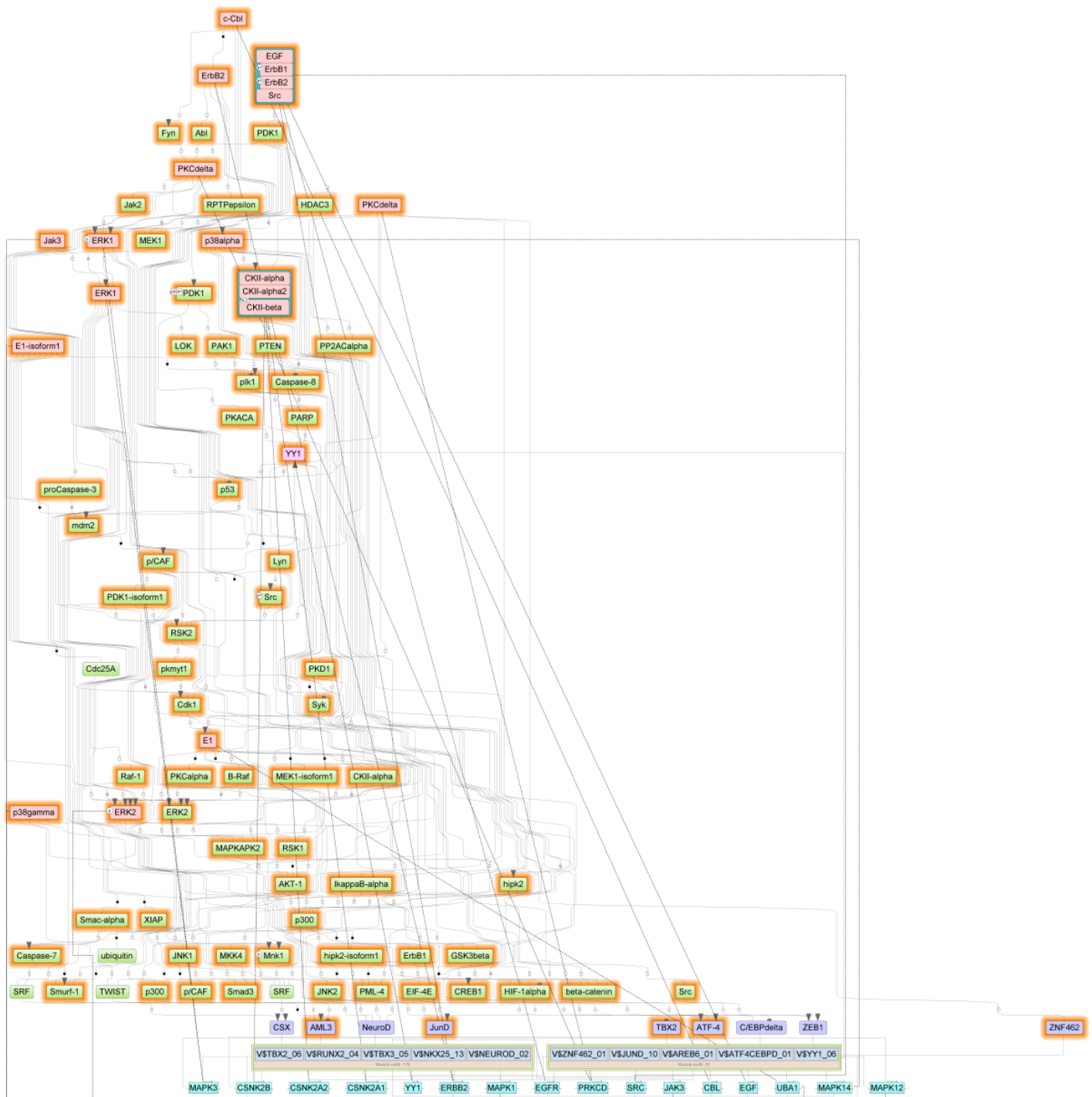
*Figure 5. Diagram of intracellular regulatory signal transduction pathways of the most frequently mutated genes in Experiment: short-term survival. Master regulators are indicated by red rectangles, transcription factors are blue rectangles, and green rectangles are intermediate molecules, which have been added to the network during the search for master regulators from selected TFs. Orange frames highlight molecules presented in original mapping.*

