

Genotypes of cancer stem cells characterized by epithelial-to-mesenchymal transition and proliferation related functions

Chueh-Lin Hsu⁺¹, Feng-Hsiang Chung⁺¹, Chih-Hao Chen¹, Tzu-Ting Hsu¹, Szu-Mam Liu¹, Dao-Sheng Chung², Ya-Fen Hsu³, Chien-Long Chen⁴, Nianhan Ma^{*1}, and Hoong-Chien Lee^{*1,5,6}

¹Institute of Systems Biology and Bioinformatics, Department of Biomedical Science and Engineering, National Central University, Zhongli, Taiwan 32001

²Department of Radiation Oncology, Landseed Hospital, Taoyuan, Taiwan 324

³Department of Surgery, Landseed Hospital, Taoyuan, Taiwan 324

⁴Department of Nephrology, Landseed Hospital, Taoyuan, Taiwan 324

⁵Department of Physics, Chung Yuan Christian University, Zhongli, Taiwan 32023

⁶Center for Dynamical Biomarkers and Translational Medicine, National Central University, Zhongli, Taiwan 32001

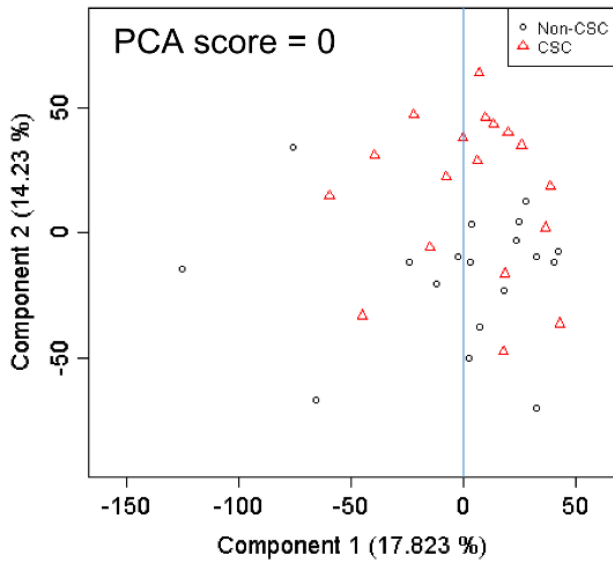
⁺CLH and FHC made equal contributions.

^{*}Corresponding authors. HCL: hcleee12345@gmail.com; NM: nianhan.ma@gmail.com.

Supplementary Information

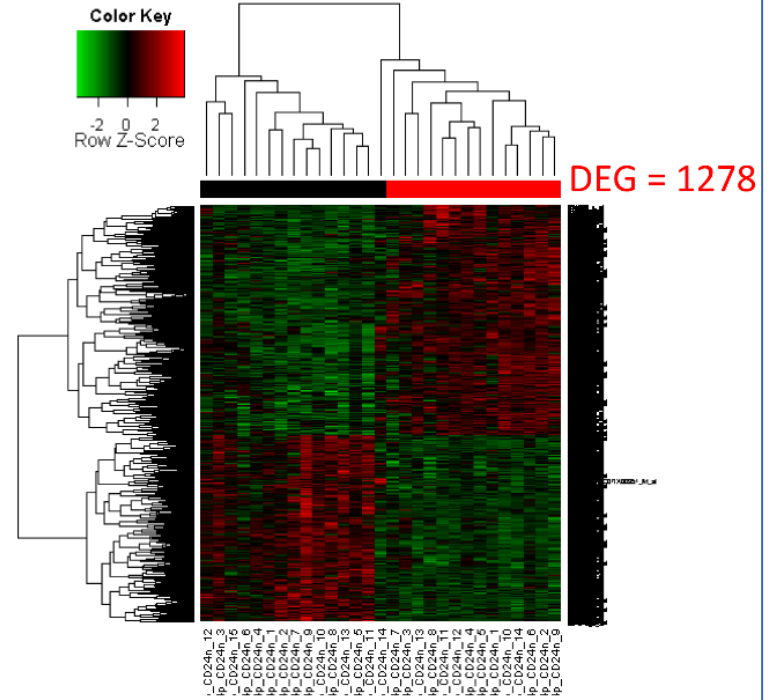
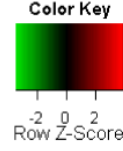
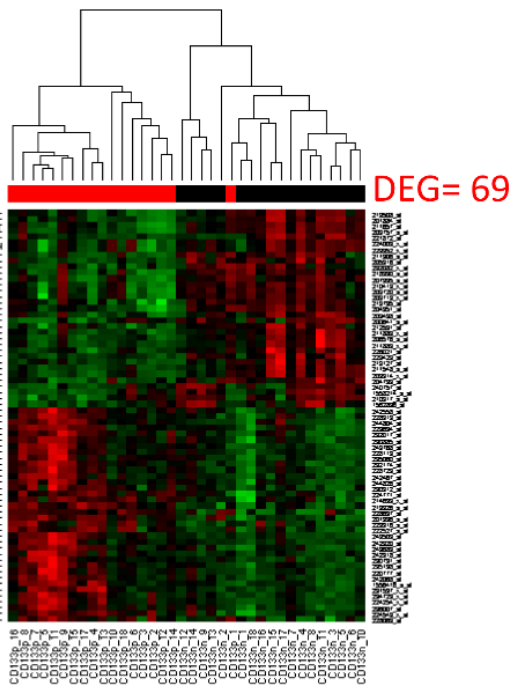
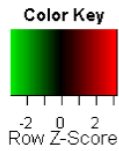
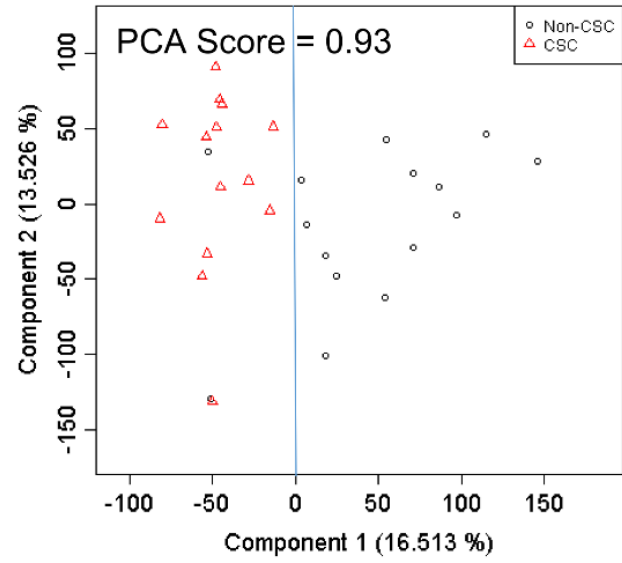
E-MEXP-993 (fail)

