

Supporting Information

Sheffer et al. 10.1073/pnas.0902232106

SI Methods

The data we work on was acquired within the framework of a Program Project Grant funded by the NCI, aimed at discovery of reliable molecular markers or targets for the detection, diagnosis, and treatment of colon cancer.

Expression Data. The expression data are composed of 390 Affymetrix GeneChip Human Genome U133A arrays. Seventeen repeats were removed from the analysis, and 2 outliers with high rates of absent values ($>14,500$ per sample), 9 normal colon samples whose expression profiles were markedly different from the other normals, 2 mislabeled samples, 1 microadenoma, 2 high grade adenomas, 30 noncolon normal samples and cell-lines, and 28 metastasis samples that exhibited high levels of stromal contamination (1). After all these filtering steps 299 samples were left: 180 primary carcinomas, 46 polyps, 43 normal colon epithelia, 21 liver metastases, and 9 lung metastases (see Table S3). The data were subjected to the following preprocessing steps: The Affymetrix MAS5 algorithm was applied on the .cel files (we tried also RMA and had reasons to work with MAS 5.0, see Fig. S3). Exact location for each gene was determined using UCSC hg18 known/RefSeq gene tables (2–4). The Affymetrix alignments of each probe set were compared with its chromosomal location. Probe sets (1,745) that did not have gene symbol/chromosomal location/alignments (including Affymetrix markers) or showed disagreement between their chromosomal location and alignments were removed (Affymetrix chromosomal locations were used for 432 probe sets whose genes did not appear in the University of California Santa Cruz (UCSC) known/RefSeq gene tables). Data were thresholded to $T = 10$ (probe sets with expression level less than T were assigned the value T). Probe sets (3,538) that had either no present calls or had expression value T in all samples, were removed. Probe sets (11,159) that represented unique genes were kept. When there were multiple probe sets per gene, the probe set that did not have $_s$, $_x$ in its identifier or had the highest expression levels was chosen. Data were subjected to \log_2 transformation.

SNP Data. The SNP data are composed of 154 SNP-based Affymetrix 50 k XbaI GeneChip Mapping Array (5). Seven mislabeled tumor samples were removed from the analysis, and 8 outlier samples (5 normal colon, 1 normal liver, and 2 primary carcinomas) and 9 cell lines, leaving 130 samples: 62 primary carcinomas, 43 normal colon epithelia, 4 normal liver, 3 normal lung, 8 liver metastases, and 10 lung metastases (see Table S3). The preprocessing steps taken for the SNP data were as follows: The data were normalized using an algorithm developed by L. Hertzberg and O. Zuk (unpublished data), providing copy number ratio values of each allele separately, for 57,768 SNPs, using the 21 normal colon samples of female subjects from our data set as a reference set. Both alleles were summed for each SNP, data were thresholded to 1 to avoid negative or small numbers (6) and \log_2 transformation was performed. *Log-copy number ratios* $CR_{n,s}$ were calculated by subtracting from the \log -transformed copy number of SNP n in sample s the \log -transformed copy number of SNP n in the matching normal sample. If there was no matching normal, the median \log -transformed copy numbers of SNP n in all normal samples were used. For chromosome X, all normal samples of the same gender were used. Smoothing was applied using the segmentation method gain and loss analysis of DNA (GLAD) (7). The $CR_{n,s}$ values were the input to GLAD and the output was a set of

segments per sample, where each segment was assigned a copy number, which was applied to all SNPs in that segment of the specific sample, producing new $CR_{n,s}$ values. Since GLAD tends to misidentify outliers as separate segments (6, 8), segments with fewer than 8 markers were joined to the neighboring segment with the closest copy number, assigning the new segment the weighted mean copy number. This step was performed recursively until there were no more such small segments.

