

Discovery of Cancer Driver Long Noncoding RNAs across 1112 Tumour Genomes: New Candidates and Distinguishing Features.

Authors

Andrés Lanzós^{1,2,3}, Joana Carlevaro-Fita^{1,2,3}, Loris Mularoni⁴, Ferran Reverter^{1,2,3}, Emilio Palumbo^{1,2,3}, Roderic Guigó^{1,2,3}, Rory Johnson^{1,2,3,5*}

1. Centre for Genomic Regulation (CRG), The Barcelona Institute of Science and Technology, Dr. Aiguader 88, Barcelona 08003, Spain

2. Universitat Pompeu Fabra (UPF), Barcelona, Spain.

3. Institut Hospital del Mar d'Investigacions Mèdiques (IMIM), 08003 Barcelona, Spain.

4. Research Unit on Biomedical Informatics, Department of Experimental and Health Sciences, Universitat Pompeu Fabra, Dr. Aiguader 88, Barcelona, Spain

5. Present address: Department of Clinical Research, University of Bern, Murtenstrasse 35, 3010 Bern, Switzerland.

* Correspondence to rory.johnson@dkf.unibe.ch

Supplementary Figure Legends

Supplementary Figure S1: Graphical representation of the ExInAator workflow. The objective of this algorithm is the creation of a local background region with identical trinucleotide content to the exons of a given gene. As an example we take a hypothetical single-exon lncRNA gene. (A) Upstream and downstream regions within the window size are defined as background regions. All the trinucleotides in the exon region are counted. The mutations (indicated by *) in the central nucleotide positions of each trinucleotide are counted. Colours indicate the four nucleotides in the exon. Background region nucleotides are not coloured. For brevity, only four trinucleotide combinations are shown. (B) Next, an iterative sampling is performed. A set of trinucleotides, identical to the exonic trinucleotide set, are randomly sampled from the background region. This operation is repeated, without replacement of the central nucleotide (red "X"), until it is no longer possible to extract the appropriate set of trinucleotides. (C) The total number of nucleotide positions and mutations are counted for the exon region and the sampled background region. These values are used for the contingency table analysis shown in Figure 1.

Supplementary Figure S2: Ratio of exonic to background mutation density.

Supplementary Figure S3: Quantile-quantile ("QQ") plots for all analyses. Grey dots represent the P values obtained in the two simulations. The black and coloured dots correspond to the P values calculated in the real dataset.

Supplementary Figure S4: Histogram of predicted candidates in Breast across the 100 iterations of shuffling analysis. First, all lncRNAs randomly placed in a new position within the whole genome, while maintaining exon-intron structure. Then ExInAator was run, with same settings as described, using real mutations from Breast. Shown are the number of predicted candidates at a cutoff of FDR<0.1. The number of predicted candidates using real gene locations is 3, greater than 97% of simulations.

Supplementary Figure S5: Quantile-quantile (“QQ”) plot of “intronic gene” simulation in Breast. Artificial “intronic” gene models were constructed, where the exons of every gene were replaced by equally-sized, randomly-selected fragments of introns from the same gene. These genes were used as the input for ExInAator and run (with same settings as for real data) on Breast samples.

Supplementary Figure S6: Expression of SAMMSON in Stomach (in which is detected by ExInAator), and skin, sun-exposed or not (in which has been previously implicated in Melanoma).

Supplementary Figure S7: Exon-level mutational density of all lncRNAs candidates against the background regions.

Supplementary Figure S8: Intersection between ExInAator, MiTranscriptome and Du et al candidates.

Supplementary Figure S9: Performance of ExInAator compared MutSig, OncodriveFM and OncodriveClust in the Pancancer Alexandrov dataset. (A) Percentage of CGC candidates amongst predicted drivers for each method across different cutoffs ($-\log_{10} Q$ value) in all the Alexandrov cancer files. (B) As for (A), but displayed by ranked number of top predicted drivers.

Supplementary Figure S10: The overlap of predicted protein-coding drivers between methods. The percent of predicted drivers of each method on Pancancer Alexander mutations that are also predicted by at least one of the other methods displayed by (A) cutoff or by (B) ranked number of candidates.

Supplementary Figure S11: QQ plots of each method on the Pancancer Alexandrov dataset. Grey dots represent the P values obtained in a simulation. The black and coloured dots correspond to the P values calculated in the real dataset.

Supplementary Figure S12: Tumour specificity of lncRNA and protein coding genes. The figure shows the percent of genes of each biotype that are mutated in one or more distinct tumour samples.

Supplementary Figure S13: Screenshots of the 12 candidate’s mutations analysed with IGV.

Supplementary Figure S14: Replication timing of cancer lncRNAs, ExInAator predictions and other lncRNAs. Cumulative distribution are shown for S50 values in HeLa cells assigned to lncRNA gene sets. S50 ratio is defined as the fraction of the S phase at which 50% of the DNA is replicated (see Materials and Methods). Low ratios indicate early replication timing. Genes are grouped by: literature-reported cancer CRL lncRNAs; all ExInAator candidates (both CRL and not); non-CRL novel ExInAator candidates; non-CRL non-candidates, being all other GENCODE lincRNAs. Candidates here were defined at a threshold of $Q \leq 0.2$. Dashed vertical lines indicate the mean value of each group.

Supplementary File Legends

Supplementary File S1: CRL lncRNAs.

Supplementary File S2: Non-CRL lncRNAs.

Supplementary File S3: The full set of lncRNA drivers predicted by ExInAator in CRL.

Supplementary File S4: The full set of lncRNA drivers predicted by ExInAator in non-CRL.

Supplementary File S5: The full set of protein-coding drivers predicted by ExInAator.

Supplementary File S6: The full set of results for lncRNAs and protein-coding genes described in this study.

Supplementary File S7: Candidates identified by Du et al. that are also considered in our study.

Supplementary File S8: Candidates identified by MiTranscriptome that are also considered in our study.

Supplementary File S9: The full set of CGC drivers analysed in this study.

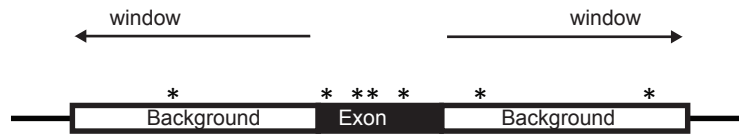
Supplementary File S10: The set of IDs for protein-coding genes used in this study.

Supplementary File S11: Protein-coding overlap with CGC across all cancers and both datasets.

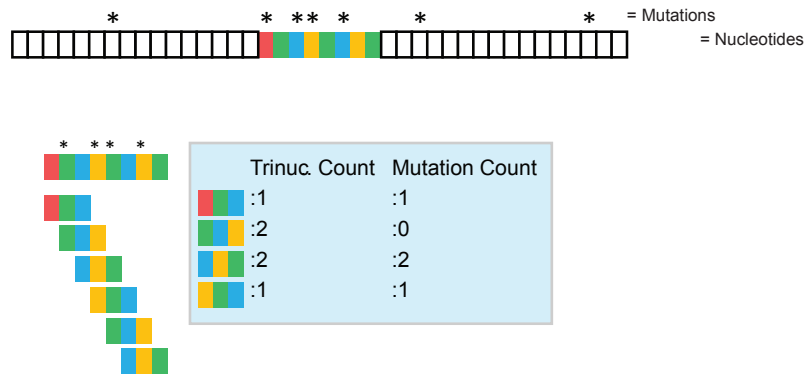
Supplementary File S12: Effects on ExInAator predictions of various filtering schemes.

Supplementary Figure S1

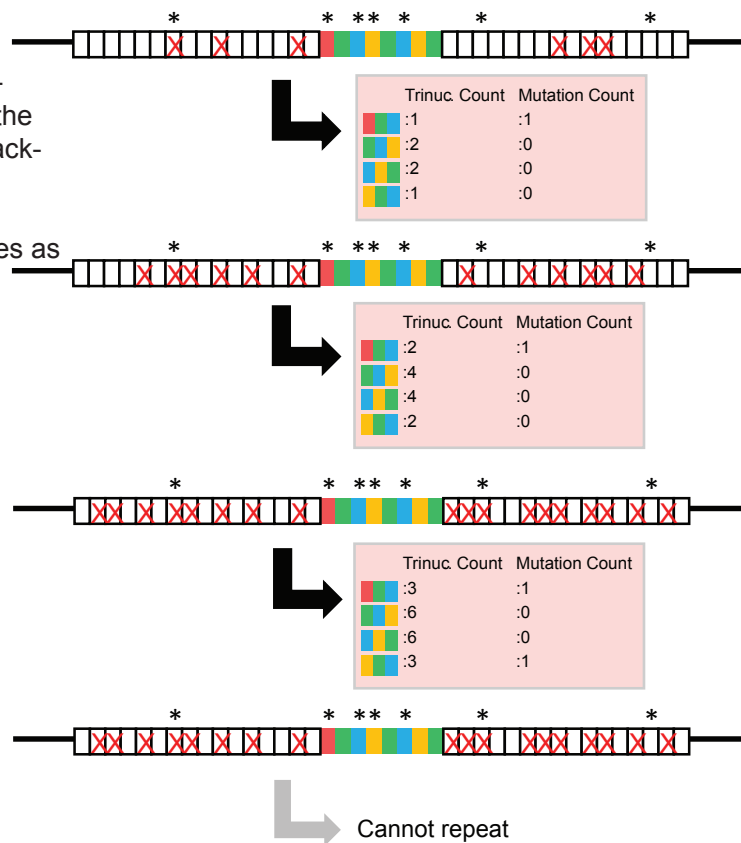
Example: Single-exon lncRNA



A. Count exon trinucleotides and their central mutations



B. Sample, without replacement, equal proportions of the same trinucleotides from background regions. Count the mutations associated with these. Repeat as many times as possible.

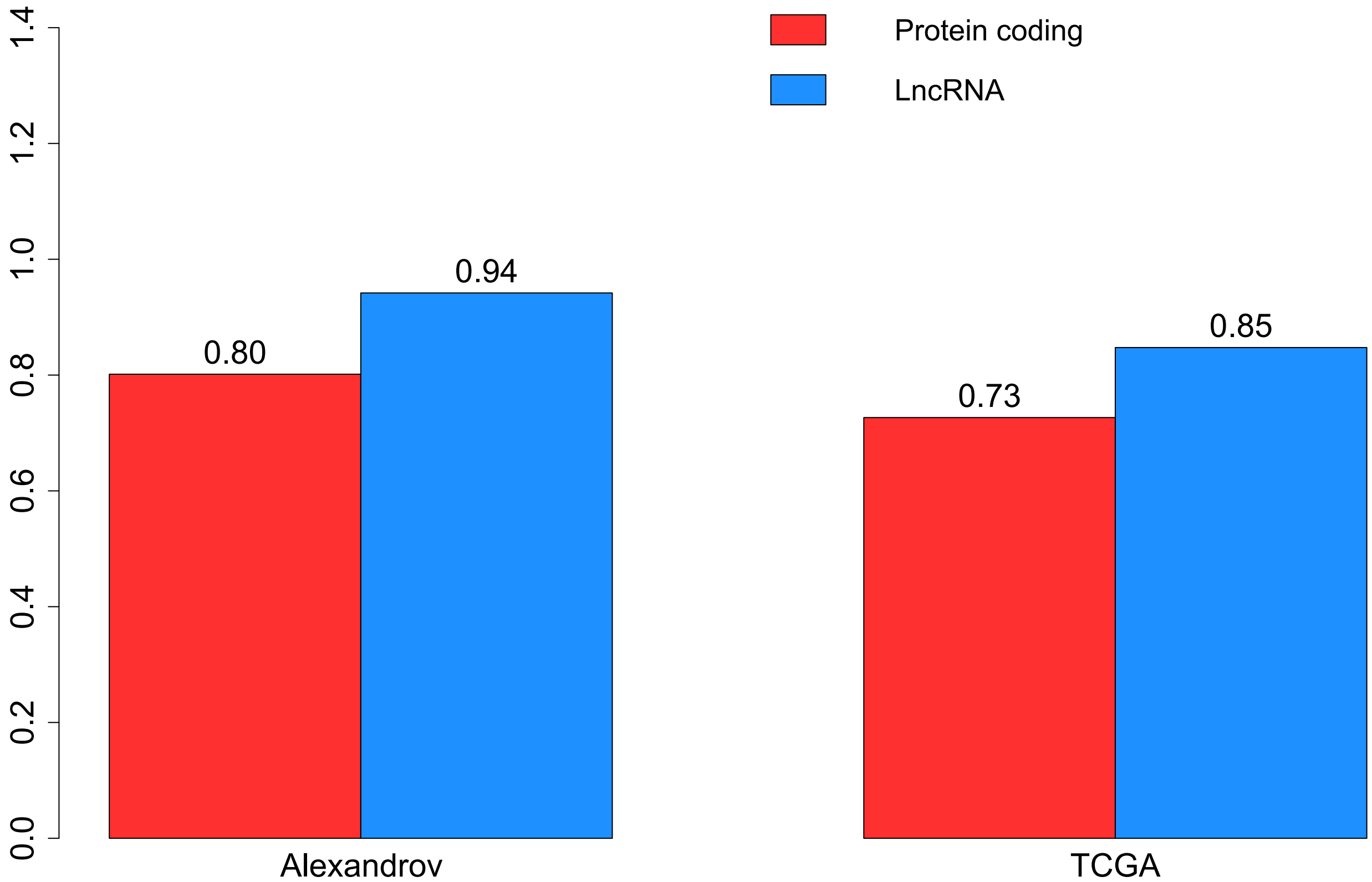


C. Compare mutation counts between exon and background regions.

	Exon	Background																														
	<table border="1"> <thead> <tr> <th>Trinucleotide</th> <th>Trinuc. Count</th> <th>Mutation Count</th> </tr> </thead> <tbody> <tr> <td>Red-Blue-Green</td> <td>:1</td> <td>:1</td> </tr> <tr> <td>Green-Blue-Yellow</td> <td>:2</td> <td>:0</td> </tr> <tr> <td>Blue-Yellow-Green</td> <td>:2</td> <td>:2</td> </tr> <tr> <td>Yellow-Green-Blue</td> <td>:1</td> <td>:1</td> </tr> </tbody> </table>	Trinucleotide	Trinuc. Count	Mutation Count	Red-Blue-Green	:1	:1	Green-Blue-Yellow	:2	:0	Blue-Yellow-Green	:2	:2	Yellow-Green-Blue	:1	:1	<table border="1"> <thead> <tr> <th>Trinucleotide</th> <th>Trinuc. Count</th> <th>Mutation Count</th> </tr> </thead> <tbody> <tr> <td>Red-Blue-Green</td> <td>:3</td> <td>:1</td> </tr> <tr> <td>Green-Blue-Yellow</td> <td>:6</td> <td>:0</td> </tr> <tr> <td>Blue-Yellow-Green</td> <td>:6</td> <td>:0</td> </tr> <tr> <td>Yellow-Green-Blue</td> <td>:3</td> <td>:1</td> </tr> </tbody> </table>	Trinucleotide	Trinuc. Count	Mutation Count	Red-Blue-Green	:3	:1	Green-Blue-Yellow	:6	:0	Blue-Yellow-Green	:6	:0	Yellow-Green-Blue	:3	:1
Trinucleotide	Trinuc. Count	Mutation Count																														
Red-Blue-Green	:1	:1																														
Green-Blue-Yellow	:2	:0																														
Blue-Yellow-Green	:2	:2																														
Yellow-Green-Blue	:1	:1																														
Trinucleotide	Trinuc. Count	Mutation Count																														
Red-Blue-Green	:3	:1																														
Green-Blue-Yellow	:6	:0																														
Blue-Yellow-Green	:6	:0																														
Yellow-Green-Blue	:3	:1																														
Total	N = 6 M = 4	n = 18 m = 2																														

