

Supporting Information

Nicolau et al. 10.1073/pnas.1102826108

SI Text

1. Microarray Data Analysis. We provide details for the microarray data analysis of the *Nederlands Kanker Instituut (NKI)* data (1) consisting of 295 tumors, the *Breast Cancer Normal (BCN)* data (2) consisting of 13 normal breast tissue samples, and the validation data sets *Ullevål University Hospital (ULL)* (3) consisting of 46 tumors of ductal histological type that had been in the study for longer than 10 mo and *HERSCH* (4) consisting of 188 primary breast tumors.

1.1. Data preprocessing. Data were retrieved, missing values imputed, then data were collapsed by UniGene cluster ID build 219, and genes present in both the tumor cohort and the normal data set were retained.

For *NKI*, data consisted of 24,479 GeneBank accession IDs on 295 tumor samples, all of which had at least 70% data. Missing data were imputed using a *knn* algorithm (5) with $k = 10$. Data were also transformed from the original \log_{10} values to \log_2 . Data were then collapsed (mean) by UniGene to the mean. The resulting data set consisted of 18,970 UniGene clusters.

For *BCN*, data from 13 normal tissue samples (nine nonneoplastic tissue from cancer patients, four reduction mammoplasty tissue) were retrieved with quality filters for each spot: (i) spot regression correlation $r > 0.6$, or (ii) channel 1 mean intensity/median background intensity > 1.5 , or (iii) channel 2 normalized (mean intensity/median background intensity) > 1.5 . Clones with 70% data were retained: 32,644 clone IDs. Missing data were imputed using a *knn* algorithm (5) with $k = 10$. Data were then collapsed by UniGene to 18,971 UniGene clusters. Of these, 12,237 UniGene IDs were in common with the *NKI* data set, and 17,441 were in common with the *ULL* data set (see below).

For *ULL*, data from 46 tumors were retrieved with quality filters for each spot: (i) spot regression correlation $r > 0.6$, or (ii) channel 1 mean intensity/median background intensity > 1.5 , or (i) channel 2 normalized (mean intensity/median background intensity) > 1.5 . Only clones with 70% good data were retained: 31,667 clone IDs. Missing data were imputed using a *knn* algorithm (5) with $k = 10$. Data were then combined with normal tissue data *BCN* and collapsed by UniGene to 17,441 UniGene clusters.

For *HERSCH*, data from 188 primary tumors were retrieved with quality filters for each spot: (i) spot regression correlation $r > 0.6$, or (ii) channel 1 mean intensity/median background intensity > 1.5 , or (iii) channel 2 normalized (mean intensity/median background intensity) > 1.5 . Only clones with 70% good data were retained: 32,644 clone IDs. Missing data were imputed using a *knn* algorithm (5) with $k = 10$. Data were then combined with normal tissue data *BCN* and collapsed by UniGene to 18,896 UniGene clusters.

1.2. Disease-Specific Genomic Analysis (DSGA). For *NKI* and *BCN*, data from tumors and normal tissue were combined along the common 12,237 UniGenes, and columns were normalized to have the magnitude of the mean vector magnitude of 13 normal tissue samples. The *Healthy State Model (HSM)* was constructed from normal tissue data $\{\vec{N}_1, \dots, \vec{N}_{13}\}$ as follows: *FLAT* construction (2) is a method to de-sparse the data in high dimensions by substituting for each normal tissue vector \vec{N}_i , its fit \hat{N}_i to a linear model in the other normal tissue vectors:

$$\hat{N}_i = \sum_{\substack{1 \leq j \leq 13 \\ j \neq i}} \beta_j \vec{N}_j.$$

This was shown to decrease noise in simulated data and help identify a good dimension reduction for *Principal Component*

Analysis (PCA). We use a method described in ref. 2 to compute the Wold invariant (6) designed to measure a version of signal-to-noise ratio:

$$W(l) = \left(\frac{\lambda_l^2}{\lambda_{l+1}^2 + \dots + \lambda_{13}^2} \right) \frac{(n-l-1)(13-l)}{(n+13-2l)}.$$

Fig. S1 plots $W(l)$ vs. the dimension l and shows a jump at $l = 10$, indicating that signal-to-noise ratio is higher at dimension 10, thereby justifying *PCA* dimension reduction of the *FLAT* normal data to 10. This produced the 10 dimensional *HSM*. Linear models are then used to compute the fitted tumor data matrix to the *HSM* (normal component *Nc.mat*) and the residuals (disease component *Dc.mat*). Along with tumor data, a leave-one-out procedure gives an estimate of the deviation of normal tissue data from the model of the healthy state *HSM*. Details of this procedure are found in ref. 2.

The validation data sets *ULL* and *HERSCH* were similarly transformed using the same normal data set *BCN*.

