# Methylome analysis reveals dysregulated developmental and viral pathways in breast cancer

Mohammad OE Abdallah[1], Ubai K Algizouli[1], Maram Abbas Suliman[1], Rawya Abdulaziz Abdulrahman[1], Mahmoud Koko[1], Ghimja fessahaye, Jamal Haleem Shakir[2], Ahmed H. Fahal[3]  Ahmed M Elhassan[1], Muntaser E Ibrahim[1]  and  Hiba S Mohamed [1]*


* Correspondence: hibasalah@iend.org

Institute of Endemic Disease, University of Khartoum, P. O. Box 102, Khartoum, Sudan

**Abbreviations**

MOA (melsiddieg@gmail.com); UKA (ubai@hotmail.com)

MAS (maramo_83@hotmail.com); RAA (rawya83@hotmail.com)

 MK (mahmoud.altayeb@gmail.com); GS (ghimjaf@yahoo.com)

JHS (jamalhaleem@yahoo.com);AHF (ahfahal@hotmail.com)

AME (ahmedelhassan85@gmail.com);MEI (mibrahim@iend.org);

HSM (hibasalah@iend.org)

## Abstract

### Background

Breast cancer (BC) ranks among the most common cancers in Sudan and worldwide with hefty toll on female health and human resources. Recent studies have uncovered a common BC signature characterized by low frequency of oncogenic mutations and high frequency of epigenetic silencing of major BC tumor suppressor genes. Therefore, we conducted a genome-wide methylome study to characterize aberrant DNA methylation in breast cancer.

### Results

Differential methylation analysis between primary tumor samples and normal samples from healthy adjacent tissues yielded 20188 differentially methylated positions (DMPs), which is further divided into 13633 hypermethylated sites corresponding to 5339 genes and 6555 hypomethylated sites corresponding to 2811 genes. Moreover, bioinformatics analysis revealed epigenetic dysregulation of major developmental pathways including hippo signaling pathway.  We also uncovered many clues to a possible role for EBV infection in BC

### Conclusion

Our results clearly show the utility of epigenetic assays in interrogating breast cancer tumorigenesis, and pinpointing specific developmental and viral pathways dysregulation that might serve as potential biomarkers or targets for therapeutic interventions.

**Keywords**

Methylome, Breast Cancer, Epigenetics, DNA Methylation, HM450, Epigenome Reference.

**Background**

Breast cancer (BC) is the most common cancer among females in Sudan [1–3], and is still a leading cause of high morbidity and mortality across the world. According to a recent report from the national cancer registry[2], BC had an incidence rate of 25.1 per 100.000, more than twice the incidence rate of the second commonest cancer. Furthermore, Sudanese BC patients tend to present at young age, at late stage, and with advanced disease compared to their counterparts in other countries [4]. Another study [5] reported a young age of presentation for locally advanced BC. Therefore, there is an urgent need for serious epidemiologic and molecular studies in order to trace the underlying mechanisms behind BC, and for developing better early detection methods as well as a nationwide educational effort to tackle this ravaging disease.

Epigenetics has emerged as a new, rapidly growing field in biology, with significant implications for cancer research. Epigenetic modifications include DNA methylation, and histone modifications, although they both do not alter DNA sequence per se, they influence chromatin remodeling and thus offer a dynamic and flexible way of controlling gene expression.

DNA methylation of cytosine residues occurs predominantly at CpG sites, and is mediated by three DNA methyltransferases (DNMTs). DNMT1, which maintains DNA

methylation during cell replication, and a pair of DNMT3s - DNMT3a and DNMT3b - which is responsible for de novo DNA methylation. Epigenetic reprogramming through genome-wide alteration of DNA methylation (methylome) is critical for control of development and differentiation in normal cells and tissues, however, faulty epigenetic reprogramming, as in aberrant DNA methylation, can be a major driver of multiple types of cancer including BC [6, 7]

Methylome analysis has proved to be very pertinent to the study of the different aspects of cancer tumorigenesis. The vast majority of methylation changes occurs in a tissue-specific manner [8], which makes methylome profiling a very sensitive and specific method for delineating dysregulated epigenetic pathways at the tissue level, as in cancer, which usually arises from a single tissue. Moreover, DNA methylation is a stable epigenetic mark that is ideal for development of biomarker assays, which can offer a rapid, cost effective, a and minimally invasive diagnostic/prognostic tests [9, 10]. Additionally, methylome analysis has been effectively used in tumor subtype classification [11–15]. Furthermore, genome-wide methylome assays have also proved to be very useful in detecting and profiling viral epigenetic signature in cancer [16–18].

The aim of the present study was to investigate genome-wide DNA methylation profile of breast cancer in Sudanese patients utilizing Illumina Infinium HumanMethylation450 BeadChips (HM450) methylation assay. This array-based assay is widely used in epigenetics studies, and is a reliable, cost effective, high throughput method. We conducted methylome analysis comparing primary BC tissue samples against normal samples from adjacent healthy tissues. The results of this study provide a valuable insight into the epigenetics of BC in Sudanese patients.

**Results**

Genome-wide DNA Differential methylation Analysis

Each of three approaches - listed in Materials and Methods- produced a list of differentially methylated sites: Limma, 39940; Wilcoxon, 34099; Nimbl, 22251 (0.2 median beta value difference, Benjamini-Hochberg adjusted p-value ≤0.05). Here we only report the results for final set obtained from Nimbl-compare module, which represents the intersection of the three methods. The final set consisted of 20188 differentially methylated CpG sites, which is further divided into 13633 hypermethylated sites corresponding to 5339 genes and 6555 hypomethylated sites corresponding to 2811 genes. Nimbl unique approach ensured detection of differentially methylated positions (DMPs) that have the largest effect size as illustrated in **Fig 1.** A volcano plot showing the demarcation of differentially methylated sites by both statistical significance and effect size is shown in **Fig 1.B** Hierarchical clustering of the top 250 differentially methylated sites sorted by F value (low intragroup variability and higher intergroup variability) is shown in **Fig 2**. The resulting heatmap and dendogram showed clear separation of tumor samples from normal samples

Genomic Distribution of Differentially Methylated CpG sites

Differentially hypermethylated and hypomethylated sites displayed similar distribution with regard to gene elements as defined by HM450 –TSS1500, TSS200, First Exon, gene body, and 3UTR – **Fig 3.A**. However, they showed an asymmetric distribution with regard to CpG island relation with most of the hypermethylated sites mapping to CpG islands, whereas most of the hypomethylated sites mapped to open sea **Fig 3.B.**

Of the 13633 hypermethylated sites, 24.37% (N=3323) mapped to Dnase hypersensitive sites compared with only 8.67% (N =568) of hypomethylated sites. Interestingly, while a greater percentage of hypermethylated compared to hypomethylated sites overlapped differentially methylated regions (DMR), [54.83% (N=1612), 11.47% (N=46)] respectively, hypomethylated sites were more concentrated in cancer DMR (CDMR), with 49.63%  compared with 14.66% in hypermethylated sites, hypomethylated sites were more concentrated in cancer DMR (CDMR), with 49.63% compared with 14.66% in hypermethylated sites. The genomic distribution of hypermethylation and hypomethylation sites at each chromosome is shown in Additional File 1 and Additional File 2.

