

Figure S1.

Probe signal (y-axis) as a function of distance to enzyme cut PmeI (x-axis). The green lines indicate +/- 600 bases from the enzyme cut.

Figure S2.

These figures show the results of bisulfite validation for 3 loci.

A) a gene, two de-methylated probes in the same experiment.

B) AluSq element, microarray methylation profiling performed twice, each time demonstrating the methylation of the region.

C) AluY element, methylation profiling performed 5 times, each time demonstrating the de-methylation of the region.

Each figure is composed of a picture generated using UCSC genome browser to plot the location of the probe in the genomic context. Below, a textual summary that shows the location of the probe (●), CpG island (*) repetitive elements (T-LTR, A-Alu, 1-Line1, 2-Line2, =-other, S-Sine). Vertical bars indicate the location of an enzyme recognition site for Acil, HhaI and McrBC enzymes. The region in the text is analogous to the one shown above in the UCSC genome browser graphic. The region surrounding the probe (highlighted in YELLOW on top of the McrBC endonuclease cleavage track) has been bisulfite treated, amplified and sequenced. The outcome of sequencing is shown in the bottom portion of the figure. The sequenced reads were aligned to the genomic sequence using clustalW. Subsequently the alignments were collapsed to only show locations of Cs in CpGs. An unmethylated CpG is indicated as “○” and methylated as “@”.

Figure S3A

Ordering within plots. Per-experiment average methylation levels of the most informative subset of L1P and the least informative probes near DNA transposons and AluSq regions. Experiments are not ordered. The faint blue line indicates the average values of significant L1P probes in normal, non-tumor adjacent, tumor and sperm experiments (from top to bottom).

Figure S3B

Ordering within plots. Per-experiment average methylation levels of the most informative subset of L1P and the least informative probes near DNA transposons and AluSq regions. Experiments ordered within their groups based on L1P - mean(AluSq + DNA) probe values. The faint blue line indicates the average values of significant L1P probes in normal, non-tumor adjacent, tumor and sperm experiments (from top to bottom)

Figure S4A.

Example of a per-experiment plot showing average methylation levels of 4 categories of genomic compartments. See Supplementary Section 3.

Figure S4B.

Example of per-category plot showing distributions of average methylation levels of 4 categories of genomic compartments. See Supplementary Section 3

Figure S5.

Normal, Non-Tumor Adjacent, Tumor, Sperm variance in methylation levels ordered by repetitive element.

Figure S6.

Normal, Non-Tumor Adjacent,, Tumor, Sperm variance in methylation levels ordered by repetitive element, MaLR.

Figure S7.

Normal, Non-Tumor Adjacent,, Tumor, Sperm variance in methylation levels ordered by repetitive element, Alu.

Figure S8.

Normal, Non-Tumor Adjacent, Tumor, Sperm variance in methylation levels ordered by repetitive element, ERV.

Figure S9.

Normal, Non-Tumor Adjacent, Tumor, Sperm variance in methylation levels ordered by repetitive element, SVA.

Figure S10.

Normal, Non-Tumor Adjacent, Tumor, Sperm variance in methylation levels ordered by repetitive element, L1P.

Figure S11.

Composite figure explaining the repeat CpG content using alignment plots and bin plots. The top portion, enveloped using the gray lines describes the specific repetitive element (alignment).

A) All sequences annotated by repeat masker as belonging to AluYb were filtered to output only those for which there is a probe on the microarray. From all these sequences, the ones that are shorter than the median of all sequences were further excluded to focus the subsequent analysis on longer (possibly full length) elements. The alignment of 60 of those sequences selected at random was then created using clustalw2 (1.0.11) program and standard parameters. The output of clustalw2 was subsequently parsed using a series of custom python and R scripts to create the figure. The nucleotides are colored to highlight the GC rich regions using the following scheme: G - yellow, C - green, A - light gray, T - pink

B)

Summary of the alignment in A), histogram indicates relative conservation of CpG at a given location. Per alignment locus, the number of CpGs was counted and normalized by the number of sequences aligned (i.e. 60)

C)

Summary of the alignment in A), (inverted) stacked histogram that summarizes a distribution of nucleotides at a given location. The bars are stacked from the most abundant nucleotide on the top, to the least abundant on the bottom. The colors are the same as in the part A). The most abundant nucleotide is marked in the middle of its bar.

The bottom portion characterizes the genomic context of the repetitive element family (bin plot).

The first 4 sub-plots, characterize all repetitive elements of a particular class in the human genome. The bins of plot D summarize the distribution of CpG counts in all sequences of all repetitive elements from a given lineage (the central bin marked in red) and 1,500 bases up~ and downstream from the repeat in 100 base increments per bin. The distribution of CpG in the repeat bin and external bins are presented in the form of a standard box and whisker plot, where the thick line inside the box indicates a median, the box is drawn around 25th and 75th percentiles, and the outliers are indicated as dots.

Plots E and F keep the binning structure of the sequence as in plot D, and show the average number of potential enzyme cuts among all the sequences per bin normalized to the size of the bin. Gray lines indicate the standard deviation.

Plot G is pertinent to the red bin of Plot D, it shows the distribution of sizes of all genomic repeats of a given family which were included in the central bin of plot D.

The plot H-K are analogous to D-G, with the only difference being that now the subset of repeats with a probe on the microarray are included in this analysis.

The plot L shows the repeat associated probe's location with respect to the repetitive element.

The supplementary figures S12 through S21 contains the bin plots for all families and subfamilies of probed repetitive elements discussed in the main text of the article. 54 alignment plots for these families are available in a zipped archive in the Supplementary Materials.

Figure S12.

CpG content summary for lineages of probed MaLR repetitive elements and 1,500 nucleotides up and down stream from the repeats. The flanking sequence is displayed in bins of 100 bases long.

Figure S13.

CpG content summary for lineages of probed Alu repetitive elements and 1,500 nucleotides up and down stream from the repeats. The flanking sequence is displayed in bins of 100 bases long.

Figure S14.

CpG content summary for lineages of probed ERV repetitive elements and 1,500 nucleotides up and down stream from the repeats. The flanking sequence is displayed in bins of 100 bases long.

Figure S15.

CpG content summary for lineages of probed SVA repetitive elements and 1,500 nucleotides up and down stream from the repeats. The flanking sequence is displayed in bins of 100 bases long.

Figure S16.

CpG content summary for lineages of probed L1P repetitive elements and 1,500 nucleotides up and down stream from the repeats. The flanking sequence is displayed in bins of 100 bases long.

Figure S17.

CpG content side-by-side comparison of an older versus younger lineage within repetitive family, MaLR: MLT1C - MST1A

Figure S18.

CpG content side-by-side comparison of an older versus younger lineage within repetitive family, SVA: SVA_F - SVA_B

Figure S19.

CpG content side-by-side comparison of an older versus younger lineage within repetitive family, L1P: L1PA17 - L1PA4

Figure S20.

CpG content side-by-side comparison of an older versus younger lineage within repetitive family, L1P: L1PA3 - L1HS

Figure S21.

CpG content side-by-side comparison of an older versus younger lineage within repetitive family, Alu: AluYd8 - AluYb9

Figure S22.

ASCIIMAPs showing the location of probes interrogating MaLR elements.

The text provides a summary of a region using ASCII characters (generated using a tool ASCIIMap). The ASCIIMap tracks show the location of the probe (●). The location of a CpG island is marked underneath (*) as are the locations of repetitive elements in the area (**1**-Line1, **2**-Line2, **T**-LTR(MaLR), **S**- SINE, **=**-other, **A**-Alu, etc.). The vertical bars (|) indicate the presence of an enzyme recognition site for Acil, HhaI and McrBC enzymes respectively. The resolution of 1 character is about 15 nucleotides. The regions shown are approximately 1kb in length. The genomic coordinates of each regions are given using HG16 build of the human genome. Note that in some cases a probe appears to map inside a region designated as a repetitive element. The probe design algorithm always ensures that the sequence is unique before designing a probe.

The apparent paradox that a repetitive element may have parts that are unique sequences can be explained by considering the age of the repetitive elements for which the probe is designed. For example, an element of the family MLT1C, 85 MYO: over a span of millions of years since it appeared in its original form in the genome, its sequence have deteriorated from its consensus so much that although the element can still be classified as MLT1C now (based on the overall structure and certain sequence patterns), its sequences acquired enough random mutations that the probe algorithm can recognize certain parts within this MLT1C as unique in the genome. This can be best illustrated in the Supplementary Figures S12 through S16, in the “repeat associated probes per bin” panels. For repetitive element families that are younger, i.e. the elements that haven’t had evolutionary time to acquire mutations differentiating them from their respective consensus, the probe designer most likely designs the probe within the 100 bases flanking region of the repetitive element. Conversely, for the older repetitive elements (20, 30, 40+ MYO), the probe designer is able to find regions that have uniquely diverged from the global consensus of the repeat family. The archive of alignments of repetitive elements within their families included in the Supplementary Materials further illustrates how similar or divergent the sequences are within a lineage of repetitive elements.

Figure S1.

Probe signal as a function of distance to enzyme cut PmeI

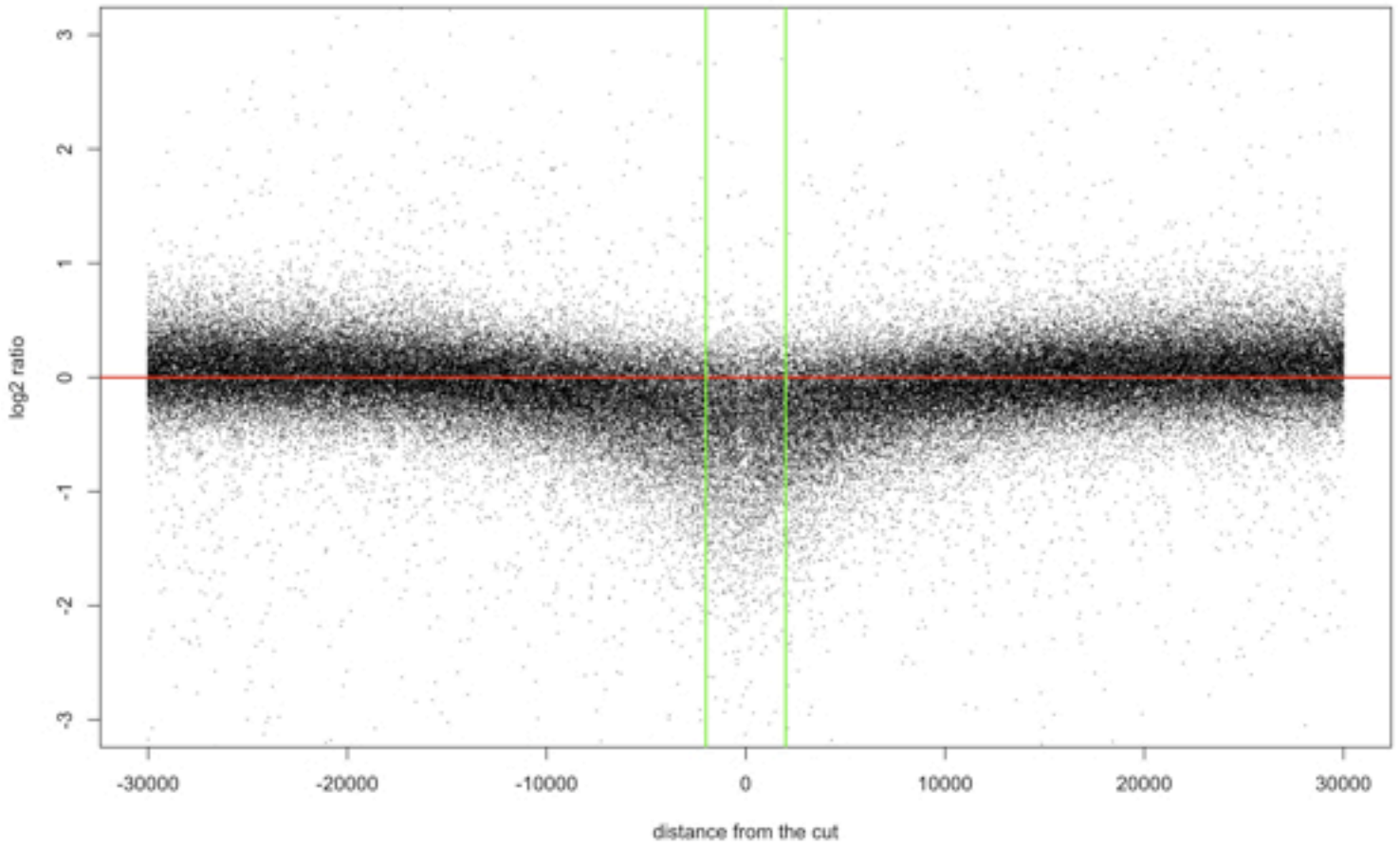


Figure S3A.

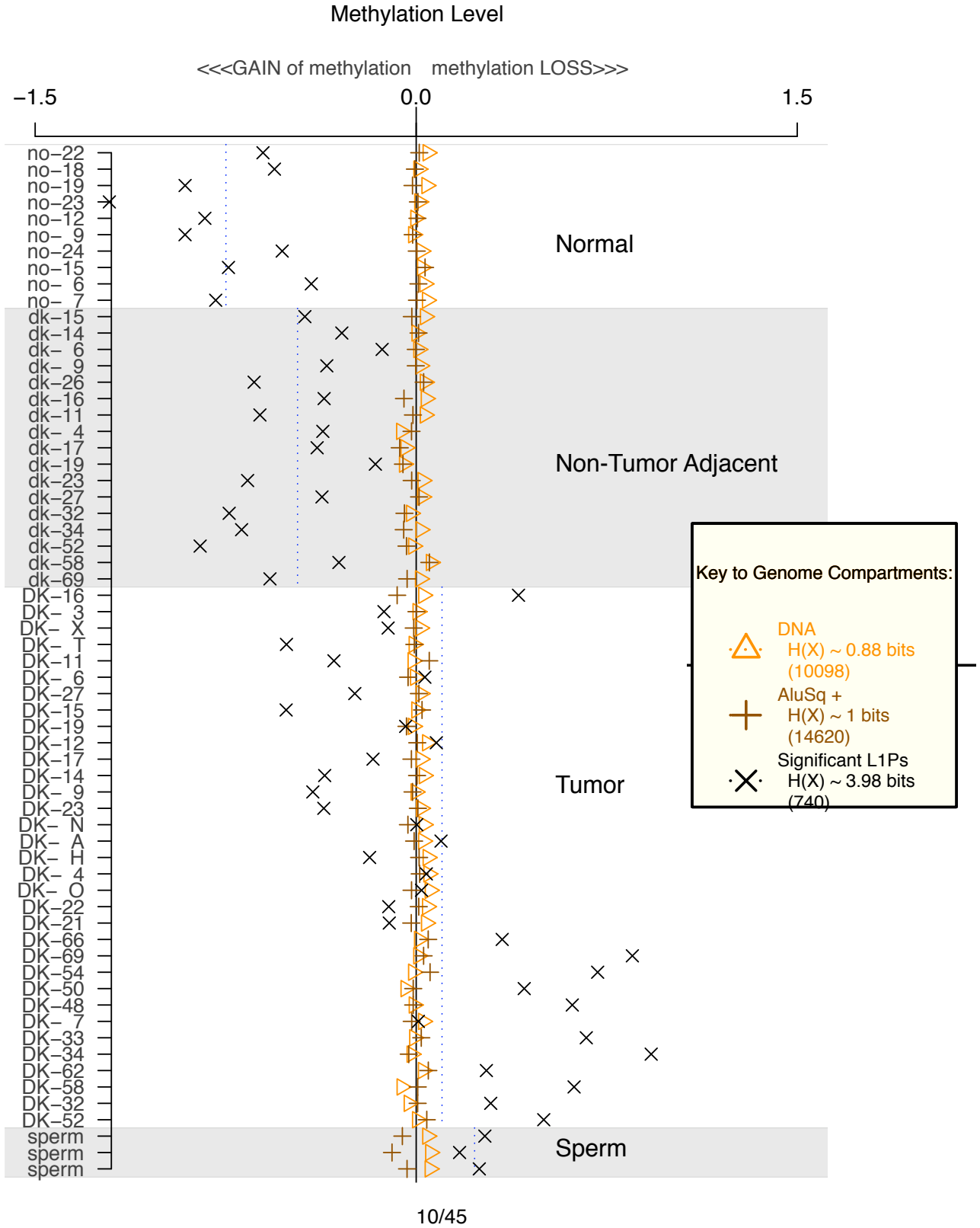


Figure S4A.

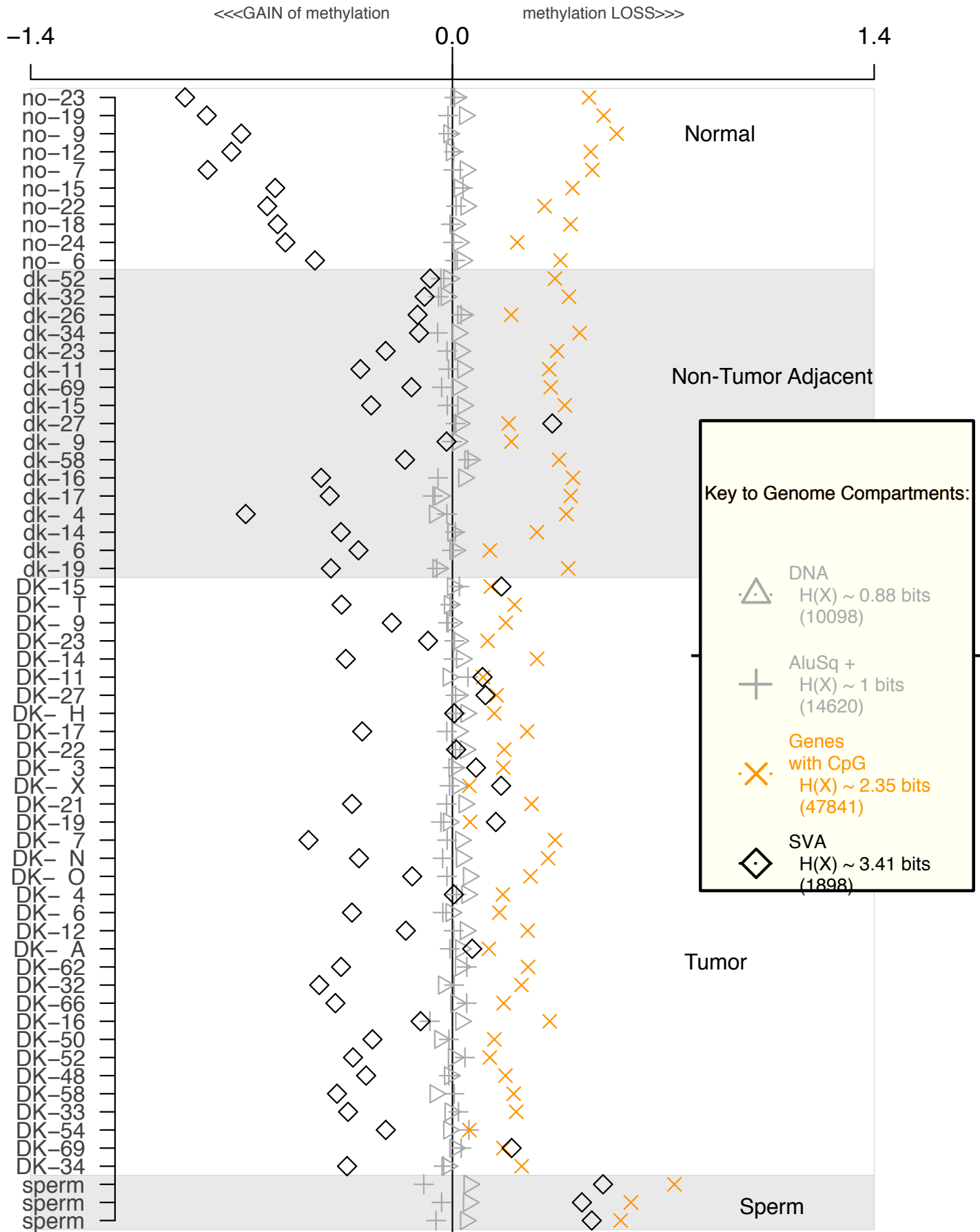


Figure S4B.