For gene thresholding, the 12,237 genes in the disease component matrix *Dc.mat* of tumors were reduced to 262 through the following method of testing for significance in deviation from the null hypothesis space. For each gene we computed the 5th and 95th percentiles of values in the disease components of the 295 tumors, and we recorded the larger of the two in absolute value and denoted the collection of these gene-by-gene deviations from normal by *MaxAbs595*. A histogram of these values is seen in Fig. S2. We then computed the 85th and 98th percentiles of *MaxAbs595* and denoted these as *relaxed* threshold and *stringent* threshold, respectively. A total of 1,836 genes exceeded the relaxed threshold, and 245 genes exceeded the stringent threshold. Genes were retained for further analysis if they passed the relaxed threshold and if they were also highly correlated ($r > 0.6$) to at least three genes that passed the stringent threshold. A total of 262 genes satisfied the condition. This method ensures that genes are retained in the analysis if they not only (i) deviate significantly from the null hypothesis space *HSM* but (ii) do so in groups of highly correlated genes. We denote the reduced matrix of disease component of *NKI* data: *nkiDc.mat*. The result of clustering the *nkiDc.mat* array and gene mean-centered can be found in supplementary folder Dataset S1: *nkiDc.AGmc.cdt*. It can be explored with *TreeView* (7), and all of the known clusters of genes can be observed, but because this is not germane to our present study we forgo any in-depth analysis of this clustering.

We did not follow the same thresholding procedure for the validation data sets *ULL* and *HERSCH*; rather, we found that of the 262 genes retained in the *NKI* data set, 255 genes were present in the *ULL* data and 221 in the *HERSCH* data.

1.3. Progression Analysis of Disease (PAD) on NKI. We give details of *PAD* on the reduced and *DSGA*-transformed *NKI* data matrix: *nkiDc.mat* of 295 tumors and 262 genes. First, this was combined with the leave-one-out matrix that estimates normal tissue: *bcnL1.mat*. The *Mapper* filter function was computed on each column vector, as explained in the main text (Eq. 2). The image space was then fragmented into 15 intervals, with 80% overlap. Two outputs of mapper were obtained: the first, which included all of the bins, can be found in Fig. 3 (main text). The second provides the tighter streamlined subset of *Mapper* output, by excluding all bins with only one data point in them. The two outputs appear side by side in Fig. S3.

1.4. Comparison with clustering. Although *Mapper* incorporates clustering at the local level, the final output captures a wide

Histogram of larger absolute value of
5th and 95th percentiles for each gene
NK/ DSGA Disease Component

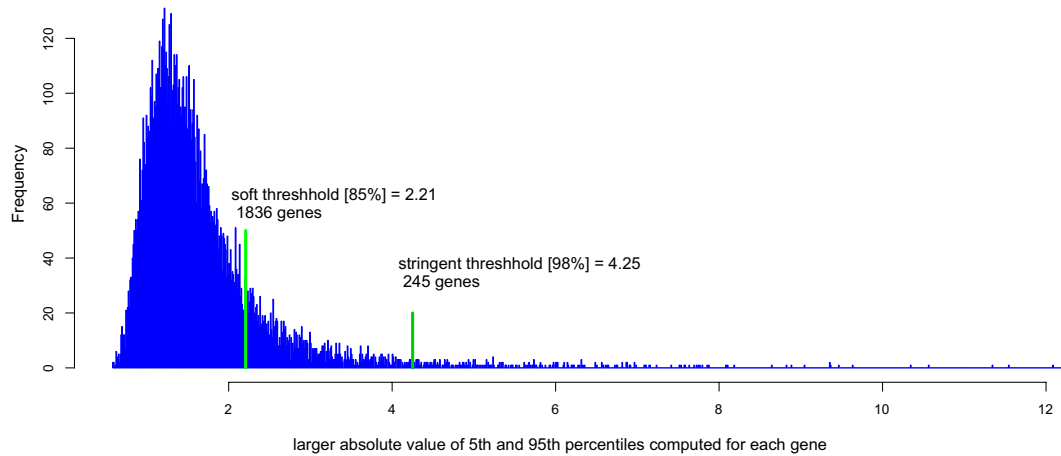


Fig. S2. For each gene, the 95th and 5th percentiles of expression levels in the disease component is computed. The larger of the two in absolute value denoted as Q_{gene} gives an estimate of the extent of deviation from normal for the gene. This deviation can be positive, indicating overexpression relative to normal levels, or negative, indicating underexpression relative to normal levels. The figure shows a histogram of the collection Q_{gene} of deviations from normal for the set of all genes. There are 1,836 genes for which this value exceeds the 85th percentile (*lax*-threshold genes) and 245 genes for which it exceeds the 95th percentile (*stringent*-threshold genes).

