

**Figure S1. Identification of breast cancer epitypes in the discovery cohort.** (A) Bootstrap clustering of 188 breast tumors using 2,108 CpGs and a seven-group cluster solution. The heatmap displays correlations between pairs of tumors. Row and column annotation bars represent the seven groups (epitypes). (B) Plot tracking group assignments of individual samples from bootstrap clustering of 188 breast tumors using 2,108 CpGs and solutions with three to seven sample clusters. Each column corresponds to a tumor sample. (C) Number of CpGs with differential methylation levels between 188 breast tumor samples and 96 normal samples as a function of number of tumors with differential methylation levels. 2,108 CpGs were identified as having differential methylation levels in at least 10 tumors as compared to normal tissue samples. (D) The seven bootstrap clusters (epitypes) are robust across a wide range of different CpG sets. The heatmap shows the overlap between bootstrap analyses performed using different number of CpGs from (C) and a seven-group cluster solution. (E) Principal component analysis of technical factors, clinicopathological factors, and DNA methylation epitypes. The full set of CpGs with no missing values was used (n=356,331). The heatmap shows the significance of the association between a principal component (columns) and a specific factor (rows) [21]. In the heatmap a more intense color means that more variation is explained by a specific factor. (F) Distribution of tumors across three bisulphite conversion plates used in the sample labeling across the seven breast cancer epitypes. (G) Distribution of tumors from different BeadChips across the seven breast cancer epitypes.

**Figure S2. Reproducibility of breast cancer epitypes in the TCGA validation cohort.** (A) Global hypermethylation scores for all CpGs (left) and all CpGs in promoters and CpG islands (right) across the epitypes for the validation cohort. The number of samples in each epitype is indicated at the top. (B) Principal component analysis of technical factors, clinicopathological factors, and DNA methylation epitypes for the 669 tumors in the validation cohort. The full set of

CpGs with no missing values was used (n=454,553). The heatmap shows the significance of the association between a principal component (columns) and a specific factor (rows) [21]. In the heatmap a more intense color means that more variation is explained by a specific factor. Epitype corresponds to the classification of the validation cohort using the classifier for the seven epitypes constructed using the discovery cohort. **(C)** Distribution of tumors in different TCGA sample batches across the seven breast cancer epitypes. The total number of tumors per epitype is indicated at the top. **(D)** Distribution of tumors from different TCGA BeadChips across the seven breast cancer epitypes. **(E)** Overlap of epitype classification (based on the seven epitypes identified using 2,108 CpGs and the discovery cohort) of tumors in the TCGA validation cohort (vertical axis) with unsupervised bootstrap tumor groups identified in the TCGA validation cohort (horizontal axis) using an eight group solution and the same 2,108 CpGs. Matrix cells indicate the percentage of tumors in an unsupervised bootstrap group found in an epitype. The number of samples in each unsupervised bootstrap group is indicated at the top.

**Figure S3. Activity of breast cancer specific gene expression modules across the epitypes. (A)** Expression levels of eight breast cancer specific gene modules in the discovery cohort (the microarray data in the discovery cohort contained expression levels for 12 of 21 genes in the mitotic checkpoint module, 13 of 22 in the mitotic progression module, 1 of 8 in the basal module, 7 of 9 in the steroid response module 39 of 71 in the immune response module, 20 of 32 in the stroma module, 4 of 18 in the lipid module, and 5 of 6 in the early response module). **(B)** Expression levels of eight breast cancer specific gene modules in the validation cohort (The TCGA RNAseq data in the validation cohort contained expression levels for all genes in the modules, except for two genes in the immune response module).

**Figure S4. Epitype-specific hypermethylation patterns.** (A) The average methylation level of the CpGs methylated in ET7 in each tumor in the discovery cohort (left) and validation cohort (right) stratified by epitype. (B) The average methylation level of the CpGs methylated in ET5 in each tumor in the discovery cohort (left) and validation cohort (right) stratified by epitype. The number of tumors in each epitype is indicated at the top of each of the plots.

