# Supplementary Tables and Figures for "Genome-Wide Prediction of cis-Regulatory Regions Using Supervised Deep Learning Methods"

Yifeng Li[*1,2], Wenqiang Shi[†1], and Wyeth W. Wasserman[‡1]

[1]Centre for Molecular Medicine and Therapeutics, BC Children's Hospital Research Institute, Department of Medical Genetics, University of British Columbia
[2]Digital Technologies Research Centre, National Research Council Canada

Table S1: Numbers of labelled regions in our data.

| Cell | A-E | I-E | A-P | I-P | A-X | I-X | UK | Total |
|---|---|---|---|---|---|---|---|---|
| A549 | 387 | 40,387 | 10,907 | 128,998 | 8,998 | 14,712 | 81,217 | 285,606 |
| GM12878 | 2,878 | 28,156 | 10,816 | 73,891 | 8,226 | 19,078 | 80,004 | 223,049 |
| HelaS3 | 1847 | 32,179 | 10,759 | 79,009 | 9,123 | 22,071 | 81,502 | 236,485 |
| HepG2 | 1465 | 34,556 | 11,467 | 96,184 | 9,931 | 19,071 | 79,417 | 252,091 |
| HUVEC | 1226 | 35,143 | 11,254 | 101,861 | 9,739 | 18,249 | 80,333 | 257,805 |
| K562 | 894 | 34,392 | 10,076 | 82,829 | 9,033 | 20,261 | 78,081 | 235,566 |
| MCF7 | 249 | 36,873 | 10,733 | 81,510 | 10,829 | 13,653 | 82,663 | 236,510 |

Table S2: Number of features used for each cell type and number of common features between any two cell types in our labelled data.

| Cell | A549 | GM12878 | HelaS3 | HepG2 | HMEC | HUVEC | K562 | MCF7 |
|---|---|---|---|---|---|---|---|---|
| A549 | 45 | 38 | 26 | 38 | 15 | 19 | 38 | 25 |
| GM12878 | - | 101 | 50 | 54 | 16 | 21 | 69 | 26 |
| HelaS3 | - | - | 74 | 43 | 16 | 22 | 56 | 23 |
| HepG2 | - | - | - | 72 | 16 | 21 | 57 | 28 |
| HMEC | - | - | - | - | 16 | 16 | 16 | 9 |
| HUVEC | - | - | - | - | - | 24 | 24 | 13 |
| K562 | - | - | - | - | - | - | 135 | 29 |
| MCF7 | - | - | - | - | - | - | - | 38 |

[*]E-mail address: `yifeng@cmmt.ubc.ca`, and `yifeng.li@nrc-cnrc.gc.ca`.
[†]E-mail address: `shi@cmmt.ubc.ca`.
[‡]Corresponding author, E-mail address: `wyeth@cmmt.ubc.ca`.

Table S3: Confusion matrices of classifying CRE-seq and MPRA validated regions using DECRES.

| Original Class\DECRES Class | A-E | A-P | BG | Total |
|---|---|---|---|---|
| CRE-seq Positive Combined Predicted Enhancer in K562 | 254 | 0 | 132 | 386 |
| CRE-seq Negative Combined Predicted Enhancer in K562 | 433 | 1 | 378 | 812 |
| Combined Predicted Repressed Region in K562 | 0 | 0 | 298 | 298 |
| Total | 687 | 1 | 808 | 1496 |
| MPRA Positive Enhancer in K562 | 120 | 0 | 2 | 122 |
| MPRA Negative Enhancer in K562 | 179 | 0 | 15 | 194 |
| Total | 299 | 0 | 17 | 316 |
| MPRA Positive Enhancer in HepG2 | 182 | 1 | 3 | 186 |
| MPRA Negative Enhancer in HepG2 | 217 | 5 | 45 | 267 |
| Total | 399 | 6 | 48 | 453 |

Table S4: Numbers of predicted and cell-specific *cis*-regulatory regions in the whole human genome. Columns 2-3: results of two-class prediction. Columns 4-7: results of three class prediction.

| Cell | A-E+A-P | Specific | A-E | Specific | A-P | Specific |
|---|---|---|---|---|---|---|
| GM12878 | 90,192 | 27,904 | 70,905 | 27,185 | 19,287 | 6,090 |
| HelaS3 | 100,102 | 22,268 | 92,509 | 22,866 | 7,593 | 979 |
| HepG2 | 114,873 | 35,986 | 105,007 | 41,109 | 9,866 | 2,135 |
| HMEC | 104,621 | 28,774 | 88,803 | 26,262 | 15,818 | 3,043 |
| HUVEC | 110,347 | 16,415 | 97,069 | 35,073 | 13,278 | 582 |
| K562 | 133,940 | 30,835 | 122,321 | 56,374 | 11,619 | 371 |

Table S5: Numbers of predicted A-Es and A-Ps on the 102,021 BDT loci. AiA: Active in the FANTOM Enhancer Atlas. IiA: Inactive in the FANTOM Enhancer Atlas. NiA: Not included in the FANTOM Enhancer Atlas. Specific: Predicted cell-specific A-Es.

| Cell | A-E | A-P | BG |
|---|---|---|---|
| GM12878 | 11,910 | 1,975 | 88,136 |
| HelaS3 | 12,743 | 226 | 89,052 |
| HepG2 | 10,761 | 488 | 90,772 |
| HMEC | 13,356 | 1,267 | 87,398 |
| HUVEC | 17,192 | 758 | 84,071 |
| K562 | 13,936 | 288 | 87,797 |

Table S6: Overlap between our predicted A-Es and the FANTOM enhancer atlas.

| Cell\Element | AiA | IiA | NiA | Specific |
|---|---|---|---|---|
| GM12878 | 2,088 | 5,144 | 4,678 | 4,199 |
| HelaS3 | 1,459 | 6,435 | 4,849 | 1,073 |
| HepG2 | 1,178 | 4,888 | 4,695 | 1,495 |
| HMEC | 1,250 | 6,644 | 5,462 | 1,916 |
| HUVEC | 879 | 8,716 | 7,597 | 3,560 |
| K562 | 750 | 7,383 | 5,803 | 2,965 |

Table S7: Transcription factors whose binding motifs are enriched in specific cells.

| Cell | Transcription Factor | Functionality |
|---|---|---|
| GM12878 | RUNX1 | RUNX1 and other Runt-related factors play crucial roles in haematopoiesis. Translocation of RUNX1 leads to severe acute myeloid leukemia [11]. |
| GM12878 | REL/NF-$\kappa$B Factors | Play a critical role in immune response to infection [4]. |
| HelaS3 | C/EBP-Related Factors | C/EBP-related factors, a subfamily of the basic leucine zipper factors (bZIP), regulating genes involved in immune and inflammatory responses, are expressed in cervix [1]. |
| HelaS3 | CREB-Related Factors (Another bZIP Subfamily) | Possibly up-regulate Bcl-2 expression in apoptotic HeLa cells induced by trichosanthin [15]. |
| HelaS3 | TEAD1 | Plays a role of apoptotic resistance in Hela cells [9]. |
| HelaS3 | AP-2 Factors | Can act as tumor suppressor, and malfunction of AP-2 was found in cervical cancer cells [2]. |
| HelaS3 | HOX Factors | a subgroup of HOX genes involve in cervical carcinoma [8]. |
| HepG2 | HNF1A, HNF1B, FOXA1, FOXA2, HNF4A, and HNF4G | These factors are hepatocyte nuclear factors that regulate liver-specific genes [3, 16]. |
| HepG2 | C/EBPs | Have a pivotal role in liver development and function [13]. |
| HMEC | TEAD Family | Members of this family regulate epithelial-mesenchymal transition [17]. |
| HUVEC | GATA Factors | Involve in networks of key determinants of vascular endothelial cell identity [7]. |
| HUVEC | SOX Factors | Regulate vascular cell development and growth (vasculogenesis) [12]. |
| HUVEC | CREB-Related Factors | Involve in several specific pathways in HUVEC cell [14]. |
| K562 | GATA-Type, NFYA, and NFYB USF1 | The NFY factors cooperate with GATA1 to mediate erythroid-specific transcription, and coassociate with FOS binding to both promoters and enhancers in K562 [5]. The NFY factors are also found cooperate with USF1 and USF2 to active HOXB4 for hematopoiesis [18]. |
| K562 | STAT5 | Maintains the high-level cell proliferation of K562 [10]. |

Figure S1: Mean performance and standard deviation of 10-fold cross-validations using the MLP model on our labelled data of eight cell types. A-E: Active Enhancer, A-P: Active Promoter, A-X: Active Exon, I-E: Inactive Enhancer, I-P: Inactive Promoter, I-X: Inactive Exon, UK: Unknown or Uncharacterized, BG: I-E+I-P+A-X+I-X+UK.



Figure S2: Cumulative DECRES membership probabilities (scores) of enhancers that were tested as positives and negatives by CRE-seq or MPRA. These enhancers were predicted as A-Es by DECRES (See Table S3).



Figure S3: Comparing the mean auPRCs over 100 resampling and retraining on our labelled regions using different feature sets. "Experimental" means our experimentally derived next generation sequencing feature set. "Sequence" means the set of 351 sequence properties used in [6]. "Experimental+Sequence" means the combination of these two sets. The $p$-values in each legend were obtained using two-tailed Student's t-test to compare "Experimental"-based results with "Experimental+Sequence"-based and "Sequence"-based results, respectively.

Figure S4: Feature importance and box plots of top features in the 3-class (A-E versus A-P versus BG) scenario. A: Feature importance discovered by randomized DFS (RDFS) and random forest (RF) on HelaS3, HepG2 and K562 cells. RF's feature importance scores were normalized to [0,1] for better comparison with RDFS. B: For the top 10 features of the 3-class models generated for four well-characterized cell lines, box plots depict the range of observed feature values (log2 scale) for 7 sequence classes.

Figure S5: Feature importance, classification performance, and top features in the 2-class (A-E+A-P versus BG) scenario. A: Feature importance discovered by randomized DFS (RDFS) and random forest (RF). The random forest's feature importance scores were normalized to [0,1] for better comparison with randomized DFS. B: auPRC versus the number of features incorporated into the RDFS and RF. The annotated points indicate where a line with slope 0.5 intersects a fitted curve). C: For the top 6 features of the 2-class models generated for four well-characterized cell lines, box plots depict the range of observed feature values (log2 scale) for 7 sequence classes.

Figure S6: Box plots of all features for A549.

Figure S7: Box plots of all features for GM12878.

Figure S8: Box plots of all features for HelaS3.

Figure S9: Box plots of all features for HepG2.

Figure S10: Box plots of all features for HMEC.



Figure S11: Box plots of all features for HUVEC.

Figure S12: Box plots of all features for MCF7.

Figure S13: Box plots of all features for K562.

Figure S14: An example of predicted regulatory regions in the UCSC Genome Browser.



Figure S15: Functional and motif analysis of the DECRES genome-wide predictions on cell line GM12878. A: Distance from the predicted A-Es to gene TSSs. B: Distance from the predicted A-Ps to gene TSSs. C,D,E: Top 20 enriched biological processes, pathways, and diseases, respectively, in the predicted cell-specific CRRs. F: Enriched *de novo* motifs in the predicted cell specific CRRs. Column 4: families of best-matched TFs. Column 5: best match scores.