## Comparison to Reference Epigenome

We utilized data from the recently released Human epigenome reference data [19] to annotate the set of deferentially methylated CpG sites. We mapped hyper and hypo DMPs in the promoter region from our data against two reference epigenome breast cell lines: HMEC(Human mammary primary epithelial cells), and vHMEC (Human mammary primary epithelial cell variant)[20, 21]].We examined the change in chromatin states - from the 15-chromatin states model [19] - that accompany the acquisition or loss of DNA methylation in the context of transitioning from normal to tumor states. Our results revealed a noticeable gain of repressive marks for the hypermethylated DMPs, which increased from 55.5% in HMEC cells to 78.7% in vHMEC cells. Interestingly, we also found a slight increase in the percentage of repressive marks in the hypomethylated DMPs, which increased from 54.3% to 61.6%. Notably, in both cases, most of the upsurge in repressive regions were concentrated in Polycomb-repressed regions **Fig**

**3.C, and 3.D.**

In addition, we observed a marked drop of all active chromatin states except for weak transcription and distal enhancer activity between the HMEC and vHMEC cells for the hypermethylated group. On the other hand, the hypomethylated group showed multiple notable shifts: From quiescent to Polycomb repression, from weak transcription to strong transcription, and from distal enhancers to genic enhancer (intronic enhancers).

## Candidate biomarkers discovery

Nimbl method was used for detection and prioritization of candidate biomarkers with greatest inter-group variability, and lowest intra-group variability[22]. Using this approach, we were able to identify a number of new as well as previously well-known BC biomarkers. Among the genes that showed significant promoter hypermethylation, we identified PAX6 [23, 24], WT1 [25], SOX1 [26], and TP73 [27, 28], all of them have been previously associated with BC. We also identified a set of previously uncharacterized biomarkers like PCDHGA1, HOXC4, and TBX15. To validate our candidate genes we interrogated our candidate gene list against BC methylome data from the Cancer Genome Atlas Network: http://cancergenome.nih.gov/ as compiled by MethHC [29] web portal. All the genes from our data were also significantly hypermethylated in the TCGA dataset. **Fig** 4 shows promoter hypermethylation of the TP73 gene.

## Pathway and Network Analysis

Results from the ReactomeFI for the EDG network uncovered a massive network of 1310 nodes (genes) and 5097 edges (interactions), while the EUG list produced a smaller network of 763 nodes and 2265 edges. Furthermore, loading the NCI (National

Cancer Institute) cancer gene index identified 781, and 470, neoplasia related genes from the EDG, and EUG networks respectively, of which 332 EDG genes, and 222 EUG genes were associated with breast cancer in the cancer gene index.

Pathway enrichment analysis on the EUG network. Identified hippo signaling, Wnt signaling, and many extracellular matrix and metastasis promoting pathways as summarized in **Table 1**. Performing the pathway enrichment analysis on the breast cancer EUG subnetwork also identified hippo signaling and pathways of extracellular matrix in addition to pathways involved in immune response against viruses **Table 2**. Interestingly, breast cancer subnetwork showed significant enrichment for Epstein-Barr virus infection (FDR <0.001).

Pathway analysis on the EDG network identified Neuroactive ligand-receptor interactions, G-protein signaling, RAP1 signaling, RAS signaling, Potassium channel signaling and many other pathways as summarized in **Table 3.** While the smaller EDG breast cancer subnetwork showed significant enrichment for a multitude of pathways including all the pathways that were enriched in the EDG network in addition to many cancer related and immune response pathways**.** Interestingly, the EDG sub network was also significant for direct p53 effectors. The complete list of enriched pathways for the EDG breast cancer subnetwork is shown in **Additional file 3: Table S1.**

 **Discussion**

Reference Human Epigenome Annotations

The recent release of the human reference epigenome data by the Roadmap project ushered in a new era of epigenomics. The current study utilized this new wealth of

information to interpret epigenetic data. We successfully mapped hyper and hypo DMPs to chromatin states from normal and premalignant reference breast cells (HMEC and vHMEC respectively). Despite the fact that vHMEC is a premalignant and not a primary tumor cell, we argue that vHMEC is a suitable model for the epigenetic changes that accompany BC tumorigenesis because the vast majority of epigenetic changes tend to occur early during BC tumorigenesis [30–33].

Notably, our data revealed a strong Polycomb repression in both hypermethylated and hypomethylated CpG sites. These findings, are in accordance with the emerging evidence that DMPs are enriched for Polycomb repression in primary breast tumors [34] and triple negative BC [35]. Moreover, Various elements of the Polycomb repressive complexes are well known to be overexpressed in BC [36, 37] and are required for stem cell state in mammary tumors [38, 39]. Interestingly, Reyngold et al, found that unlike primary tumors, genes methylated in metastatic lesions seem to lack Polycomb repressive marks [40].

## Network-Based pathway enrichment analysis

Our network-based pathway enrichment results for the EUG network revealed many upregulated pathways that have been previously associated with BC tumorigenesis. Hippo signaling, which appeared as the top significantly enriched pathway in our results, has recently emerged as an important regulator of BC growth, migration, invasiveness, stemness, as well as drug resistance[41]. Wang et al, demonstrated that overexpression of YAP enhanced BC formation and growth. Hiemer et al, found that both TAZ and YAP -key effectors of the Hippo pathway -are crucial to promote and maintain TGFβ-induced tumorigenic phenotypes in breast cancer cells [42]. In addition,

YAP was demonstrated to mediate drug resistance to RAF and MEK targeted cancer therapy [43, 44]. Interestingly, we also reported an upregulated Wnt signaling pathway, which has been linked to BC growth and malignant behavior [45]. Jinhua et al found that Wnt signaling pathway is required for triple- negative breast cancer development [46]. Recent studies have suggested long lasting reduced Wnt signaling as the mechanism by which early pregnancy protects against BC [47].

Regarding the EDG network, Neuroactive ligand-receptor interaction, in addition to GPCR, RAS and Rap1 signaling were among the most significantly enriched pathways. Recent studies have found Neuroactive ligand-receptor interaction related genes to be hypermethylated in colorectal and EBV associated gastric cancers [20,21]. Elements of RAS signaling like RASSF has been frequently found to be hypermethylated in BC [49], moreover, Qin et al, has demonstrated that resveratrol is able to demethylate RASSF1 promoter through decreased DNMT1 and DNMT3b in mammary tumors [50, 51]. Notably, we reported the apparent silencing of multiple pro-tumor pathways in our results like GPCR and RAP1 signaling, the precise significance of this findings remains unclear. In addition, we also noticed the bivalent enrichment of multiple pathways (where different elements of the same pathway are both up and down regulated). Interpreting such perturbations is tricky, and predicting the net outcome of those perturbations might not be readily obvious given the crosstalk between different pathways.