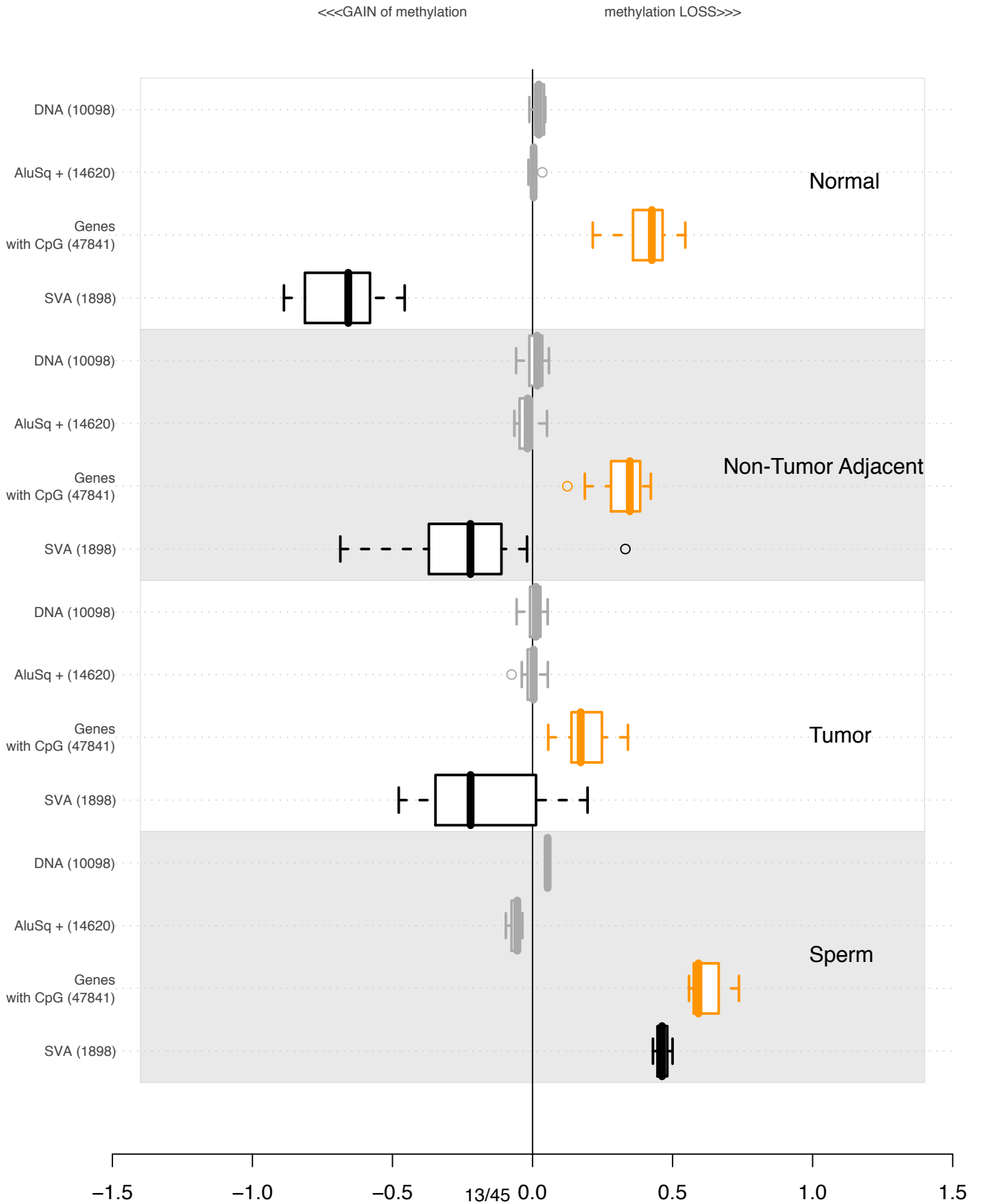


Figure S5.

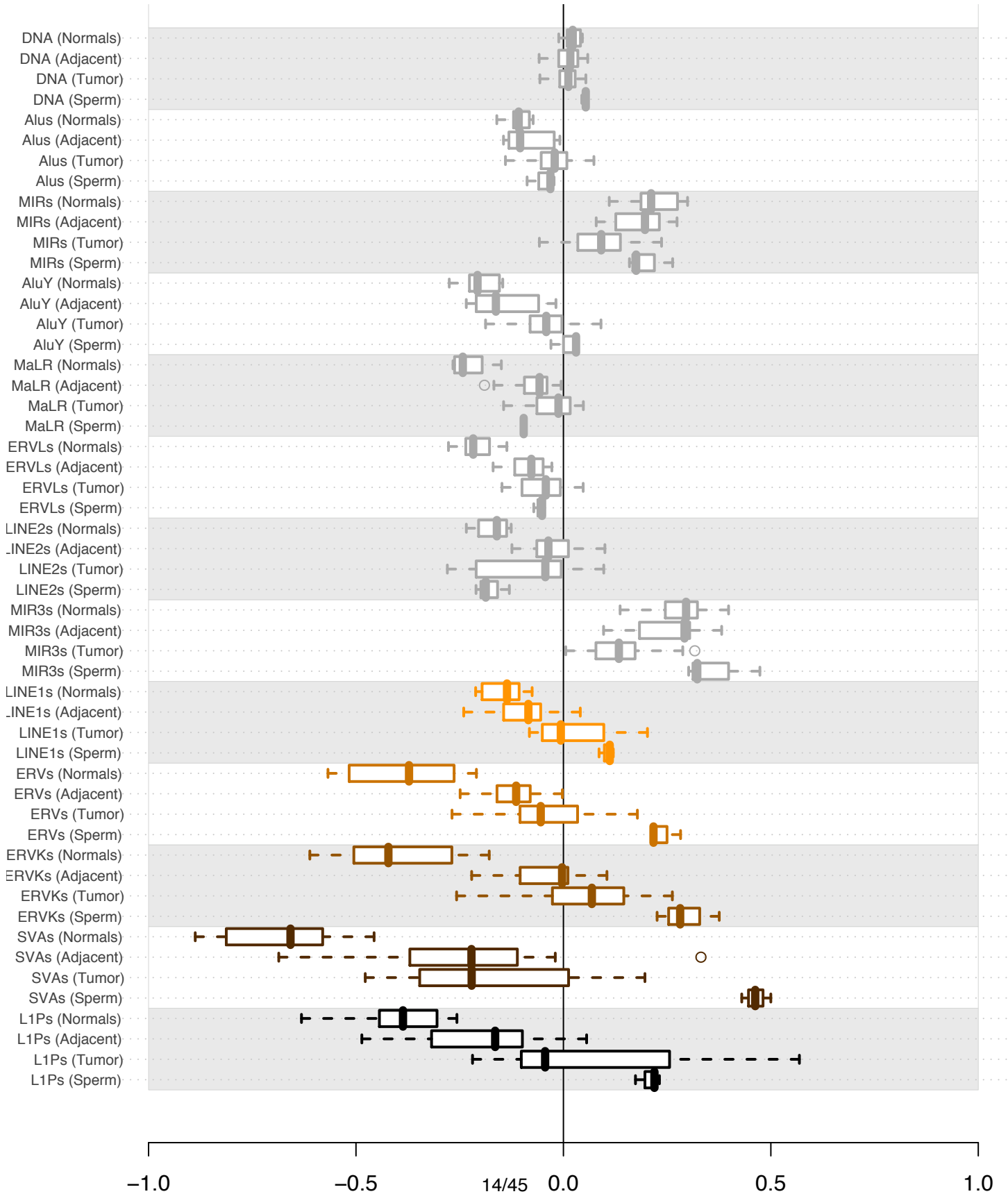


Figure S6.

<<<GAIN of methylation

methylation LOSS>>>

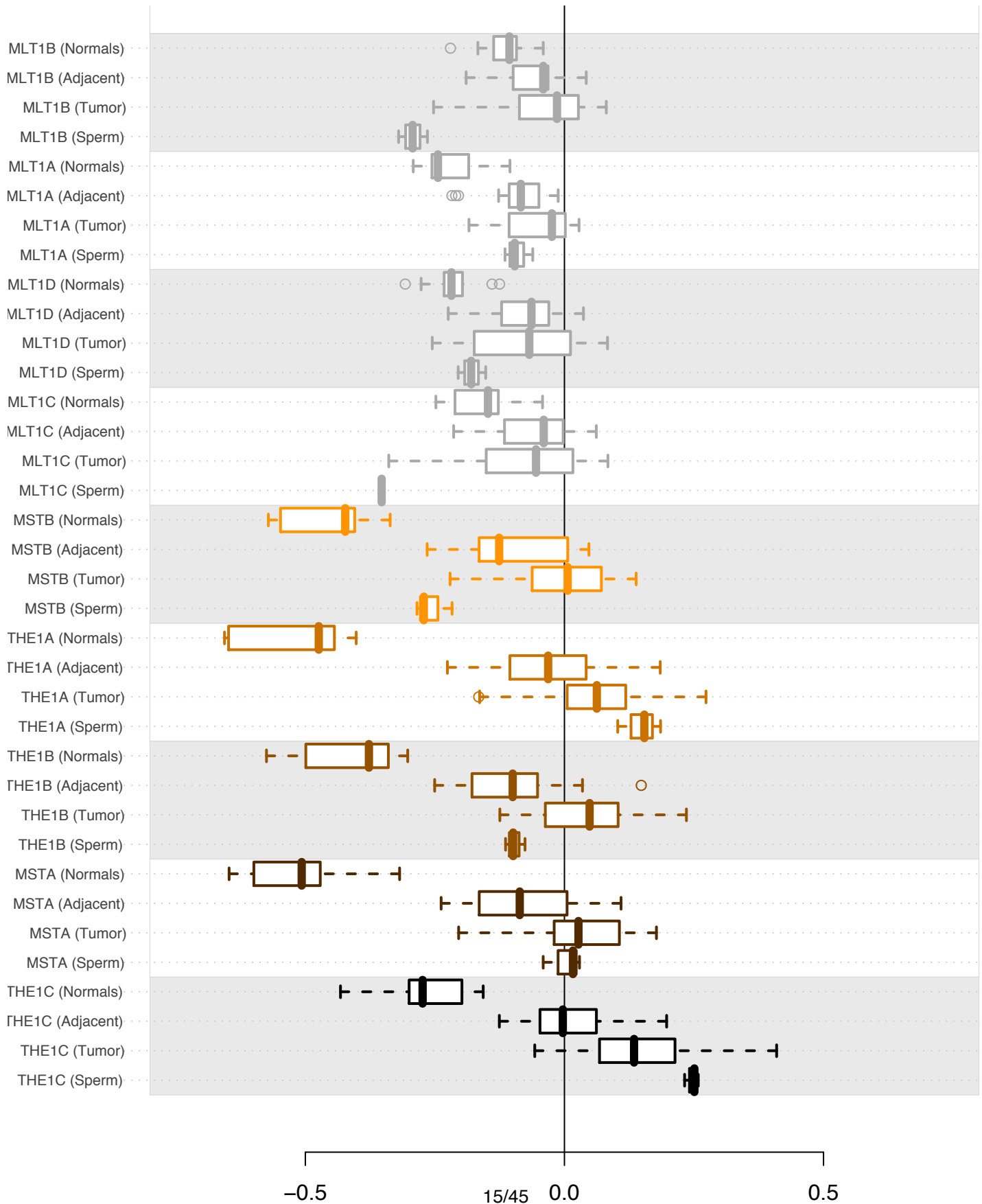


Figure S7.

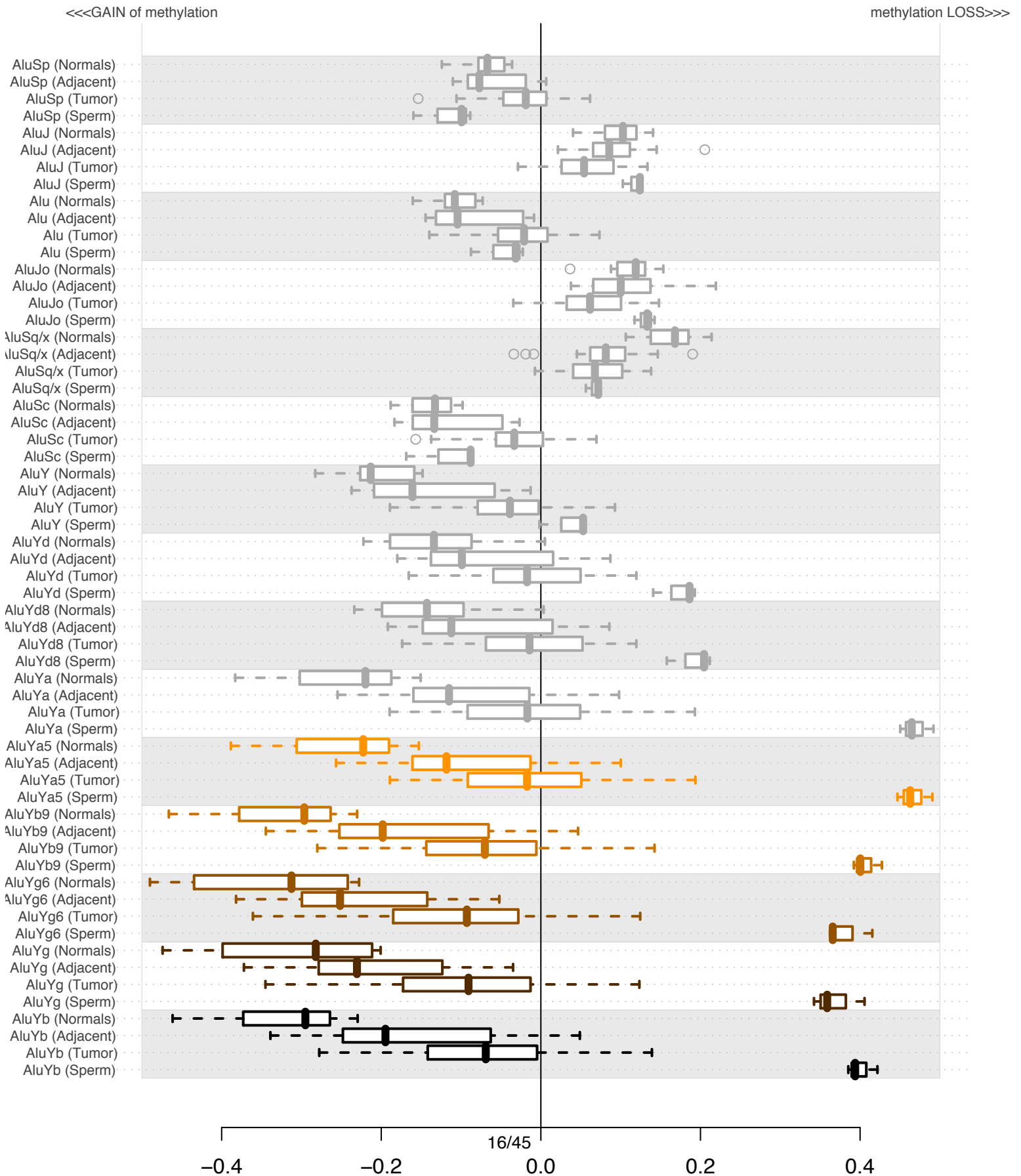


Figure S8.

<<<GAIN of methylation methylation LOSS>>>

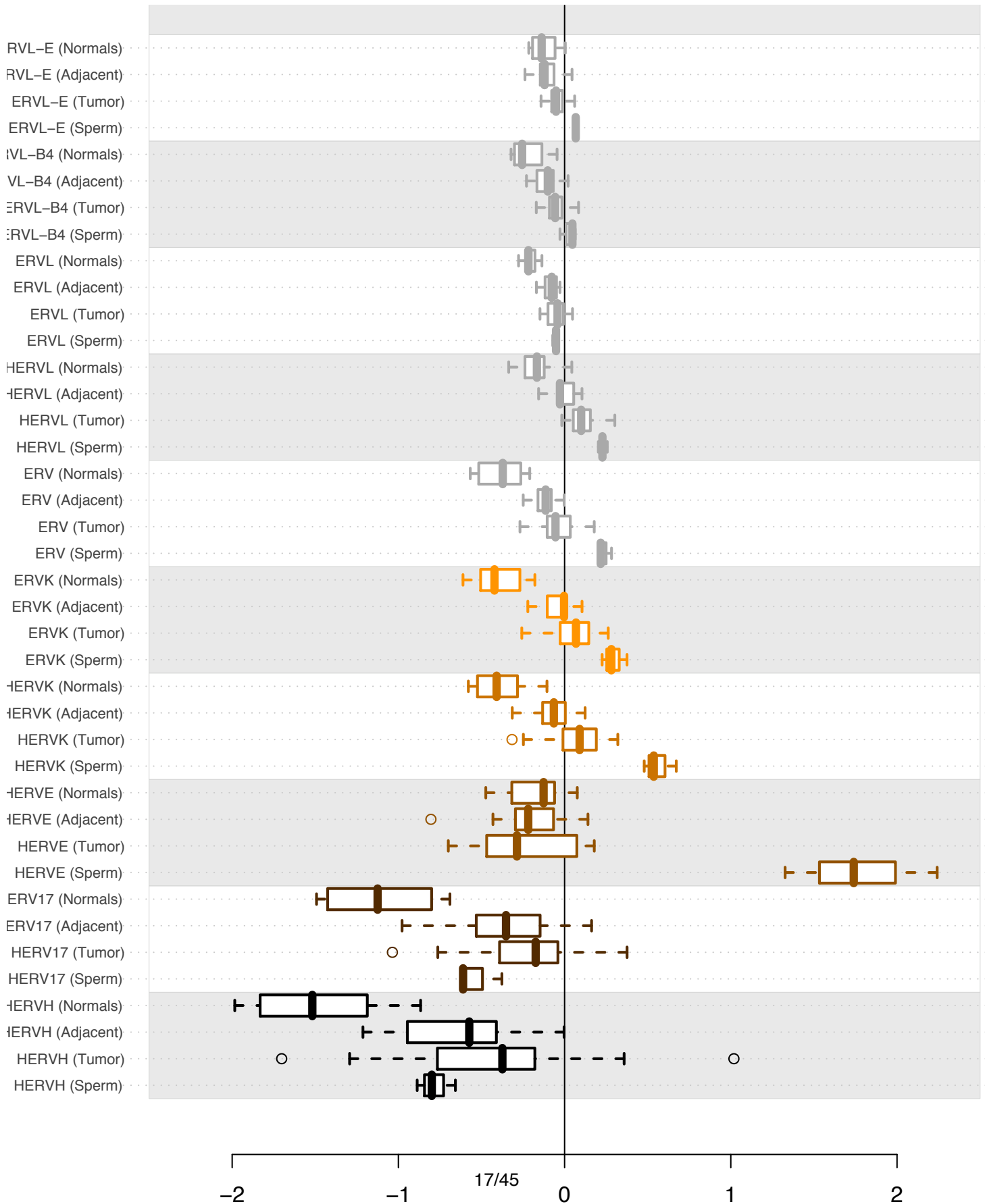


Figure S9.