14

Figure S16 D table:

| Motif | p-value | Best Match | Family | Score |
|---|---|---|---|---|
| | 1e-150 | FOS, JUN | FOS-related factors, Jun-related factors | 0.96 |
| | 1e-136 | Erg | Ets-related factors | 0.93 |
| | 1e-104 | REL | NFkB-related factors | 0.96 |
| | 1e-95 | RUNX1 | Runt-related factors | 0.95 |
| | 1e-66 | Irf5 | Interferon-regulatory factors | 0.80 |
| | 1e-38 | POU2F2 | POU domain factors | 0.93 |
| | 1e-38 | SP1 | Three-zinc finger Kruppel-related factors | 0.85 |
| | 1e-29 | Gm397.2 | BetaBetaAlpha-zinc finger | 0.72 |
| | 1e-28 | DCE_S_III | Unknown | 0.75 |
| | 1e-28 | Gcm1_1 | Glial cells missing (GCM) | 0.72 |
| | 1e-24 | Pax2 | Paired domain only | 0.75 |
| | 1e-24 | Bhlhb2_1 | Helix-Loop-Helix | 0.86 |
| | 1e-22 | Srf_2 | MADS | 0.75 |
| | 1e-20 | KLF5 | Three-zinc finger Kruppel-related factors | 0.71 |
| | 1e-20 | NFYB | Heteromeric CCAAT-binding factors | 0.96 |
| | 1e-18 | ETS::E-box | ETS::bHLH | 0.65 |
| | 1e-18 | USF2 | bHLH-ZIP factors | 0.68 |
| | 1e-16 | MEF2A | MADS, Regulators of differentiation | 0.86 |
| | 1e-15 | CREB1 | CREB-related factors | 0.72 |
| | 1e-14 | Foxo1 | Forkhead box factors | 0.69 |
| | 1e-14 | Mycn | bHLH-ZIP factors | 0.78 |
| | 1e-12 | Hnf1b | POU domain factors | 0.68 |
| | 1e-12 | Irf6.2 | Interferon-regulatory factors | 0.61 |

D



Figure S16: Functional and motif analysis of the DECRES NiA enhancers on cell line GM12878. A,B,C: Top 20 enriched biological processes, pathways, and diseases, respectively, in the NiA enhancers. D: Enriched *de novo* motifs in the NiA enhancer regions. Column 4: the families of best-matched TFs. Column 5: best match scores.

Figure S17 C table:

| Motif | p-value | Best Match | Family | Score |
|---|---|---|---|---|
| | 1e-2223 | CTCF | More than 3 adjacent zinc finger factors | 0.90 |
| | 1e-2062 | JUND, FOS | Jun-related factors, Fos-related factors | 0.97 |
| | 1e-855 | bZIP_cEBP-like_subclass | C/EBP-related factors | 0.91 |
| | 1e-637 | TFAP2A | AP-2 | 0.93 |
| | 1e-279 | Klf1, Klf4 | Three-zinc finger Kruppel-related factors | 0.91 |
| | 1e-263 | Nr2f2_1 | Hormone-nuclear receptor | 0.86 |
| | 1e-199 | USF2, USF1 | bHLH-ZIP factors | 0.92 |
| | 1e-162 | TEAD1 | TEF-1-related factors | 0.89 |
| | 1e-161 | NFIC::TLX1 | Nuclear factor 1::NK-related factors | 0.81 |
| | 1e-135 | Six1 | Homeobox | 0.91 |
| | 1e-87 | GATA3 | GATA-type zinc fingers | 0.98 |
| | 1e-79 | Tcfap2b_1 | Helix-Loop-Helix | 0.67 |
| | 1e-79 | Foxo1 | Forkhead box factors | 0.75 |
| | 1e-77 | Chop/Ddit3 | C/EBP-related factors | 0.93 |
| | 1e-67 | Forkhead_class | Forkhead box factors | 0.94 |
| | 1e-67 | NFIC | Nuclear factor 1 | 0.86 |
| | 1e-51 | Hoxa9 | HOX-related factors | 0.96 |
| | 1e-40 | at_AC_acceptor | Unknown | 0.61 |
| | 1e-33 | GFY::Staf | Unknown::Zinc finger | 0.79 |
| | 1e-30 | Hoxc12 | HOX-related factors | 0.77 |

C

Figure S17: Functional and motif analysis of the DECRES genome-wide predictions on cell line HelaS3. A,B: Top enriched biological processes and diseases (no pathways enriched for HelaS3), respectively, in the predicted cell-specific CRRs. C: Enriched *de novo* motifs in the predicted cell specific CRRs. Column 4: families of best-matched TFs. Column 5: best match scores.

15

**Figure S18**

**GO Biological Process** — -log10(Binomial p value)

- response to decreased oxygen levels — 29.33
- response to oxygen levels — 29.07
- response to hypoxia — 27.58
- apoptotic signaling pathway — 25.10
- transforming growth factor beta receptor signaling pathway — 22.46
- cellular response to decreased oxygen levels — 19.77
- cellular response to oxygen levels — 17.88
- regulation of carbohydrate metabolic process — 16.84
- regulation of transforming growth factor beta production — 16.81
- cellular response to hypoxia — 16.73
- regulation of glucose metabolic process — 16.21
- smooth muscle cell migration — 16.20
- regulation of transforming growth factor beta2 production — 15.46
- extrinsic apoptotic signaling pathway — 15.44
- response to reactive oxygen species — 15.03
- muscle cell migration — 14.44
- regulation of transcription from RNA polymerase II promoter in response to stress — 14.18
- regulation of cellular carbohydrate metabolic process — 14.07
- regulation of DNA-dependent transcription in response to stress — 13.93
- positive regulation of cellular catabolic process — 13.11

*(A)*

**MSigDB Pathway** — -log10(Binomial p value)

- Chronic myeloid leukemia — 27.71
- PDGFR-beta signaling pathway — 20.81
- Keratinocyte Differentiation — 16.50
- Endocytosis — 15.78
- IL6-mediated signaling events — 15.55
- Direct p53 effectors — 15.26
- HIF-1-alpha transcription factor network — 15.23
- Pancreatic cancer — 15.00
- mTOR signaling pathway — 14.29
- RXR and RAR heterodimerization with other nuclear receptor — 14.15
- Role of Calcineurin-dependent NFAT signaling in lymphocytes — 13.33
- Genes involved in Signaling by TGF-Beta Receptor Complex — 13.28
- Validated targets of C-MYC transcriptional repression — 13.21
- Bladder cancer — 12.96
- Cell cycle — 12.79
- Regulation of nuclear SMAD2/3 signaling — 12.68
- Cell Cycle: G1/S Check Point — 12.56
- C-MYB transcription factor network — 12.25
- Genes involved in Transcriptional Regulation of White Adipocyte Differentiation — 12.23
- NFkB activation by Nontypeable Hemophilus influenzae

*(B)*

**Disease Ontology** — -log10(Binomial p value)

- DNA virus infectious disease — 46.73
- dsDNA virus infectious disease — 38.74
- neck neoplasm — 38.30
- neck cancer — 37.96
- neck carcinoma — 32.11
- thyroid neoplasm — 28.79
- head and neck squamous cell carcinoma — 28.27
- malignant neoplasm of thyroid — 27.97
- (+)ssRNA virus infectious disease — 25.69
- hepatitis B — 25.22
- Herpesviridae infectious disease — 24.56
- papillary epithelial neoplasm — 23.46
- upper respiratory tract disease — 22.41
- colon adenocarcinoma — 21.38
- transitional cell carcinoma — 21.04
- large intestine adenocarcinoma — 20.81
- gastric adenocarcinoma — 19.95
- stomach carcinoma — 19.61
- neoplasm of body of uterus — 19.14
- thyroid carcinoma — 18.99

*(C)*

**(D)**

| Motif | p-value | Best Match | Family | Score |
|---|---|---|---|---|
| | 1e-413 | FOS, JUN | Fos-related factors, Jun-related factors | 0.97 |
| | 1e-79 | KLF5 | Three-zinc finger Kruppel-related factors | 0.95 |
| | 1e-61 | CEBPA | C/EBP-related factors | 0.91 |
| | 1e-46 | bZIP_CREB/G-box-like_subclass | CREB-related factors | 0.84 |
| | 1e-45 | FLI1 | Ets-related factors | 0.91 |
| | 1e-32 | CTCF | More than 3 adjacent zinc finger factors | 0.87 |
| | 1e-30 | FOXD1 | Forkhead box factors | 0.82 |
| | 1e-30 | Sp4.2 | BetaBetaAlpha-zinc finger | 0.82 |
| | 1e-26 | REL | NFkB-related factors | 0.92 |
| | 1e-24 | SP1 | Three-zinc finger Kruppel-related factors | 0.73 |
| | 1e-24 | Hoxc9 | HOX-related factors | 0.73 |
| | 1e-22 | NFIC | Nuclear factor 1 | 0.78 |
| | 1e-22 | Zbtb12.2 | BetaBetaAlpha-zinc finger | 0.64 |
| | 1e-20 | TEAD1 | TEF-1-related factors | 0.79 |
| | 1e-19 | Foxo1 | Forkhead box factors | 0.72 |
| | 1e-19 | Hand1::Tcfe2a | Helix-Loop-Helix | 0.62 |
| | 1e-18 | ETS::E-box | Ets::bHLH | 0.79 |
| | 1e-16 | CEBPA | C/EBP-related factors | 0.69 |
| | 1e-16 | TBP | TBP-related factors | 0.73 |
| | 1e-16 | NFIC::TLX1 | Nuclear factor 1::NK-related factors | 0.65 |
| | 1e-16 | Tcfap2e.2 | Helix-Loop-Helix | 0.76 |
| | 1e-15 | Gata4 | GATA-type zinc fingers | 0.88 |
| | 1e-14 | NFIC | Nuclear factor 1 | 0.60 |

Figure S18: Functional and motif analysis of the DECRES NiA enhancers on cell line HelaS3. A,B,C: Top 20 enriched biological processes, pathways, and diseases, respectively, in the NiA enhancers. D: Enriched *de novo* motifs in the NiA enhancer regions. Column 4: families of best-matched TFs. Column 5: best match scores.

---

**Figure S19**

**GO Biological Process** — -log10(Binomial p value)

- steroid metabolic process — 81.07
- sterol metabolic process — 57.32
- plasma lipoprotein particle assembly — 41.62
- lipid homeostasis — 41.27
- regulation of plasma lipoprotein particle levels — 38.69
- regulation of steroid metabolic process — 38.19
- protein-lipid complex assembly — 37.78
- cholesterol homeostasis — 36.11
- sterol homeostasis — 34.71
- regulation of steroid biosynthetic process — 31.67
- acylglycerol metabolic process — 30.33
- triglyceride homeostasis — 30.29
- regulation of RNA stability — 30.16
- regulation of mRNA stability — 30.05
- neutral lipid metabolic process — 29.08
- regulation of glucose transport — 28.76
- nephron epithelium morphogenesis — 28.30
- triglyceride metabolic process — 27.44
- plasma lipoprotein particle clearance — 24.69
- body morphogenesis — 22.64

*(A)*

**MSigDB Pathway** — -log10(Binomial p value)

- FOXA2 and FOXA3 transcription factor networks — 97.68
- PPAR signaling pathway — 51.80
- Genes involved in Lipoprotein metabolism — 46.09
- Genes involved in PPARA Activates Gene Expression — 38.49
- Stabilization and expansion of the E-cadherin adherens junction — 22.34
- Thyroid cancer — 18.52
- E-cadherin signaling in the nascent adherens junction — 16.12
- LKB1 signaling events — 15.01
- Genes involved in Regulation of Hypoxia-inducible Factor (HIF) by Oxygen — 8.84

*(B)*

**Disease Ontology** — -log10(Binomial p value)

- fatty liver — 55.78
- inborn errors lipid metabolism — 47.95
- familial hyperlipidemia — 45.15
- inborn errors of amino acid metabolism — 35.64
- hypercholesterolemia — 34.25
- familial hyperlipoproteinemia — 23.42
- familial hypercholesterolemia — 21.12