## EBV signature

We previously reported a strong association between EBV and BC in Sudanese patients [52] , we also reported frequent epigenetic silencing of major tumor suppressor

genes coupled with low frequency of known tumor associated mutations in the same population [53]. In this study, we have demonstrated genome-wide epigenetic alterations consistent with our original proposition that epigenetic changes are the primary driver of BC tumorigenesis in Sudanese patients.

A myriad of recent studies point toward a common theme in EBV associated cancers characterized by genome-wide epigenetic changes coupled with a paucity of mutations. EBV infection is now known to play significant role in epithelial cancers like nasopharyngeal and gastric carcinomas mainly through genome-wide epigenetic changes [54–56]. Li et al observed a unique epiphenotype of EBV associated carcinomas suggesting a predominant role for EBV infection in the ensuing epigenetic dysregulation of those cancers [17]. Another study attributed the genome-wide promoter methylation in EBV driven gastric cancer to the induced expression of DNA methyltransferase-3b (DNMT3b) [57].

Our data mirrored the overall unique pattern of EBV infection characterized by sweeping epigenetic changes accompanied by low mutation frequency. Furthermore, we also showed that the EUG network was significantly enriched for EBV infection pathway **Fig 5**. In addition, results from MSig perturbations obtained from GREAT web tool (which predicts functions of cis regulatory elements) [58], showed significant enrichment for a set of downregulated genes which had been previously correlated with increased expression of EBV EBNA1 protein in NPC, in the hypermethylated CpG sites group, data not shown.. For the hypomethylated CpG group, we found genes upregulated in B2264-19/3 cells (primary B lymphocytes) within 30-60 min after activation of LMP1 to be significantly enriched in MSig oncogenic signature. These findings taken together

provide the first bioinformatics evidence of a possible active role for EBV infection in BC tumorigenesis in Sudanese patients.

## Conclusions

In conclusion, our study uncovered interesting epigenetic patterns characterized by increased acquisition of Polycomb repressive marks, as revealed from comparison to human reference epigenome breast cells. We identified many potential BC biomarkers like TP73, and TBX15. We also identified many significantly enriched developmental pathways including Hippo and Wnt signaling pathways. Additionally, our bioinformatics analysis indicated a possible role for EBV infection in BC tumorigenesis.

## Materials and Methods

### Ethical considerations

Ethical approval for this study was obtained from the Institute of Endemic Diseases, University of Khartoum Ethical Committee. Written informed consent was obtained from all participants; all clinical investigations were conducted according to the principles expressed in the Declaration of Helsinki:

http://www.wma.net/en/30publications/10policies/b3/index.html

### Samples, DNA extraction, and Illumina Infinium HumanMethylation 450 (HM450) array

Genomic DNA was extracted from eight samples of primary breast tumors and eight normal samples from adjacent healthy. All samples were independently reviewed by histopathologists. DNA was extracted from tissues using Promega genomic DNA purification kit [59] following the standard protocol as described by the manufacturer.

DNA methylome profiling was performed using Illumina Infinium HumanMethylation 450 (HM450) [60]BeadChip array by Beijing Genomics Institute (BGI)

## Data Preprocessing

For quality control, any array probes with p detection value less than 0.05 or missing beta values were removed. In addition, array sites corresponding to sex chromosomes or mapping to SNPs were filtered out. Peak-based correction[61] (PBC) was used to normalize the final dataset and to correct for probe type bias. Density plots of beta values for individual samples are shown in Additional File 4**: Fig S3**

## Genome-wide DNA Differential methylation Analysis

A trilateral approach consisting of two statistical methods augmented by one numerical method was used for the differential methylation analysis: Moderated T test from R limma [62] package ; Wilcoxon test (Non Parametric test) from R stat package; and Nimbl [22] (Numerical Identification of Methylation Biomarker Lists) which is a Matalab package designed to identify and prioritize differentially methylated sites.

Nimbl core module identify potential biomarkers by calculating a score based on the inter-group and intra-group variability:

Score = beta_valdist – (mediandiff – beta_valdist)

Where beta_valdist is the distance in beta values between non-overlapping groups and mediandiff is the absolute difference of the medians of each group [22]. It then assigns

high scores for CpG sites that achieve higher discrimination between groups while maintaining low within-group variability. Nimbl-compare module was also used to extract the final set of CpG sites that were identified by all three methods. Hierarchical clustering analysis was performed using the top 250 differentially methylated sites sorted by F value.

## Reference Epigenome Annotations

Bed files of chromatin states for both HMEC and vHMEC cells were obtained from Roadmap web portal: http://egg2.wustl.edu/roadmap/web_portal/, further analysis was performed in GALAXY web-based platform [63–65] and R statistical software.

## Network and Pathway Analysis

Differential methylation analysis produced two lists of differentially methylated genes (hyper and hypo) and their enrichment of differentially methylated sites in their gene regions, i.e., promoter region, gene body, 3UTR, etc. The aggregated gene list was sorted by the count of methylated sites in the promoter area, first exon, and gene body regions. Subsequently all epigenetically upregulated genes (EUG) were combined in a single group, i.e., genes bearing methylation marks that promote gene expression - hypomethylation in the promoter area, and first exon or hypermethylation in the gene body region- in a single group. Then we compiled a second group of epigenetically downregulated genes (EDG), i.e., genes bearing methylation marks that inhibits gene expression, i.e., hypermethylation of the promoter area, and the first exon or hypomethylation of the gene body region. We excluded other gene-based regions that are not well correlated with gene expression from further analysis.

We utilized ReactomeFI[66], a Cytoscape [67] app to perform network and pathways analysis. Projecting the lists of EDG and EUG groups through the ReactomeFI functional network produced two corresponding networks. To extract breast cancer specific subnetworks from EUG and EDG groups we loaded NCI cancer index from within the ReactomeFI app, and we selected nodes that corresponded to malignant breast cancer.

## Abbreviations

BC, breast cancer; DMP, differentially methylated position; DMR, differentially methylated region; CDMR, cancer differentially methylated site; TSS, transcription start site; UTR, untranslated region; MSig, mutation signature

## Competing interests

Authors declare no competing interests.

## Authors' contributions

HSM conceived and design the study and contributed to manuscript writing and data interpretation. MEI contributed to study design and manuscript writing. MOA performed the data analysis, contributed to interpretation and prepared the manuscript draft. JHS and AHF recruited patients and provided samples. MK contributed to data analysis. AME performed the histopathology, UKA, MAS, RA, and GS contributed to sample collection, DNA extraction and purification. All authors read and approved the final manuscript

**Acknowledgments**

**Funding**

# Authors' information

[1]Institute of Endemic Disease, University of Khartoum, P. O. Box 102. Khartoum Sudan

[2] Khartoum teaching hospital. Khartoum. Sudan

[3]Faculty of Medicine. University of Khartoum. Khartoum. Sudan

# References

1. Elamin A, Ibrahim ME, Abuidris D, Mohamed KEH, Mohammed SI: **Part I: cancer in Sudan-burden, distribution, and trends breast, gynecological, and prostate cancers**. *Cancer Med* 2015:n/a–n/a.