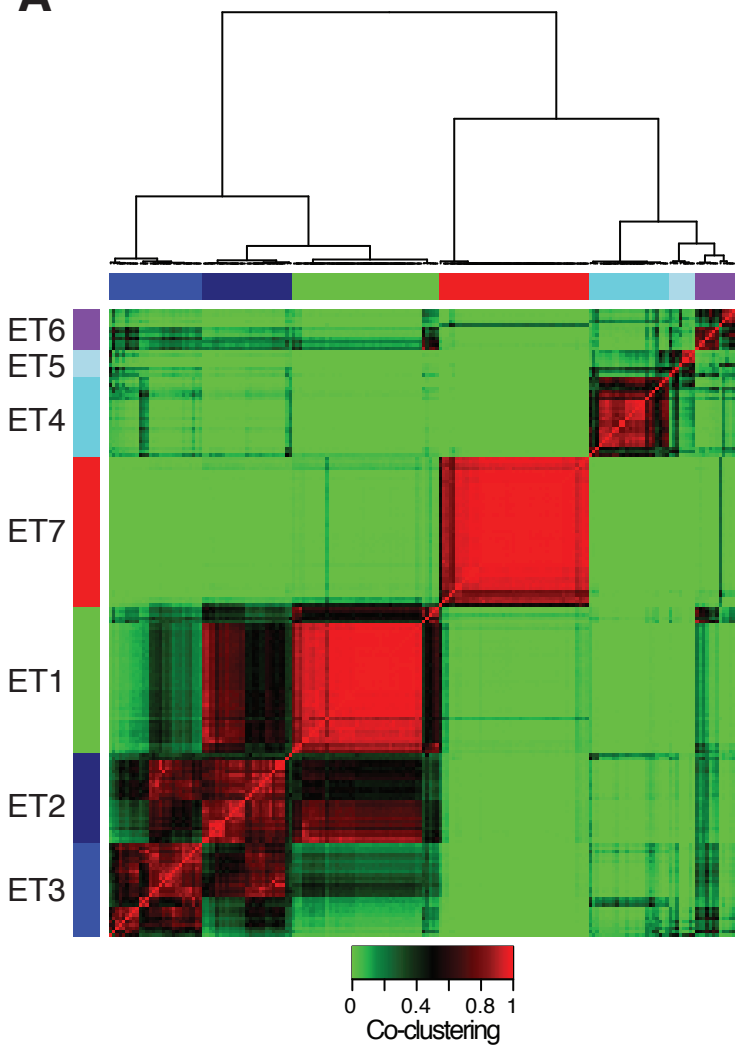
**Figure S5. Epitype-specific hypomethylation patterns.** (A) The average methylation level of the CpGs demethylated in ET4, demethylated in luminal tumors and demethylated in ET7 in each tumor in the discovery cohort (top) and validation cohort (bottom) stratified by epitype. The number of tumors in each epitype is indicated at the top of each of the plots. (B) For a given set of CpGs (All/Demethylated in ET4/Demethylated in luminal/Demethylated in ET7) a bar indicates the fraction of CpGs assigned to the respective chromatin state in HMEC. (C) Pearson correlation between gene expression and methylation levels for the CpGs selected as demethylated in ET7, luminal and ET4. The numbers on top indicate the number of CpGs matched to a gene with gene expression data. (D) The average gene expression levels across 661 breast tumors in the validation cohort stratified by epitype and 106 normal breast tissue samples from TCGA for the CpGs demethylated in ET4, luminal and ET7 and matched to a gene with gene expression data. (E) The average methylation levels across different subpopulations of blood cells for the CpGs demethylated in ET4, luminal and ET7. (F) For a given set of CpGs (All/Demethylated in ET4/Demethylated in luminal) a bar indicates the fraction of CpGs situated in different types of DNA repeats. (G) For a given set of CpGs (All/Demethylated in ET4/Demethylated in luminal) a bar indicates the fraction of CpGs with more or less than 5Mb to the nearest chromosome end.

**Figure S6. Amplifications and IntClust groups across epitypes.** (A) The frequency of samples in

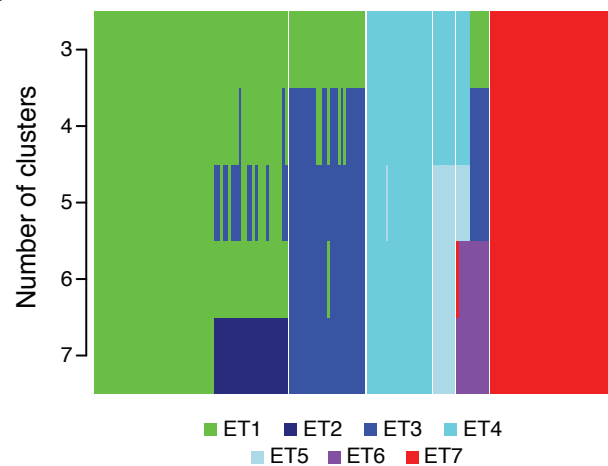
each epitype harboring specific amplifications for the discovery cohort (left) and the validation cohort (right). **(B)** The number of samples in each epitype classified into the ten IntClust groups for the discovery cohort (left) and the validation cohort (right). The number of tumors in each epitype is indicated at the top of each of the plots.