<<<GAIN of methylation methylation LOSS>>>

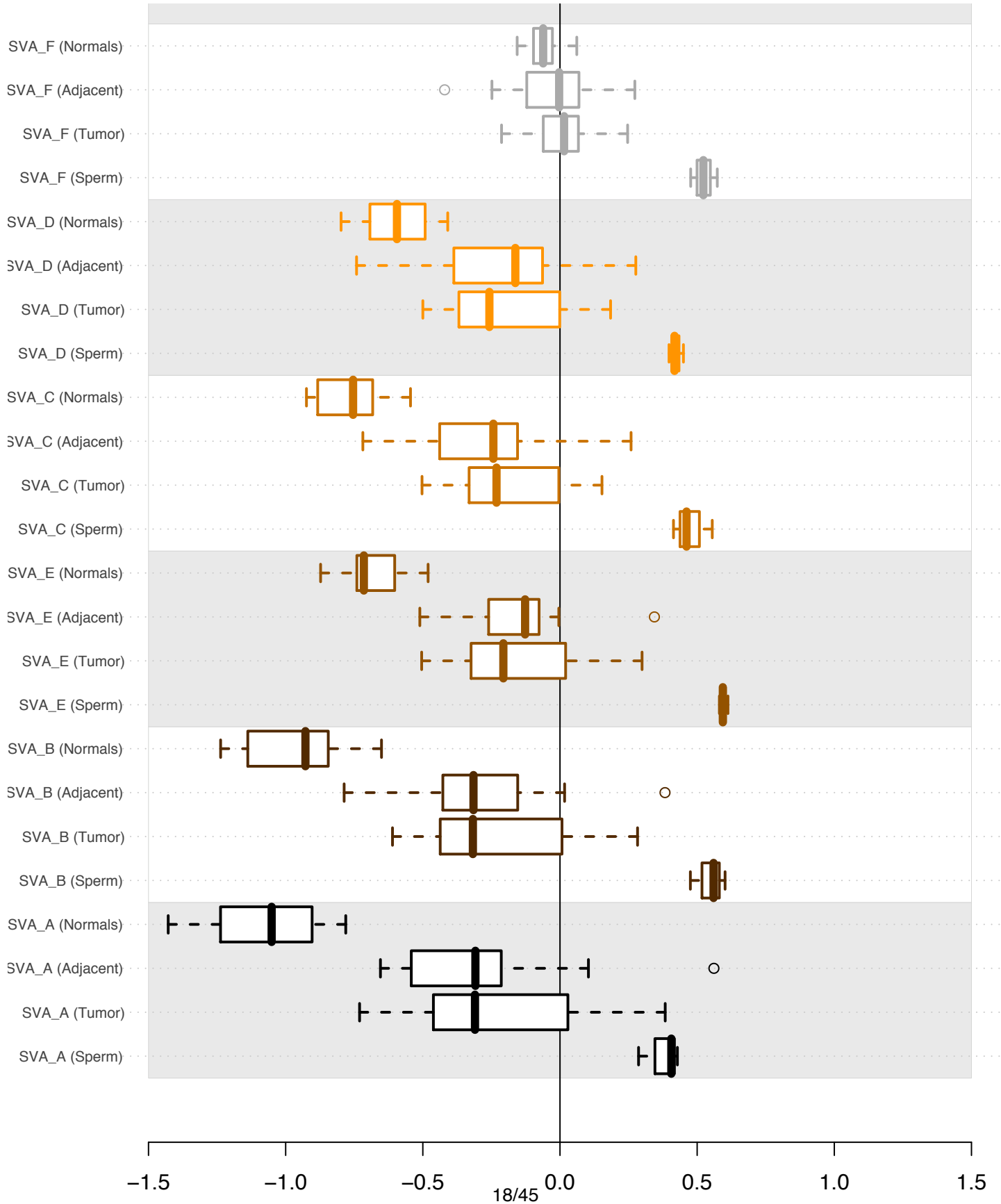
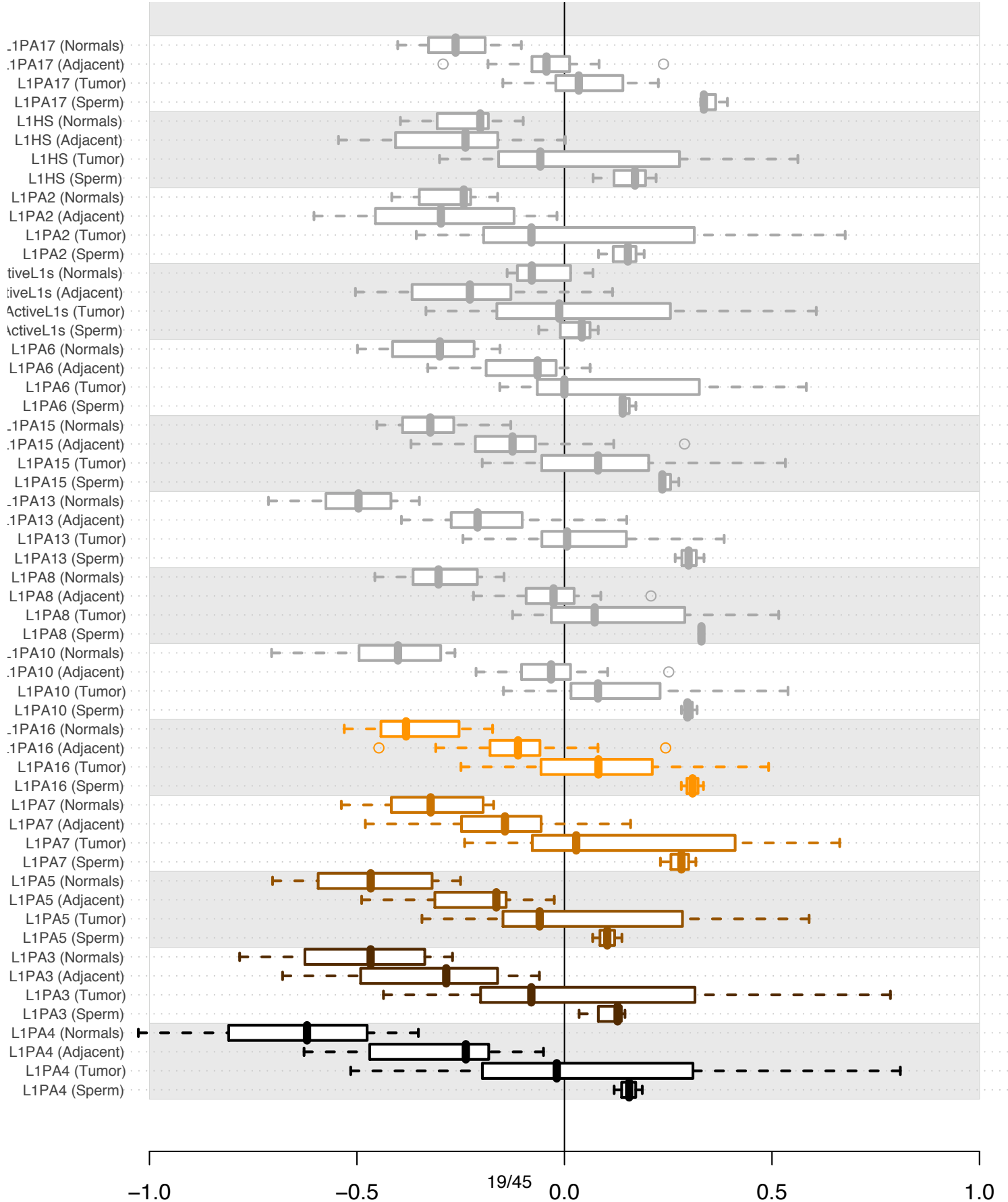


Figure S10.

<<<GAIN of methylation

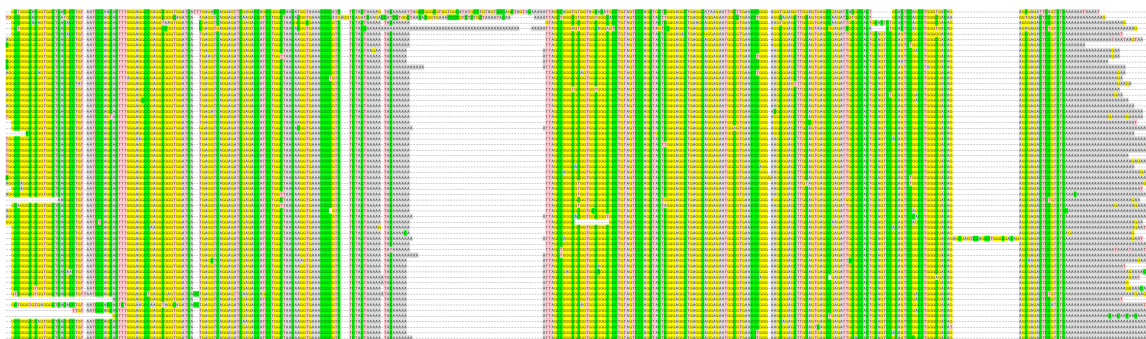
methylation LOSS>>>



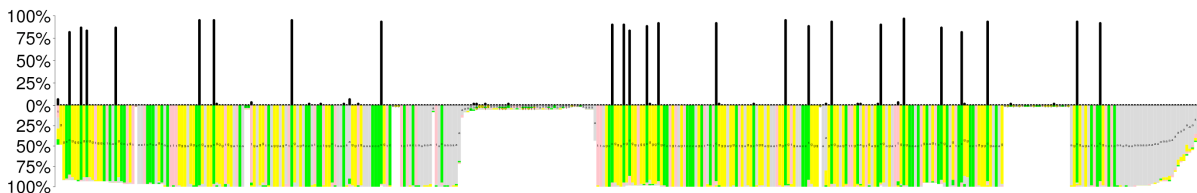
alignment, CG content and conservation of a random 60 AluYb(s)

Figure S11

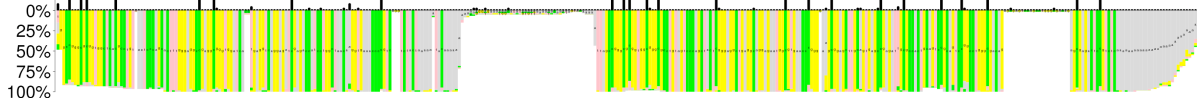
A



B

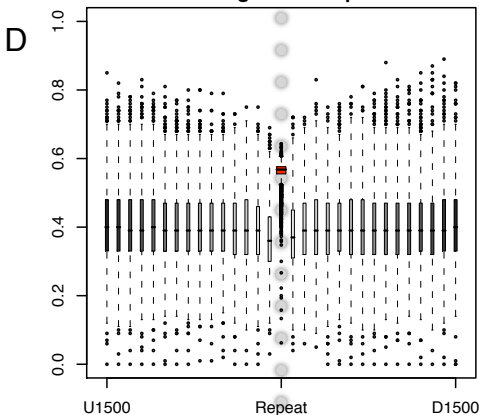


C

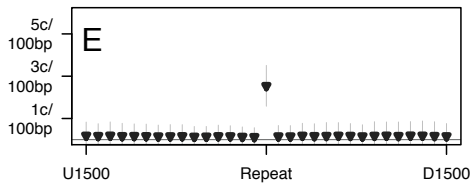


ALL GENOMIC

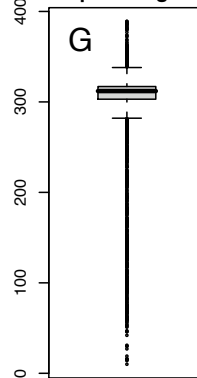
CpG content relative to binsize
3365 all genomic repeats.



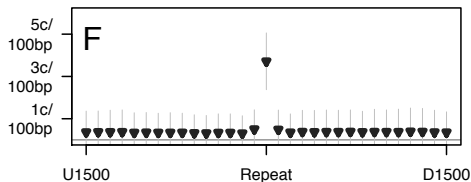
all potential Acil + HHal cuts per 100 base



Distribution of repeat lengths



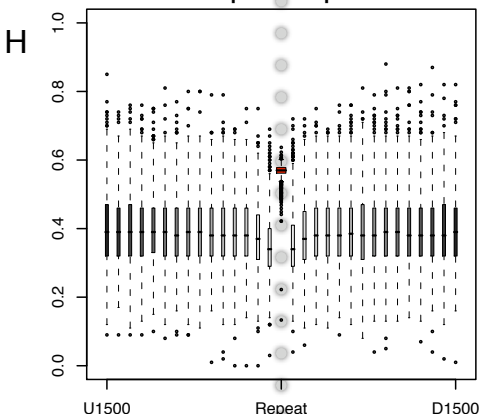
all potential McrBC cuts per 100 bases



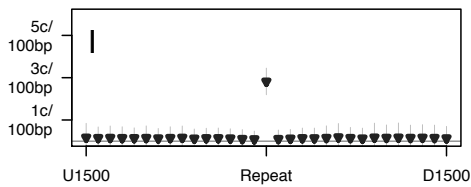
AluYb

PROBED

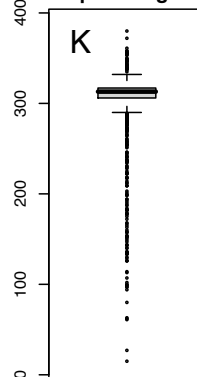
CpG content relative to binsize
1414 probed repeats.



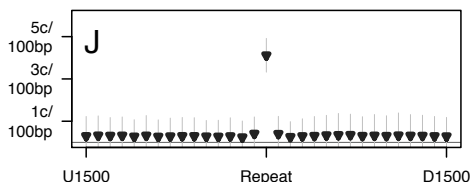
all potential Acil + HHal cuts per 100 base



Distribution of repeat lengths



all potential McrBC cuts per 100 bases



Repeat associated probes per bin

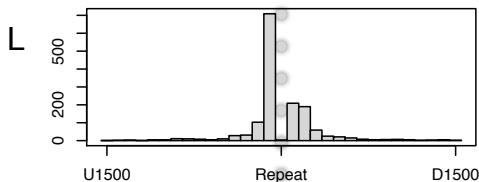
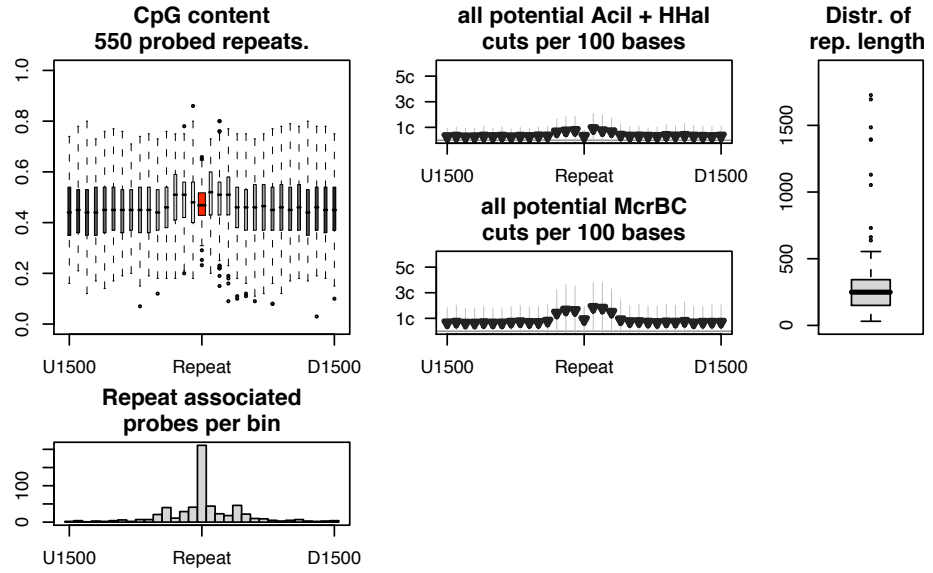
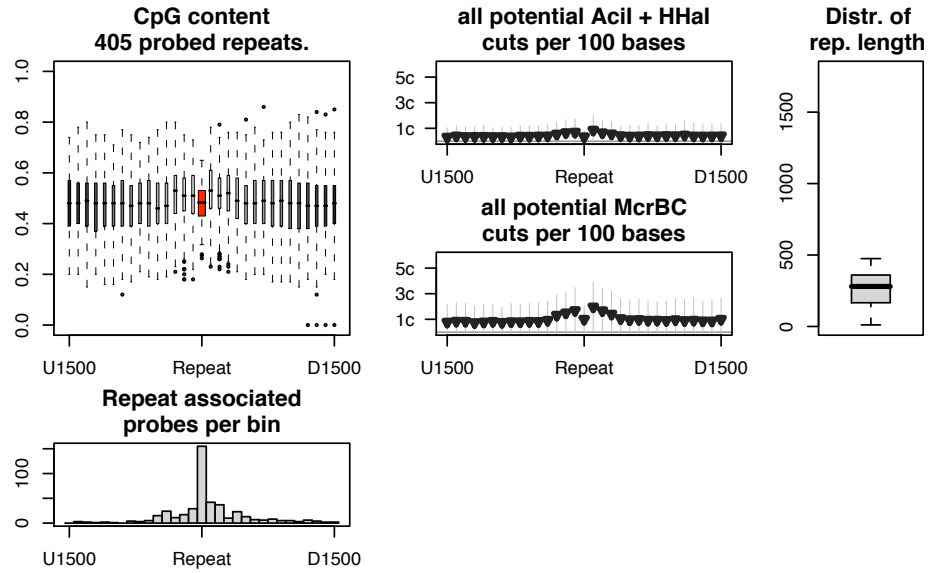


Figure S12

MLT1A
~77 MYO



MLT1B
~80 MYO



MLT1C
~85 MYO

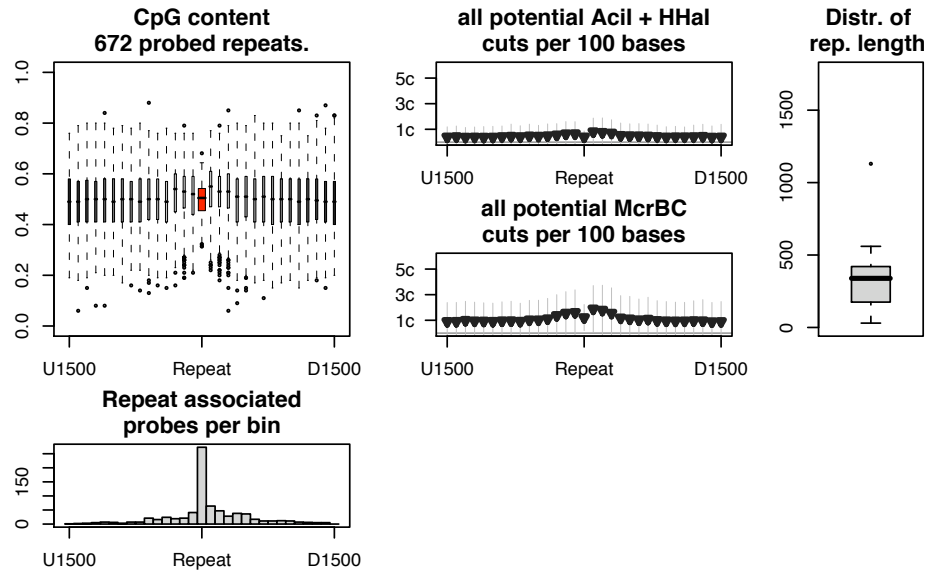
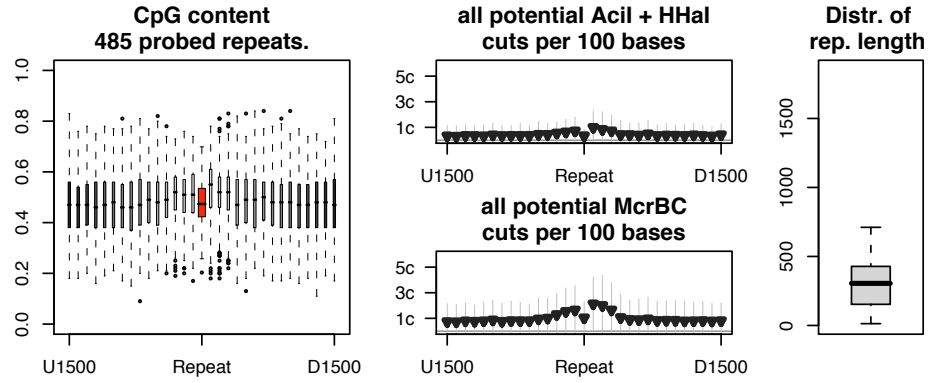
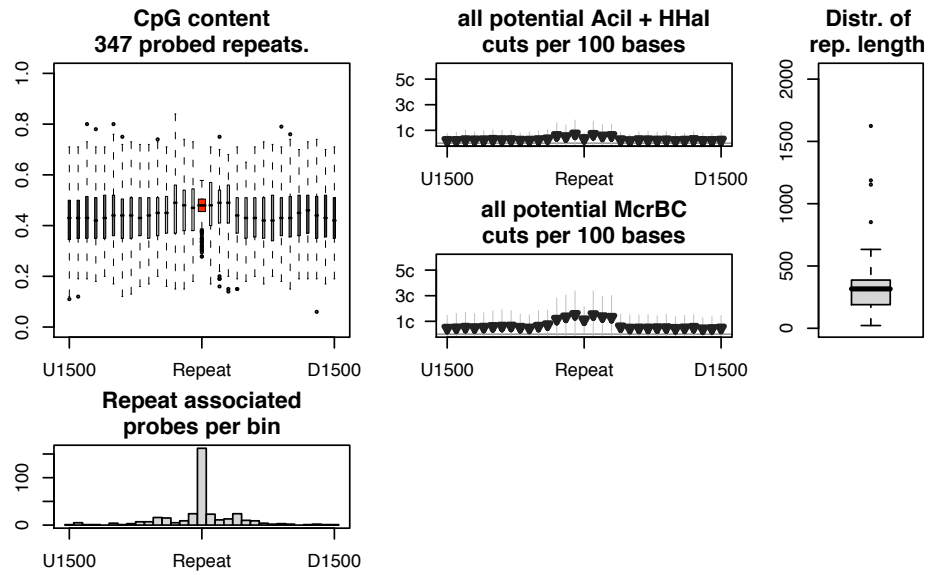