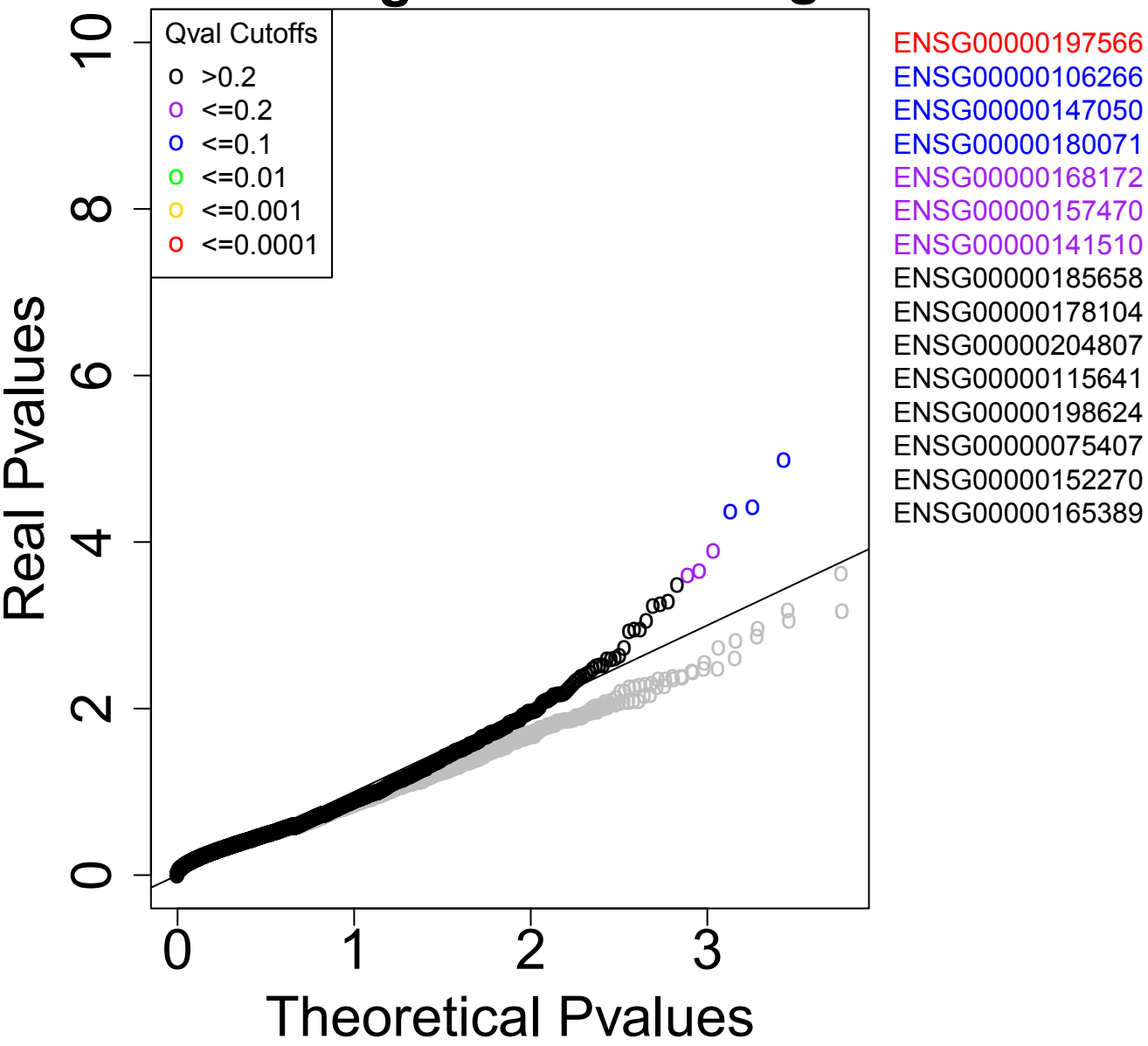
Supplementary Figure S2

Exon/Background Mutation/Kb

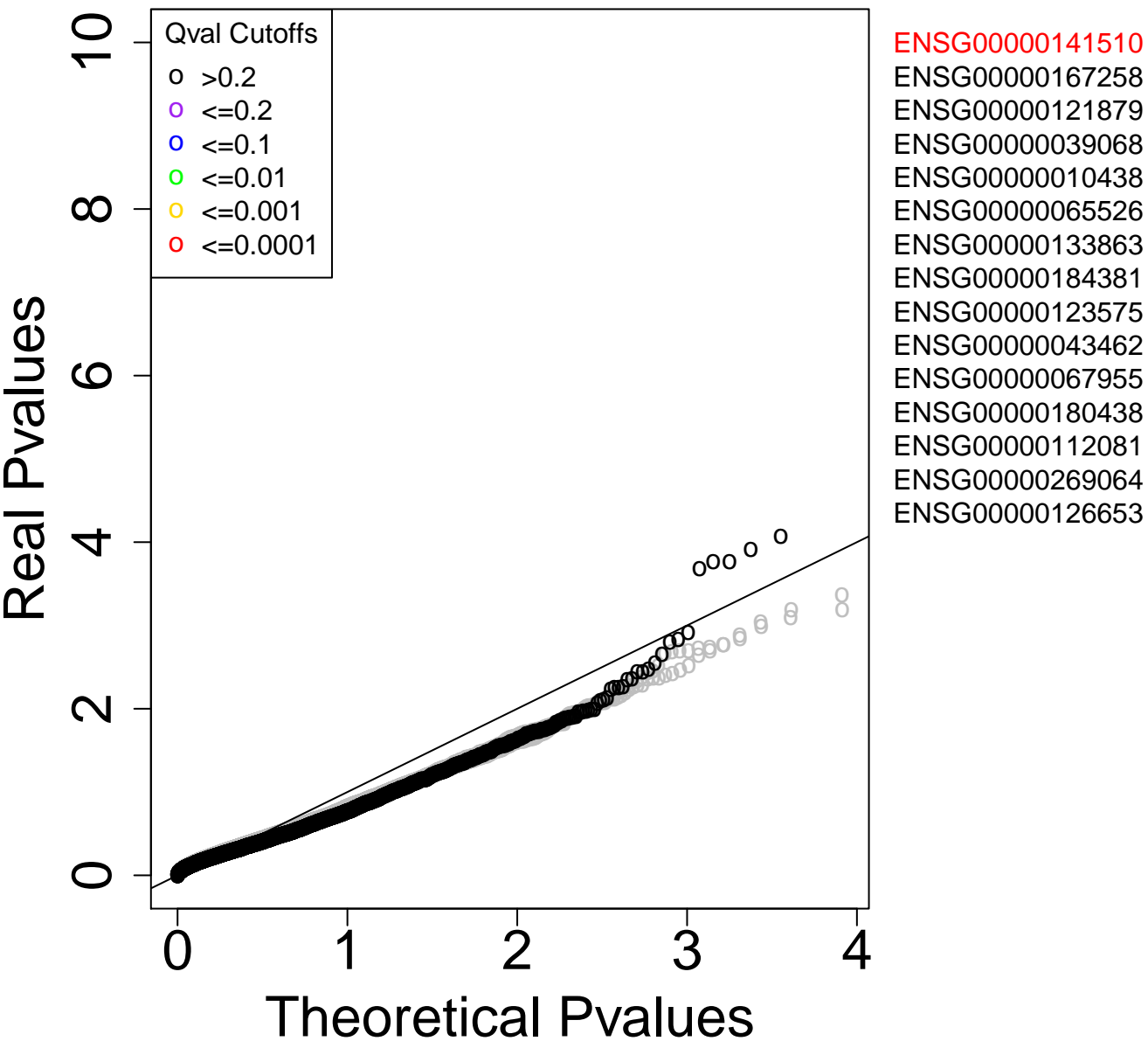


Supplementary Figure S3

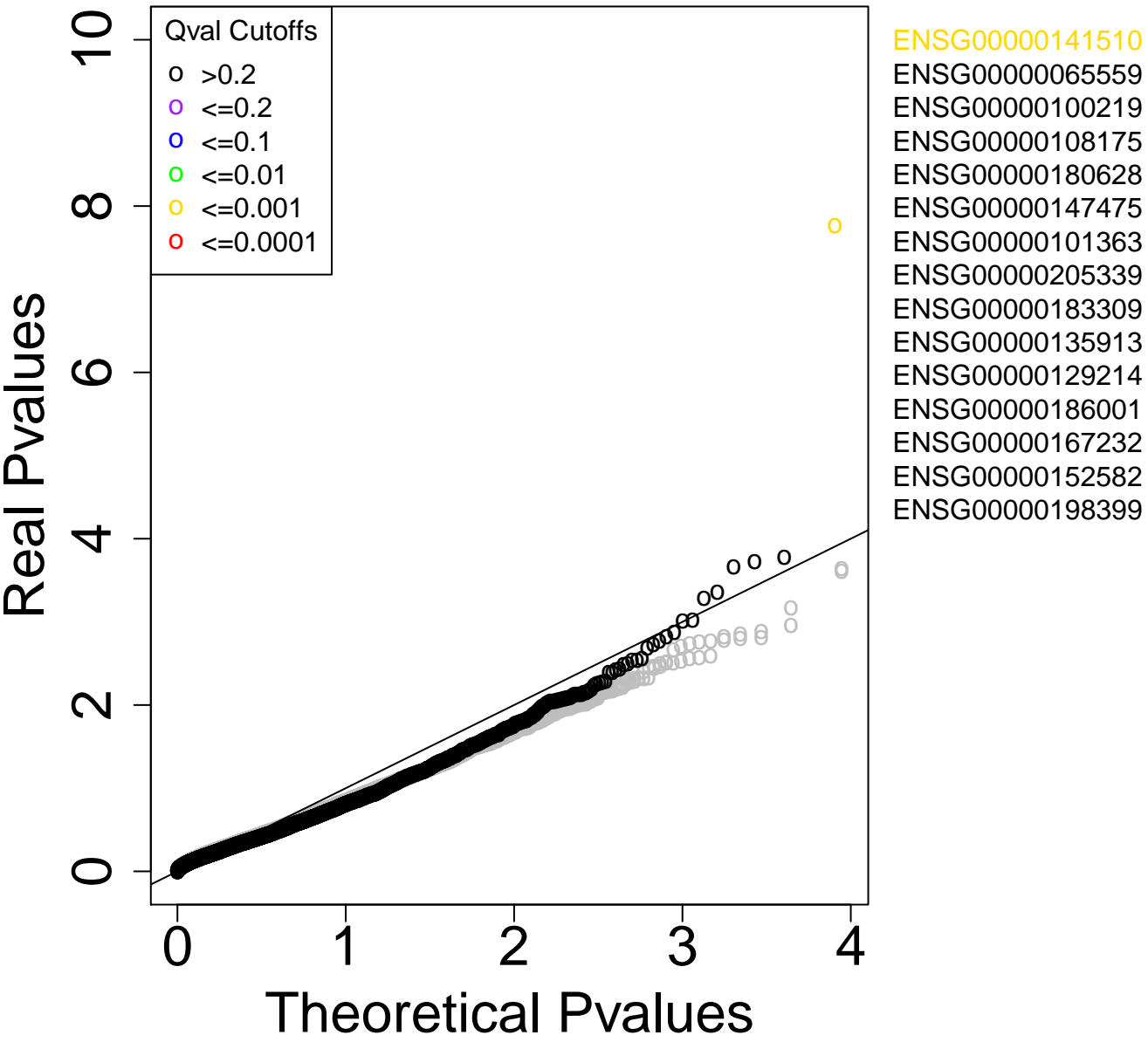
Coding BLCA – 5407 genes



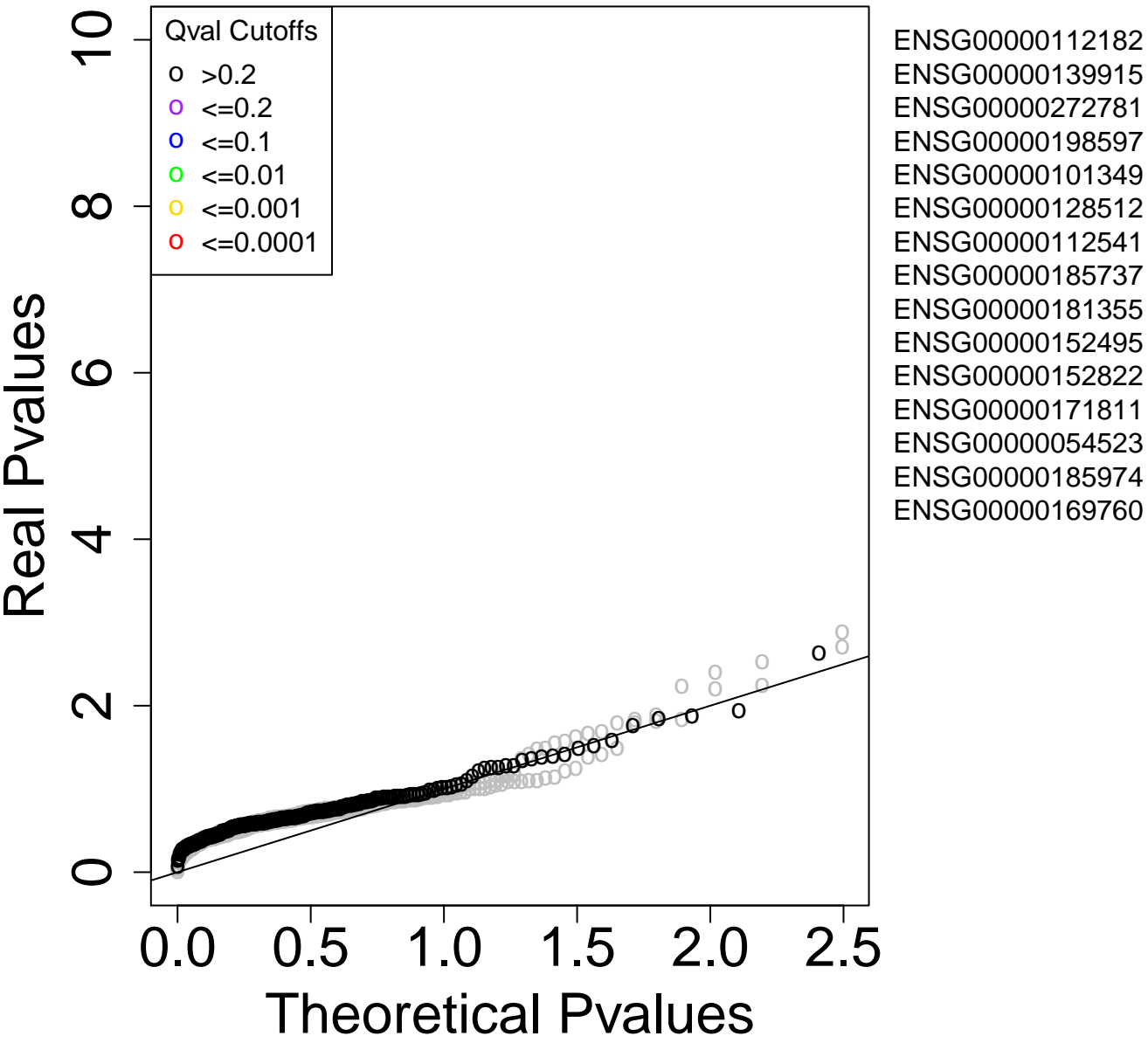
Coding BRCA – 7115 genes



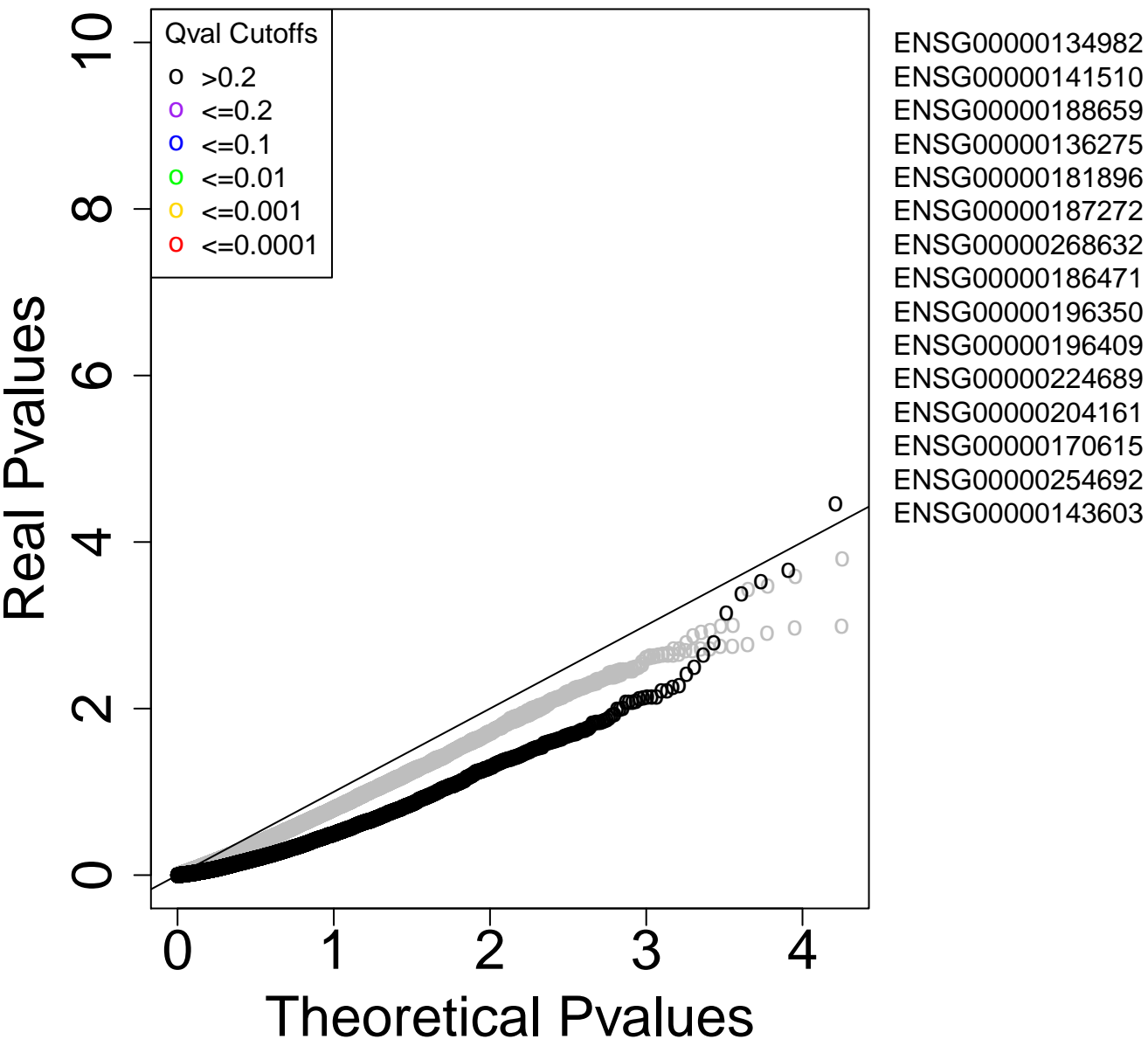
Coding Breast – 8045 genes



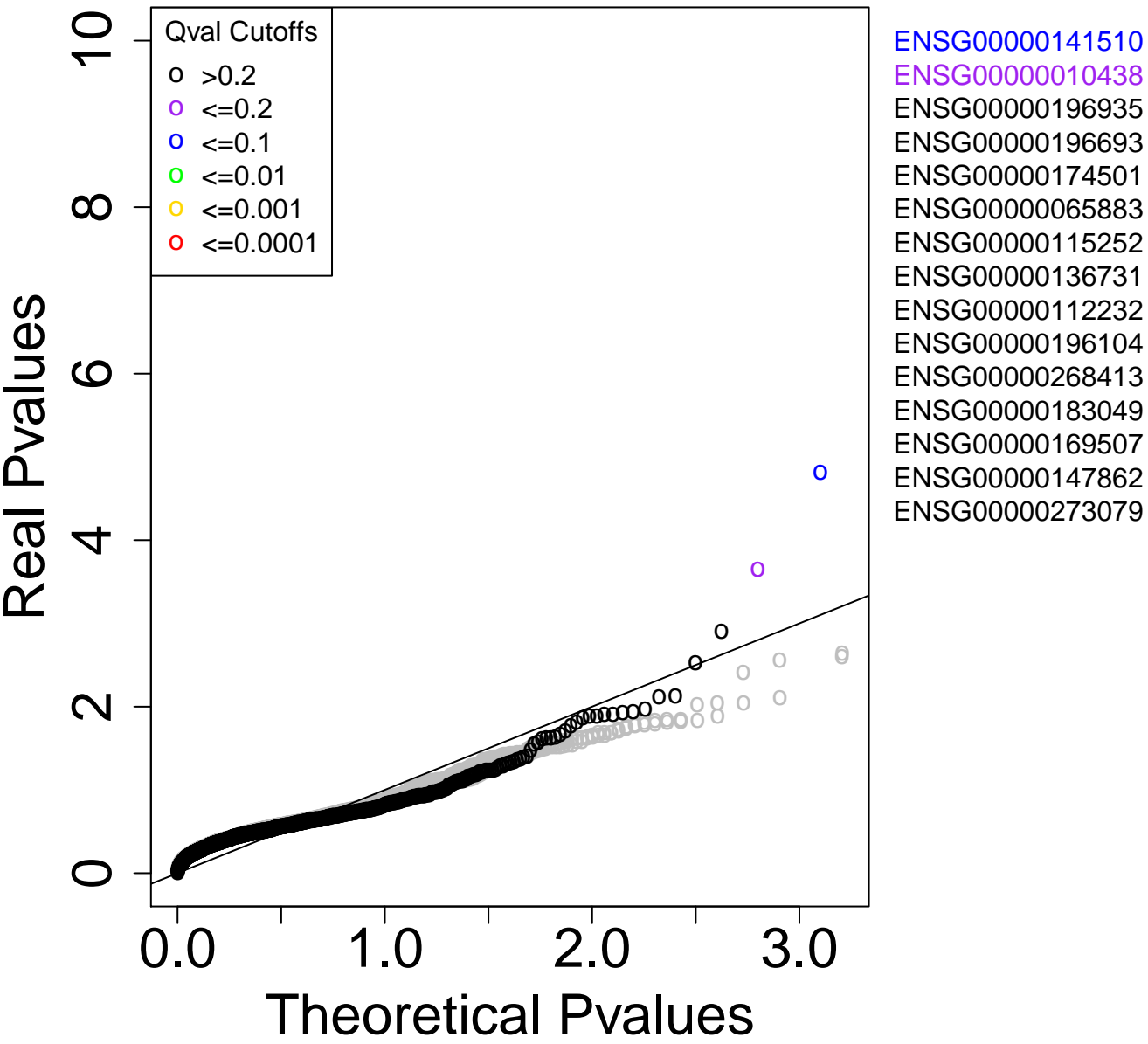
Coding CLL – 256 genes



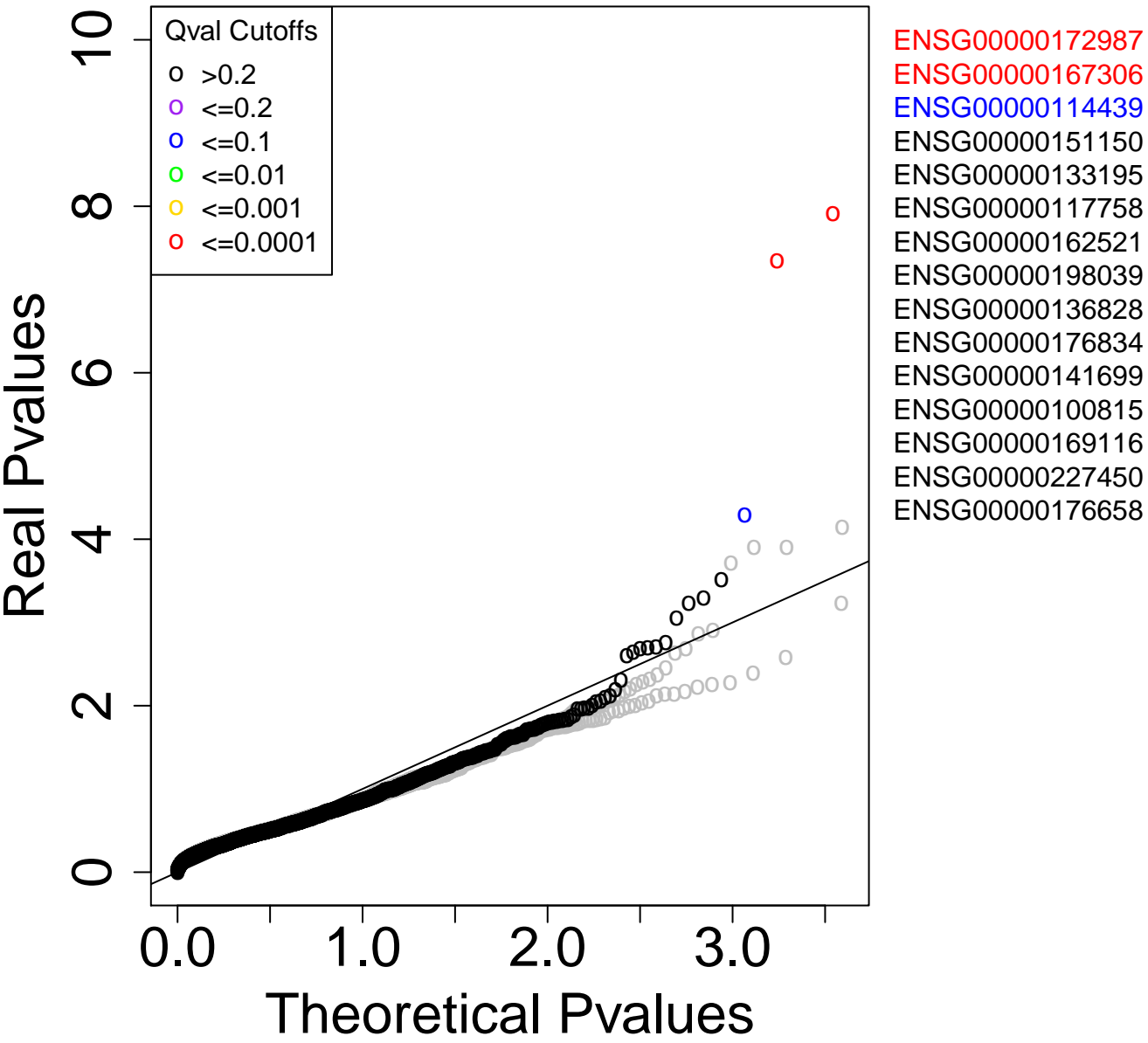