**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 9. 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 |
|---|---|---|---|
| PSMA7 | proteasome 20S subunit alpha 7 | 3 | 258 |
| PTPRE | protein tyrosine phosphatase receptor type E | 1 | 312 |
| PLK1 | polo like kinase 1 | 5 | 457 |
| PPP2CA | protein phosphatase 2 catalytic subunit alpha | 3 | 476 |
| JAK3 | Janus kinase 3 | 3 | 513 |
| IL2RB | interleukin 2 receptor subunit beta | 4 | 544 |

*Table 10. 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.*
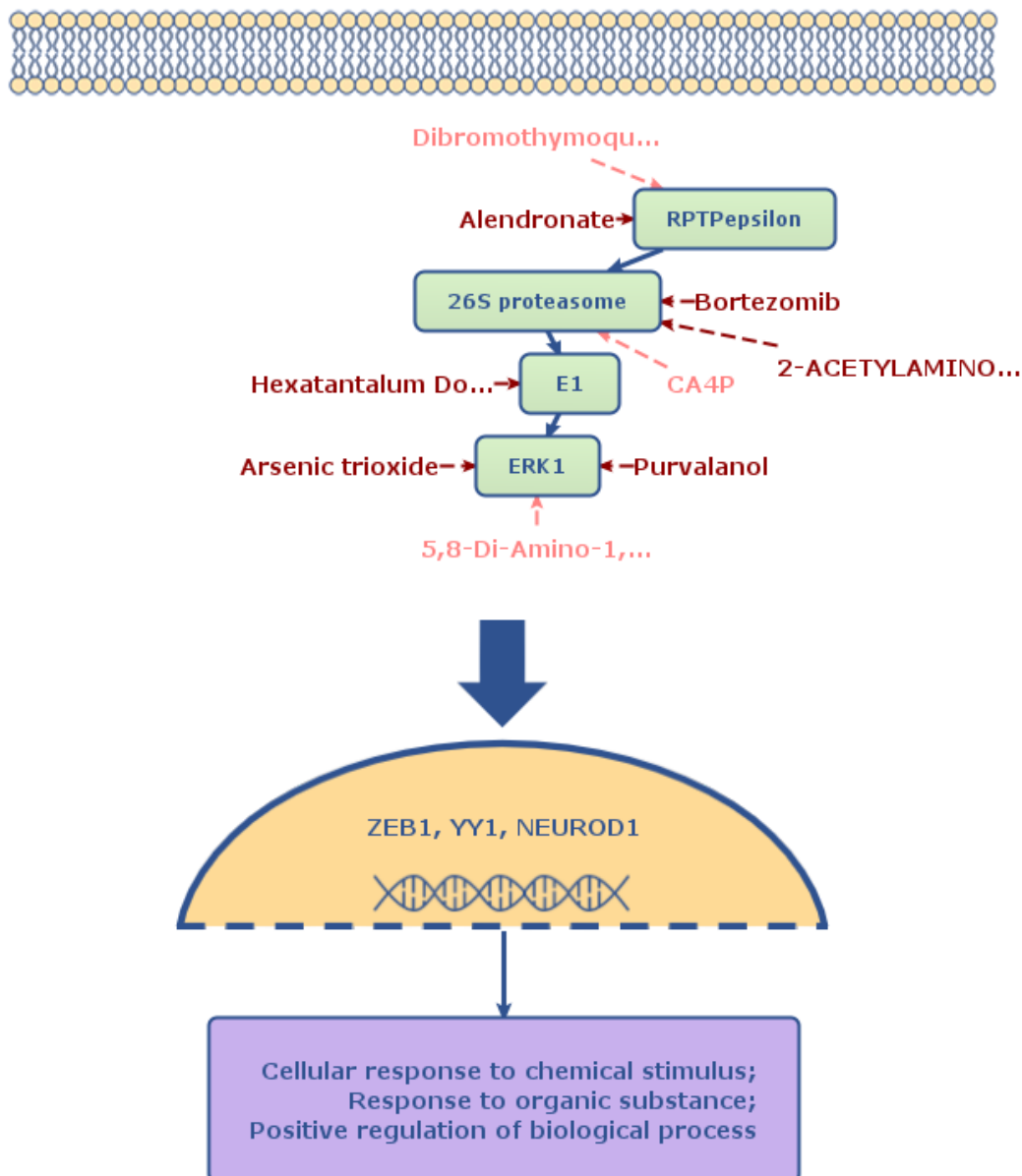
**See full table →**