Figure S12

MLT1D
~100 MYO



MSTA
~60 MYO



MSTB
~75 MYO

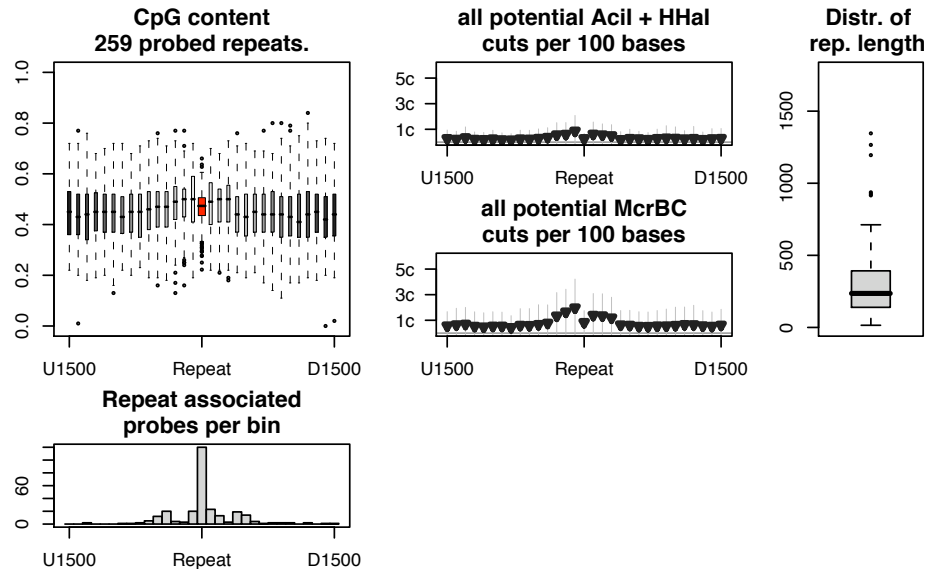
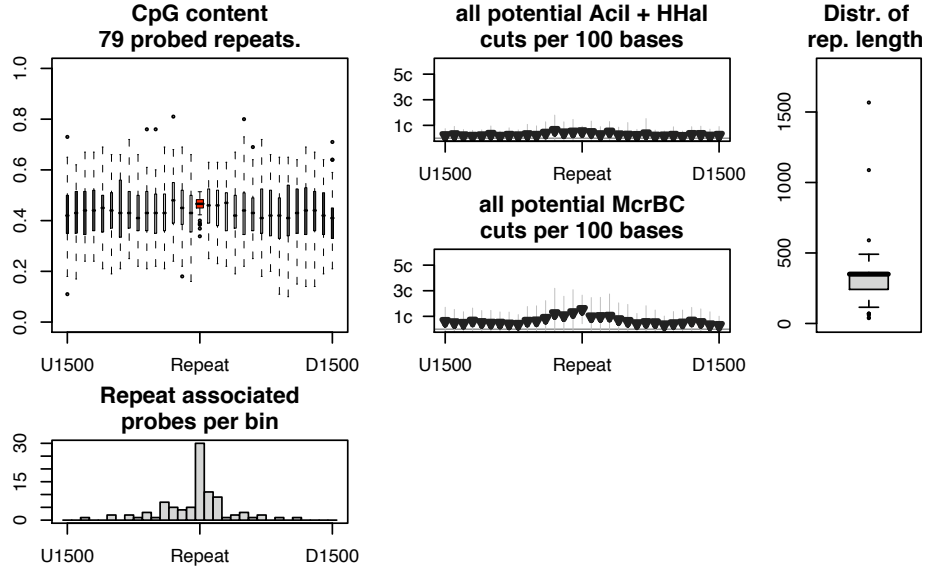
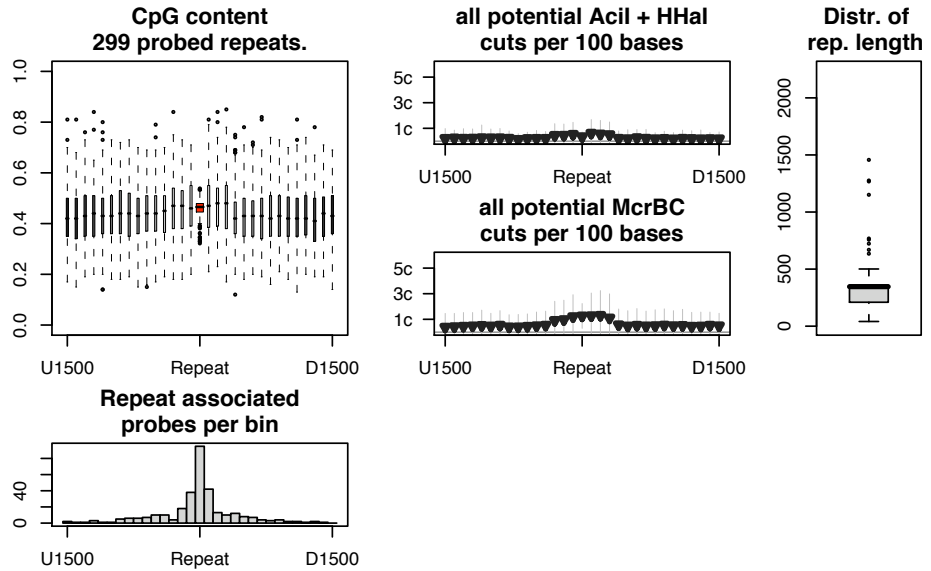


Figure S12

THE1A
~45 MYO



THE1B
~50 MYO



THE1C
~55 MYO

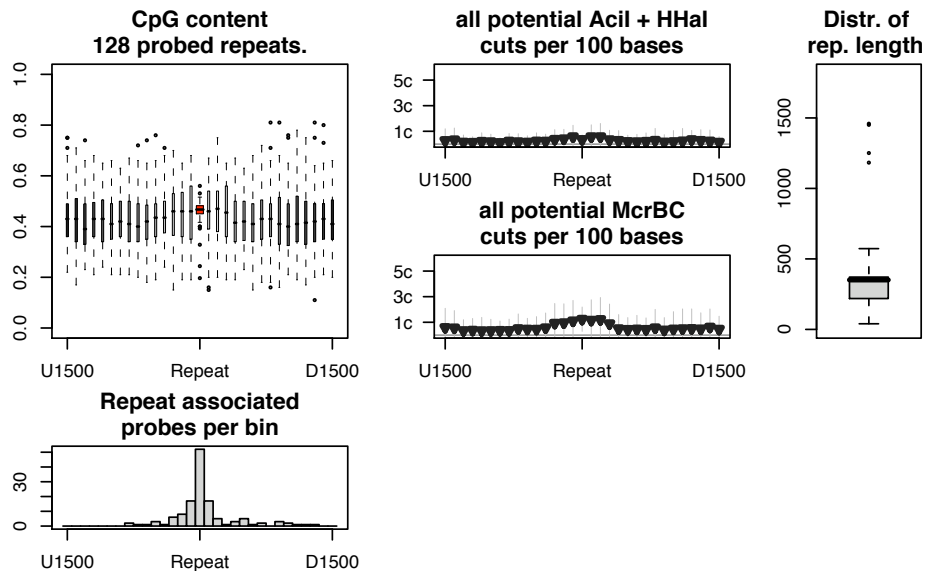
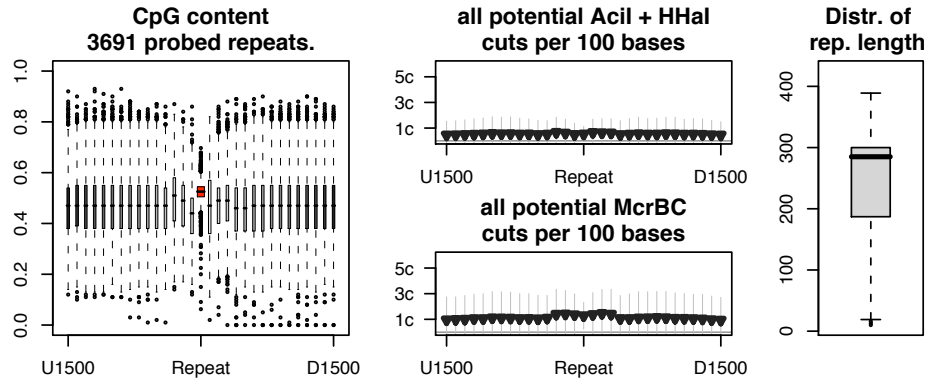
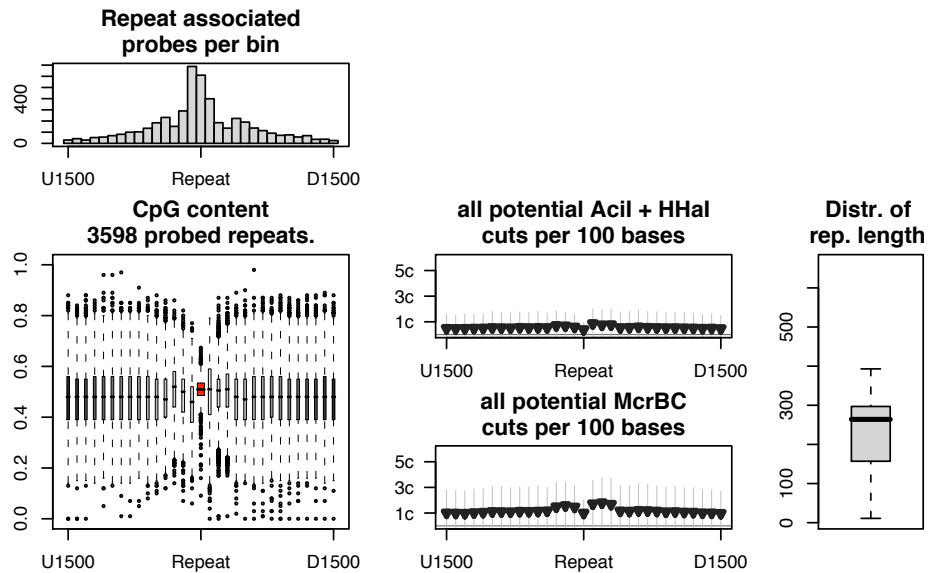


Figure S13

AluJb
~61 MYO



AluJo
~61 MYO



AluSc
~29 MYO

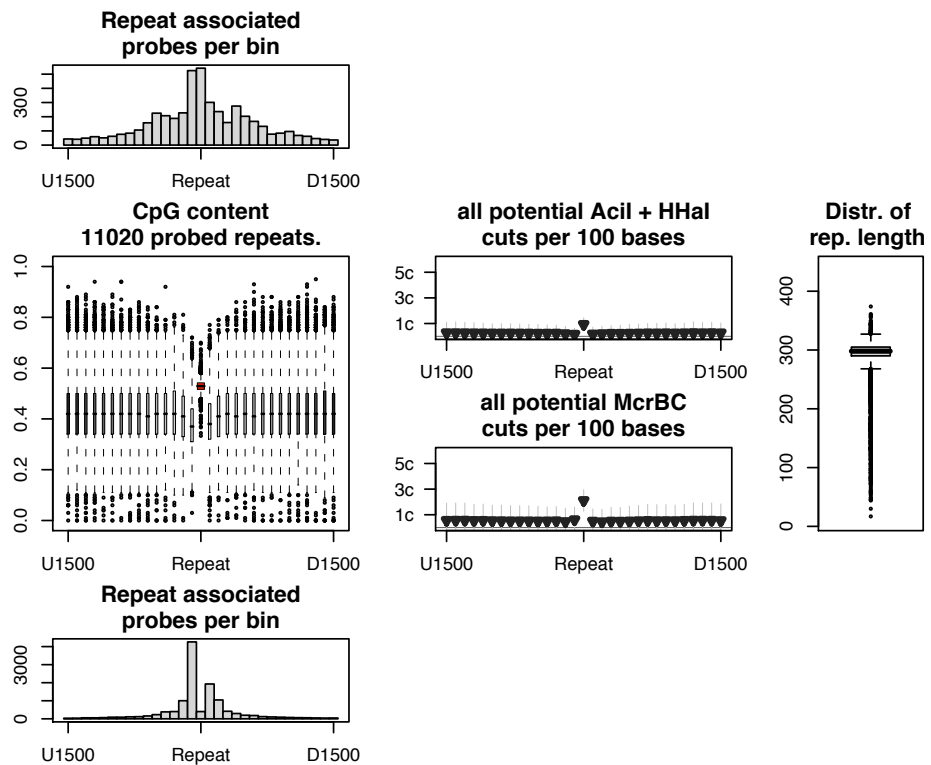
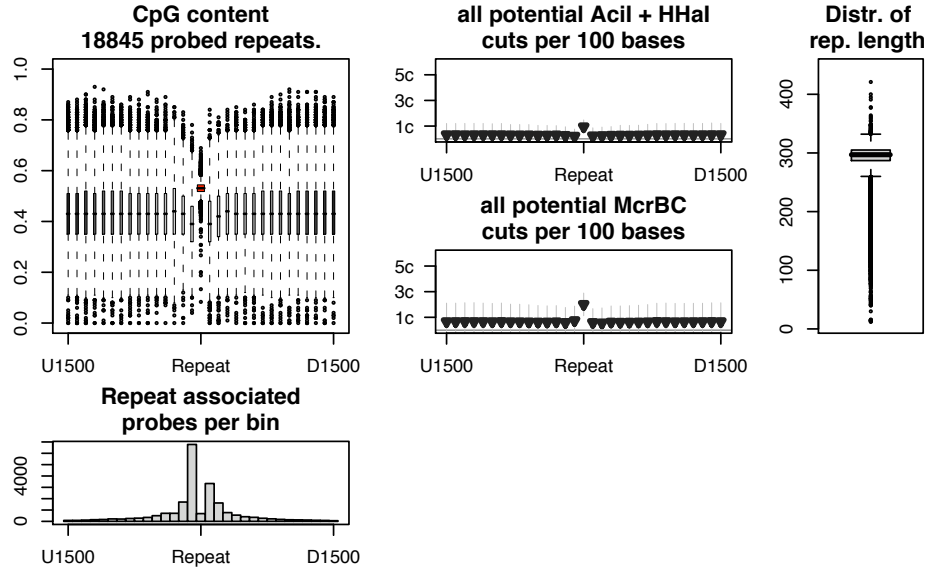
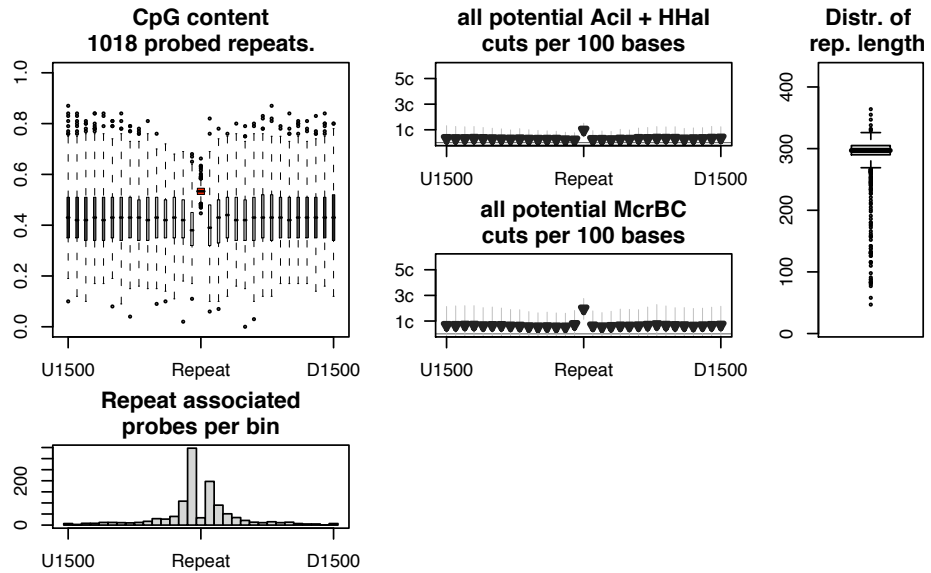


Figure S13

AluSg
~32 MYO



AluSg1
~30 MYO



AluSp
~33 MYO

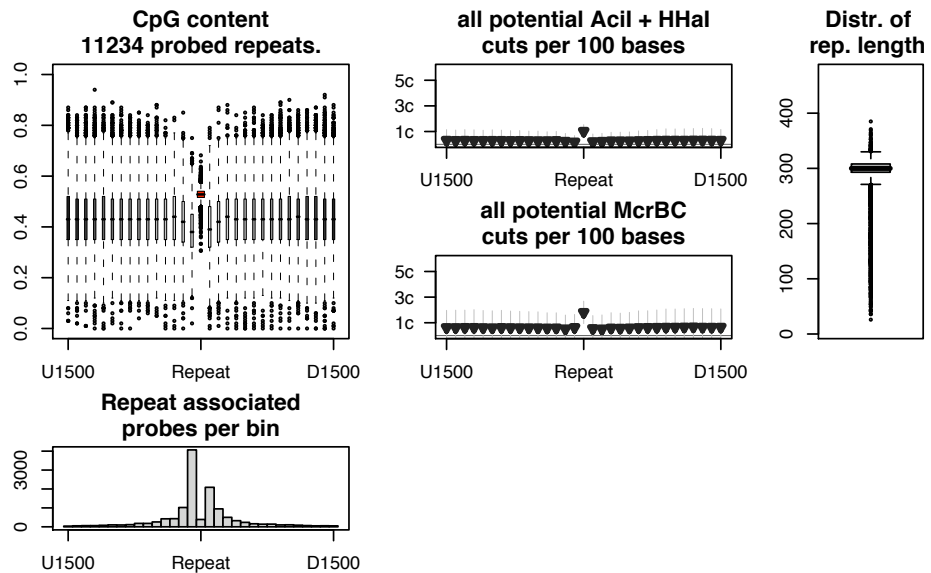
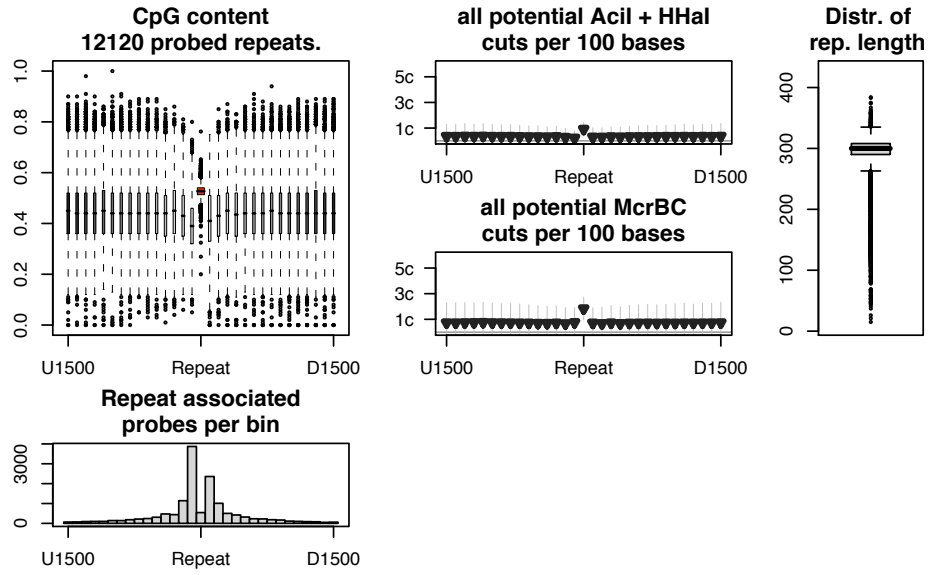
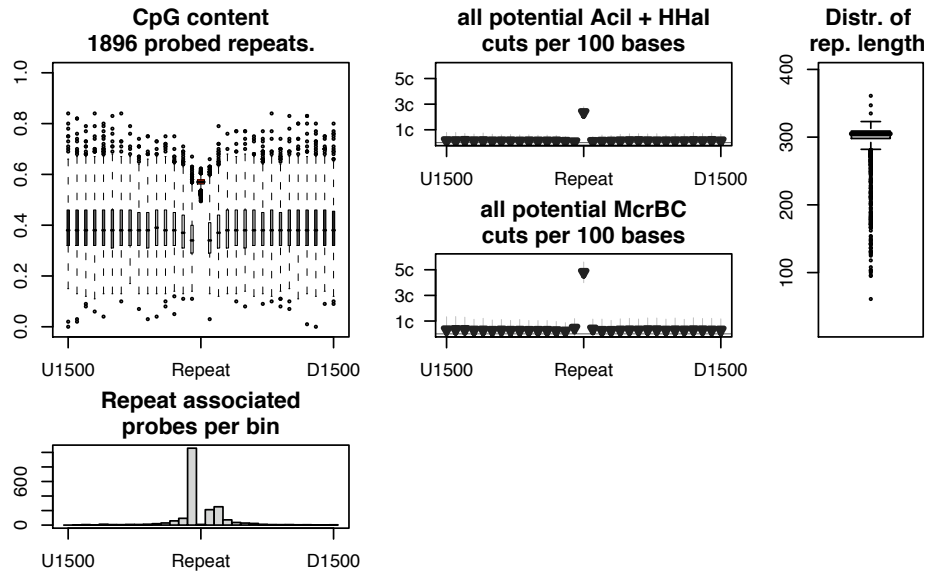


Figure S13

AluSq
~33 MYO



AluYa
~11 MYO



AluYa5
~6.5 MYO

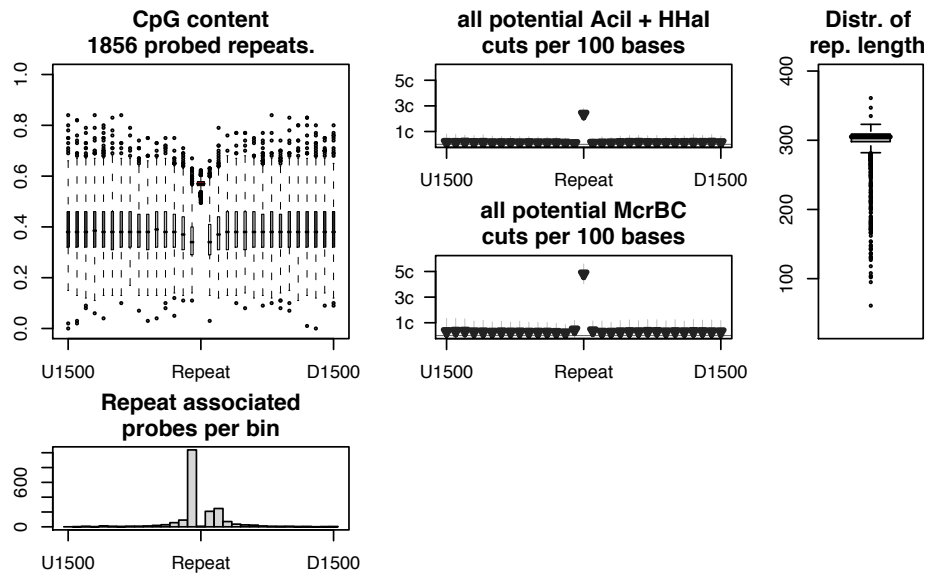


Figure S13.