Figure S1

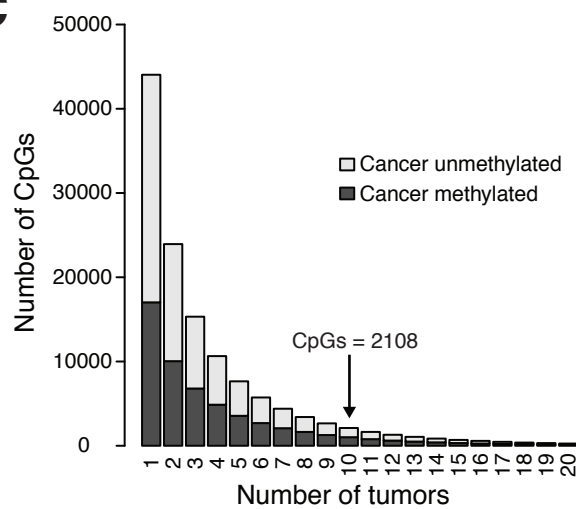
**A**



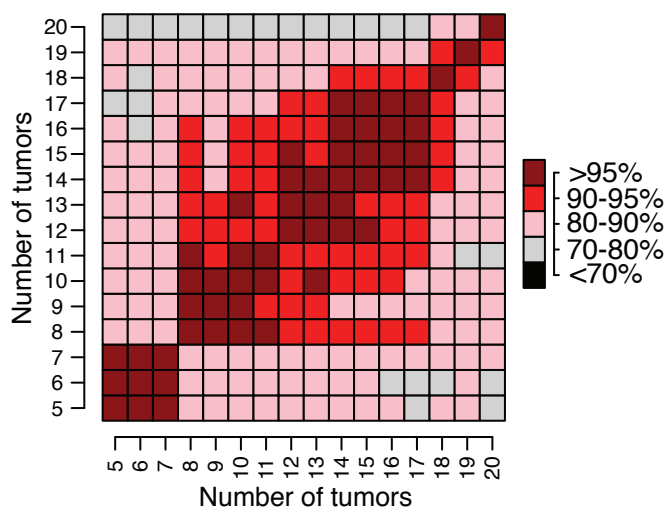
**B**



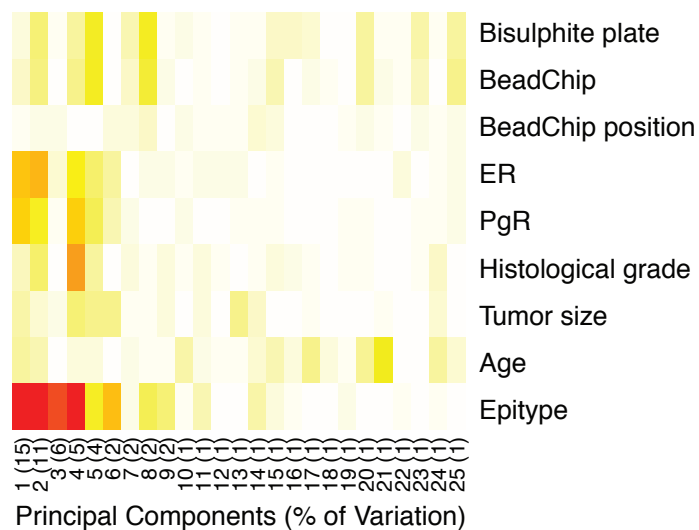
**C**



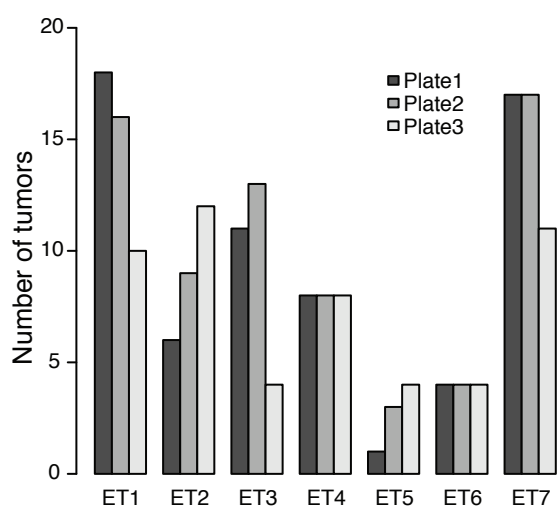