Identification of CINons. Our analysis was based on genomic identification of significant targets in cancer (GISTIC) a statistical algorithm developed by Beroukhi et al. (10). The input to GISTIC were data from 55 aneuploid tumor samples (marked in red in Fig. 1, see text). The thresholds $\theta^{amp} = 0.0974$, $\theta^{del} = -0.1121$ were used, corresponding to 0.1% quantiles of the distribution of all $CR_{n,s}$ from the normal samples. G_i scores (for each SNP) were calculated for amplifications and deletions separately, and each score was assigned a FDR q -value (9) that assesses its statistical significance, taking into consideration its amplitude and frequency, over the null distribution where the $CR_{n,s}$ in each sample are randomly permuted. The q -values were generated as follows: the SNPs were sorted according to their p -values (p_i), and the q_i -values of the procedure were corrected, in descending order, according to $q_i = \min(\min(q_i, q_{i+1}), 1)$. When the contiguous set of deleted/amplified SNPs extended over the centromere, it was divided into 2 CINons, one on each chromosomal arm. This configuration identified all broad regions and a few focal regions, and is referred to as *configuration 1*. Driver mutations are presumed to be located in the peaks [of the $-\log(q\text{-value})$] measured across the CINons, i.e., at the SNPs having the lowest q -value). Since some of the broad regions exhibit fairly uniformly low q -values and no clear prominent peak, a different configuration, referred to as *configuration 2*, was used to identify the peak regions. In this configuration, SNPs located on GLAD segments that consist of at least 75% of the SNPs in a chromosomal arm of a specific sample, were assigned copy number zero which means that these SNPs, in this sample, are removed from the subsequent analysis. These modified $CR_{n,s}$ values were then used as an input to the second run of GISTIC, keeping the same thresholds as in configuration 1. Focal/peak regions were looked for, by taking the SNPs with the lowest q -values in each of the CINons found by this second run. In this way, it was possible to identify some new focal regions and some peak regions that belonged to the broad CINons identified earlier. The peak region in 7q31.1 was manually added to the list of peak CINons as it appeared significant by eye but had the second lowest q -values. A sample s was considered to have an aberrant CINon if $\text{median}(CR_{i,s}) < \theta^{del}$ or $\text{median}(CR_{i,s}) > \theta^{amp}$, for all SNPs i located within this CINon (shown in Table 1). Leave-one-out statistics was used originally in GISTIC to determine robust borders of a CINon; this part was not performed in our analysis since some of the peak CINons are based on only 1 or 2 samples, and by removing these samples we would have missed the peak. The same analysis was performed with GISTIC over all 80 tumor samples (by adding the 25 near-diploid samples), producing almost identical results.

Correlation Between Expression and Copy Number. Denote by K the group of all colon tissue samples (including normal colon, primary carcinoma, and metastasis to lung and liver) that had both measures of expression and copy numbers (79 samples, Table S3). For each probe set n , the Pearson correlation

coefficient was calculated between the $ER_{n,s}$ values of these 79 samples, and the corresponding values $MCR_{n,s}$ that represent the copy number ratios measured in sample s from SNPs located near the probe set n . The value of each of these $MCR_{n,s}$ was determined from the median of the $CR_{i,s}$ ratios of all of the SNPs i located within a window near probe set n . The window used extends from 1,000 bps before the transcription start position to 1,000 bps after the transcription end position of the gene represented by the probe set n . If no SNPs were found in this window, the mean of the first SNP before the transcription start position and the first SNP after transcription end position was taken. We calculated the correlation only for those probe sets for which both expression and copy number shifted (versus normals) in the same direction, i.e., $mean(ER_{n,s}) * mean(MCR_{n,s}) > 0$, where we average over the samples in K . Each correlation was assigned a P -value for testing the hypothesis of no correlation against the alternative, of a nonzero correlation. FDR of 25% was then used to filter the most significantly correlated probe sets.

CINon Expression Table. To assess whether a sample harbors amplification or deletion, we compared the $CE_{i,s}$ and the $CC_{i,s}$ (see *Methods*) to the upper and lower 0.1% quantiles of these values in the normal samples.

Annotation. Annotations of genes and noncoding small RNA were taken from UCSC tables, hg18 (2–4, 10, 11). Genes were considered to belong to a CINone if located in the interval between the first gene before and the first gene after the CINon's boundaries.

Outcome. Outcome labels were assigned according to the following rules: recurrence after more than 60 months was considered a good outcome, recurrence after less than or equal to 60 months was considered a poor outcome. In case there was no recurrence or recurrence information was not available, the follow-up status of the patient was tested; if the patient had died of the disease it was considered a poor outcome, otherwise, if the patient had a follow-up interval larger than 60 it was considered a good

outcome. Follow-up interval of less than or equal to 60 was considered an unknown outcome. For the Kaplan-Meier survival analysis, only primary tumors were taken.

Pathway Analysis. The list of probe sets that showed significant correlation between their expression levels and their copy numbers, that were also located within the broad CINons found earlier, including the focal CINon 1p, were analyzed using DAVID (12, 13) for enrichment of Biocarta and KEGG pathways. The background was chosen as Affymetrix HG-U133A. Pathways that passed 25% FDR according to DAVID were selected.