| Gene symbol | Gene Description | Druggability score | Total rank |
|---|---|---|---|
| PSMC5 | proteasome 26S subunit, ATPase 5 | 1.28 | 258 |
| PSMD5 | proteasome 26S subunit, non-ATPase 5 | 1.28 | 258 |
| PSMA7 | proteasome 20S subunit alpha 7 | 2.48 | 258 |
| PSMD4 | proteasome 26S subunit, non-ATPase 4 | 1.28 | 258 |
| PSMC2 | proteasome 26S subunit, ATPase 2 | 1.28 | 258 |
| PSMC3 | proteasome 26S subunit, ATPase 3 | 1.28 | 258 |

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:

- E1
- RPTPepsilon
- ERK1
- 26S proteasome

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: Bortezomib, Alendronate, CA4P, Arsenic trioxide, Hexatantalum Dodecabromide, Dibromothymoquinone, Purvalanol, 2-ACETYLAMINO-4-METHYL-PENTANOIC ACID [1-(1-FORMYL-PENTYLCARBAMOYL)-3-METHYL-BUTYL]-AMIDE and 5,8-Di-Amino-1,4-Dihydroxy-Anthraquinone, 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 11. 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) |
|------|-------------|-----------|------------------------|---------|-------------------------------|
| Regorafenib | KIT, KDR, ABL1, PDGFRB, FGFR1, TEK, RET... | 4 | 11 | Colorectal Neoplasms, Gastrointestinal Stromal Tumors, Neoplasms, Rectal Neoplasms | small molecule, approved |
| Nintedanib | FGFR3, SRC, KDR, LYN, FGFR1 | 20 | 6 | Idiopathic Pulmonary Fibrosis, Pulmonary Fibrosis | small molecule, approved |
| Sorafenib | KIT, KDR, PDGFRB, FGFR1, BRAF, RAF1, RET | 22 | 4 | Carcinoma, Hepatocellular, Carcinoma, Renal Cell, Liver Neoplasms, Neoplasms, Noma, Thrombosis | small molecule, approved, investigational |
| Sunitinib | KIT, KDR, PDGFRB, PDGFRA | 28 | 6 | Carcinoma, Renal Cell, Gastrointestinal Neoplasms, Gastrointestinal Stromal Tumors, Intestinal Neoplasms, Lung Neoplasms, Neoplasms, Neuroendocrine Tumors... | small molecule, approved, investigational |
| Dasatinib | KIT, SRC, ABL1, PDGFRB, YES1, FYN, ABL2 | 34 | 3 | Leukemia, Leukemia, Lymphoid, Leukemia, Myelogenous, Chronic, BCR-ABL Positive, Leukemia, Myeloid, Precursor Cell Lymphoblastic Leukemia-Lymphoma | small molecule, approved, investigational |

## Repurposing drugs

Table 12. 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) |
|------|-------------|-----------|---------|-------------------------------|
| Tofacitinib | JAK3, JAK2, JAK1 | 85 | Arthritis, Arthritis, Rheumatoid | small molecule, approved |
| Anti-thymocyte Globulin (Rabbit) | ITGB1, ITGAV, ITGAL, ITGB3, CD4 | 91 | Anemia, Anemia, Aplastic, Leukemia, Liver Diseases | biotech, approved |
| XL184 | KIT, KDR, TEK | 113 | Neoplasms, Thyroid Neoplasms | small molecule, investigational |
| Alendronate | PTPRS, PTPRE | 116 | Arteriosclerosis, Arthritis, Arthritis, Rheumatoid, Bone Demineralization, Pathologic, Bone Diseases, Bone Diseases, Metabolic, Cystic Fibrosis... | small molecule, approved |
| Tirofiban | ITGB3, ITGA2B | 125 | Acute Coronary Syndrome, Coronary Artery Disease, Coronary Disease, Myocardial Infarction, No-Reflow Phenomenon, ST Elevation Myocardial Infarction | small molecule, approved |