PCA



GSE7513 (success)

PCA



Supplementary Figure 1 High PCA-scoring datasets yielded larger differentially expressed gene (DEG) sets than low-scoring datasets

The dataset E-MEXP-933 had a PCA score of zero and yielded 69 DEGs. The dataset Breast_CD44_GSE7513 had a PCA score of 0.93 and yielded 1278 DEGs. Genes were selected using LIMMA [51], with false discovery rate < 0.05 and fold change > 2 .

Suppl. Fig. 2

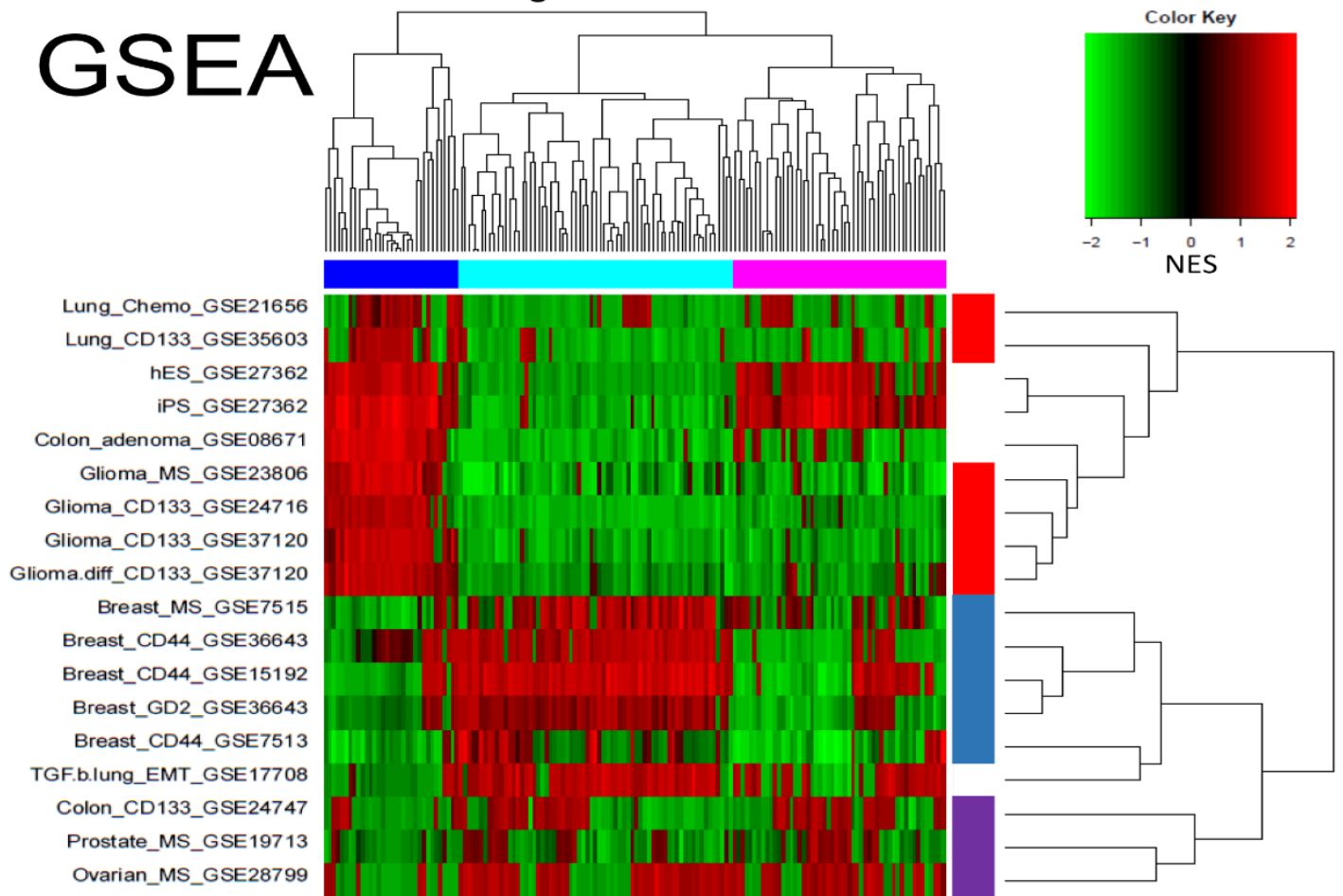
	TGF-b-lung_EMT_GSE17708 (2279)	iPS_GSE27362 (3060)	hES_GSE27362 (3106)	colon_adenoma_GSE08671 (1927)	Prostate_MS_GSE19713 (618)	Ovarian_MS_GSE28799 (2100)	Lung_Chemo_GSE21656 (123)	Lung_CD133_GSE35603 (5591)	Glioma-diff_CD133_GSE37120 (1291)	Glioma_MS_GSE23806 (2049)	Glioma_CD133_GSE37120 (1400)	Glioma_CD133_GSE24716 (606)	Colon_CD133_GSE24747 (247)	Breast_MS_GSE7515 (2877)	Breast_GD2_GSE36643 (465)	Breast_CD44_GSE7513 (1278)	Breast_CD44_GSE36643 (1137)	Breast_CD44_GSE15192 (3093)
DEG overlap(%)																		
Breast_CD44_GSE15192 (3093)	100	62	28	63	24	9.3	9.2	11	16	11	15	6.5	21	17	13	7.2	6.7	26
Breast_CD44_GSE36643 (1137)	23	100	14	64	13	2.8	3.3	3.6	6.3	4	6.4	1.6	9.9	6.6	5.2	2.9	2.6	10
Breast_CD44_GSE7513 (1278)	12	16	100	20	9	2.8	5.3	4.9	4.9	4.5	7.6	3.3	8.6	8.6	4.4	3.3	3.5	9
Breast_GD2_GSE36643 (465)	9.4	26	7.3	100	4.1	1.2	1	1.3	3.2	1.4	2.6	2.4	3.7	3.4	2.3	1	1	3
Breast_MS_GSE7515 (2877)	22	32	20	25	100	8.1	11	8.3	22	8.3	12	5.7	18	20	17	6.2	6.6	20
Colon_CD133_GSE24747 (247)	0.7	0.6	0.6	0.7	0.7	100	4	2.9	0.9	1.9	0.7	1.6	0.6	0.5	0.9	0.6	0.6	1
