# AR and DUSP9 are promising druggable targets for treating Hypertension that control activity of ETS1, RAD21 and BCL6 transcription factor on of highly methylated genes in blood tissue

Demo User geneXplain GmbH info@genexplain.com Data received on 08/04/2022 ; Run on 09/07/2023 ; Report generated on 09/07/2023

Genome Enhancer release 3.2 (TRANSFAC®, TRANSPATH® and HumanPSD<sup>™</sup> release 2023.1)



### Abstract

In the present study we applied the software package "Genome Enhancer" to a data set that contains *epigenomics* data obtained from *blood* tissue. The study is done in the context of *Hypertension*. 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 highly methylated genes: ETS1, RAD21 and BCL6. The subsequent network analysis suggested

- TSC22D3
- MKP-4
- AR-isoform1
- AR

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: Imatinib, Flavopiridol, Moexipril and 2,5,7-Trihydroxynaphthoquinone.

## 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 highly methylated genes for the transcription factors (TFs) involved in their regulation and, thus, important for the process under study; (2) re-constructing the signaling pathways that activate these TFs and identifying master regulators at the top of such pathways. For the first step, the database TRANSFAC® [6] is employed together with the TF binding site identification

algorithms Match [7] and CMA [8]. The second step involves the signal transduction database TRANSPATH® [9] and special graph search algorithms [10] 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<sup>™</sup> 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<sup>™</sup> 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
df_hp	Epigenomics
df_norm	Epigenomics

Hypertension	Control
LCL_0001	LCL_0007
■ df_hp = LCL 0005	➡ df_norm
df_hp	df_norm
tot_0016 ₩df_hp	df_norm
ELCL_0019 ₩df_hp	LCL_0017
ELCL_0021 #df_hp	ELCL_0018
LCL_0023	LCL_0022
ELCL_0028	LCL_0029
LCL_0032	
ELCL_0037	LCL_0066
LCL_0041	LCL_0074
LCL_0044	LCL_0076
•• df_hp = LCL_0048	E LCL_0077
■ df_hp = LCL_0050	← df_norm ⊑ LCL_0083
■ df_hp ■ LCL 0052	"₩ df_norm ⊑ LCL 0087
# df_hp	df_norm
# df_hp	df_norm
ELCL_0055 #df_hp	LCL_0091
ELCL_0057	ELCL_0095
LCL_0060	LCL_0003
E_LCL_0063	LCL_0097
ELCL_0072	LCL_0102
LCL_0082	LCL_0107
LCL_0084	
LCL_0085	
	LCL_0128
	ELCL_0129
■ df_hp = LCL_0093	- ₩ df_norm
■ df_hp = ICI 0101	<sup>™</sup> # df_norm ⊏ ICL 0136
# df_hp	<b>#</b> df_norm

ELCL\_0105 LCL\_0111 LCL\_0112 LCL\_0113 LCL\_0120 LCL\_0125 LCL\_0131 LCL\_0141 LCL\_0162 ELCL\_0171 LCL\_0175 LCL\_0177 LCL\_0189 ELCL\_0190 LCL\_0192 ELCL\_0193 LCL\_0194 ELCL\_0198 ELCL\_0199 LCL\_0200 LCL\_0205 LCL\_0211 ELCL\_0213 LCL\_0214 LCL\_0217 ELCL\_0233 ELCL\_0235 ELCL\_0236 LCL\_0237 ELCL\_0241 ELCL\_0246 ELCL\_0255 LCL\_0256 LCL\_0257 ELCL\_0262 ELCL\_0265 LCL\_0266 LCL\_0272 ELCL\_0280 ELCL\_0285 ELCL\_0298 LCL\_0299 LCL\_0310 ELCL\_0314 LCL\_0322

E,	LGL_0144 df_norm
E,	LCL_0145
E.	LCL_0147
	df_norm LCL 0148
<b>-</b> #	df_norm
E#	df_norm
E <sub>#</sub>	LCL_0155 df_norm
E,	LCL_0158 df_norm
E,	LCL_0159
E,	LCL_0135
Е.	ar_norm LCL_0164
- #	df_norm LCL 0165
-#	df_norm
E <sub>#</sub>	df_norm
E,	LCL_0168 df_norm
E,	LCL_0174 df_norm
E,	LCL_0176
E.,	LCL_0179
Е.	dt_norm LCL_0183
-*	df_norm ICI 0188
F#	df_norm
E,	df_norm
E <sub>#</sub>	LCL_0197 df_norm
E,	LCL_0204 df_norm
E,	LCL_0209
E,	LCL_0219
Е.,	LCL_0220
Е.	LCL_0221
F	df_norm LCL_0223
-*	df_norm ICI 0224
5	df_norm
F#	df_norm
E,	LCL_0234 df_norm
E <sub>#</sub>	LCL_0244 df_norm
E,	LCL_0245 df.norm
E,	LCL_0253
E,	LCL_0258
E.	LCL_0260
F	df_norm LCL_0264
-#	df_norm ICI 0274
<b>-</b> #	df_norm
E#	df_norm
E#	df_norm
E#	LCL_0291 df_norm
E,	LCL_0293 df_norm
E,	LCL_0294 df norm
E,	LCL_0295
E,	LCL_0301
fe.	LCL_0303
-# F	df_norm LCL 0306
-#	df_norm

ELCL\_0324 LCL\_0327 ELCL\_0338 ELCL\_0350 LCL\_0351 LCL\_0352 ELCL\_0355 ELCL\_0356 LCL\_0363 LCL\_0366 LCL\_0369 ELCL\_0372 ELCL\_0374 LCL\_0392 ELCL\_0393 ELCL\_0398 ELCL\_0399 ELCL\_0401 LCL\_0403 ELCL\_0406 ELCL\_0408 ELCL\_0410 LCL\_0417 LCL\_0420 ELCL\_0430 ELCL\_0434 LCL\_0441 ELCL\_0444 LCL\_0445 ELCL\_0447 ELCL\_0460 LCL\_0468 LCL\_0473 ELCL\_0478 ELCL\_0483 ELCL\_0484 LCL\_0486 ELCL\_0507 ELCL\_0509 ELCL\_0511 LCL\_0512 ELCL\_0514 ELCL\_0515 ELCL\_0518 LCL\_0519