*(C)*

**(D)**

| Motif | p-value | Best Match | Family | Score |
|---|---|---|---|---|
| | 1e-3527 | HNF4A, HNF4G | RXR-related receptors (NR2) | 0.97 |
| | 1e-2261 | FOXA1, Foxa2 | Forkhead box factors | 0.98 |
| | 1e-969 | REST | Factors with multiple dispersed zinc fingers | 0.82 |
| | 1e-571 | REST | Factors with multiple dispersed zinc fingers | 0.78 |
| | 1e-480 | CTCF | More than 3 adjacent zinc finger factors | 0.92 |
| | 1e-455 | Esrrb | Steroid hormone receptors (NR3) | 0.95 |
| | 1e-453 | HNF1B, HNF1A | POU domain factors | 0.89 |
| | 1e-366 | bZIP_cEBP-like_subclass | C/EBP-related factors | 0.92 |
| | 1e-337 | Arid5a.1 | ARID-related factors | 0.67 |
| | 1e-283 | Nuclear_Receptor_class | Unknown | 0.78 |
| | 1e-206 | TEAD1 | TEF-1-related factors | 0.90 |
| | 1e-161 | TCF7L2 | TCF-7-related factors | 0.91 |
| | 1e-144 | Tcf7l2.2 | TCF-7-related factors | 0.76 |
| | 1e-110 | ELF1 | Ets-related factors | 0.92 |
| | 1e-108 | JUND, FOS | Jun-related factors, Fos-related factors | 0.97 |
| | 1e-78 | Zfp691.1 | BetaBetaAlpha-zinc finger | 0.60 |
| | 1e-74 | SP1 | Three-zinc finger Kruppel-related factors | 0.92 |
| | 1e-67 | Smad3.1 | SMAD factors | 0.55 |
| | 1e-65 | Zfp161.2 | BetaBetaAlpha-zinc finger | 0.71 |
| | 1e-61 | Atoh1 | Tal-related factors | 0.75 |
| | 1e-55 | Nkx2-5(var.2) | NK-related factors | 0.59 |
| | 1e-46 | DCE_S_II | Unknown | 0.63 |

Figure S19: Functional and motif analysis of the DECRES genome-wide predictions on cell line HepG2. A,B,C: Top enriched biological processes, pathways and diseases, respectively, in the predicted cell-specific CRRs. D: Enriched *de novo* motifs in the predicted cell specific CRRs. Column 4: families of best-matched TFs. Column 5: best match scores.

Figure S20: Functional and motif analysis of the DECRES NiA enhancers on cell line HepG2. A,B,C: Top 20 enriched biological processes, pathways, and diseases, respectively, in the NiA enhancers. D: Enriched *de novo* motifs in the NiA enhancer regions. Column 4: families of best-matched TFs. Column 5: best match scores.



Figure S21: Functional and motif analysis of the DECRES genome-wide predictions on cell line HMEC. A,B,C: Top enriched biological processes,pathways and diseases, respectively, in the predicted cell-specific CRRs. D: Enriched *de novo* motifs in the predicted cell specific CRRs. Column 4: families of best-matched TFs. Column 5: best match scores.

**Figure S22**

GO Biological Process — −log10(Binomial p value)

| Term | value |
|---|---|
| apoptotic signaling pathway | 26.39 |
| regulation of apoptotic signaling pathway | 23.83 |
| regulation of extrinsic apoptotic signaling pathway | 21.36 |
| placenta development | 19.36 |
| transforming growth factor beta receptor signaling pathway | 18.35 |
| innate immune response-activating signal transduction | 18.33 |
| activation of innate immune response | 18.04 |
| positive regulation of innate immune response | 18.03 |
| negative regulation of cytokine production | 16.92 |
| pattern recognition receptor signaling pathway | 16.87 |
| response to reactive oxygen species | 15.34 |
| cell-substrate junction assembly | 14.92 |
| platelet degranulation | 14.49 |
| muscle cell migration | 14.13 |
| toll-like receptor 5 signaling pathway | 13.92 |
| toll-like receptor 10 signaling pathway | 13.92 |
| platelet-derived growth factor receptor signaling pathway | 13.68 |
| toll-like receptor TLR1:TLR2 signaling pathway | 13.63 |
| toll-like receptor TLR6:TLR2 signaling pathway | 13.63 |
| regulation of cell junction assembly | 13.59 |

MSigDB Pathway — −log10(Binomial p value)

| Term | value |
|---|---|
| Focal adhesion | 33.65 |
| ErbB1 downstream signaling | 26.64 |
| Direct p53 effectors | 23.89 |
| PDGFR-beta signaling pathway | 23.47 |
| Chronic myeloid leukemia | 22.79 |
| Pancreatic cancer | 21.98 |
| Small cell lung cancer | 21.70 |
| Keratinocyte Differentiation | 19.69 |
| a6b1 and a6b4 Integrin signaling | 19.40 |
| NFkB activation by Nontypeable Hemophilus influenzae | 18.74 |
| Adherens junction | 18.61 |
| ECM-receptor interaction | 17.92 |
| Glioma | 17.90 |
| Signaling events mediated by focal adhesion kinase | 17.10 |
| ERK1/ERK2 MAPK Pathway | 15.10 |
| Dorso-ventral axis formation | 14.73 |
| Hematopoietic cell lineage | 14.60 |
| Validated transcriptional targets of AP1 family members Fra1 and Fra2 | 14.58 |
| Renal cell carcinoma | 14.44 |
| p73 transcription factor network | 14.33 |

Disease Ontology — −log10(Binomial p value)

| Term | value |
|---|---|
| DNA virus infectious disease | 60.24 |
| dsDNA virus infectious disease | 52.00 |
| neck neoplasm | 40.60 |
| neck cancer | 40.31 |
| neck carcinoma | 34.75 |
| cancer of urinary tract | 31.43 |
| Herpesviridae infectious disease | 30.72 |
| head and neck squamous cell carcinoma | 30.51 |
| lichen planus | 28.92 |
| neoplasm of body of uterus | 28.74 |
| upper respiratory tract disease | 28.24 |
| lung adenocarcinoma | 27.94 |
| hepatitis B | 27.93 |
| (+)ssRNA virus infectious disease | 26.90 |
| smooth muscle tumor | 26.47 |
| nephrosis | 26.06 |
| transitional cell carcinoma | 25.78 |
| Adenoviridae infectious disease | 25.24 |
| malignant neoplasm of prostate | 24.35 |
| malignant neoplasm of male genital organ or tract | 24.31 |

D:

| Motif | p-value | Best Match | Family | Score |
|---|---|---|---|---|
|  | 1e-288 | FOS, JUN | Fos-related factors, Jun-related factors | 0.95 |
|  | 1e-70 | Erg | Ets-related factors | 0.92 |
|  | 1e-33 | KLF5 | Three-zinc finger Kruppel-related factors | 0.92 |
|  | 1e-30 | TEAD1 | TEF-1-related factors | 0.87 |
|  | 1e-30 | Atf1_1 | Leucine zipper | 0.92 |
|  | 1e-24 | FOXO3 | Forkhead box factors | 0.68 |
|  | 1e-23 | RUNX1 | Runt-related factors | 0.70 |
|  | 1e-21 | Zbtb3_2 | BetaBetaAlpha-zinc finger | 0.69 |
|  | 1e-21 | Klf1 | Three-zinc finger Kruppel-related factors | 0.75 |
|  | 1e-20 | RFX5 | RFX-related factors | 0.73 |
|  | 1e-19 | CEBPA | C/EBP-related factors | 0.86 |
|  | 1e-19 | Rhox11_2 | Paired-related HD factors | 0.69 |
|  | 1e-19 | Gm397_2 | BetaBetaAlpha-zinc finger | 0.77 |
|  | 1e-19 | Osr1_2 | BetaBetaAlpha-zinc finger | 0.61 |
|  | 1e-18 | Gcm1_1 | Glial Cells Missing (GCM) | 0.64 |
|  | 1e-17 | Gmeb1_1 | Sand | 0.66 |
|  | 1e-15 | INR | Unknown | 0.61 |
|  | 1e-14 | FOXI1 | Forkhead box factors | 0.65 |
|  | 1e-14 | Srf_2 | MADS | 0.87 |
|  | 1e-14 | DCE_S_II | Unknown | 0.70 |
|  | 1e-14 | Nr2e3 | RXR-related receptors (NR2) | 0.60 |
|  | 1e-13 | TBP | TBP-related factors | 0.72 |
|  | 1e-13 | Hoxc9 | HOX-related factors | 0.72 |

Figure S22: Functional and motif analysis of the DECRES NiA enhancers on cell line HMEC. A,B,C: Top 20 enriched biological processes, pathways, and diseases, respectively, in the NiA enhancers. D: Enriched *de novo* motifs in the NiA enhancer regions. Column 4: families of best-matched TFs. Column 5: best match scores.
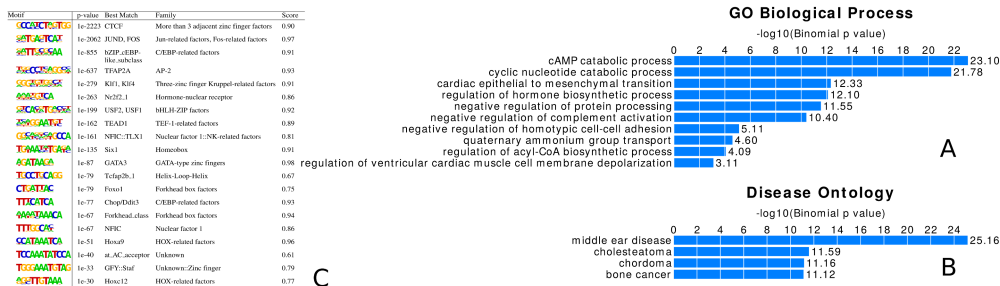
---

**Figure S23**

GO Biological Process — −log10(Binomial p value)

| Term | value |
|---|---|
| blood vessel development | 114.09 |
| blood vessel morphogenesis | 108.14 |
| angiogenesis | 108.02 |
| endocytosis | 55.65 |
| regulation of Rho protein signal transduction | 51.42 |
| platelet activation | 50.72 |
| response to transforming growth factor beta stimulus | 49.58 |
| cellular response to transforming growth factor beta stimulus | 48.63 |
| transforming growth factor beta receptor signaling pathway | 47.83 |
| phagocytosis | 41.56 |
| platelet degranulation | 35.88 |
| regulation of epithelial cell migration | 35.67 |
| regulation of Rho GTPase activity | 33.78 |
| regulation of cell shape | 33.68 |
| vascular endothelial growth factor receptor signaling pathway | 32.73 |
| regulation of endothelial cell migration | 31.74 |
| cell-substrate adhesion | 31.72 |
| cytoskeleton-dependent intracellular transport | 31.37 |
| cell-matrix adhesion | 30.58 |
| Fc-gamma receptor signaling pathway involved in phagocytosis | 28.83 |

MSigDB Pathway — −log10(Binomial p value)

| Term | value |
|---|---|
| Genes involved in Hemostasis | 80.29 |
| Genes involved in Platelet activation, signaling and aggregation | 59.39 |
| Focal adhesion | 52.40 |
| Signaling events mediated by VEGFR1 and VEGFR2 | 48.81 |
| Genes involved in Signaling by TGF-beta Receptor Complex | 46.72 |
| Signaling events mediated by focal adhesion kinase | 46.48 |
| Genes involved in Signaling by Rho GTPases | 45.97 |
| TGF-beta receptor signaling | 40.77 |
| Genes involved in Transcriptional activity of SMAD2/SMAD3:SMAD4 heterotrimer | 37.39 |
| Integrin-linked kinase signaling | 34.76 |
| Netrin-mediated signaling events | 34.57 |
| Genes involved in Response to elevated platelet cytosolic Ca2+ | 34.51 |
| Integrin Signaling Pathway | 34.04 |
| Leukocyte transendothelial migration | 33.94 |
| PDGFR-beta signaling pathway | 32.03 |
| VEGF, Hypoxia, and Angiogenesis | 31.12 |
| VEGF and VEGFR signaling network | 30.97 |
| Genes involved in Integrin cell surface interactions | 30.55 |
| Regulation of RhoA activity | 30.51 |
| Genes involved in Sema4D in semaphorin signaling | 30.40 |