Glioma_CD133_GSE24716 (606)	1.8	1.8	2.5	1.3	2.3	9.7	100	18	2.1	14	2.1	2.4	1.5	1.6	2.4	1.5	1.2	1.7
Glioma_CD133_GSE37120 (1400)	4.9	4.4	5.4	3.9	4	16	41	100	5.3	47	4.6	5.7	4.2	4.5	5.1	3	2.9	4.3
Glioma_MS_GSE23806 (2049)	11	11	7.9	14	16	7.3	7.1	7.8	100	7.4	9.6	6.5	11	14	22	4.3	5.1	11
Glioma-diff_CD133_GSE37120 (1291)	4.5	4.6	4.5	3.9	3.7	9.7	29	44	4.7	100	4.1	2.4	3.9	4.4	5	2.9	2.7	3.9
Lung_CD133_GSE35603 (5591)	27	32	33	32	24	16	19	18	26	18	100	13	30	29	27	12	12	34
Lung_Chemo_GSE21656 (123)	0.3	0.2	0.3	0.7	0.2	0.8	0.5	0.5	0.4	0.2	0.3	100	0.3	0.3	0.4	1	1.3	0.4
Ovarian_MS_GSE28799 (2100)	14	18	14	17	13	5.3	5.3	6.4	11	6.3	11	5.7	100	21	12	5.2	4.9	11
Prostate_MS_GSE19713 (618)	3.3	3.6	4.2	4.5	4.3	1.2	1.7	2	4.3	2.1	3.2	1.6	6.2	100	4.5	1.3	1.3	2.7
colon_adenoma_GSE08671 (1927)	8.1	8.9	6.7	9.7	11	6.9	7.8	7.1	21	7.4	9.3	5.7	11	14	100	3.7	4.4	12
hES_GSE27362 (3106)	7.2	7.8	8.1	6.9	6.7	7.3	7.4	6.7	6.5	7.1	6.8	24	7.7	6.6	5.9	100	81	7.1
iPS_GSE27362 (3060)	6.7	7	8.4	6.7	7	7.7	6.3	6.4	7.6	6.4	6.8	32	7.2	6.6	6.9	80	100	7.3
TGF-b-lung_EMT_GSE17708 (2279)	19	20	16	15	16	8.9	6.4	7.1	12	6.8	14	6.5	12	10	14	5.2	5.4	100

Supplementary Figure 2 Percentage overlap matrix of DEG sets selected from 18 datasets

DEGs for each dataset was selected using LIMMA, with FDR < 0.05 and FC >2. An entry in the upper (lower) triangle of the overlap matrix is the ratio of the number of common genes to number of genes in DEG sets of corresponding dataset in the top row (column on left).