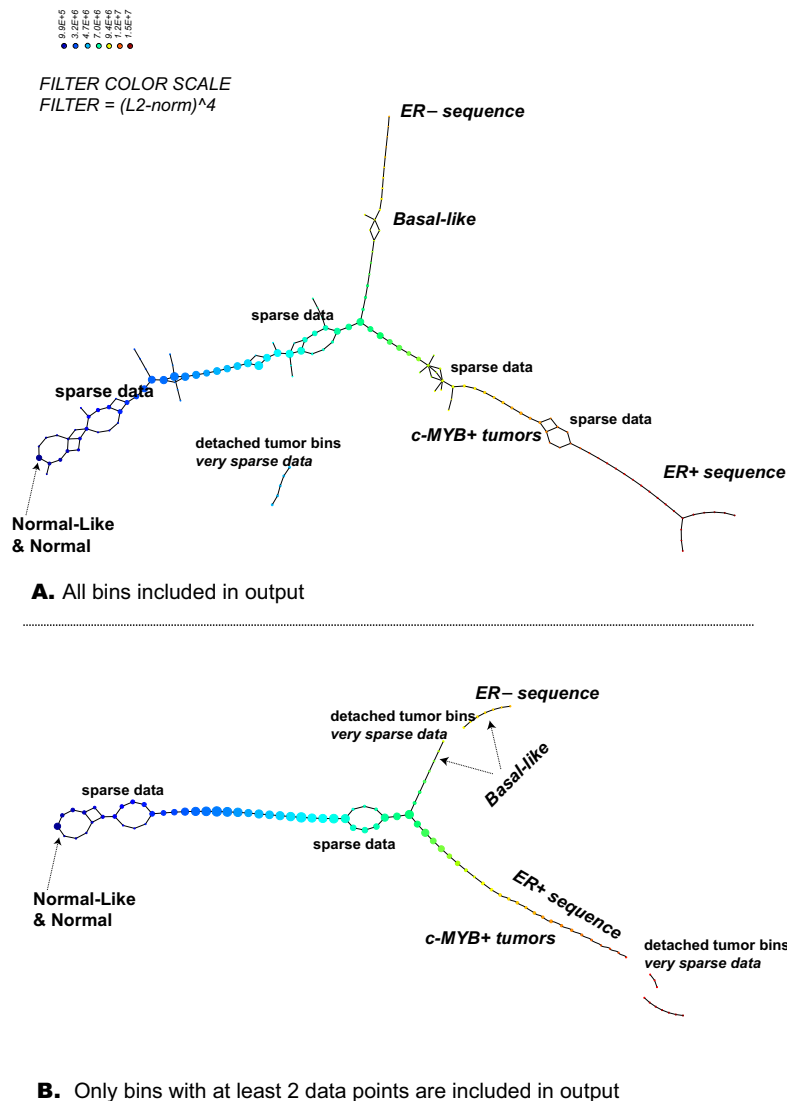


Fig. S3. (A) Complete output of the analysis. Each colored disk represents a bin containing several data points or patients. Thus, individual patients (data points) are not visible, and we only see bins containing collections of very similar points. This step provides a simplification of the original set of data points, because instead of showing a multitude of individual points, it shows a much smaller collection of bins, each bin containing a collection of very similar points. The size of bins relates to the number of data points contained in them. Thus, bins containing many data points appear as large discs, whereas bins with few points are drawn much smaller. When two bins have patients in common, an edge connects them. Thus, the bins provide a granularity to the overall set of data points, and the connections between bins, the edges that connect them, capture a rough shape of the data. Each data point has assigned to it a value of the *Mapper* xti filter function, and the bins are colored by the average value of this function for the points in the bin. The legend with assigned colors is seen at the top. In this particular example, each data point is a tumor sample, its gene expression transformed by *DSGA* to measure deviation from the *HSM*. The filter function is the overall amount of deviation from the *HSM*. Thus, red bins contain patients whose overall molecular profiles deviate a lot from normal, whereas blue bins contain patients whose profile is very close to normal. Sometimes data points are quite sparse, and this sparseness is visible in the output as well. When the data points become some what sparse, we see the graph fan out in a slight web-like feature. When data becomes really sparse, pieces of the graph become completely disconnected. Areas of local data sparseness are indicated in the figure. Finally, some bins are very small, containing only a few data points. To get a more streamlined, simplified picture, we can choose to ignore bins that are very small. This is similar to ignoring outliers. (B) Same output, but with bins containing single points not shown. Notice that this more streamlined version loses some of the sparseness information (for example that the long *ER*⁺ arm no longer exhibits the sparseness at the halfway point) and accentuates some sparseness areas by causing breaks in some places (for example the *Basal* arm now appears in two pieces).