2. Saeed IE, Weng H-Y, Mohamed KH, Mohammed SI: **Cancer incidence in Khartoum, Sudan: first results from the Cancer Registry, 2009-2010.** *Cancer Med* 2014, **3**:1075–84.

3. Elgaili EM, Abuidris DO, Rahman M, Michalek AM, Mohammed SI: **Breast cancer burden in central Sudan.** *Int J Womens Health* 2010, **2**:77–82.

4. Awadelkarim KD, Arizzi C, Elamin EOM, Hamad HMA, De Blasio P, Mekki SO, Osman I, Biunno I, Elwali NE, Mariani-Costantini R, Barberis MC: **Pathological, clinical and prognostic characteristics of breast cancer in Central Sudan versus Northern Italy: implications for breast cancer in Africa.** *Histopathology* 2008, **52**:445–56.

5. Ahmed AAM: **Clinicopathological profile of female sudanese patients with locally advanced breast cancer.** *Breast Dis* 2014.

6. Katz T a, Huang Y, Davidson NE, Jankowitz RC: **Epigenetic reprogramming in breast cancer: From new targets to new therapies.** *Ann Med* 2014(April):1–12.

7. Liu H, Li X, Dong C: **Epigenetic and metabolic regulation of breast cancer stem cells**. *J Zhejiang Univ Sci B* 2015, **16**:10–17.

8. Lokk K, Modhukur V, Rajashekar B, Märtens K, Mägi R, Kolde R, Kolt Ina M, Nilsson TK, Vilo J, Salumets A, Tõnisson N: **DNA methylome profiling of human tissues identifies global and tissue-specific methylation patterns.** *Genome Biol* 2014, **15**:R54.

9. Fackler MJ, Lopez Bujanda Z, Umbricht C, Teo WW, Cho S, Zhang Z, Visvanathan K, Jeter S, Argani P, Wang C, Lyman JP, de Brot M, Ingle JN, Boughey J, McGuire K, King T a, Carey L a, Cope L, Wolff AC, Sukumar S: **Novel methylated biomarkers and a robust assay to detect circulating tumor DNA in metastatic breast cancer.** *Cancer Res* 2014, **74**:2160–70.

10. Bock C: **Epigenetic biomarker development**. :1–14.

11. Bediaga NG, Acha-Sagredo A, Guerra I, Viguri A, Albaina C, Ruiz Diaz I, Rezola R, Alberdi MJ, Dopazo J, Montaner D, Renobales M, Fernández AF, Field JK, Fraga MF, Liloglou T, de Pancorbo MM: **DNA methylation epigenotypes in breast cancer molecular subtypes.** *Breast Cancer Res* 2010, **12**:R77.

12. Rhee J-K, Kim K, Chae H, Evans J, Yan P, Zhang B-T, Gray J, Spellman P, Huang TH-M, Nephew KP, Kim S: **Integrated analysis of genome-wide DNA methylation and gene expression profiles in molecular subtypes of breast cancer.** *Nucleic Acids Res* 2013, **41**:8464–74.

13. Park SY, Kwon HJ, Choi Y, Lee HE, Kim S-W, Kim JH, Kim IA, Jung N, Cho N-Y, Kang GH: **Distinct patterns of promoter CpG island methylation of breast cancer subtypes are associated with stem cell phenotypes**. *Mod Pathol* 2011, **25**:185–196.

14. List M, Hauschild A, Tan Q, Kruse TA: **Classification of Breast Cancer Subtypes by combining Gene Expression and DNA Methylation Data**. 2014, **11**.

15. Stefansson O a., Moran S, Gomez A, Sayols S, Arribas-Jorba C, Sandoval J, Hilmarsdottir H, Olafsdottir E, Tryggvadottir L, Jonasson JG, Eyfjord J, Esteller M: **A DNA methylation-based definition of biologically distinct breast cancer subtypes**. *Mol Oncol* 2014:1–14.

16. Lechner M, Fenton T, West J, Wilson G, Feber A, Henderson S, Thirlwell C, Dibra HK, Jay A, Butcher L, Chakravarthy AR, Gratrix F, Patel N, Vaz F, O'Flynn P, Kalavrezos N, Teschendorff AE, Boshoff C, Beck S: **Identification and functional validation of HPV-mediated hypermethylation in head and neck squamous cell**

**carcinoma.** *Genome Med* 2013, **5**:15.

17. Li L, Zhang Y, Guo BB, Chan FKLL, Tao Q: **Oncogenic induction of cellular high CpG methylation by Epstein-Barr virus in malignant epithelial cells.** *Chin J Cancer* 2014:604–608.

18. Gulley ML: **Genomic assays for Epstein – Barr virus-positive gastric adenocarcinoma**. 2015, **47**:e134–12.

19. Consortium RE, Kundaje A, Meuleman W, Ernst J, Bilenky M, Yen A, Heravi-Moussavi A, Kheradpour P, Zhang Z, Wang J, Ziller MJ, Amin V, Whitaker JW, Schultz MD, Ward LD, Sarkar A, Quon G, Sandstrom RS, Eaton ML, Wu Y-C, Pfenning AR, Wang X, Claussnitzer M, Liu Y, Coarfa C, Harris RA, Shoresh N, Epstein CB, Gjoneska E, Leung D, et al.: **Integrative analysis of 111 reference human epigenomes**. *Nature* 2015, **518**:317–330.

20. Berman H, Zhang J, Crawford YG, Gauthier ML, Fordyce CA, McDermott KM, Sigaroudinia M, Kozakiewicz K, Tlsty TD: **Genetic and epigenetic changes in mammary epithelial cells identify a subpopulation of cells involved in early carcinogenesis.** *Cold Spring Harb Symp Quant Biol* 2005, **70**:317–27.

21. Dumont N, Crawford YG, Sigaroudinia M, Nagrani SS, Wilson MB, Buehring GC, Turashvili G, Aparicio S, Gauthier ML, Fordyce CA, Mcdermott KM, Tlsty TD: **Research article Human mammary cancer progression model recapitulates methylation events associated with breast premalignancy**. 2009, **11**:1–17.

22. Wessely F, Emes RD: **Identification of DNA methylation biomarkers from Infinium arrays.** *Front Genet* 2012, **3**(August):161.

23. Urrutia G, Laurito S, Marzese DM, Gago F, Orozco J, Tello O, Branham T, Campoy EM, Roqué M: **Epigenetic variations in breast cancer progression to lymph node metastasis.** *Clin Exp Metastasis* 2015, **32**:99–110.