Coding CRC – 16265 genes



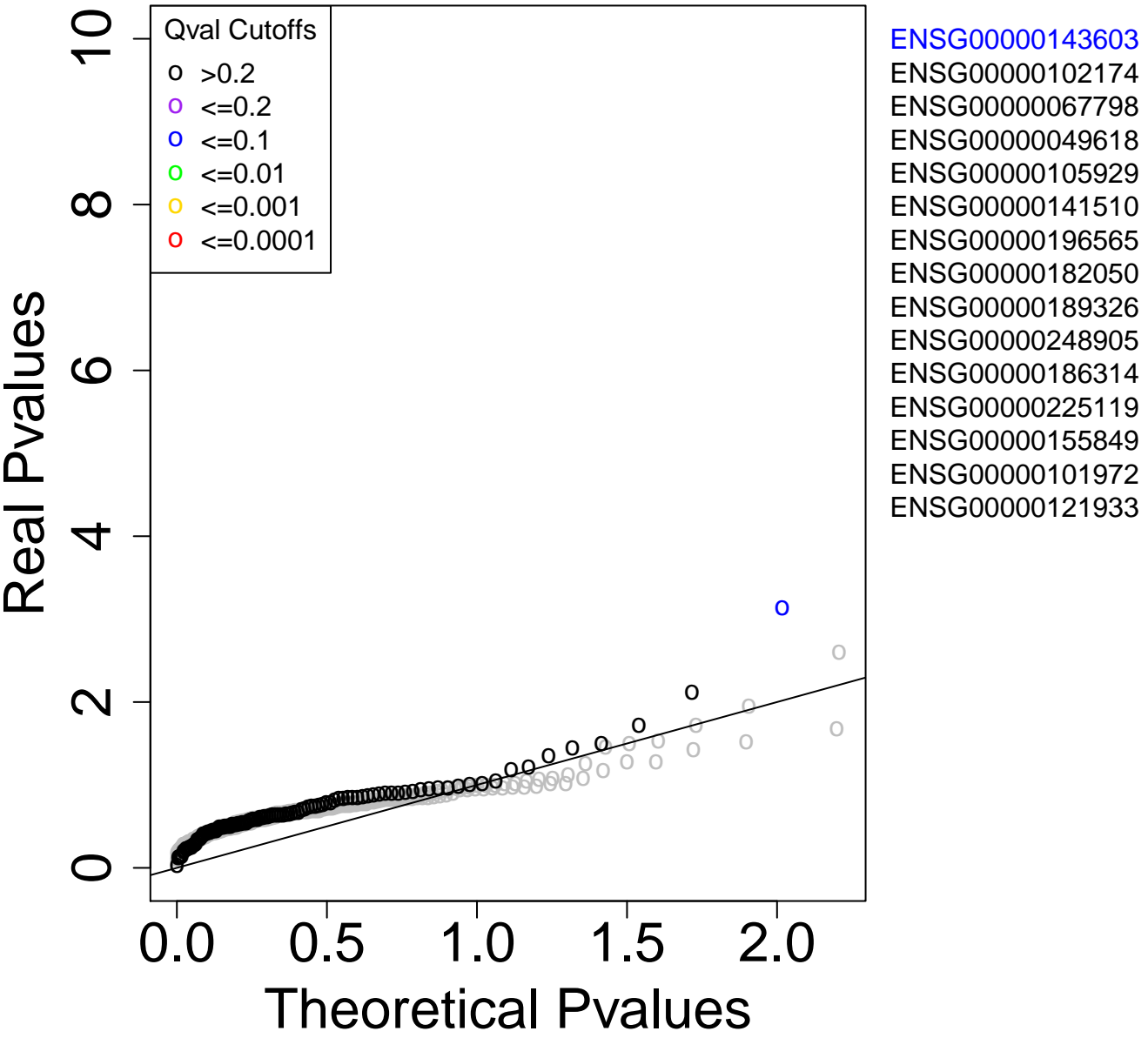
Coding GBM – 1262 genes



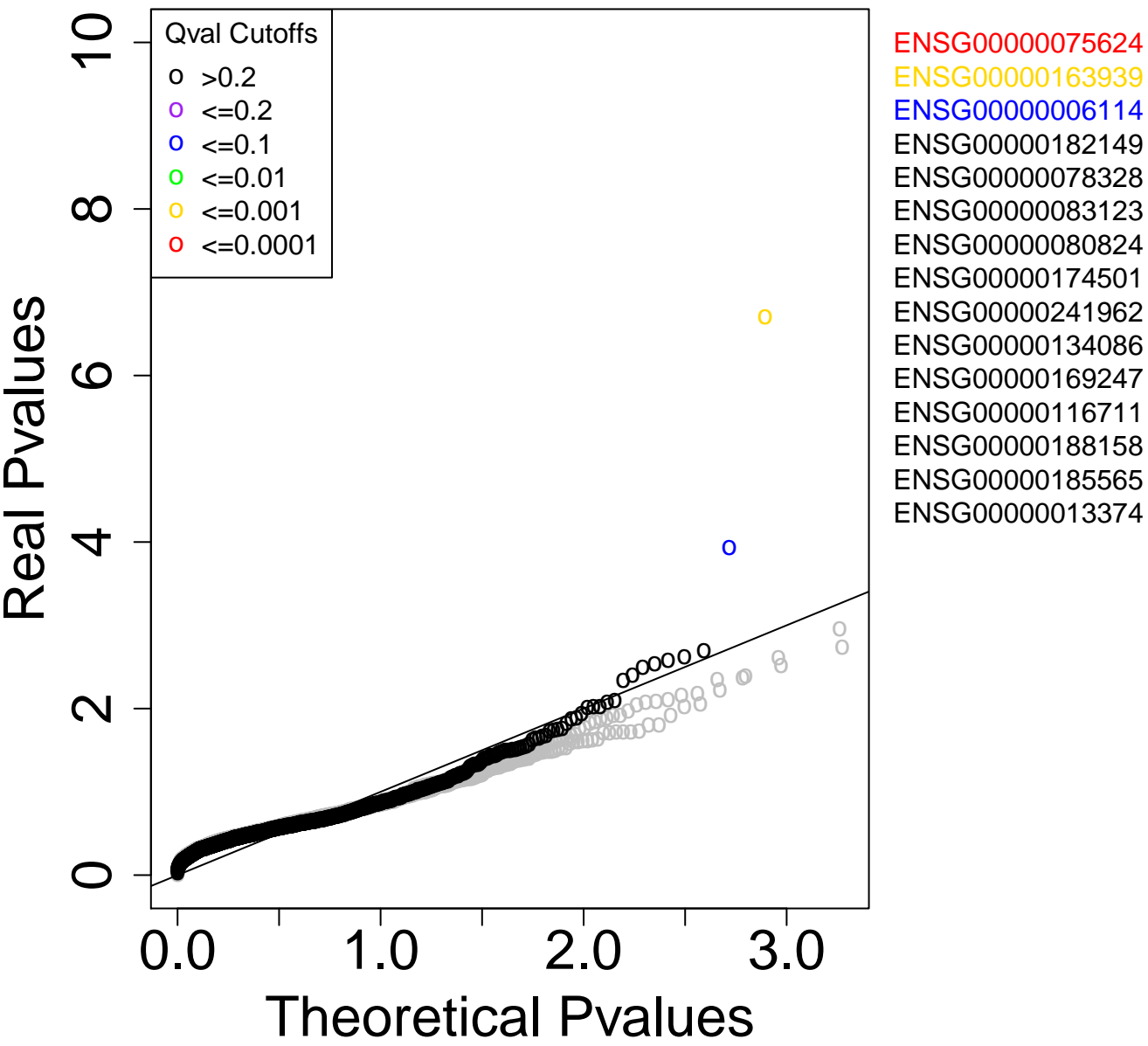
Coding HNSC – 3482 genes



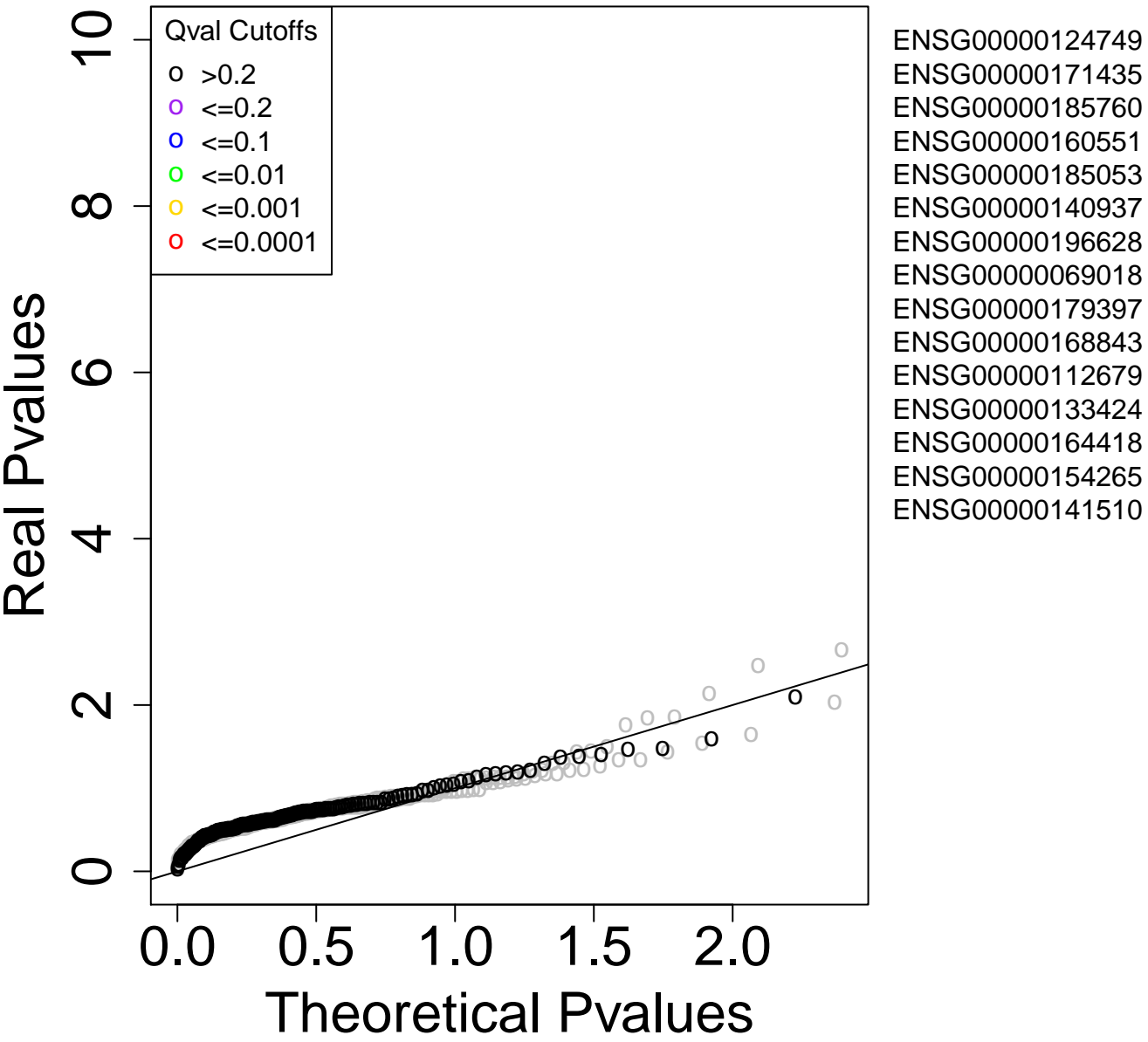
Coding KICH – 104 genes



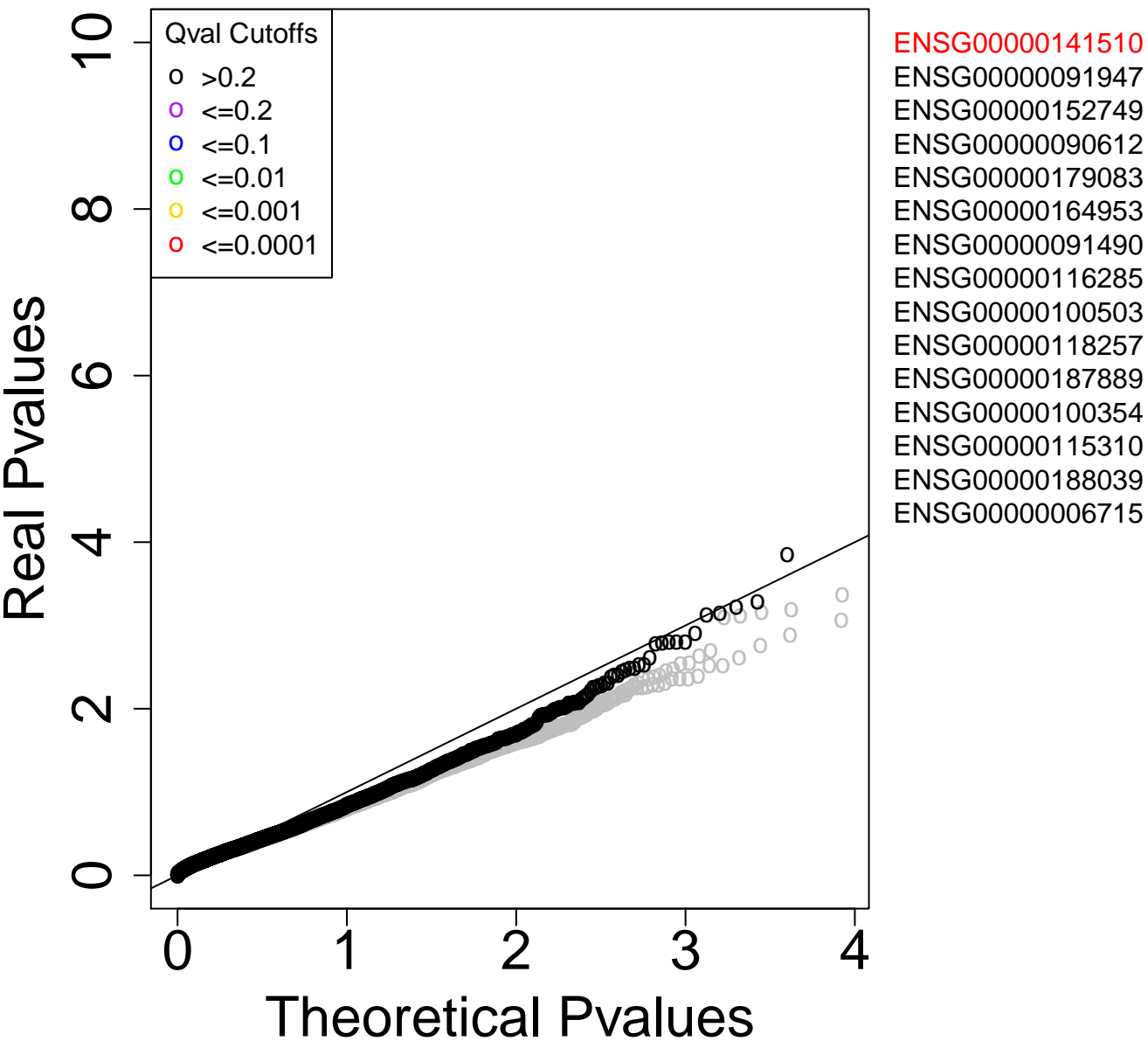
Coding KIRC – 1566 genes



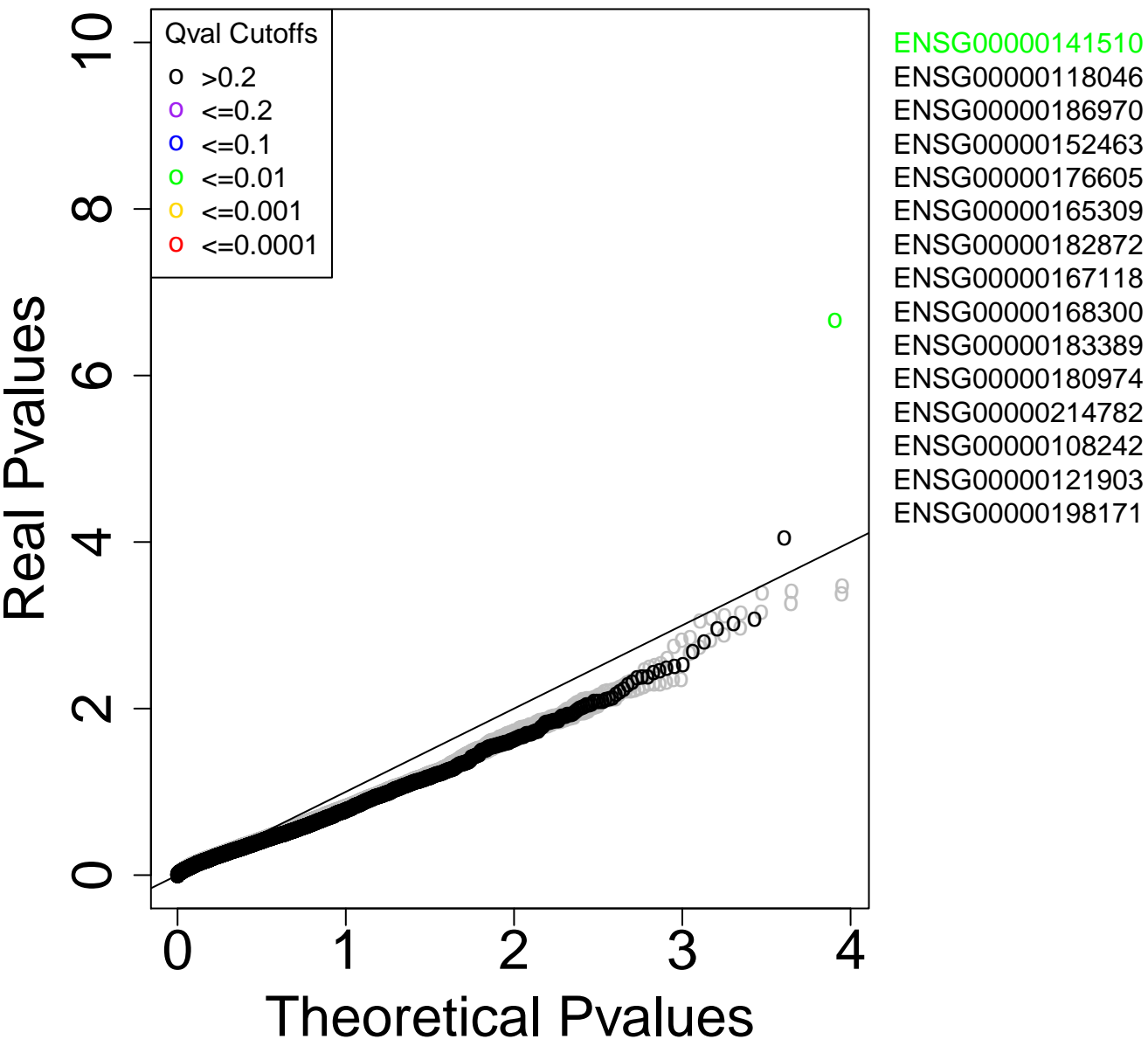
Coding LGG – 168 genes



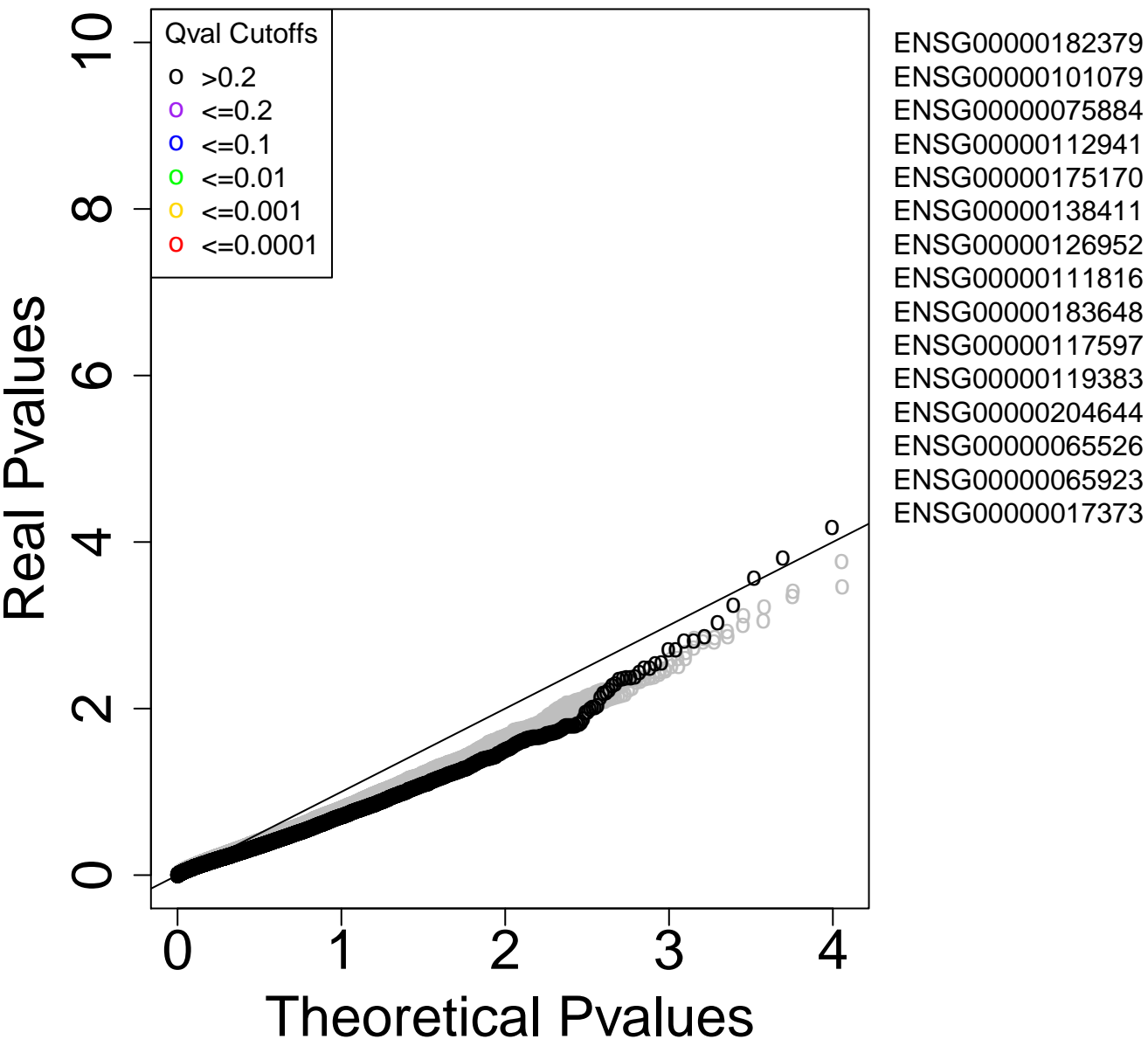
Coding Liver – 7982 genes