**D**



**E**



**F**



**G**

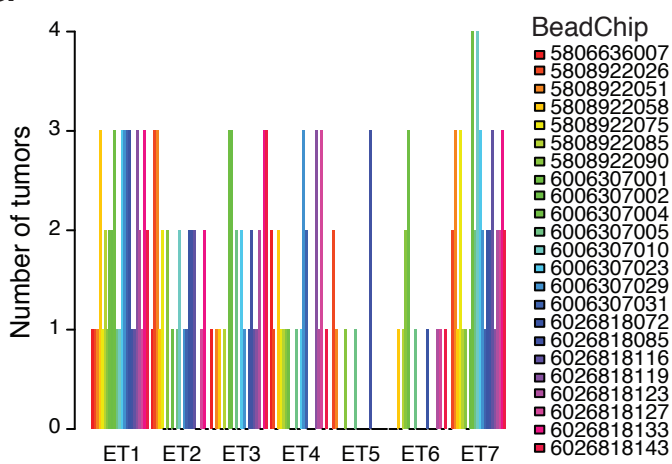
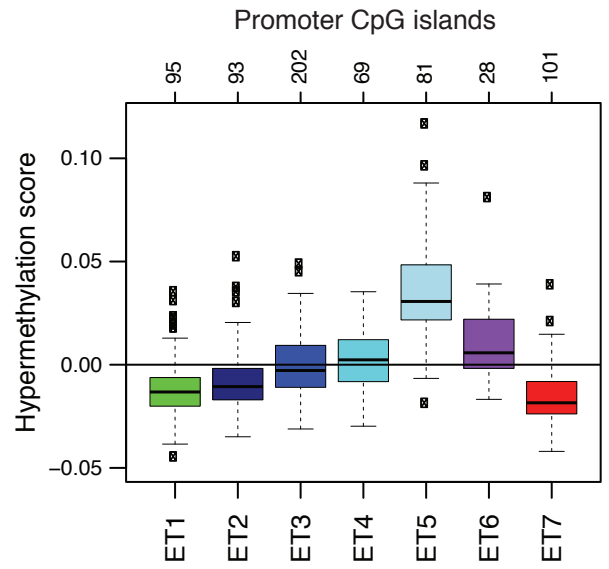
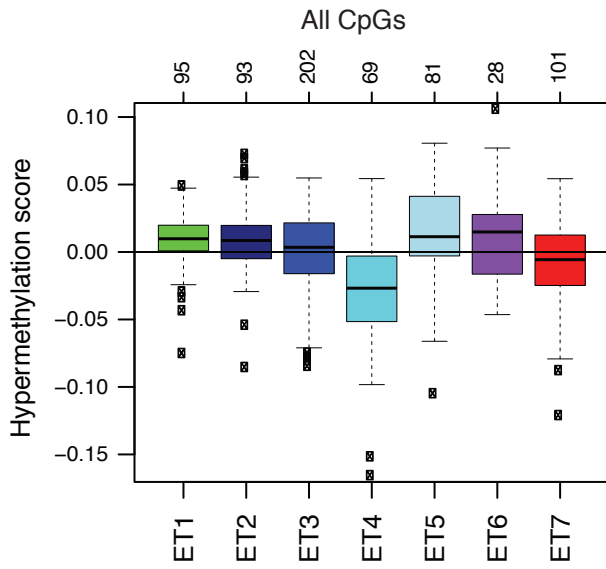
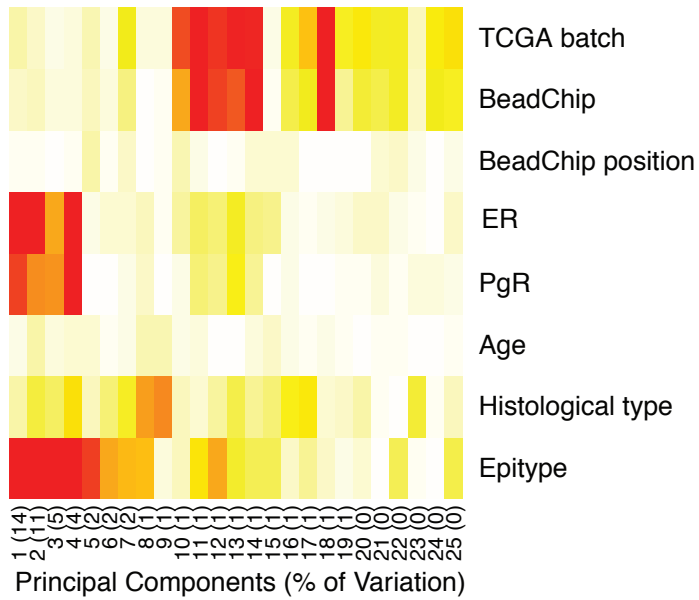