E,	LCL_0307
_	dt_norm LCL 0315
H	df_norm
E <sub>#</sub>	df_norm
E,	LCL_0318 df norm
Е.	LCL_0319
-	df_norm LCL 0325
-#	df_norm
E <sub>#</sub>	df_norm
E,	LCL_0329 df_norm
E,	LCL_0336
Е.,	LCL_0337
•	df_norm LCL 0341
<b>-</b> #	df_norm
E <sub>#</sub>	df_norm
E <sub>#</sub>	LCL_0361 df_norm
E,	LCL_0365
Е.	LCL_0375
-	df_norm LCL 0388
Ħ	df_norm
E <sub>#</sub>	df_norm
E <sub>#</sub>	LCL_0427 df_norm
E,	LCL_0428
E.,	LCL_0429
E.	LCL_0437
-*	df_norm ICI 0442
F#	df_norm
E <sub>#</sub>	df_norm
E,	LCL_0476 df_norm
E,	LCL_0496
Е.,	LCL_0498
	df_norm LCL_0499
-#	df_norm
F#	df_norm
E <sub>#</sub>	df_norm
E <sub>#</sub>	LCL_0504 df_norm
E,	LCL_0505
E.	LCL_0506
f	LCL_0531
-#	df_norm
<b>-</b> #	df_norm
E,	df_norm
E,	LCL_0568 df_norm
E <sub>#</sub>	LCL_0578
Е.	LCL_0585
۴ F	df_norm LCL_0586
*#	df_norm
5#	df_norm
E#	LCL_0599 df_norm
E <sub>#</sub>	LCL_0605 df_norm
E <sub>#</sub>	LCL_0607
E.	LCL_0616
f	df_norm LCL_0618
-#	df_norm

1

ELCL\_0522 ELCL\_0528 ELCL\_0532 LCL\_0534 LCL\_0537 ELCL\_0547 ELCL\_0552 LCL\_0554 LCL\_0542 ELCL\_0569 ELCL\_0570 ELCL\_0571 LCL\_0581 LCL\_0584 ELCL\_0590 ELCL\_0597 LCL\_0604 LCL\_0610 ELCL\_0611 LCL\_0612 ELCL\_0620 LCL\_0630 LCL\_0631 ELCL\_0637 ELCL\_0641 LCL\_0642 ELCL\_0659 ELCL\_0660 ELCL\_0665 ELCL\_0673 LCL\_0677 ELCL\_0678 ELCL\_0687 ELCL\_0691 LCL\_0694 ELCL\_0697 ELCL\_0701 ELCL\_0702 ELCL\_0707 ELCL\_0711 LCL\_0723 ELCL\_0726 ELCL\_0727 LCL\_0732 LCL\_0733

E,	LCL_0625
E.	LCL_0626
- #	df_norm ICI 0627
<b>F</b> #	df_norm
E,	df_norm
E#	LCL_0634 df_norm
E <sub>#</sub>	LCL_0638 df_norm
E,	LCL_0640
E.,	LCL_0643
Е.	LCL_0644
F.	df_norm LCL_0645
-#	df_norm LCL 0655
-#	df_norm
F#	df_norm
E <sub>#</sub>	df_norm
E#	LCL_0676 df_norm
E,	LCL_0690 df_norm
E,	LCL_0692 df norm
E,	 LCL_0693
E,	LCL_0700
Е.	LCL_0705
F.	df_norm LCL_0709
-# E	df_norm LCL 0717
5# -	df_norm
F#	df_norm
E#	df_norm
E	LCL_0724 df_norm
E#	LCL_0728 df_norm
E,	LCL_0731 df_norm
E,	LCL_0734 df.norm
E,	_ LCL_0738
E,	LCL_0739
	LCL_0745
Ē	df_norm LCL_0749
-# F	df_norm LCL 0754
-#	df_norm
F#	df_norm
E	df_norm
E,	LCL_0762 df_norm
E#	LCL_0764 df_norm
E#	LCL_0765 df_norm
E,	LCL_0769 df_norm
E <sub>#</sub>	LCL_0780
E,	LCL_0784
E.	LCL_0788
f.	LCL_0789
=#	df_norm LCL_0793
•# E	df_norm LCL 0797
<b>-</b> #	df_norm
F#	df_norm

ELCL\_0737 ELCL\_0741 LCL\_0742 ELCL\_0743 LCL\_0755 ELCL\_0773 LCL\_0774 LCL\_0775 ELCL\_0778 ELCL\_0783 ELCL\_0785 LCL\_0795 LCL\_0796 ELCL\_0800 ELCL\_0808 ELCL\_0810 ELCL\_0814 ELCL\_0817 ELCL\_0836 ELCL\_0837 LCL\_0866 LCL\_0869 ELCL\_0871 ELCL\_0873 LCL\_0877 ELCL\_0893 ELCL\_0903 ELCL\_0905 ELCL\_0907 LCL\_0911 ELCL\_0912 ELCL\_0916 ELCL\_0923 ELCL\_0935 ELCL\_0937 ELCL\_0942 ELCL\_0950 ELCL\_0958 ELCL\_0959 ELCL\_0960 ELCL\_0961 ELCL\_0963 ELCL\_0967 ELCL\_0975 ELCL\_0978