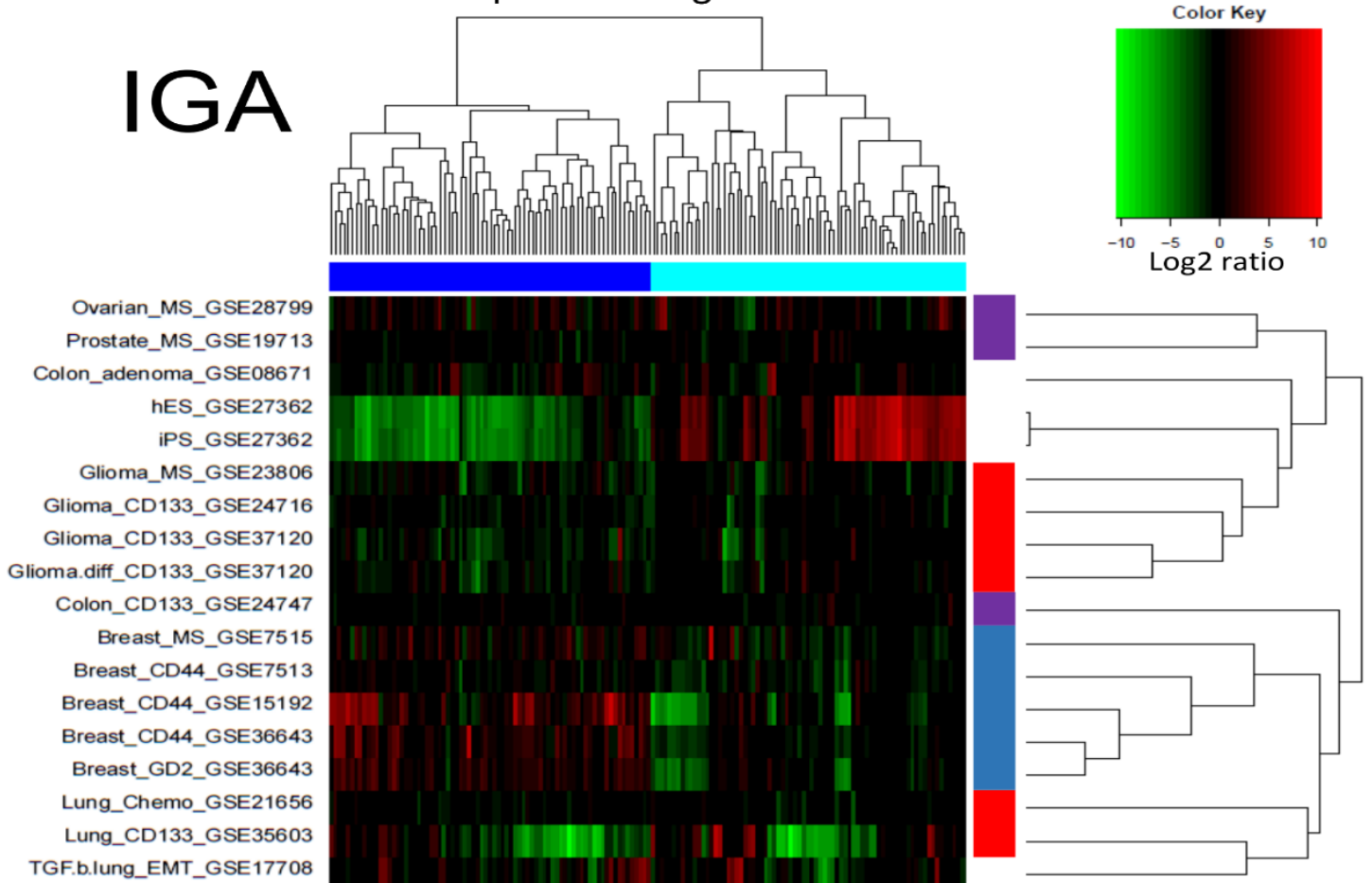
NOM p-val < 0.05
152 gene-sets

GSEA



Top-152 var. genes

IGA

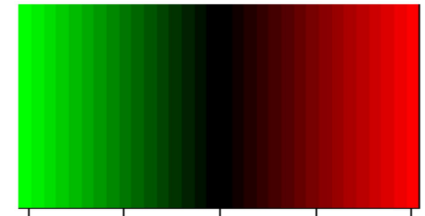


Supplementary Figure 3 Two-way hierarchical clusters of 18 datasets in IGA and GSEA

Top, in GSEA, dataset represented by NESs of 152 most significantly enriched/ depleted molecular signatures (nominal $p < 0.05$ by GSEA algorithm). Bottom, in IGA, dataset represented by log₂-ratios of 152 genes with the highest variance across datasets. Clustering was based on Pearson correlation of log₂-ratios in IGA and NESs in GSEA (Methods).

Suppl. Fig. 4

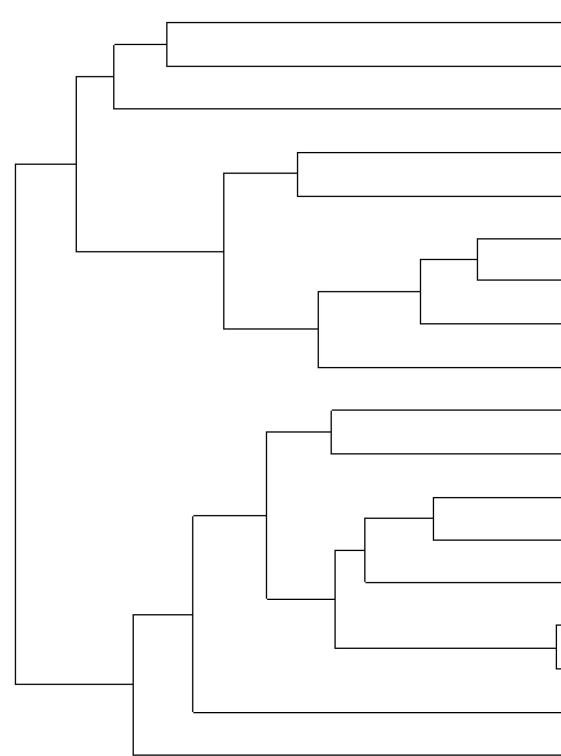
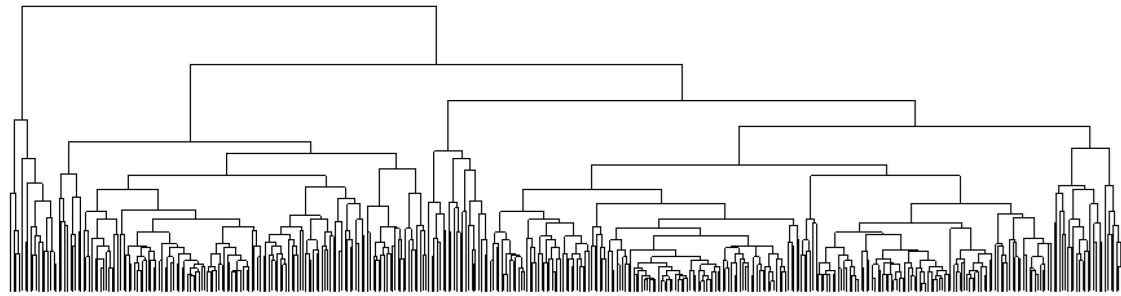
Color Key