Figure S2

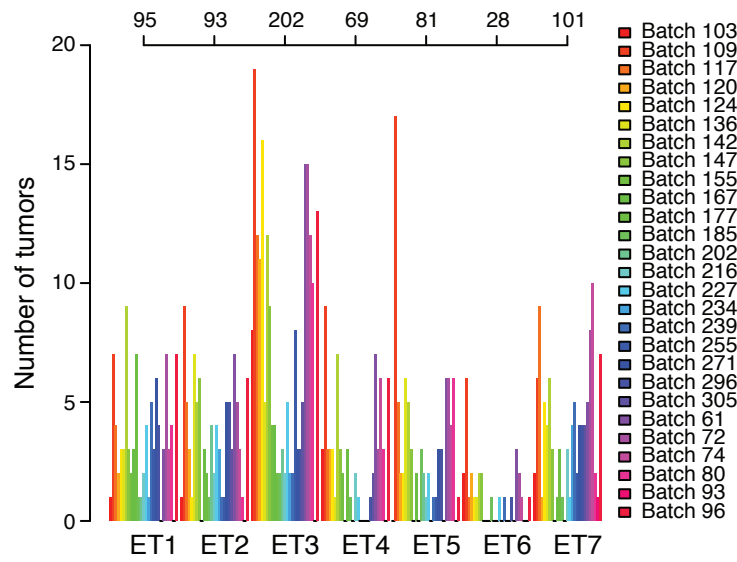
**A**



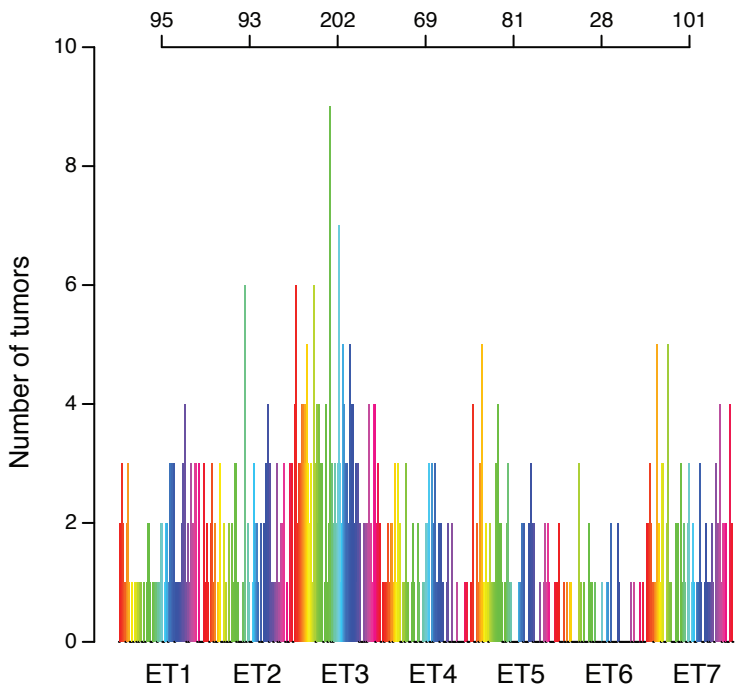
**B**



**C**



**D**



**E**

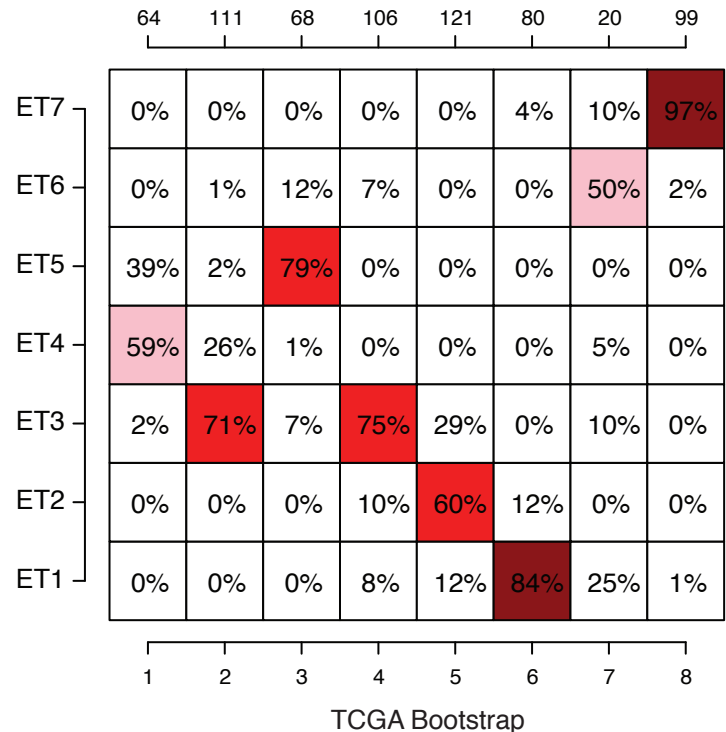
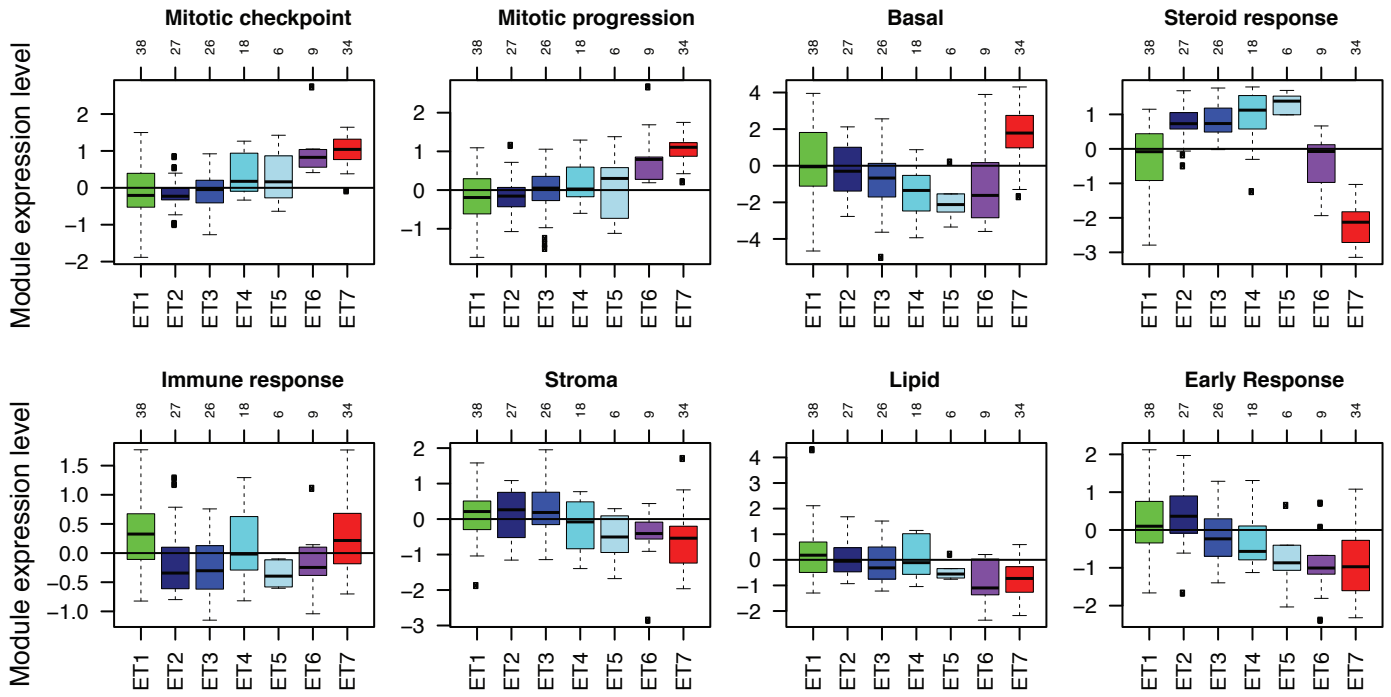