Es.	LCL_0805
<b>H</b>	df_norm
E,	LCL_0809
F	LCL 0811
"#	df_norm
E,	LCL_0813
<b>F</b> =	LCL 0821
<b>H</b>	df_norm
E,	LCL_0823
<b>F</b> =	LCL 0824
<b>H</b>	df_norm
E,	LCL_0826
E=	LCL_0827
"H	df_norm
E,	LCL_0828 df_norm
Е.	LCL_0831
	df_norm
E,	df_norm
E.,	LCL_0835
_	df_norm
F#	df_norm
E.,	LCL_0847
	at_norm
F#	df_norm
E,	LCL_0853
	ICI 0854
F#	df_norm
E,	LCL_0858
<b>F</b> .	LCL 0863
*	df_norm
E,	LCL_0865 df norm
E.	LCL_0870
"H	df_norm
E,	LCL_0874 df_norm
E.,	LCL_0879
_"	
F#	df_norm
E,	LCL_0883
<b>F</b> .	LCL 0884
Ħ	df_norm
E,	LCL_0887
E=	LCL_0888
"H	df_norm
E,	LCL_0889 df_norm
Е.	LCL_0890
	df_norm
E.	df_norm
Е,	LCL_0898
_	LCL 0904
<b>H</b>	df_norm
E,	LCL_0909
F	LCL_0914
Ħ	df_norm
E,	LCL_0917 df_norm
E	LCL_0920
	df_norm
E,	df_norm
E.	LCL_0922
-	df_norm
F#	df_norm
E,	LCL_0928
_	LCL 0943
Ħ	df_norm
E,	LCL_0954 df norm
F	LCL_0964
-#	df_norm

E.,	LCL_0980
•	df_hp LCL_0984
Ē	df_hp LCL 0987
-#	df_hp
F#	df_hp
E,	df_hp
E,	LCL_1002 df_hp
E,	LCL_1006 df_hp
E,	LCL_1007
Е.,	LCL_1025
	at_np LCL_1027
-# F	df_hp LCL 1028
-#	df_hp
F#	df_hp
E,	LCL_1041 df_hp
E,	LCL_1048 df_hp
E,	LCL_1055 df hp
E,	LCL_1066
Е.	LCL_1069
F	df_hp LCL_1078
-#	df_hp ICL 1086
F#	df_hp
E,	df_hp
E,	LCL_1094 df_hp
E,	LCL_1096 df_hp
E,	LCL_1098 dfhp
E.	LCL_1100
E.	LCL_1102
	at_np LCL_1104
THE	df_hp LCL 1105
-#	df_hp
F#	df_hp
E,	df_hp
E,	LCL_1115 df_hp
E,	LCL_1119 df_hp
E,	LCL_1120
Е.,	LCL_1121
	dt_hp LCL_1123
-#	df_hp LCL 1124
Ħ	df_hp
<b>F</b> #	df_hp
E <sub>#</sub>	df_hp
E,	LCL_1133 df_hp
E,	LCL_1134 df_hp
E,	LCL_1135
E,	LCL_1137
E.	LCL_1142
f E	ar_hp LCL_1143
=#	df_hp LCL 1145
Ħ	df_hp
F#	df_hp

E,	LCL_0968 df_norm
E,	LCL_0969
E,	LCL_0982
E.	LCL_0990
E.,	LCL_0994
E.	LCL_1005
E.	df_norm LCL_1021
Е.	df_norm LCL_1032
F	df_norm LCL_1034
Ē	df_norm LCL 1035
-# =	df_norm LCL 1039
*	df_norm
Ħ	df_norm
F#	df_norm
F.	df_norm
E,	df_norm
E,	LCL_1079 df_norm
E,	LCL_1087 df_norm
E,	LCL_1088 df_norm
E,	LCL_1097 df_norm
E,	LCL_1101 df_norm
E,	LCL_1122 df_norm
E,	LCL_1144 df_norm
E#	LCL_1162 df_norm
E,	LCL_1176 df_norm
E,	LCL_1185 df_norm
E#	LCL_1193 df_norm
E,	LCL_1194
E,	LCL_1202
	_



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: Hypertension versus Control.

### 3.1. Identification of target genes

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

Table 2. Top ten highly methylated genes in Hypertension vs. Control. See full table  $\rightarrow$ 

ID	Gene description	Gene symbol	Gene schematic representation	Number of methylation sites	Methylation sites in exons	Methylation sites in 5' region
ENSG00000204956	protocadherin gamma subfamily A, 1	PCDHGA1		41	5	0
ENSG00000250349	novel proline rich Gla (G- carboxyglutamic acid) 1 (PRRG1) and tetraspanin 7 (TSPAN7) protein	ENSG00000250349		38	2	2
ENSG0000081853	protocadherin gamma subfamily A, 2	PCDHGA2		37	5	0
ENSG00000204970	protocadherin alpha 1	PCDHA1		34	0	0
ENSG00000254245	protocadherin gamma subfamily A, 3	PCDHGA3		34	5	0
ENSG00000204969	protocadherin alpha 2	PCDHA2		33	3	0
ENSG00000254221	protocadherin gamma subfamily B, 1	PCDHGB1		31	5	0
ENSG00000255408	protocadherin alpha 3	PCDHA3		30	2	0
ENSG00000204967	protocadherin alpha 4	PCDHA4		28	1	0
ENSG00000204965	protocadherin alpha 5	PCDHA5		27	3	1

### 3.2. Functional classification of genes

A functional analysis of highly methylated genes was done by mapping the genes to several known ontologies, such as Gene Ontology (GO), disease ontology (based on HumanPSD<sup>™</sup> 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.

### Highly methylated genes in Hypertension vs. Control:

300 top methylated genes were taken for the mapping.

GO (biological process)

biological\_process Gene Ontology treemap



Figure 2. Enriched GO (biological process) of highly methylated genes in Hypertension vs. Control. Full classification  $\rightarrow$ 

#### TRANSPATH® Pathways (2023.1)





#### HumanPSD(TM) disease (2023.1)



Neurodevelopmental Disorders Signs and Symptoms Neurologic Manifestations

Genetic Diseases, X-Linked Behavior and Behavior Mechanisms Intellectual Disability

Neurobehavioral Manifestations Heredodegenerative Disorders, Nervous System

Mental Retardation, X-Linked Cholestasis, Intrahepatic Liver Cirrhosis, Biliary

Figure 4. Enriched HumanPSD(TM) disease (2023.1) of highly methylated genes in Hypertension vs. Control. The size of the bars correspond to the number of biomarkers of the given disease found among the input set.