Mutation Status. Mutation status of *p53*, *APC*, *kRAS* were obtained and also methylation status for *APC* was collected. *APC* was considered as mutated if it was either mutated or methylated.

MIN Status. Microsatellite instability (MSI) status of a tumor was determined by the National Cancer Institute (NCI) according to a set of 5 microsatellite markers demonstrating instability (14). When 2 or more markers are positive the tumor is considered as MSI-High (MIN).

T-Test for Genes That Are Related to the Putative TSG and Oncogenes. Z-test was used to test the hypothesis that a value of gene i in sample j belongs to the distribution of the corresponding values of this gene in the normal samples. This was done to all of the genes in the expression data. FDR of 5% was used on the whole set of calculated p -values, resulting in a matrix of '1's, '0's, '-1's where '1' is placed for all upregulated values that passed FDR, '-1' for all downregulated values that passed FDR and '0' for the rest. This matrix represents the genes that were differentially expressed in each cancer sample versus normal tissues.

For each putative tumor suppressor or oncogenes, a t test was performed comparing the samples with '1' against samples with '0' (for *CCDC68* the samples with '-1' were compared against samples with '0'). FDR of 5% was used to select the genes that best separate the 2 groups of samples. Pathway enrichment analysis was performed using DAVID as described above.

1. Tsafir D, et al. (2005) Sorting points into neighborhoods (SPIN): data analysis and visualization by ordering distance matrices. *Bioinformatics* 21:2301–2308.
2. Karolchik D, et al. (2003) The UCSC Genome Browser Database. *Nucleic Acids Res* 31:51–54.
3. Hsu F, et al. (2006) The UCSC Known Genes. *Bioinformatics* 22:1036–1046.
4. Pruitt KD, Tatusova T, Maglott DR (2005) NCBI Reference Sequence (RefSeq): A curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res* 33:D501–504.
5. Bacolod MD, et al. (2008) The signatures of autozygosity among patients with colorectal cancer. *Cancer Res* 68:2610–2621.
6. Beroukhim R, et al. (2007) Assessing the significance of chromosomal aberrations in cancer: Methodology and application to glioma. *Proc Natl Acad Sci USA* 104:20007–20012.
7. Hupe P, et al. (2004) Analysis of array CGH data: From signal ratio to gain and loss of DNA regions. *Bioinformatics* 20:3413–3422.
8. Lai WR, Johnson MD, Kucherlapati R, Park PJ (2005) Comparative analysis of algorithms for identifying amplifications and deletions in array CGH data. *Bioinformatics* 21:3763–3770.
9. Benjamini Y, Hochberg Y (1995) Controlling the False Discovery Rate—a Practical and Powerful Approach to Multiple Testing. *J R Stat Soc Series B Stat Methodol* 57:289–300.
10. Griffiths-Jones S (2004) The microRNA Registry. *Nucleic Acids Res* 32:D109–111.
11. Weber MJ (2005) New human and mouse microRNA genes found by homology search. *FEBS J* 272:59–73.
12. Dennis G, Jr., et al. (2003) DAVID: Database for annotation, visualization, and integrated discovery. *Genome Biol* 4:P3.
13. Hosack DA, et al. (2003) Identifying biological themes within lists of genes with EASE. *Genome Biol* 4:R70.
14. Boland CR, et al. (1998) A National Cancer Institute Workshop on Microsatellite Instability for cancer detection and familial predisposition: Development of international criteria for the determination of microsatellite instability in colorectal cancer. *Cancer Res* 58:5248–5257.

