**A**

Frequency (y-axis), No. CpGs per gene (x-axis)

mode = 15
median = 15
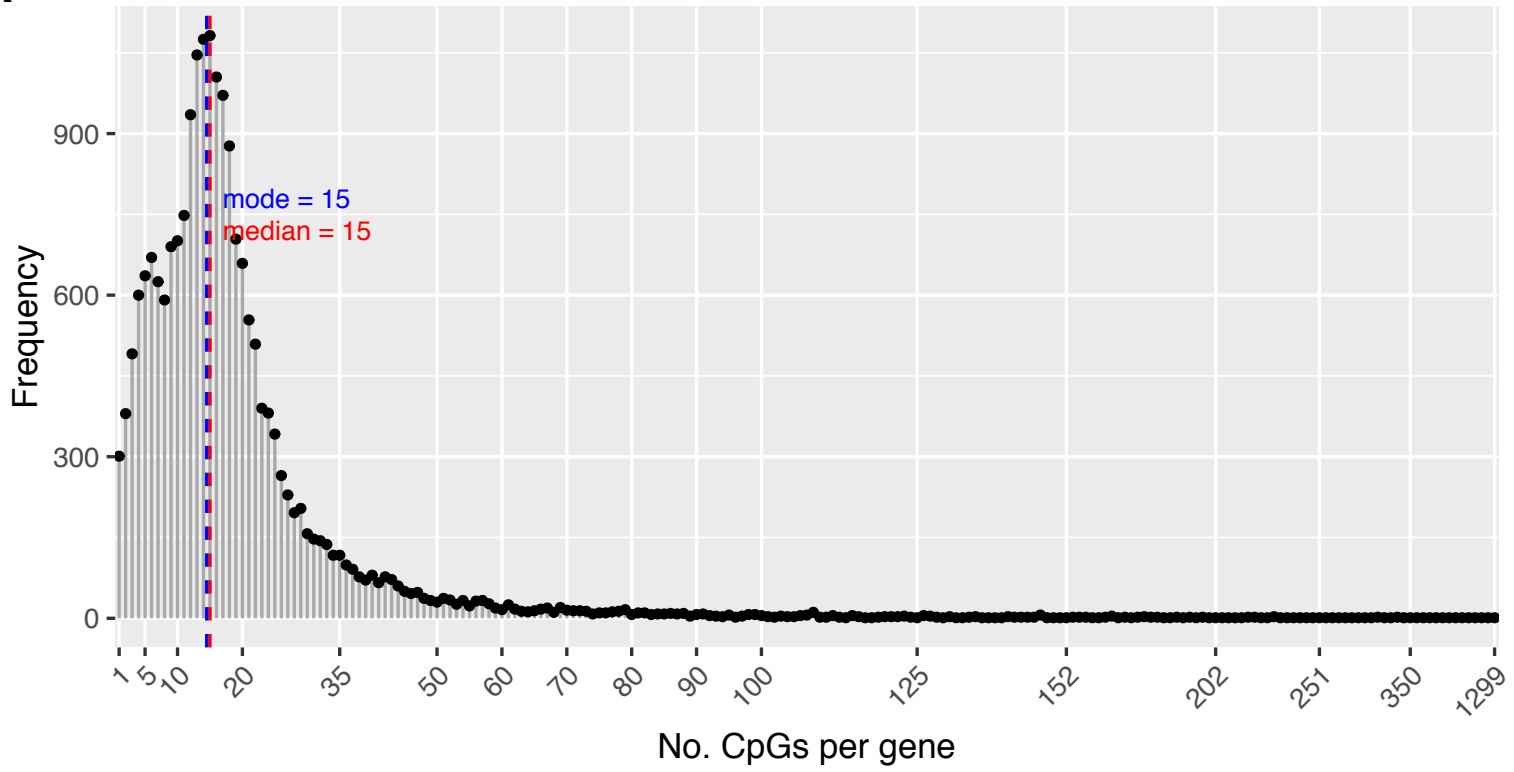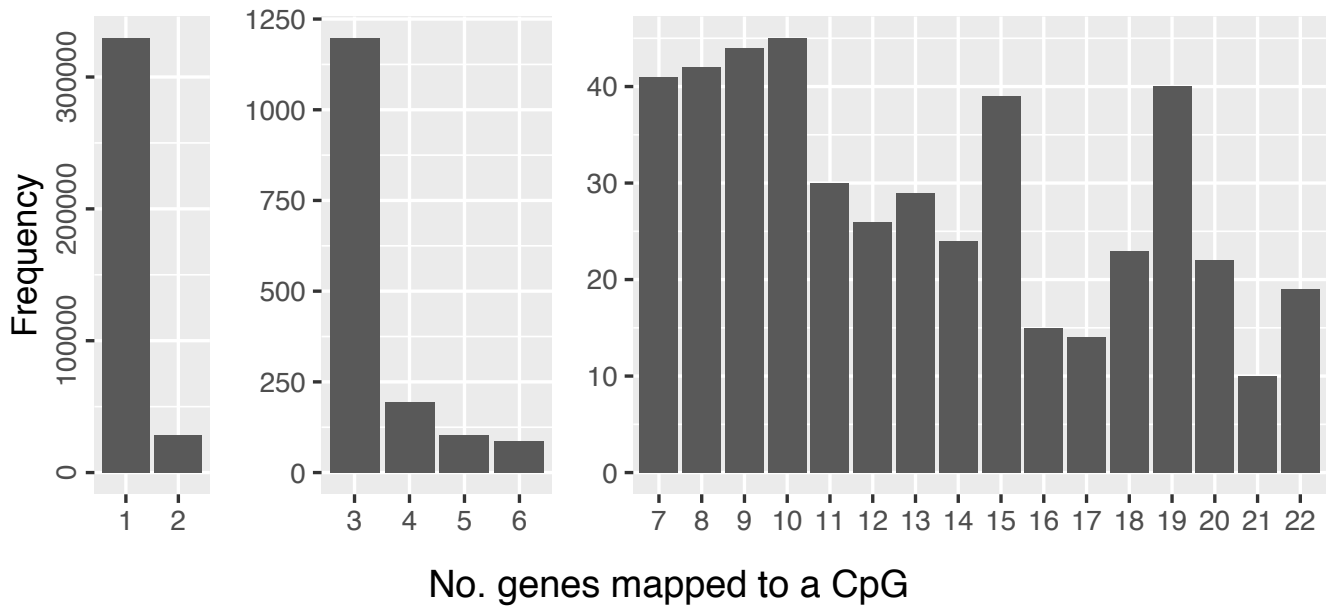
**B**

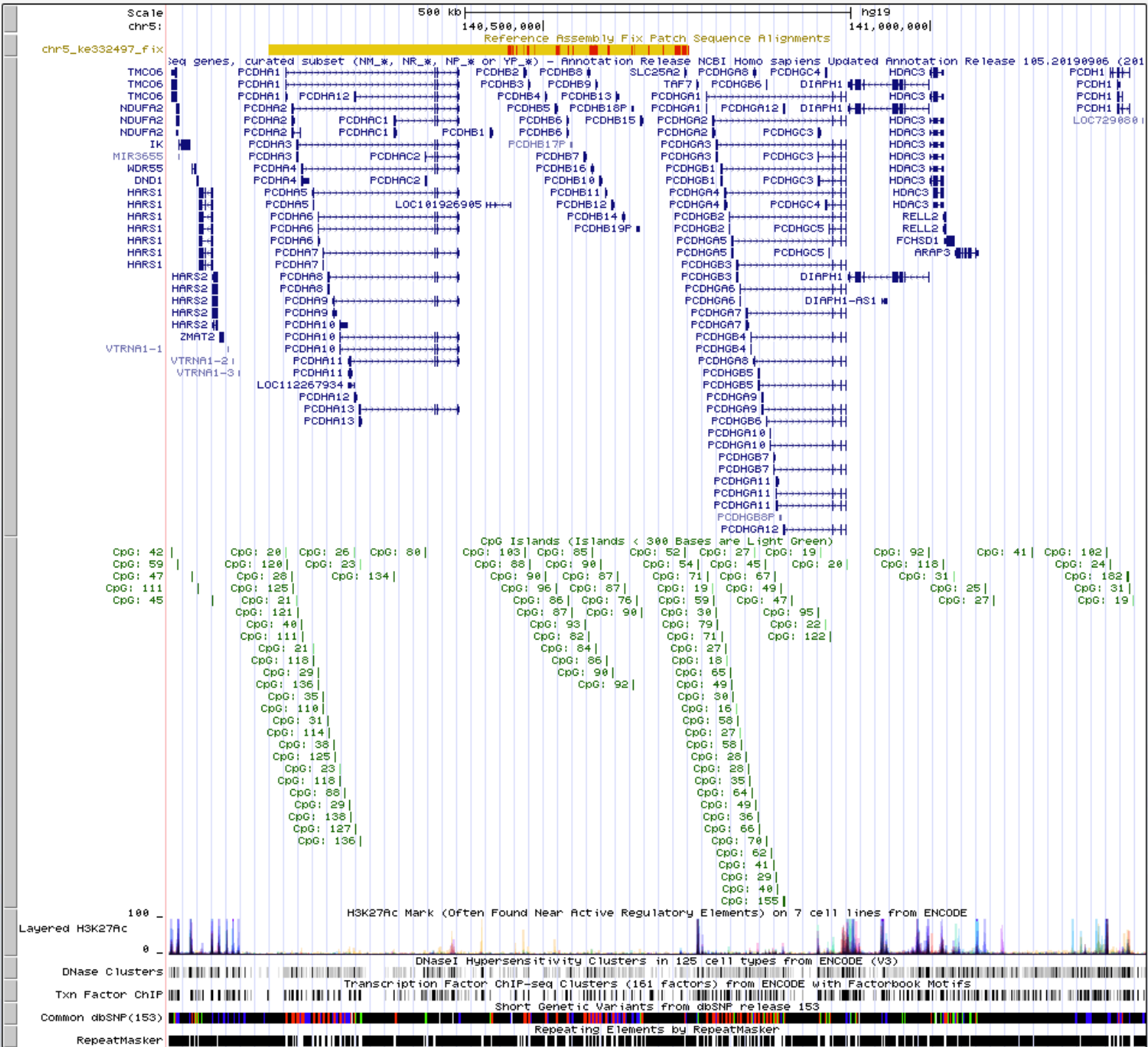Frequency (y-axis), No. genes mapped to a CpG (x-axis)