#### Full classification –

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



Highly methylated genes in Hypertension vs. Control hits

-- Highly methylated genes in Hypertension vs. Control -log10(P-value)

### 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 Epigenomics data from the track(s) "Methylation track" to predict positions of potential **enhancers** regulating the highly methylated genes revealed by comparative epigenomics 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 areas of the highly methylated 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 (highly methylated genes in Hypertension vs. Control).

To build the most specific composite modules we choose top methylated genes as the input of CMA algorithm. The obtained CMA model is then applied to compute CMA score for all highly methylated genes in Hypertension vs. Control.

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

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



Model score (-p\*log10(pval)): 30.20 Wilcoxon p-value (pval): 4.38e-57 Penalty (p): 0.536 Average yes-set score: 9.54 Average no-set score: 8.19 AUC: 0.73 Separation point: 9.04 False-positive: 28.63% False-negative: 33.39% The AUC of the model achieves value significantly higher than expected for a random set of regulatory regions Z-score = 3.00



📕 No-set 📕 Yes-set — Separation point

Table 3. List of top ten highly methylated genes in Hypertension vs. Control 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  $\rightarrow$ 

Ensembl IDs	Gene symbol	Gene description	CMA score	Factor names
ENSG00000188385	JAKMIP3	Janus kinase and microtubule interacting protein 3	16.01	hHR21(h), GTF3C2(h), NF-1C(h), BCL-6(h), c- Ets-1(h), DBP(h), ZBTB7C(h)
ENSG00000102057	KCND1	potassium voltage-gated channel subfamily D member 1	14.51	NF-1C(h), ZBTB7C(h), GCMa(h), c-Ets-1(h), BCL-6(h), DBP(h), hHR21(h)
ENSG00000154358	OBSCN	obscurin, cytoskeletal calmodulin and titin-interacting RhoGEF	14.39	hHR21(h), BCL-6(h), GTF3C2(h), NF-1C(h), ZBTB7C(h), DBP(h), GCMa(h)
ENSG00000196593	ANKRD20A19P	ankyrin repeat domain 20 family member A19, pseudogene	14.35	DBP(h), GCMa(h), ZBTB7C(h), c-Ets-1(h), GTF3C2(h), hHR21(h), BCL-6(h)
ENSG00000229797		novel transcript	14.35	hHR21(h), NF-1C(h), DBP(h), c-Ets-1(h), GTF3C2(h), BCL-6(h), GCMa(h)
ENSG00000151322	NPAS3	neuronal PAS domain protein 3	14.23	DBP(h), GCMa(h), c-Ets-1(h), ZBTB7C(h), hHR21(h), BCL-6(h), NF-1C(h)
ENSG00000261308	FIGNL2	fidgetin like 2	14.22	DBP(h), GCMa(h), ZBTB7C(h), c-Ets-1(h), GTF3C2(h), hHR21(h), BCL-6(h)
ENSG00000228156		novel transcript, antisense to NOL4L	14.15	NF-1C(h), hHR21(h), GTF3C2(h), BCL-6(h), ZBTB7C(h), GCMa(h), DBP(h)
ENSG00000113504	SLC12A7	solute carrier family 12 member 7	14.01	BCL-6(h), DBP(h), GCMa(h), ZBTB7C(h), c-Ets- 1(h), GTF3C2(h), NF-1C(h)
ENSG00000101438	SLC32A1	solute carrier family 32 member 1	13.94	ZBTB7C(h), DBP(h), c-Ets-1(h), hHR21(h), GCMa(h), GTF3C2(h), NF-1C(h)

On the basis of the enhancer models we identified transcription factors potentially regulating the *target genes* of our interest. We found 8 transcription factors controlling expression of highly methylated genes in Hypertension vs. Control (see Table 4).

Table 4. Transcription factors of the predicted enhancer model potentially regulating the highly methylated genes (highly methylated genes in Hypertension vs. Control). **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**  $\rightarrow$ 

ID	Gene symbol	Gene description	Regulatory score	Yes-No ratio
MO000059013	ETS1	ETS proto-oncogene 1, transcription factor	2.82	1.27
MO000042938	RAD21	RAD21 cohesin complex component	2.19	1.21
MO000026319	BCL6	BCL6 transcription repressor	1.96	3
MO000028669	DBP	D-box binding PAR bZIP transcription factor	1.8	1.12
MO000057363	GTF3C2	general transcription factor IIIC subunit 2	1.74	1.21
MO000026306	GCM1	glial cells missing transcription factor 1	1.61	4
MO000024750	NFIC	nuclear factor I C	0	1.43
MO000225595	ZBTB7C	zinc finger and BTB domain containing 7C	0	1.25

The following diagram represents the key transcription factors, which were predicted to be potentially regulating highly methylated genes in the analyzed pathology: ETS1, RAD21 and BCL6.



## 3.4. Finding master regulators in networks

In the second step of the upstream analysis common regulators of the revealed TFs were identified. These master regulators appear to be the key candidates for therapeutic targets as they have a master effect on regulation of intracellular pathways that activate the pathological process of our study. The identified master regulators are shown in Table 5.

Table 5. Master regulators that may govern the regulation of highly methylated genes in Hypertension vs. Control. **Total rank** is the sum of the ranks of the master molecules sorted by keynode score, CMA score, epigenomics data. See full table  $\rightarrow$ 

ID	Master molecule	Gene symbol	Gene description	Total rank
MO000097095	AR-isoform1(h){aceK630} {aceK632}{aceK633}	AR	androgen receptor	32
MO000021454	AR(h)	AR	androgen receptor	33
MO000058849	AR-isoform1(h)	AR	androgen receptor	33
MO000078001	AR(h){sumo}	AR	androgen receptor	33
MO000078150	AR(h){ub}	AR	androgen receptor	33
MO000162563	AR-isoform2(h)	AR	androgen receptor	33
MO000342259	AR-isoform3(h)	AR	androgen receptor	33
MO000342260	AR-isoform4(h)	AR	androgen receptor	33
MO000058267	RSK(h){p}	RPS6KA1, RPS6KA2, RPS6KA3	ribosomal protein S6 kinase A1, ribosomal protein S6 kinase A2, ribosomal protein S6 kinase A3	71
MO000019070	XIAP(h)	XIAP	X-linked inhibitor of apoptosis	80

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 highly methylated genes, and selected master regulators, which are responsible for the regulation of these TFs.