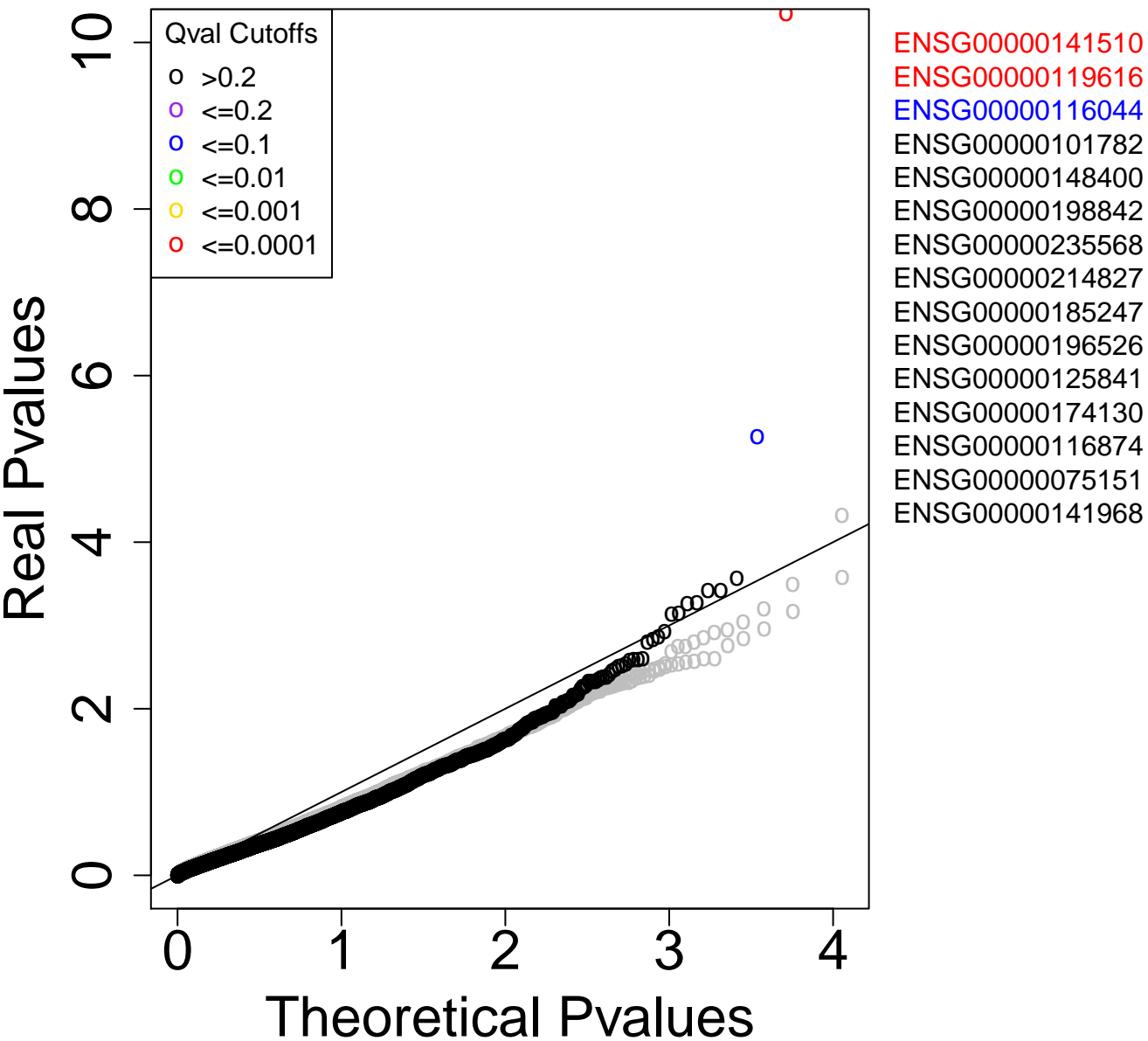
Coding LUAD – 8082 genes



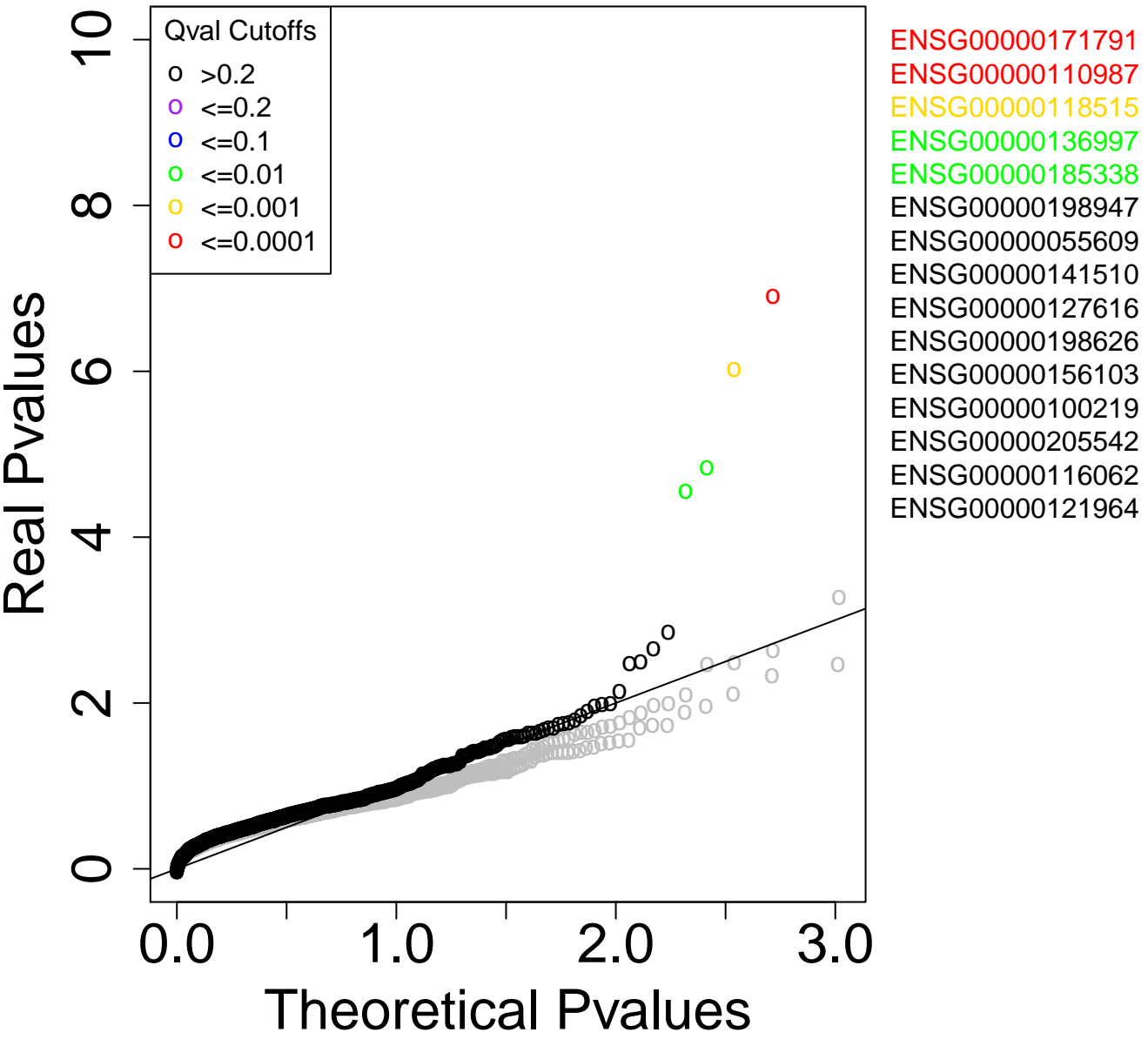
Coding Lung_adeno – 9893 genes



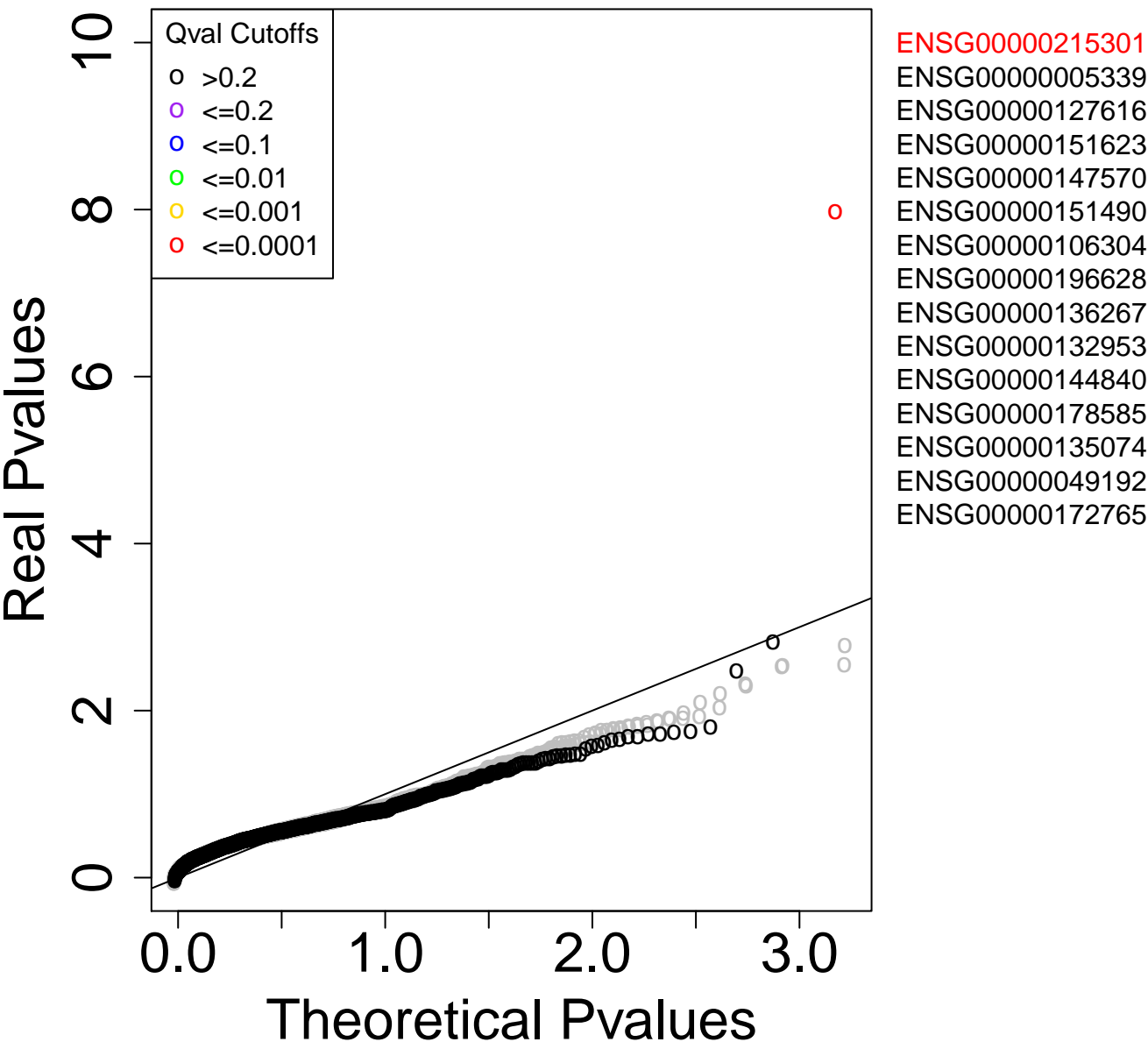
Coding LUSC – 10329 genes



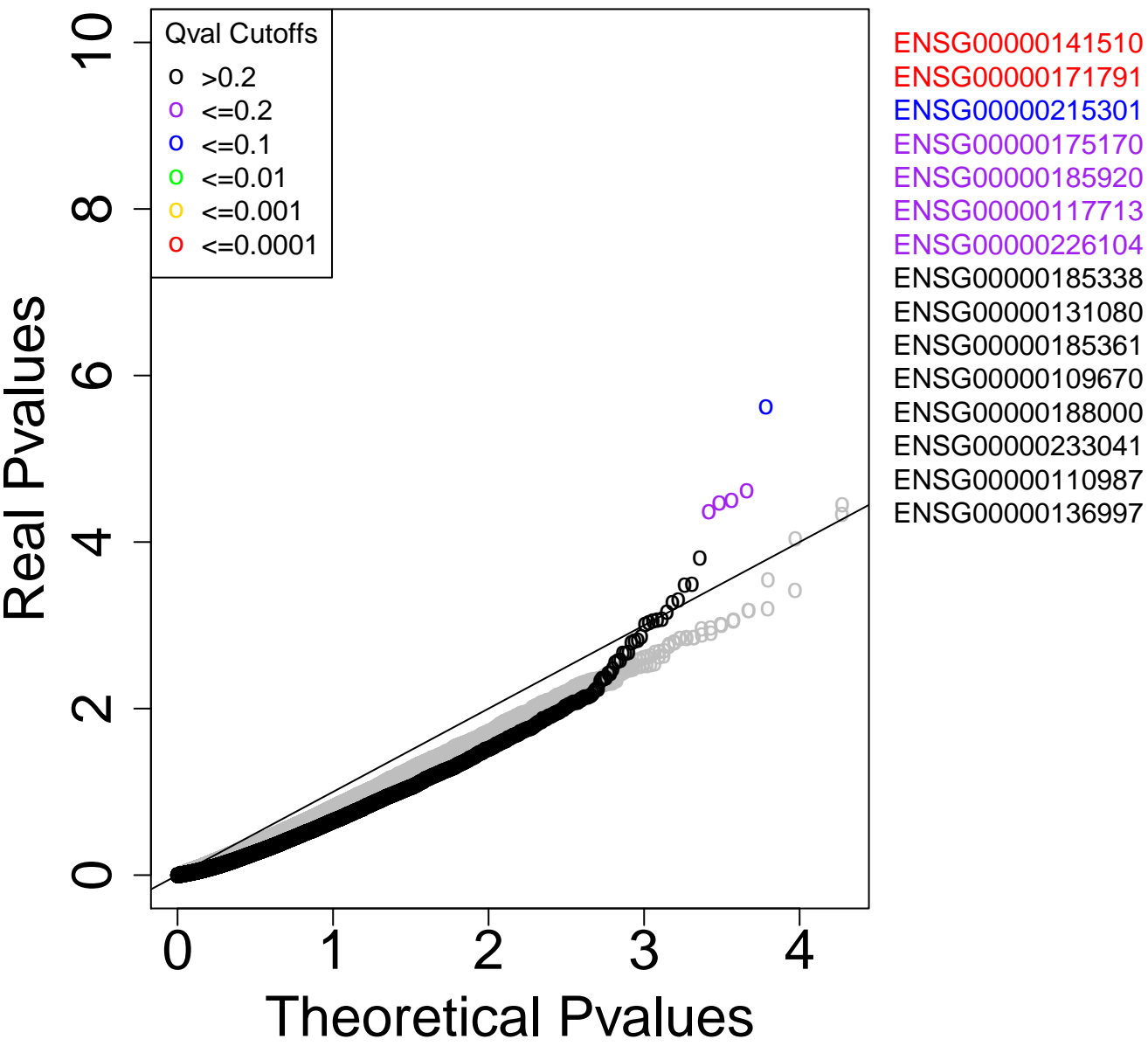
Coding Lymphoma_B-cell – 1038 genes



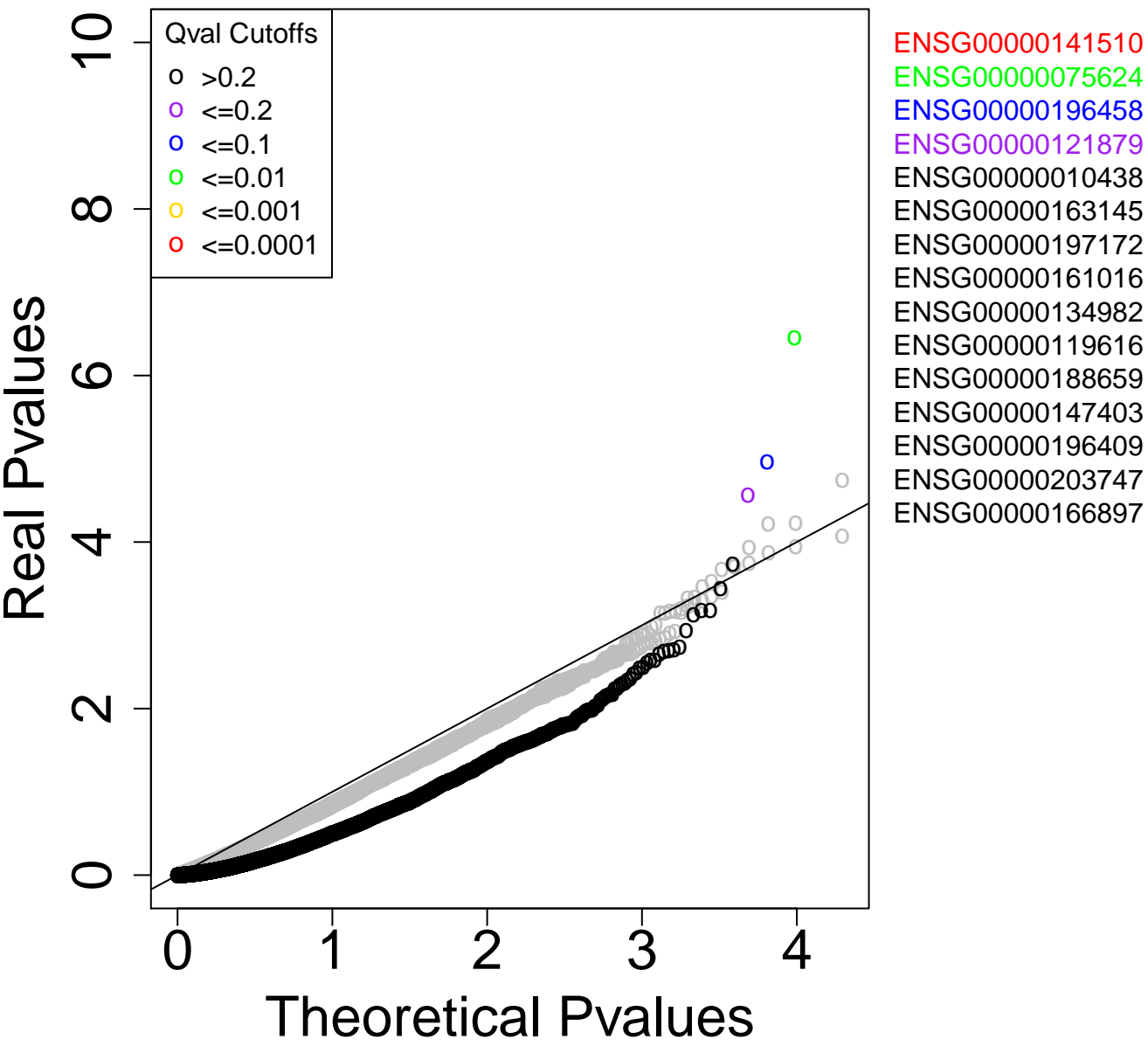
Coding Medulloblastoma – 1488 genes



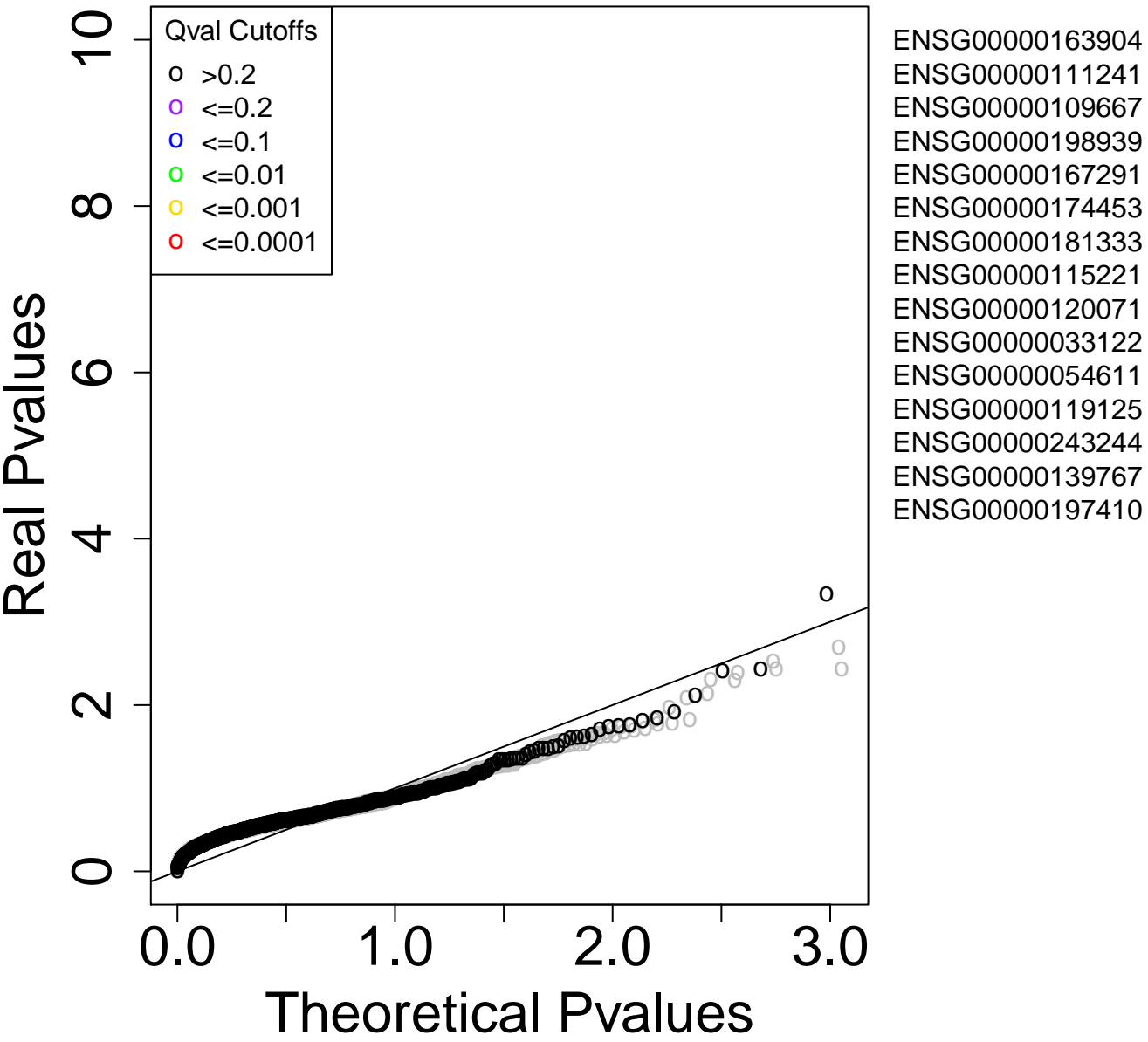
Ranking Pancancer_Alexandrov – 18261 genes