PAM analysis c-MYB+ group vs. Normal

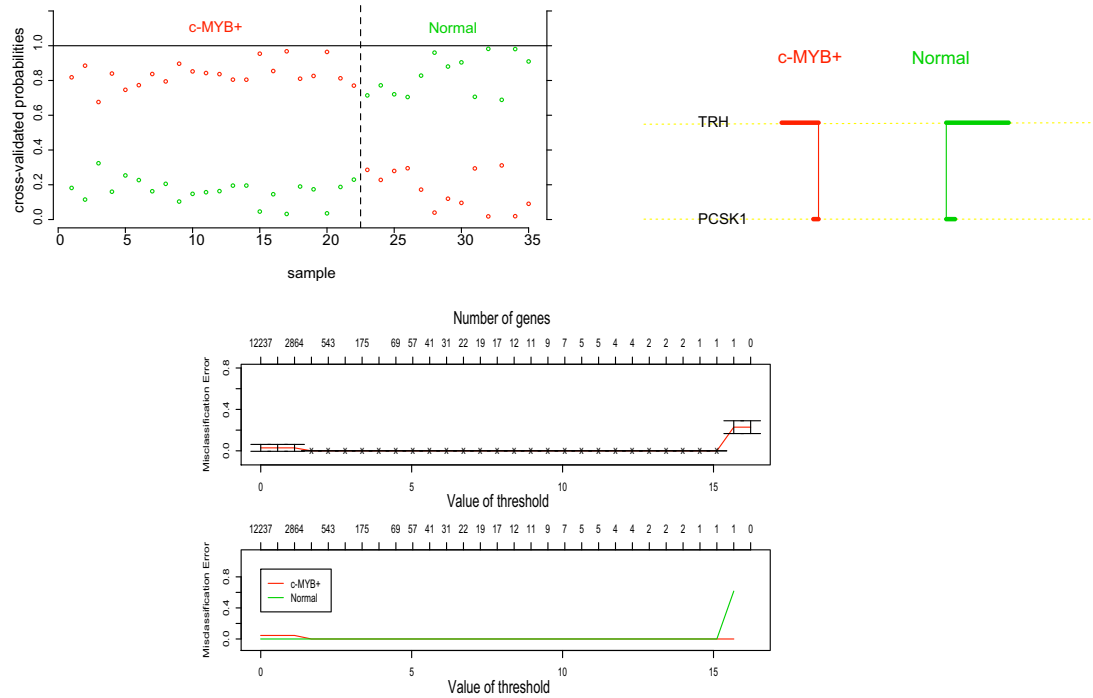


Fig. S5. Output of PAM analysis on the *c-MYB*⁺ group vs. *Normal* data. Two genes provide class prediction with *error rate* = 0: *TRH*, *TSH-releasing hormone*, and *PCSK1*, *proprotein convertase subtilisin kexin type 1*. The centroids, cross-validation probabilities, and misclassification error plots are shown.

Table S1. Genes significantly up-regulated and down-regulated in *MYB*+ vs. the rest of *ER*+ sequence

UniGene build 219	Gene symbol	q value	Gene information	MYB level vs. rest of <i>ER</i> + sequence
Hs.654446	<i>MYB</i>	0	MYB v-myb myeloblastosis viral oncogene homolog(avian)	Up
Hs.88417	<i>SUSD3</i>	0	SUSD3 Sushi domain containing 3	Up
Hs.414028	<i>C9orf116</i>	0	C9orf116 Chromosome 9 ORF 116	Up
Hs.532634	<i>IFI27</i>	5.15	IFI27 IFN, α -inducible protein 27 Hs.532634	Down
Hs.477891	<i>CPB1</i>	5.15	CPB1 Carboxypeptidase B1 (tissue) Hs.477891	Down
Hs.49760	<i>ORC6L</i>	5.15	ORC6L Origin recognition complex, subunit 6 like (yeast) Hs.49760	Down
Hs.517307	<i>MX1</i>	5.15	MX1 Myxovirus (influenza virus) resistance 1, IFN-inducible protein p78 (mouse) Hs.517307	Down
Hs.77367	<i>CXCL9</i>	5.15	CXCL9 Chemokine (C-X-C motif) ligand 9 Hs.77367	Down
Hs.501778	<i>TRIM22</i>	5.15	TRIM22 Tripartite motif-containing 22 Hs.501778	Down
Hs.521459	<i>ADAMDEC1</i>	5.15	ADAMDEC1 ADAM-like, decysin 1 Hs.521459	Down
Hs.458485	<i>ISG15</i>	5.15	ISG15 ISG15 ubiquitin-like modifier Hs.458485	Down
Hs.109225	<i>VCAM1</i>	5.15	VCAM1 Vascular cell adhesion molecule 1 Hs.109225	Down
Hs.17518	<i>RSAD2</i>	5.15	RSAD2 Radical S-adenosyl methionine domain containing 2 Hs.17518	Down
Hs.7155	<i>CMPK2</i>	5.15	CMPK2 Cytidine monophosphate (UMP-CMP) kinase 2, mitochondrial Hs.7155	Down
Hs.20315	<i>IFIT1</i>	6.51	IFIT1 IFN-induced protein with tetratricopeptide repeats 1 Hs.20315	Down
Hs.306777	<i>GSDMB</i>	6.51	GSDMB Gasdermin B Hs.306777	Down
Hs.715518	<i>STAT1</i>	6.51	STAT1 Signal transducer and activator of transcription 1, 91kDa Hs.715518	Down
Hs.709313	<i>B2M</i>	6.51	B2M Beta-2-microglobulin Hs.709313	Down
Hs.584823	<i>PLA2G7</i>	6.51	PLA2G7 Phospholipase A2, group VII (platelet-activating factor acetylhydrolase, plasma) Hs.584823	Down
Hs.181244	<i>HLA-A</i>	6.51	HLA-A Major histocompatibility complex, class I, A Hs.181244	Down
Hs.473341	<i>SAMSN1</i>	6.51	SAMSN1 SAM domain, SH3 domain and nuclear localization signals 1 Hs.473341	Down
Hs.523847	<i>IFI6</i>	6.51	IFI6 IFN, α -inducible protein 6 Hs.523847	Down
Hs.504641	<i>CD163</i>	6.51	CD163 CD163 molecule Hs.504641	Down
Hs.250615	<i>CYP2A6</i>	15.08	CYP2A6 Cytochrome P450, family 2, subfamily A, polypeptide 6 Hs.250615	Down
Hs.655652	<i>LILRB2</i>	15.08	LILRB2 Leukocyte Ig-like receptor, subfamily B (with TM and ITIM domains), member 2 Hs.655652	Down
Hs.459265	<i>ISG20</i>	15.08	ISG20 IFN stimulated exonuclease gene 20kDa Hs.459265	Down
Hs.926	<i>MX2</i>	15.08	MX2 Myxovirus (influenza virus) resistance 2 (mouse) Hs.926	Down
Hs.525157	<i>TNFSF13B</i>	15.08	TNFSF13B Tumor necrosis factor (ligand) superfamily, member 13b Hs.525157	Down
Hs.86859	<i>GRB7</i>	15.08	GRB7 Growth factor receptor-bound protein 7 Hs.86859	Down
Hs.352018	<i>TAP1</i>	15.08	TAP1 Transporter 1, ATP-binding cassette, subfamily B (MDR/TAP) Hs.352018	Down
Hs.32763	<i>GRIA2</i>	15.08	GRIA2 Glutamate receptor, ionotropic, AMPA 2 Hs.32763	Down
Hs.654585	<i>PSMB9</i>	15.08	PSMB9 Proteasome (prosome, macropain) subunit, β type, 9 (large multifunctional peptidase 2) Hs.654585	Down
Hs.718626	<i>KIF20A</i>	15.08	KIF20A Kinesin family member 20A Hs.718626	Down
Hs.474787	<i>IL2RB</i>	15.08	IL2RB Interleukin 2 receptor, β Hs.474787	Down
Hs.650174	<i>HLA-E</i>	15.08	HLA-E Major histocompatibility complex, class I, E Hs.650174	Down