-10 -5 0 5 10

Log2 ratio

FGS-398 var. genes

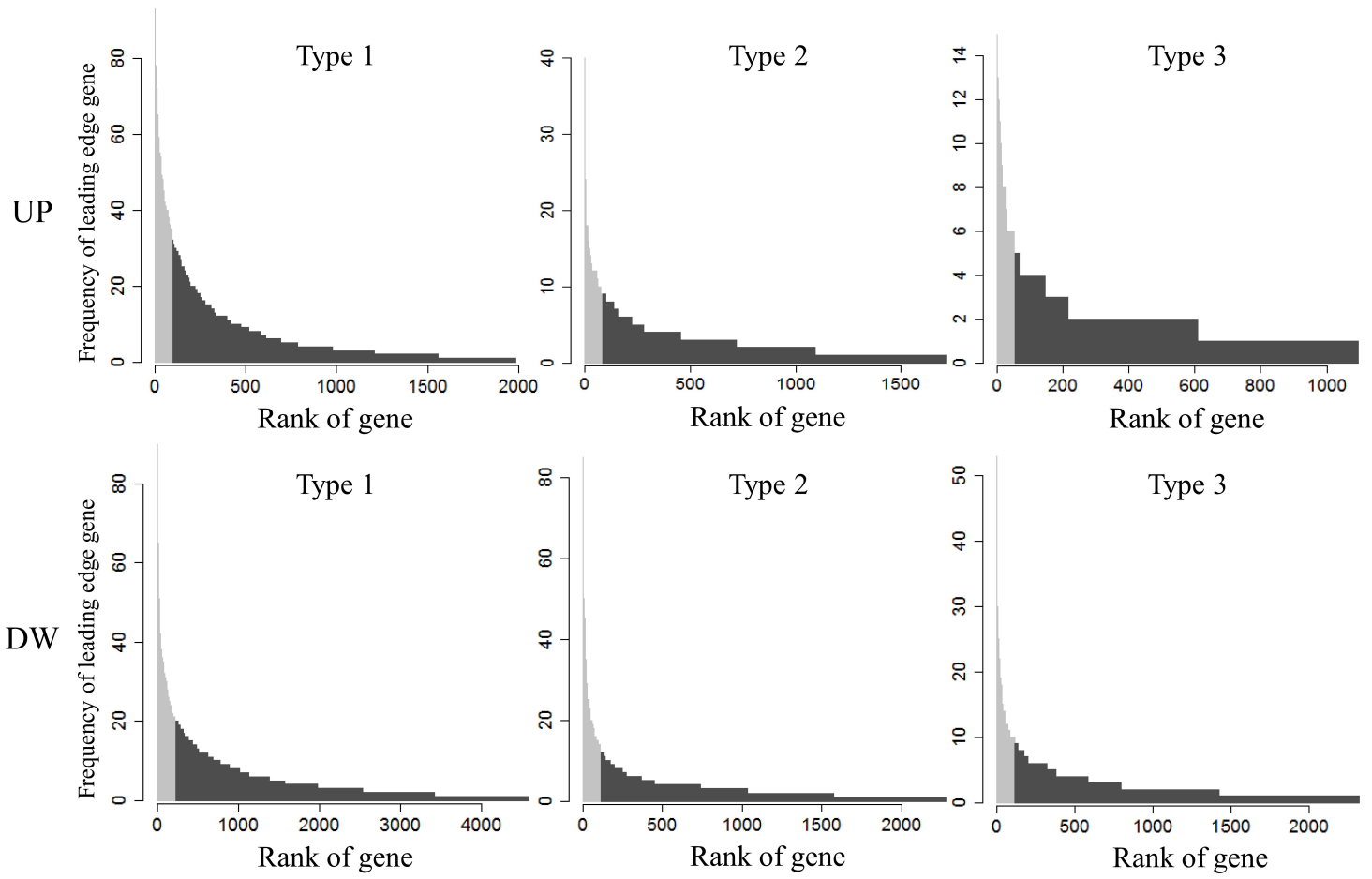


Prostate_MS_GSE19713
Ovarian_MS_GSE28799
Colon_CD133_GSE24747
TGF.b.lung_EMT_GSE17708
Breast_MS_GSE7515
Breast_GD2_GSE36643
Breast_CD44_GSE36643
Breast_CD44_GSE15192
Breast_CD44_GSE7513
Colon_adenoma_GSE08671
Glioma_MS_GSE23806
Glioma.diff_CD133_GSE37120
Glioma_CD133_GSE37120
Glioma_CD133_GSE24716
iPS_GSE27362
hES_GSE27362
Lung_CD133_GSE35603
Lung_Chemo_GSE21656

Supplementary Figure 4 Heatmap for 18 datasets against 398 genes selected from 152 signatures

The 398 genes are the highest occurrence frequency (top 5%) genes in the 152 signatures used in the heatmap of Fig. 2. Contents of the three cluster formed by the 18 datasets are identical to those in Fig. 2.