**Fig. S1**. **Array design bias for the Illumina HumanMethylation 450K BeadChip. (A)** Frequency plot of the numbers of CpGs measuring methylation across each gene for the 450K array (probe number bias). The most extreme value is 1299 CpGs measuring methylation for a single gene. The median is 15 and the mode is 15. **(B)** Split bar chart showing the numbers of genes annotated to each CpG (multi-gene bias). While the majority of CpGs are annotated to only one gene, there is still a large number annotated to 2 or more genes.
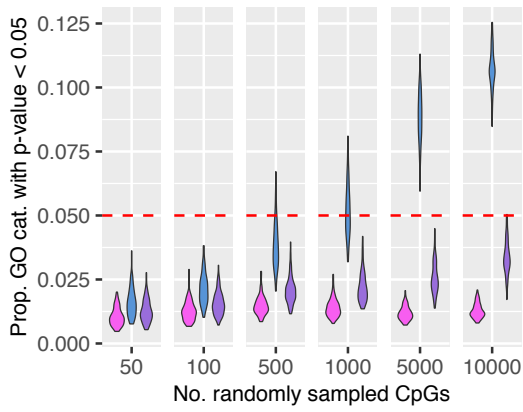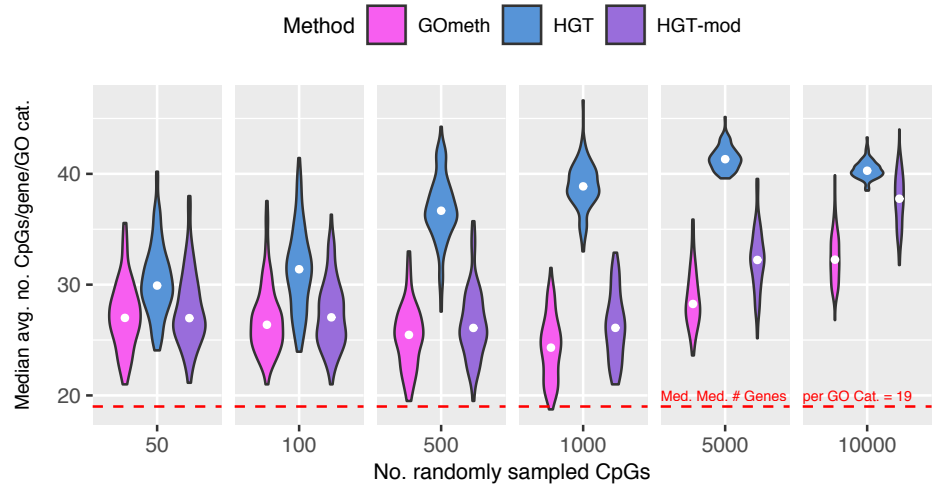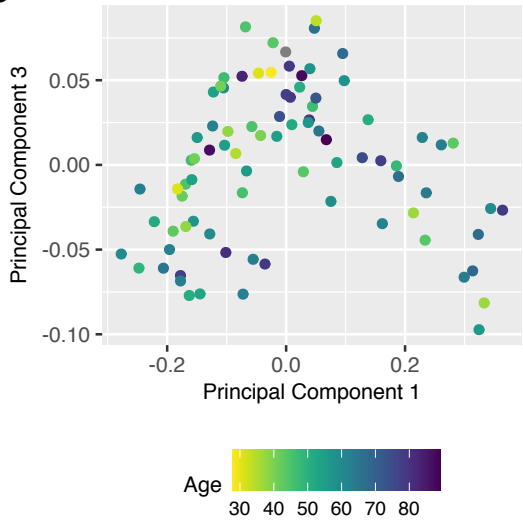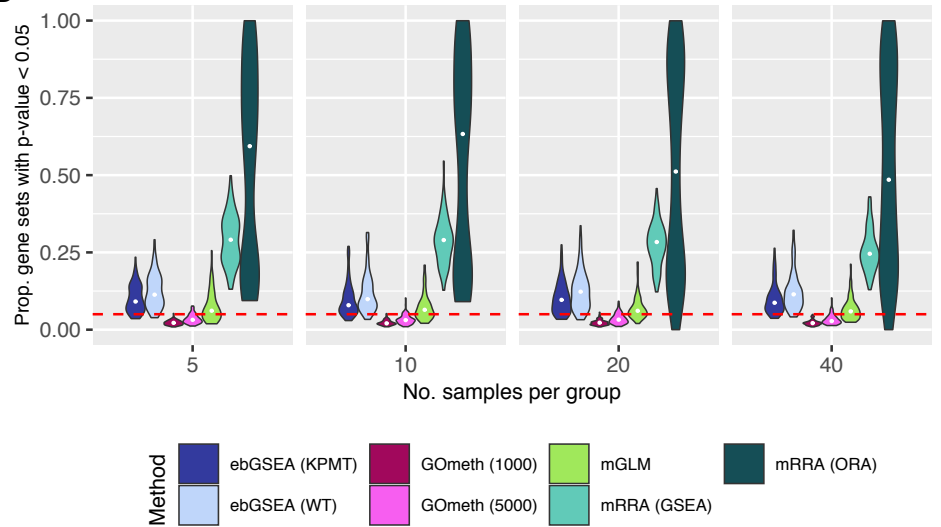
**Fig. S2.** UCSC genome browser snapshot of the protocadherin gamma cluster.

**Fig. S3. Evaluation of false discovery rate control for 450K data. (A)** Type I error rates across 100 simulations for varying numbers of randomly sampled CpGs. **(B)** Median average numbers of CpGs per gene for GO categories with an unadjusted p-value < 0.05. The hypergeometric test is biased towards GO categories with more CpGs per gene on average. GOmeth = adjust for probe-number and multi-gene bias; HGT = hypergeometric test; HGT-mod = adjust for probe-number bias only. **(C)** Multidimensional scaling plot of normal samples from TCGA BRCA data, coloured by age. **(D)** False discovery rate control of seven gene set testing methods using normal samples from TCGA BRCA data. Two groups were generated by randomly sampling *n* samples per group, followed by differential methylation analysis and subsequent gene set testing. This was repeated 100 times at each sample size. The proportion of gene sets with unadjusted p-value < 0.05 across the 100 null simulations is shown for each method, at each sample size. Methods with good false discovery rate control should have r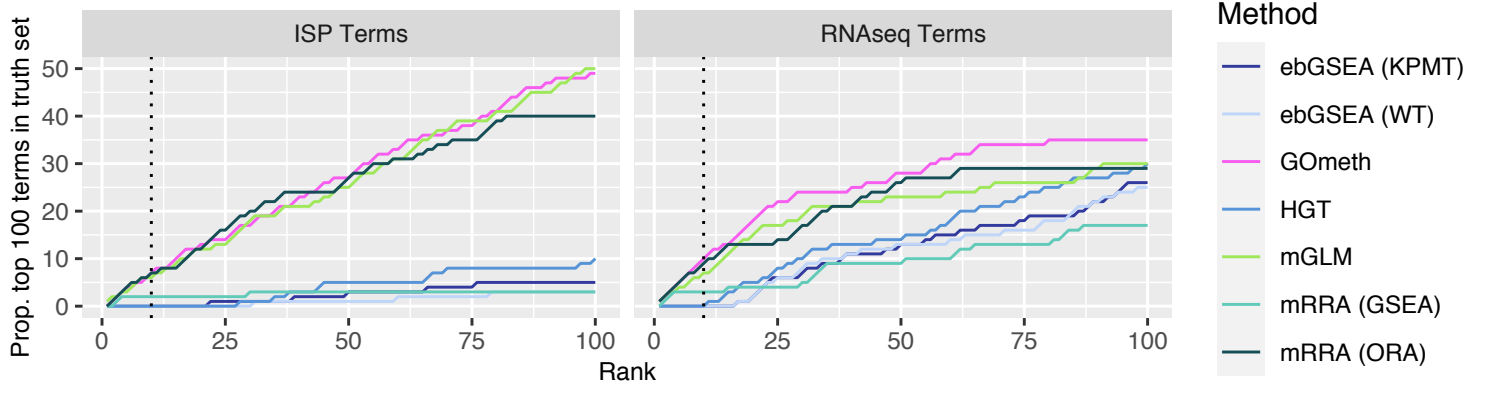elatively tight distributions around the red dashed line at 0.05. ebGSEA (KPMT) = ebGSEA using Known Population Median Test; ebGSEA (WT) = ebGSEA using Wilcoxon Test; GOmeth (1000) = GOmeth using top 1000 ranked probes; GOmeth (5000) = GOmeth using top 5000 ranked probes; mGLM = methylglm; mRRA (GSEA) = methylRRA using gene set enrichment analysis; mRRA (ORA) = methylRRA using over-representation analysis. **(E)** Gene set testing was performed on the results of the three blood cell type comparisons: CD4 T-cells vs. CD8 T-cells, monocytes vs. neutrophils and B-cells vs. NK cells, using the *MethylGSA* methods: mGLM, mRRA (GSEA) and mRRA (ORA), with several combinations of minimum and maximum gene set size parameters. When the minimum gene set size is set to less than five, mGLM and mRRA (ORA) highly rank very small gene sets containing the number of genes equal to the minimum size parameter. mRRA is also biased towards highly ranking very large gene sets, if they are not filtered out. mGLM = methylglm; mRRA (GSEA) = methylRRA using gene set enrichment analysis; mRRA (ORA) = methylRRA using over-representation analysis.
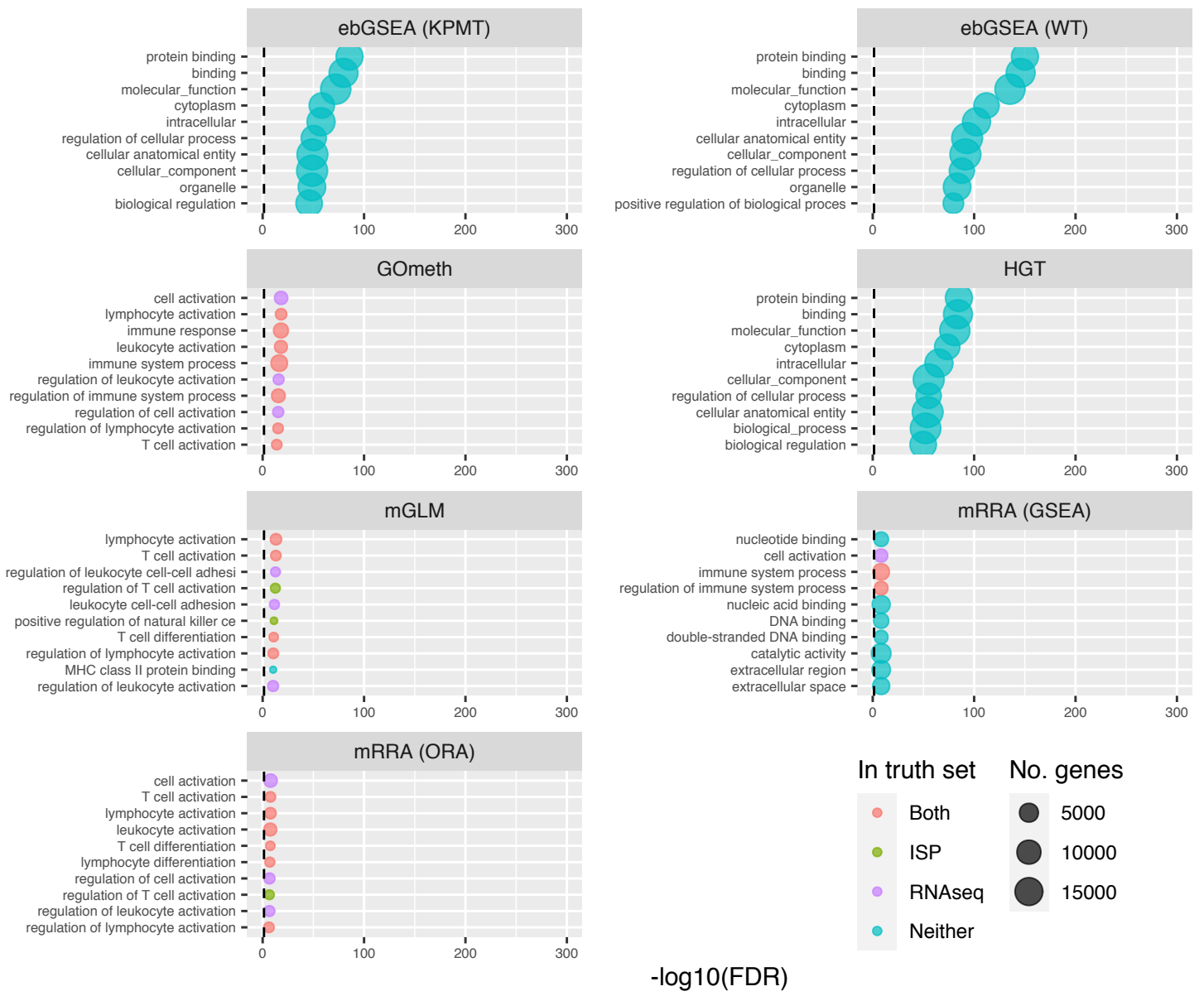
**Fig. S4. Characteristics of gene expression datasets and analysis. (A)** Flow sorted blood cell types RNA-seq data: Library sizes (length-scaled TPM) for each of the samples after Salmon mapping and quantification. **(B)** Flow sorted blood cell types RNA-seq data: Distribution of counts per million (CPM) after filtering and normalisation. **(C)** Flow sorted blood cell types RNA-seq data: Multidimensional scaling plot of the filtered and normalised RNAseq data. **(D)** Flow sorted blood cell types RNA-seq data: Numbers of differentially expressed genes with an adjusted p-value < 0.05, for each cell type comparison: CD4 T-cells vs. CD8 T-cells, monocytes vs. neutrophils and B-cells vs. NK cells. The blue bar is the number of significantly down-regulated genes, and the red bar is the number that are significantly up-regulated; e.g. ~2,300 genes are down-regulated and ~2500 are up-regulated in B-cells, compared to NK cells. (E) Log2 intensity distributions of the raw data for all B-cell development samples. (F) Multidimensional scaling plot of the normalised and filtered B-cell development Affymetrix array data. (G) Numbers of differentially expressed genes with an adjusted p-value < 0.01, for Stage 1 versus Stage 2 of B-cell development. The blue bar is the number of significantly down-regulated genes, and the red bar is the number that are significantly up-regulated; e.g. ~400 genes are down-regulated and ~500 are up-regulated in Stage 1, compared to Stage 2.