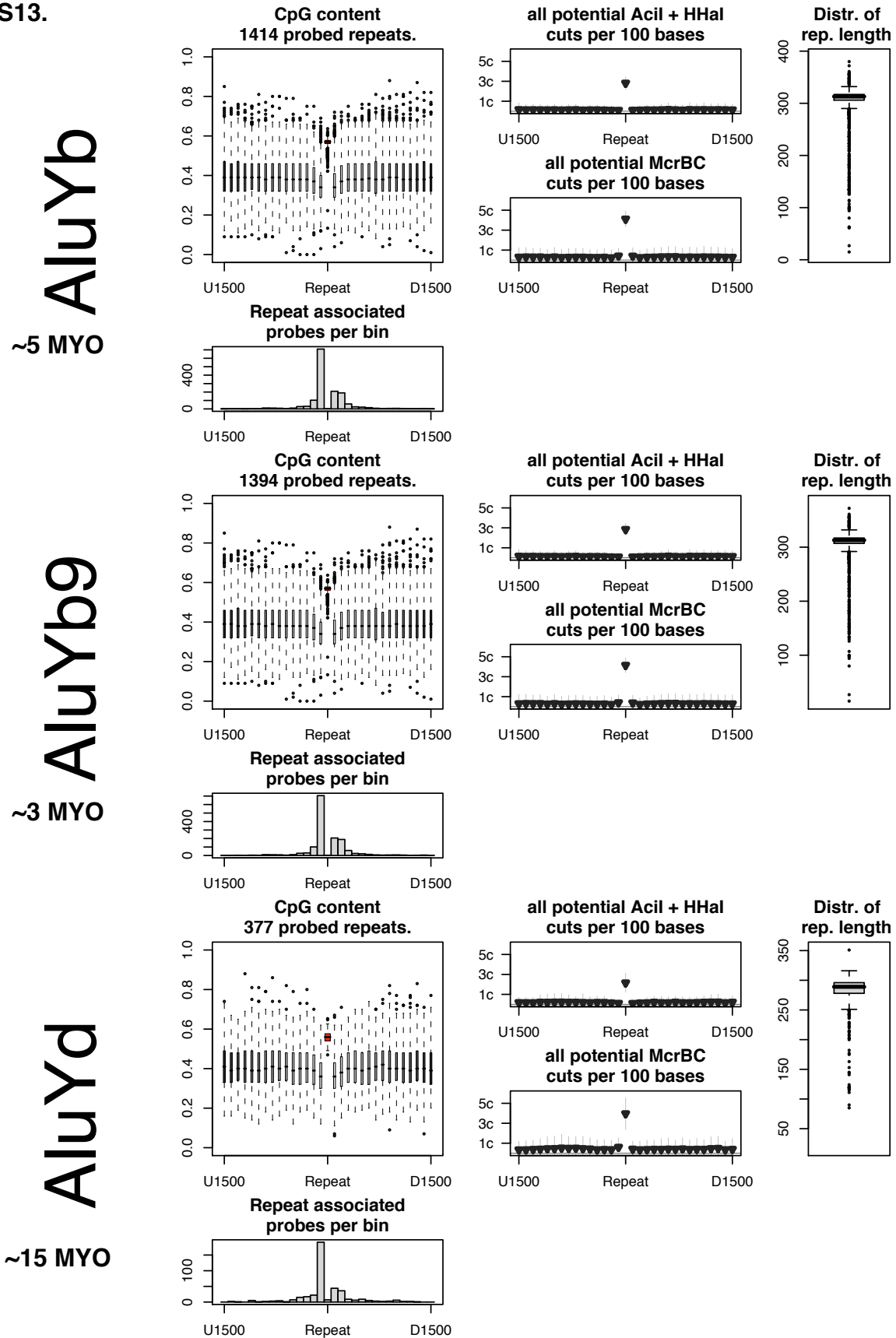
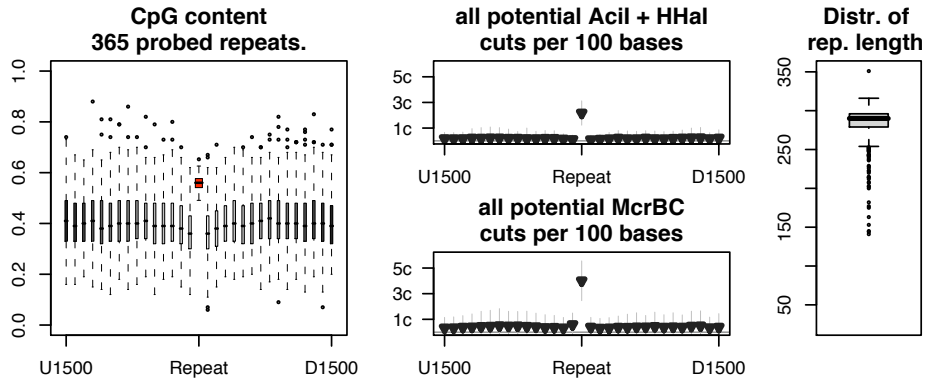
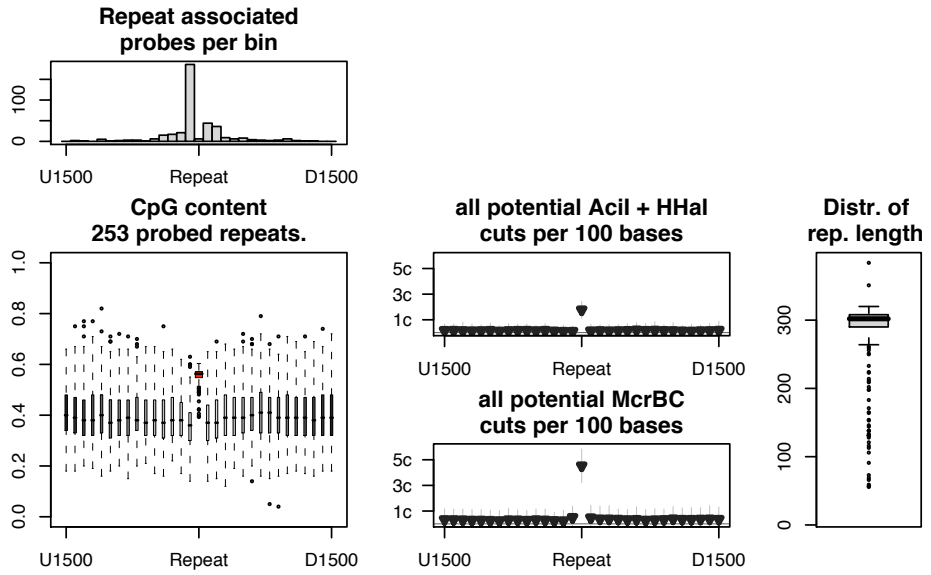


Figure S13

AluYd8
~15 MYO



AluYg
~5 MYO



AluYg6
~3 MYO

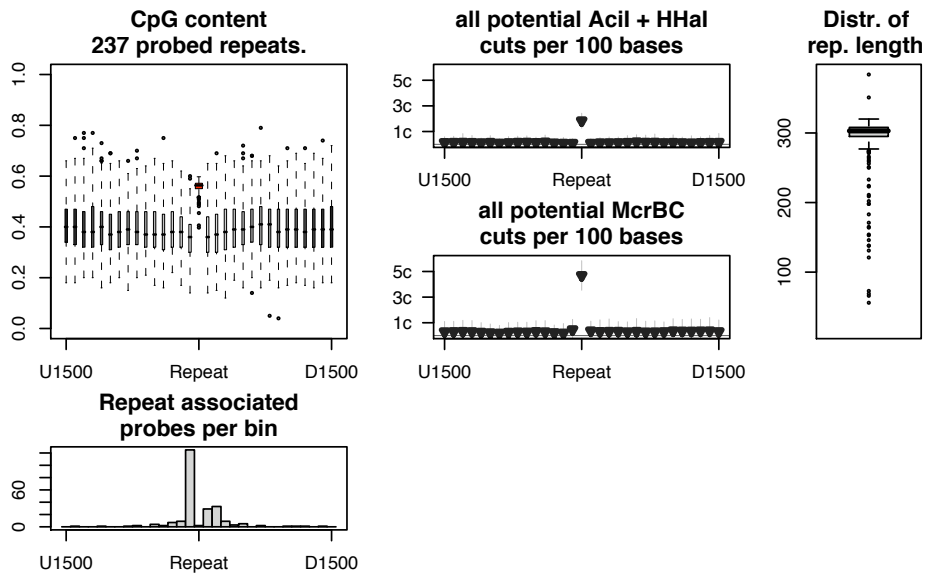
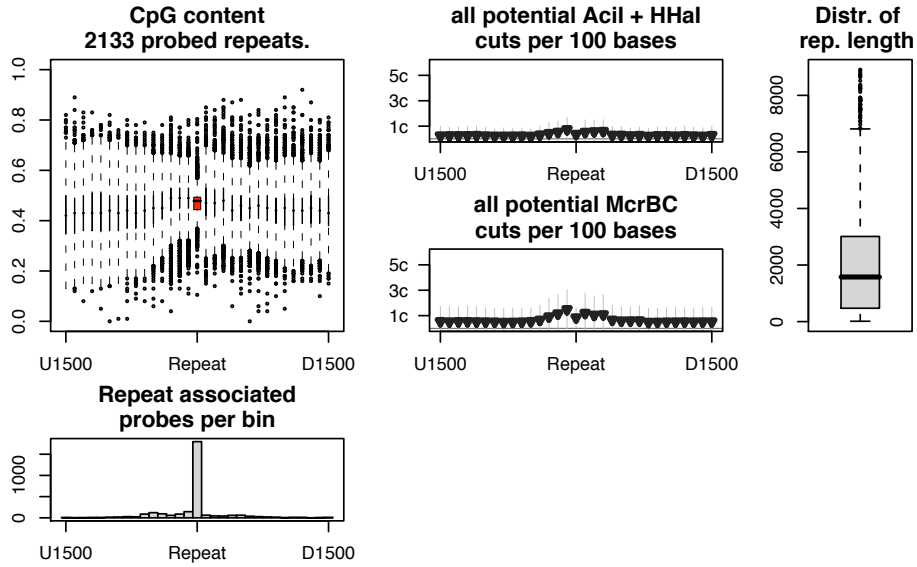
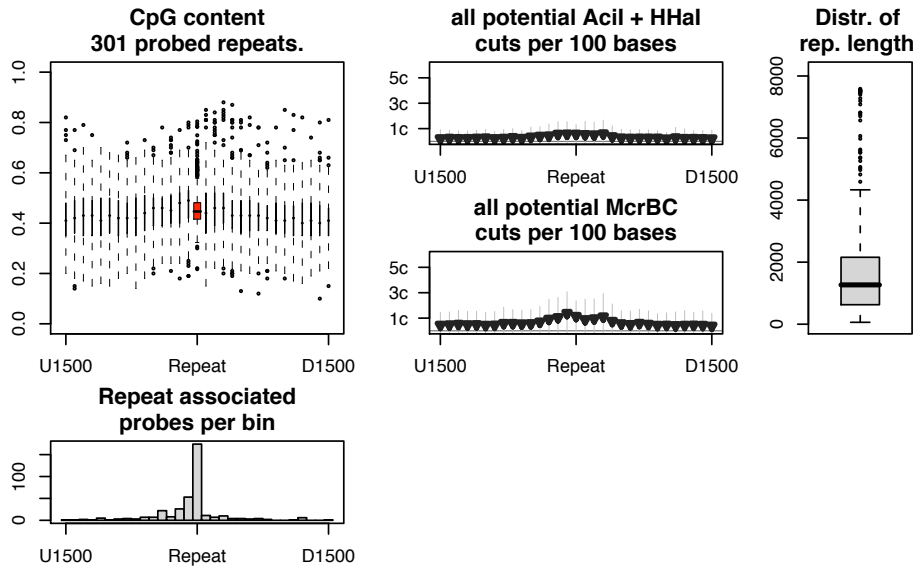


Figure S14

ERV



ERVK



ERVL-B4

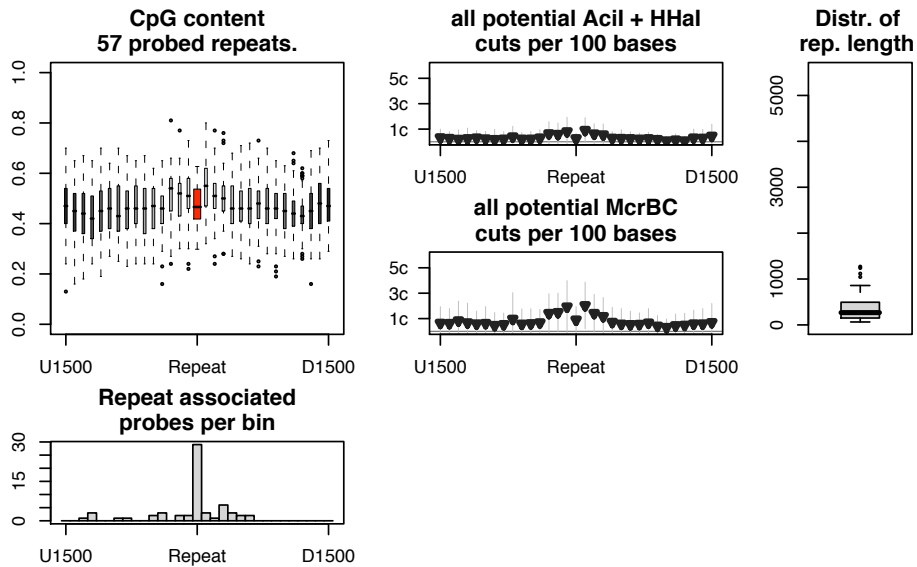
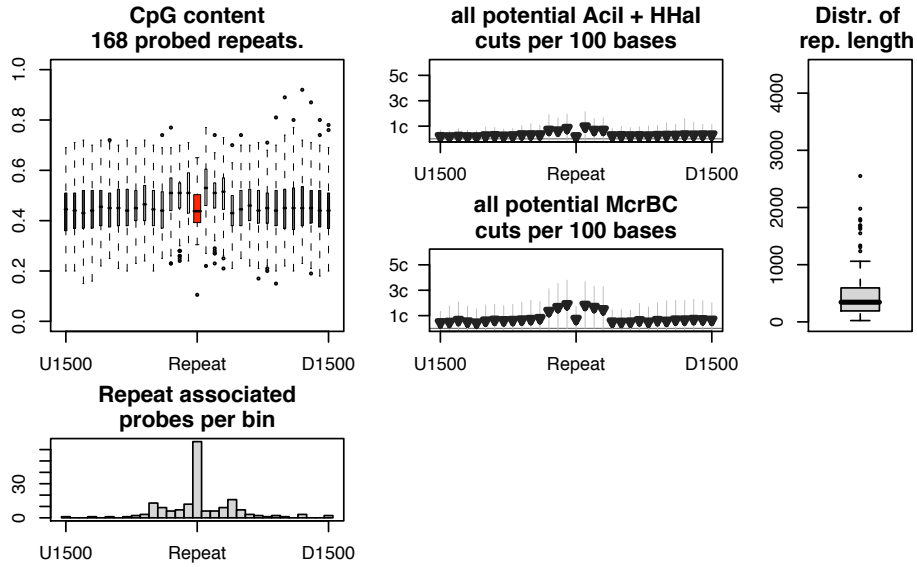
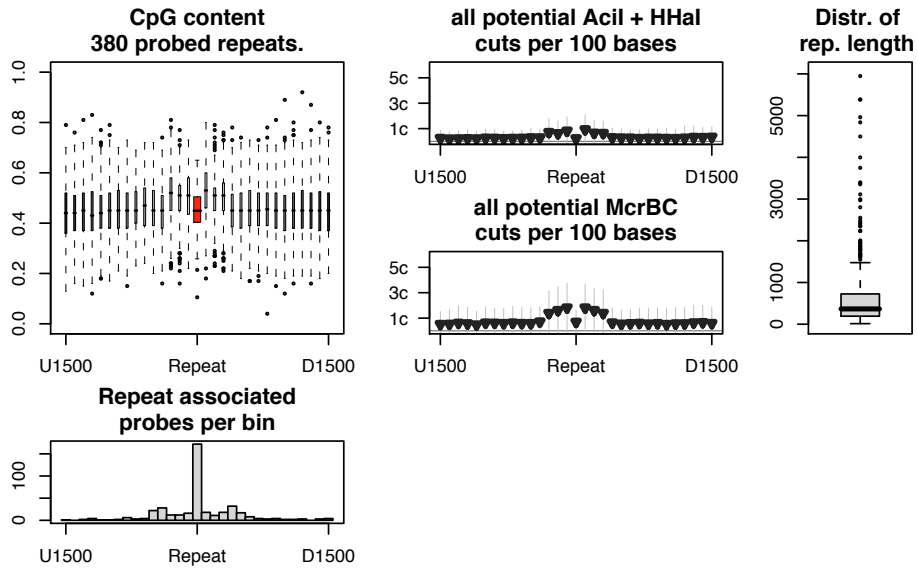


Figure S14

ERV1-E



ERV1



HERV17

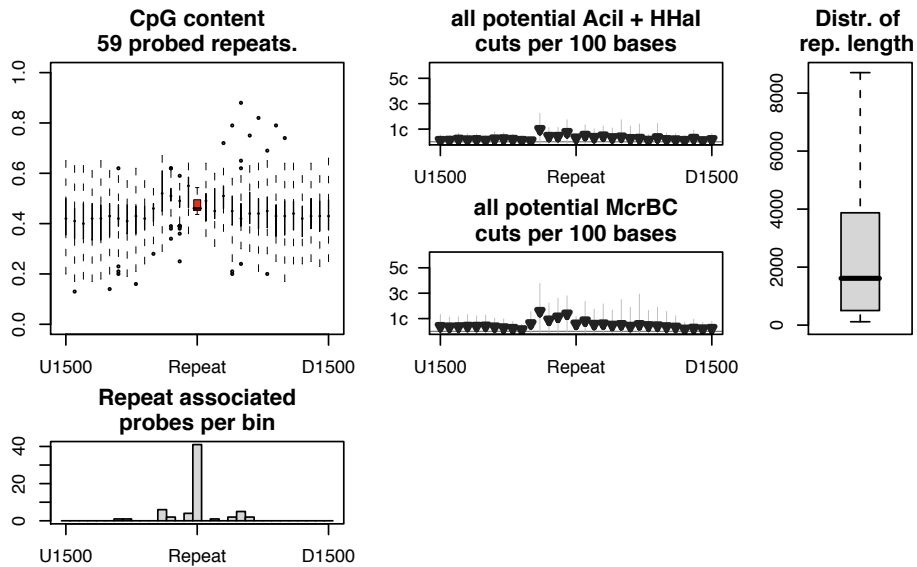
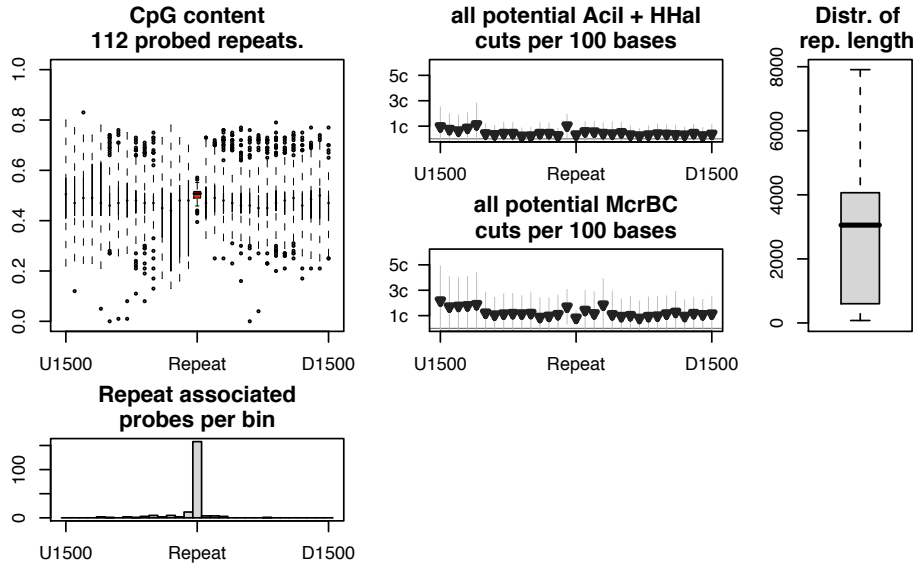
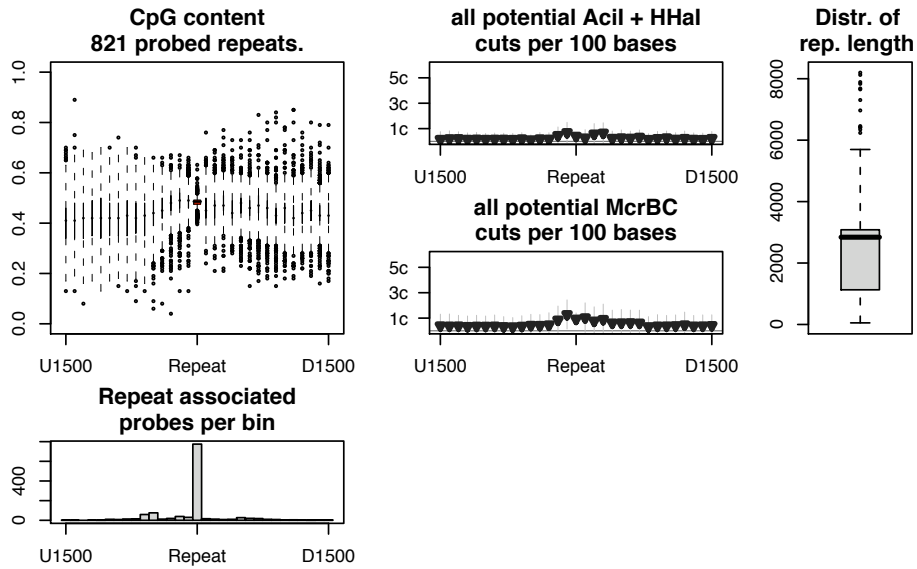