Figure 5. Diagram of intracellular regulatory signal transduction pathways of highly methylated genes in Hypertension vs. Control. 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. **See full diagram**  $\rightarrow$ 

# 4. Finding prospective drug targets

The identified master regulators that may govern pathology associated genes were checked for druggability potential using HumanPSD<sup>M</sup> [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<sup>M</sup> 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<sup>™</sup> 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 Methods section for the details) such master regulator proteins were not used in further analysis of drug prediction.

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

Table 6. Prospective drug targets selected from full list of identified master regulators filtered by Druggability score from HumanPSD<sup>TM</sup> 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**  $\rightarrow$ 

Gene symbol	Gene Description	Druggability score	<b>Total rank</b>
AR	androgen receptor	70	33
TSC22D3	TSC22 domain family member 3	3	99
TSPYL2	TSPY like 2	1	187
F8	coagulation factor VIII	4	289
LCK	LCK proto-oncogene, Src family tyrosine kinase	45	291
IKBKG	inhibitor of nuclear factor kappa B kinase regulatory subunit gamma	1	301

Table 7. 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.

Gene symbol	Gene Description	Druggability score	Total rank
AR	androgen receptor	28.53	33
DUSP9	dual specificity phosphatase 9	35.02	140
EFNB1	ephrin B1	1.08	250
LCK	LCK proto-oncogene, Src family tyrosine kinase	11.04	291
HDAC4	histone deacetylase 4	1.34	311
PTPRH	protein tyrosine phosphatase receptor type H	16.75	313

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:

- TSC22D3
- MKP-4
- AR-isoform1
- AR

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, Dexamethasone, Norelgestromin, Estrone, Progesterone, Drospirenone and Ulipristal, should be considered as a prospective research initiative for further drug repurposing and drug development. These drugs were selected as top matching treatments to the most prospective drug targets of the studied pathology, however, these results should be considered with special caution and are to be used for research purposes only, as there is not enough clinical information for adapting these results towards immediate treatment of patients.

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

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

# 5. Identification of potential drugs

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

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

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

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

- Target activity score (depends on ranks of all targets that were found for the selected drug);
- Disease activity score (weighted sum of number of clinical trials on disease(s) under study where the selected drug is
  known to be applied or PASS Disease activity score cheminformatically predicted property of the compound to be active
  against the studied disease(s));
- Clinical validity score (applicable only for drugs predicted on the basis of literature curation in HumanPSD<sup>™</sup> database (Tables 8 and 9), reflects the number of the highest clinical trials phase on which the drug was tested for any pathology).

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

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

are calculated as a quotient of maximum drug rank and the drug rank of the given drug multiplied by 100.

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

## Drugs approved in clinical trials



Table 8. 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<sup>TM</sup> database) See full table  $\rightarrow$ 

Name	Target names	Drug score	Disease activity score	Disease trial phase
Imatinib	MAPK10, RPS6KA3, BMPR1A, MARK3, PRKAA2, KDR, INSR, LTK, EGFR, ACVR2A, PRKG1, ERBB2, CAMK2B, PRKD1, PKMYT1, IGF1R, BMX, LCK, PIM2, EPHB2, RPS6KA2, MAPK4, CSNK1G2, PAK3, BTK, CAMKK1, MAPK11, TGFBR2, FYN, FES, RET, STK3	94	7	Phase 3: Hypertension, Astrocytoma, Bone Marrow Diseases, COVID-19, Cerebral Infarction, Familial Primary Pulmonary Hypertension, Fibrosarcoma, Fibrosis, Gastrointestinal Stromal Tumors, Glioblastoma, Graft vs Host Disease, Hematologic Diseases, Hypertension, Pulmonary, Idiopathic Pulmonary Fibrosis, Infarction, Ischemia, Ischemic Stroke, Leukemia, Leukemia, Biphenotypic, Acute, Leukemia, Lymphoid, Leukemia, Myelogenous, Chronic, BCR-ABL Positive, Leukemia, Myeloid, Leukemia, Myeloid, Chronic-Phase, Lung Diseases, Lymphoma, Lymphoma, Non-Hodgkin, Mucositis, Neoplasms, Nephrogenic Fibrosing Dermopathy, Nerve Sheath Neoplasms, Neurilemmoma, Neurofibrosarcoma, Pneumonia, Precursor Cell Lymphoblastic Leukemia-Lymphoma, Precursor T-Cell Lymphoblastic Leukemia-Lymphoma, Pulmonary Arterial Hypertension, Pulmonary Fibrosis, Recurrence, Sarcoma, Severe Acute Respiratory Syndrome, Stroke
Sorafenib	MAPK10, RPS6KA3, BMPR1A, MARK3, PRKAA2, KDR, INSR, LTK, EGFR, ACVR2A, PRKG1, ERBB2, CAMK2B, PRKD1, PKMYT1, IGF1R, BMX, LCK, PIM2, EPHB2, RPS6KA2, MAPK4, CSNK1G2, PAK3, PRKCZ, BTK, CAMKK1, MAPK11, TGFBR2, FYN, FES, RET, STK3	92	3	Phase 2: Hypertension, Acute Disease, Adenocarcinoma, Adenocarcinoma of Lung, Adenocarcinoma, Follicular, Adenoma, Adenoma, Liver Cell, Adrenal Cortex Neoplasms, Brain Abscess, Brain Neoplasms, Breast Neoplasms, Breast Neoplasms, Male, Carcinoid Tumor, Carcinoma, Carcinoma, Ductal, Carcinoma, Nearotocarinoma, Carcinoma, Islet Cell, Carcinoma, Medullary, Carcinoma, Neuroendocrine, Carcinoma, Non-Small-Cell Lung, Carcinoma, Neuroina, Carcinoma, Renal Cell, Carcinoma, Small Cell, Carcinoma, Squamous Cell, Carcinoma, Renal Cell, Carcinoma, Small Cell, Carcinoma, Squamous Cell, Carcinoma, Transitional Cell, Carcinoma, Colonic Neoplasms, Colorectal Neoplasms, Desmoplastic Small Round Cell Tumor, Digestive System Neoplasms, Disease Progression, Endocrine Gland Neoplasms, Esophageal Neoplasms, Galoractal Neoplasms, Fibroma, Fibrosarcoma, Fibrosis, Gallbladder Neoplasms, Gastrinoma, Glasarconma, Glucagonoma, Head and Neck Neoplasms, Glioblastoma, Gliosarcoma, Glucagonoma, Head and Neck Neoplasms, Hemangiosarcoma, Hepatitis, Hepatitis A, Hepatitis B, Hepatitis C, Hepatoblastoma, Hepatopulmonary Syndrome, Histiocytoma, Histiocytoma, Benign Fibrous, Histiocytoma, Malignant Fibrous, Hypertension, Portal, Hypopharyngeal Neoplasms, Immunoblastic Lymphadenopathy, Insulinoma, Intestinal Neoplasms, Keloid, Kidney Diseases, Kidney Neoplasms, Klatskin Tumor, Laryngeal Diseases, Laryngeal Neoplasms, Ieeiomyosarcoma, Leukemia, Leukemia, Biphenotytic, Acute, Leukemia, Myeloomocytic, Chronic, BCR-ABL Positive, Leukemia, Myeloid, Leukemia, Myeloid, Acute, Leukemia, Myelomonocytic, Chronic, Leukemia, Myeloid, Leukemia, Myeloid, Acute, Leukemia, Myelomonocytic, Chronic, Leukemia, Myeloomonocytic, Juvenile, Leukemia, Promyelocytic, Acute, Leukemia, Ture, Byelos, Lymphadenopathy, Lymphoma, Non-Hodgkin, Lymphoma, T-Cell, Lymphoma, T-Cell, Cutaneous, Lymphoma, Non-Hodgkin, Lymphoma, T-Cell, Lymphoma, T-Cell, Cutaneous, Lymphoma, Non-Hodgkin, Lymphoma, T-Cell, Lymphoma, T-Cell, Cutaneous, Lymphoma, Non-Hodgkin, Lymphoma, T-Cell, Lymphoma, See