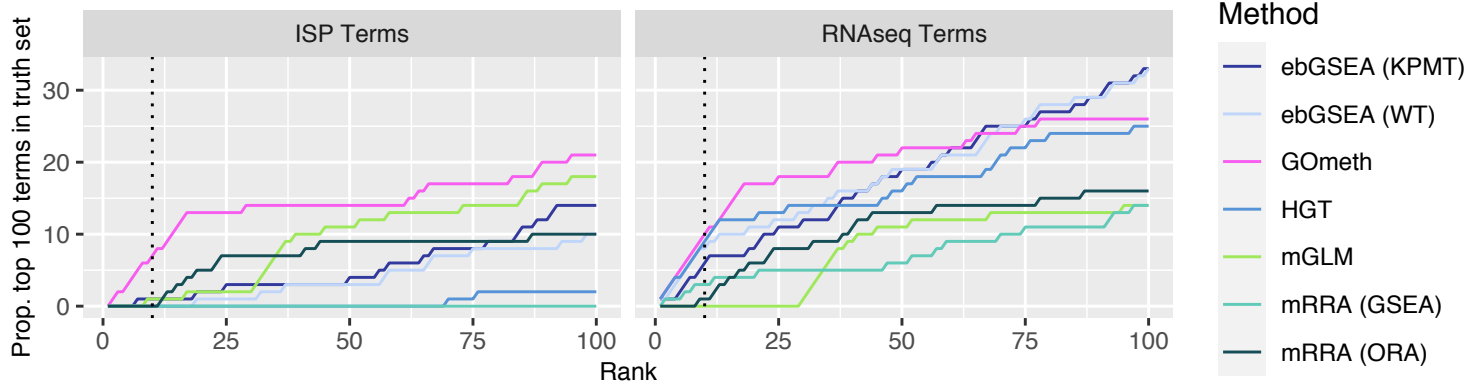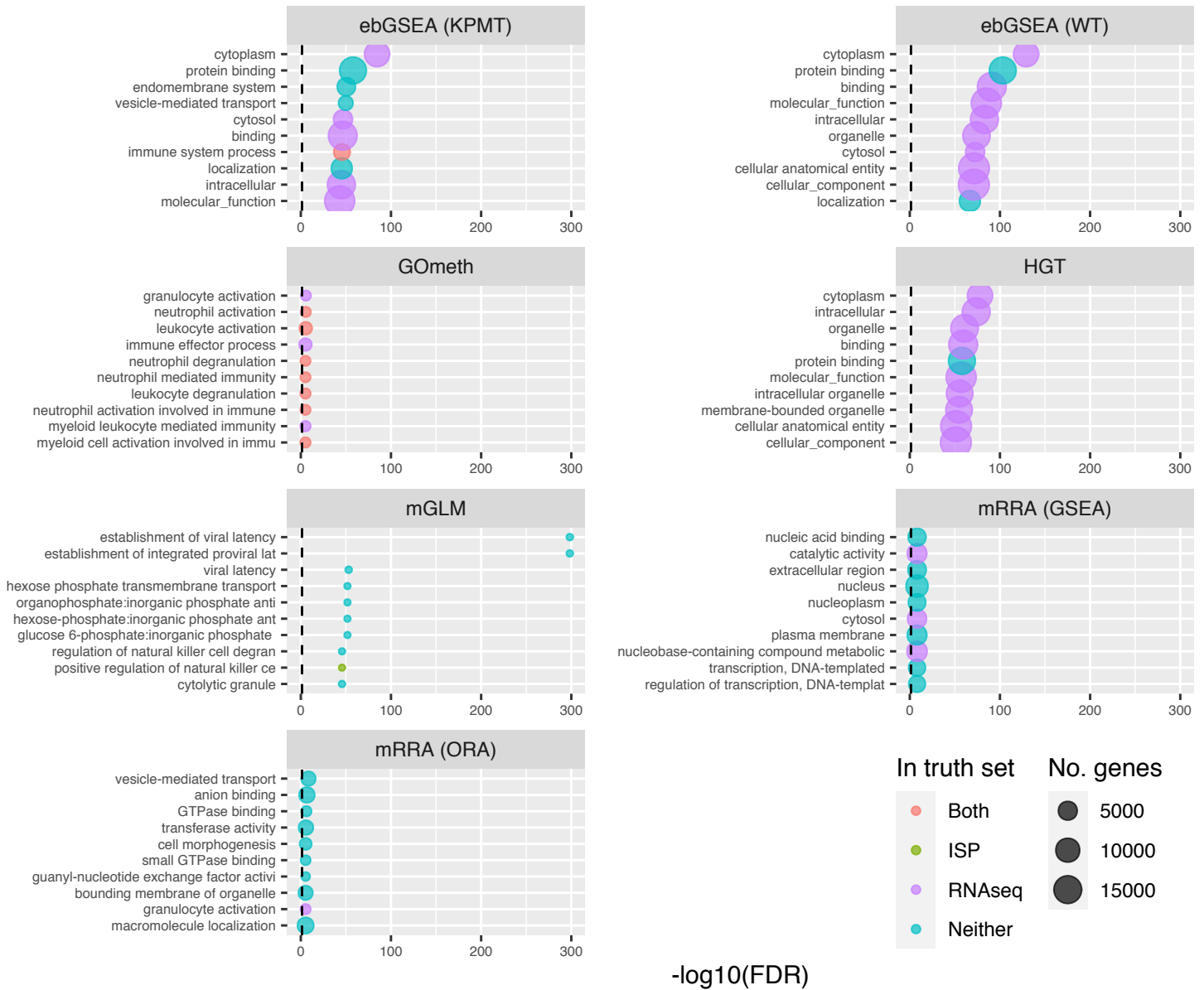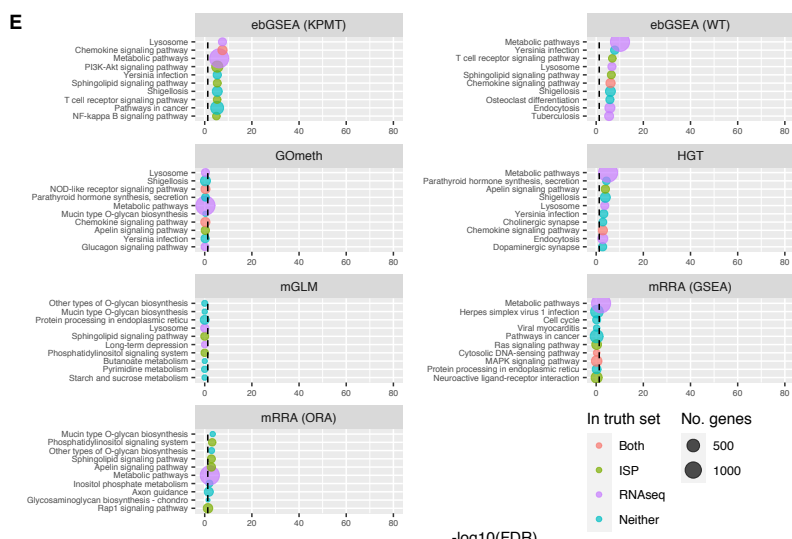**Fig. S5. Comparison of gene set testing performance for CD4 T cells vs. CD8 T-cells (GO).**
**(A)** Cumulative number of GO terms, as ranked by various methods, that are present in each truth set for the CD4 T-cells vs. CD8 T-cells comparison. ISP Terms = immune-system process child terms truth set; RNAseq Terms = top 100 terms from RNAseq analysis of the same cell types. **(B)** Bubble plots of the top 10 GO terms as ranked by various gene set testing methods for the CD4 T-cells vs. CD8 T-cells comparison. The size of the bubble indicates the relative number of genes in the set. The colour of the bubble indicates whether the term is present in either RNAseq (purple) or ISP (green) truth sets, both (red) or neither (blue). ebGSEA (KPMT) = ebGSEA using Known Population Median Test; ebGSEA (WT) = ebGSEA using Wilcoxon Test; GOmeth = GOmeth using either FDR < 0.05 (for contrasts with <5000 significant CpGs) or top 5000 most significant probes; HGT = hypergeometric test; mGLM = methylglm; mRRA (GSEA) = methylRRA using gene set enrichment analysis; mRRA (ORA) = methylRRA using over-representation analysis.
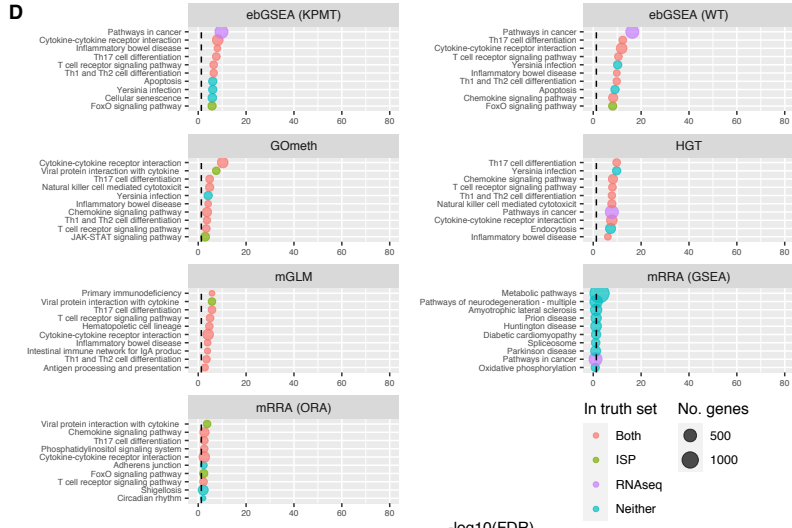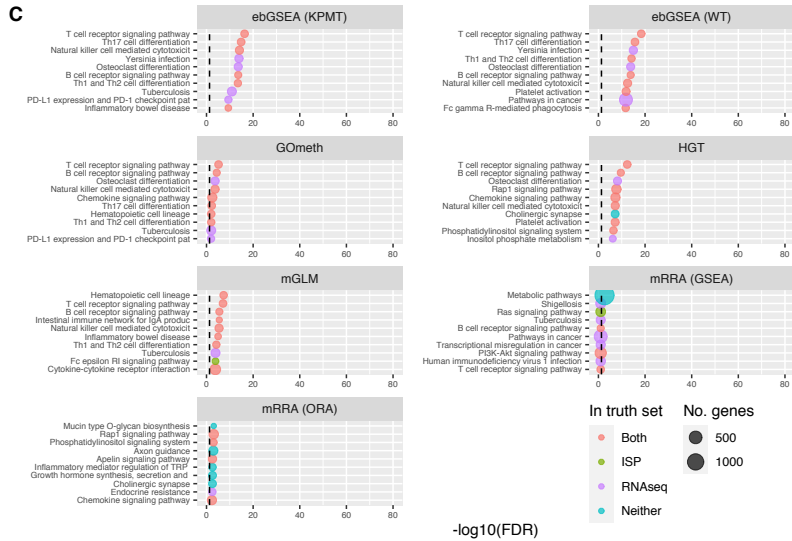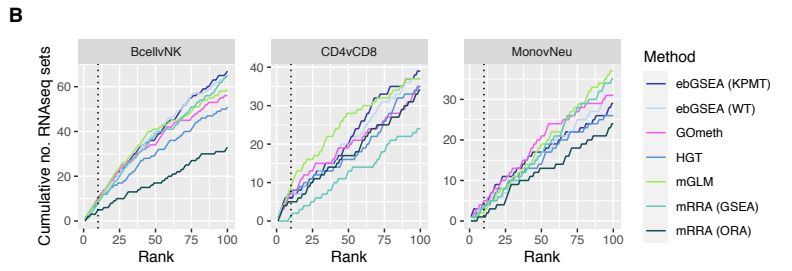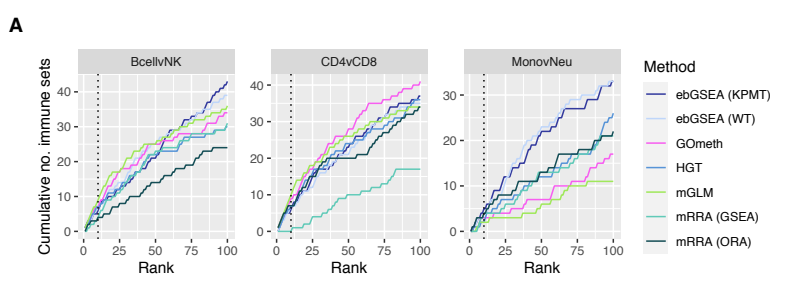
**Fig. S6. Comparison of gene set testing performance for monocytes vs. neutrophils (GO).**
**(A)** Cumulative number of GO terms, as ranked by various methods, that are present in each truth set for the monocytes vs. neutrophils comparison. ISP Terms = immune-system process child terms truth set; RNAseq Terms = top 100 terms from RNAseq analysis of the same cell types. **(B)** Bubble plots of the top 10 GO terms as ranked by various gene set testing methods for the monocytes vs. neutrophils comparison. The size of the bubble indicates the relative number of genes in the set. The colour of the bubble indicates whether the term is present in either RNAseq (purple) or ISP (green) truth sets, both (red) or neither (blue). ebGSEA (KPMT) = ebGSEA using Known Population Median Test; ebGSEA (WT) = ebGSEA using Wilcoxon Test; GOmeth = GOmeth using either FDR < 0.05 (for contrasts with <5000 significant CpGs) or top 5000 most significant probes; HGT = hypergeometric test; mGLM = methylglm; mRRA (GSEA) = methylRRA using gene set enrichment analysis; mRRA (ORA) = methylRRA using over-representation analysis.
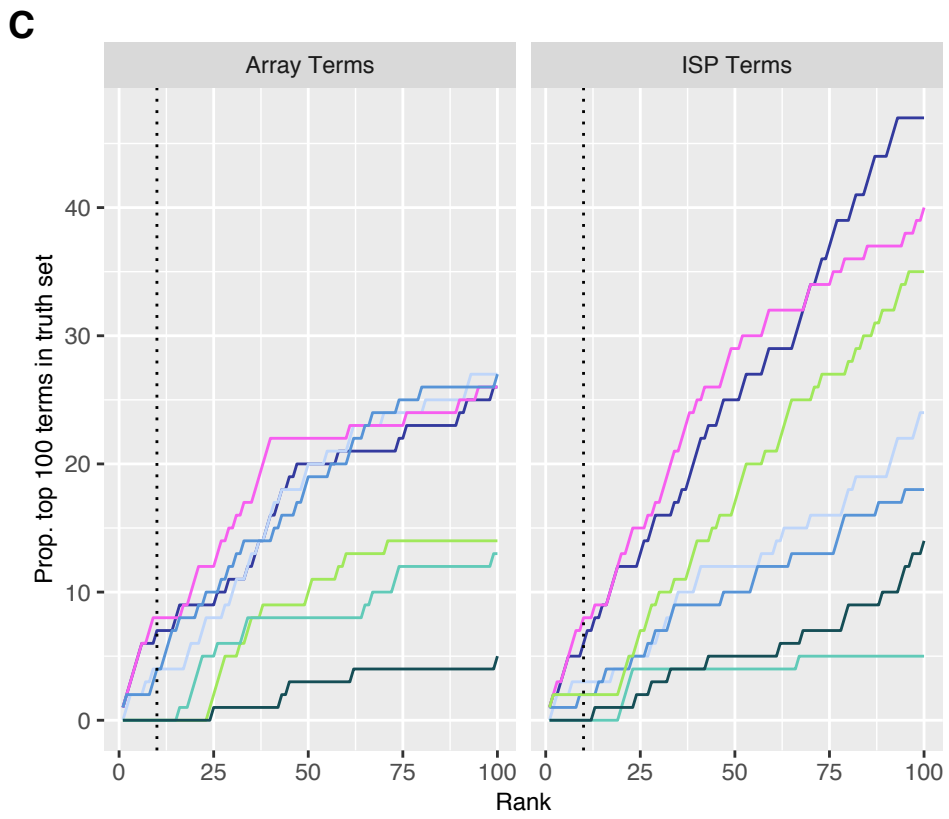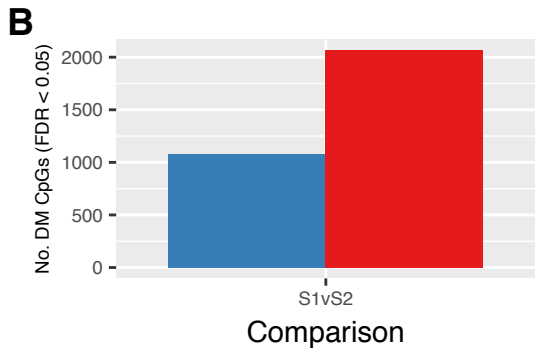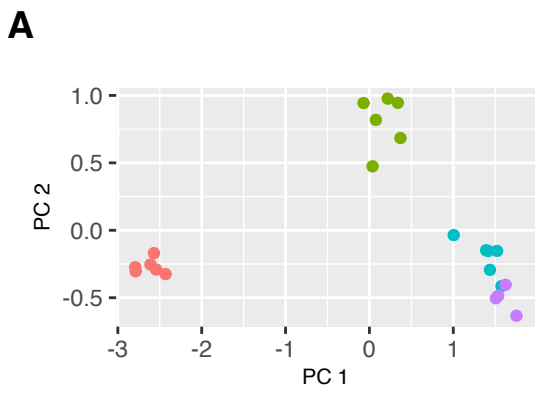
**Fig. S7. Comparison of gene set testing performance for sorted blood cell types data (KEGG). (A)** Cumulative number of KEGG pathways, as ranked by various methods, that are present in the immune truth set for all comparisons. Immune sets = all pathways belonging to the following categories: Immune system, Immune disease, Signal transduction, Signaling molecules and interaction. **(B)** Cumulative number of KEGG pathways, as ranked by various methods, that are present in the RNAseq truth set for all comparisons. RNAseq Terms = top 100 KEGG pathways from RNAseq analysis of the same cell types. **(C)** Bubble plots of the top 10 KEGG pathways as ranked by various gene set testing methods for the B-cell vs. NK comparison. The size of the bubble indicates the relative number of genes in the set. The colour of the bubble indicates whether the term is present in either RNAseq (purple) or ISP (green) truth sets, both (red) or neither (blue). **(D)** Bubble plots of the top 10 KEGG pathways as ranked by various gene set testing methods for the CD4 T-cells vs. CD8 T-cells comparison. **(E)** Bubble plots of the top 10 KEGG pathways as ranked by various gene set testing methods for the monocytes vs. neutrophils comparison. ebGSEA (KPMT) = ebGSEA using Known Population Median Test; ebGSEA (WT) = ebGSEA using Wilcoxon Test; GOmeth = GOmeth using either FDR < 0.05 (for contrasts with <5000 significant CpGs) or top 5000 most significant probes; HGT = hypergeometric test; mGLM = methylglm; mRRA (GSEA) = methylRRA using gene set enrichment analysis; mRRA (ORA) = methylRRA using over-representation analysis.