Disease Ontology — −log10(Binomial p value)

| Term | value |
|---|---|
| colon adenocarcinoma | 29.52 |
| neoplasm in vascular tissue | 25.42 |
| Adenoviridae infectious disease | 22.07 |
| papillary adenocarcinoma | 19.70 |
| hemangioma | 17.72 |
| Wiskott-Aldrich syndrome | 15.73 |
| osteonecrosis | 9.62 |
| angiosarcoma | 8.96 |
| proliferative diabetic retinopathy | 8.57 |
| splenic disease | 8.39 |
| chronic myeloproliferative disease | 8.25 |
| lymphoid leukemia | 5.21 |
| hepatitis E | 2.63 |

D:

| Motif | p-value | Best Match | Family | Score |
|---|---|---|---|---|
|  | 1e-999 | JUNB, JUND, FOS, FOSL2 | Jun-related factors, Fos-related factors | 0.97 |
|  | 1e-689 | FEV, FLI1 | Ets-related factors | 0.93 |
|  | 1e-142 | ELF5 | Ets-related factors | 0.69 |
|  | 1e-64 | SOX9 | SOX-related factors | 0.93 |
|  | 1e-48 | GATA3, Gata4, Gata1 | GATA-type zinc fingers | 0.96 |
|  | 1e-44 | Foxa2, FOXA1 | Forkhead box factors | 0.74 |
|  | 1e-43 | SPIB | Ets-related factors | 0.65 |
|  | 1e-35 | Mycn | bHLH-ZIP factors | 0.96 |
|  | 1e-31 | Foxa2 | Forkhead box factors | 0.75 |
|  | 1e-30 | EHF | Ets-related factors | 0.55 |
|  | 1e-29 | Mafb | Maf-related factors | 0.77 |
|  | 1e-24 | Hic1_1 | BetaBetaAlpha-zinc finger | 0.67 |
|  | 1e-23 | Gata5_1 | GATA-type zinc fingers | 0.70 |
|  | 1e-23 | EHF | Ets-related factors | 0.70 |
|  | 1e-21 | NFIC | Nuclear factor 1 | 0.86 |
|  | 1e-19 | Nkx3-1_1 | NK-related factors | 0.71 |
|  | 1e-18 | E2F6 | E2F-related factors | 0.60 |
|  | 1e-15 | Tcfap2b_2 | Helix-Loop-Helix | 0.72 |
|  | 1e-13 | MZF1 | More than 3 adjacent zinc finger factors | 0.63 |
|  | 1e-12 | SMAD2::SMAD3::SMAD4 | SMAD factors | 0.62 |

Figure S23: Functional and motif analysis of the DECRES genome-wide predictions on cell line HUVEC. A,B,C: Top enriched biological processes, pathways and diseases, respectively, in the predicted cell-specific CRRs. D: Enriched *de novo* motifs in the predicted cell specific CRRs. Column 4: families of best-matched TFs. Column 5: best match scores.
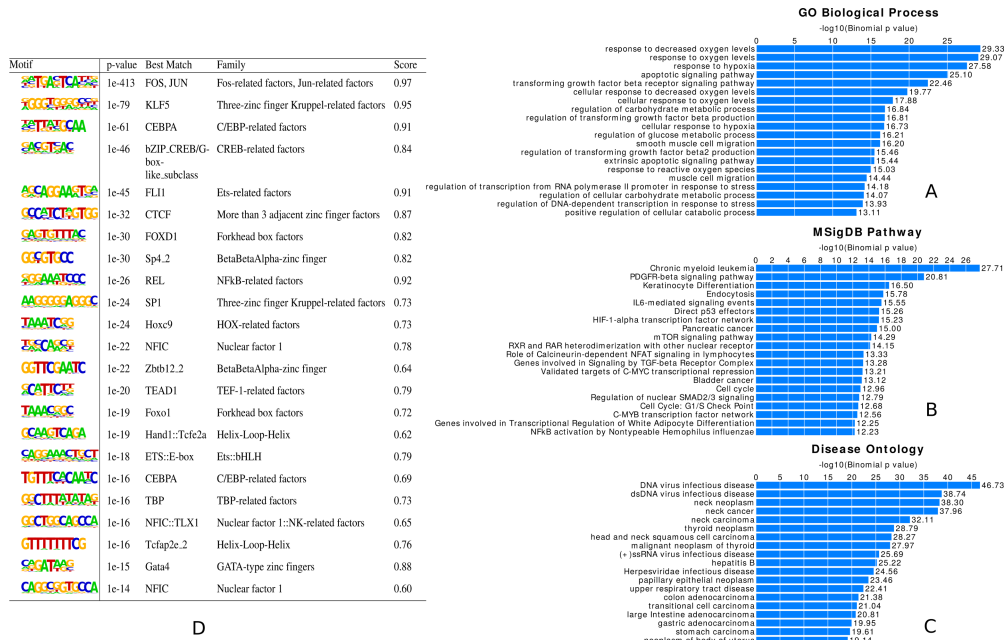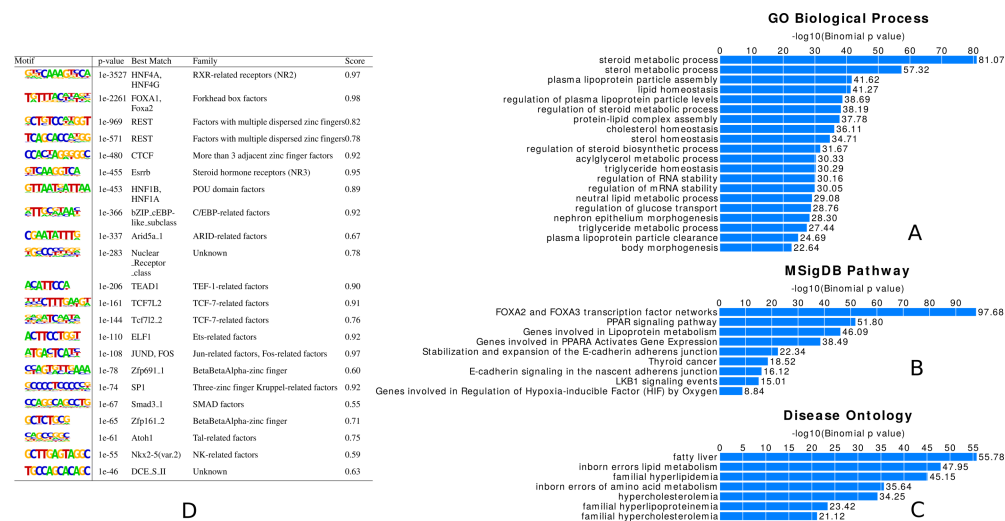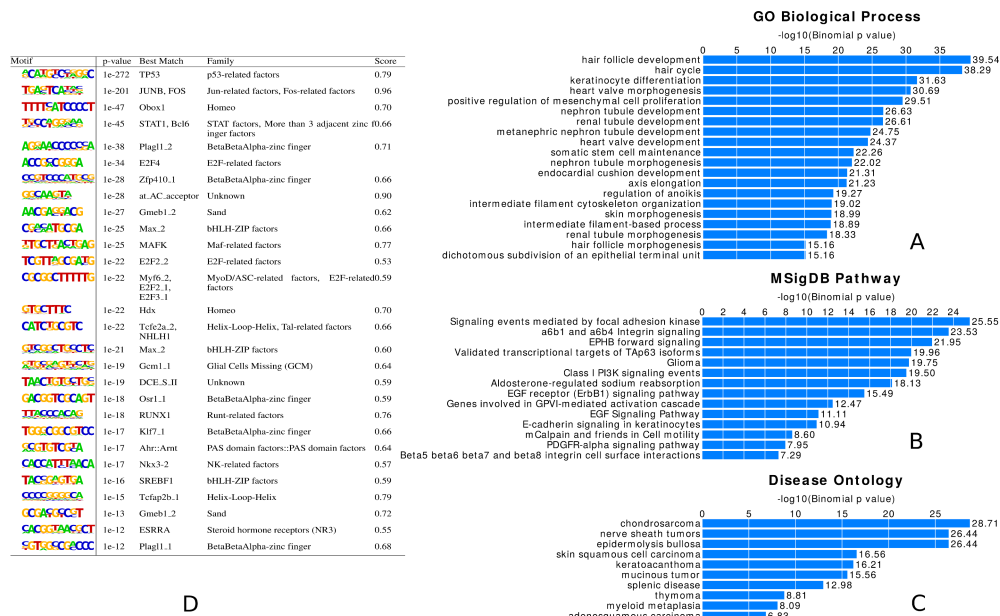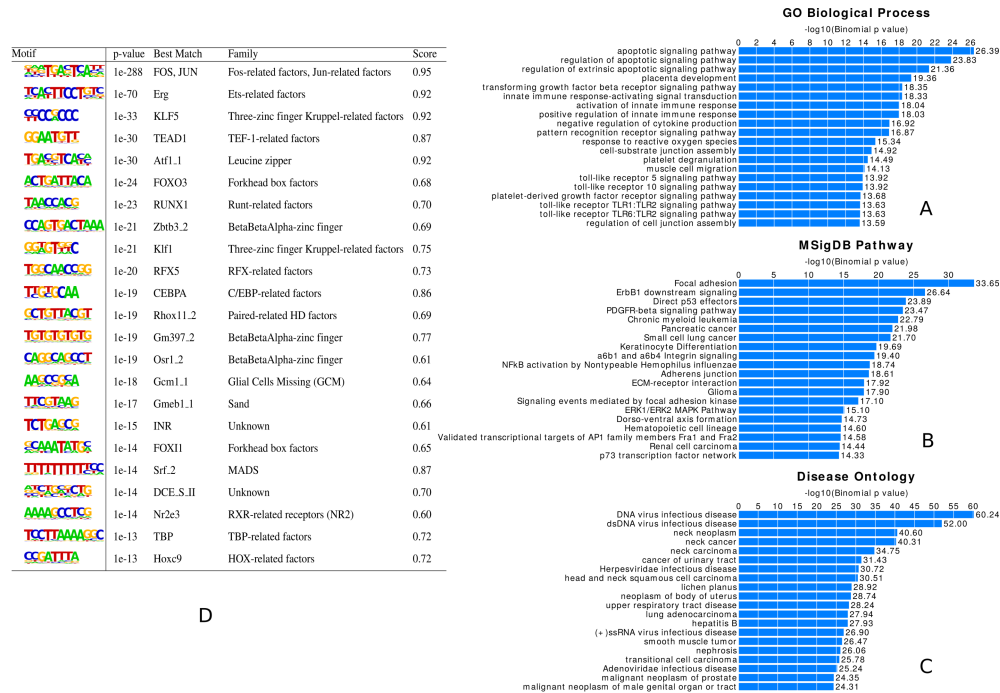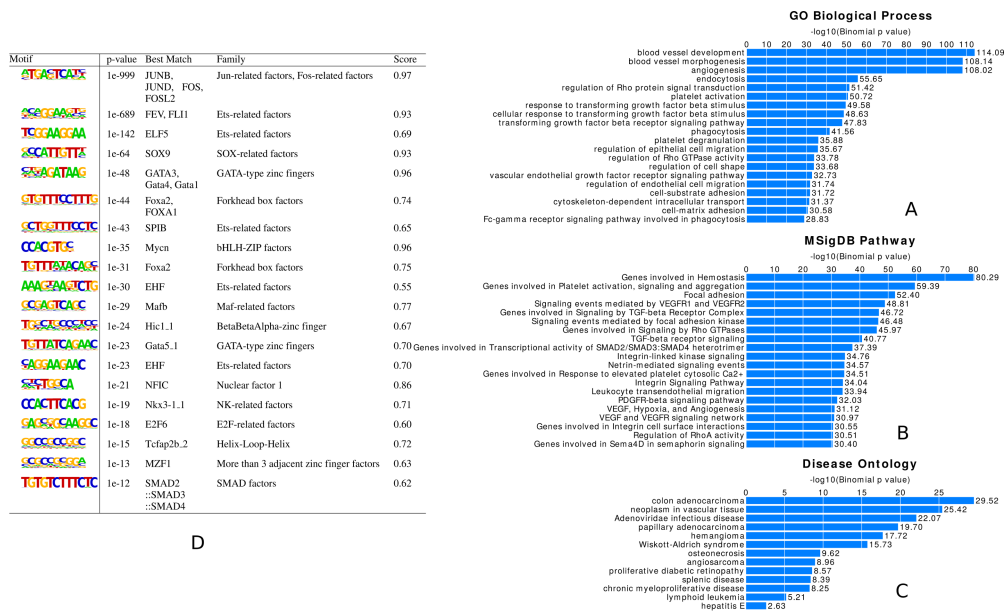
18