				Carcinoma, Somatostatinoma, Squamous Cell Carcinoma of Head and Neck, Stomach Neoplasms, Syndrome, Testicular Neoplasms, Thrombosis, Thyroid Cancer, Papillary, Thyroid Carcinoma, Anaplastic, Thyroid Diseases, Thyroid Neoplasms, Tongue Neoplasms, Triple Negative Breast Neoplasms, Ureteral Neoplasms, Urethral Neoplasms, Urinary Bladder Neoplasms, Uterine Cervical Neoplasms, Uveal Neoplasms, Vaccinia, Vipoma, Wilms Tumor
Drospirenone	PGR, AR	92	5	Phase 3: Hypertension, Acne Vulgaris, Cysts, Dysbiosis, Essential Hypertension, Hot Flashes, Hypersensitivity, Insulin Resistance, Neural Tube Defects, Polycystic Ovary Syndrome, Spinal Dysraphism, Syndrome
Pazopanib	MAPK10, RPS6KA3, BMPR1A, MARK3, PRKAA2, KDR, INSR, LTK, EGFR, ACVR2A, PRKG1, ERBB2, CAMK2B, PRKD1, PKMYT1, IGF1R, BMX, LCK, PIM2, RPS6KA2, MAPK4, CSNK1G2, PAK3, BTK, CAMKK1, MAPK11, TGFBR2, FYN, FES, RET, STK3	91	2	Phase 2: Hypertension, Adenocarcinoma, Anemia, Brain Abscess, Brain Neoplasms, Breast Neoplasms, Breast Neoplasms, Male, Carcinoid Tumor, Carcinoma, Carcinoma, Islet Cell, Carcinoma, Ovarian Epithelial, Carcinoma, Renal Cell, Carcinoma, Non-Small- Cell Lung, Carcinoma, Ovarian Epithelial, Carcinoma, Renal Cell, Carcinoma, Squamous Cell, Carcinoma, Transitional Cell, Carcinoma, Central Nervous System Neoplasms, Cholangiocarcinoma, Chondrosarcoma, Chondrosarcoma, Mesenchymal, Corneal Neovascularization, Dermatofibrosarcoma, Desmoplastic Small Round Cell Tumor, Dilatation, Pathologic, Drug-Related Side Effects and Adverse Reactions, Edema, Endocrine Gland Neoplasms, Epistaxis, Fallopian Tube Neoplasms, Fibroma, Fibromatosis, Aggressive, Fibrosarcoma, Gallbladder Neoplasms, Gastrinnea, Gastrointestinal Neoplasms, Gastrointestinal Stromal Tumors, Glioblastoma, Glioma, Gliosarcoma, Glomus Tumor, Glucagonoma, Granular Cell Tumor, Hemangioendothelioma, Hemangioendothelioma, Epithelioid, Hemangiopericytoma, Hemangiosarcoma, Hemorrhage, Histiocytoma, Histiocytoma, Benign Fibrous, Histiocytoma, Malignant Fibrous, Hypersensitivity, Inflammatory Breast Neoplasms, Insulinoma, Intestinal Neoplasms, Kidney Neoplasms, Leiomyosarcoma, Leukemia, Leukemia, Myeloid, Leukemia, Myeloid, Acute, Liposarcoma, Lung Neoplasms, Lymphedema, Macular Degeneration, Malignant Carcinoid Syndrome, Melanoma, Mesothelioma, Mesothelioma, Malignant, Mixed Tumor, Multiple Endocrine Neoplasia, Multiple Endocrine Neoplasis, Neoplasms, Fibrous Tissue, Neoplasms, Germ Cell and Embryonal, Neovascularization, Pathologic, Nerve Sheath Neoplasms, Nervous System Neoplasms, Neuroendocrine Tumors, Neurofibrosarcoma, Osteosarcoma, Ovarian Neoplasms, Pancreatic Neoplasms, Sarcoma, Aveolar Soft Part, Sarcoma, Clear Cell, Sarcoma, Synovial, Small Cell Lung Carcinoma, Solitary Fibrous Tumors, Somatostatinoma, Squamous Cell Carcinoma of Head and Neck, Stomach Neoplasms, Syndrome, Telangiectasia, Hereditary Hemorrhagic, Telangiectasis, Thyroid Carcinoma, Anaplastic, Th
estradiol benzoate	BDNF, AR	87	5	Phase 3: Hypertension, Alzheimer Disease, Amenorrhea, Anorexia, Anorexia Nervosa, Arnold-Chiari Malformation, Atrophic Vaginitis, Atrophy, Breast Neoplasms, Congenital Abnormalities, Endometriosis, Essential Hypertension, Fibroma, Genetic Diseases, Inborn, Gonadal Dysgenesis, Hemorrhage, Hot Flashes, Hypogonadism, Infertility, Infertility, Female, Leiomyoma, Menopause, Premature, Menorrhagia, Metrorrhagia, Migraine Disorders, Myofibroma, Myoma, Neoplasms, Pain, Premature Birth, Primary Ovarian Insufficiency, Syndrome, Turner Syndrome, Uterine Hemorrhage, Vaginitis