Table 13. *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)*

| Name | Target names | Drug rank | Target activity score |
|---|---|---|---|
| Camptothecin | HIF1A, CASP3 | 127 | 0.31 |
| Topotecan | HIF1A, CASP3 | 128 | 0.31 |
| LE-SN38 | HIF1A, CASP3 | 131 | 0.29 |
| MGI-114 | CASP8, CASP9 | 157 | 0.15 |
| Ouabain | STAT3, CASP8, HIF1A, CASP3, RELA | 159 | 0.14 |

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)*

| Name | Target names | Drug rank | Target activity score |
|---|---|---|---|
| 2,6-Dihydroanthra/1,9-Cd/Pyrazol-6-One | MAPK10, RPS6KA3, IRAK4, CDK6, CAMK2G, CSNK1A1, PAK2... | 25 | 10.23 |
| Iodophenyl | STK10, ROCK2, MARK3, PAK2, GSK3B, SLK, VRK1... | 27 | 16.14 |
| Rbt205 Inhibitor | RPS6KA3, CDK6, CAMK2G, GRK2, MAP3K10, PRKCQ, PRKACA... | 31 | 15.93 |
| Uracil mustard | KDR, ABL1, JAK3, EPHA4, PDGFRA, INSR, MST1R... | 42 | 5.83 |
| 3,5-Diaminophthalhydrazide | RPS6KA3, IRAK4, CAMK2G, CSNK1A1, PRKCQ, PRKCA, GSK3B... | 53 | 5.67 |

As the result of drug search we propose the following drugs as most promising candidates for treating the pathology under study: Regorafenib, Tofacitinib, Camptothecin and 2,6-Dihydroanthra/1,9-Cd/Pyrazol-6-One. These drugs were selected for acting on the following targets: RET, JAK3, HIF1A and CHEK2, 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 *genomics* data. The study is done in the context of *Colorectal Neoplasms*. The data were pre-processed, statistically analyzed and genes carrying sequence variations 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:

**Regorafenib, Tofacitinib, Camptothecin and 2,6-Dihydroanthra/1,9-Cd/Pyrazol-6-One**

These drugs were selected for acting on the following targets: RET, JAK3, HIF1A and CHEK2, 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:

**E1, RPTPepsilon, ERK1 and 26S proteasome**

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: Bortezomib, Alendronate, CA4P, Arsenic trioxide, Hexatantalum Dodecabromide, Dibromothymoquinone, Purvalanol, 2-ACETYLAMINO-4-METHYL-PENTANOIC ACID [1-(1-FORMYL-PENTYLCARBAMOYL)-3-METHYL-BUTYL]-

AMIDE and 5,8-Di-Amino-1,4-Dihydroxy-Anthraquinone. 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 genes carrying sequence variations 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.

- E1
- RPTPepsilon
- ERK1
- 26S proteasome

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 2021.1 (geneXplain GmbH, Wolfenbüttel, Germany) (https://genexplain.com/transfac).
The master regulator search uses the TRANSPATH® database (BIOBASE), release 2021.1 (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 2021.1 (https://genexplain.com/humanpsd).
The Ensembl database release Human100.38 (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 it's 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 = In_disease * In_transpath * Gene mutation weight, where

In_disease = 2.0 for genes assigned to selected diseases,
In_transpath = 1.5 for genes mapped to Transpath pathways,
and In_disease = 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-score$_{PSD}$*),
  2. ranking by "Disease activity score" (*D-score$_{PSD}$*).
"Target activity score" ( *T-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$_{PSD}$*):

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

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:

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

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
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, Gloriozova 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 - Detailed report. Composite modules and master regulators (the most frequently mutated genes in Experiment: short-term survival).
2. Supplementary table 2 - 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.