Table S1. Cont.

UniGene build 219	Gene symbol	q value	Gene information	MYB level vs. rest of ER ⁺ sequence
Hs.143961	<i>CCL18</i>	15. 08	CCL18 Chemokine (C-C motif) ligand 18 (pulmonary and activation-regulated) Hs.143961	Down
Hs.81337	<i>LGALS9</i>	15. 08	LGALS9 Lectin, galactoside-binding, soluble, 9 Hs.81337	Down
Hs.474217	<i>CDC45L</i>	15. 08	CDC45L CDC45 cell division cycle 45-like (S. cerevisiae) Hs.474217	Down
Hs.301921	<i>CCR1</i>	15. 08	CCR1 Chemokine (C-C motif) receptor 1 Hs.301921	Down
Hs.16362	<i>P2RY6</i>	15. 08	P2RY6 Pyrimidinergic receptor P2Y, G protein coupled, 6 Hs.16362	Down
Hs.419259	<i>REC8</i>	15. 08	REC8 REC8 homolog (yeast) Hs.419259	Down
Hs.591742	<i>IL7R</i>	15. 08	IL7R Interleukin 7 receptor Hs.591742	Down
Hs.647962	<i>ZIC1</i>	18.67	ZIC1 Zic family member 1 (odd-paired homolog, Drosophila) Hs.647962	Down
Hs.43388	<i>RTP4</i>	18. 67	RTP4 Receptor (chemosensory) transporter protein 4 Hs.43388	Down
Hs.376208	<i>LTB</i>	18. 67	LTB Lymphotoxin β (TNF superfamily, member 3) Hs.376208	Down
Hs.14623	<i>IFI30</i>	18. 67	IFI30 IFN, γ-inducible protein 30 Hs.14623	Down
Hs.660866	<i>CTSL2</i>	18. 67	CTSL2 Cathepsin L2 Hs.660866	Down
Hs.278658	<i>KRT86</i>	18. 67	KRT86 Keratin 86 Hs.278658	Down
Hs.1051	<i>GZMB</i>	18. 67	GZMB Granzyme B (granzyme 2, cytotoxic T lymphocyte-associated serine esterase 1) Hs.1051	Down
Hs.1594	<i>CENPA</i>	18. 67	CENPA Centromere protein A Hs.1594	Down
Hs.161985	<i>TMPRSS4</i>	18. 67	TMPRSS4 Transmembrane protease, serine 4 Hs.161985	Down
Hs.153752	<i>CDC25B</i>	18. 67	CDC25B Cell division cycle 25 homolog B (S. pombe) Hs.153752	Down
Hs.446352	<i>ERBB2</i>	18. 67	ERBB2 V-erb-b2 erythroblastic leukemia viral oncogene homolog 2, neuro/glioblastoma derived oncogene homolog (avian) Hs.446352	Down
Hs.497599	<i>WARS</i>	18. 67	WARS Tryptophanyl-tRNA synthetase Hs.497599	Down
Hs.182231	<i>TRH</i>	18. 67	TRH TSH-releasing hormone Hs.182231	Down
Hs.521903	<i>LY6E</i>	20.44	LY6E Lymphocyte antigen 6 complex, locus E Hs.521903	Down
Hs.370036	<i>CCR7</i>	20. 44	CCR7 Chemokine (C-C motif) receptor 7 Hs.370036	Down

Table S2. Genes significantly up-regulated and down-regulated in MYB+ vs. Normal tissue