24. Wang D, Yang P-N, Chen J, Zhou X-Y, Liu Q-J, Li H-J, Li C-L: **Promoter hypermethylation may be an important mechanism of the transcriptional inactivation of ARRDC3, GATA5, and ELP3 in invasive ductal breast carcinoma.** *Mol Cell Biochem* 2014, **396**:67–77.

25. Ruike Y, Imanaka Y, Sato F, Shimizu K, Tsujimoto G: **Genome-wide analysis of aberrant methylation in human breast cancer cells using methyl-DNA immunoprecipitation combined with high-throughput sequencing.** *BMC Genomics* 2010, **11**:137.

26. Conway K, Edmiston SN, May R, Kuan PF, Chu H, Bryant C, Tse C-K, Swift-Scanlan T, Geradts J, Troester MA, Millikan RC: **DNA methylation profiling in the**

**Carolina Breast Cancer Study defines cancer subclasses differing in clinicopathologic characteristics and survival.** *Breast Cancer Res* 2014, **16**:450.

27. Marzese DM, Hoon DSB, Chong KK, Gago FE, Orozco JI, Tello OM, Vargas-Roig LM, Roqué M: **DNA methylation index and methylation profile of invasive ductal breast tumors.** *J Mol Diagn* 2012, **14**:613–22.

28. Moarii M, Pinheiro A, Sigal-Zafrani B, Fourquet A, Caly M, Servant N, Stoven V, Vert J-P, Reyal F: **Epigenomic Alterations in Breast Carcinoma from Primary Tumor to Locoregional Recurrences**. *PLoS One* 2014, **9**:e103986.

29. Huang W-Y, Hsu S-D, Huang H-Y, Sun Y-M, Chou C-H, Weng S-L, Huang H-D: **MethHC: a database of DNA methylation and gene expression in human cancer.** *Nucleic Acids Res* 2014, **5712121**:1–6.

30. Fleischer T, Frigessi A, Johnson KC, Edvardsen H, Touleimat N, Klajic J, Riis MLH, Haakensen VD, Wärnberg F, Naume B, Helland Å: **Genome-wide DNA methylation profiles in progression to in situ and invasive carcinoma of the breast with impact on gene transcription and prognosis**. 2014:1–13.

31. Tommasi S, Karm DL, Wu X, Yen Y, Pfeifer GP: **Methylation of homeobox genes is a frequent and early epigenetic event in breast cancer.** *Breast Cancer Res* 2009, **11**:R14.

32. Reinhardt D, Cruickshanks HA, Mjoseng HK, Rmcphersonedacuk RCM, Lentini A, Donnchadunicanigmmedacuk DSD, Saripenningsedacuk SP, Steveandertonedacuk SMA, Benson M, Richardmeehanigmmedacuk RRM: **Rapid reprogramming of epigenetic and transcriptional profiles in mammalian culture systems**. 2015.

33. Kalari S, Pfeifer GP: *Identification of Driver and Passenger DNA Methylation in Cancer by Epigenomic Analysis*. 1st edition. *Volume 70*. Elsevier Inc.; 2010.

34. Hon GC, Hawkins RD, Caballero OL, Lo C, Lister R, Pelizzola M, Valsesia A, Ye Z, Kuan S, Edsall LE, Camargo AA, Stevenson BJ, Ecker JR, Bafna V, Strausberg RL, Simpson AJ, Ren B: **Global DNA hypomethylation coupled to repressive chromatin domain formation and gene silencing in breast cancer**. 2012:246–258.

35. Stirzaker C, Zotenko E, Song JZ, Qu W, Nair SS, Locke WJ, Stone A, Armstong NJ, Robinson MD, Dobrovic A, Avery-kiejda KA, Peters KM, French JD, Stein S, Korbie DJ, Trau M, Forbes JF, Scott RJ, Brown MA, Francis GD, Clark SJ: **prognostic value**. *Nat Commun* 2015, **6**:1–11.

36. Gonzalez ME, Moore HM, Li X, Toy K a, Huang W, Sabel MS, Kidwell KM, Kleer CG: **EZH2 expands breast stem cells through activation of NOTCH1 signaling.** *Proc Natl Acad Sci U S A* 2014, **111**:3098–103.

37. Cho J-H, Dimri M, Dimri GP: **A Positive Feedback Loop Regulates the Expression of Polycomb Group Protein BMI1 via WNT Signaling Pathway.** *J Biol Chem* 2012, **288**:3406–3418.

38. van Vlerken LE, Kiefer CM, Morehouse C, Li Y, Groves C, Wilson SD, Yao Y, Hollingsworth RE, Hurt EM: **EZH2 is required for breast and pancreatic cancer stem cell maintenance and can be used as a functional cancer stem cell reporter.** *Stem Cells Transl Med* 2013, **2**:43–52.

39. Polytarchou C, Iliopoulos D, Struhl K: **An integrated transcriptional regulatory circuit that reinforces the breast cancer stem cell state**. *Proc Natl Acad Sci* 2012, **109**:14470–14475.

40. Reyngold M, Turcan S, Giri D, Kannan K, Walsh L a., Viale A, Drobnjak M, Vahdat LT, Lee W, Chan T a.: **Remodeling of the methylation landscape in breast cancer metastasis**. *PLoS One* 2014, **9**:1–10.

41. Shi P, Feng J, Chen C: **Hippo pathway in mammary gland development and breast cancer**. *Acta Biochim Biophys Sin (Shanghai)* 2014, **47**(December 2014):53–59.

42. Hiemer SE, Szymaniak AD, Varelas X: **The transcriptional regulators TAZ and YAP direct transforming growth factor β-induced tumorigenic phenotypes in breast cancer cells.** *J Biol Chem* 2014, **289**:13461–74.

43. Lin L, Sabnis AJ, Chan E, Olivas V, Cade L, Pazarentzos E, Asthana S, Neel D, Yan JJ, Lu X, Pham L, Wang MM, Karachaliou N, Cao MG, Manzano JL, Ramirez JL, Torres JMS, Buttitta F, Rudin CM, Collisson E a, Algazi A, Robinson E, Osman I, Muñoz-Couselo E, Cortes J, Frederick DT, Cooper Z a, McMahon M, Marchetti A, Rosell R, et al.: **The Hippo effector YAP promotes resistance to RAF- and MEK-targeted cancer therapies**. *Nat Genet* 2015, **47**:250–256.

44. Keren-Paz A, Emmanuel R, Samuels Y: **YAP and the drug resistance highway**. *Nat Genet* 2015, **47**:193–194.

45. Chiang K-C, Yeh C-N, Chung L-C, Feng T-H, Sun C-C, Chen M-F, Jan Y-Y, Yeh T-S, Chen S-C, Juang H-H: **WNT-1 inducible signaling pathway protein-1 enhances growth and tumorigenesis in human breast cancer**. *Sci Rep* 2015, **5**:8686.