Figure S14

HERVE



HERVH



HERVK

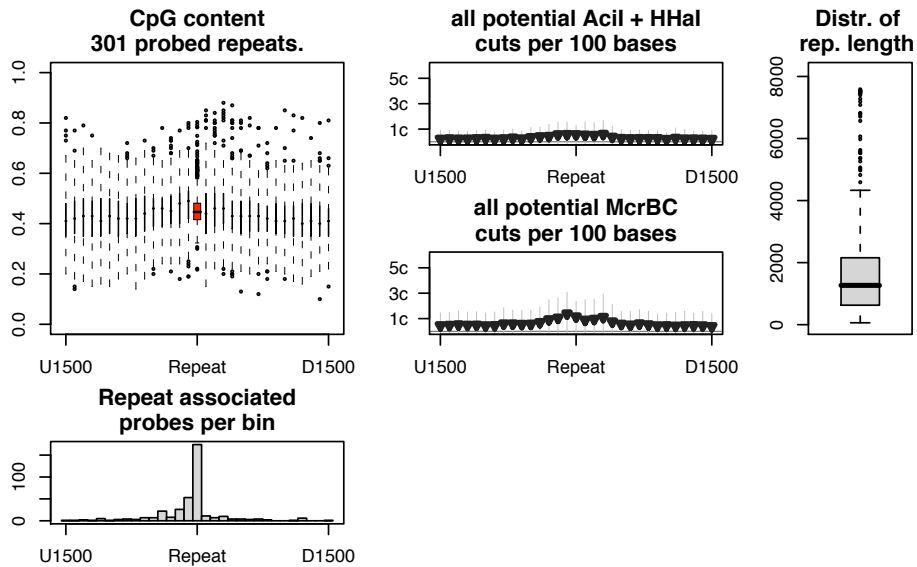


Figure S14

HERVL

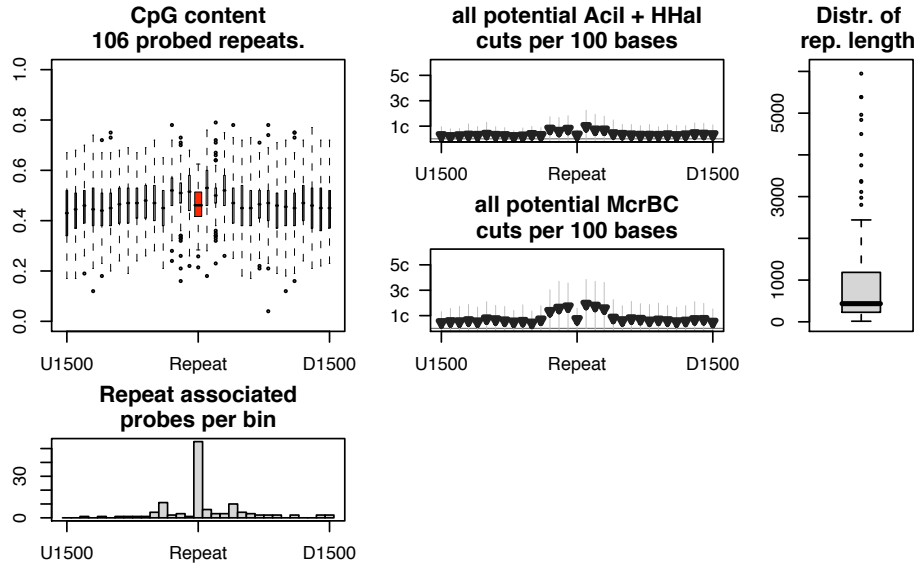
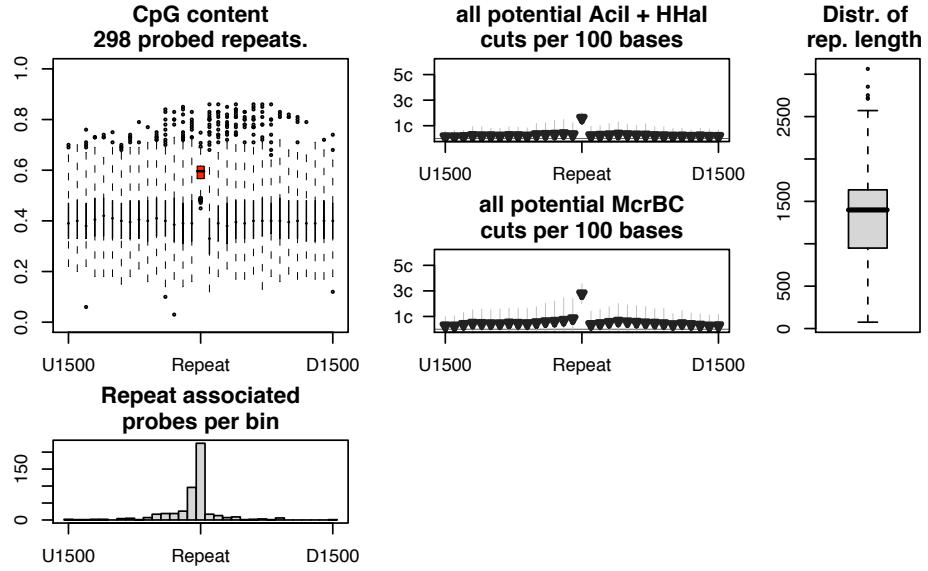
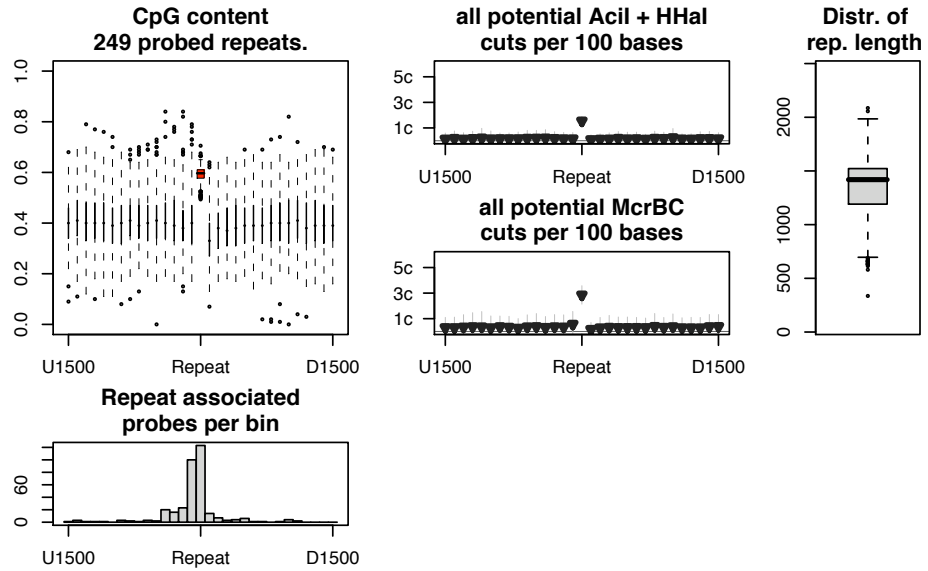


Figure S15

SVA_A
~10-17 MYO



SVA_B
~11-12 MYO



SVA_C
~10-11 MYO

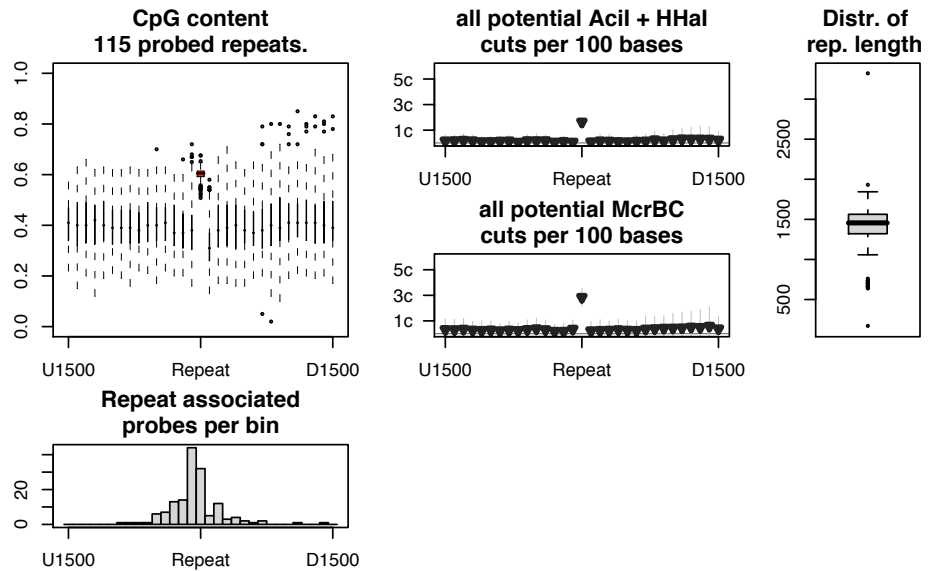
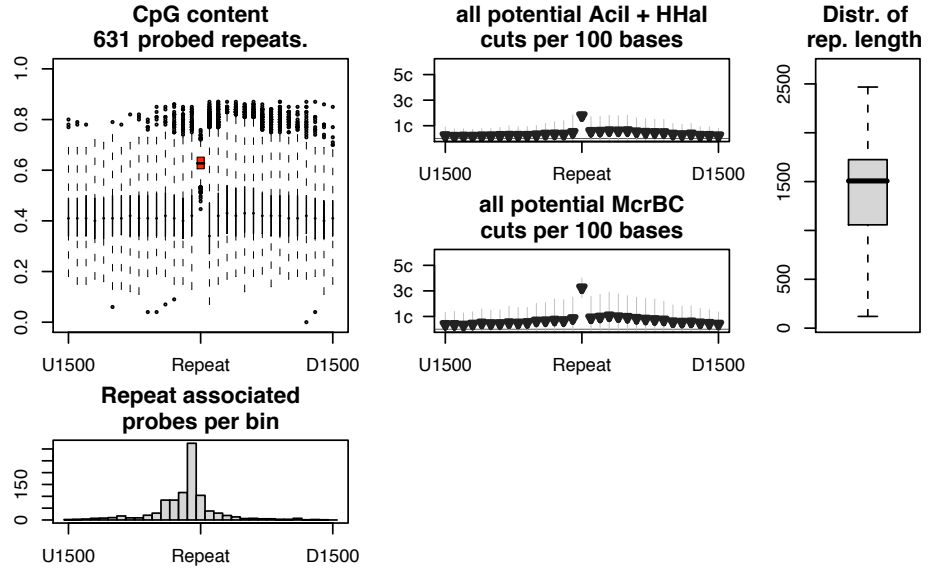
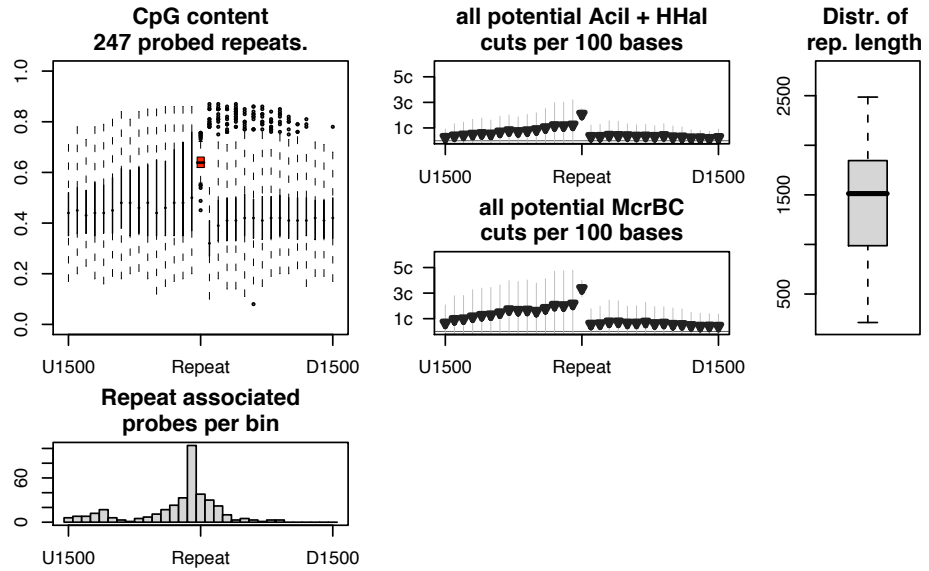


Figure S15

SVA_D
~9.5 MYO



SVA_E
~2.4-4.5 MYO



SVA_F
~2.7-3.6 MYO

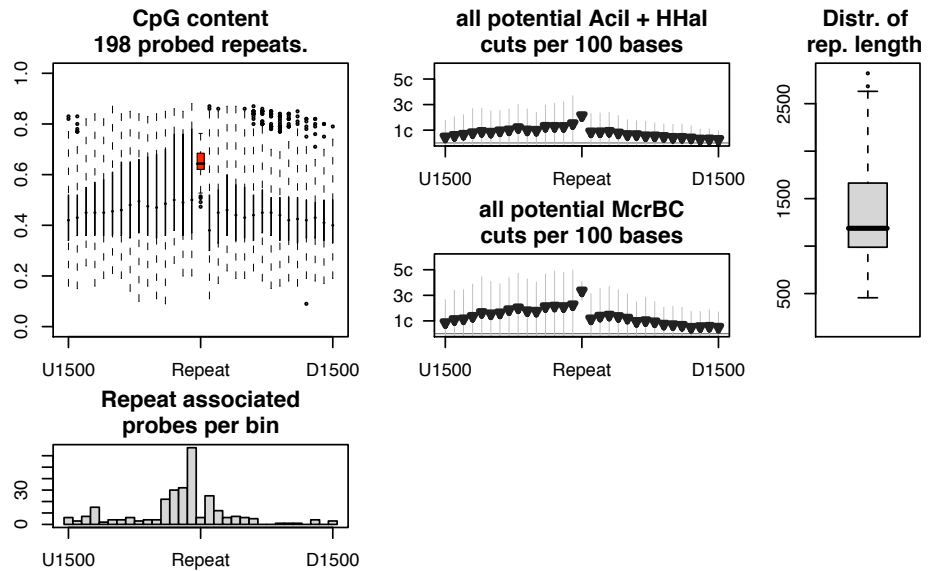
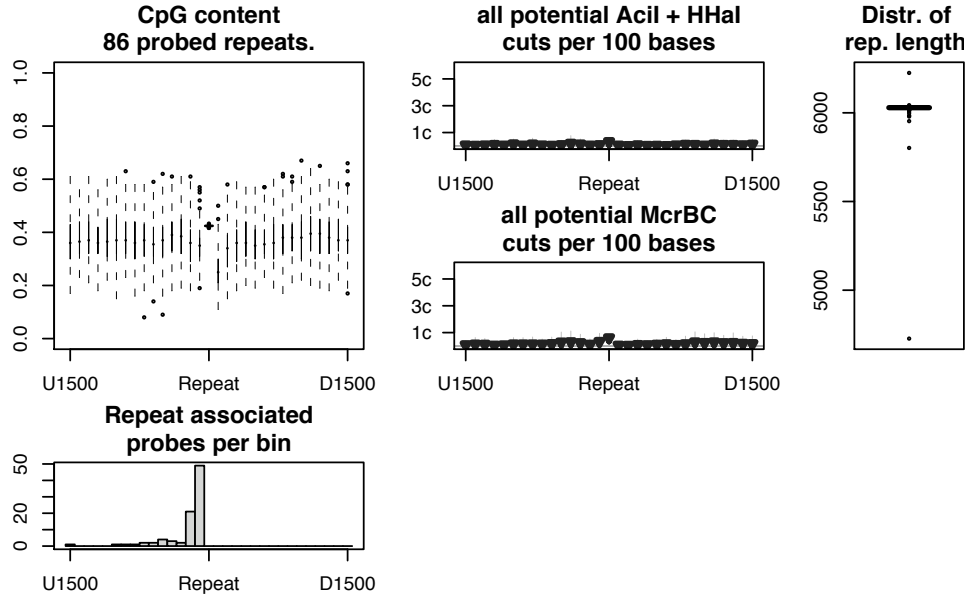
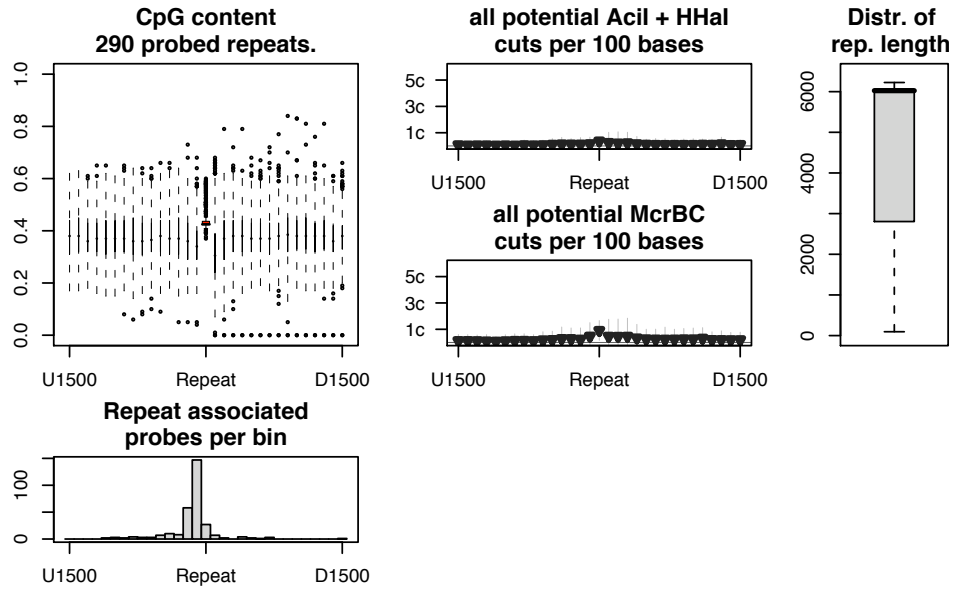