**Fig. S8. Comparison of gene set testing performance on Gene Ontology (GO) categories for B-cell development data. (A)** Multidimensional scaling plot of 450k array B-cell development data. **(B)** Numbers of differentially methylated CpGs with an adjusted p-value < 0.05, for Stage1 versus Stage 2. The blue bar is the number of significant CpGs that are less methylated in Stage 1 relative to Stage 2, and the red bar is the number that are more methylated; e.g. ~1000 CpGs are less methylated and ~2000 are more methylated in Stage 1, compared to Stage 2. **(C)** Cumulative number of GO terms, as ranked by various methods, that are present in each truth set for Stage 1 versus Stage 2. ISP Terms = immune-system process child terms truth set; Array Terms = top 100 terms from Affymetrix array analysis of the same B-cell development stages. **(D)** Bubble plots of the top 10 GO terms as ranked by various gene set testing methods. The size of the bubble indicates the relative number of genes in the set. The colour of the bubble indicates whether the term is present in either Array (purple) or ISP (green) truth sets, both (red) or neither (blue). ebGSEA (KPMT) = ebGSEA using Known Population Median Test; ebGSEA (WT) = ebGSEA using Wilcoxon Test; GOmeth = GOmeth using either FDR < 0.05 (for contrasts with <5000 significant CpGs) or top 5000 most significant probes; HGT = hypergeometric test; mGLM = methylglm; mRRA (GSEA) = methylRRA using gene set enrichment analysis; mRRA (ORA) = methylRRA using over-representation analysis.
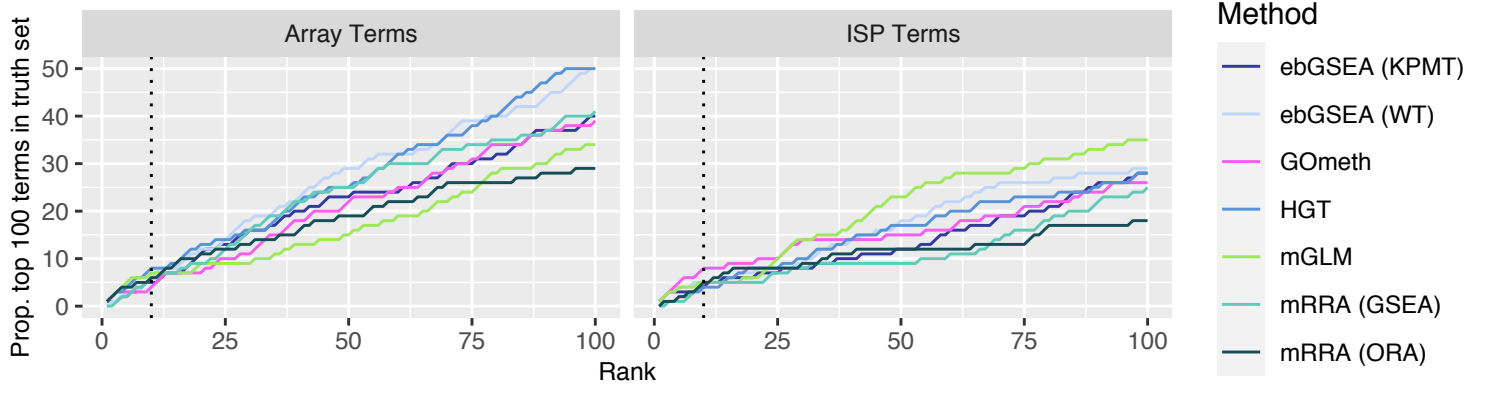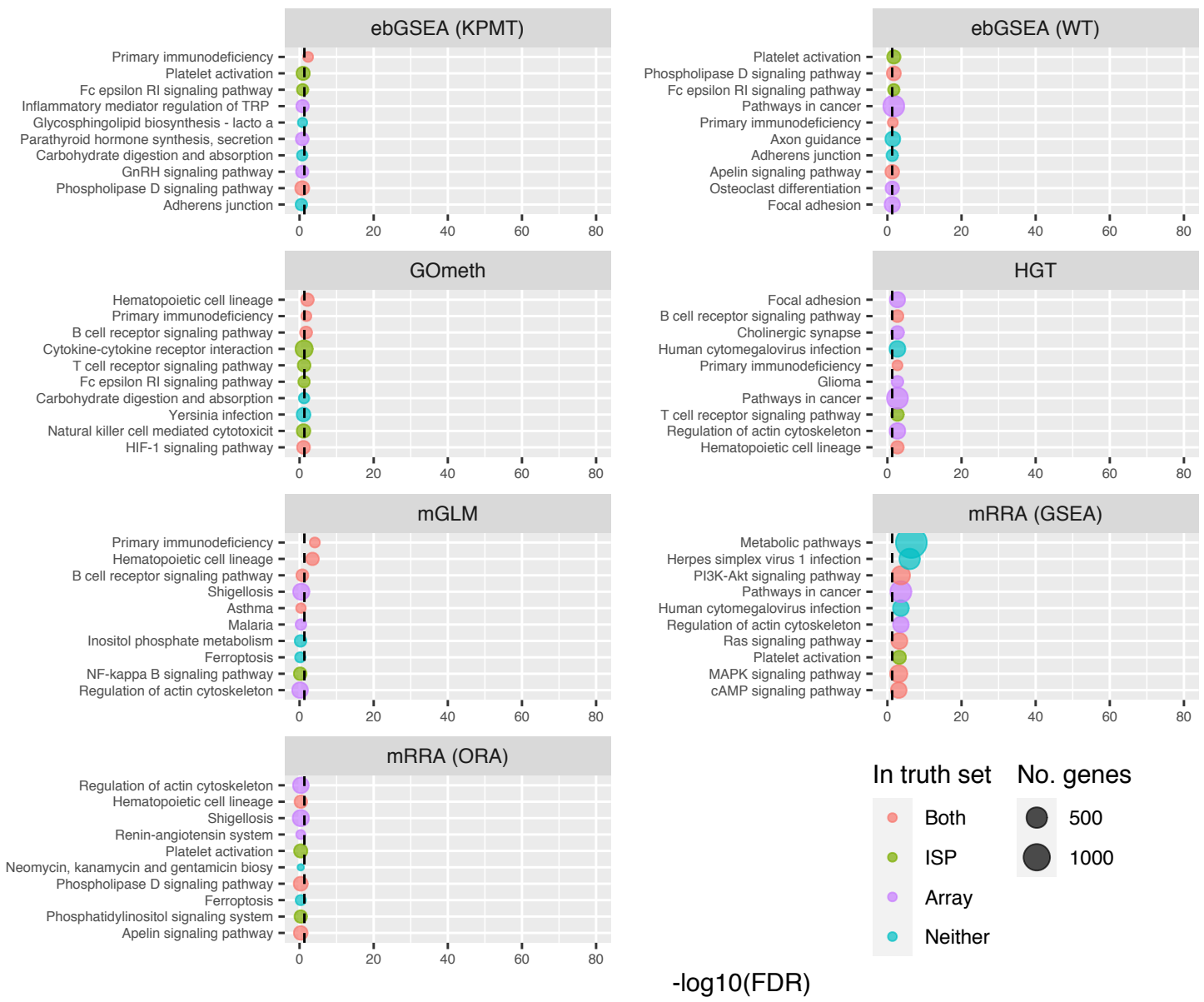
**Fig. S9. Comparison of gene set testing performance for Stage 1 vs. Stage 2 of B-cell development (KEGG). (A)** Cumulative number of KEGG terms, as ranked by various methods, that are present in each truth set for Stage 1 versus Stage 2. ISP Terms = immune-system process child terms truth set; Array Terms = top 100 terms from Affymetrix array analysis of Stage 1 vs. Stage 2. **(B)** Bubble plots of the top 10 KEGG pathways as ranked by various gene set testing methods. The size of the bubble indicates the relative number of genes in the set. The colour of the bubble indicates whether the term is present in either Array (purple) or ISP (green) truth sets, both (red) or neither (blue). ebGSEA (KPMT) = ebGSEA using Known Population Median Test; ebGSEA (WT) = ebGSEA using Wilcoxon Test; GOmeth = GOmeth using either FDR < 0.05 (for contrasts with <5000 significant CpGs) or top 5000 most significant probes; HGT = hypergeometric test; mGLM = methylglm; mRRA (GSEA) = methylRRA using gene set enrichment analysis; mRRA (ORA) = methylRRA using over-representation analysis.