SI References for CGH, SNP studies

1. Andersen CL, et al. (2007) Frequent occurrence of uniparental disomy in colorectal cancer. *Carcinogenesis* 28:38–48.
2. Bartos JD, et al. (2007) aCGH local copy number aberrations associated with overall copy number genomic instability in colorectal cancer: Coordinate involvement of the regions including BCR and ABL. *Mutat Res* 615:1–11.
3. Camps J, et al. (2006) Genome-wide differences between microsatellite stable and unstable colorectal tumors. *Carcinogenesis* 27:419–428.
4. Camps J, et al. (2008) Chromosomal breakpoints in primary colon cancer cluster at sites of structural variants in the genome. *Cancer Res* 68:1284–1295.
5. Cardoso J, Boer J, Morreau H, Fodde R (2007) Expression and genomic profiling of colorectal cancer. *Biochim Biophys Acta* 1775:103–137.
6. Diep CB, et al. (2004) Genome characteristics of primary carcinomas, local recurrences, carcinomatoses, and liver metastases from colorectal cancer patients. *Mol Cancer* 3:6.
7. Douglas EJ, et al. (2004) Array comparative genomic hybridization analysis of colorectal cancer cell lines and primary carcinomas. *Cancer Res* 64:4817–4825.
8. Grade M, et al. (2006) Aneuploidy-dependent massive deregulation of the cellular transcriptome and apparent divergence of the Wnt/beta-catenin signaling pathway in human rectal carcinomas. *Cancer Res* 66:267–282.
9. Habermann JK, et al. (2007) Stage-specific alterations of the genome, transcriptome, and proteome during colorectal carcinogenesis. *Genes Chromosomes Cancer* 46:10–26.
10. Jones AM, et al. (2007) Analysis of copy number changes suggests chromosomal instability in a minority of large colorectal adenomas. *J Pathol* 213:249–256.
11. Kim MY, et al. (2006) Recurrent genomic alterations with impact on survival in colorectal cancer identified by genome-wide array comparative genomic hybridization. *Gastroenterology* 131:1913–1924.
12. Lassmann S, et al. (2007) Array CGH identifies distinct DNA copy number profiles of oncogenes and tumor suppressor genes in chromosomal- and microsatellite-unstable sporadic colorectal carcinomas. *J Mol Med* 85:289–300.
13. Martin ES, et al. (2007) Common and distinct genomic events in sporadic colorectal cancer and diverse cancer types. *Cancer Res* 67:10736–10743.
14. Tanami H, et al. (2005) Involvement of cyclin D3 in liver metastasis of colorectal cancer, revealed by genome-wide copy-number analysis. *Lab Invest* 85:1118–1129.