Figure S3

**A**

Discovery cohort



**B**

Validation cohort

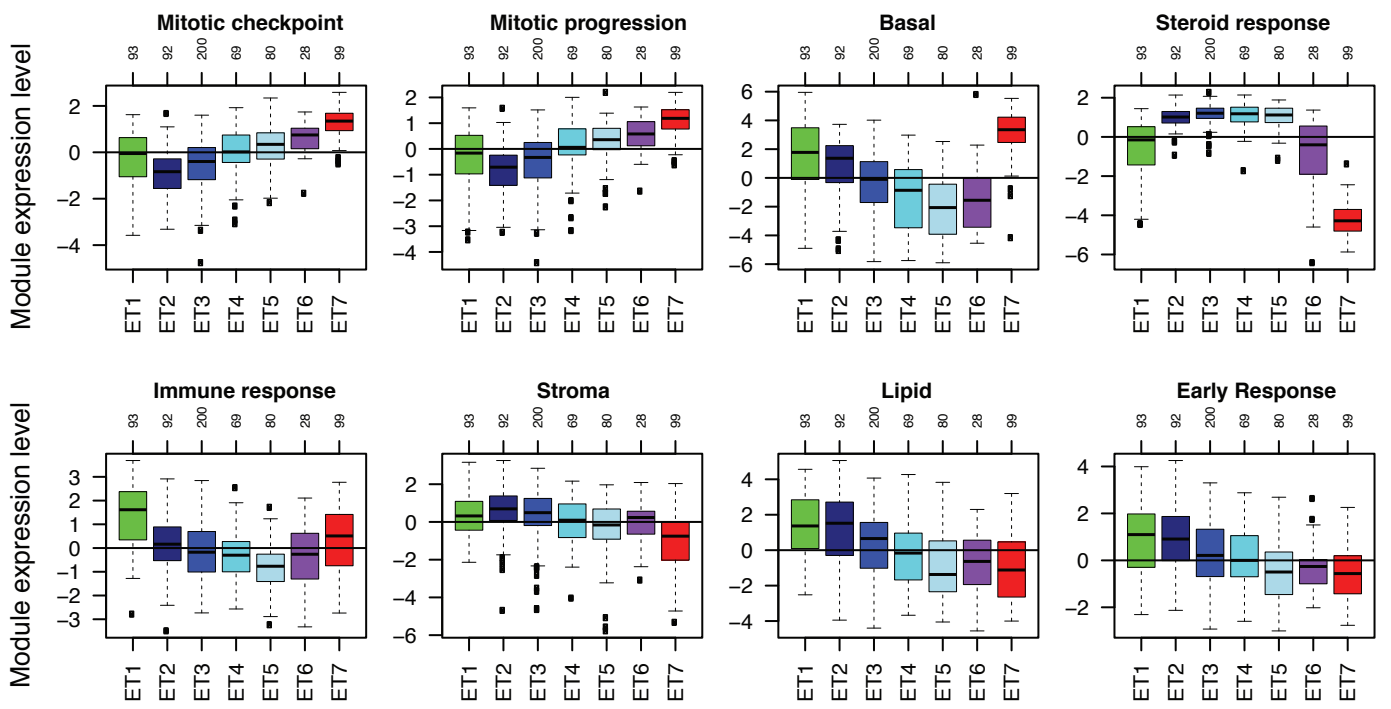
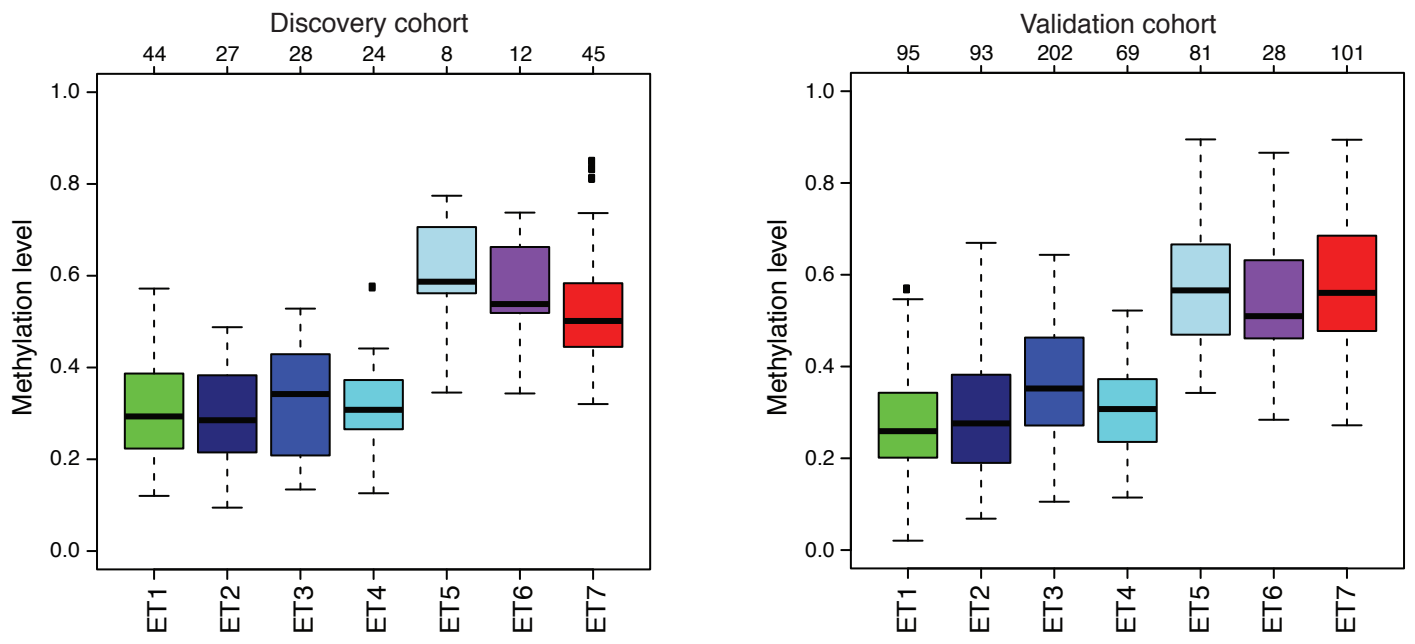