## GO Biological Process
*-log10(Binomial p value)*

- apoptotic signaling pathway — 44.81
- nucleotide-binding domain, leucine rich repeat containing receptor signaling pathway — 35.92
- negative regulation of inflammatory response — 31.86
- regulation of cell junction assembly — 31.26
- regulation of translation — 30.82
- transforming growth factor beta receptor signaling pathway — 30.44
- leukocyte migration — 29.51
- regulation of apoptotic signaling pathway — 28.05
- regulation of inflammatory response — 27.90
- cellular response to oxygen levels — 27.64
- regulation of angiogenesis — 26.94
- negative regulation of protein kinase activity — 25.44
- cell-substrate adhesion — 25.15
- regulation of focal adhesion assembly — 23.41
- positive regulation of angiogenesis — 23.24
- phagocytosis — 22.14
- activation of innate immune response — 21.84
- innate immune response-activating signal transduction — 21.69
- pattern recognition receptor signaling pathway — 21.69
- platelet degranulation — 21.65

**A**

## MSigDB Pathway
*-log10(Binomial p value)*

- PDGFR-beta signaling pathway — 41.04
- Chronic myeloid leukemia — 39.02
- Focal adhesion — 32.37
- Genes involved in Platelet activation, signaling and aggregation — 29.53
- Genes involved in Sema4D in semaphorin signaling — 28.53
- Signaling events mediated by focal adhesion kinase — 28.32
- TGF-beta receptor signaling — 28.10
- Genes involved in Signaling by TGF-beta Receptor Complex — 27.87
- IL6-mediated signaling events — 26.60
- Pancreatic cancer — 26.59
- TNF receptor signaling pathway — 25.71
- NFkB activation by Nontypeable Hemophilus influenzae — 25.06
- ErbB1 downstream signaling — 24.20
- Genes involved in Signaling by Interleukins — 23.19
- Integrin Signaling Pathway — 23.02
- Genes involved in Transcriptional activity of SMAD2/SMAD3:SMAD4 heterotrimer — 22.84
- Genes involved in Integrin cell surface interactions — 22.83
- Fc gamma R-mediated phagocytosis — 22.67
- AP-1 transcription factor network — 22.61
- HIF-1-alpha transcription factor network — 22.48

**B**

## Disease Ontology
*-log10(Binomial p value)*

- DNA virus infectious disease — 73.21
- dsDNA virus infectious disease — 63.86
- epithelial ovarian cancer — 47.24
- (+)ssRNA virus infectious disease — 37.85
- Adenoviridae infectious disease — 37.27
- Herpesviridae infectious disease — 32.05
- transitional cell carcinoma — 32.02
- hepatitis B — 31.37
- head and neck squamous cell carcinoma — 29.46
- laryngeal squamous cell carcinoma — 29.39
- large intestine adenocarcinoma — 28.26
- colon adenocarcinoma — 27.94
- neoplasm of body of uterus — 27.85
- laryngeal neoplasm — 27.03
- lichen planus — 26.89
- carcinoma of larynx — 26.89
- laryngeal disease — 26.03
- colon carcinoma — 25.71
- colonic neoplasm — 24.57
- myelofibrosis — 23.84

**C**

**D**

| Motif | p-value | Best Match | Family | Score |
|---|---|---|---|---|
| | 1e-382 | JUNB, FOS | Jun-related factors, Fos-related factors | 0.96 |
| | 1e-230 | Erg | Ets-related factors | 0.96 |
| | 1e-42 | bZIP CREB/G-box-like subclass | CREB-related factors | 0.91 |
| | 1e-39 | GC-box, KLF5 | Three-zinc finger Kruppel-related factors | 0.89 |
| | 1e-32 | FOXD1 | Forkhead box factors | 0.68 |
| | 1e-32 | NFKB1 | NFkB-related factors | 0.79 |
| | 1e-31 | RUNX1 | Runt-related factors | 0.73 |
| | 1e-31 | Gata1 | GATA-type zinc fingers | 0.94 |
| | 1e-27 | Nkx2-4 | NK-related factors | 0.57 |
| | 1e-27 | Mafb | Maf-related factors | 0.67 |
| | 1e-24 | Crx | Paired-related HD factors | 0.74 |
| | 1e-24 | Myog, Tcf12 | MyoD/ASC-related factors, E2A-related factors | 0.93 |
| | 1e-23 | NFATC2 | NFAT-related factors | 0.63 |
| | 1e-23 | Hand1::Tcfe2a, NFIC | Tal-related factors::E2A-related factors, Nuclear factor 1 | 0.69 |
| | 1e-21 | ARNT::HIF1A | PAS domain factors::PAS domain factors | 0.75 |
| | 1e-21 | Mybl1_1 | Myb/SANT domain factors | 0.71 |
| | 1e-21 | Rfxdc2_1 | RFX-related factors | 0.78 |
| | 1e-19 | Srf_2 | MADS | 0.87 |
| | 1e-18 | Gm397_2 | BetaBetaAlpha-zinc finger | 0.72 |
| | 1e-17 | Isgf3g_2 | Interferon-regulatory factors | 0.67 |
| | 1e-16 | Zfp161_2 | BetaBetaAlpha-zinc finger | 0.69 |
| | 1e-16 | at_AC_acceptor | Unknown | 0.61 |
| | 1e-13 | Rxra | RXR-related receptors (NR2) | 0.55 |
| | 1e-12 | Esrra_2 | Hormone-nuclear receptor | 0.68 |
| | 1e-12 | Irf3_2 | Interferon-regulatory factors | 0.61 |

Figure S24: Functional and motif analysis of the DECRES NiA enhancers on cell line HUVEC. A,B,C: Top 20 enriched biological processes, pathways, and diseases, respectively, in the NiA enhancers. D: Enriched *de novo* motifs in the NiA enhancer regions. Column 4: families of best-matched TFs. Column 5: best match scores.

## GO Biological Process
*-log10(Binomial p value)*

- erythrocyte homeostasis — 44.21
- myeloid cell homeostasis — 40.56
- regulation of gene expression, epigenetic — 36.26
- histone lysine methylation — 33.23
- myeloid leukocyte activation — 32.74
- protein methylation — 29.69
- histone methylation — 29.51
- neutrophil activation — 26.13
- leukocyte degranulation — 25.41
- granulocyte differentiation — 23.28
- negative regulation of transforming growth factor beta receptor signaling pathway — 22.74
- regulated secretory pathway — 18.90
- natural killer cell activation — 18.08
- tooth mineralization — 15.05
- negative regulation of Notch signaling pathway — 13.01
- nitric oxide metabolic process — 10.38
- natural killer cell differentiation — 9.88
- peptidyl-lysine methylation — 6.17
- positive regulation of interleukin-2 biosynthetic process — 4.63

**A**

## MSigDB Pathway
*-log10(Binomial p value)*

- IL3-mediated signaling events — 31.99
- GMCSF-mediated signaling events — 25.33
- TPO Signaling Pathway — 24.43
- Fc epsilon RI signaling pathway — 23.78
- T Cell Receptor Signaling Pathway — 22.76
- IL2 signaling events mediated by STAT5 — 19.52
- IL2-mediated signaling events — 18.82
- PDGF Signaling Pathway — 18.69
- Role of Calcineurin-dependent NFAT signaling in lymphocytes — 16.71
- Fc Epsilon Receptor I Signaling in Mast Cells — 15.84
- Effects of calcineurin in Keratinocyte Differentiation — 10.13
- Genes involved in Iron uptake and transport — 7.54

**B**

## Disease Ontology
*-log10(Binomial p value)*

- juvenile myelomonocytic leukemia — 10.54
- vascular dementia — 9.09
- toxoplasmosis — 5.36
- coccidiosis — 5.27