UniGene build 219	Gene symbol	q value	Gene information	MYB level vs. normal
Hs.414028	<i>C9orf116</i>	0	C9orf116 Chromosome 9 ORF 116 Hs.414028	Up
Hs.406050	<i>DNALI1</i>	0	DNALI1 Dynein, axonemal, light intermediate chain 1 Hs.406050	Up
Hs.163484	<i>FOXA1</i>	0	FOXA1 Forkhead box A1 Hs.163484	Up
Hs.76704	[Hs.76704]	0	NA Transcribed locus Hs.76704	Up
Hs.654446	<i>MYB</i>	0	MYB V-myb myeloblastosis viral oncogene homolog (avian) Hs.654446	Up
Hs.88417	<i>SUSD3</i>	0	SUSD3 Sushi domain containing 3 Hs.88417	Up
Hs.494496	<i>FBP1</i>	0	FBP1 Fructose-1,6-bisphosphatase 1 Hs.494496	Up
Hs.448520	<i>SLC7A2</i>	0	SLC7A2 Solute carrier family 7 (cationic amino acid transporter, y+ system), member 2 Hs.448520	Up
Hs.534847	<i>C4A</i>	0	C4A Complement component 4A (Rodgers blood group) Hs.534847	Up
Hs.496240	<i>AR</i>	0	AR Androgen receptor Hs.496240	Up
Hs.631650	<i>GLT8D2</i>	0	GLT8D2 Glycosyltransferase 8 domain containing 2 Hs.631650	Up
Hs.91109	<i>PRR15</i>	0	PRR15 Proline rich 15 Hs.91109	Up
Hs.387057	<i>THSD4</i>	0	THSD4 Thrombospondin, type I, domain containing 4 Hs.387057	Up
Hs.98265	<i>ST6GAL2</i>	0	ST6GAL2 ST6 β -galactosamide α -2,6-sialyltransferase 2 Hs.98265	Up
Hs.208124	<i>ESR1</i>	0	ESR1 Estrogen receptor 1 Hs.208124	Up
Hs.111779	<i>SPARC</i>	0	SPARC Secreted protein, acidic, cysteine-rich (osteonectin) Hs.111779	Up
Hs.480819	<i>TBC1D9</i>	0	TBC1D9 TBC1 domain family, member 9 (with GRAM domain) Hs.480819	Up
Hs.437638	<i>XBP1</i>	0	XBP1 X-box binding protein 1 Hs.437638	Up
Hs.444414	<i>AFF3</i>	0	AFF3 AF4/FMR2 family, member 3 Hs.444414	Up
Hs.524134	<i>GATA3</i>	0	GATA3 GATA binding protein 3 Hs.524134	Up
Hs.467733	<i>GREB1</i>	0	GREB1 GREB1 protein Hs.467733	Up
Hs.458573	<i>PDGFRL</i>	0	PDGFRL Platelet-derived growth factor receptor-like Hs.458573	Up
Hs.210995	<i>CA12</i>	0	CA12 Carbonic anhydrase XII Hs.210995	Up
Hs.523468	<i>SCUBE2</i>	0	SCUBE2 Signal peptide, CUB domain, EGF-like 2 Hs.523468	Up
Hs.654370	<i>FAP</i>	0	FAP Fibroblast activation protein, α Hs.654370	Up
Hs.489142	<i>COL1A2</i>	0	COL1A2 Collagen, type I, α 2 Hs.489142	Up
Hs.416108	<i>CRKRS</i>	0	CRKRS Cdc2-related kinase, arginine/serine-rich Hs.416108	Up
Hs.371147	<i>THBS2</i>	0	THBS2 Thrombospondin 2 Hs.371147	Up
Hs.519601	<i>ID4</i>	0	ID4 Inhibitor of DNA binding 4, dominant negative helix-loop-helix protein Hs.519601	Up
Hs.100686	<i>AGR3</i>	0	AGR3 Anterior gradient homolog 3 (<i>Xenopus laevis</i>) Hs.100686	Up
Hs.435655	<i>ASPN</i>	0	ASPN Asporin Hs.435655	Up
Hs.425777	<i>UBE2L6</i>	0	UBE2L6 Ubiquitin-conjugating enzyme E2L 6 Hs.425777	Up
Hs.659093	[Hs.659093]	0	NA Transcribed locus Hs.659093	Up
Hs.93764	<i>CPA4</i>	0	CPA4 Carboxypeptidase A4 Hs.93764	Up
Hs.719277	<i>SLC39A6</i>	0	SLC39A6 Solute carrier family 39 (zinc transporter), member 6 Hs.719277	Up
Hs.604376	[Hs.604376]	0	NA Transcribed locus Hs.604376	Up
Hs.95612	<i>DSC2</i>	0	DSC2 Desmocollin 2 Hs.95612	Up
Hs.8059	<i>SYT4</i>	0	SYT4 Synaptotagmin IV Hs.8059	Up
Hs.1925	<i>DSG3</i>	0	DSG3 Desmoglein 3 (pemphigus vulgaris antigen) Hs.1925	Up
Hs.8786	<i>CHST2</i>	0	CHST2 Carbohydrate (N-acetylglucosamine-6-O) sulfotransferase 2 Hs.8786	Up
Hs.24950	<i>RGSS5</i>	0	RGSS5 Regulator of G protein signaling 5 Hs.24950	Up
Hs.19492	<i>PCDH8</i>	0	PCDH8 Protocadherin 8 Hs.19492	Up
Hs.520339	<i>COL10A1</i>	0	COL10A1 Collagen, type X, α 1 Hs.520339	Up
Hs.5210	<i>GMFG</i>	0.46	GMFG Glia maturation factor, γ Hs.5210	Up
Hs.497636	<i>LAMB3</i>	0.46	LAMB3 Laminin, β 3 Hs.497636	Up
Hs.6360	<i>TMCC2</i>	0.46	TMCC2 Transmembrane and coiled-coil domain family 2 Hs.6360	Up
Hs.34526	<i>CXCR6</i>	0.46	CXCR6 Chemokine (C-X-C motif) receptor 6 Hs.34526	Up
Hs.504115	<i>TRIM29</i>	0.85	TRIM29 Tripartite motif-containing 29 Hs.504115	Up

Table S2. Cont.