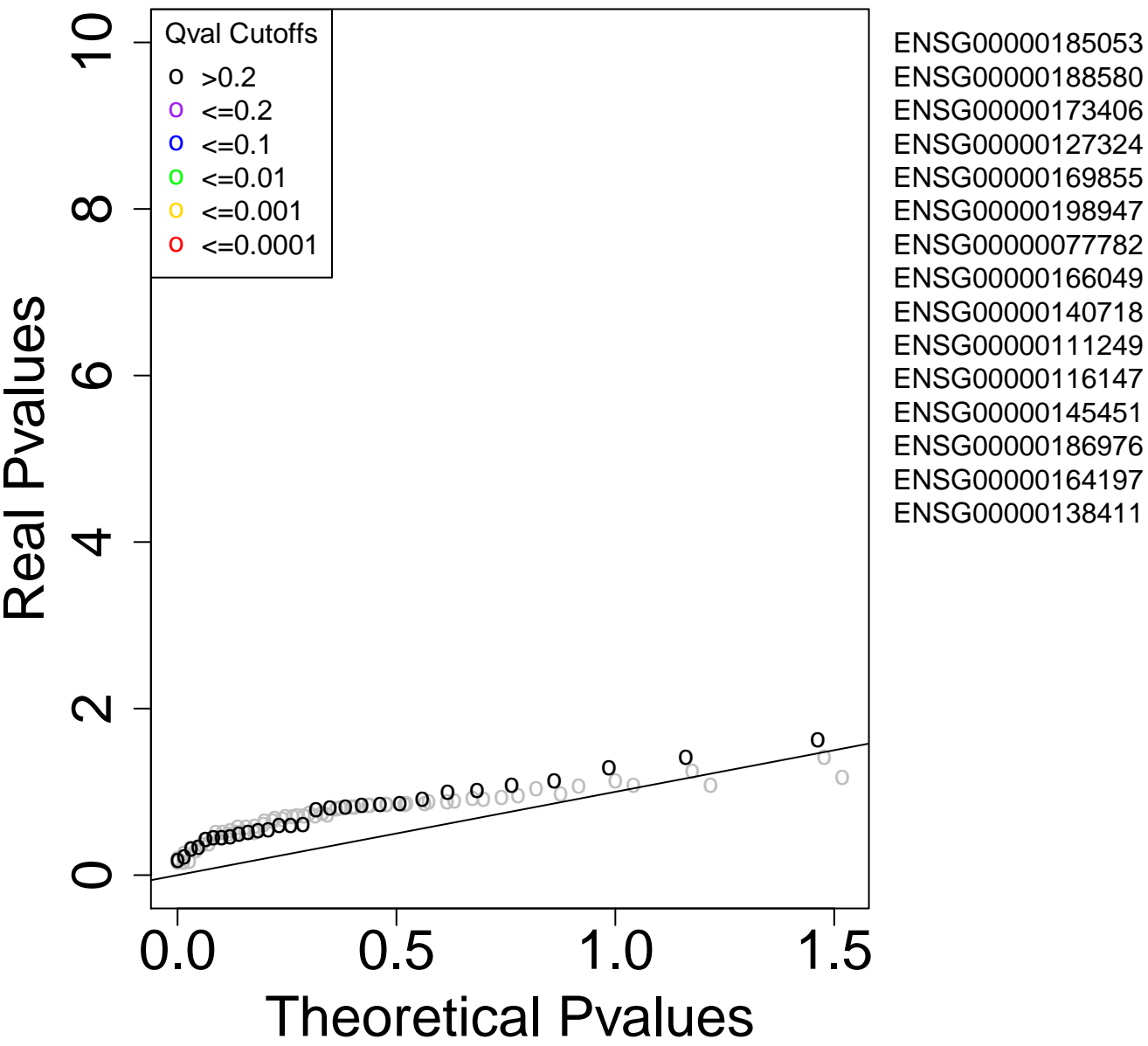
Coding Pancancer_TCGA – 19313 genes



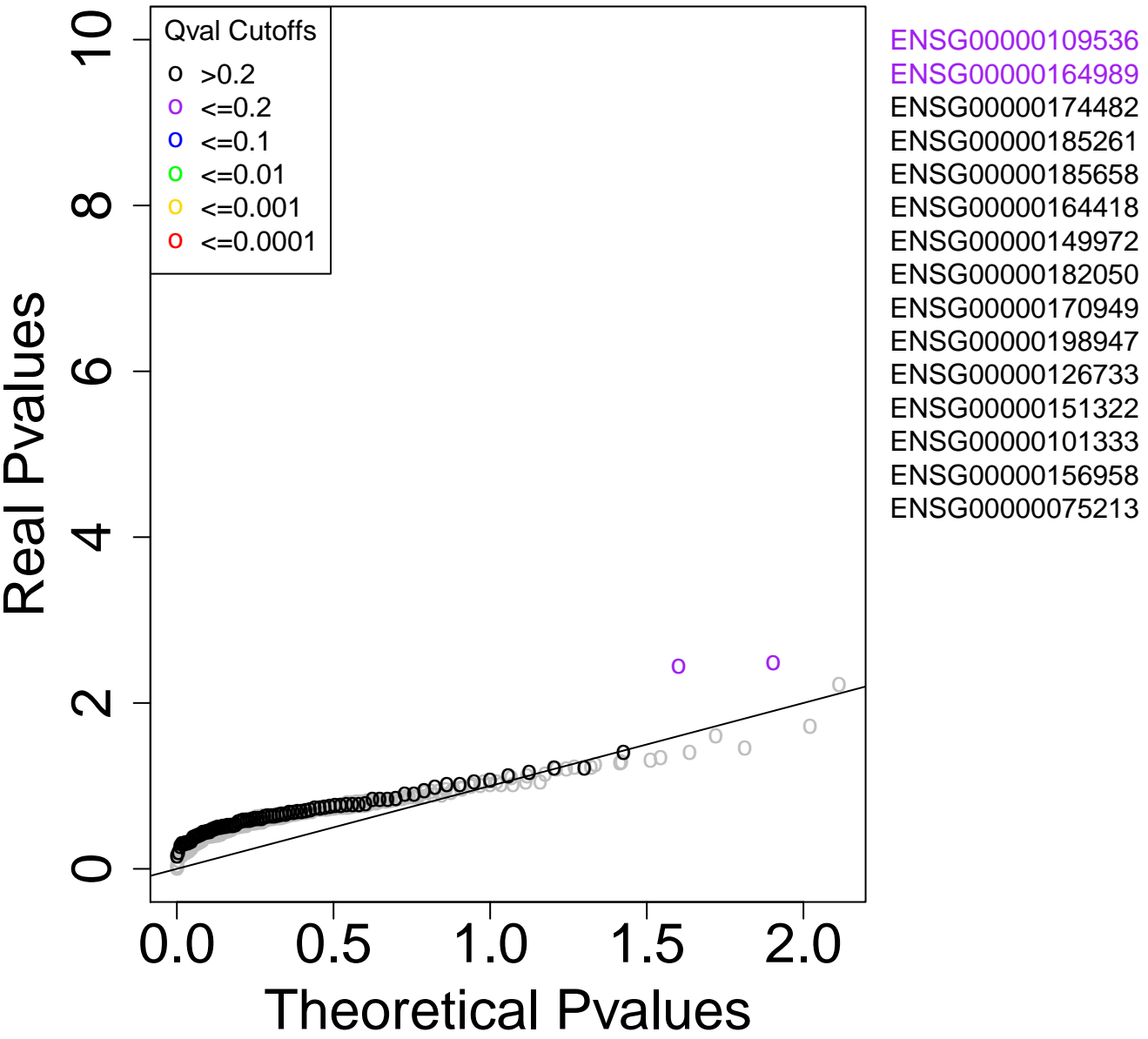
Coding Pancreas – 960 genes



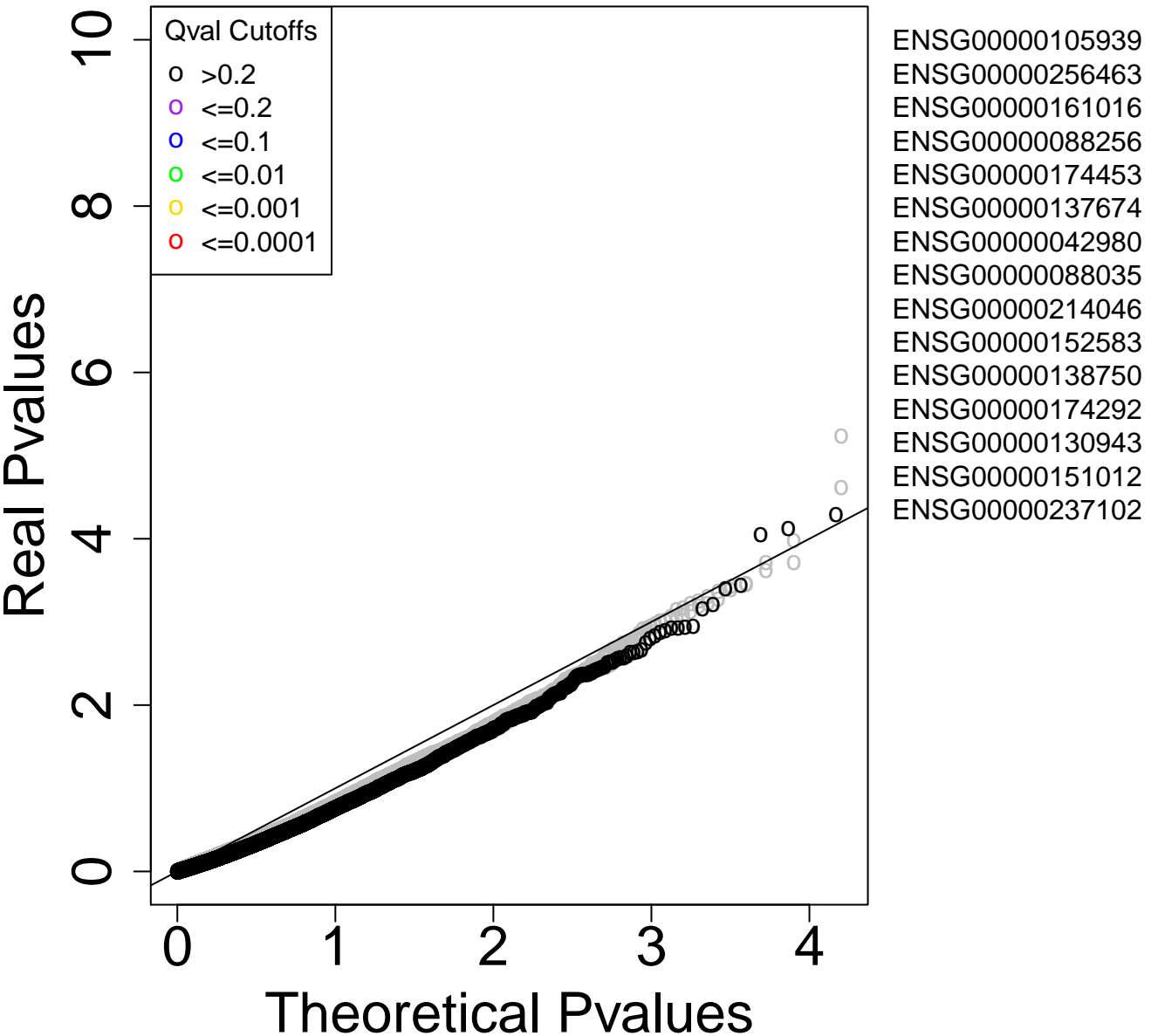
oding Pilocytic_astrocytoma – 29 genes



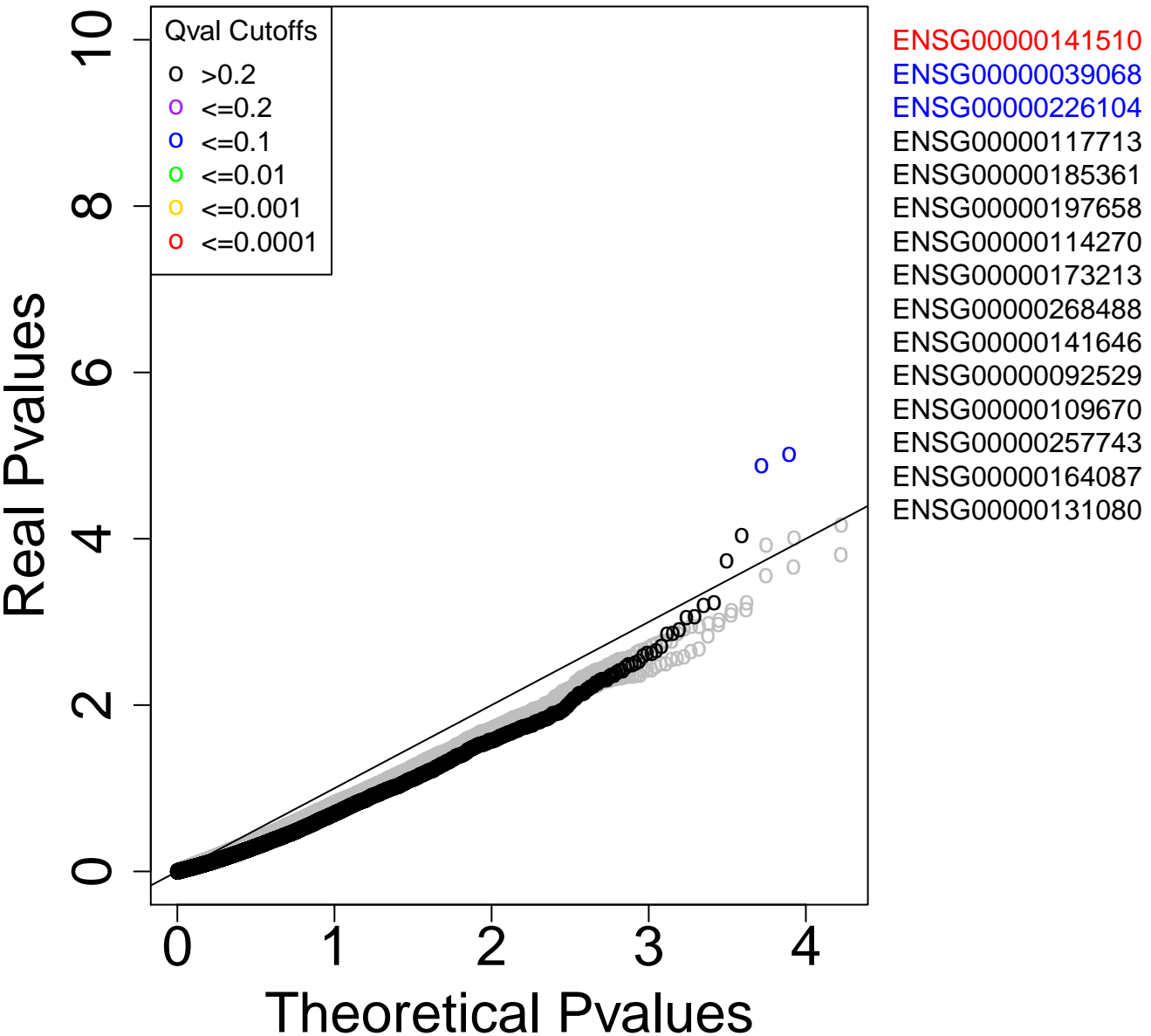
Coding PRAD – 80 genes



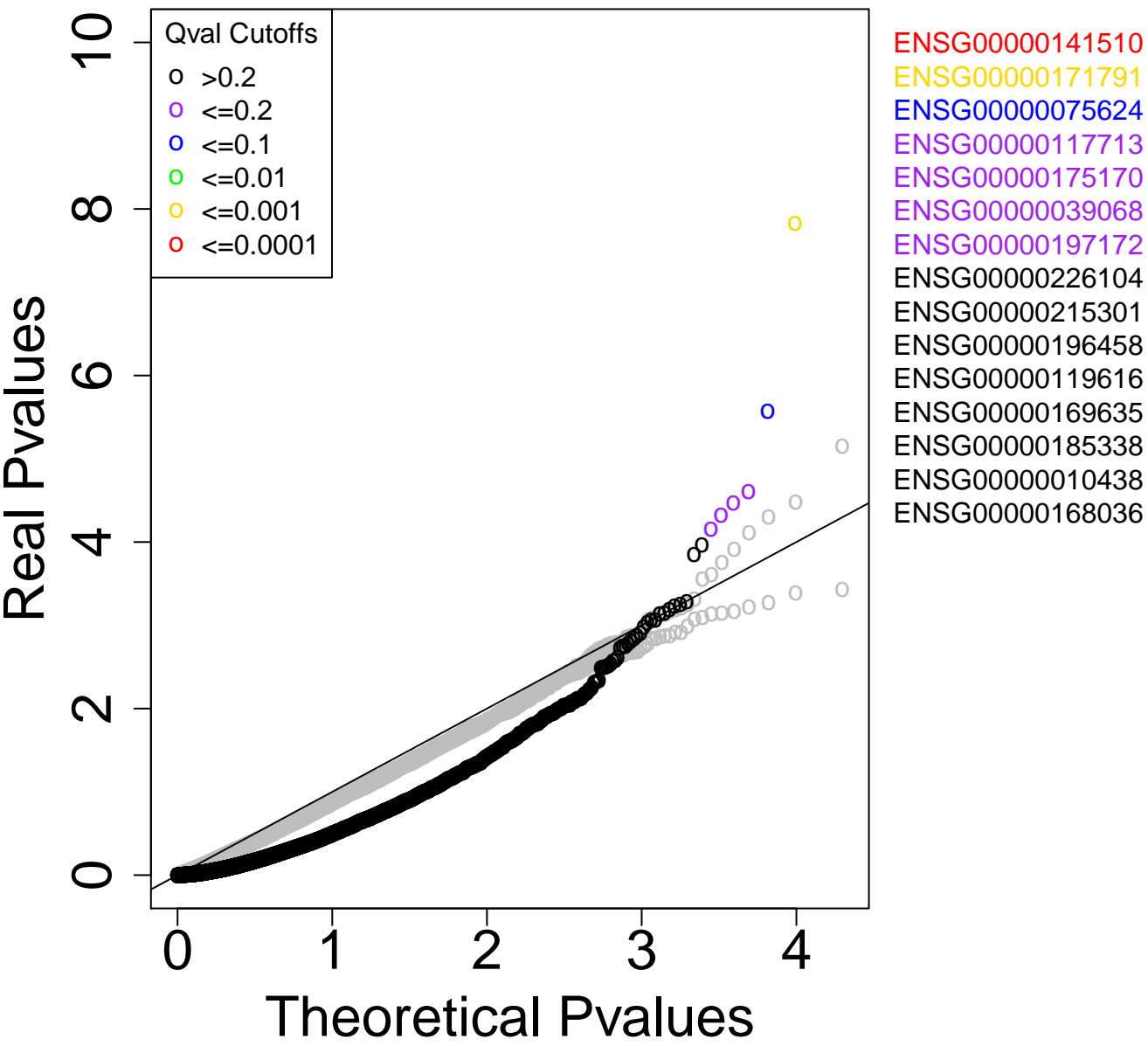
Coding SKCM – 14721 genes



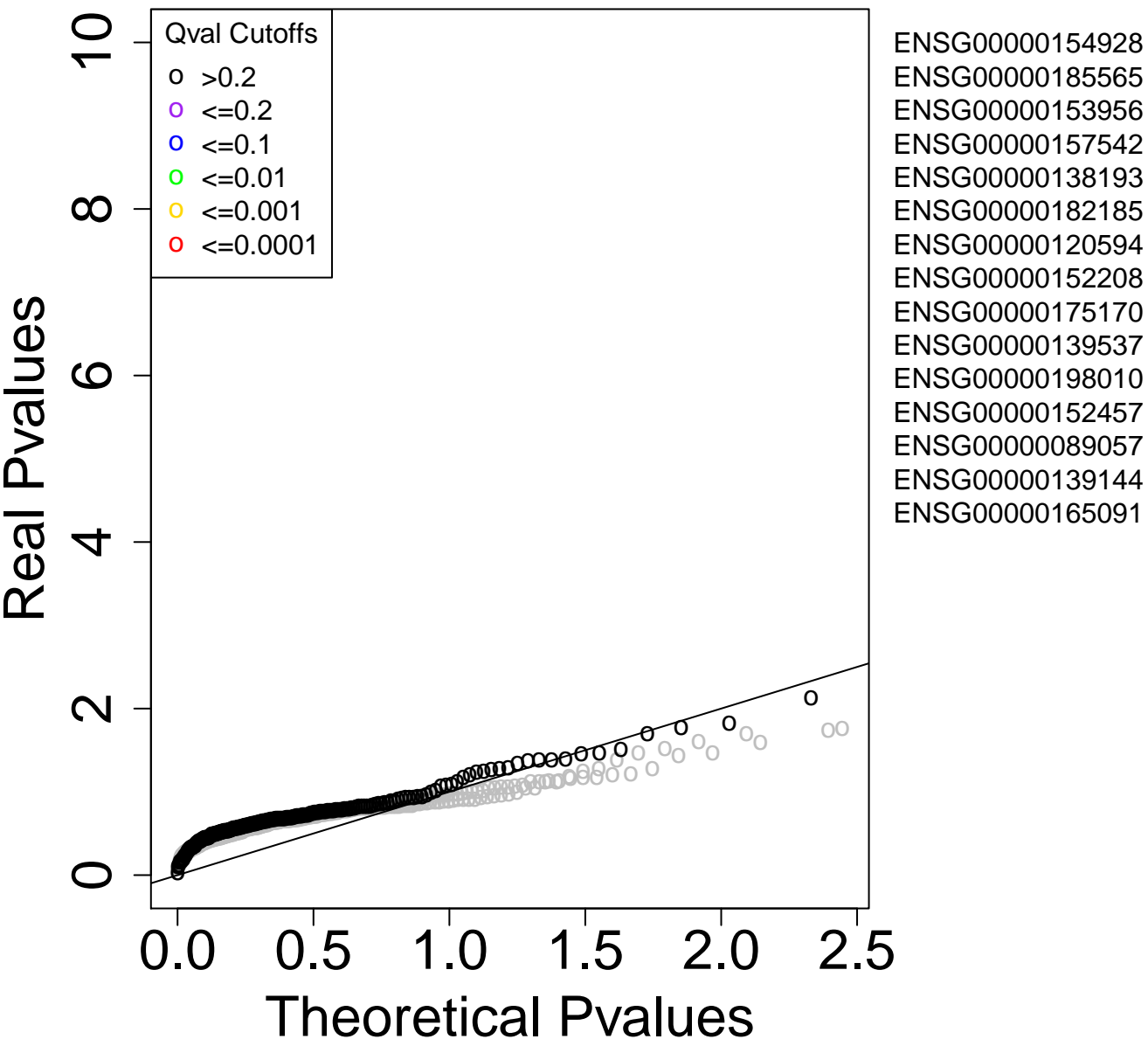
Coding Stad – 15662 genes



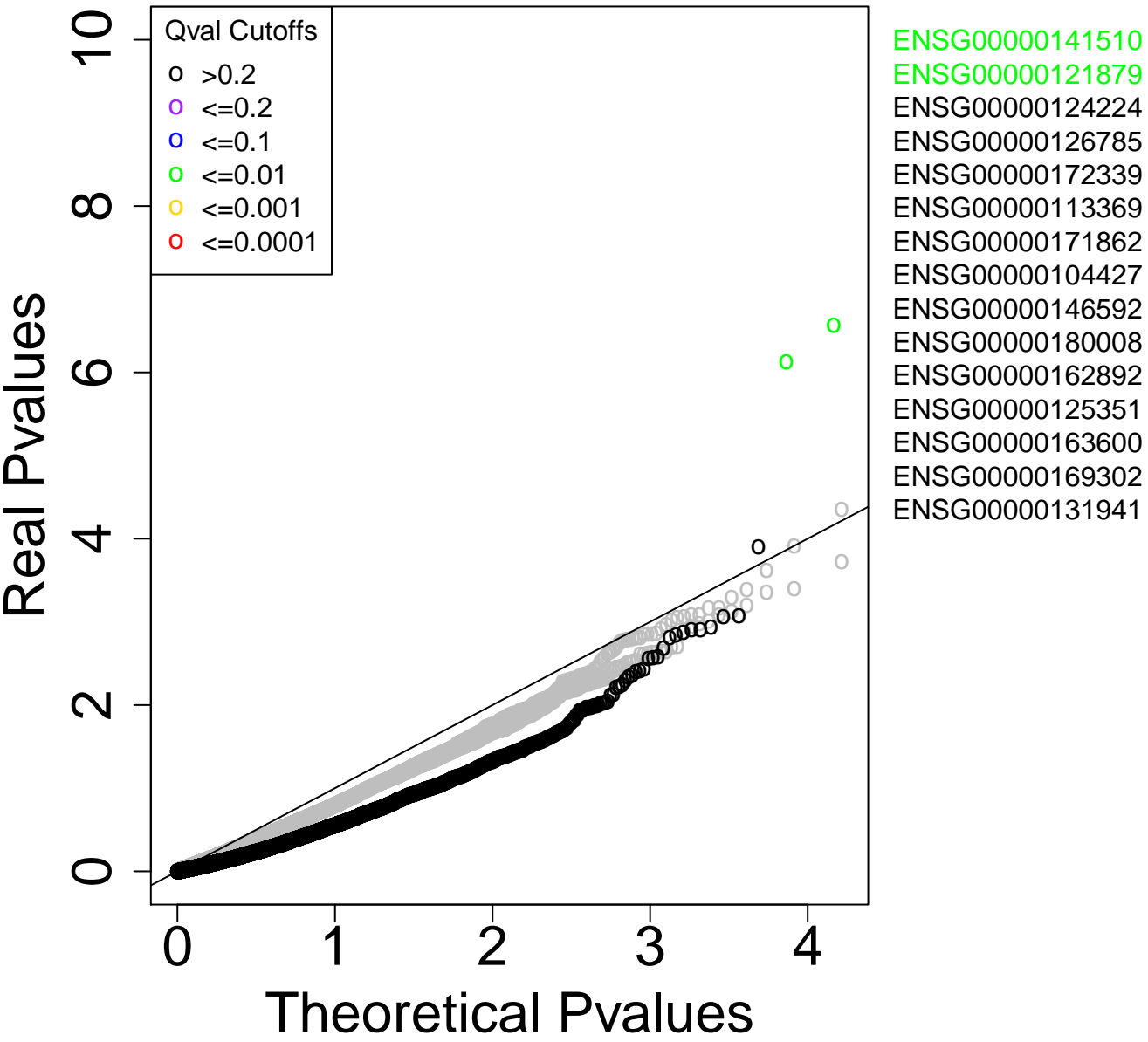