46. Xu J, Prosperi JR, Choudhury N, Olopade OI, Goss KH: **β-Catenin Is Required for the Tumorigenic Behavior of Triple-Negative Breast Cancer Cells**. *PLoS One* 2015, **10**:e0117097.

47. Meier-abt F, Bentires-alj M, Rochlitz C: **Breast Cancer Prevention : Lessons to be Learned from Mechanisms of Early Pregnancy – Mediated Breast Cancer Protection**. 2015:803–808.

48. Naumov V a, Generozov E V, Zaharjevskaya NB, Matushkina DS, Larin AK, Chernyshov S V, Alekseev M V, Shelygin Y a, Govorun VM: **Genome-scale analysis of DNA methylation in colorectal cancer using Infinium HumanMethylation450 BeadChips.** *Epigenetics* 2013, **8**:921–34.

49. Zhu W, Qin W, Hewett JE, Sauter ER: **Quantitative evaluation of DNA hypermethylation in malignant and benign breast tissue and fluids.** *Int J Cancer* 2010, **126**:474–82.

50. Zhu W, Qin W, Zhang K, Rottinghaus GE, Chen Y-C, Kliethermes B, Sauter ER: **Trans-resveratrol alters mammary promoter hypermethylation in women at increased risk for breast cancer.** *Nutr Cancer* 2012, **64**:393–400.

51. Qin W, Zhang K, Clarke K, Weiland T, Sauter ER: **Methylation and miRNA effects of resveratrol on mammary tumors vs. normal tissue.** *Nutr Cancer* 2014, **66**:270–7.

52. Yahia ZA, Adam AA, Elgizouli M, Hussein A, Masri MA, Kamal M, Mohamed HS, Alzaki K, Elhassan AM, Hamad K, Ibrahim ME: **Epstein Barr virus: a prime candidate of breast cancer aetiology in Sudanese patients.** *Infect Agent Cancer* 2014, **9**:9.

53. Masri MA, Abdel Seed NM, Fahal AH, Romano M, Baralle F, El Hassam AM, Ibrahim ME: **Minor role for BRCA2 (exon11) and p53 (exon 5-9) among Sudanese breast cancer patients.** *Breast Cancer Res Treat* 2002, **71**:145–7.

54. Lin D-C, Meng X, Hazawa M, Nagata Y, Varela AM, Xu L, Sato Y, Liu L-Z, Ding L-W, Sharma A, Goh BC, Lee SC, Petersson BF, Yu FG, Macary P, Oo MZ, Ha CS, Yang H, Ogawa S, Loh KS, Koeffler HP: **The genomic landscape of nasopharyngeal carcinoma.** *Nat Genet* 2014, **46**:866–871.

55. Uozaki H, Fukayama M: **Epstein-Barr virus and gastric carcinoma--viral carcinogenesis through epigenetic mechanisms.** *Int J Clin Exp Pathol* 2008, **1**:198–216.

56. Niller HH, Szenthe K, Minarovits J: **Epstein-Barr virus-host cell interactions: an epigenetic dialog?** *Front Genet* 2014, **5**(October):367.

57. Zhao J, Liang Q, Cheung K-F, Kang W, Lung RWM, Tong JHM, To KF, Sung JJY, Yu J: **Genome-wide identification of Epstein-Barr virus-driven promoter methylation profiles of human genes in gastric cancer cells.** *Cancer* 2012:1–9.

58. McLean CY, Bristor D, Hiller M, Clarke SL, Schaar BT, Lowe CB, Wenger AM, Bejerano G: **GREAT improves functional interpretation of cis-regulatory regions.** *Nat Biotechnol* 2010, **28**:495–501.

59. **Wizard® Genomic DNA Purification Kit Protocol**

[https://www.promega.com/resources/protocols/technical-manuals/0/wizard-genomic-dna-purification-kit-protocol/]

60. **Infinium HumanMethylation450 BeadChip Kit** [http://www.illumina.com/products/methylation_450_beadchip_kits.html]

61. Dedeurwaerder S, Defrance M, Calonne E, Denis H, Sotiriou C, Fuks F: **Evaluation of the Infinium Methylation 450K technology.** *Epigenomics* 2011, **3**:771–84.

62. Smyth GK: **Linear models and empirical bayes methods for assessing differential expression in microarray experiments.** *Stat Appl Genet Mol Biol* 2004, **3**:Article3.

63. Giardine B, Riemer C, Hardison RC, Burhans R, Elnitski L, Shah P, Zhang Y, Blankenberg D, Albert I, Taylor J, Miller W, Kent WJ, Nekrutenko A: **Galaxy: A platform for interactive large-scale genome analysis**. *Genome Res* 2005, **15**:1451–1455.

64. Blankenberg D, Von Kuster G, Coraor N, Ananda G, Lazarus R, Mangan M, Nekrutenko A, Taylor J: **Galaxy: a web-based genome analysis tool for experimentalists.** *Curr Protoc Mol Biol* 2010, **Chapter 19**:Unit 19.10.1–21.

65. Goecks J, Nekrutenko A, Taylor J: **Galaxy: a comprehensive approach for supporting accessible, reproducible, and transparent computational research in the life sciences.** *Genome Biol* 2010, **11**:R86.

66. Fi R, Plugin C: **Reactome FI Cytoscape Plugin 4 Reactome FI Cytoscape Plugin 4 - ReactomeWiki**. 2014:1–23.

67. Shannon P, Markiel A, Ozier O, Baliga NS, Wang JT, Ramage D, Amin N, Schwikowski B, Ideker T: **Cytoscape: A software Environment for integrated models of biomolecular interaction networks**. *Genome Res* 2003, **13**:2498–2504.

**Illustrations and figures**

**Figure 1: Genome-wide DNA Differential methylation Analysis of study samples**.

(A) shows differentially methylated CpG sites (defined as median beta value difference equal to or more than 0.2) identified using three methods: Limma (L; 34099 sites),

Wilcoxon (W; 39940 sites), and Nimbl, (N; 22251 sites). The color code shows sites identified by each method alone and in combination. A final set which represents the intersection of three approaches (L + W + N; black dots) consisted of 20188 sites was obtained by Nimble-compare module and used for analysis in this study. (B) A volcano plot showing the demarcation of differentially methylated sites by both statistical significance and effect size. The sites targeted in this study are those with high effect size (median beta value difference equal to or more than 0.2) and low p-value (equal to or more than 0.01, shown as –log10). The dotted lines show these cut-offs. Targeted sites for analysis are those in outer upper rectangular area of the plot.

**Figure 2: Hierarchical clustering of highly differentially methylated positions**. Differentially methylated positions (DMPs) were sorted by F value (low intragroup variability and higher intergroup variability) and the top 250 sites were tested for clustering between study samples. Hierarchical clustering heatmap and dendogram are depicted in this figure, showing a clear separation of tumor samples from normal samples (top dendogram, control samples above green bar, tumor samples above orange bar). DMS median p-value heatmap shows a contrasting state of differential methylation between tumor and control samples indicating both gain and loss of differential methylation states in tumor tissues.