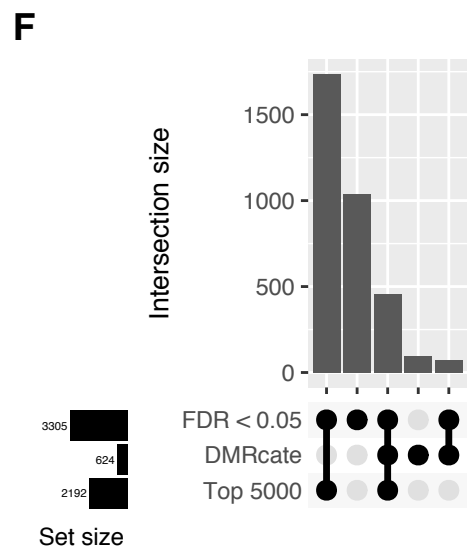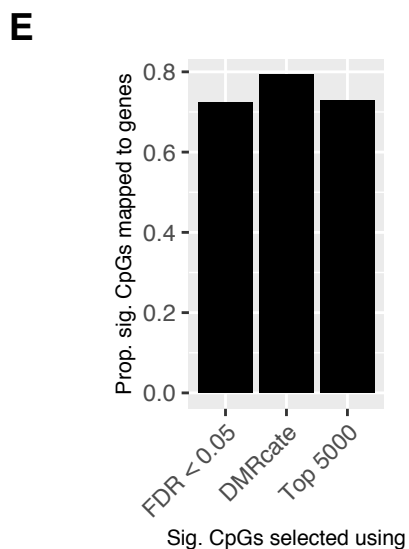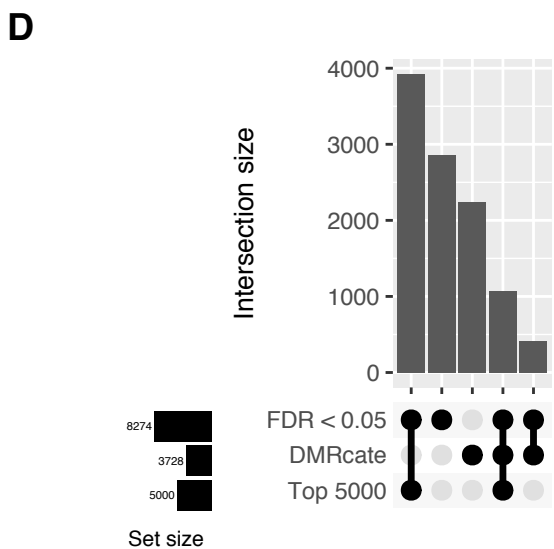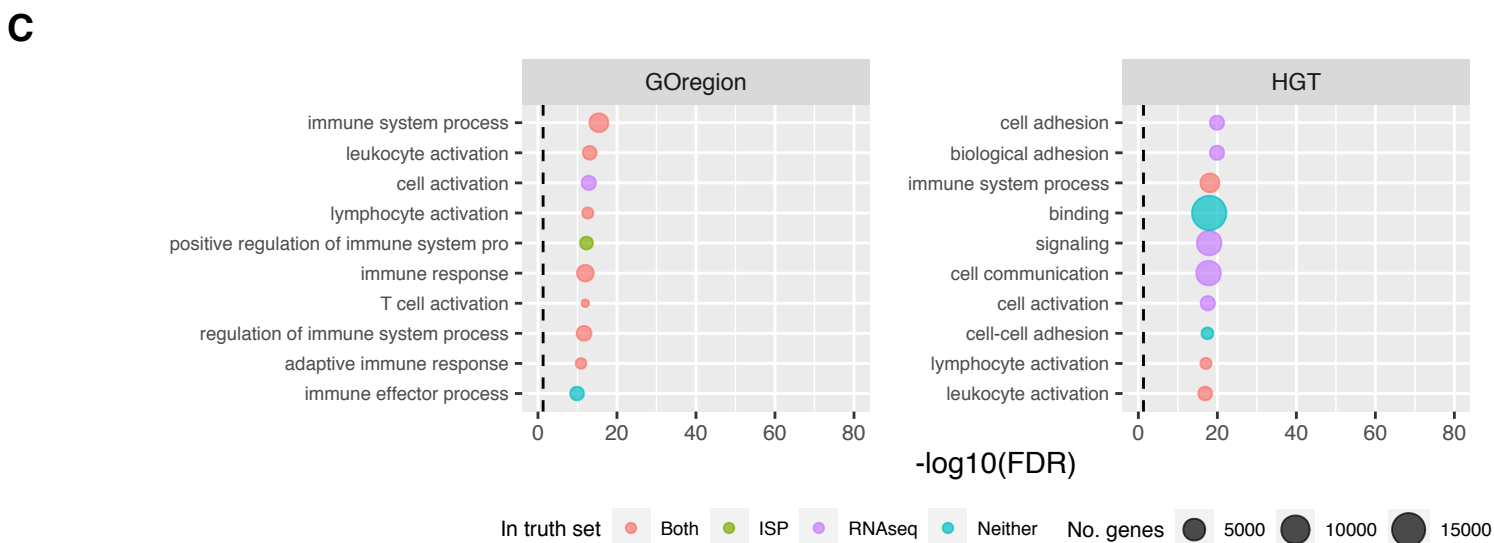
**Fig. S10. Evaluation of the performance of GOregion for CD4 T-cells vs. CD8 T-cells. (A)** Bias plot showing that genes that have more CpGs measuring methylation are more likely to have a differentially methylated region. This plot is produced from EPIC array sorted blood cell type data, comparing CD4 T-cells vs. CD8 T-cells. **(B)** Cumulative number of GO terms, as ranked by GOregion and a simple hypergeometric test (HGT), that are present in each truth set for the CD4 T-cells vs. CD8 T-cells comparison. ISP Terms = immune-system process child terms truth set; RNAseq Terms = top 100 terms from RNAseq analysis of the same cell types. **(C)** Bubble plots of the top 10 GO terms as ranked by GOregion and a simple HGT for the CD4 T-cells vs. CD8 T-cells comparison. The size of the bubble indicates the relative number of genes in the set. The colour of the bubble indicates whether the term is present in either RNAseq (purple) or ISP (green) truth sets, both (red) or neither (blue). **(D)** Upset plot showing the characteristics of th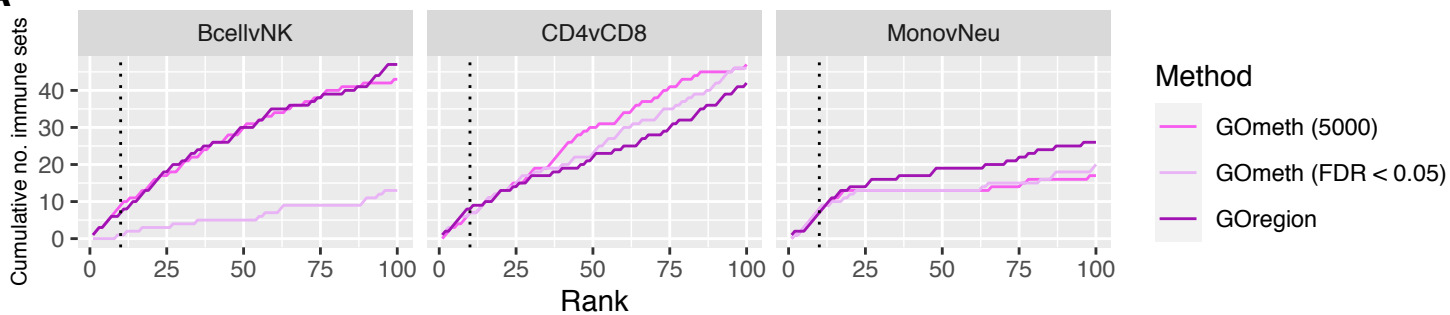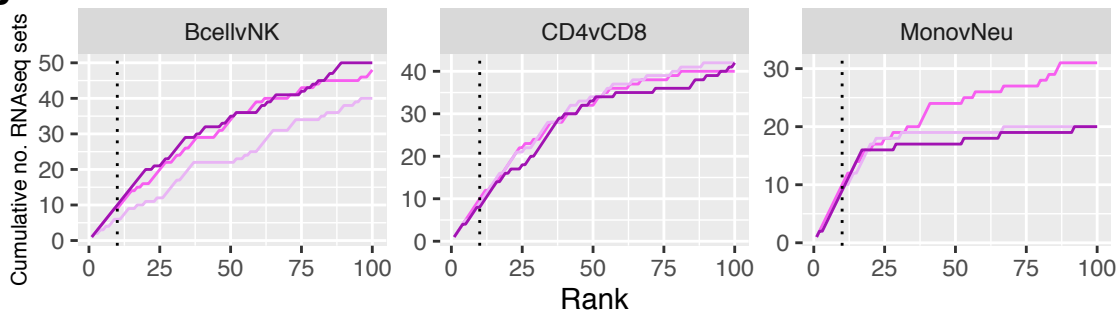e CpGs selected as "significant" for the CD4 T-cells vs. CD8 T-cells comparison by a probe-wise differential methylation analysis using a significance cut off (FDR < 0.05), the top 5000 CpGs as ranked by the probe-wise analysis (Top 5000) or the CpGs underlying the filtered *DMRcate* regions (DMRcate). The probe-wise analysis with FDR < 0.05 identified over 8,000 CpGs as "significant" and had the most unique CpGs. Although DMRcate identified the fewest "significant" CpGs (~3,700), over 2,000 were unique to that approach. **(E)** Proportion of "significant" CpGs that are annotated to genes as identified by the three different strategies. **(F)** Upset plot showing the characteristics of the genes that "significant" CpGs are annotated to, as identified by the three different strategies, for the CD4 T-cells vs. CD8 T-cells comparison. CpGs identified by the probe-wise analysis with FDR < 0.05 map to over 3,000 genes. The CpGs identified by DMRcate mapped to ~600 genes, several of which are unique to this approach.