Coding Superpancancer – 19594 genes



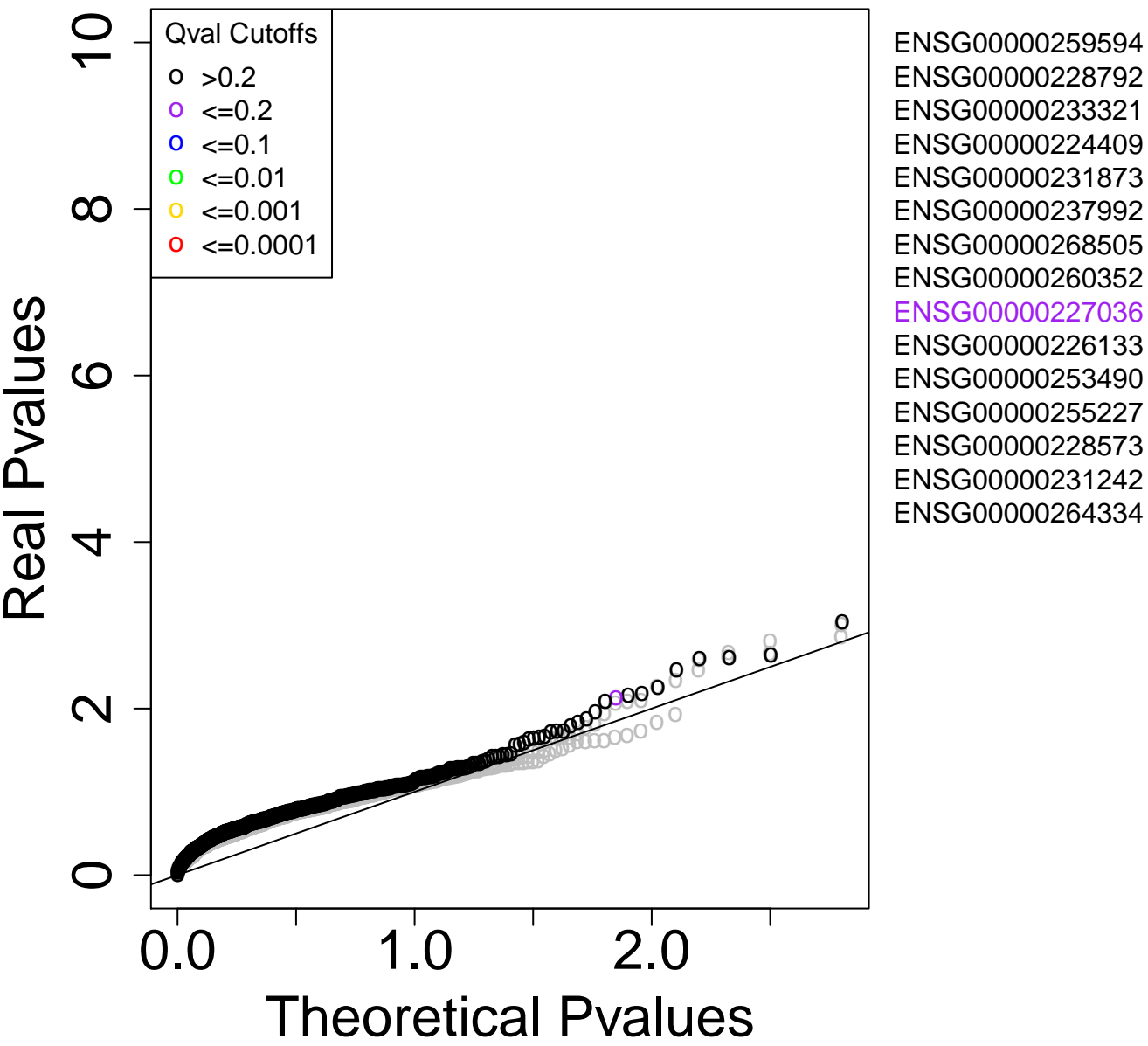
Coding THCA – 214 genes



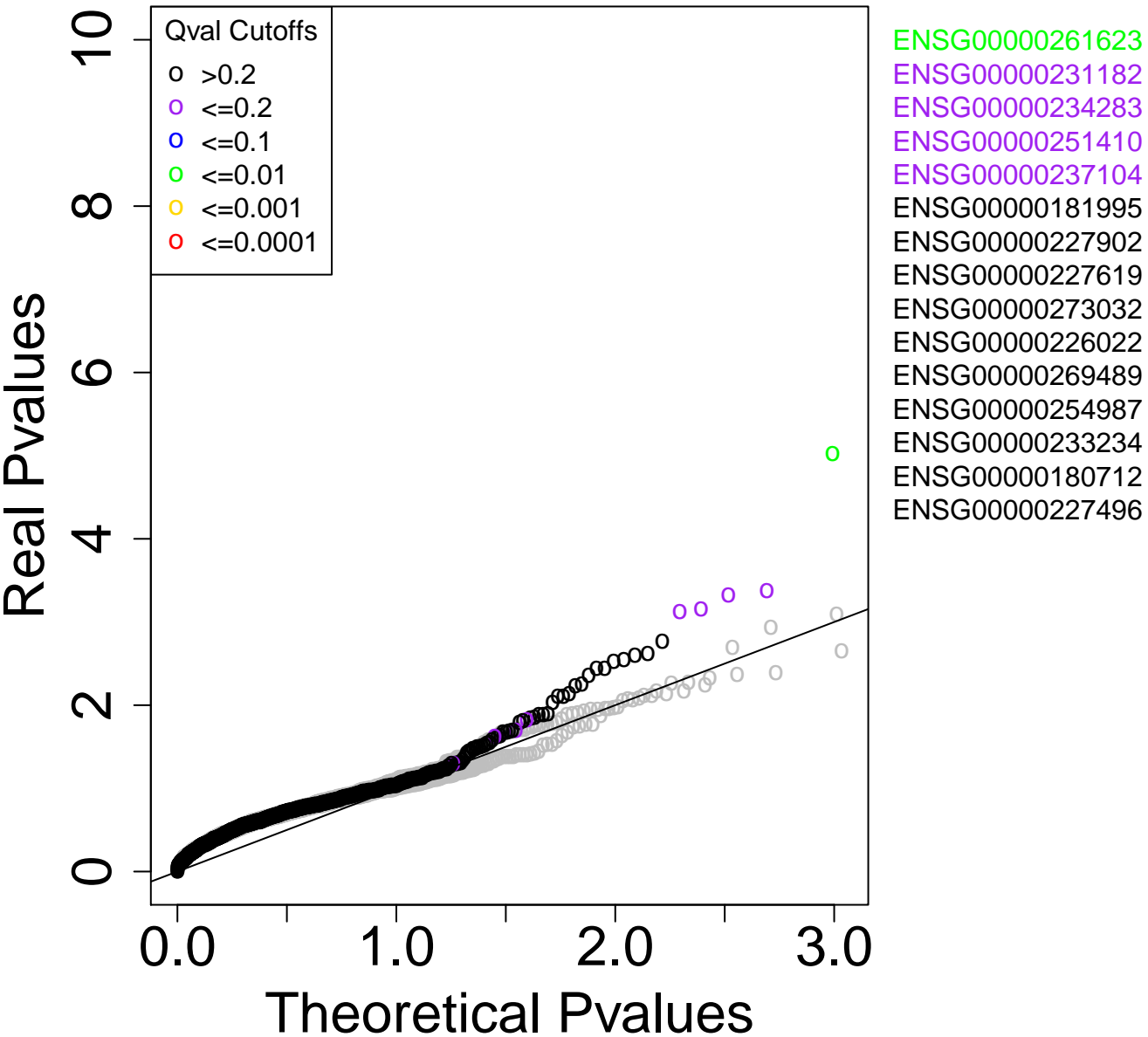
Coding UCEC – 14611 genes



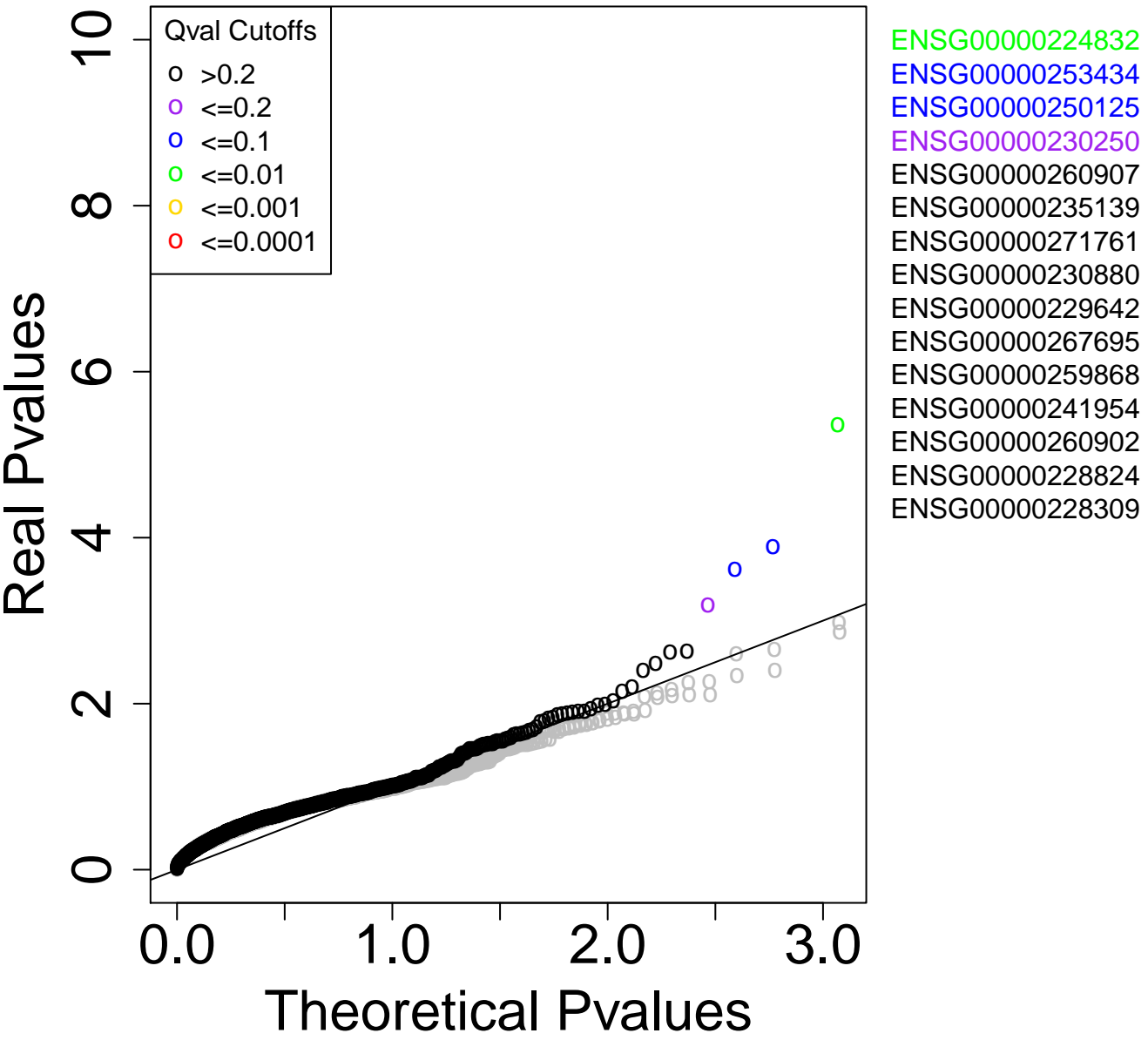
LncRNA BLCA – 637 genes



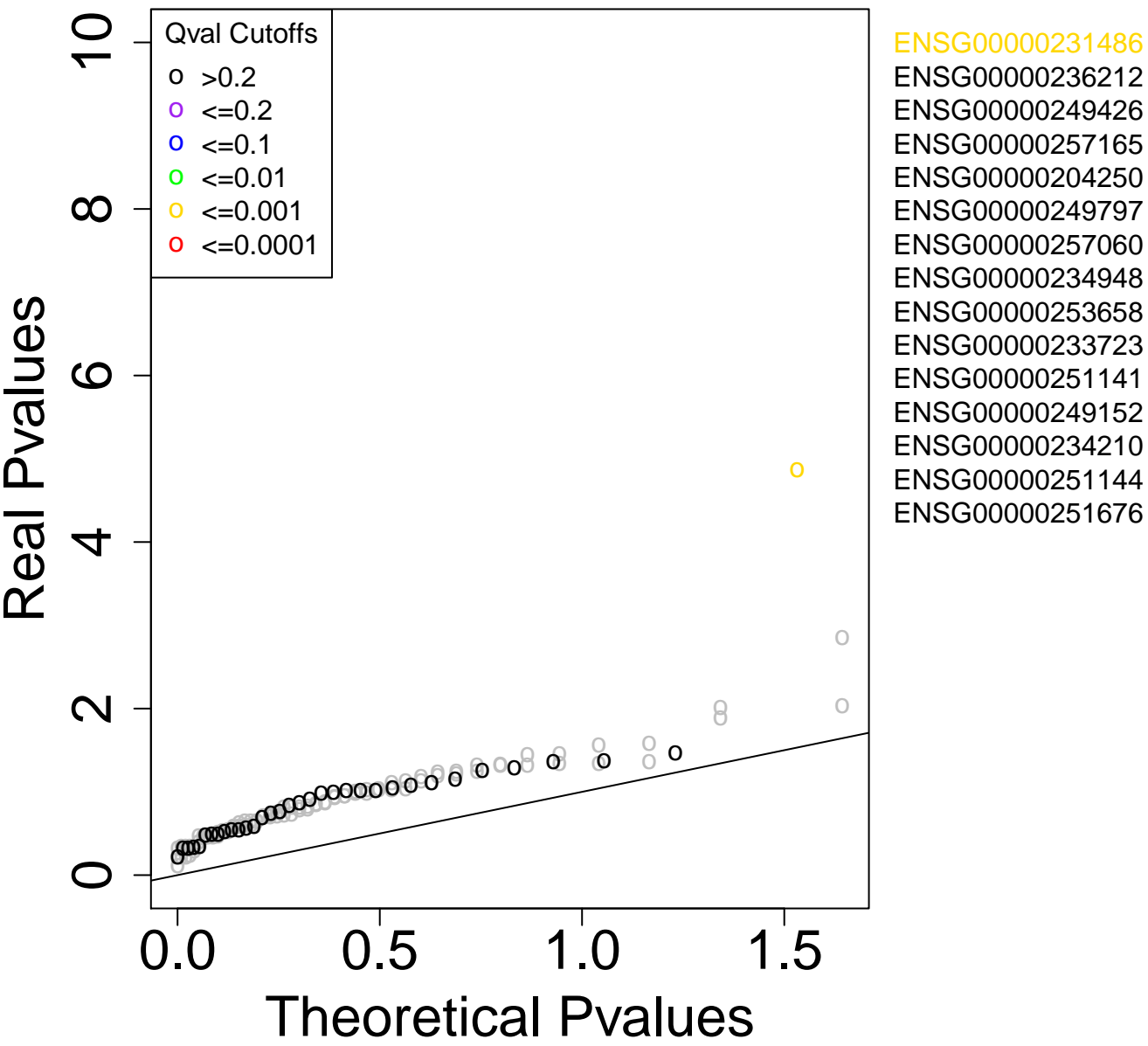
LncRNA BRCA – 985 genes



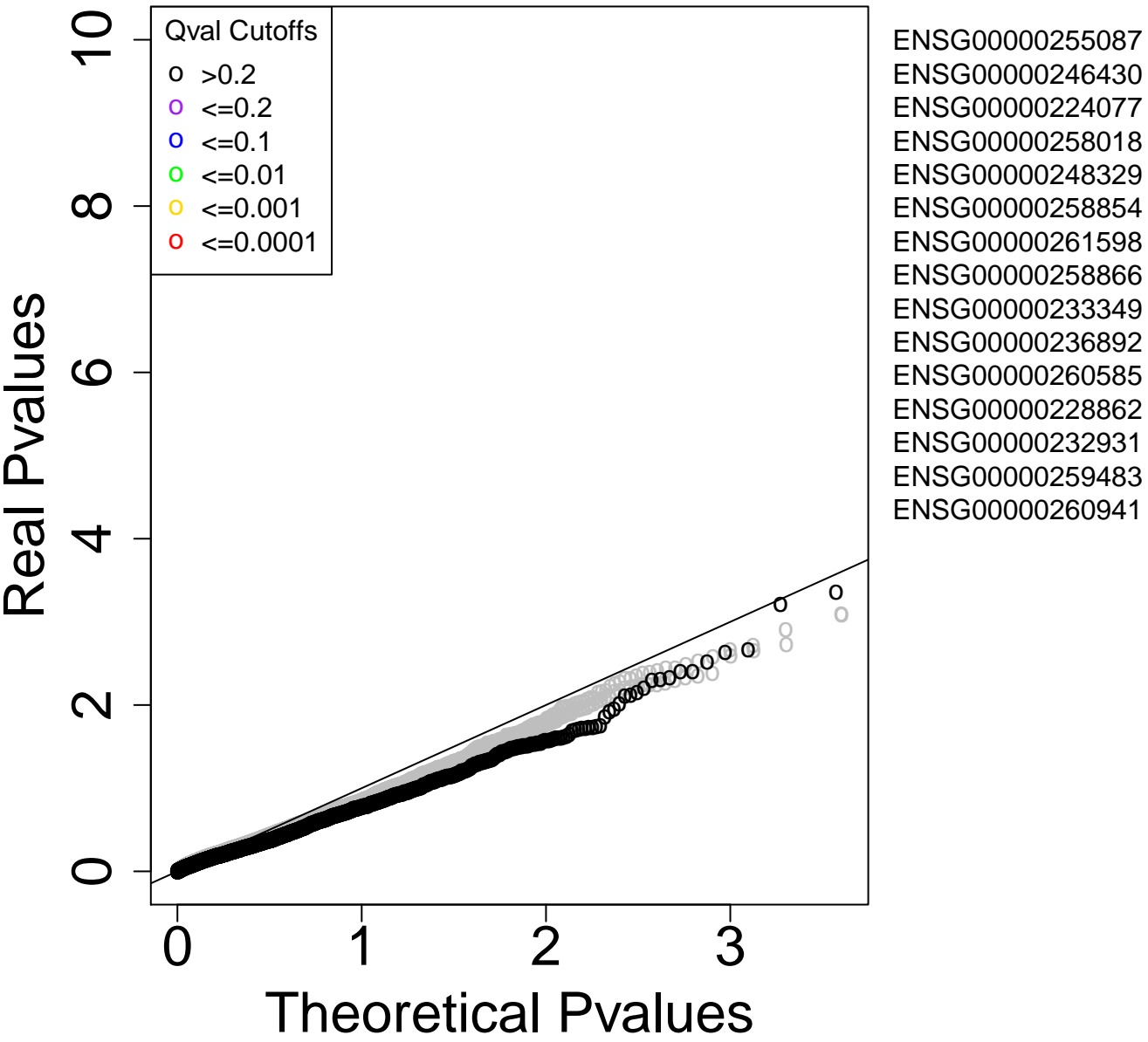
LncRNA Breast – 1169 genes



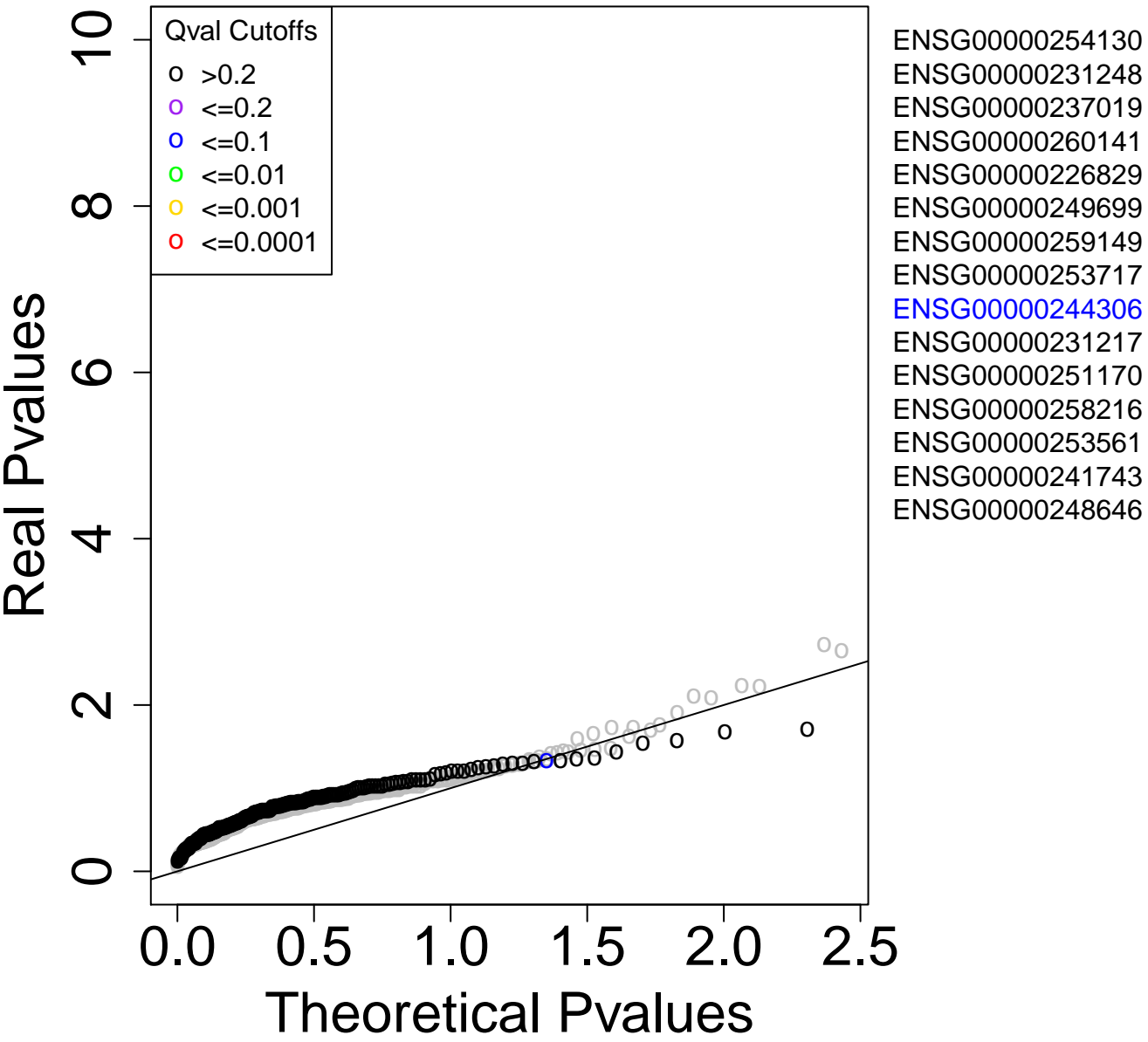
LncRNA CLL – 34 genes



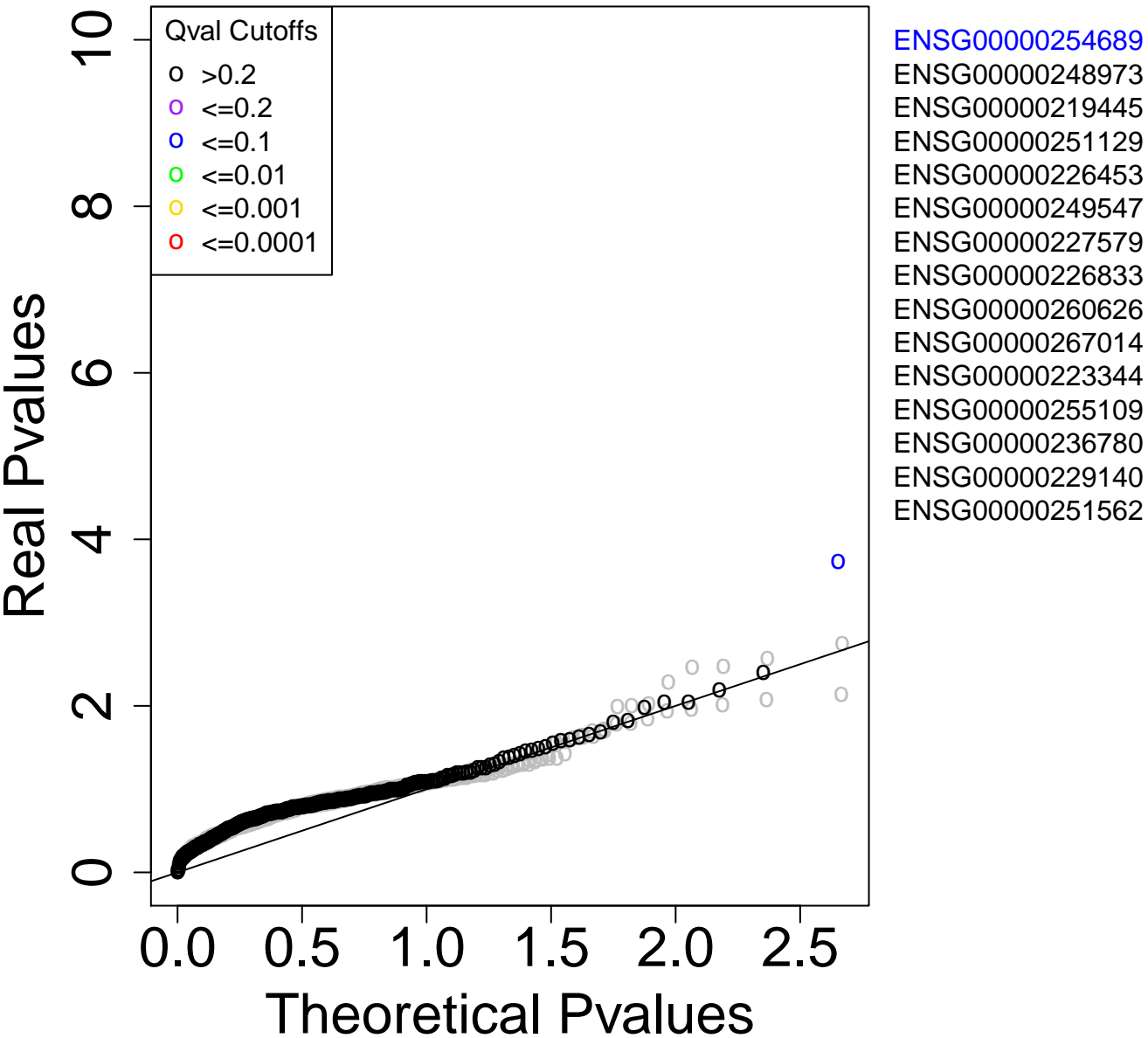