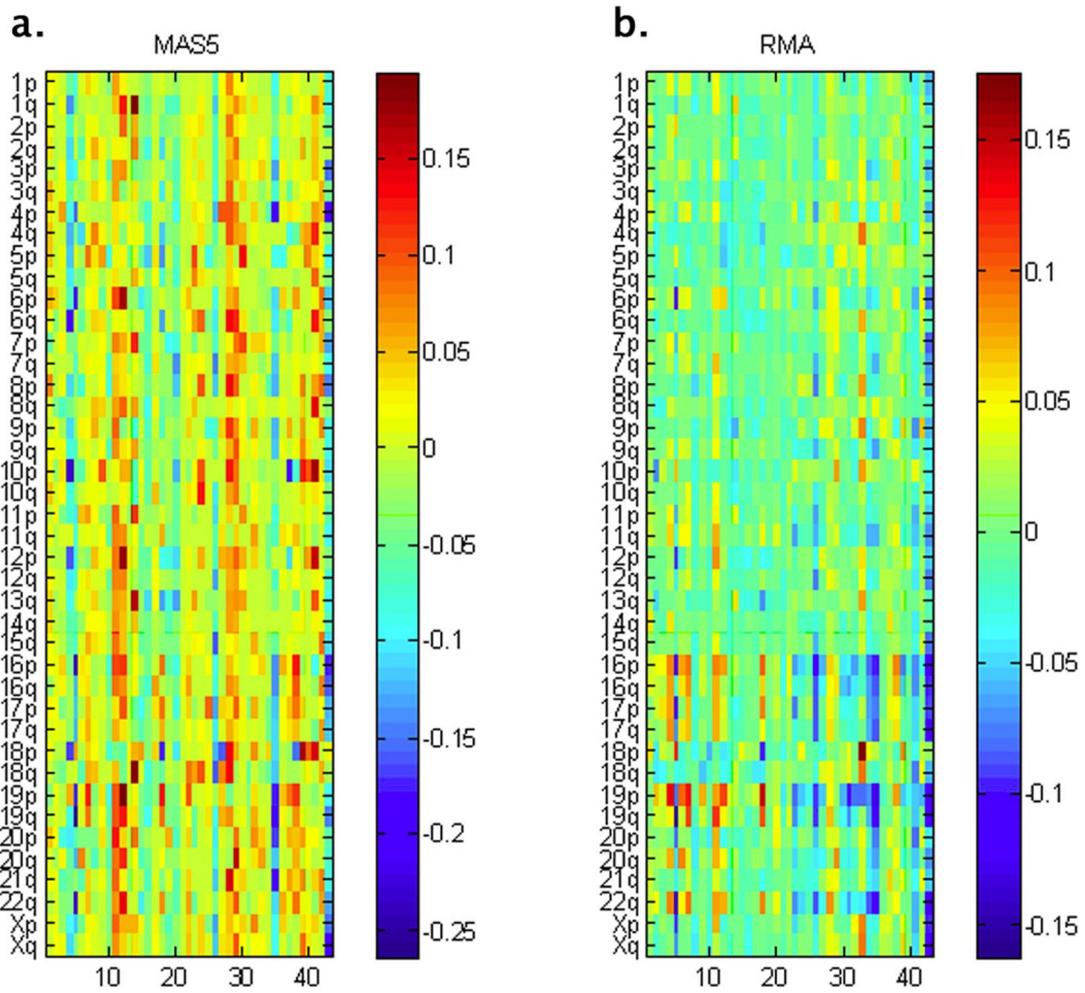


Fig. S3. MAS5 vs. RMA. The following matrices show the difference between MAS5 and RMA in the normal colon samples. The rows in each matrix represent chromosomal arms and the columns represent normal colon samples. Each entry ($M_{i,s}$) is calculated as the median of all probe sets $ER_{n,s}$ located on chromosomal arm i and sample s . The matrix in *a* represents the values after MAS5 normalization, as described in *SI Methods*. RMA values, as shown in *b* had the same preprocess stages, except they were not log₂ transformed, and the fold-change values were calculated compared to all normal samples without taking into consideration the batch effect of the different protocols that were used (see *SI Methods*). For some reason RMA shows greater difference between chromosomes 1–15 and 16–X, that is not shown by the MAS5 algorithm.

Table S1. Table of correlated genes that reside on focal/peak CINons

Id	Symbol	CINon	Chr	Band	Q-value
222258.s.at	SH3BP4	1	2	q37.2	6.43E-15
202142.at	COPS8	1	2	q37.3	2.26E-05
221575.at	SCLY	1	2	q37.3	1.65E-07
221548.s.at	ILKAP	1	2	q37.3	1.37E-02
218301.at	RNPEPL1	1	2	q37.3	5.43E-01
210264.at	GPR35	1	2	q37.3	3.67E-06
218106.s.at	MRPS10	2	6	p21.1	4.74E-06
217931.at	CNPY3	2	6	p21.1	1.34E-10
218061.at	MEA1	2	6	p21.1	2.36E-10
214383.x.at	KLHDC3	2	6	p21.1	1.75E-01
36084.at	CUL7	2	6	p21.1	1.57E-10
207011.s.at	PTK7	2	6	p21.1	2.95E-12
202401.s.at	SRF	2	6	p21.1	7.28E-01
213204.at	PARC	2	6	p21.1	5.77E-01
39817.s.at	C6orf108	2	6	p21.1	1.25E-07
220554.at	SLC22A7	2	6	p21.1	7.47E-02
213485.s.at	ABCC10	2	6	p21.1	3.10E-12
47608.at	TJAP1	2	6	p21.1	2.05E-04
209317.at	POLR1C	2	6	p21.1	9.72E-23
219380.x.at	POLH	2	6	p21.1	9.89E-01
221050.s.at	GTPBP2	2	6	p21.1	1.17E-01
218385.at	MRPS18A	2	6	p21.1	3.60E-01
210512.s.at	VEGFA	2	6	p21.1	2.70E-24
201802.at	SLC29A1	2	6	p21.1	2.09E-16
200064.at	HSP90AB1	2	6	p21.1	7.79E-34
203927.at	NFKBIE	2	6	p21.1	1.78E-11
209056.s.at	CDC5L	2	6	p21.1	1.39E-08
203510.at	MET	5	7	q31.2	4.50E-30
207614.s.at	CUL1	6	7	q36.1	2.59E-07
203358.s.at	EZH2	6	7	q36.1	6.54E-10
202033.s.at	RB1CC1	9	8	q11.23	7.38E-03
221504.s.at	ATP6V1H	9	8	q11.23	8.33E-06
216241.s.at	TCEA1	9	8	q11.23	7.62E-08
212449.s.at	LYPLA1	9	8	q11.23	4.56E-05
218027.at	MRPL15	9	8	q11.23	1.52E-02
219231.at	TGS1	9	8	q12.1	1.85E-18
202625.at	LYN	9	8	q12.1	2.40E-08
214003.x.at	RPS20	9	8	q12.1	1.57E-15
205372.at	PLAG1	9	8	q12.1	5.43E-01
218642.s.at	CHCHD7	9	8	q12.1	3.99E-01
218516.s.at	IMPAD1	9	8	q12.1	2.43E-01
222087.at	PVT1	10	8	q24.21	2.35E-19
219299.at	TRMT12	10	8	q24.13	2.88E-02
209510.at	RNF139	10	8	q24.13	5.00E-01
209218.at	SQLE	10	8	q24.13	8.30E-17
201985.at	KIAA0196	10	8	q24.13	2.11E-02
202241.at	TRIB1	10	8	q24.13	7.91E-01
214532.x.at	POU5F1	10	8	q24.21	1.51E-06
202431.s.at	MYC	10	8	q24.21	1.21E-25
219099.at	C12orf5	11	12	p13.32	1.96E-06
204146.at	RAD51AP1	11	12	p13.32	6.32E-10
218258.at	POLR1D	13	13	q12.2	5.10E-19
204328.at	TMC6	16	17	q25.3	5.35E-06
201079.at	SYNGR2	16	17	q25.3	2.55E-08
202338.at	TK1	16	17	q25.3	2.45E-04
219394.at	PGS1	16	17	q25.3	5.00E-01
220370.s.at	USP36	16	17	q25.3	3.62E-03
206724.at	CBX4	16	17	q25.3	5.48E-26
222116.s.at	TBC1D16	16	17	q25.3	3.10E-17
217796.s.at	NPLOC4	16	17	q25.3	2.21E-02
210428.s.at	HGS	16	17	q25.3	2.13E-09
200654.at	P4HB	16	17	q25.3	6.59E-06
211716.x.at	ARHGDI1A	16	17	q25.3	9.43E-01
204970.s.at	MAFG	16	17	q25.3	1.95E-03
202148.s.at	PYCR1	16	17	q25.3	3.53E-24