**C**

**D**

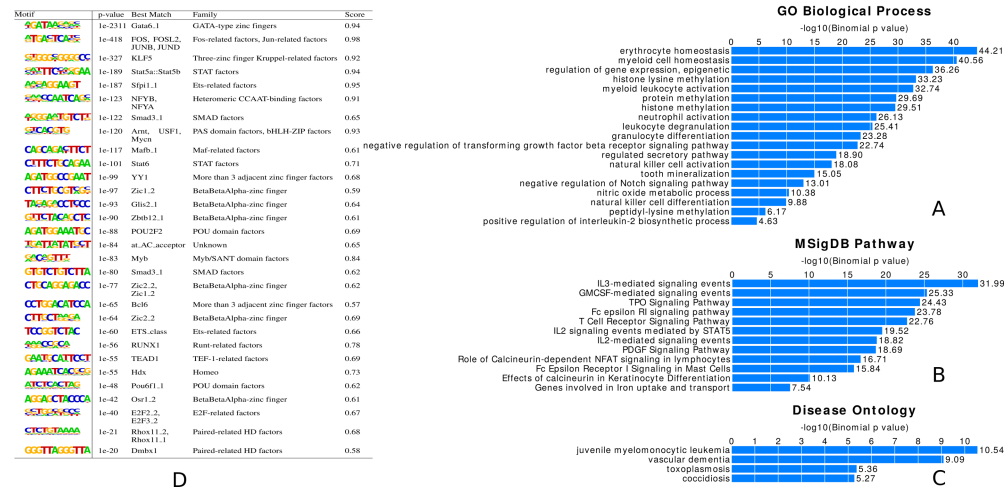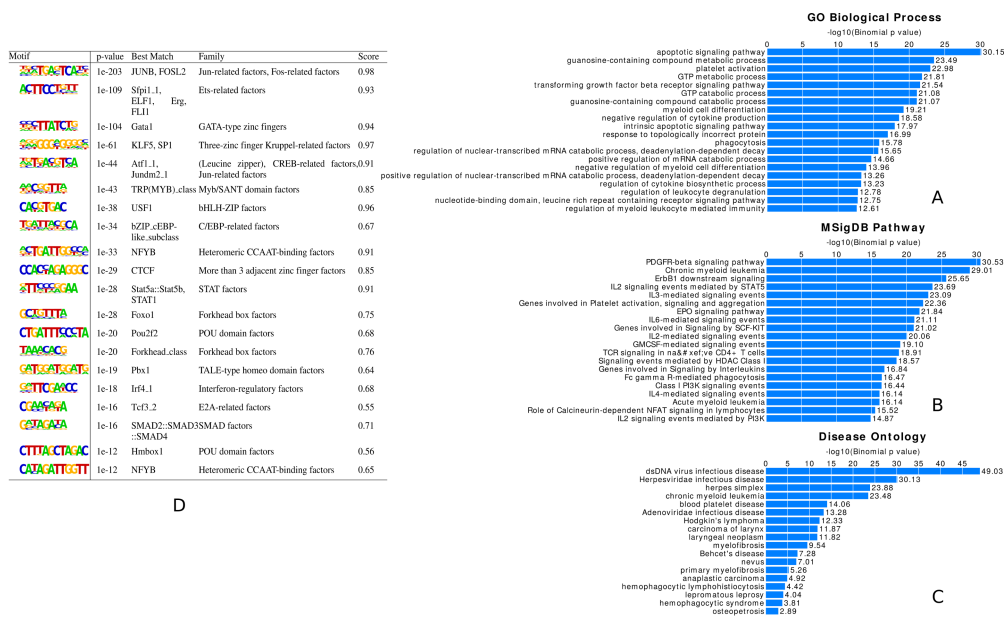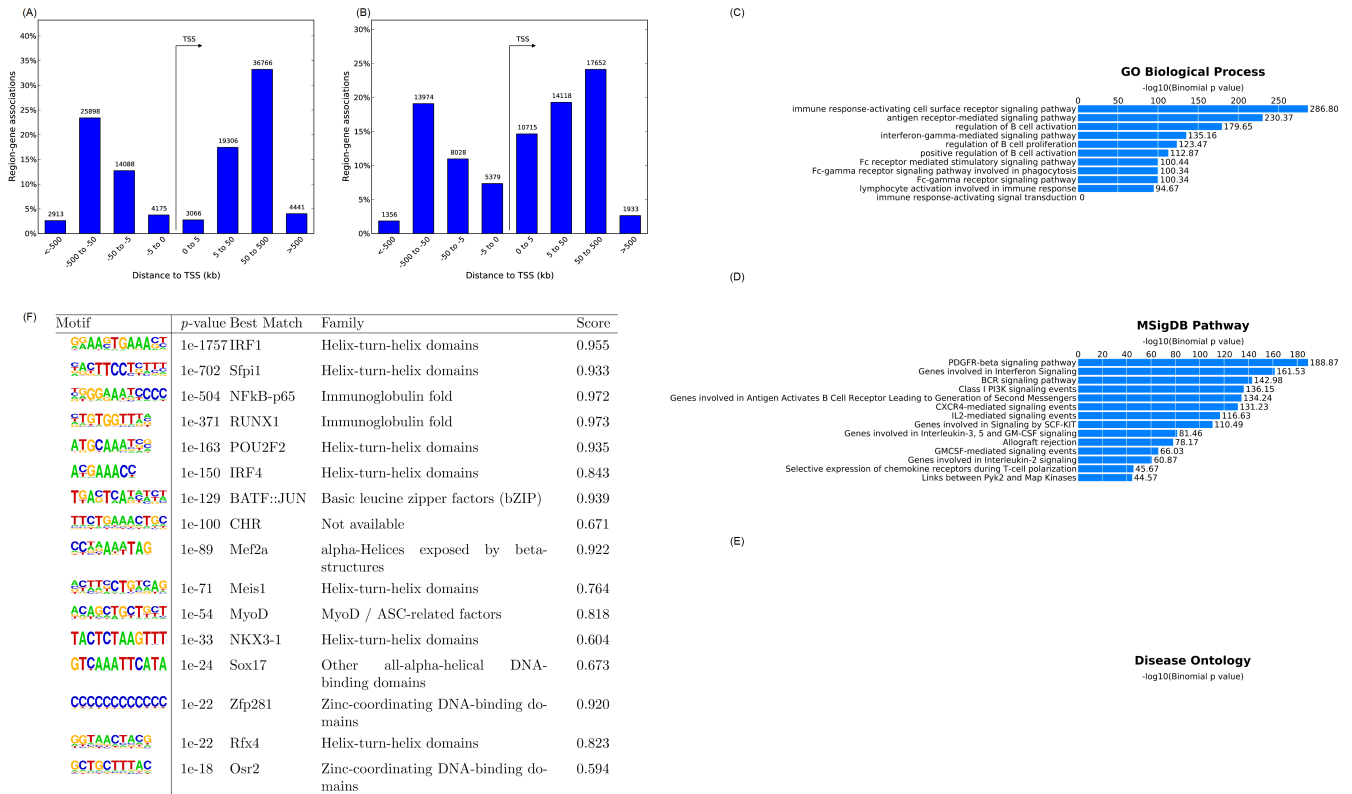| Motif | p-value | Best Match | Family | Score |
|---|---|---|---|---|
| | 1e-2311 | Gata6_1 | GATA-type zinc fingers | 0.94 |
| | 1e-418 | FOS, FOSL2, JUNB, JUND | Fos-related factors, Jun-related factors | 0.98 |
| | 1e-327 | KLF5 | Three-zinc finger Kruppel-related factors | 0.92 |
| | 1e-189 | Stat5a::Stat5b | STAT factors | 0.94 |
| | 1e-187 | Sfpi1_1 | Ets-related factors | 0.95 |
| | 1e-123 | NFYB, NFYA | Heteromeric CCAAT-binding factors | 0.91 |
| | 1e-122 | Smad3_1 | SMAD factors | 0.65 |
| | 1e-120 | Arnt, USF1, Mycn | PAS domain factors, bHLH-ZIP factors | 0.93 |
| | 1e-117 | Mafb_1 | Maf-related factors | 0.61 |
| | 1e-101 | Stat6 | STAT factors | 0.71 |
| | 1e-99 | YY1 | More than 3 adjacent zinc finger factors | 0.68 |
| | 1e-97 | Zic1_2 | BetaBetaAlpha-zinc finger | 0.59 |
| | 1e-93 | Glis2_1 | BetaBetaAlpha-zinc finger | 0.64 |
| | 1e-90 | Zbtb12_1 | BetaBetaAlpha-zinc finger | 0.61 |
| | 1e-88 | POU2F2 | POU domain factors | 0.69 |
| | 1e-84 | at_AC_acceptor | Unknown | 0.65 |
| | 1e-83 | Myb | Myb/SANT domain factors | 0.84 |
| | 1e-80 | Smad3_1 | SMAD factors | 0.62 |
| | 1e-77 | Zic2_2, Zic1_2 | BetaBetaAlpha-zinc finger | 0.62 |
| | 1e-65 | Bcl6 | More than 3 adjacent zinc finger factors | 0.57 |
| | 1e-64 | Zic2_2 | BetaBetaAlpha-zinc finger | 0.69 |
| | 1e-60 | ETS_class | Ets-related factors | 0.66 |
| | 1e-56 | RUNX1 | Runt-related factors | 0.78 |
| | 1e-55 | TEAD1 | TEF-1-related factors | 0.69 |
| | 1e-55 | Hdx | Homeo | 0.73 |
| | 1e-48 | Pou6f1_1 | POU domain factors | 0.62 |
| | 1e-42 | Osr1_2 | BetaBetaAlpha-zinc finger | 0.61 |
| | 1e-40 | E2F2_2, E2F3_2 | E2F-related factors | 0.67 |
| | 1e-21 | Rhox11_2, Rhox11_1 | Paired-related HD factors | 0.68 |
| | 1e-20 | Dmbx1 | Paired-related HD factors | 0.58 |

Figure S25: Functional and motif analysis of the DECRES genome-wide predictions on cell line K562. A,B,C: Top enriched biological processes, pathways and diseases, respectively, in the predicted cell-specific CRRs. D: Enriched *de novo* motifs in the predicted cell specific CRRs. Column 4: families of best-matched TFs. Column 5: best match scores.
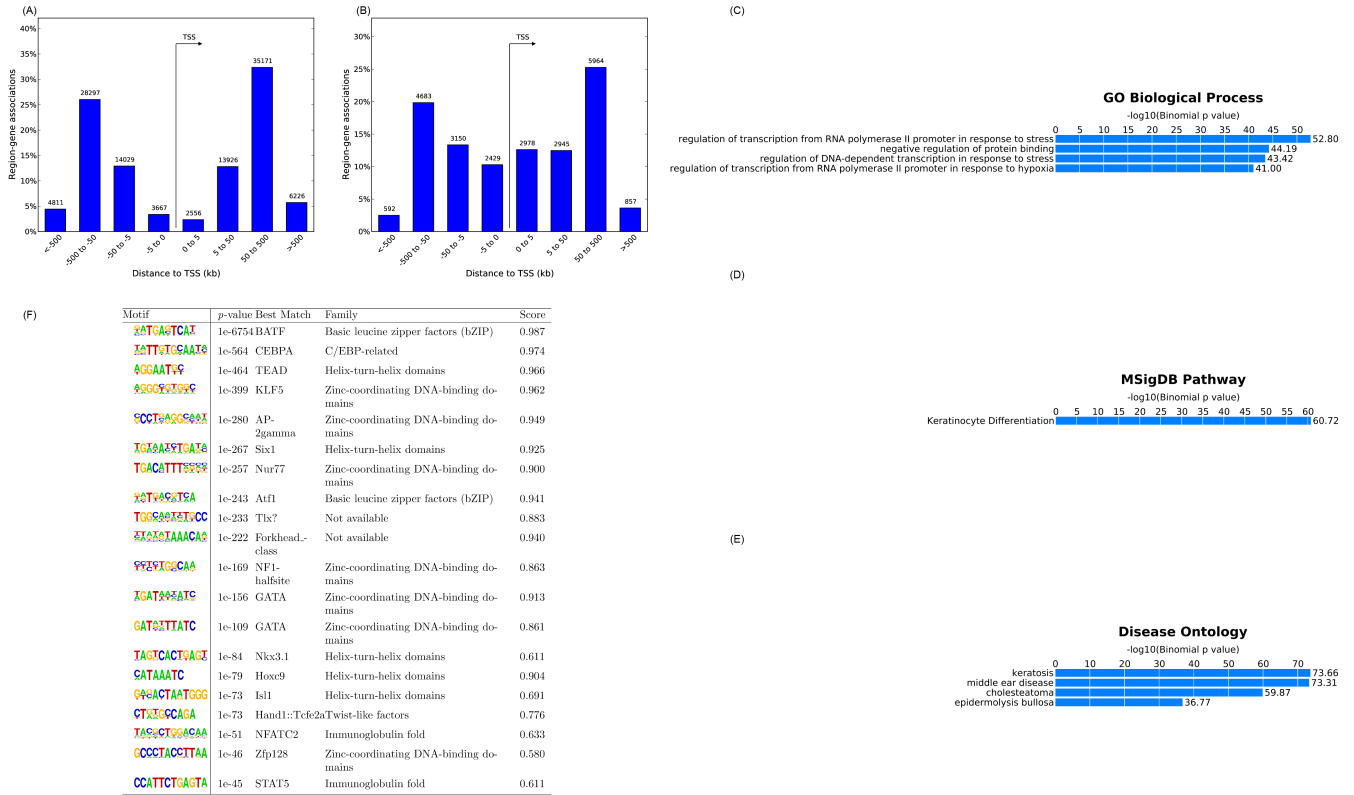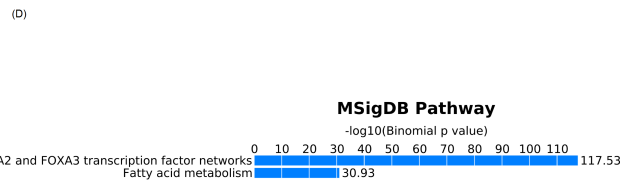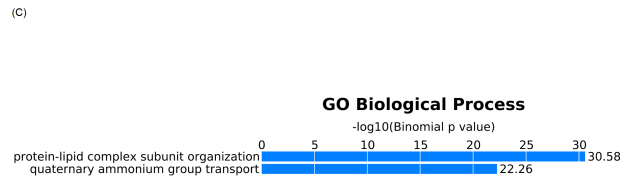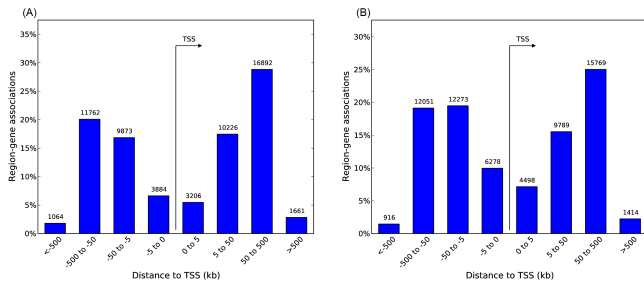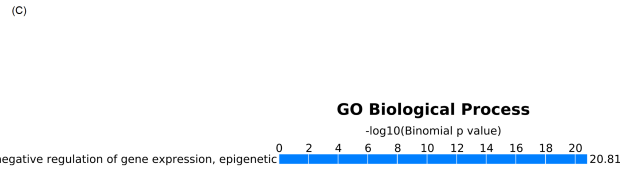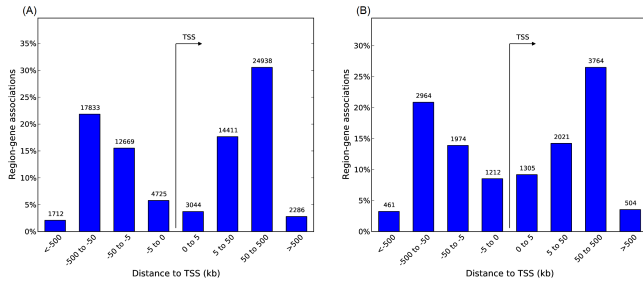
**GO Biological Process**