Suppl. Fig. 5



Supplementary Figure 5 Selection of highest occurrence frequency Up-regulated genes (URGs) and down-regulated genes (DRGs) for the three genotypes T1, T2, and T3

Genes in the PNS and NNS signatures (Fig. 5) are separately ranked by frequency of occurrence. Those having the top 5% frequencies in PNS (NNS) are selected as URGs/DRGs. The size of URG/DRG is 83/191 for T1, 70/91 for T2, 46/95 for T3.

Supplementary Table 1 Cluster gene sets (CGS) selected from three signature clusters in heatmap of Fig. 2.

Cluster	Gene list
Cluster 1 94 genes	CDC6,CCNA2,CDK2,CDC7,MCM3,MAD2L1,BUB1,AURKA,MCM2,CCNB1,UBE2C,SMC1A,HDAC2,CDK4,KIF11,DBF4,RRM1,MCM4,BIRC5,BUB1B,PKMYT1,CCNB2,SKP2,PCNA,HMGB2,PLK1,CDC20,RAN,MCM6,FEN1,RAD54L,E2F1,MCM5,ESPL1,POLD1,KIAA0101,RAD51,MCM7,TTK,CHEK1,E2F3,CDC25A,ZWINT,SSRP1,TOP2A,MYC,PRKDC,H2AFZ,CDC25C,KHDRBS1,BUB3,SMC3,KIF2C,NOLC1,RF,C4,TYMS,CENPF,NASP,POLA1,ANAPC10,TFDP1,TIPIN,STMN1,RACGAP1,TPX2,POLE,LSM4,CDH1,EED,NCAPH,USP1,NUSAP1,CDK6,SCARB1,NME1,SNRPA,CCT3,FBXO5,PRMT5,SMC4,CDC23,KIF22,TCPI1,CCT5,CDCA5,CCNE2,CDKN3,NEK2,TIMELESS,ERBB3,RPA3,CCNE1,CSE1L,C1QBP
Cluster 2 164 genes	CD44,NRP1,DAB2,PPAP2B,IL1R1,SERPINE1,COL3A1,TGFBI,PDGFRA,TPM1,THBS1,IGFBP3,COL1A2,FN1,S100A4,RRAS,VIM,CYP1B1,LAMB1,TIMP1,SRPX,SPARC,TIMP2,EMP3,CCL2,TNFRSF1A,GAS6,MYLK,TXNIP,ACTA2,BCL6,LUM,COL1A1,CLU,DKK3,C1S,TGFBI,LAMC1,CD59,MMP2,CD63,CDKN1A,GBP2,AKR1C3,COL6A1,MMP1,CYR61,CEBPD,PRSS23,B2M,DLC1,RAB31,SEMA3C,COL6A3,FTL,NNMT,TSC22D3,CTSB,MGLL,CDH11,OSMR,TRAM2,COL5A1,NRCAM,IGFBP5,PDGFRB,AKR1C1,FBN1,PLOD2,HTRA1,PMP22,FTH1,BTG1,ITGAV,POSTN,EPH2,ADM,TSC22D1,MARCKS,COL4A1,CALD1,GABARAPL1,PDE4DIP,LMO2,LOXL2,COL6A2,GPNMB,ITGA5,HSPB8,VEGFC,CCDC80,TRIM22,APOE,TPM2,PHLDA1,NUAK1,FAP,LPXN,IL6,ID2,NR3C1,MT1H,S100A6,ITGB5,RGL1,LOX,TAGLN,THBS2,TIMP3,SYNE1,SERPING1,FADS1,CDH2,SEMA5A,LOXL1,SNAI2,CPM,EMP1,ACTN1,CD151,DCN,CCND1,GAS1,LAMB2,EPHX1,GJA1,IL6ST,COL5A2,COL8A1,ACSL1,PCDH9,ZCCHC6,GSN,LRP1,MT1F,ANPEP,CD14,EFEMP1,ENPP2,CDH13,GEM,SMAD3,AHNAK,PAPSS2,SAT1,ANXA1,GLIPR1,ANXA5,SPOCK1,SHC1,ZBTB20,RAB13,RHOBTB3,TGFB3,MAN1A1,LTBP2,ARL6IP5,IFITM3,IFI16,ENG,SAT6,SEPP1,HLA-F,YPEL5
Cluster 3 157 genes	CDH1,CLDN7,KRT18,S100P,SPINT2,CLDN4,CD24,KRT19,SCNN1A,AGR2,MAP7,GATA3,ERBB3,PPL,SLC9A3R1,CA2,TPD52,MUC1,ELF3,VAV3,MAL2,KRT8,SOX4,GRB7,LAMB3,SLPI,MYB,SSH3,LAMC2,PKP3,TJP3,FXRD3,CLU,AQP3,TACSTD2,CLDN1,EGFR,DNAJC12,SORL1,KRT17,TRIM29,MYO5C,AKR1C1,FBP1,CDH3,MOSC1,PDK4,SELENBP1,MAOA,MT2A,KRT15,TPD52L1,FOXA1,S100A2,LAD1,KCNK1,SFN,DSP,JUP,TFAP2C,KRT5,PERP,BTG1,ID2,DDR1,TFPI2,CA12,SDC4,NEBL,ANXA3,RTKN,B3GNT3,AKR1C3,SERPIN5,FGFR3,MT1X,RND1,KRT14,DST,MATN2,KIAA1324,KRT6B,BIK,OCN,SFRP1,EHF,HOOK1,GPR160,MLPH,ITGB4,MMP7,ST14,MEST,AZGP1,KRT7,GPM6B,GABRP,SOX9,F2RL1,ARHGFE5,CLGN,ARL4C,LLGL2,GALNT3,F11R,ABCG2,TGM2,AKR1C2,IGFBP2,STC1,RAPGEF5,CECAM6,RNF43,PODXL,ETV4,NMU,LPL,LCN2,KLF5,TMEM30B,RAB25,CYB5A,SLC27A2,SYTL2,PRSS23,PHLDA1,SYT17,ANXA9,CAMK2N1,TFE3,SCUBE2,UCP2,BLVRB,WWC1,CNKSR1,CCND1,RHOD,ST6GAL1,ATF3,CHD7,CELSR1,HIPK2,LIMA1,MYO6,GPX3,MYLK,RHOB,STOM,S100A14,FLNB,SPINT1,DFNA5,TMEM45B,MBP,TGFA,KIAA0040,EPN3