Id	Symbol	CINon	Chr	Band	Q-value
218908.at	ASPCR1	16	17	q25.3	2.93E-12
212968.at	RFNG	16	17	q25.3	1.42E-03
217782.s.at	GP51	16	17	q25.3	1.38E-04
64438.at	FLJ22222	16	17	q25.3	3.75E-01
218130.at	C17orf62	16	17	q25.3	8.34E-02
219862.s.at	NARF	16	17	q25.3	7.31E-11
203064.s.at	FO XK2	16	17	q25.3	8.81E-07
209076.s.at	WDR45L	16	17	q25.3	7.22E-12
208804.s.at	SFRS6	19	20	q13.11	3.95E-01
213837.at	L3MBTL	19	20	q13.11	5.88E-01
218709.s.at	IFT52	19	20	q13.12	1.84E-06
201710.at	MYBL2	19	20	q13.12	4.78E-19
209020.at	C20orf111	19	20	q13.12	4.26E-09
208429.x.at	HNF4A	19	20	q13.12	8.59E-01
219633.at	C20orf121	19	20	q13.12	1.77E-03
221471.at	SERINC3	19	20	q13.12	2.87E-05
217718.s.at	YWHAB	19	20	q13.12	6.00E-07
201870.at	TOMM34	19	20	q13.12	7.63E-26
205411.at	STK4	19	20	q13.12	2.38E-03
202071.at	SDC4	19	20	q13.12	1.30E-08
217770.at	PIGT	19	20	q13.12	1.22E-04
202954.at	UBE2C	19	20	q13.12	4.42E-30
205388.at	TNNC2	19	20	q13.12	3.52E-01
217592.at	ZSWIM1	19	20	q13.12	7.28E-01
202075.s.at	PLTP	19	20	q13.12	7.38E-03
89948.at	C20orf67	19	20	q13.12	4.16E-01
78330.at	ZNF335	19	20	q13.12	2.69E-02
203936.s.at	MMP9	19	20	q13.12	5.36E-14
219447.s.at	SLC35C2	19	20	q13.12	1.36E-11
55692.at	ELMO2	19	20	q13.12	2.98E-08
217875.s.at	PMEPA1	20	20	q13.31	9.22E-18
204092.s.at	AURKA	20	20	q13.2	9.17E-23
202190.at	CSTF1	20	20	q13.31	6.30E-09
217737.x.at	C20orf43	20	20	q13.31	4.62E-09
209590.at	BMP7	20	20	q13.31	2.31E-05
201558.at	RAE1	20	20	q13.31	3.07E-26
213405.at	RAB22A	20	20	q13.32	1.17E-04
202549.at	VAPB	20	20	q13.32	8.96E-03
221500.s.at	STX16	20	20	q13.32	2.49E-13
89476.r.at	NPEPL1	20	20	q13.32	3.95E-03
220607.x.at	TH1L	20	20	q13.32	2.12E-22
217801.at	ATP5E	20	20	q13.32	2.17E-07
217851.s.at	SLMO2	20	20	q13.32	2.27E-01
204554.at	PPP1R3D	20	20	q13.33	1.39E-04
211038.s.at	CROCCL1	24	1	p36.13	3.36E-03
201155.s.at	MFN2	24	1	p36.22	3.43E-16
212326.at	VPS13D	24	1	p36.22	7.79E-16
202481.at	DHRS3	24	1	p36.22	7.65E-03
217992.s.at	EFHD2	24	1	p36.21	1.62E-01
212146.at	PLEKHM2	24	1	p36.21	5.00E-01
218934.s.at	HSPB7	24	1	p36.13	6.53E-12
221813.at	FBXO42	24	1	p36.13	2.37E-01
202675.at	SDHB	24	1	p36.13	2.03E-10
209791.at	PADI2	24	1	p36.13	1.98E-31
221656.s.at	ARHGEF10L	24	1	p36.13	1.31E-01
212394.at	KIAA0090	24	1	p36.13	1.25E-01
216381.x.at	AKR7A3	24	1	p36.13	2.07E-11
202139.at	AKR7A2	24	1	p36.13	2.25E-13
37012.at	CAPZB	24	1	p36.13	7.53E-04
37005.at	NBL1	24	1	p36.13	6.80E-02
203649.s.at	PLA2G2A	24	1	p36.13	9.02E-06
218309.at	CAMK2N1	24	1	p36.12	3.30E-08
218246.at	C1orf166	24	1	p36.12	3.39E-10
209018.s.at	PINK1	24	1	p36.12	4.05E-11
201935.s.at	EIF4G3	24	1	p36.12	2.37E-10
201749.at	ECE1	24	1	p36.12	1.20E-05