-log10(Binomial p value)

| Term | Value |
|---|---|
| apoptotic signaling pathway | 30.15 |
| guanosine-containing compound metabolic process | 23.49 |
| platelet activation | 22.98 |
| GTP metabolic process | 21.81 |
| transforming growth factor beta receptor signaling pathway | 21.54 |
| GTP catabolic process | 21.08 |
| guanosine-containing compound catabolic process | 21.07 |
| myeloid cell differentiation | 19.21 |
| negative regulation of cytokine production | 18.58 |
| intrinsic apoptotic signaling pathway | 17.97 |
| response to topologically incorrect protein | 16.99 |
| regulation of nuclear-transcribed mRNA catabolic process, deadenylation-dependent decay | 15.78 |
| positive regulation of mRNA catabolic process | 15.65 |
| negative regulation of myeloid cell differentiation | 14.66 |
| positive regulation of nuclear-transcribed mRNA catabolic process, deadenylation-dependent decay | 13.96 |
| regulation of cytokine biosynthetic process | 13.26 |
| regulation of leukocyte degranulation | 13.23 |
| nucleotide-binding domain, leucine rich repeat containing receptor signaling pathway | 12.78 |
| regulation of myeloid leukocyte mediated immunity | 12.75 |
| | 12.61 |

A

**MSigDB Pathway**

-log10(Binomial p value)

| Term | Value |
|---|---|
| PDGFR-beta signaling pathway | 30.53 |
| Chronic myeloid leukemia | 29.01 |
| ErbB1 downstream signaling | 25.65 |
| IL2 signaling events mediated by STAT5 | 23.69 |
| IL3-mediated signaling events | 23.09 |
| Genes involved in Platelet activation, signaling and aggregation | 22.36 |
| EPO signaling pathway | 21.84 |
| IL6-mediated signaling events | 21.11 |
| Genes involved in Signaling by SCF-KIT | 21.02 |
| IL2-mediated signaling events | 20.06 |
| GMCSF-mediated signaling events | 19.10 |
| TCR signaling in naïve CD4+ T cells | 18.91 |
| Signaling events mediated by HDAC Class I | 18.57 |
| Genes involved in Signaling by Interleukins | 16.84 |
| Fc gamma R-mediated phagocytosis | 16.47 |
| Class I PI3K signaling events | 16.44 |
| IL4-mediated signaling events | 16.14 |
| Acute myeloid leukemia | 16.14 |
| Role of Calcineurin-dependent NFAT signaling in lymphocytes | 15.52 |
| IL2 signaling events mediated by PI3K | 14.87 |

B

**Disease Ontology**

-log10(Binomial p value)

| Term | Value |
|---|---|
| dsDNA virus infectious disease | 49.03 |
| Herpesviridae infectious disease | 30.13 |
| herpes simplex | 23.88 |
| chronic myeloid leukemia | 23.48 |
| blood platelet disease | 14.06 |
| Adenoviridae infectious disease | 13.28 |
| Hodgkin's lymphoma | 12.33 |
| carcinoma of larynx | 11.87 |
| laryngeal neoplasm | 11.82 |
| myelofibrosis | 9.54 |
| Behcet's disease | 7.28 |
| nevus | 7.01 |
| primary myelofibrosis | 5.26 |
| anaplastic carcinoma | 4.92 |
| hemophagocytic lymphohistiocytosis | 4.42 |
| lepromatous leprosy | 4.04 |
| hemophagocytic syndrome | 3.81 |
| osteopetrosis | 2.89 |

C

| Motif | p-value | Best Match | Family | Score |
|---|---|---|---|---|
| | 1e-203 | JUNB, FOSL2 | Jun-related factors, Fos-related factors | 0.98 |
| | 1e-109 | Sfpi1.1, ELF1, Erg, FLI1 | Ets-related factors | 0.93 |
| | 1e-104 | Gata1 | GATA-type zinc fingers | 0.94 |
| | 1e-61 | KLF5, SP1 | Three-zinc finger Kruppel-related factors | 0.97 |
| | 1e-44 | Atf1.1, Jundm2.1 | (Leucine zipper), CREB-related factors, Jun-related factors | 0.91 |
| | 1e-43 | TRP(MYB).class | Myb/SANT domain factors | 0.85 |
| | 1e-38 | USF1 | bHLH-ZIP factors | 0.96 |
| | 1e-34 | bZIP.cEBP-like.subclass | C/EBP-related factors | 0.67 |
| | 1e-33 | NFYB | Heteromeric CCAAT-binding factors | 0.91 |
| | 1e-29 | CTCF | More than 3 adjacent zinc finger factors | 0.85 |
| | 1e-28 | Stat5a::Stat5b, STAT1 | STAT factors | 0.91 |
| | 1e-28 | Foxo1 | Forkhead box factors | 0.75 |
| | 1e-20 | Pou2f2 | POU domain factors | 0.68 |
| | 1e-20 | Forkhead.class | Forkhead box factors | 0.76 |
| | 1e-19 | Pbx1 | TALE-type homeo domain factors | 0.64 |
| | 1e-18 | Irf4.1 | Interferon-regulatory factors | 0.68 |
| | 1e-16 | Tcf3.2 | E2A-related factors | 0.55 |
| | 1e-16 | SMAD2::SMAD3::SMAD4 | SMAD factors | 0.71 |
| | 1e-12 | Hmbox1 | POU domain factors | 0.56 |
| | 1e-12 | NFYB | Heteromeric CCAAT-binding factors | 0.65 |

D

Figure S26: Functional and motif analysis of the DECRES NiA enhancers on cell line K562. A,B,C: Top 20 enriched biological processes, pathways, and diseases, respectively, in the NiA enhancers. D: Enriched *de novo* motifs in the NiA enhancer regions. Column 4: families of best-matched TFs. Column 5: best match scores.

(A)

(B)

(C)

**GO Biological Process**

(D)

**MSigDB Pathway**

(E)

**Disease Ontology**

(F)

| Motif | $p$-value | Best Match | Family | Score |
|---|---|---|---|---|
| | 1e-1757 | IRF1 | Helix-turn-helix domains | 0.955 |
| | 1e-702 | Sfpi1 | Helix-turn-helix domains | 0.933 |
| | 1e-504 | NFkB-p65 | Immunoglobulin fold | 0.972 |
| | 1e-371 | RUNX1 | Immunoglobulin fold | 0.973 |
| | 1e-163 | POU2F2 | Helix-turn-helix domains | 0.935 |
| | 1e-150 | IRF4 | Helix-turn-helix domains | 0.843 |
| | 1e-129 | BATF::JUN | Basic leucine zipper factors (bZIP) | 0.939 |
| | 1e-100 | CHR | Not available | 0.671 |
| | 1e-89 | Mef2a | alpha-Helices exposed by beta-structures | 0.922 |
| | 1e-71 | Meis1 | Helix-turn-helix domains | 0.764 |
| | 1e-54 | MyoD | MyoD / ASC-related factors | 0.818 |
| | 1e-33 | NKX3-1 | Helix-turn-helix domains | 0.604 |
| | 1e-24 | Sox17 | Other all-alpha-helical DNA-binding domains | 0.673 |
| | 1e-22 | Zfp281 | Zinc-coordinating DNA-binding domains | 0.920 |
| | 1e-22 | Rfx4 | Helix-turn-helix domains | 0.823 |
| | 1e-18 | Osr2 | Zinc-coordinating DNA-binding domains | 0.594 |

Figure S27: Functional and motif analysis of the Combined genome-wide predictions on cell line GM12878. A: Distance from the predicted A-Es to gene TSSs. B: Distance from the predicted A-Ps to gene TSSs. C,D,E: Top 20 enriched biological processes, pathways, and diseases, respectively, in the predicted cell-specific CRRs. F: Enriched *de novo* motifs in the predicted cell specific CRRs. Column 4: families of best-matched TFs. Column 5: best match scores.

**(A)** — Distance to TSS (kb) histogram with values: 4811, 28297, 14029, 3667, 2556, 13926, 35171, 6226

**(B)** — Distance to TSS (kb) histogram with values: 592, 4683, 3150, 2429, 2978, 2945, 5964, 857

**(C) GO Biological Process**
-log10(Binomial p value)

| | |
|---|---|
| regulation of transcription from RNA polymerase II promoter in response to stress | 52.80 |
| negative regulation of protein binding | 44.19 |
| regulation of DNA-dependent transcription in response to stress | 43.42 |
| regulation of transcription from RNA polymerase II promoter in response to hypoxia | 41.00 |

**(D) MSigDB Pathway**
-log10(Binomial p value)

| | |
|---|---|
| Keratinocyte Differentiation | 60.72 |

**(E) Disease Ontology**
-log10(Binomial p value)

| | |
|---|---|
| keratosis | 73.66 |
| middle ear disease | 73.31 |
| cholesteatoma | 59.87 |
| epidermolysis bullosa | 36.77 |

**(F)**

| Motif | p-value | Best Match | Family | Score |
|---|---|---|---|---|
| | 1e-6754 | BATF | Basic leucine zipper factors (bZIP) | 0.987 |
| | 1e-564 | CEBPA | C/EBP-related | 0.974 |
| | 1e-464 | TEAD | Helix-turn-helix domains | 0.966 |
| | 1e-399 | KLF5 | Zinc-coordinating DNA-binding domains | 0.962 |
| | 1e-280 | AP-2gamma | Zinc-coordinating DNA-binding domains | 0.949 |
| | 1e-267 | Six1 | Helix-turn-helix domains | 0.925 |
| | 1e-257 | Nur77 | Zinc-coordinating DNA-binding domains | 0.900 |
| | 1e-243 | Atf1 | Basic leucine zipper factors (bZIP) | 0.941 |
| | 1e-233 | Tlx? | Not available | 0.883 |
| | 1e-222 | Forkhead_class | Not available | 0.940 |
| | 1e-169 | NF1-halfsite | Zinc-coordinating DNA-binding domains | 0.863 |
| | 1e-156 | GATA | Zinc-coordinating DNA-binding domains | 0.913 |
| | 1e-109 | GATA | Zinc-coordinating DNA-binding domains | 0.861 |
| | 1e-84 | Nkx3.1 | Helix-turn-helix domains | 0.611 |
| | 1e-79 | Hoxc9 | Helix-turn-helix domains | 0.904 |
| | 1e-73 | Isl1 | Helix-turn-helix domains | 0.691 |
| | 1e-73 | Hand1::Tcfe2a | Twist-like factors | 0.776 |
| | 1e-51 | NFATC2 | Immunoglobulin fold | 0.633 |
| | 1e-46 | Zfp128 | Zinc-coordinating DNA-binding domains | 0.580 |
| | 1e-45 | STAT5 | Immunoglobulin fold | 0.611 |

Figure S28: Functional and motif analysis of the Combined genome-wide predictions on cell line HelaS3. A: Distance from the predicted A-Es to gene TSSs. B: Distance from the predicted A-Ps to gene TSSs. C,D,E: Top 20 enriched biological processes, pathways, and diseases, respectively, in the predicted cell-specific CRRs. F: Enriched *de novo* motifs in the predicted cell specific CRRs. Column 4: families of best-matched TFs. Column 5: best match scores.