	ABLIM1,AKT3,ANAPC10,ANXA1,APC,ATR,AVEN,BARD1,BCAT1,BCL2,BCL2L1,BIRC3,BIRC5,BNIP2,BNIP3, BRAF*,BUB1,CALD1,CASP3,CAV1*,CCNA2*,CCND1,CCNE1,CCNE2,CDC16,CDC23,CDC25B,CDC25C,CDC4
T3-DRG	2EP3,CDC6,CDK6*,CDK8,CDKN1A,CDKN1C,CENPF,CFLAR,CHEK1,CHUK,COL4A2,CRIM1,CUL1,DBF4,E2F
20* + 75	3,EGFR*,FLNA,FN1,FYN,GAB2*,GPX1,GRB2*,GSK3B*,GSTP1,HDAC3,HDAC9,HRAS*,HTATIP2,IGF1R,IKBK
genes	B,ITGAV,ITGB1,KRAS*,KRT7,MAD2L1,MAP2K1,MAPK1,MT2A,MYC,MYH9,NCK2,NRAS*,PAK2,PEA15*,PF KP,PIK3CA*,PIK3CB*,PIK3CD,PIK3R1*,PIK3R3,PTPN11,PTPN14,RB1,SAMD4A,SHC1*, SKP2,SOS1*,SPTBN1,STAT1,TAX1BP1,TCF4,THBS1,TIPIN,TPM1*,TUBB2A,YWHAH,ZFP161

Annotation of T1, T2 and T3 genes in Supplementary Table 2 Among the T1-URGs, CDC6* is an oncogenic cell-cycle gene ¹; MCM2*, MCM3*, MCM4, MCM5*, MCM6, and MCM7 are marker genes for proliferation ²; and MSH2* is a DNA repair enzyme gene ³. Among the T1-DRGs, CD44* and CD24* and markers of breast CSC ⁴ that indicate enhanced invasive capability ⁵; TGFBI promotes metastasis ⁶; APOE* (when expressed) indicates enhanced invasive and metastasis capabilities ⁷; and FN1* indicates metastatic outgrowth ⁸. All of the above suggest enhanced proliferation and suppressed metastasis. Among the T2-URGs, IGFBP3* and IGFBP4* express the presence of cytokines for mesenchymal stem cells ⁹; FBN1*, MMP2*, PDGFRA*, and TWIST1 indicate EMT activity ¹⁰; VIM* promotes cell morphogenesis during EMT ¹¹; TGFBI promotes metastasis ⁶; and FN1 indicates metastatic outgrowth ⁸. Among the T2-DRGs, CDH1*, ERBB3*, KRT14*, KRT18*, KRT5, KRT8*, MUC1*, and OCLN* are epithelial markers ¹²⁻¹⁵. All of these are indications of EMT activity. Among the T3-URGs, APOE* indicates enhanced invasive and metastasis capabilities ⁷; ABCA1* indicates resistance to chemo drugs ¹⁶; APOE*, ABCA1*, NR1H3*, and PLTP* indicate cell-migration related ¹⁷ cholesterol metabolism ¹⁸; AGR2* indicates tumorigenesis and induces metastasis ¹⁹. Among the T3-DRGs, BRAF*, EGFR*, KRAS*, MAP2K1, and PIK3CA* induce EMT leading to invasion ^{20,21}; KRAS* also promotes proliferations and suppresses diversification and apoptosis ²²; PEA15* regulates cell's invasive ability ²³. The signals from T3 are mixed; together it indicates suppression of proliferation, drug resistance, tumorigenesis, metastasis but suppression of EMT.