Id	Symbol	CINon	Chr	Band	Q-value
203911_at	RAP1GAP	24	1	p36.12	1.30E-07
214230_at	CDC42	24	1	p36.12	1.15E-06
219103_at	DDEF1	24	1	p36.12	6.68E-09
202292_x.at	LYPLA2	24	1	p36.11	5.06E-04
202528_at	GALE	24	1	p36.11	3.63E-01
202772_at	HMGCL	24	1	p36.11	6.10E-20
202838_at	FUCA1	24	1	p36.11	1.44E-37
217779_s.at	PNRC2	24	1	p36.11	1.44E-08
202553_s.at	SXF2	24	1	p36.11	3.35E-08
209007_s.at	C1orf63	24	1	p36.11	5.43E-01
217766_s.at	TMEM50A	24	1	p36.11	7.31E-11
57082_at	LDLRAP1	24	1	p36.11	3.42E-01
221269_s.at	SH3BGR13	24	1	p36.11	5.96E-05
218547_at	DHDDS	24	1	p36.11	8.68E-30
208668_x.at	HMG2	24	1	p36.11	8.38E-01
203379_at	RPS6KA1	24	1	p36.11	9.44E-15
212152_x.at	ARID1A	24	1	p36.11	9.03E-04
218799_at	ATPBD1B	24	1	p36.11	5.44E-03
209453_at	SLC9A1	24	1	p36.11	3.66E-18
219278_at	MAP3K6	24	1	p36.11	1.01E-01
212111_at	STX12	24	1	p35.3	1.34E-17
201756_at	RPA2	24	1	p35.3	9.92E-01
205309_at	SMPDL3B	24	1	p35.3	3.73E-09
218671_s.at	ATPIF1	24	1	p35.3	2.91E-06
219235_s.at	PHACTR4	24	1	p35.3	1.72E-12
218977_s.at	TRSPAP1	24	1	p35.3	1.21E-01
201696_at	SFRS4	24	1	p35.3	2.09E-10
202898_at	SDC3	24	1	p35.2	9.02E-03
204054_at	PTEN	29	10	q23.31	7.15E-07
211285_s.at	UBE3A	32	15	q11.2	6.24E-01
202604_x.at	ADAM10	33	15	q22.1	3.65E-01
217828_at	SLTM	33	15	q22.1	9.14E-10
218761_at	RNF111	33	15	q22.1	1.76E-01
209120_at	NR2F2	34	15	q26.2	3.39E-03
220180_at	CCDC68	39	18	q21.2	1.58E-30
218145_at	TRIB3	41	20	p13	2.83E-35
221827_at	RBCK1	41	20	p13	1.25E-17
212073_at	CSNK2A1	41	20	p13	1.85E-01
201052_s.at	PSMF1	41	20	p13	3.55E-06
219958_at	C20orf46	41	20	p13	8.49E-08
202897_at	SIRPA	41	20	p13	2.62E-02
208821_at	SNRPB	41	20	p13	4.15E-13
200875_s.at	NOL5A	41	20	p13	1.48E-24
203459_s.at	VPS16	41	20	p13	2.20E-03
213795_s.at	PTPRA	41	20	p13	1.73E-04
215544_s.at	UBOX5	41	20	p13	8.77E-08
204447_at	ProSAPiP1	41	20	p13	6.43E-15
218159_at	C20orf116	41	20	p13	3.44E-19
209171_at	ITPA	41	20	p13	1.25E-11
50314_i.at	C20orf27	41	20	p13	7.79E-16
212437_at	CENPB	41	20	p13	1.92E-01
201853_s.at	CDC25B	41	20	p13	1.36E-28
218809_at	PANK2	41	20	p13	2.14E-04
204668_at	RNF24	41	20	p13	2.97E-03
210357_s.at	SMOX	41	20	p13	4.23E-18
219570_at	C20orf23	43	20	p12.1	3.63E-01
217792_at	SNX5	44	20	p11.23	4.39E-01
37254_at	ZNF133	44	20	p11.23	1.64E-06
219951_s.at	C20orf12	44	20	p11.23	2.58E-02
205218_at	POLR3F	44	20	p11.23	2.35E-01
201582_at	SEC23B	44	20	p11.23	1.21E-01