Figure S16

ActiveL1S
~5 MYO



L1HS
~3.1 MYO



L1PA2
~7.6 MYO

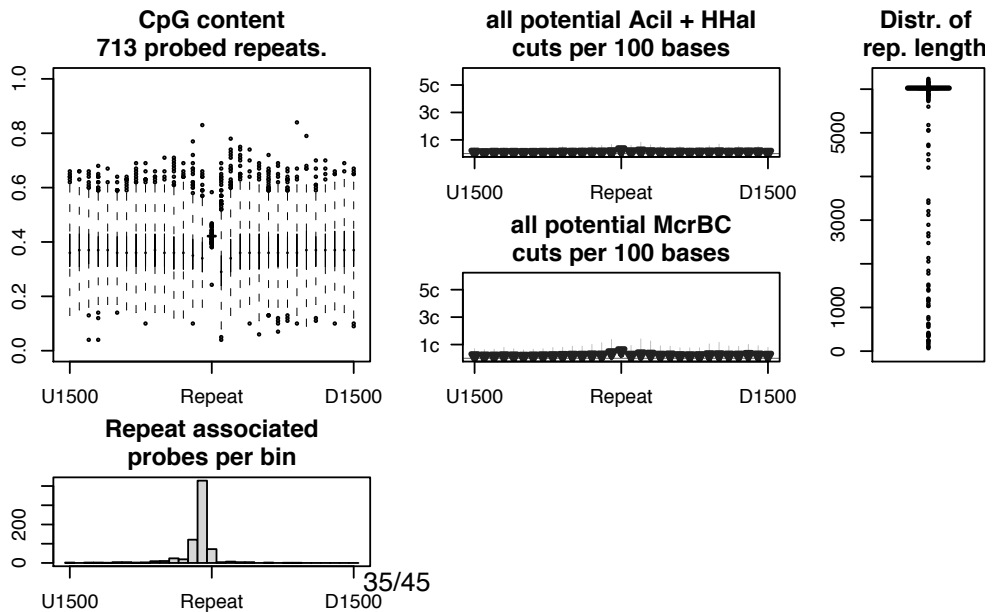
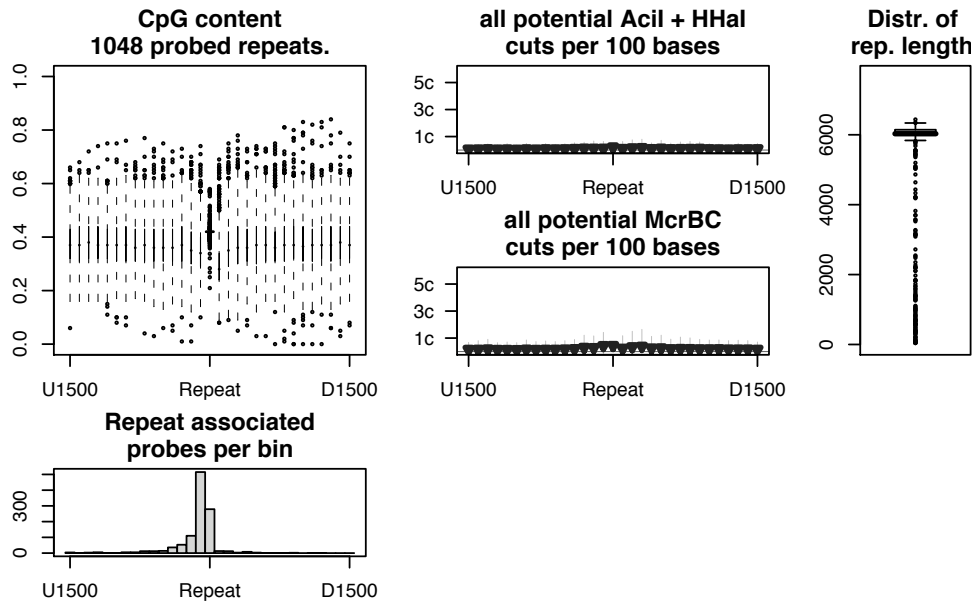
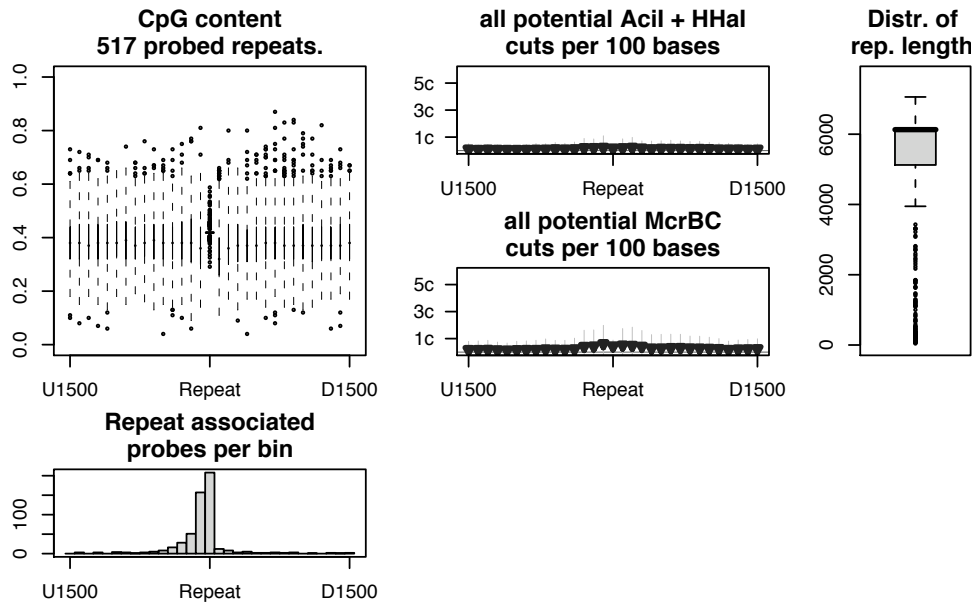


Figure S16

L1PA3
~12.5 MYO



L1PA4
~18.0 MYO



L1PA5
~20.4 MYO

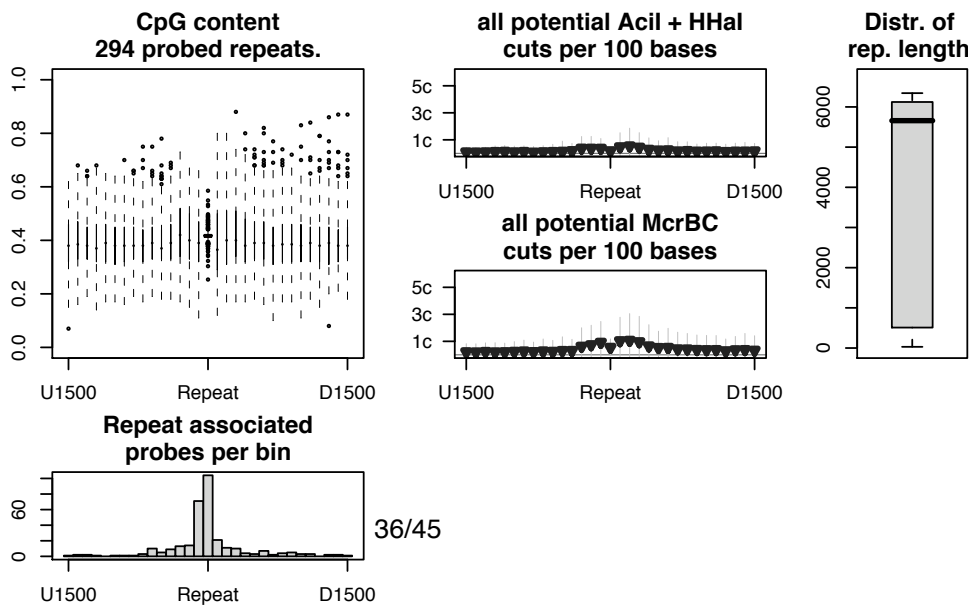
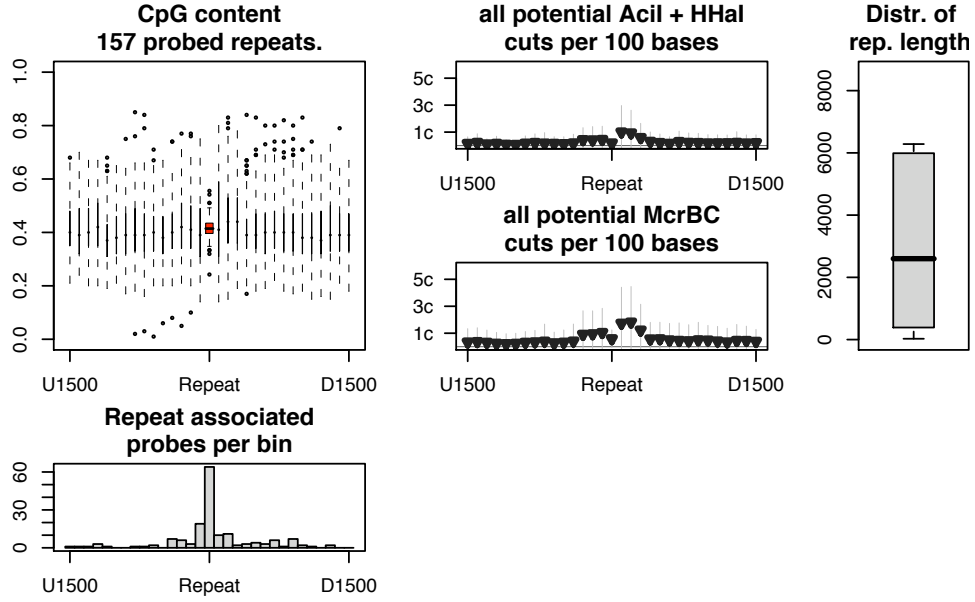
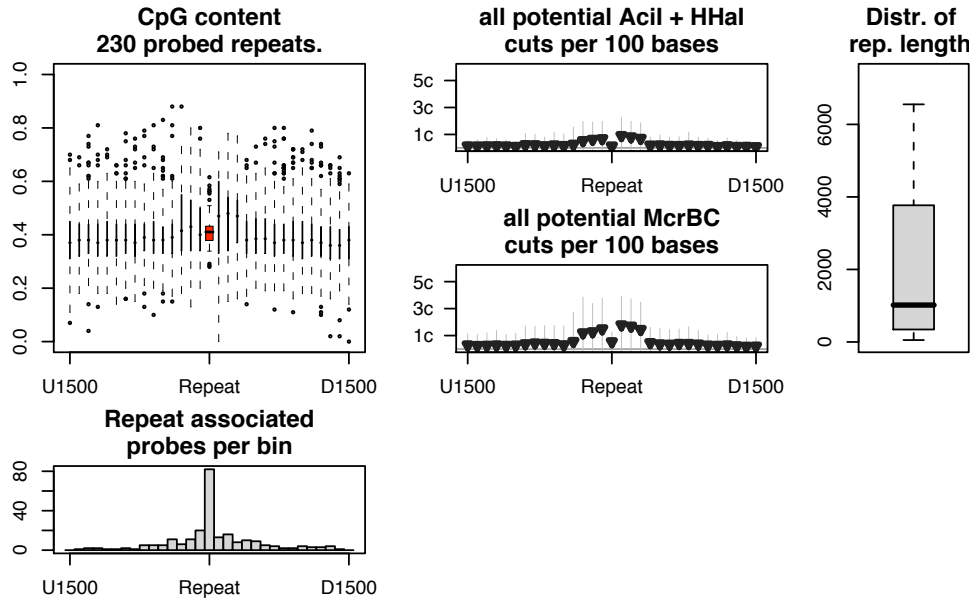


Figure S16

L1PA6
~26.8 MYO



L1PA7
~31.4 MYO



L1PA8
~40.9 MYO

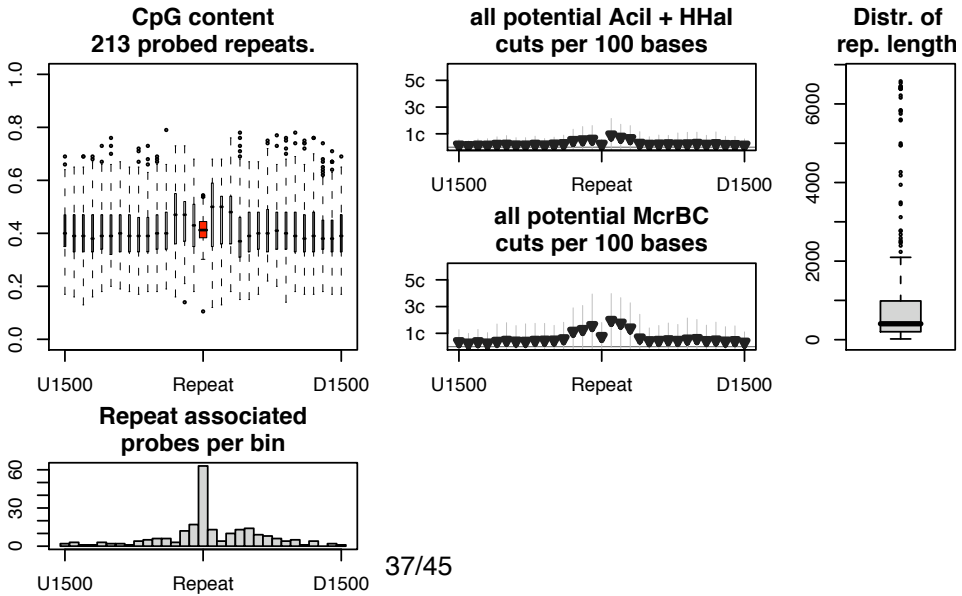
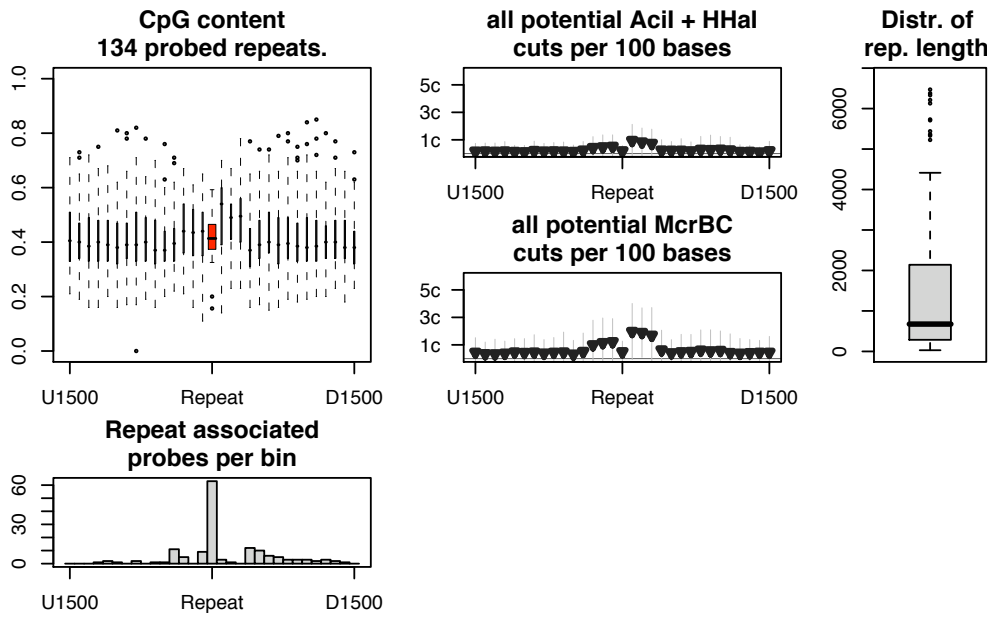
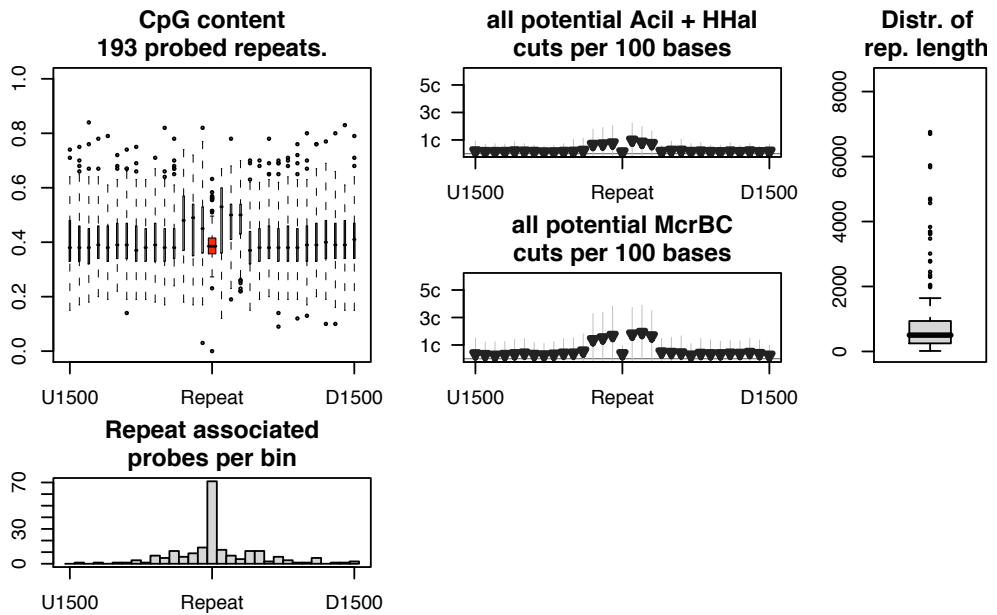


Figure S16

L1PA10
~46.4 MYO



L1PA13
~60 MYO



L1PA15
~70.5 MYO

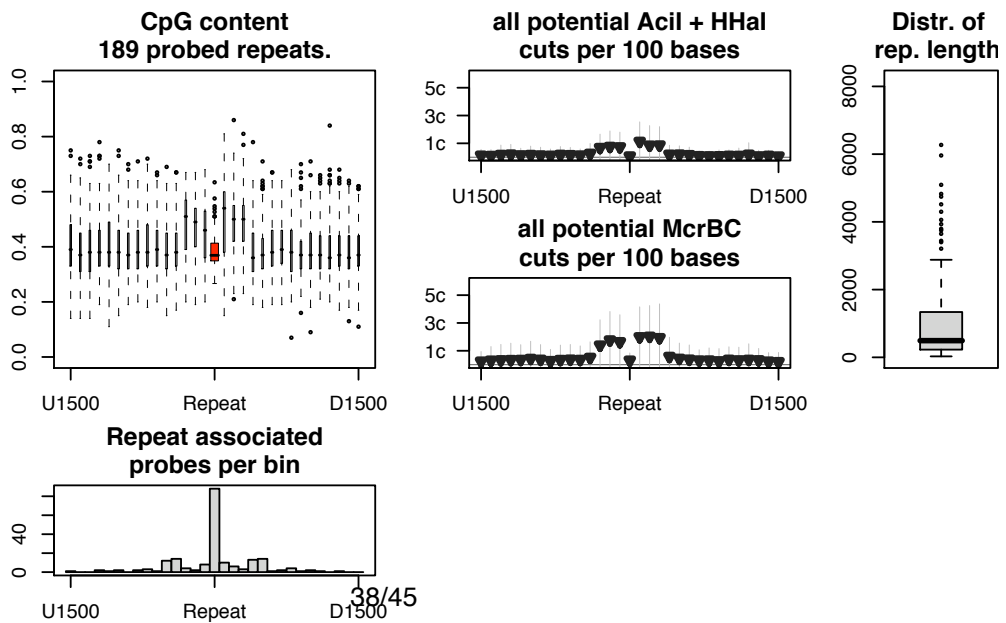
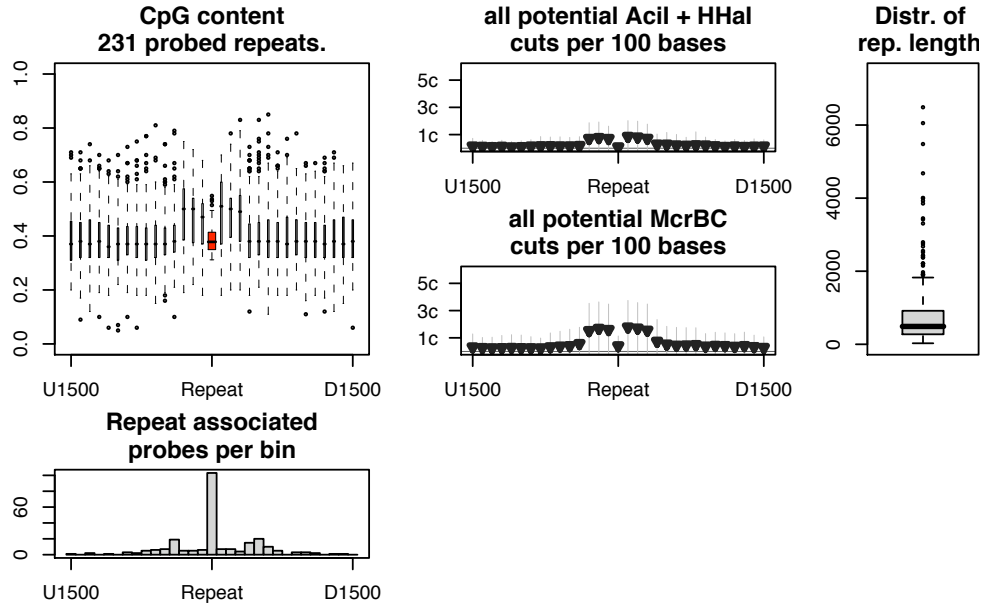