**Figure 3: Genomic and epigenomic distribution of differentially methylated positions (DMPs)**. This figure details the number of DMPs in relation to gene elements, CpG islands and chromatin states. (A) Distribution of hyper and hypo methylated CpG sites in relation to gene elements. TSS: transcription start site, UTR: untranslated region. (B) Distribution of hyper and hypo methylated CpG sites in relation to CpG Islands. N_:

north, S_: south (C) Distribution of Hypomethylated CpG sites in relation to chromatin states. (D) Distribution of Hypermethylated CpG sites in relation to chromatin states. Fourteen chromatin states are shown:

**Figure 4: Hypermethylation of the TP73 gene**. Differential methylation Beta-values for 8 tumor and 8 control samples at methylation array sites of TP73 gene are shown. The figure contains three tracks: the genomic location of the TP73 and its different RefSeq transcripts are shown in the 'Chromosome' and 'RefSeq genes' tracks respectively; the 'Methylation' track shows the methylation level in each tumor sample (red dot) and control sample (blue dot). The overall discordance in methylation Beta-values between tumor samples (red line in the methylation track) and control samples (blue line) is notable specially at TSS both for the long and short RefSeq transcripts (genomic areas around 3.56 mb and 3.6 mb, respectively). Tumor samples show relatively high beta-values compared to controls at these sites indicating differential promoter hypermethylation. TSS: Transcription Start Site.


**Figure 5: Tumor Epigenetically Upregulated Genes (EUG) in Epstein-Barr Virus Infection pathway**. Many genes bearing methylation marks that promote gene expression  (hypomethylation in the promoter area and first exon or hypermethylation in the gene body region) – referred to in this study Epigenetically Upregulated Genes – were found to be integral parts of EBV Infection KEGG pathway (highlighted  red and grey boxes). This group of genes showed significant enrichment for Epstein-Barr Virus Infection Pathway (red boxes are highly enriched nodes). Epstein-Barr Virus Infection

KEGG Pathway was obtained from KEGG pathways database

[http://www.kegg.jp/pathway/hsa05169].

**Tables and captions**

**Table 1 Pathway enrichment results for EUG network**

| Pathway | Number of genes in the Geneset | Number of genes in the Network | FDR |
|---|---|---|---|
| Hippo signaling pathway | 154 | 31 | <1.000e-03 |
| Arrhythmogenic right ventricular cardiomyopathy (ARVC) | 74 | 20 | <5.000e-04 |
| L1CAM interactions | 94 | 21 | 2.50E-04 |
| Wnt signaling pathway | 269 | 41 | 3.33E-04 |

**Table 1: Pathway enrichment analysis results for Epigenetically Upregulated Genes (EUG) interaction network**. This pathway enrichment analysis and the interaction network were prepared using ReactomeFI Cytoscape app. The table shows the enriched pathways, the number of genes in the pathway from the total query gene set, and the number of genes in the pathway found in the interaction network. Results having p-values <0.01 and a False Detection Rate < 0.001 are shown.

**Table 2 Pathway enrichment results for EUG breast cancer subnetwork**

| Pathway | Number of genes in the Geneset | Number of genes in the Network | FDR |
|---|---|---|---|
| CXCR4-mediated signaling events | 79 | 10 | 2.50E-04 |
| AP-1 transcription factor network | 70 | 9 | 7.27E-04 |
| HIF-1-alpha transcription factor network | 66 | 9 | 4.00E-04 |
| Viral myocarditis | 59 | 9 | 2.50E-04 |
| Pathways in cancer | 327 | 22 | <1.000e-03 |
| HTLV-I infection | 260 | 18 | 3.33E-04 |
| Proteoglycans in cancer | 225 | 17 | 2.00E-04 |
| Epstein-Barr virus infection | 202 | 16 | 1.67E-04 |
| Hippo signaling pathway | 154 | 14 | 3.33E-04 |
| Natural killer cell mediated cytotoxicity | 135 | 13 | 1.43E-04 |
| Alzheimer disease-presenilin pathway | 111 | 12 | 5.00E-04 |

**Table 2: Pathway enrichment results for breast cancer related Epigenetically Upregulated Genes (EUG) subnetwork**. ReactomeFI cytoscape app was used to extract breast cancer related subnetworks from EUG set by loading NCI cancer index and performing pathway enrichment analysis on interaction networks. Nodes that corresponded to malignant breast cancer were selected. The table shows the enriched pathways, the number of genes in the pathway from the total query gene set, and the number of genes in the pathway found in the interaction network. Results having p-values <0.01 and a False Detection Rate < 0.001 are shown.

**Table 3 Pathway enrichment results for EDG network**

| Pathway | Number of genes in the Geneset | Number of genes in the Network | FDR |
|---|---|---|---|
| Neuroactive ligand-receptor interaction | 275 | 81 | <1.000e-03 |
| GPCR ligand binding | 433 | 107 | <5.000e-04 |
| PI3K-Akt signaling pathway | 346 | 81 | <3.333e-04 |
| Extracellular matrix organization | 263 | 65 | <2.500e-04 |
| Pathways in cancer | 327 | 74 | <2.000e-04 |
| Rap1 signaling pathway | 213 | 54 | <1.667e-04 |
| Regulation of actin cytoskeleton | 215 | 54 | <1.429e-04 |
| Neurotransmitter Receptor Binding And Downstream Transmission In The Postsynaptic Cell | 137 | 40 | <1.250e-04 |
| Potassium Channels | 86 | 30 | <1.111e-04 |
| Heterotrimeric G-protein signaling pathway-Gi alpha and Gs alpha mediated pathway | 147 | 41 | <1.000e-04 |
| Proteoglycans in cancer | 225 | 54 | <9.091e-05 |
| Ras signaling pathway | 227 | 54 | <8.333e-05 |
| ECM-receptor interaction | 86 | 28 | 1.54E-04 |
| Calcium signaling pathway | 181 | 45 | 2.14E-04 |
| FGF signaling pathway | 92 | 29 | 2.00E-04 |
| Focal adhesion | 206 | 48 | 2.50E-04 |
| Gastrin-CREB signalling pathway via PKC and MAPK | 207 | 48 | 2.35E-04 |
| Cell adhesion molecules (CAMs) | 143 | 37 | 2.78E-04 |
| Wnt signaling pathway | 269 | 57 | 4.21E-04 |
| MAPK signaling pathway | 259 | 55 | 4.00E-04 |
| Heterotrimeric G-protein signaling pathway-Gq alpha and Go alpha mediated pathway | 108 | 30 | 5.24E-04 |
| IL4-mediated signaling events | 63 | 21 | 7.73E-04 |
| HTLV-I infection | 260 | 54 | 7.83E-04 |
| GABAergic synapse | 90 | 26 | 7.92E-04 |
| Signaling by Type 1 Insulin-like Growth Factor 1 Receptor (IGF1R) | 86 | 25 | 8.40E-04 |

| | | | |
|---|---|---|---|
| Retrograde endocannabinoid signaling | 103 | 28 | 8.85E-04 |
| Melanoma | 71 | 22 | 9.26E-04 |

**Table 3: Pathway analysis on the Epigenetically Downregulated Genes (EDG) interaction network**. The functional interaction network was constructed using ReactomeFI cytoscape app. The table shows the enriched pathways, the number of genes in the pathway from the total query gene set, and the number of genes in the pathway found in the interaction network. Results having p-values <0.01 and a False Detection Rate < 0.001 are shown.