The q-value for each probe set calculated from the *t*-test between primary tumors and normal samples.

Table S2. Expression levels of each putative TSG and oncogene separate the samples into 2 groups

	Pathway	Number of genes in the pathway
PMEPA1	Oxidative phosphorylation	26
	Valine leucine and isoleucine degradation	12
	Citrate cycle (TCA cycle)	8
	Glyoxylate and dicarboxylate metabolism	5
	Glutathione metabolism	8
	TACI and BCMA stimulation of B cell immune responses	5
	Electron Transport Reaction in Mitochondria	4
	Fatty acid metabolism	9
POLR1D	Leukocyte transendothelial migration	35
	B cell receptor signaling pathway	21
	Toll-like receptor signaling pathway	28
	Natural killer cell mediated cytotoxicity	32
	Ribosome	23
	T cell receptor signaling pathway	23
	Hematopoietic cell lineage	21
	Cytokine–cytokine receptor interaction	47
	Chondroitin sulfate biosynthesis	7
	Cell adhesion molecules (CAMs)	26
Fc epsilon RI signaling pathway	18	
CCDC68	Long-term depression	12
	Aminoacyl-tRNA biosynthesis	7
	p53 signaling pathway	10

The table details the pathway analysis for the genes that differentiate between the 2 groups of samples (see *SI Methods*). The separations into 2 groups was as follows: *PMEPA1*: 162 samples had overexpression of this gene and 48 samples had normal-like expression, 973 genes passed the *t*-test (5% FDR) between these 2 groups; *POLR1D*: 127 samples had overexpression and 83 samples were normal-like expression, 1,602 genes passed 5% FDR; *CCDC68*: 179 samples were underexpressed and 31 samples were normal-like expression, 802 passed 5% FDR.

Table S3. Summary of the samples used in the copy number and expression analysis

Tissue	Number of samples	Number of samples in expression	Number of samples in SNP			Number of samples in Both expression and SNP		
			All	Near-euploid	Aneuploid	All	Near-euploid	Aneuploid
Polyp	46	46	—	—	—	—	—	—
Normal colon	71	43	43	—	—	15	—	—
Normal liver	4	—	4	—	—	—	—	—
Normal lung	3	—	3	—	—	—	—	—
Primary tumor	187	180	62	22	40	55	18	37
Liver metastasis	24	21	8	1	7	5	1	4
Lung metastasis	15	9	10	2	8	4	—	4
Total	350	299	130	25	55	79	19	45
Primary tumor—stage I	28	28	11	4	7	11	4	7
Primary tumor—stage II	48	47	18	7	11	17	6	11
Primary tumor—stage III	50	49	14	4	10	13	3	10
Primary tumor—stage IV	61	56	19	7	12	14	5	9