LncRNA CRC – 3750 genes



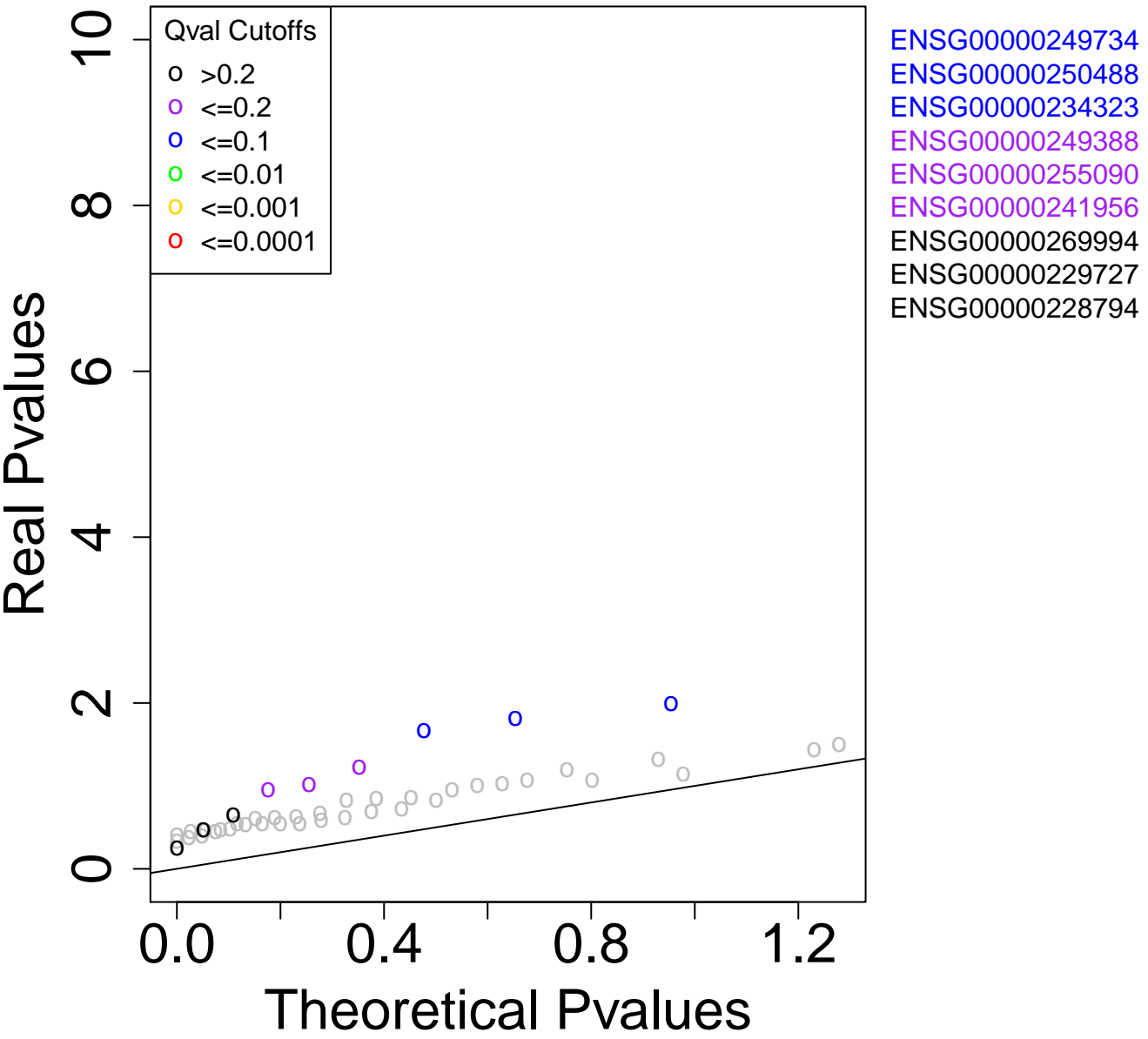
LncRNA GBM – 202 genes



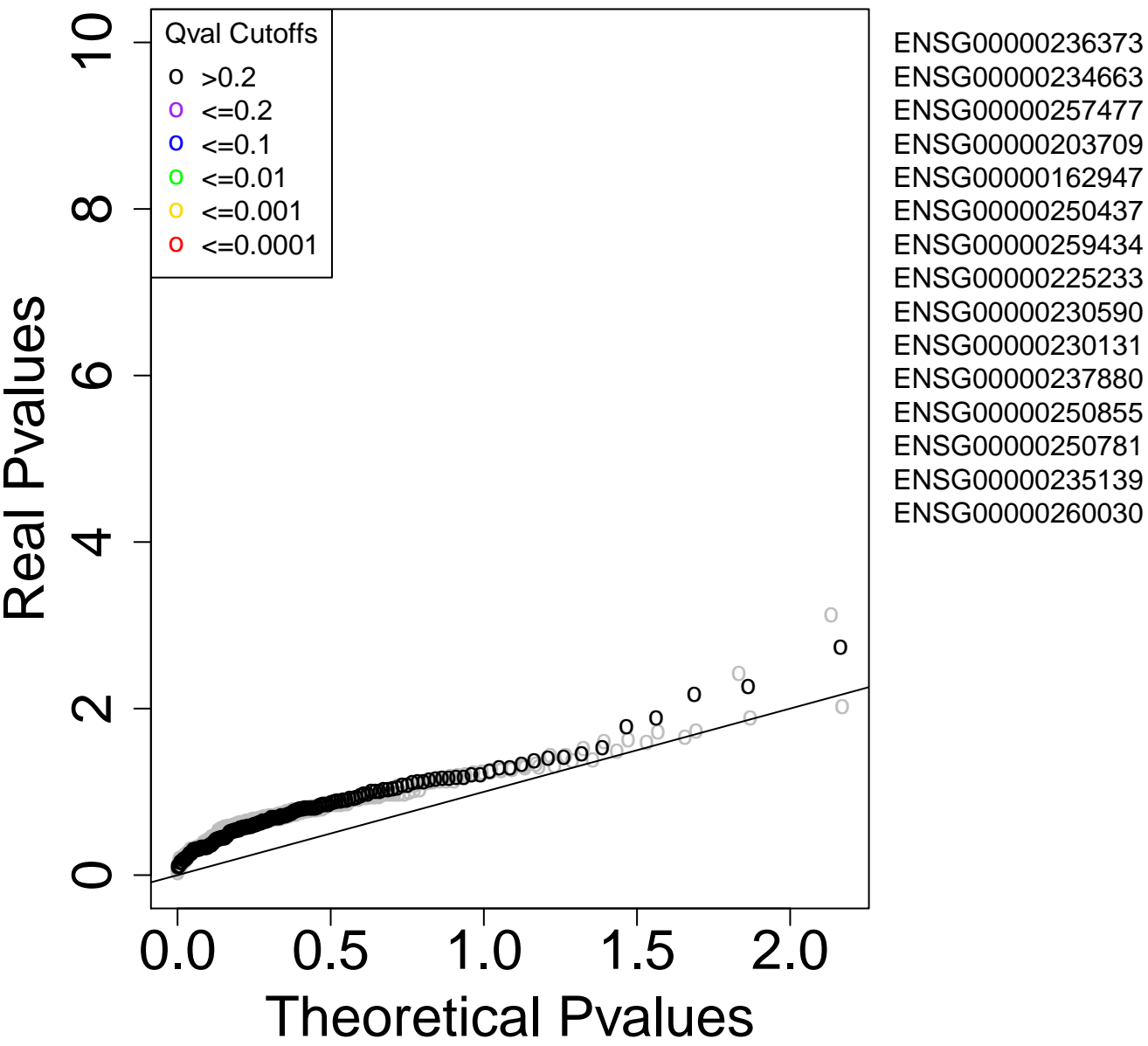
LncRNA HNSC – 450 genes



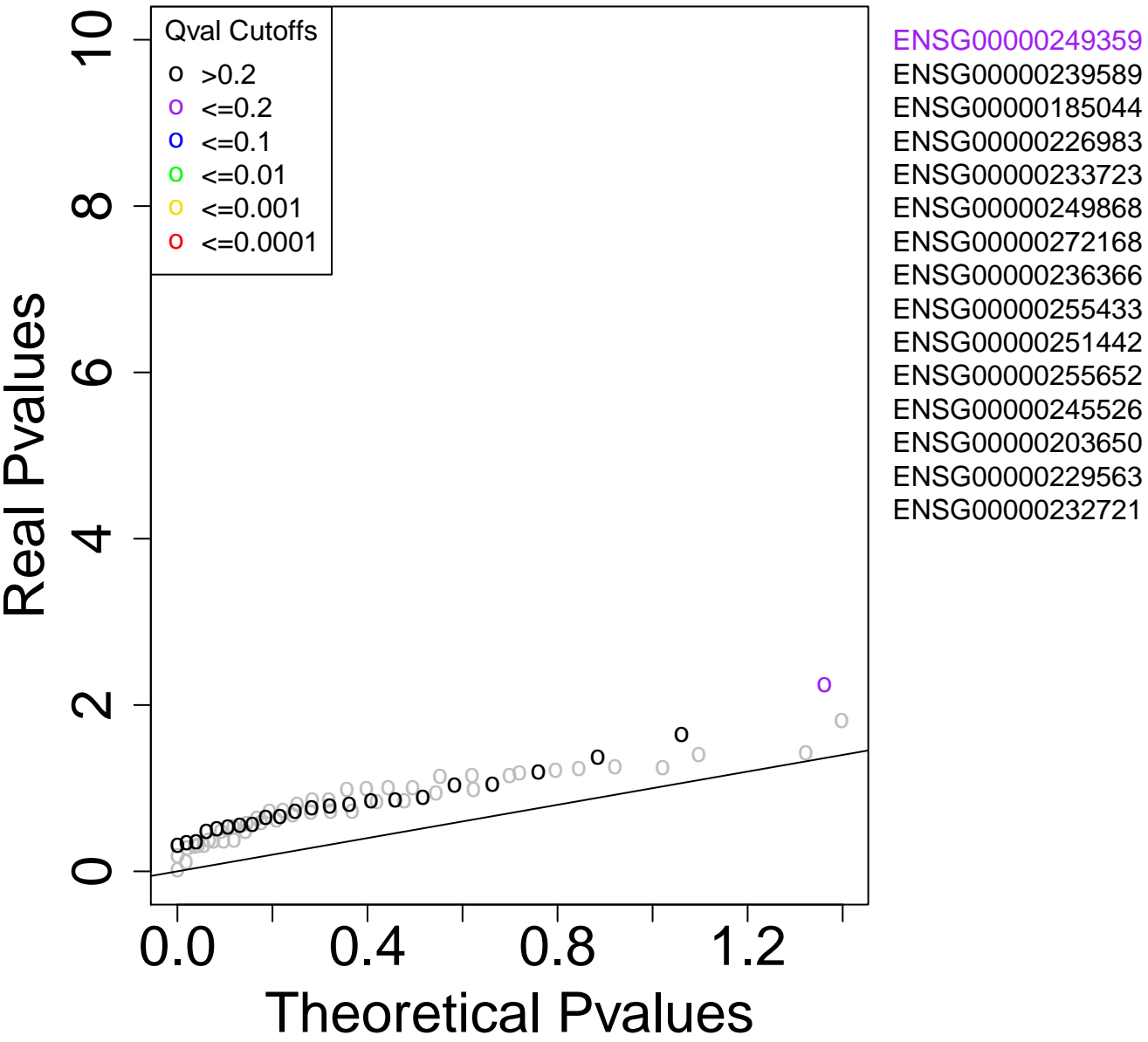
LncRNA KICH – 9 genes



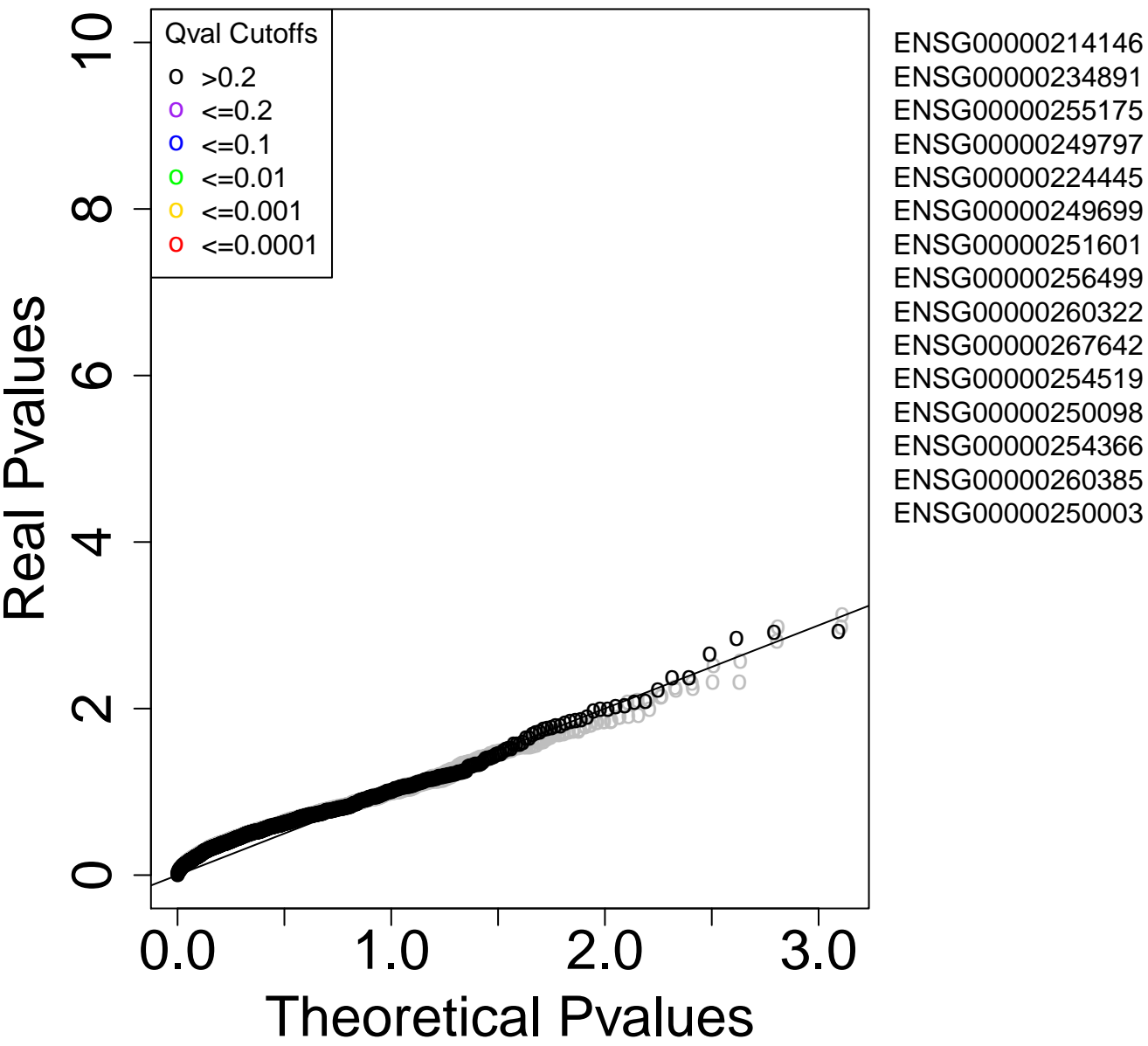
LncRNA KIRC – 146 genes



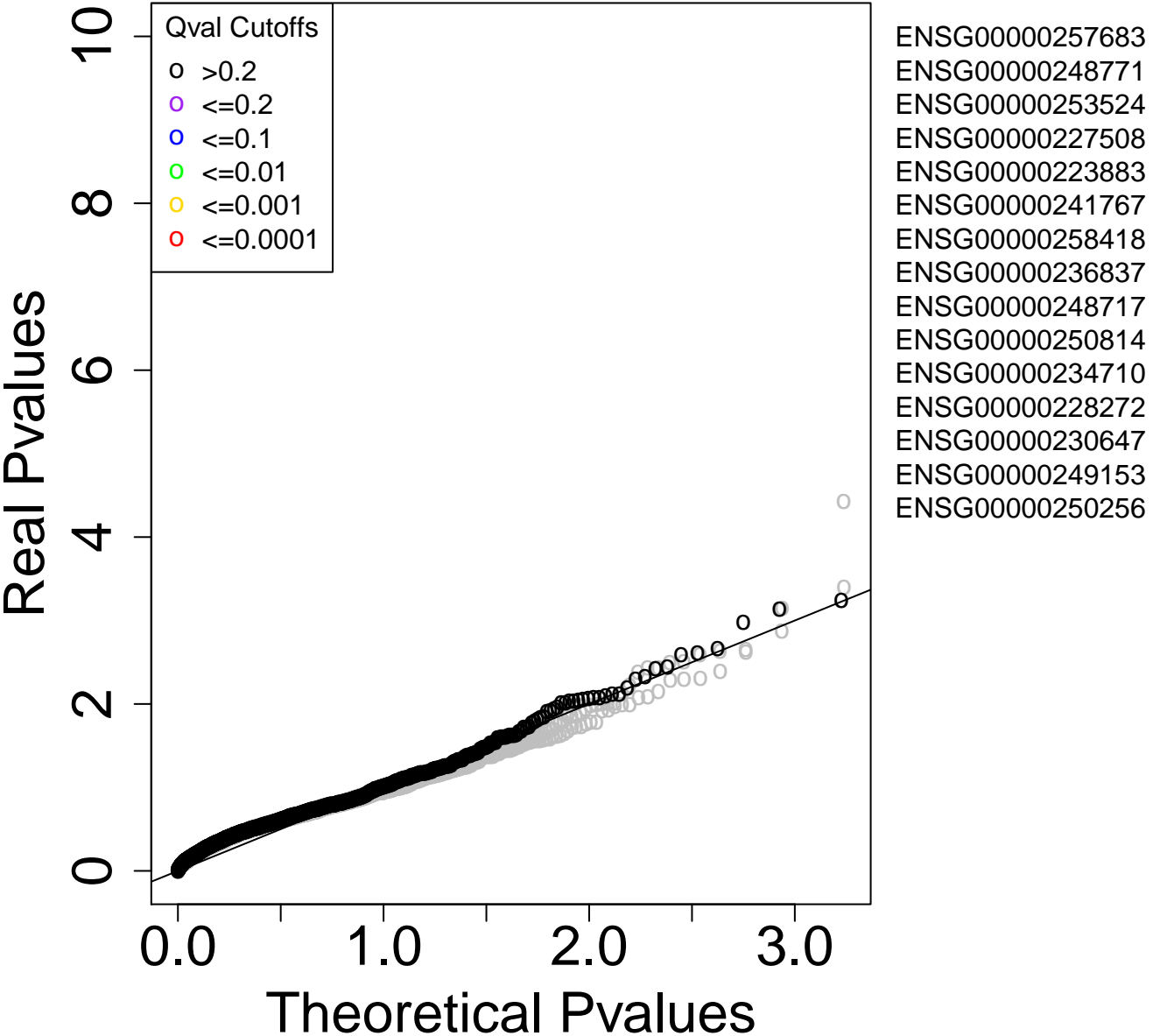
LncRNA LGG – 23 genes



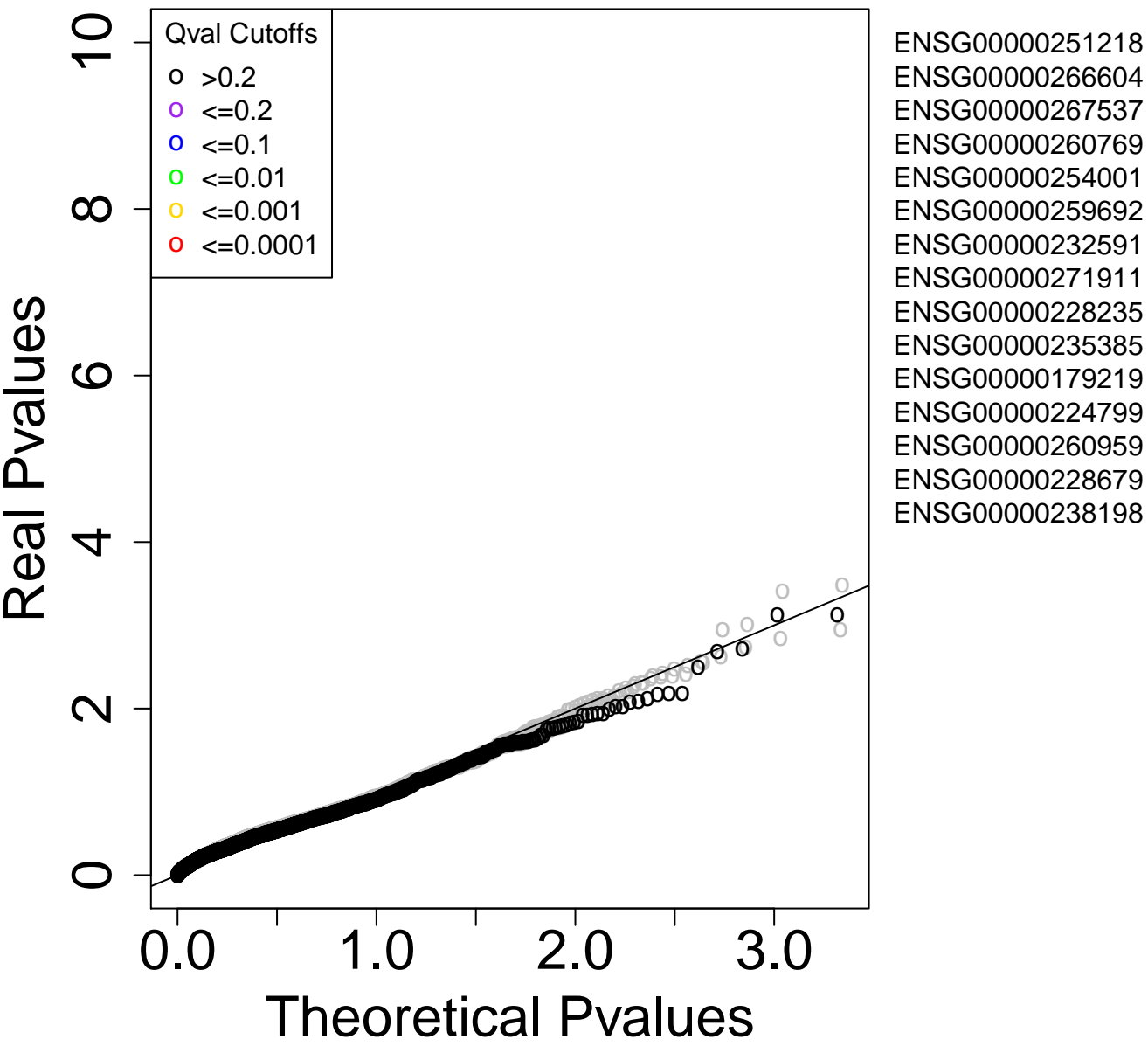
LncRNA Liver – 1237 genes



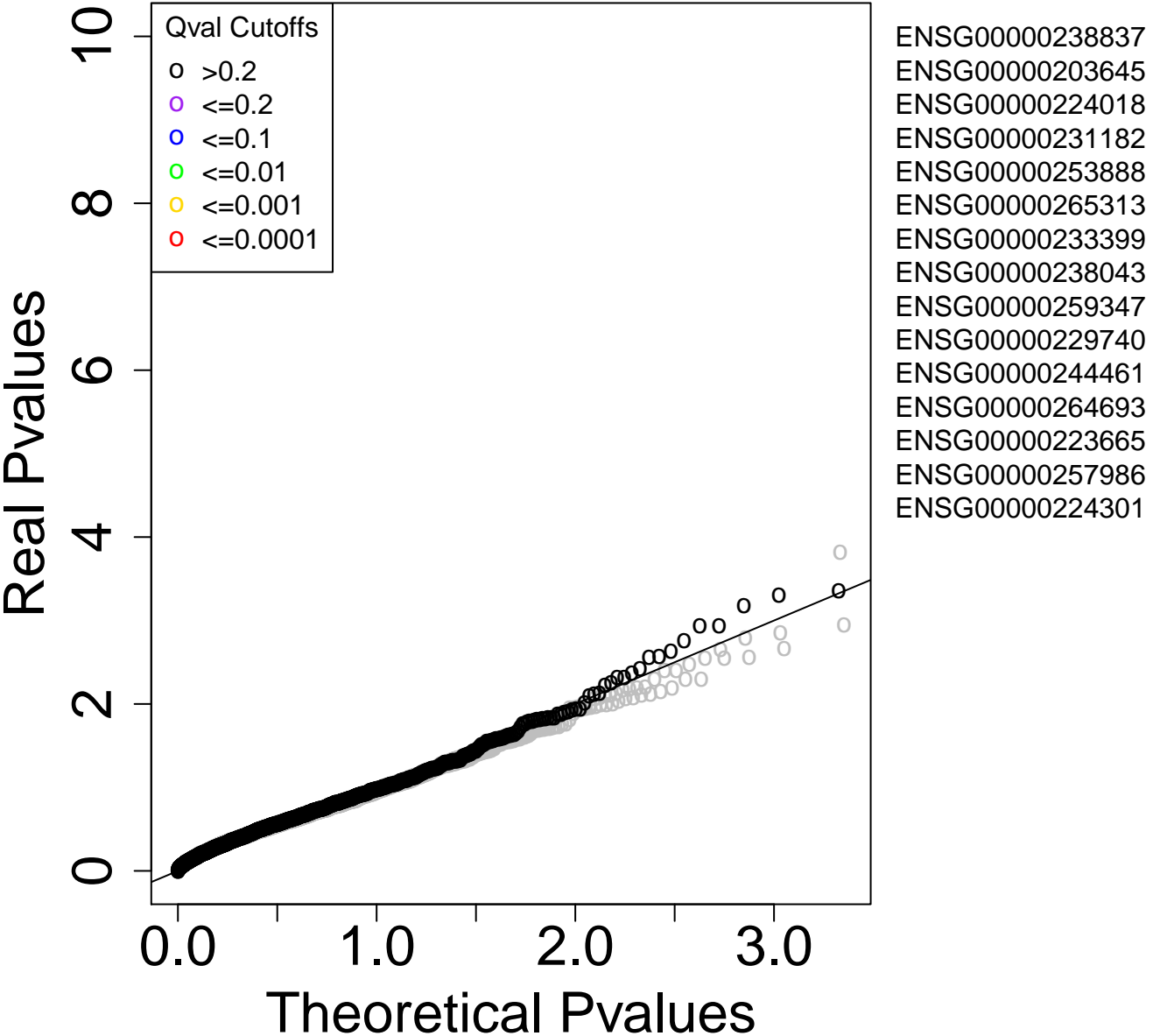