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

# <u>Repurposing drugs</u>



Table 9. Repurposed drugs used in clinical trials for other pathologies (prospective drugs against the identified drug targets on the basis of literature curation in HumanPSD<sup>m</sup> database)

See	full table $\rightarrow$		
Name	Target names	Drug score	Maximum trial phase
Flavopiridol	MAPK10, RPS6KA3, CDK6, BMPR1A, MARK3, PRKAA2, KDR, INSR, LTK, EGFR, ACVR2A, PRKG1, ERBB2, CAMK2B, PRKD1, PKMYT1, IGF1R, BMX, LCK, PIM2, EPHB2, RPS6KA2, MAPK4, CSNK1G2, PAK3, BTK, CAMKK1, XIAP, MAPK11, TGFBR2, FYN, FES, RET, STK3	86	Phase 2: Adenocarcinoma, Brain Abscess, Breast Neoplasms, Carcinoma, Hepatocellular, Carcinoma, Ovarian Epithelial, Carcinoma, Renal Cell, Embolism, Endometrial Neoplasms, Esophageal Neoplasms, Germinoma, Granuloma, Head and Neck Neoplasms, Hodgkin Disease, Hypereosinophilic Syndrome, Immunoblastic Lymphadenopathy, Kidney Neoplasms, Leukemia, Leukemia, Basophilic, Acute, Leukemia, Eosinophilic, Acute, Leukemia, Erythroblastic, Acute, Leukemia, Lymphocytic, Chronic, B-Cell, Leukemia, Lymphoid, Leukemia, Megakaryoblastic, Acute, Leukemia, Monocytic, Acute, Leukemia, Myeloid, Leukemia, Myeloid, Acute, Leukemia, Myelomonocytic, Acute, Leukemia, Prolymphocytic, Leukemia, T-Cell, Leukemia-Lymphoma, Adult T-Cell, Liver Neoplasms, Lymphadenopathy, Lymphatic Diseases, Lymphoma, Large B-Cell, Diffuse, Lymphoma, B-Cell, Marginal Zone, Lymphoma, Mantle-Cell, Lymphoma, Non-Hodgkin, Lymphoma, T-Cell, Lymphoma, T-Cell, Cutaneous, Lymphomatoid Granulomatosis, Melanoma, Multiple Myeloma, Mycoses, Mycosis Fungoides, Myelodysplastic Syndromes, Neoplasms, Neoplasms, Germ Cell and Embryonal, Neoplasms, Plasma Cell, Ovarian Neoplasms, Pancreatic Neoplasms, Peritoneal Neoplasms, Prostatic Neoplasms, Castration-Resistant, Recurrence, Sarcoma, Seminoma, Sezary Syndrome, Stomach Neoplasms, Testicular Neoplasms, Thromboembolism, Waldenstrom Macroglobulinemia
ruboxistaurin	MAPK10, RPS6KA3, BMPR1A, MARK3, PRKAA2, KDR, INSR, LTK, EGFR, ACVR2A, PRKG1, ERBB2, CAMK2B, PRKD1, PKMYT1, IGF1R, BMX, LCK, PIM2, EPHB2, PRKCG, RPS6KA2, MAPK4, CSNK1G2, PAK3, PRKCZ, BTK, CAMKK1, MAPK11, TGFBR2, FYN, FES. RET. STK3	86	Phase 3: Diabetes Mellitus, Diabetes Mellitus, Type 1, Diabetes Mellitus, Type 2, Diabetic Neuropathies, Diabetic Retinopathy, Edema, Macular Edema, Nervous System Diseases, Peripheral Nervous System Diseases, Retinal Diseases
bms-387032	MAPK10, RPS6KA3, BMPR1A, MARK3, PRKAA2, KDR, INSR, LTK, EGFR, ACVR2A, PRKG1, ERBB2, CAMK2B, PRKD1, PKMYT1, IGF1R, BMX, LCK, PIM2, EPHB2, PRKCG, RPS6KA2, MAPK4, CSNK1G2, PAK3, BTK, CAMKK1, MAPK11, TGFBR2, FYN, FES, RET, STK3	86	Phase 1: Leukemia, Lymphocytic, Chronic, B-Cell, Lymphoma, Mantle-Cell, Multiple Myeloma, Neoplasms
Dasatinib	MAPK10,	86	Phase 4: Hematologic Neoplasms, Leukemia, Leukemia, Lymphoid, Leukemia, Myelogenous,