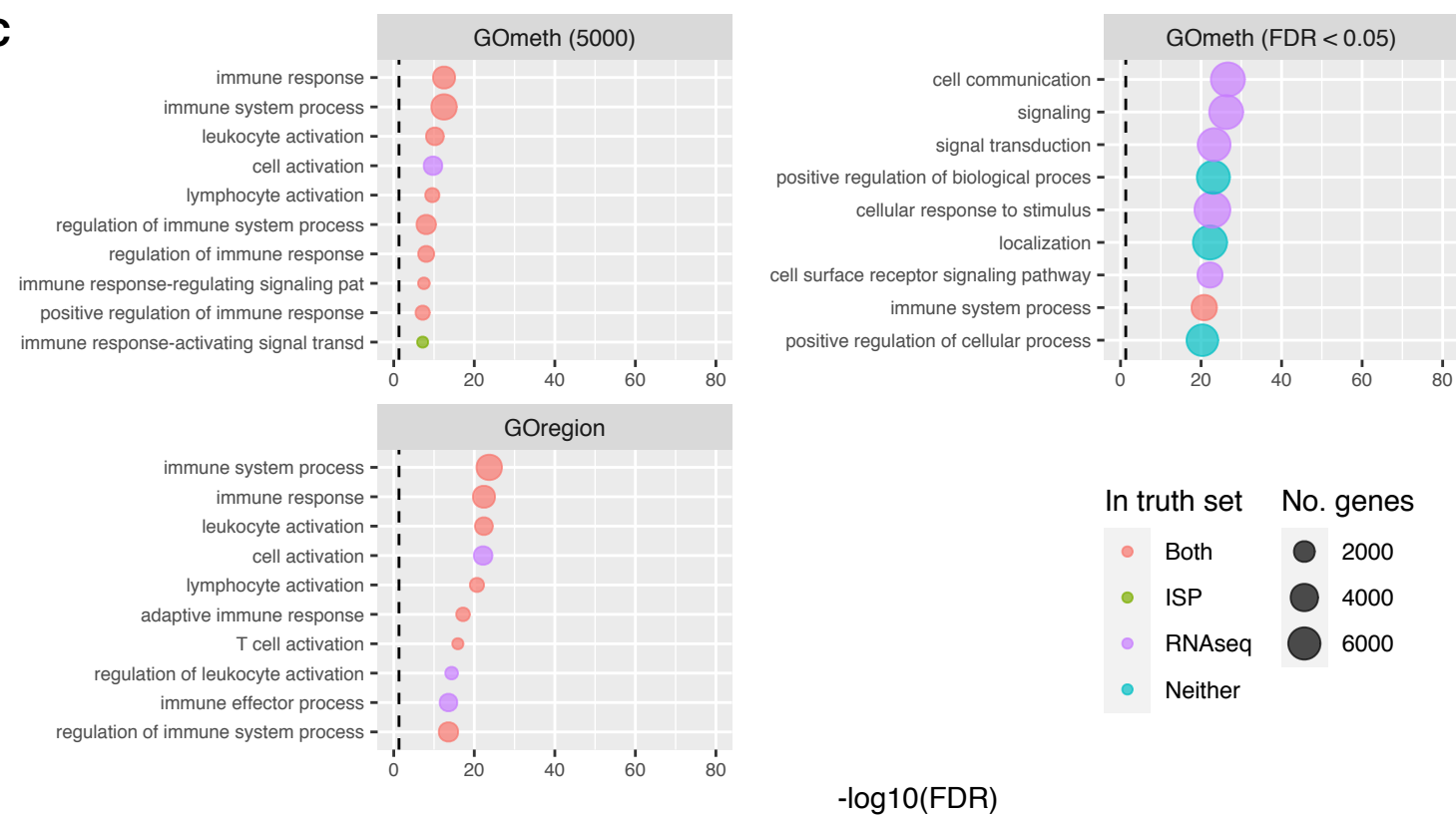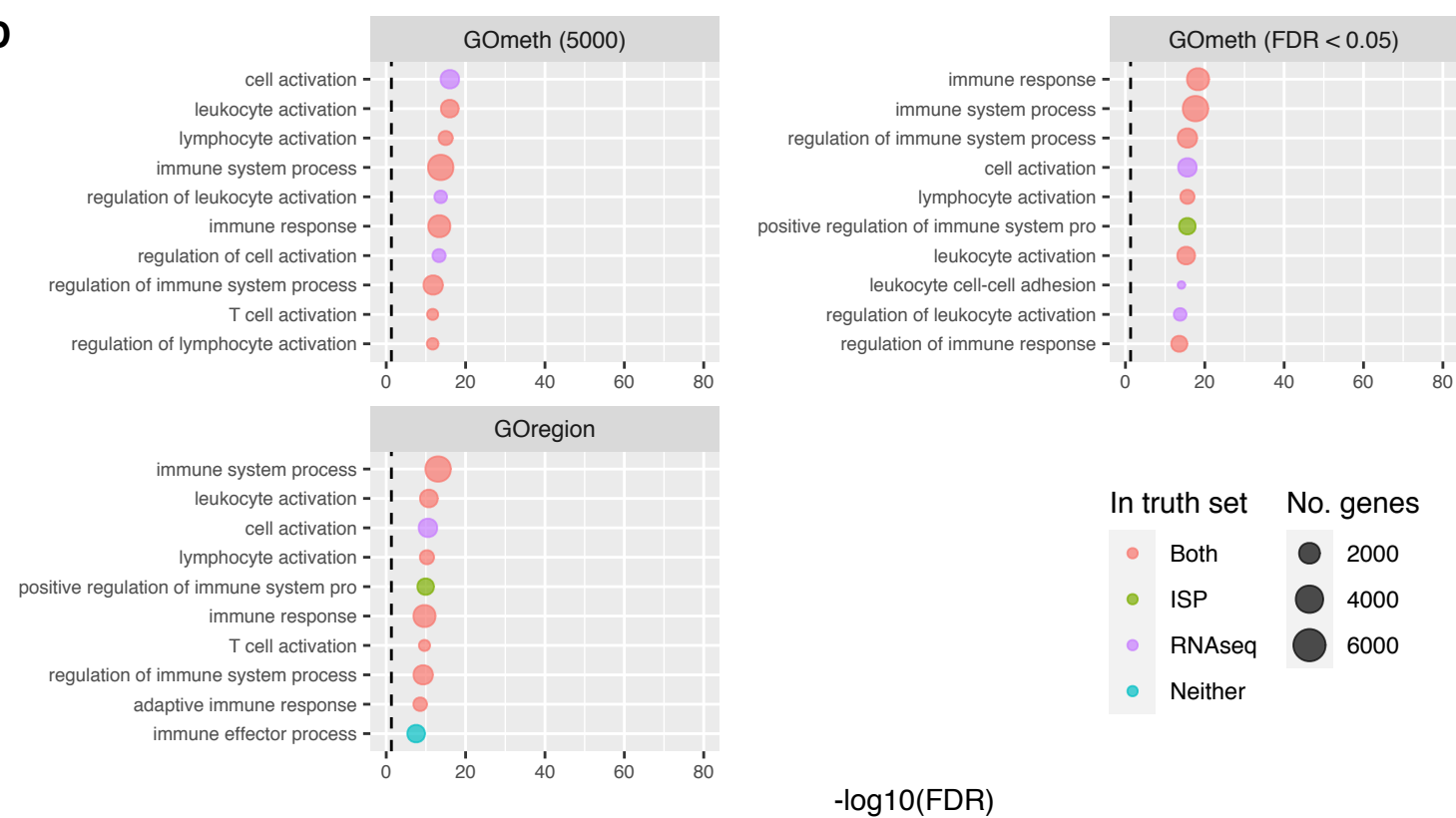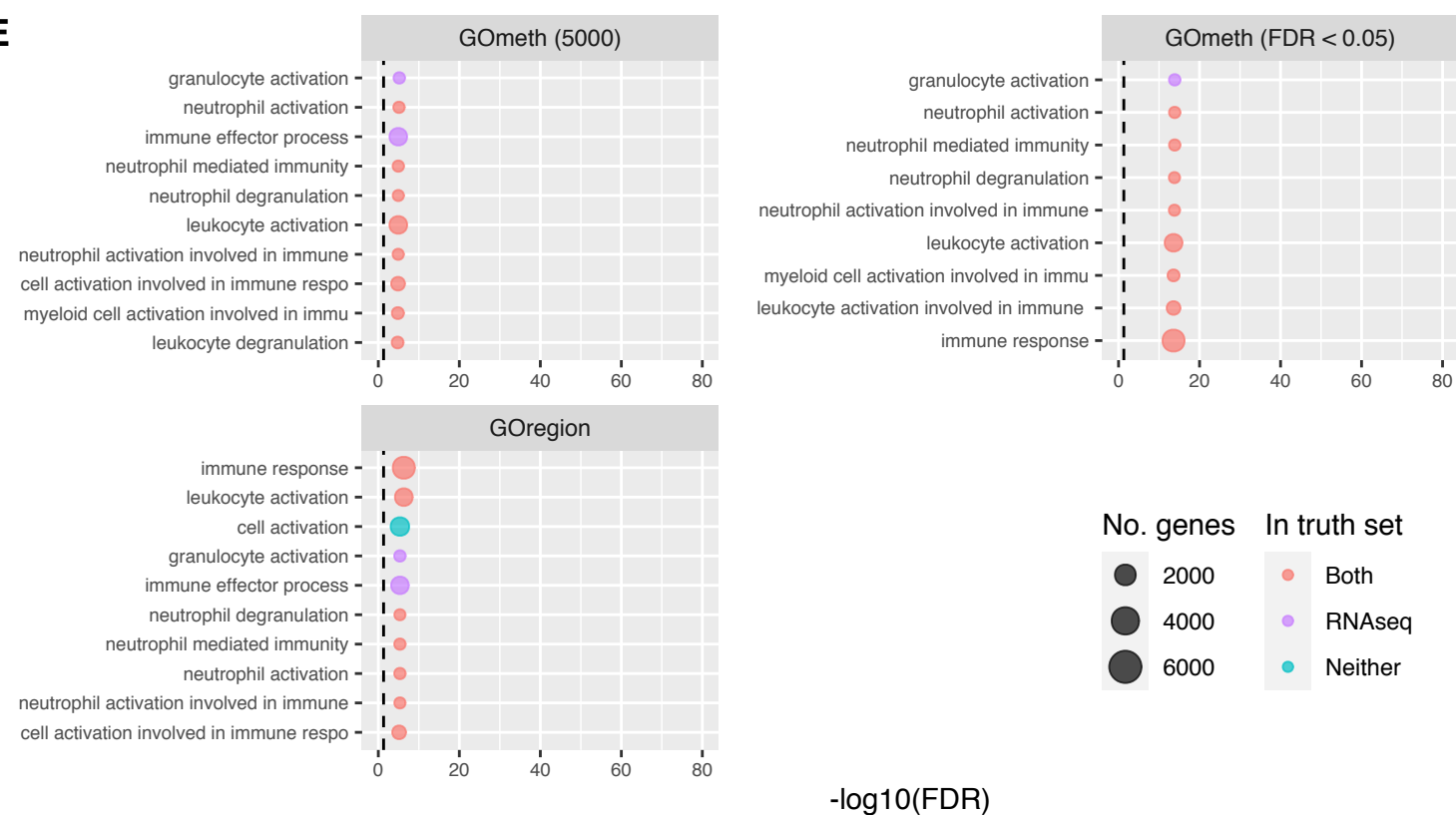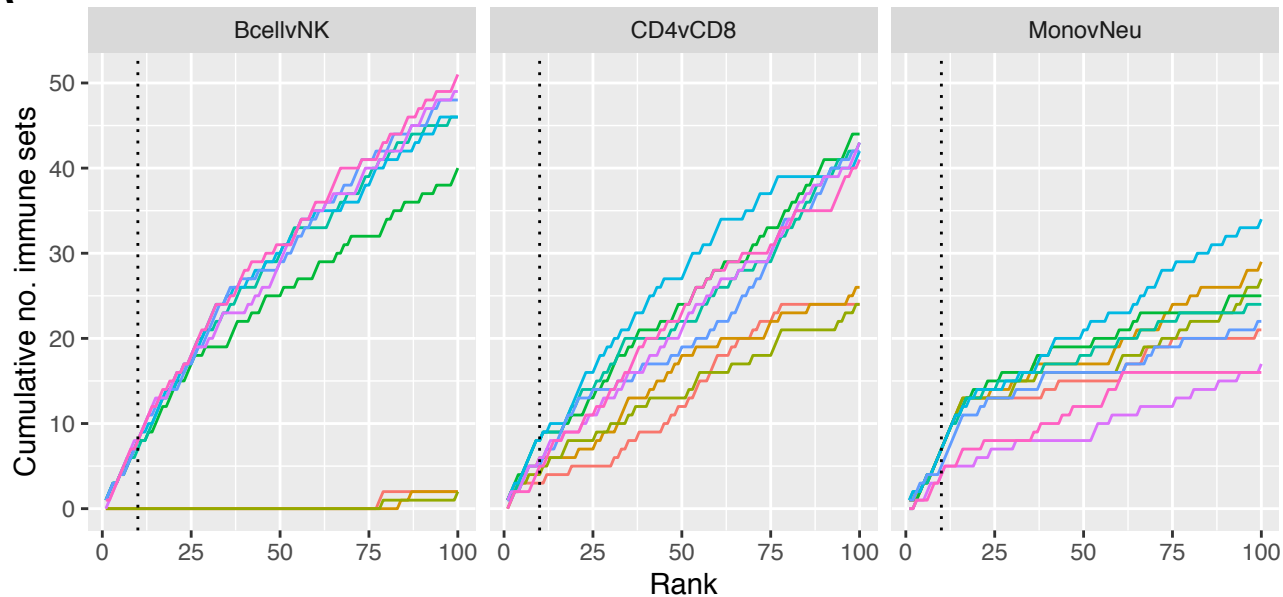
**Fig. S11. Evaluation of the performance of GOregion for monocytes vs. neutrophils. (A)** Bias plot showing that genes that have more CpGs measuring methylation are more likely to have a differentially methylated region. This plot is produced from EPIC array sorted blood cell type data, comparing monocytes vs. neutrophils. **(B)** Cumulative number of GO terms, as ranked by GOregion and a simple hypergeometric test (HGT), that are present in each truth set for the monocytes vs. neutrophils comparison. ISP Terms = immune-system process child terms truth set; RNAseq Terms = top 100 terms from RNAseq analysis of the same cell types. **(C)** Bubble plots of the top 10 GO terms as ranked by GOregion and a simple HGT for the monocytes vs. neutrophils comparison. The size of the bubble indicates the relative number of genes in the set. The colour of the bubble indicates whether the term is present in either RNAseq (purple) or ISP (green) truth sets, both (red) or neither (blue). **(D)** Upset plot showing the characteristics of the CpGs selected as "significant" for the monocytes vs. neutrophils comparison by a probe-wise differential methylation analysis using a significance cut off (FDR < 0.05), the top 5000 CpGs as ranked by the probe-wise analysis (Top 5000) or the CpGs underlying the filtered *DMRcate* regions (DMRcate). The probe-wise analysis with FDR < 0.05 identified over 21,000 CpGs as "significant" and had the most unique CpGs. Although DMRcate identified fewer CpGs (~7500), ~3500 of them were unique to that approach. **(E)** Proportion of "significant" CpGs that are annotated to genes as identified by the three different strategies. **(F)** Upset plot showing the characteristics of the genes that "significant" CpGs are annotated to, as identified by the three different strategies, for the monocytes vs. neutrophils comparison. CpGs identified by the probe-wise analysis with FDR < 0.05 map to over 6,000 genes. The CpGs identified by DMRcate mapped to ~1200, several of which are unique to this approach.