Figure S4

**A**

Methylated in ET7



**B**

Methylated in ET5

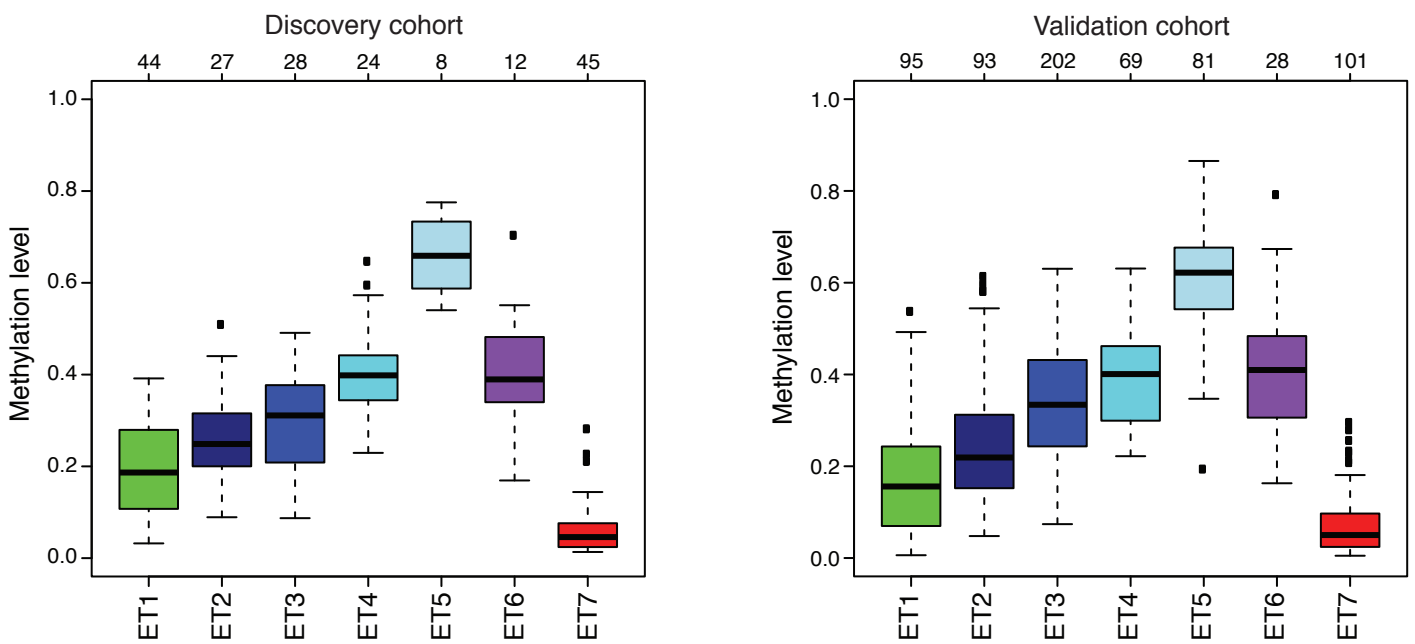
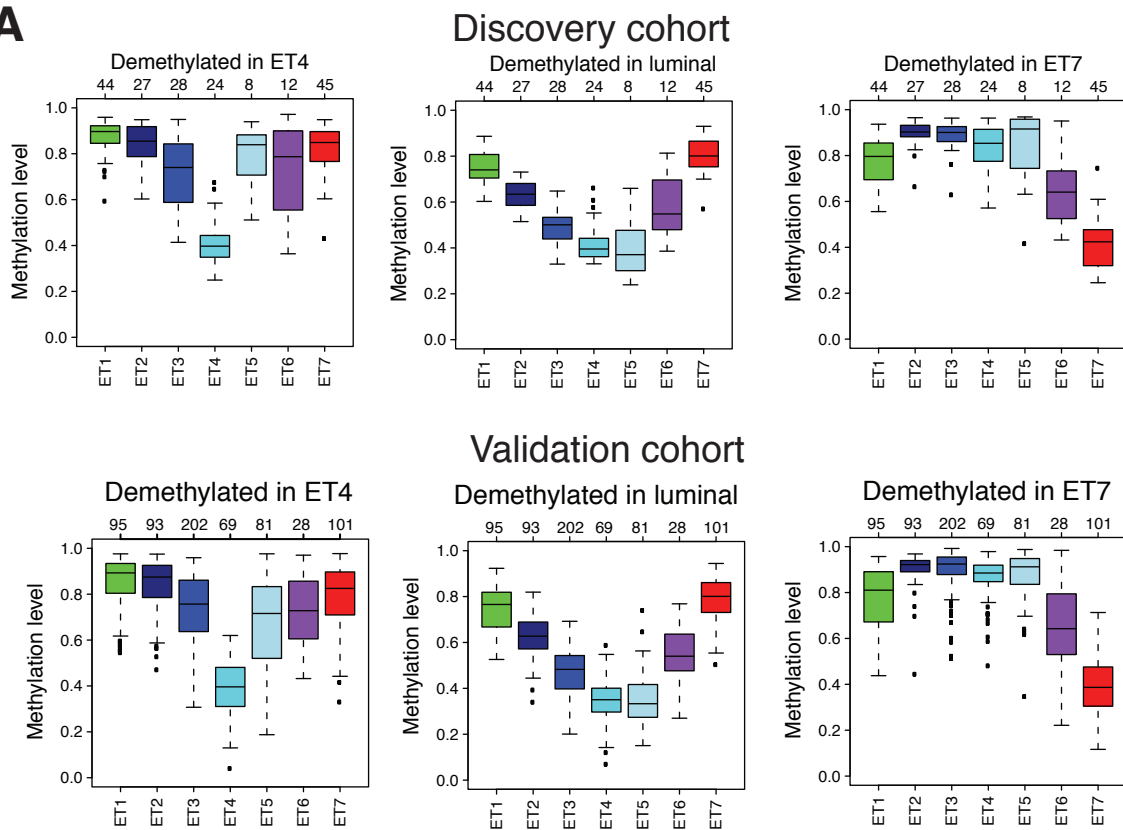