UniGene build 219	Gene symbol	q value	Gene information	MYB level vs. normal
Hs.1787	<i>PLP1</i>	0.85	PLP1 Proteolipid protein 1 Hs.1787	Up
Hs.523500	<i>CD2</i>	0.85	CD2 CD2 molecule Hs.523500	Up
Hs.131431	<i>EIF2AK2</i>	0.85	EIF2AK2 Eukaryotic translation initiation factor 2- α kinase 2 Hs.131431	Up
Hs.136348	<i>POSTN</i>	0.85	POSTN Periostin, osteoblast specific factor Hs.136348	Up
Hs.193235	<i>CPLX2</i>	0.85	CPLX2 Complexin 2 Hs.193235	Up
Hs.438	<i>MEOX1</i>	1.94	MEOX1 Mesenchyme homeobox 1 Hs.438	Up
Hs.405614	<i>CTHRC1</i>	1.94	CTHRC1 Collagen triple helix repeat containing 1 Hs.405614	Up
Hs.182231	<i>TRH</i>	0	TRH TSH-releasing hormone Hs.182231	Down
Hs.477891	<i>CPB1</i>	0	CPB1 Carboxypeptidase B1 (tissue) Hs.477891	Down
Hs.78977	<i>PCSK1</i>	0	PCSK1 Proprotein convertase subtilisin/kexin type 1 Hs.78977	Down
Hs.250615	<i>CYP2A6</i>	0	CYP2A6 Cytochrome P450, family 2, subfamily A, polypeptide 6 Hs.250615	Down
Hs.26770	<i>FABP7</i>	0	FABP7 Fatty acid binding protein 7, brain Hs.26770	Down
Hs.516874	<i>CHGB</i>	0	CHGB Chromogranin B (secretogranin 1) Hs.516874	Down
Hs.150793	<i>CHGA</i>	0	CHGA Chromogranin A (parathyroid secretory protein 1) Hs.150793	Down
Hs.77367	<i>CXCL9</i>	0	CXCL9 Chemokine (C-X-C motif) ligand 9 Hs.77367	Down
Hs.496843	<i>VGLL1</i>	0	VGLL1 Vestigial like 1 (<i>Drosophila</i>) Hs.496843	Down
Hs.268728	<i>TTYH1</i>	0	TTYH1 Tweety homolog 1 (<i>Drosophila</i>) Hs.268728	Down
Hs.416073	<i>S100A8</i>	0	S100A8 S100 calcium binding protein A8 Hs.416073	Down
Hs.473341	<i>SAMSN1</i>	0	SAMSN1 SAM domain, SH3 domain and nuclear localization signals 1 Hs.473341	Down
Hs.517307	<i>MX1</i>	0	MX1 Myxovirus (influenza virus) resistance 1, IFN-inducible protein p78 (mouse) Hs.517307	Down
Hs.532634	<i>IFI27</i>	0	IFI27 IFN, α -inducible protein 27 Hs.532634	Down
Hs.143961	<i>CCL18</i>	0	CCL18 Chemokine (C-C motif) ligand 18 (pulmonary and activation-regulated) Hs.143961	Down
Hs.458485	<i>ISG15</i>	0	ISG15 ISG15 ubiquitin-like modifier Hs.458485	Down
Hs.192859	<i>PCDH10</i>	0	PCDH10 Protocadherin 10 Hs.192859	Down
Hs.419259	<i>REC8</i>	0	REC8 REC8 homolog (yeast) Hs.419259	Down
Hs.470654	<i>CDCA7</i>	0	CDCA7 Cell division cycle associated 7 Hs.470654	Down
Hs.32763	<i>GRIA2</i>	0	GRIA2 Glutamate receptor, ionotropic, AMPA 2 Hs.32763	Down
Hs.415762	<i>LY6D</i>	0	LY6D Lymphocyte antigen 6 complex, locus D Hs.415762	Down
Hs.119689	<i>CGA</i>	0	CGA Glycoprotein hormones, α polypeptide Hs.119689	Down
Hs.278658	<i>KRT86</i>	0	KRT86 Keratin 86 Hs.278658	Down
Hs.17518	<i>RSAD2</i>	0	RSAD2 Radical S-adenosyl methionine domain containing 2 Hs.17518	Down
Hs.7155	<i>CMPK2</i>	0	CMPK2 Cytidine monophosphate (UMP-CMP) kinase 2, mitochondrial Hs.7155	Down
Hs.20315	<i>IFIT1</i>	0	IFIT1 IFN-induced protein with tetratricopeptide repeats 1 Hs.20315	Down
Hs.418167	<i>ALB</i>	0	ALB Albumin Hs.418167	Down
Hs.372578	<i>FAM65C</i>	0	FAM65C Family with sequence similarity 65, member C Hs.372578	Down
Hs.26225	<i>GABRP</i>	0	GABRP Gamma-aminobutyric acid (GABA) A receptor, ρ 1 Hs.26225	Down
Hs.151254	<i>KLK7</i>	0	KLK7 Kallikrein-related peptidase 7 Hs.151254	Down
Hs.161985	<i>TMPRSS4</i>	0	TMPRSS4 Transmembrane protease, serine 4 Hs.161985	Down
Hs.376208	<i>LTB</i>	0	LTB Lymphotoxin β (TNF superfamily, member 3) Hs.376208	Down
Hs.414629	<i>CCL13</i>	0	CCL13 Chemokine (C-C motif) ligand 13 Hs.414629	Down
Hs.521459	<i>ADAMDEC1</i>	0	ADAMDEC1 ADAM-like, decysin 1 Hs.521459	Down
Hs.79361	<i>KLK6</i>	0	KLK6 Kallikrein-related peptidase 6 Hs.79361	Down
Hs.112405	<i>S100A9</i>	0	S100A9 S100 calcium binding protein A9 Hs.112405	Down
Hs.49760	<i>ORC6L</i>	0	ORC6L Origin recognition complex, subunit 6 like (yeast) Hs.49760	Down
Hs.647962	<i>ZIC1</i>	0	ZIC1 Zic family member 1 (odd-paired homolog, <i>Drosophila</i>) Hs.647962	Down
Hs.30743	<i>PRAME</i>	0	PRAME Preferentially expressed antigen in melanoma Hs.30743	Down