(A) Region-gene associations / Distance to TSS (kb) — bars: 1064, 11762, 9873, 3884, 3206, 10226, 16892, 1661; categories: <-500, -500 to -50, -50 to -5, -5 to 0, 0 to 5, 5 to 50, 50 to 500, >500; TSS

(B) Region-gene associations / Distance to TSS (kb) — bars: 916, 12051, 12273, 6278, 4498, 9789, 15769, 1414; categories: <-500, -500 to -50, -50 to -5, -5 to 0, 0 to 5, 5 to 50, 50 to 500, >500; TSS

(C) **GO Biological Process** — -log10(Binomial p value)
- protein-lipid complex subunit organization: 30.58
- quaternary ammonium group transport: 22.26

(D) **MSigDB Pathway** — -log10(Binomial p value)
- FOXA2 and FOXA3 transcription factor networks: 117.53
- Fatty acid metabolism: 30.93

(E) **Disease Ontology** — -log10(Binomial p value)

(F)

| Motif | p-value | Best Match | Family | Score |
|---|---|---|---|---|
| | 1e-1086 | HNF4A | Zinc-coordinating DNA-binding domains | 0.968 |
| | 1e-423 | Fox:Ebox | Helix-turn-helix domains | 0.904 |
| | 1e-352 | Hnf1 | Helix-turn-helix domains | 0.949 |
| | 1e-101 | Nr5a2 | Zinc-coordinating DNA-binding domains | 0.876 |
| | 1e-87 | Six1 | Helix-turn-helix domains | 0.592 |
| | 1e-74 | MafA | Basic leucine zipper factors (bZIP) | 0.632 |
| | 1e-68 | CEBPA | C/EBP-related | 0.901 |
| | 1e-63 | Nuclear_Receptor_class | Not available | 0.786 |
| | 1e-61 | Pou6f1 | Helix-turn-helix domains | 0.576 |
| | 1e-60 | KLF5 | Zinc-coordinating DNA-binding domains | 0.875 |
| | 1e-56 | T | Basic helix-loop-helix factors (bHLH) | 0.653 |
| | 1e-49 | Egr2 | Zinc-coordinating DNA-binding domains | 0.691 |
| | 1e-48 | EWS:ERG-fusion | Helix-turn-helix domains | 0.744 |
| | 1e-48 | Tlx? | Not available | 0.658 |
| | 1e-47 | Tcfap2e | Basic domains | 0.715 |
| | 1e-43 | NFKB1 | Immunoglobulin fold | 0.733 |
| | 1e-42 | Rfx3 | Helix-turn-helix domains | 0.670 |
| | 1e-42 | at_AC_acceptor | Not available | 0.575 |
| | 1e-41 | NFATC2 | Immunoglobulin fold | 0.892 |
| | 1e-40 | DCE_S_I | Not available | 0.637 |

Figure S29: Functional and motif analysis of the Combined genome-wide predictions on cell line HepG2. A: Distance from the predicted A-Es to gene TSSs. B: Distance from the predicted A-Ps to gene TSSs. C,D,E: Top 20 enriched biological processes, pathways, and diseases, respectively, in the predicted cell-specific CRRs. F: Enriched *de novo* motifs in the predicted cell specific CRRs. Column 4: families of best-matched TFs. Column 5: best match scores.

(A)

(B)

(C)

**GO Biological Process**
-log10(Binomial p value)

| 0 | 2 | 4 | 6 | 8 | 10 | 12 | 14 | 16 | 18 | 20 |

negative regulation of gene expression, epigenetic ... 20.81

(D)

**MSigDB Pathway**
-log10(Binomial p value)

(E)

**Disease Ontology**
-log10(Binomial p value)

(F)

| Motif | p-value | Best Match | Family | Score |
|---|---|---|---|---|
| AGATAAGAⱼⱼ | 1e-2706 | Gata2 | Zinc-coordinating DNA-binding domains | 0.980 |
| TGⱼTGAⱼTCAₜⱼ | 1e-630 | Bach2 | NF-E2-like factors | 0.974 |
| AAⱼAGGAAⱼT | 1e-238 | Sfpi1 | Helix-turn-helix domains | 0.921 |
| ₐₜGGCₐGₜₜₜCC | 1e-227 | KLF5 | Zinc-coordinating DNA-binding domains | 0.921 |
| GCACCTCAGATⱼ | 1e-202 | CREB | CREB-like factors | 0.600 |
| TTCₜₜₜGAAₜ | 1e-176 | STAT5 | Immunoglobulin fold | 0.948 |
| TAGTⱼAGATGAA | 1e-175 | Six1 | Helix-turn-helix domains | 0.581 |
| TAGGGTAACTTC | 1e-173 | Spz1 | Not available | 0.667 |
| TGTAAACTGTCA | 1e-170 | Myb | Helix-turn-helix domains | 0.763 |
| CGCACₐGTGCₐC | 1e-166 | Zfp161 | Zinc-coordinating DNA-binding domains | 0.602 |
| GAGTGCCAGACA | 1e-163 | Tbox:Smad | Not available | 0.805 |
| GGGCⱼCAGGₜCA | 1e-140 | RXR | Zinc-coordinating DNA-binding domains | 0.711 |
| GTACCⱼⱼGTT | 1e-133 | Nkx3-1 | Helix-turn-helix domains | 0.546 |
| ₐCₜAⱼₜCAGAAG | 1e-130 | DCE_S_I | Not available | 0.613 |
| TTTTAGACCATA | 1e-128 | Zfp187 | Zinc-coordinating DNA-binding domains | 0.593 |
| ₐTₜAGCATGCTA | 1e-126 | Pou2f3 | Helix-turn-helix domains | 0.639 |
| GAGCCGAAGCAG | 1e-121 | DCE_S_I | Not available | 0.611 |
| CATCTTⱼGCTCC | 1e-116 | E2F2 | Helix-turn-helix domains | 0.632 |
| GGTCTACAGCTⱼ | 1e-111 | DCE_S_II | Not available | 0.614 |
| GCAGAAATCACC | 1e-109 | Hdx | Helix-turn-helix domains | 0.749 |

Figure S30: Functional and motif analysis of the Combined genome-wide predictions on cell line K562. A: Distance from the predicted A-Es to gene TSSs. B: Distance from the predicted A-Ps to gene TSSs. C,D,E: Top 20 enriched biological processes, pathways, and diseases, respectively, in the predicted cell-specific CRRs. F: Enriched *de novo* motifs in the predicted cell specific CRRs. Column 4: families of best-matched TFs. Column 5: best match scores.

# References

[1] B. Arnett, P. Soisson, B.S. Ducatman, and P. Zhang. Expression of CAAT enhancer binding protein beta (C/EBP beta) in cervix and endometrium. *Molecular Cancer*, 2:21, 2003.

[2] M. Beger, K. Butz, C. Denk, T. Williams, H.C. Hurst, and F. Hoppe-Seyler. Expression pattern of AP-2 transcription factors in cervical cancer cells and analysis of their influence on human papillomavirus oncogene transcription. *Journal of Molecular Medicine*, 79(5-6):314–320, 2001.

[3] R.H. Costa, V.V. Kalinichenko, A.X. Holterman, and X. Wang. Transcription factors in liver development, differentiation, and regeneration. *Hepatology*, 38(6):1331–1347, 2003.

[4] A. Dev, S. Iyer, B. Razani, and G. Cheng. NF-$\kappa$B and innate immunity. *Current Topics in Microbiology and Immunology*, 349:115–143, 2011.

[5] J.D. Fleming, G. Pavesi, P. Benatti, C. Imbriano, R. Mantovani, and K. Struhl. NF-Y coassociates with FOS at promoters, enhancers, repetitive elements, and inactive chromatin regions, and is stereo-positioned with growth-controlling transcription factors. *Genome Research*, 23(8):1195–1209, 2013.

[6] D. Kleftogiannis, P. Kalnis, and V.B. Bajic. DEEP: A general compuational framework for predicting enhancers. *Nucleic Acids Research*, 43(1):e6, 2015.

[7] A.K. Linnemann, H. O'Geen, S. Keles, P.J. Farnham, and E.H. Bresnick. Genetic framework for GATA factor function in vascular biology. *Proceedings of the National Academy of Sciences*, 108(33):13641–13646, 2011.

[8] R. Lopez, E. Garrido, G. Vazquez, P. Pina, C. Perez, I. Alvarado, and M. Salcedo. A subgroup of HOX Abd-B gene is differentially expressed in cervical cancer. *International Journal Gynecological Cancer*, 16(3):1289–1296, 2006.

[9] A.L. Malt, J. Cagliero, K. Legent, J. Silber, A. Zider, and D. Flagiello. Alteration of TEAD1 expression levels confers apoptotic resistance through the transcriptional up-regulation of Livin. *PLoS One*, 7(9):e45498, 2012.

[10] E.V. Mityushova, N.D. Aksenov, and I.I. Marakhova. STAT5 in regulation of chronic leukemia K562 cell proliferation: Inhibitory effect of WHI-P131. *Cell and Tissue Biology*, 4(1):63–69, 2010.

[11] T. Okuda, M. Nishimura, M. Nakao, and Y. Fujita. RUNX1/AML1: A central player in hematopoiesis. *International Journal of Hematology*, 74(3):252–257, 2001.

[12] G.V. Samant, M.O. Schupp, M. Francois, S. Moleri, R.K. Kothinti, C.Z. Chun, I. Sinha, S. Sellars, N. Leigh, K. Pramanik, M.A. Horswill, I. Remadevi, K. Li, G.A. Wilkinson, N.M. Tabatabai, M. Beltrame, P. Koopman, and R. Ramchandran. Sox factors transcriptionally regulate ROBO4 gene expression in developing vasculature in zebrafish. *The Journal of Biological Chemistry*, 286:30740–30747, 2011.

[13] H. Schrem, J. Klempnauer, and J. Borlak. Liver-enriched transcription factors in liver function and development. Part II: The C/EBPs and D site-binding protein in cell cycle control, carcinogenesis, circadian gene regulation, liver regeneration, apoptosis, and liver-specific gene regulation. *Pharmacological Reviews*, 56(2):291–330, 2004.

[14] C.C. Thornton, F. Al-Rashed, D. Calay, G.M. Birdsey, A. Bauer, H. Mylroie, B.J. Morley, A.M. Randi, D.O. Haskard, J.J. Boyle, and J.C. Mason. Methotrexate-mediated activation of an AMPK-CREB-dependent pathway: A novel mechanism for vascular protection in chronic systemic inflammation. *Annals of the Rheumatics Diseases*, 75(2):439–448, 2015.

[15] P. Wang, J. Xu, and C. Zhang. CREB, a possible upstream regulator of Bcl-2 in trichosanthin-induced HeLa cell apoptosis. *Molecular Biology Reports*, 37(4):1891–1896, 2010.

[16] Z. Wang, E.P. Bishop, and P.A. Burke. Expression profile analysis of the inflammatory response regulated by hepatocyte nuclear factor $4\alpha$. *BMC Genomics*, 12:128, 2011.

[17] H. Zhang, C. Liu, Z. Zha, B. Zhao, J. Yao, S. Zhao, Y. Xiong, Q. Lei, and K. Guan. TEAD transcription factors mediate the function of TAZ in cell growth and epithelial-mesenchymal transition. *The Journal of Biological Chemistry*, 284(20):13355–13362, 2009.

[18] J. Zhu, D.M. Giannola, Y. Zhang, A.J. Rivera, and S.G. Emerson. NF-Y cooperates with USF1/2 to induce the hematopoietic expression of HOXB4. *Blood*, 102(7):2420–2427, 2003.