**Additional Files**

**Additional File 1** (Fig S1): Genomic distribution of hypermethylation marks shown at each chromosome. Black color indicates hypermethylation sites. File format: PDF.

**Additional File 2** (Fig S2): Genomic distribution of hypomethylation marks shown at each chromosome. Black color indicates hypermethylation sites. File format: PDF.

**Additional File 3** (Table S1): Pathway enrichment results for breast cancer related Epigenetically Downregulated Genes (EDG) subnetwork. ReactomeFI cytoscape app was used to extract breast cancer related subnetworks from EUG set by loading NCI cancer index and performing pathway enrichment analysis on interaction networks. Nodes that corresponded to malignant breast cancer were selected. The table shows the enriched pathways, the number of genes in the pathway from the total query gene set, and the number of genes in the pathway found in the interaction network. Results
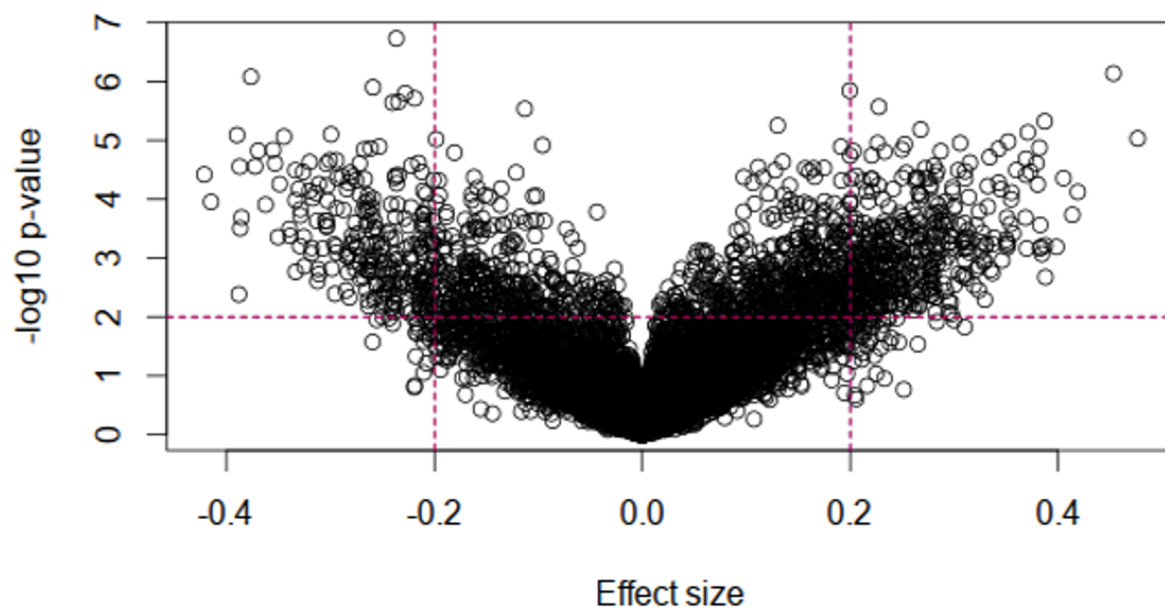
having p-values <0.01 and a False Detection Rate < 0.001 are shown. File format:
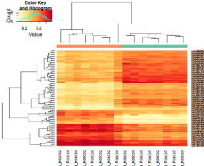
DOCX.

**Additional File 4** (Fig S3): Density plots of beta values for individual samples. Shades

of red and yellow colors represent tumor samples, whereas shades of blue and green
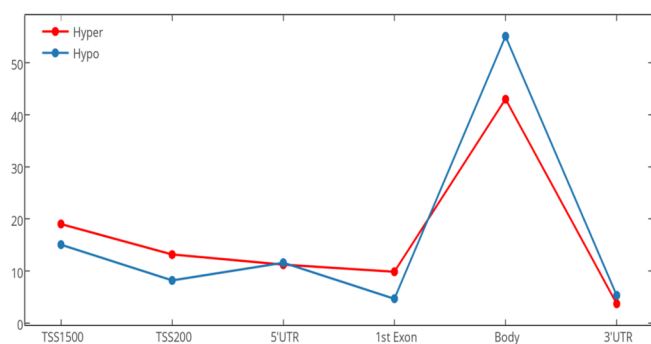
represent normal samples. File format: PDF.

**A**

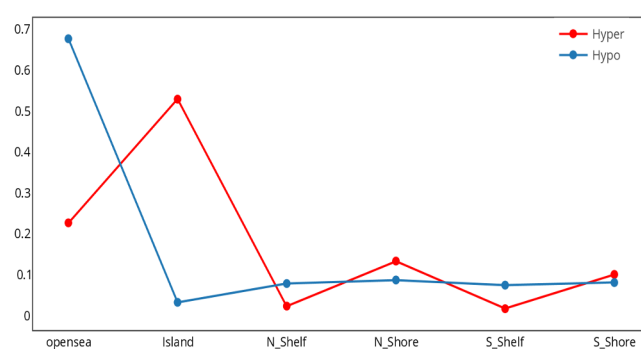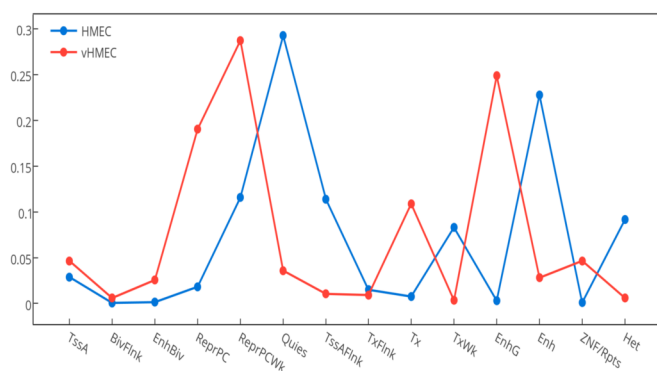Legend:
- N+L+W (20188)
- N+L (418)
- N+W (1101)
- L+W (1328)
- N (544)
- L (212)
- W (5370)

x-axis: median beta value: Control
y-axis: median beta value: Tumor

**B**

x-axis: Effect size
y-axis: -log10 p-value