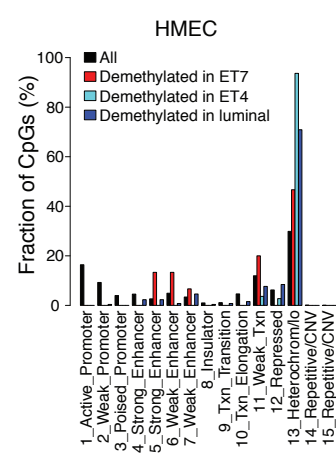


Figure S5

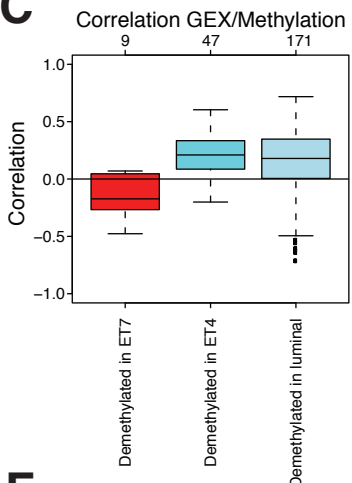
**A**



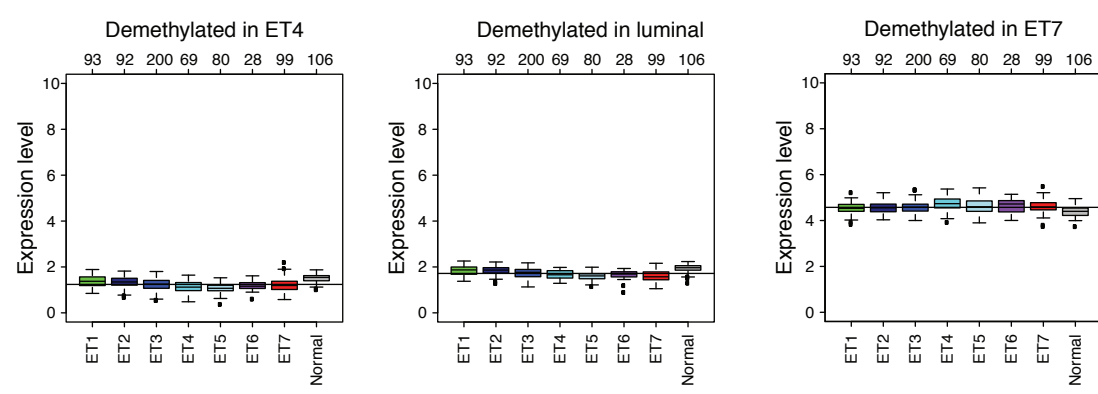
**B**



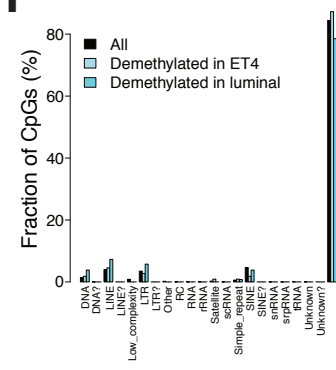
**C**



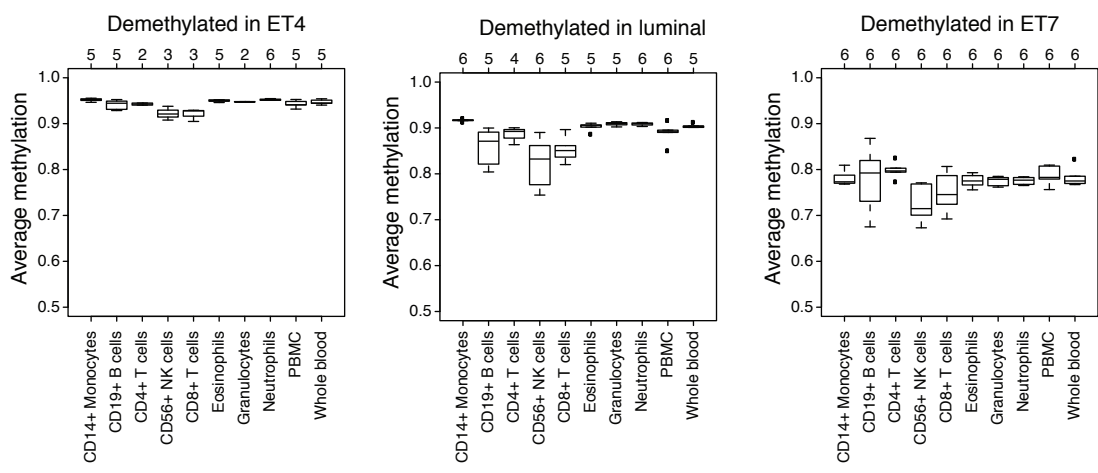
**D**



**F**



**E**



**G**