Table S3. Testing the MYB signature genes

Gene symbol	Gene name	UniGene build 219	pval MYB ⁺ group_UP Normal_LO
<i>MYC</i>	V-myc myelocytomatosis viral oncogene homolog	Hs.202453	0.24
<i>MYB</i>	V-myb myeloblastosis viral oncogene homolog	Hs.654446	4.70E-05
<i>ADA</i>	Adenosine deaminase	Hs.654536	1.60E-10
<i>CDK1</i>	Cyclin-dependent kinase 1	Hs.334562	0.00019
<i>POLD1</i>	Polymerase (DNA directed), δ 1	Hs.279413	2.60E-11
<i>PRTN3</i>	Myeloblastin proteinase 3	Hs.928	0.00014
<i>CD4</i>	T-cell surface antigen T4/Leu-3	Hs.631659	1
<i>VEGF</i>	Vascular endothelial growth factor A	Hs.73793	0.62
<i>BCL2</i>	B-cell CLL/lymphoma 2	Hs.150749	0.97
<i>KIT</i>	Proto-oncogene c-Kit mast/stem cell growth factor receptor	Hs.479754	1
<i>CD34</i>	Hematopoietic progenitor cell antigen CD34	Hs.374990	1
<i>GATA3</i>	Transacting T-cell-specific transcription factor GATA-3	Hs.524134	0.00048
<i>MPO</i>	Myeloperoxidase	Hs.458272	0.012
<i>HSP70</i>	HSPA4 heat shock 70kDa protein 4	Hs.90093	0.00064
<i>H2A.Z</i>	H2AZ histone	Hs.119192	0.00028
<i>Adora2B</i>	Adenosine receptor 2B – chicken	Hs.167046	0.01
<i>Mcm4</i>	CDC21; CDC54; MGC33310; P1-CDC21; hCdc21	Hs.460184	2.20E-05
<i>GAS41</i>	YEATS4;Yeats domeain containing 4	Hs.4029	0.00078
<i>NMU</i>	Neuromedin U	Hs.418367	0.38
<i>CCNE1</i>	Cyclin E1	Hs.244723	0.00049
<i>CCNB1</i>	cyclin B1	Hs.23960	0.021
<i>CA1</i>	Carbonic anhydrase 1	Hs.23118	0.00037
<i>PDCD4</i>	Programmed cell death 4(neoplastic transformation inhibitor)	Hs.711490	0.019
<i>COL1A1</i>	Collagen type I, α 1	Hs.172928	1
<i>COL1A2</i>	Collagen type I, α 2	Hs.489142	1
<i>CD13</i> <i>ANPEP</i>	Ananyl (membrane) aminopeptidase	Hs.1239	0.96
<i>GBX2</i>	Gastrulation brain homeobox 2	Hs.184945	0.61
<i>Actn1</i>	Actinin, α 1	Hs.509765	0.9
<i>Birc3</i>	Baculoviral IAP repeat-containing 3	Hs.127799	1
<i>Casp6</i>	caspase 6, apoptosis-related cysteine peptidase	Hs.654616	3.60E-06
<i>Cbx4</i>	Chromobox homolog 4 (Pc class homolog, <i>Drosophila</i>)	Hs.714363	0.00073
<i>Copa</i>	coatamer protein complex, subunit α	Hs.162121	0.00017
<i>Hspa8</i>	Heat shock 70kDa protein 8	Hs.702021	1.80E-05
<i>Iqgap1</i>	IQ motif containing GTPase activating protein 1	Hs.430551	0.0047
<i>Lca</i> <i>CLTA</i>	Clathrin, light chain A	Hs.522114	9.00E-07
<i>Mad111</i>	MAD1 mitotic arrest deficient-like 1 (yeast)	Hs.654838	7.30E-10
<i>Ppp3ca</i>	Protein phosphatase 3, catalytic subunit, α isozyme	Hs.435512	0.42
<i>SLC1A5</i>	Solute carrier family 1 (neutral amino acid transporter), member 5	Hs.631582	0.032
<i>Cox-2</i> <i>PTGS2</i>	Prostaglandin-endoperoxide synthase 2 (prostaglandin G/H synthase and cyclooxygenase)	Hs.196384	0.28
<i>TCRD</i> <i>TRD@</i>	T-cell receptor δ locus	Hs.74647	0.79
<i>FABP5</i>	Fatty acid binding protein 5 (psoriasis-associated)	Hs.408061	1
<i>DHRS2</i>	Dehydrogenase/reductase (SDR family) member 2	Hs.272499	0.19
<i>TGFB1</i>	Transforming growth factor, β 1	Hs.645227	0.63
<i>CTNNL1</i>	Catenin (cadherin-associated protein), α -like 1	Hs.58488	0.00059