Figure S16

L1PA16

~79.7 MYO



L1PA17

~101.1 MYO

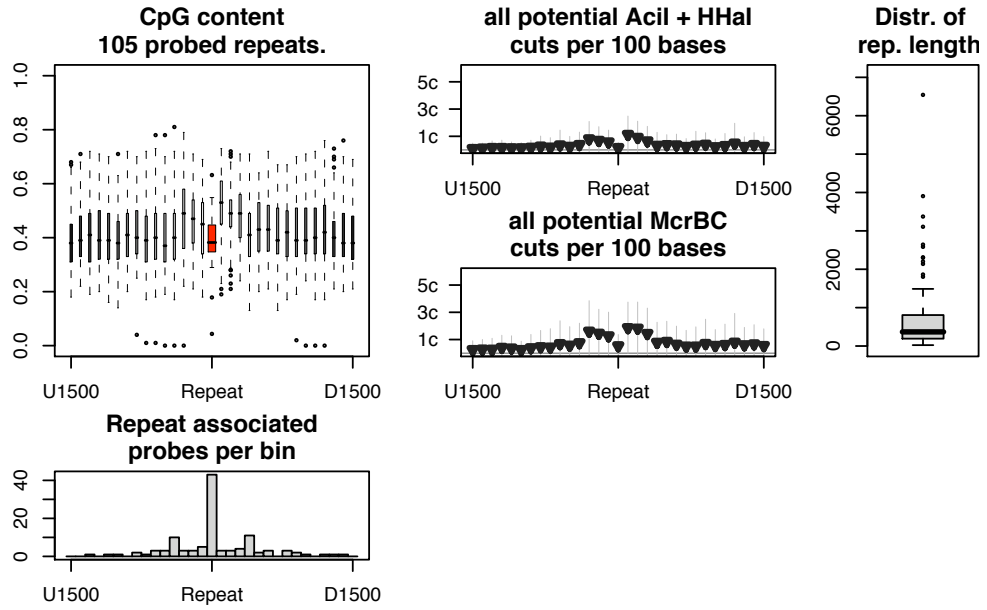


Figure S17

MLT1C

MSTA

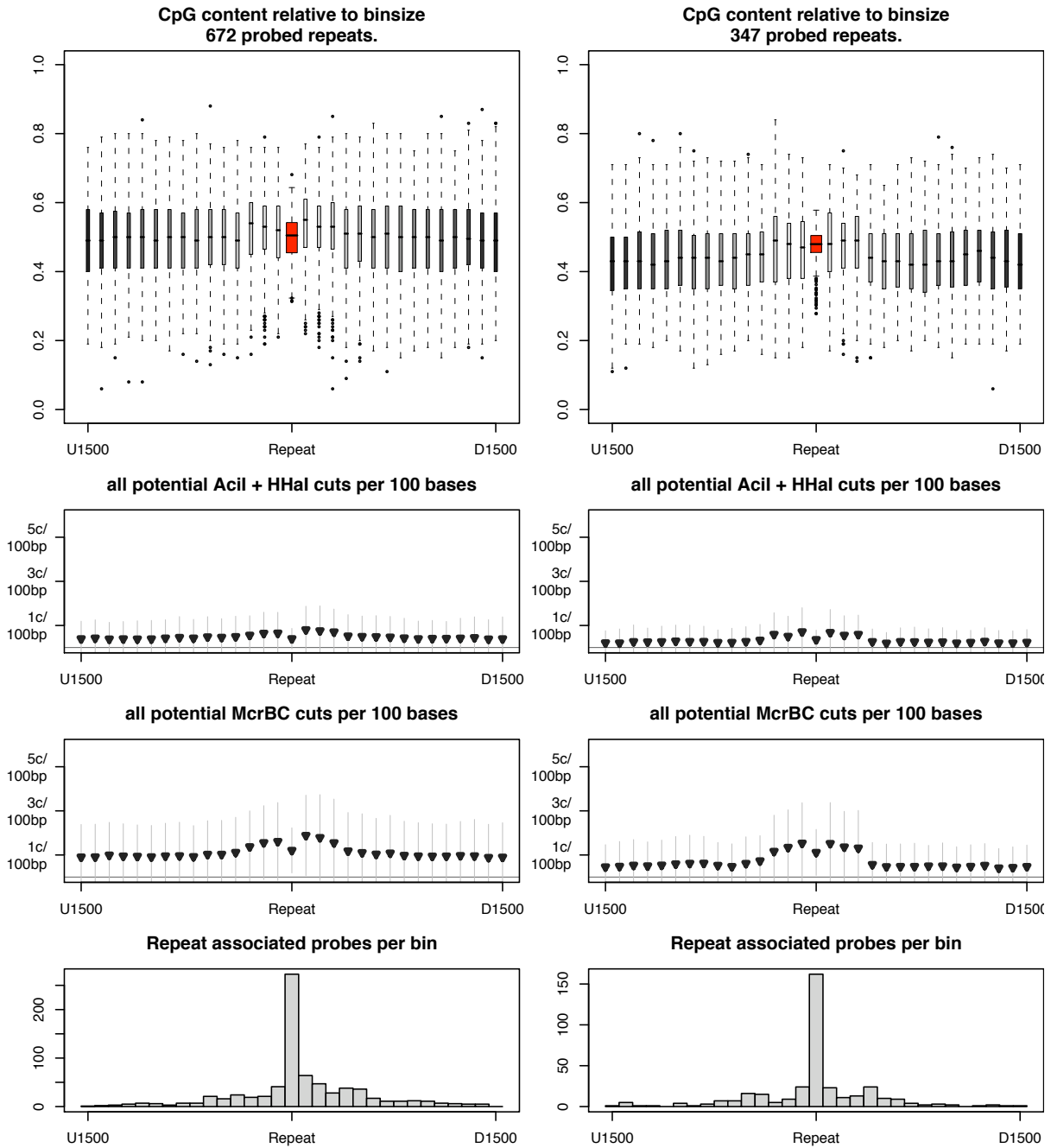


Figure S18

SVA_B

SVA_F

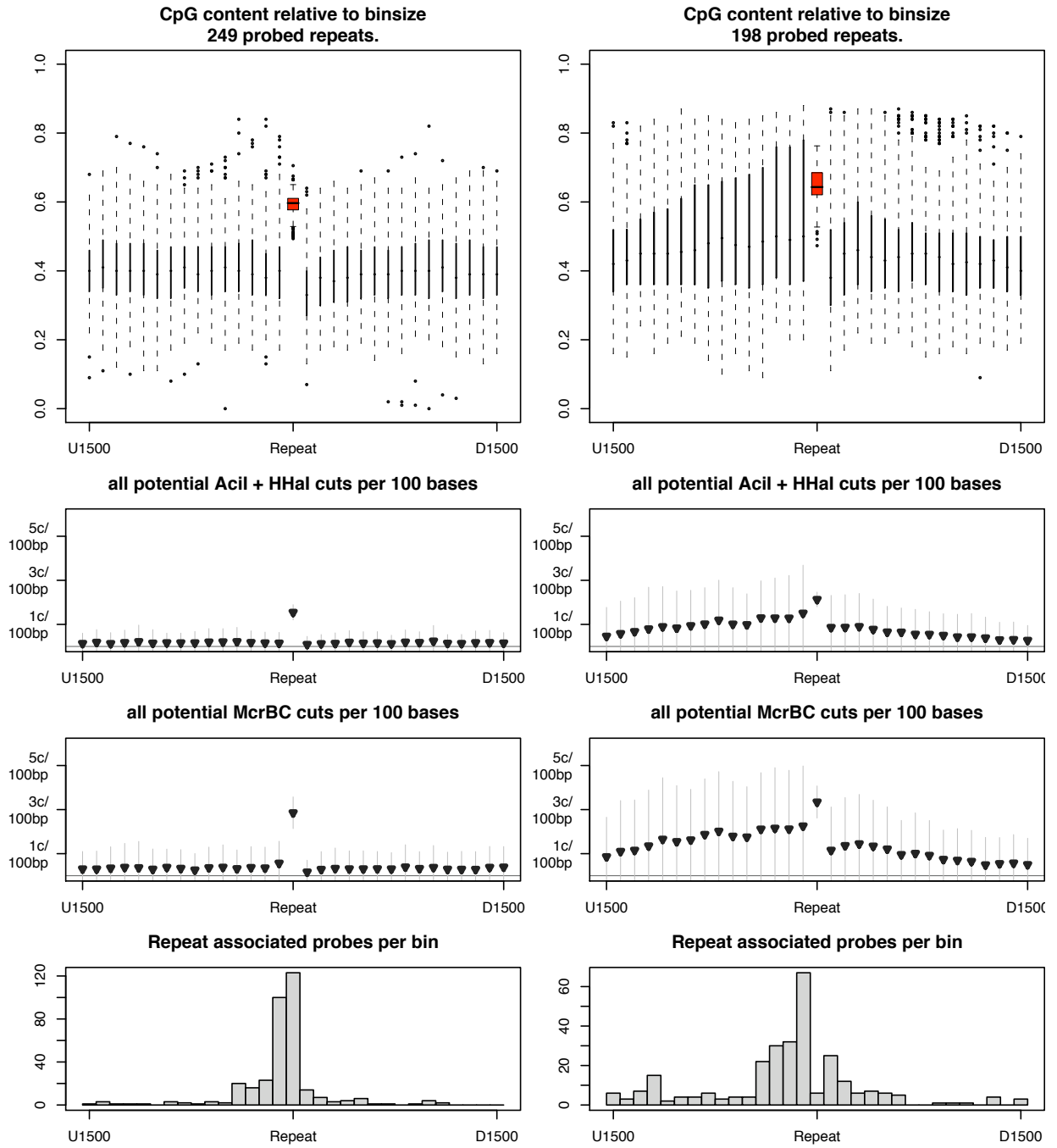


Figure S19

L1PA17

L1PA4

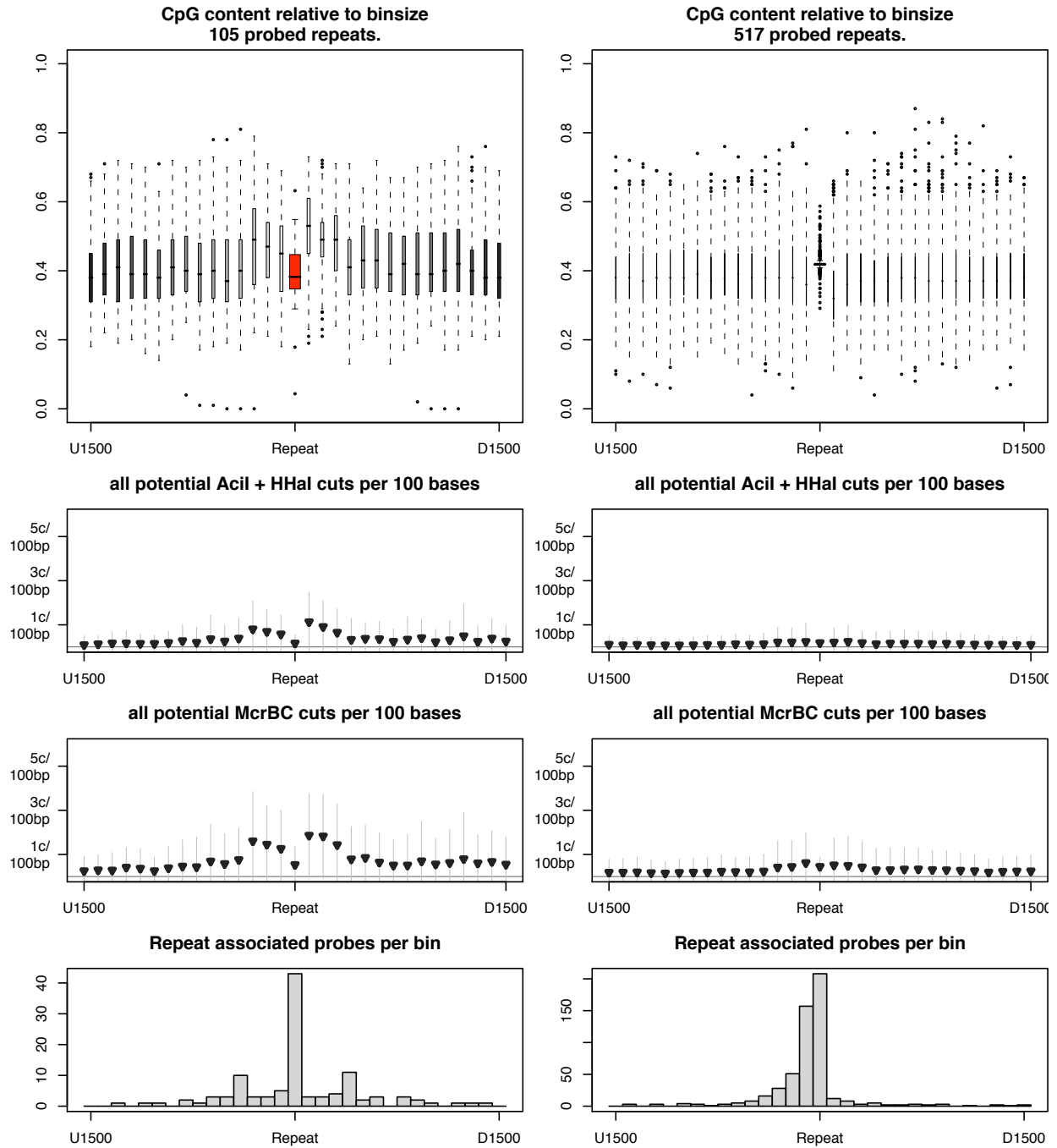


Figure S20

L1PA3

L1HS

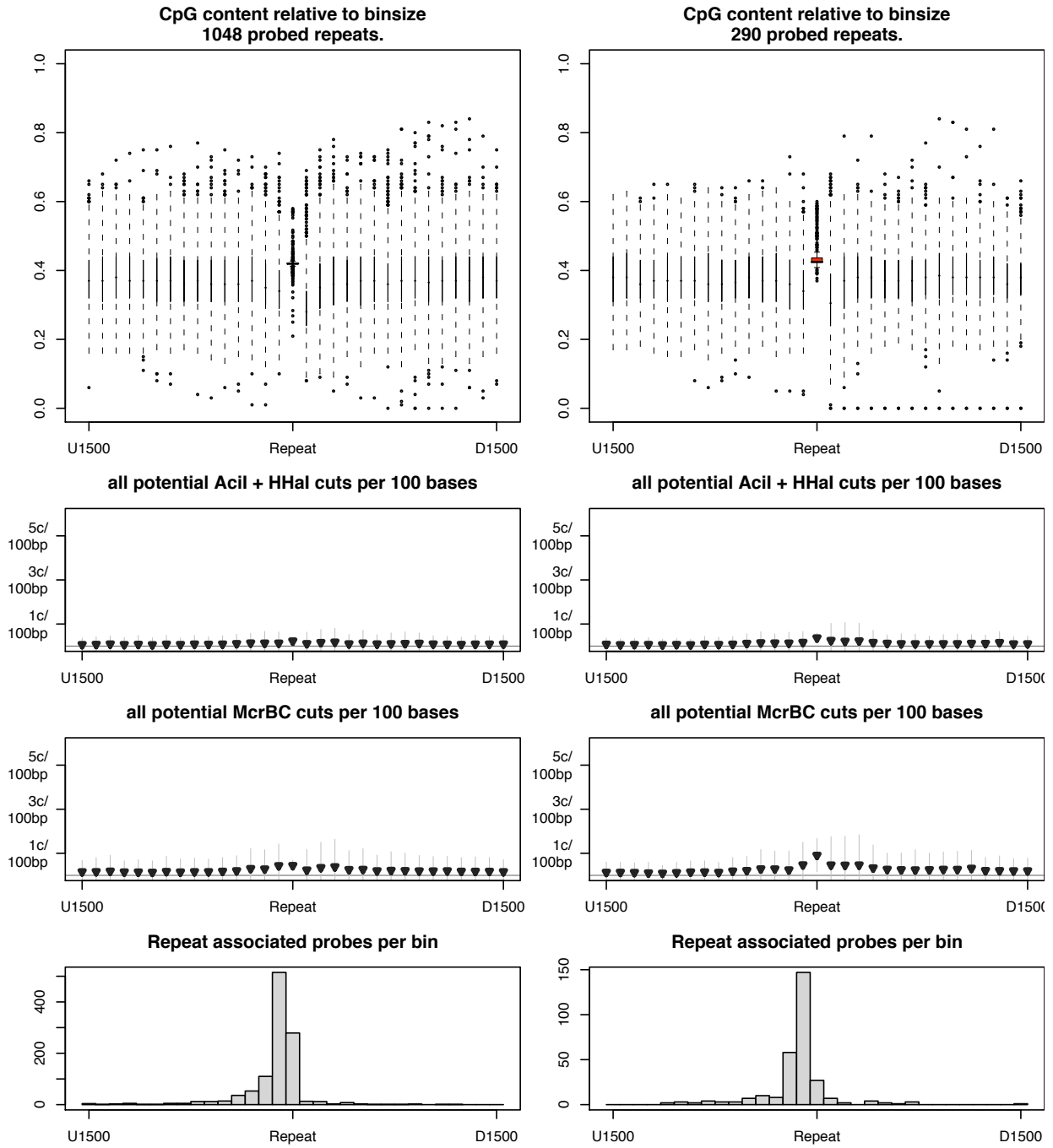


Figure S21

AluYd8

AluYb9

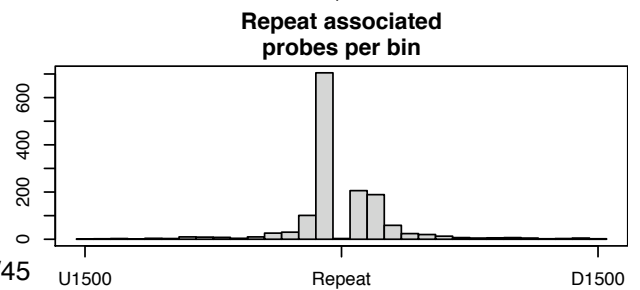
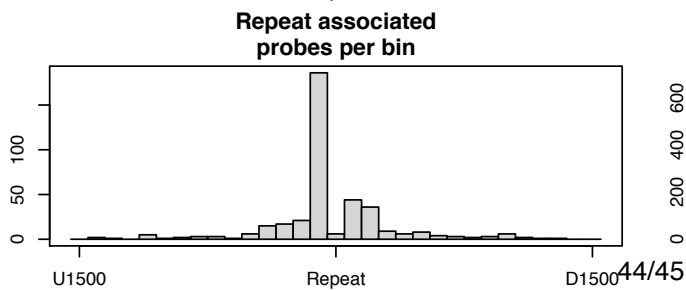
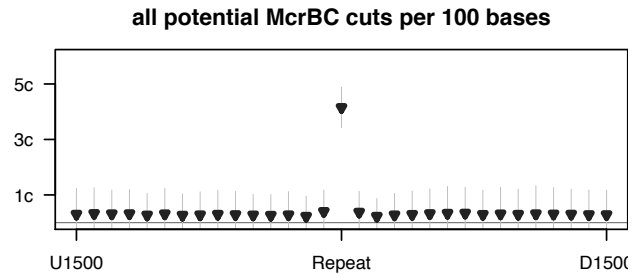
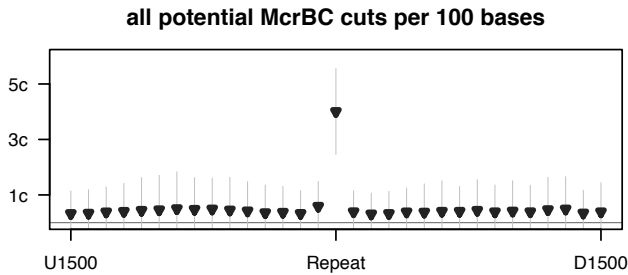
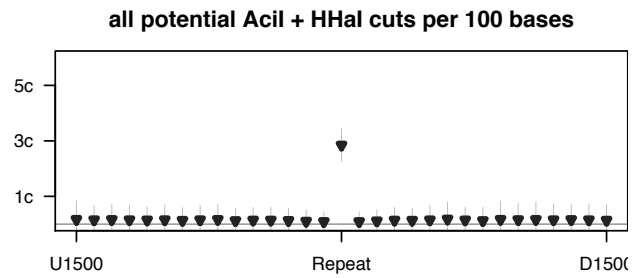
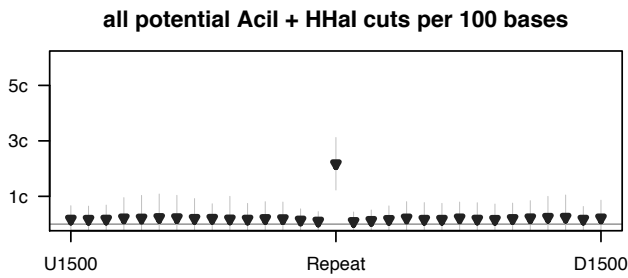
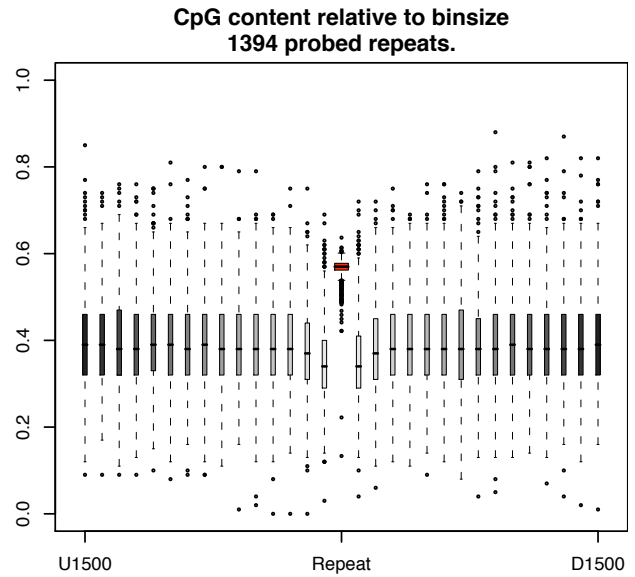
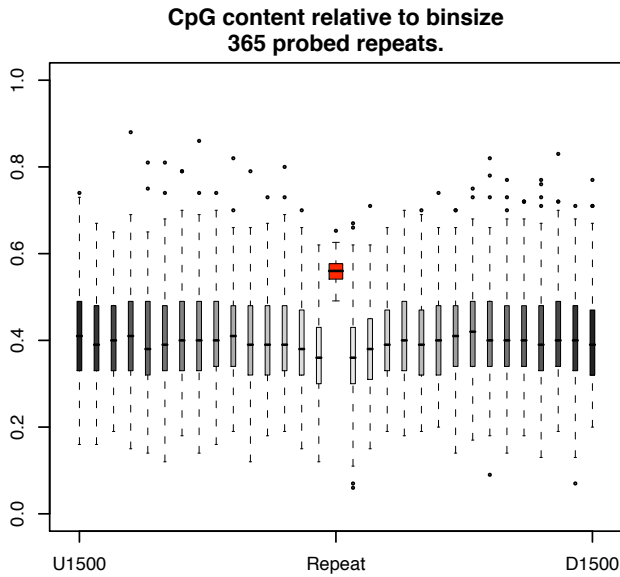


Figure S22.