	RPS6KA3, BMPR1A, MARK3, PRKAA2, KDR, INSR, LTK, EGFR, ACVR2A, PRKG1, ERBB2, CAMK2B, PRKD1, PKMYT1, IGF1R, BMX, LCK, PIM2, EPHB2, RPS6KA2, MAPK4, CSNK1G2, PAK3, PRKCZ, BTK, CAMKK1, MAPK11, TGFBR2, FYN, FES, RET, STK3		Chronic, BCR-ABL Positive, Leukemia, Myeloid, Leukemia, Myeloid, Chronic-Phase, Lymphoma, Neoplasms, Precursor Cell Lymphoblastic Leukemia-Lymphoma
:I-1033	MAPK10, RPS6KA3, BMPR1A, MARK3, PRKAA2, KDR, INSR, LTK, EGFR, ACVR2A, PRKG1, ERBB2, CAMK2B, PRKD1, PKMYT1, IGF1R, BMX, LCK, PIM2, EPHB2, RPS6KA2, MAPK4, CSNK1G2, ERBB4, PAK3, BTK, CAMKK1, MAPK11, TGFBR2, FYN, FES, RET, STK3	86	Phase 2: Breast Neoplasms, Carcinoma, Non-Small-Cell Lung, Lung Neoplasms, Neoplasms

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

Table 10. 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  $\rightarrow$ 

Name	Target names	Drug score	Target activity score
Moexipril	BDNF, ITGB3	87	0.17
Quinapril	BDNF, ITGB3	84	0.13
Nicergoline	PRKAA2, PRKG1, PKMYT1	81	8.58E-2
Darodipine	RPS6KA3, RPS6KA2	71	4.71E-2
Hesperetin	MAPK10, BDNF, MAPK4, MAPK11	71	4.86E-2



C

Table 11. 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  $\rightarrow$ 

Name	Target names	Drug score	Target activity score
2,5,7-Trihydroxynaphthoquinone	MAPK10, DUSP22, MAPK4, CDC14B, MAPK11, DUSP14, DUPD1, DUSP9, BRCA1	86	0.25
3-METHYL-1,6,8- TRIHYDROXYANTHRAQUINONE	MAPK10, DUSP22, MAPK4, CDC14B, MAPK11, DUSP14, DUPD1, DUSP9, BRCA1	86	0.21
[[N- (Benzyloxycarbonyl)Amino]Methyl]Phosphate	PTPRR, ACE2, GRIN1, DUSP22, PTPRO, PTPRF, PTPN2, PTPN13, PTPRH, PTPRD, CDC14B, DUSP14, DUPD1, DUSP9, PTPN21	85	0.56
Fenoldopam	MAPK10, GRIN1, MAPK4, MAPK11	82	0.14
histamine dihydrochloride	BMX, LCK, MAPK10, EPHB2, MAPK4, KDR, ERBB4, INSR, LTK, EGFR, BTK, MAPK11, ERBB2, FYN, RET, FES, IGF1R	82	0.62

As the result of drug search we propose the following drugs as most promising candidates for treating the pathology under study: Imatinib, Flavopiridol, Moexipril and 2,5,7-Trihydroxynaphthoquinone. These drugs were selected for acting on the following targets: LCK, ITGB3 and DUSP9, 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 *epigenomics* data obtained from *blood* tissue. The study is done in the context of *Hypertension*. The data were pre-processed, statistically analyzed and highly methylated 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:



These drugs were selected for acting on the following targets: LCK, ITGB3 and DUSP9, 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:



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, Dexamethasone, Norelgestromin, Estrone, Progesterone, Drospirenone and Ulipristal. 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 highly methylated 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.

- TSC22D3
- MKP-4
- AR-isoform1
- AR

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 highly methylated genes were analyzed using known DNAbinding motifs described in the TRANSFAC® library, release 2023.1 (geneXplain GmbH, Wolfenbüttel, Germany) (https://genexplain.com/transfac).

The master regulator search uses the TRANSPATH® database (BIOBASE), release 2023.1 (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<sup>™</sup> database, release 2023.1 (https://genexplain.com/humanpsd).

The Ensembl database release Human104.38 (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.

#### **Epigenomics data processing**

When analyzing a list of CpG sites, we compute the fold change values between the methylation status in the studied pathology and the control set. Top 10 000 CpG sites with highest logFC values are taken to further analysis. These sites are mapped to corresponding genes, which will be further compared to the list of housekeeping genes at the step of promoter analysis.

#### 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 enhancers under study as compared to a background set of promoters of housekeeping genes. 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<sup>™</sup> and predicting potential drugs using PASS program.

Method for analysis of known pharmaceutical compounds

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

- 1. ranking by "Target activity score" (T-score<sub>PSD</sub>),
- 2. ranking by "Disease activity score" (*D*-score<sub>PSD</sub>),
- 3. ranking by "Clinical validity score".

"Target activity score" (*T*-score<sub>PSD</sub>) 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<sub>PSD</sub>):

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

where *D* is the set of selected diseases, and if *D* is empty set, D-score<sub>PSD</sub>=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.

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

#### Method for prediction of pharmaceutical compounds

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

We selected compounds that satisfied the following conditions:

- 1. Toxicity below a chosen toxicity threshold (defines as Pa, probability to be active as toxic substance).
- 2. For all predicted pharmacological effects that correspond to a set of user selected disease(s) *Pa* is greater than a chosen effect threshold.

3. There are at least 2 targets (corresponding to the predicted activity-mechanisms) with predicted *Pa* greater than a chosen target threshold.

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

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

where M(s) is the set of activity-mechanisms for the given structure (which passed the chosen threshold for activitymechanisms 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
- 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.
- 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
- 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.
- 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
- 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.
- Filimonov DA, Poroikov VV. Prognosis of specters of biological activity of organic molecules. Russian chemical journal. 2006;50(2):66-75 (russ)
- 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 <a href="mailto:support@genexplain.com">support@genexplain.com</a>

## Supplementary material

- 1. Supplementary table 1 Detailed report. Composite modules and master regulators (highly methylated genes in Hypertension vs. Control).
- 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<sup>™</sup> 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.