References cited in Supplementary Table 2

1. Malumbres, M. & Barbacid, M. Cell cycle, CDKs and cancer: a changing paradigm. *Nature Reviews Cancer* **9**, 153-166, doi:10.1038/nrc2602 (2009).
2. Blow, J. J. & Hodgson, B. Replication licensing—Origin licensing: defining the proliferative state? *Trends in cell biology* **12**, 72-78, doi:10.1016/S0962-8924(01)02203-6 (2002).
3. Fujii, H. *et al.* Sphere-forming stem-like cell populations with drug resistance in human sarcoma cell lines. *International journal of oncology* **34**, 1381-1386, doi:10.3892/ijo_00000265 (2009).
4. Al-Hajj, M., Wicha, M. S., Benito-Hernandez, A., Morrison, S. J. & Clarke, M. F. Prospective identification of tumorigenic breast cancer cells. *Proceedings of the National Academy of Sciences* **100**, 3983-3988, doi:10.1073/pnas.0530291100 (2003).
5. Sheridan, C. *et al.* CD44+/CD24-breast cancer cells exhibit enhanced invasive properties: an early step necessary for metastasis. *Breast Cancer Res* **8**, R59, doi:10.1186/bcr1610 (2006).
6. Ma, C. *et al.* Extracellular matrix protein β ig-h3/TGFBI promotes metastasis of colon cancer by enhancing cell extravasation. *Genes & development* **22**, 308-321, doi:10.1101/gad.1632008 (2008).
7. Sakashita, K. *et al.* Clinical significance of ApoE expression in human gastric cancer. *Oncology reports* **20**, 1313-1319, doi:10.3892/or_00000146 (2008).
8. Soikkeli, J. *et al.* Metastatic outgrowth encompasses COL-I, FN1, and POSTN up-regulation and assembly to fibrillar networks regulating cell adhesion, migration, and growth. *The American journal of pathology* **177**, 387-403, doi:10.2353/ajpath.2010.090748 (2010).
9. Liu, C.-H. & Hwang, S.-M. Cytokine interactions in mesenchymal stem cells from cord blood. *Cytokine* **32**, 270-279, doi:10.1016/j.cyto.2005.11.003 (2005).
10. Lehmann, B. D. *et al.* Identification of human triple-negative breast cancer subtypes and preclinical models for selection of targeted therapies. *The Journal of clinical investigation* **121**, 2750, doi:10.1172/JCI45014 (2011).
11. Mendez, M. G., Kojima, S.-I. & Goldman, R. D. Vimentin induces changes in cell shape, motility, and adhesion during the epithelial to mesenchymal transition. *The FASEB Journal* **24**, 1838-1851, doi:10.1096/fj.09-151639 (2010).
12. Martin, P. *et al.* Prostate epithelial Pten/TP53 loss leads to transformation of multipotential progenitors and epithelial to mesenchymal transition. *The American journal of pathology* **179**, 422-435 (2011).

13. Fuchs, B. C. *et al.* Epithelial-to-mesenchymal transition and integrin-linked kinase mediate sensitivity to epidermal growth factor receptor inhibition in human hepatoma cells. *Cancer research* **68**, 2391-2399, doi:10.1158/0008-5472.CAN-07-2460 (2008).
14. Guaita, S. *et al.* Snail Induction of Epithelial to Mesenchymal Transition in Tumor Cells Is Accompanied by MUC1 Repression and ZEB1 Expression. *Journal of Biological Chemistry* **277**, 39209-39216, doi:10.1074/jbc.M206400200 (2002).
15. Samavarchi-Tehrani, P. *et al.* Functional genomics reveals a BMP-driven mesenchymal-to-epithelial transition in the initiation of somatic cell reprogramming. *Cell stem cell* **7**, 64-77, doi:10.1016/j.stem.2010.04.015 (2010).
16. Gillet, J.-P., Efferth, T. & Remacle, J. Chemotherapy-induced resistance by ATP-binding cassette transporter genes. *Biochimica et Biophysica Acta (BBA)-Reviews on Cancer* **1775**, 237-262, doi:10.1016/j.bbcan.2007.05.002 (2007).
17. Zuo, W. & Chen, Y.-G. Specific activation of mitogen-activated protein kinase by transforming growth factor- β receptors in lipid rafts is required for epithelial cell plasticity. *Molecular biology of the cell* **20**, 1020-1029, doi:10.1091/mbc.E08-09-0898 (2009).
18. Legry, V. *et al.* Associations between common genetic polymorphisms in the liver X receptor alpha and its target genes with the serum HDL-cholesterol concentration in adolescents of the HELENA Study. *Atherosclerosis* **216**, 166-169, doi:10.1016/j.atherosclerosis.2011.01.031 (2011).
19. Liu, D., Rudland, P. S., Sibson, D. R., Platt-Higgins, A. & Barraclough, R. Human homologue of cement gland protein, a novel metastasis inducer associated with breast carcinomas. *Cancer research* **65**, 3796-3805, doi:10.1158/0008-5472.CAN-04-3823 (2005).
20. De Roock, W., De Vriendt, V., Normanno, N., Ciardiello, F. & Tejpar, S. KRAS, BRAF, PIK3CA, and PTEN mutations: implications for targeted therapies in metastatic colorectal cancer. *The lancet oncology* **12**, 594-603, doi:10.1016/S1470-2045(10)70209-6 (2011).
21. Berg, M. & Soreide, K. EGFR and downstream genetic alterations in KRAS/BRAF and PI3K/AKT pathways in colorectal cancer—implications for targeted therapy. *Discovery medicine* **14**, 207-214 (2012).
22. Irahara, N. *et al.* NRAS mutations are rare in colorectal cancer. *Diagnostic molecular pathology: the American journal of surgical pathology, part B* **19**, 157, doi:10.1097/PDM.0b013e3181c93fd1 (2010).
23. Glading, A., Koziol, J. A., Krueger, J. & Ginsberg, M. H. PEA-15 inhibits tumor cell invasion by binding to extracellular signal-regulated kinase 1/2. *Cancer research* **67**, 1536-1544, doi:10.1158/0008-5472.CAN-06-1378 (2007).

Supplementary Table 3 Number of regulation related terms in the top-30 up-regulated and down-regulated GO terms in the three CSC genotypes.

Type	T1		T2		T3	
	Up	Down	Up	Down	Up	Down
Number of terms	0	6	3	2	0	19
Range of $-\log p$	<7.68	16.7-21.0	7.01-7.29	2.99-6.48	<2	18.3-28.0