**Fig. S12. Comparison of probe-wise and region-wise analyses in the gene set testing context for sorted blood cell types data. (A)** Cumulative number of GO terms, as ranked by various analysis strategies, that are present in the immune truth set for all comparisons. ISP Terms = immune-system process child terms truth set. **(B)** Cumulative number of GO terms, as ranked by various methods, that are present in the RNAseq truth set for all comparisons. RNAseq Terms = top 100 terms from RNAseq analysis of the same cell types. **(C)** Bubble plots of the top 10 GO terms as ranked by various analysis strategies for the B-cell vs. NK comparison. The size of the bubble indicates the relative number of genes in the set. The colour of the bubble indicates whether the term is present in either RNAseq (purple) or ISP (green) truth sets, both (red) or neither (blue). **(D)** Bubble plots of the top 10 GO terms as ranked by various analysis strategies for the CD4 T-cells vs. CD8 T-cells comparison. **(E)** Bubble plots of the top 10 GO terms as ranked by various analysis strategies for the monocytes vs. neutrophils comparison. GOmeth (5000) = GOmeth using top 5000 most significant probes; GOmeth (FDR < 0.05) = GOmeth using all probes significant at FDR < 0.05; GOregion = GOregion using all selected region coordinates.

**A**

**B**

**Fig. S13. Effect of DMR filtering on gene set testing for sorted blood cell types data. (A)**
Cumulative number of GO terms, as ranked by GOregion analysis of DMRs filtered using
various parameters, that are present in the immune truth set for all comparisons. ISP Terms =
immune-system process child terms truth set. **(B)** Cumulative number of GO terms, as ranked by
GOregion analysis of DMRs filtered using various parameters, that are present in the RNAseq
truth set for all comparisons. RNAseq Terms =  top 100 terms from RNAseq analysis of the same
cell types. $|\Delta\beta|$ = mean methylation difference across region; No. CpGs = number of CpGs
underlying region.

**Fig. S14. Evaluation of the performance of GOregion on B-cell development data. (A)** Bias plot showing that genes with more measured CpGs are more likely to have a differentially methylated region (DMR). This plot is produced from 450K array B-cell development data, comparing Stage 1 vs. Stage 2. **(B)** Cumulative number of GO terms, as ranked by GOregion and a simple hypergeometric test (HGT), that are present in each truth set for the Stage 1 versus Stage 2 comparison. ISP Terms = immune-system process child terms truth set; Array Terms = top 100 terms from Affymetrix array analysis of the same B-cell development stages. **(C)** Bubble plots of the top 10 GO terms as ranked by GOregion and a simple HGT for the Stage 1 versus Stage 2 comparison. The size of the bubble indicates the relative number of genes in the set. The colour of the bubble indicates whether the term is present in either Array (purple) or ISP (green) truth sets, both (red) or neither (blue). **(D)** Upset plot showing the characteristics of the CpGs selected as "significant" for the Stage 1 versus Stage 2 comparison by a probe-wise differential methylation analysis using a significance cut off (FDR < 0.05), the top 5000 CpGs as ranked by the probe-wise analysis (Top 5000) or the CpGs underlying the filtered *DMRcate* regions (DMRcate). The *DMRcate* analysis identified over 11,000 CpGs as "significant" and had the most unique CpGs (>7500). **(E)** Proportion of "significant" CpGs that are annotated to genes as identified by the three different strategies. **(F)** Upset plot showing the characteristics of the genes that "significant" CpGs are annotated to, as identified by the three different strategies, for the Stage 1 versus Stage 2 comparison.

**A**

Array Terms | ISP Terms

Cumulative no. truth sets

Rank

Parameters
- GOmeth (5000)
- GOmeth (FDR < 0.05)
- GOregion

**B**

GOmeth (5000)

immune response
cell activation
leukocyte activation
immune system process
lymphocyte activation
adaptive immune response
immune effector process
regulation of immune system process
cytokine production
regulation of cytokine production

GOmeth (FDR < 0.05)

immune response
immune system process
cell activation
leukocyte activation
lymphocyte activation
immune effector process
regulation of immune system process
adaptive immune response
positive regulation of immune system pro
regulation of immune response

GOregion

cell activation
leukocyte activation
immune response
immune system process
adaptive immune response
lymphocyte activation
cytokine production
regulation of cytokine production
regulation of immune system process
inflammatory response

No. genes
- 1000
- 1500
- 2000
- 2500
- 3000

In truth set
- Both
- ISP
- Array
- Neither

-log10(FDR)

**C**

Array Terms | ISP Terms

Cumulative no. truth sets

Rank

Parameters
- $|\Delta\beta| = 0$; No. CpGs = 2
- $|\Delta\beta| = 0$; No. CpGs = 3
- $|\Delta\beta| = 0$; No. CpGs = 4
- $|\Delta\beta| = 0.1$; No. CpGs = 2
- $|\Delta\beta| = 0.1$; No. CpGs = 3
- $|\Delta\beta| = 0.1$; No. CpGs = 4
- $|\Delta\beta| = 0.2$; No. CpGs = 2
- $|\Delta\beta| = 0.2$; No. CpGs = 3
- $|\Delta\beta| = 0.2$; No. CpGs = 4

**Fig. S15. Comparison of probe-wise and region-wise analyses in the gene set testing context using B-cell development data. (A)** Cumulative number of GO terms, as ranked by various analysis strategies, that are present in each truth set for Stage 1 vs. Stage 2. ISP Terms = immune-system process child terms truth set; Array Terms = top 100 terms from Affymetrix array analysis of Stage 1 vs. Stage 2. GOmeth (5000) = GOmeth using top 5000 most significant probes; GOmeth (FDR < 0.05) = GOmeth using all probes significant at FDR < 0.05; GOregion = GOregion using all selected region coordinates. **(B)** Bubble plots of the top 10 GO terms as ranked by various analysis strategies for the Stage 1 vs. Stage 2 comparison. The size of the bubble indicates the relative number of genes in the set. The colour of the bubble indicates whether the term is present in either Array (purple) or ISP (green) truth sets, both (red) or neither (blue). **(C)** Effect of DMR filtering on gene set testing. Cumulative number of GO terms, as ranked by GOregion analysis of DMRs filtered using various parameters, that are present in each truth set for the Stage 1 vs. Stage 2 comparison. ISP Terms = immune-system process child terms truth set. Array Terms = top 100 terms from Affymetrix array analysis of the same B-cell development stages. $|\Delta\beta|$ = mean methylation difference across region; No. CpGs = number of CpGs underlying region.