LncRNA LUAD – 1683 genes



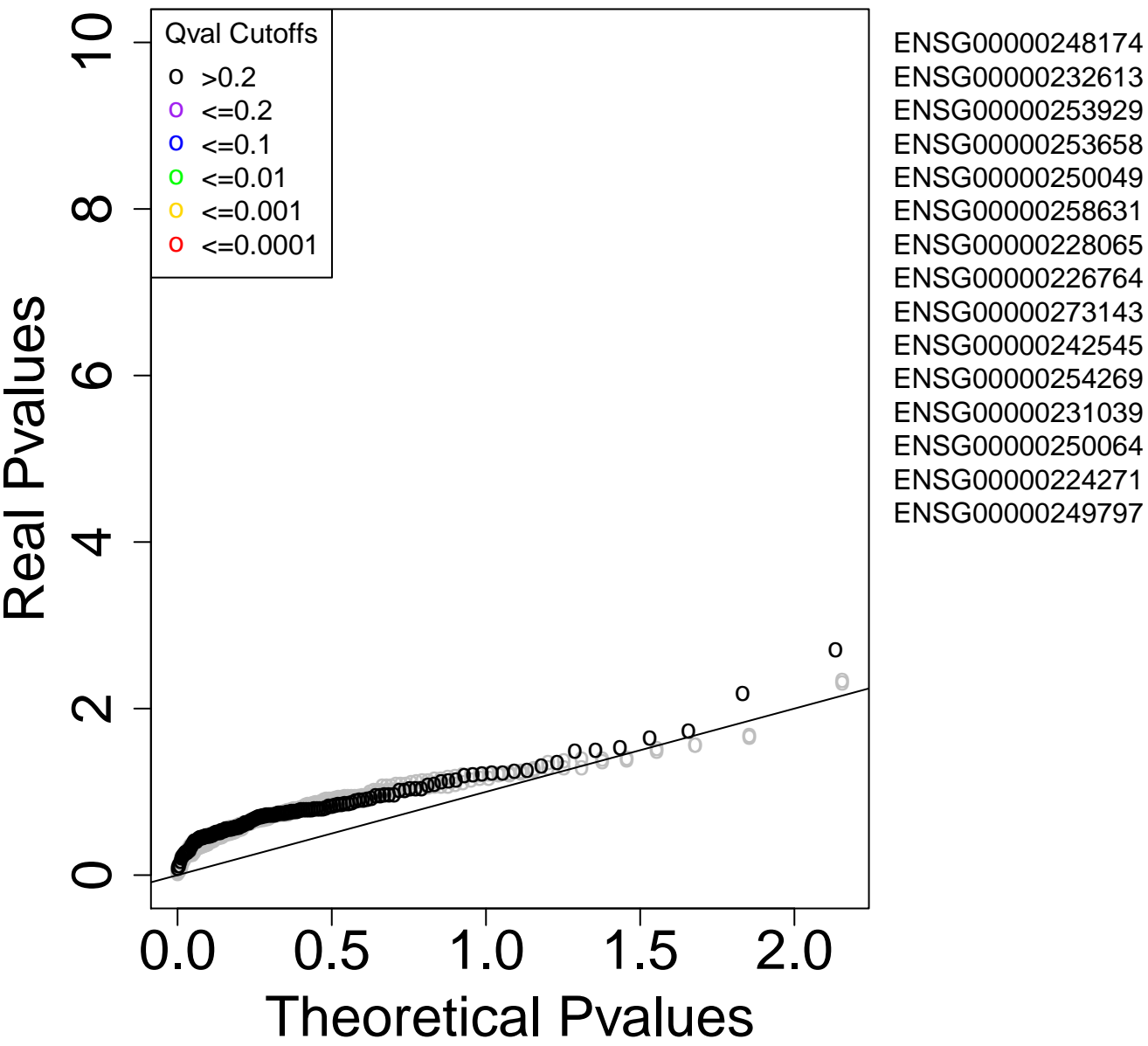
LncRNA Lung_adeno – 2073 genes



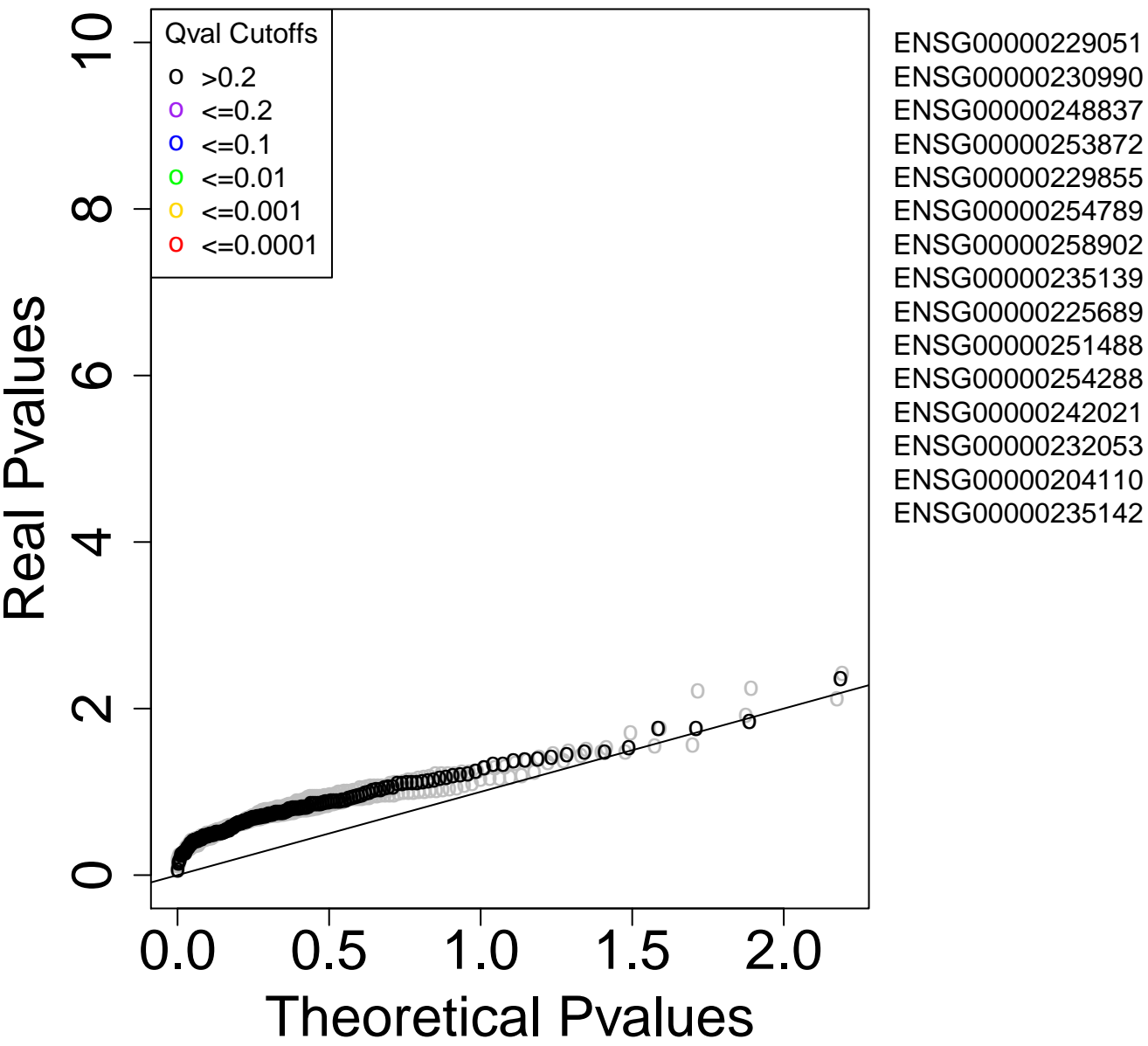
LncRNA LUSC – 2117 genes



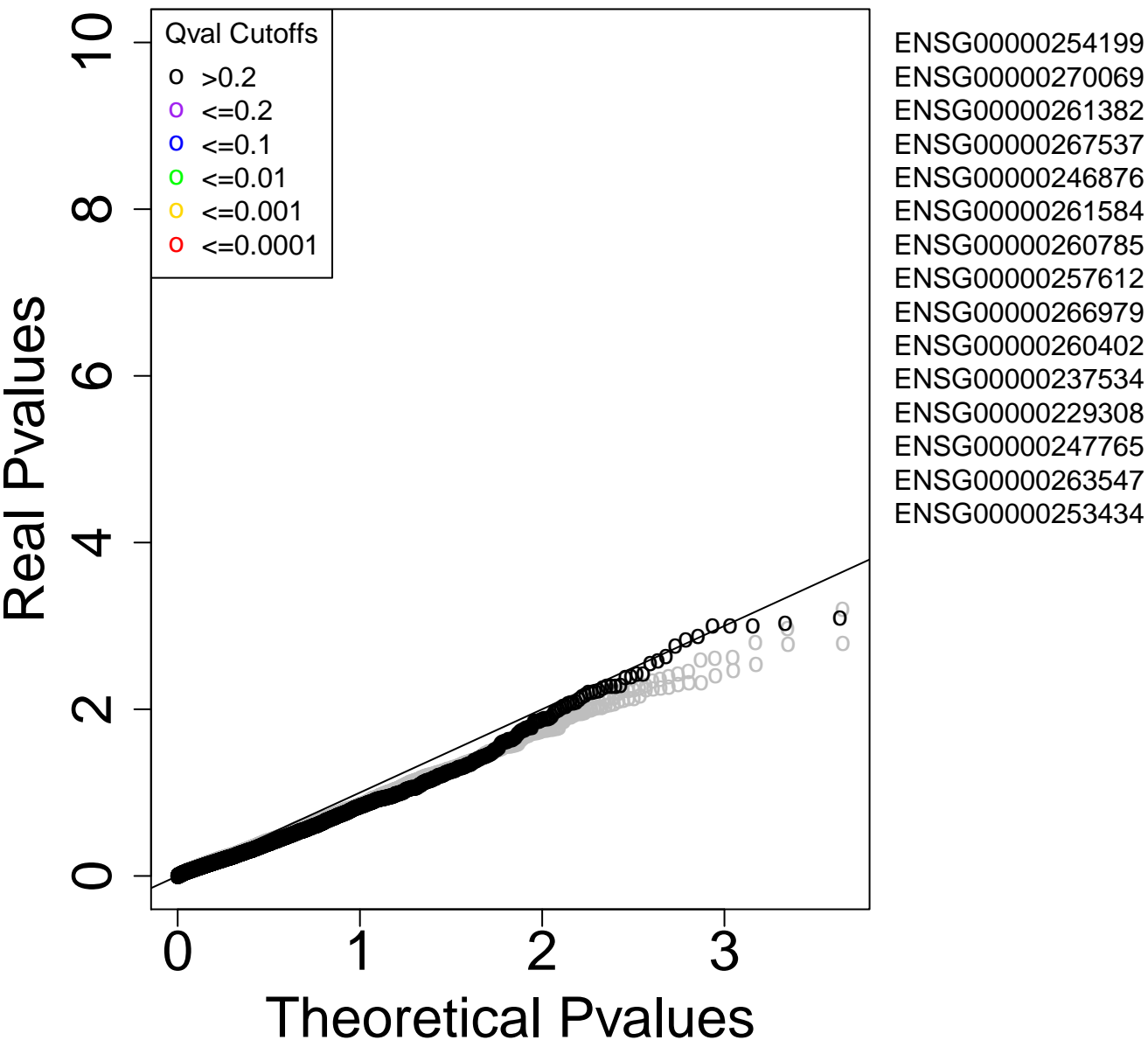
lncRNA Lymphoma_B-cell – 136 genes



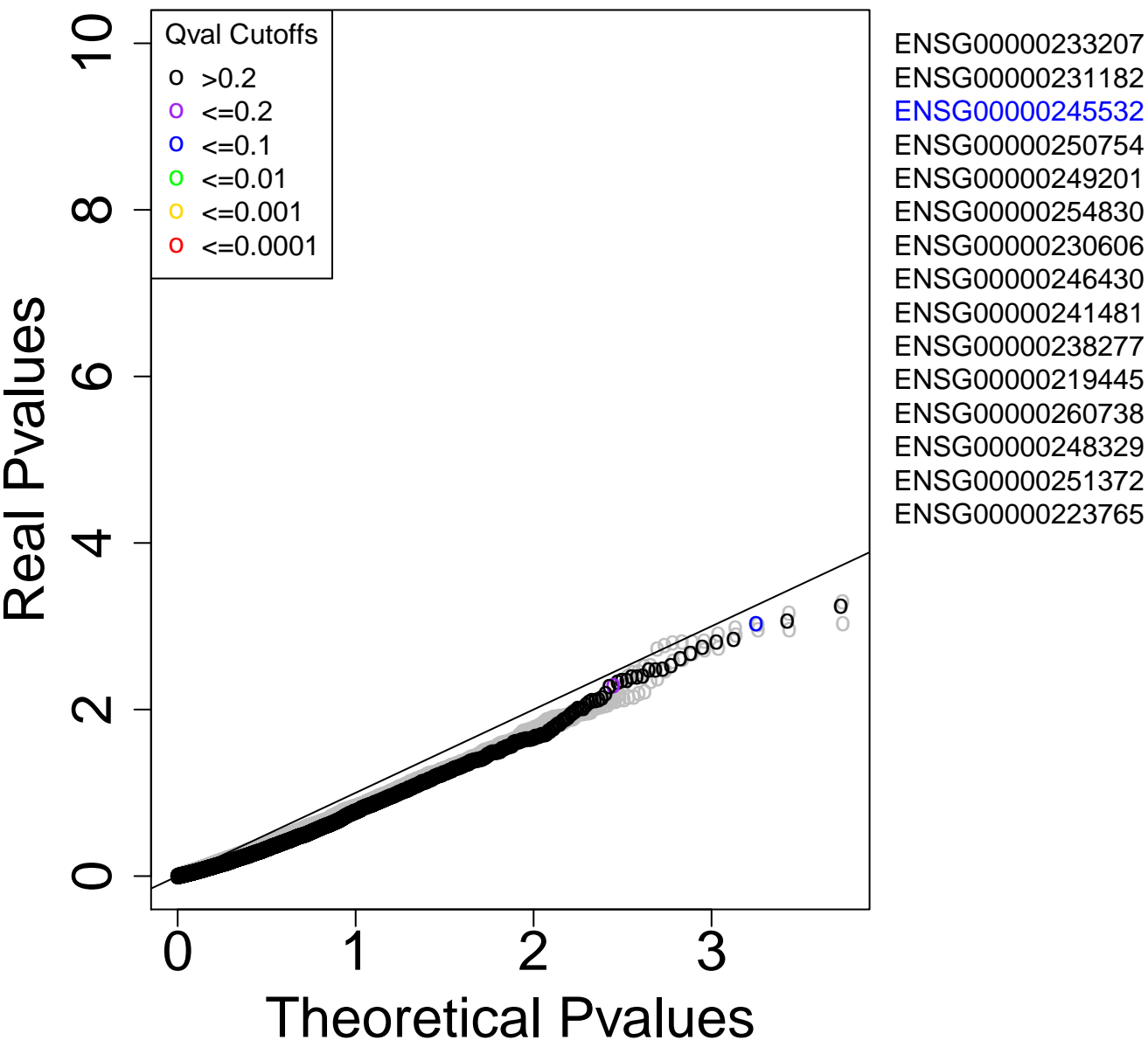
LncRNA Medulloblastoma – 154 genes



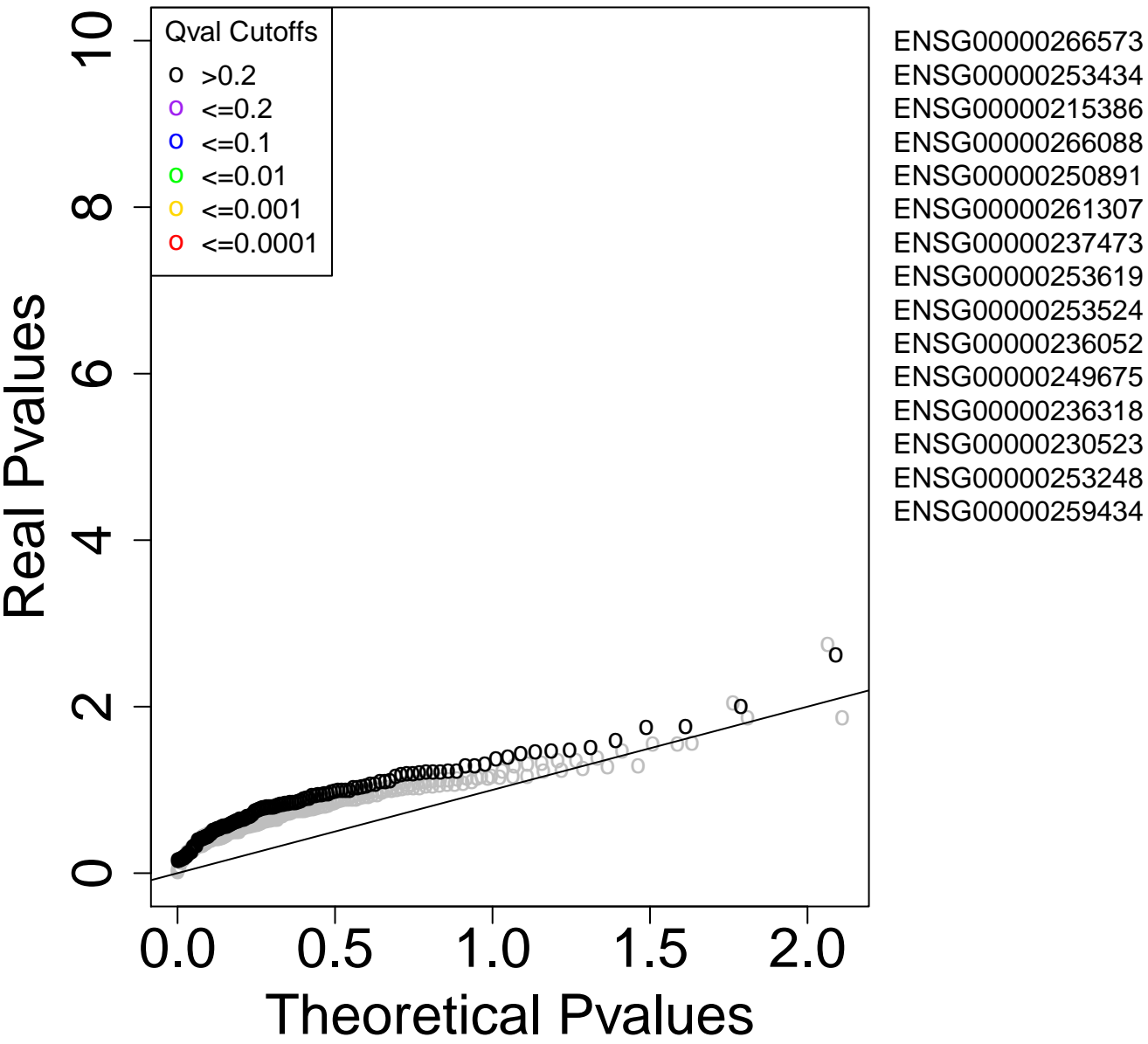
RNA Pancancer_Alexandrov – 4303 genes



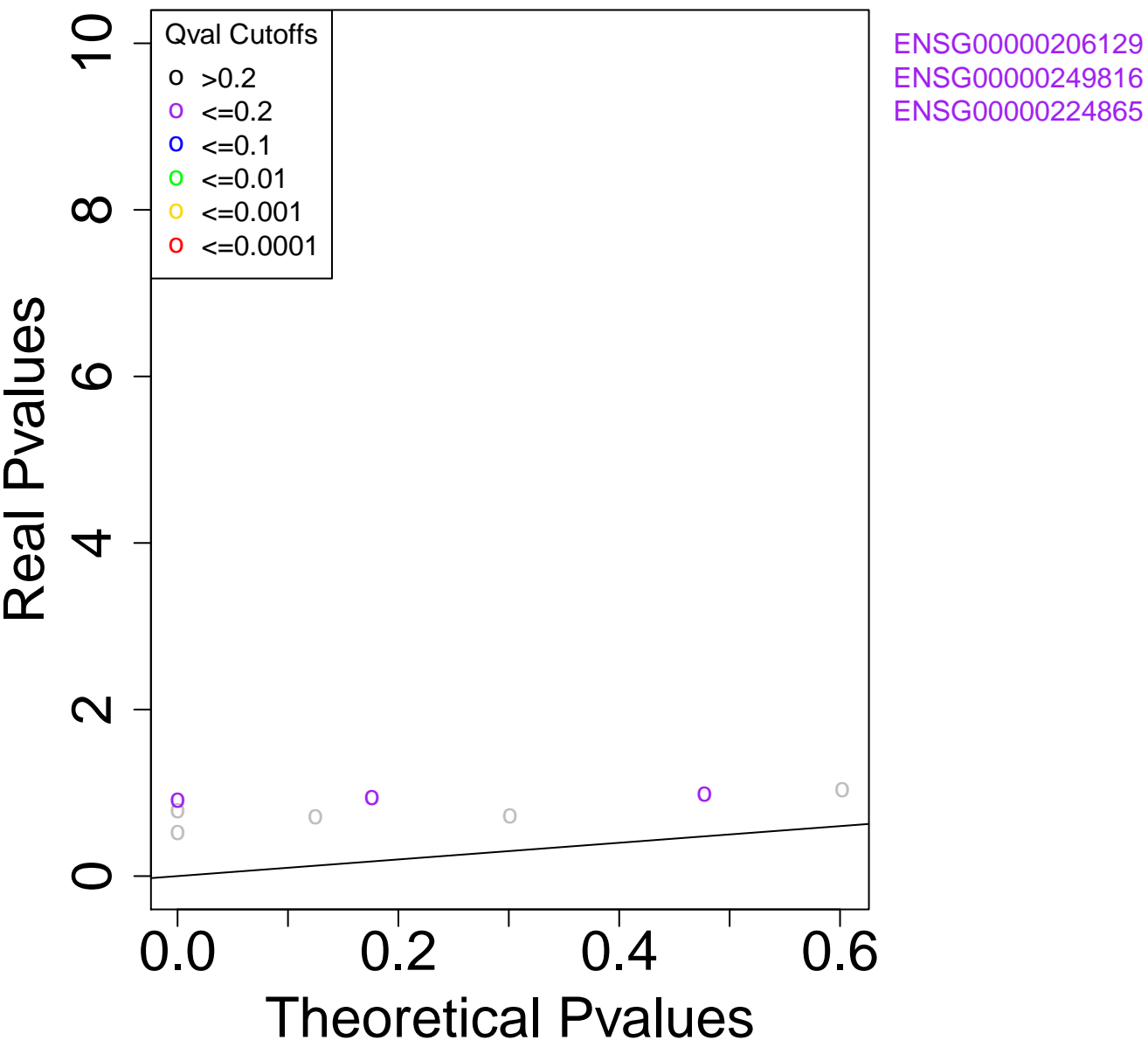
lncRNA Pancancer_TCGA – 5331 genes



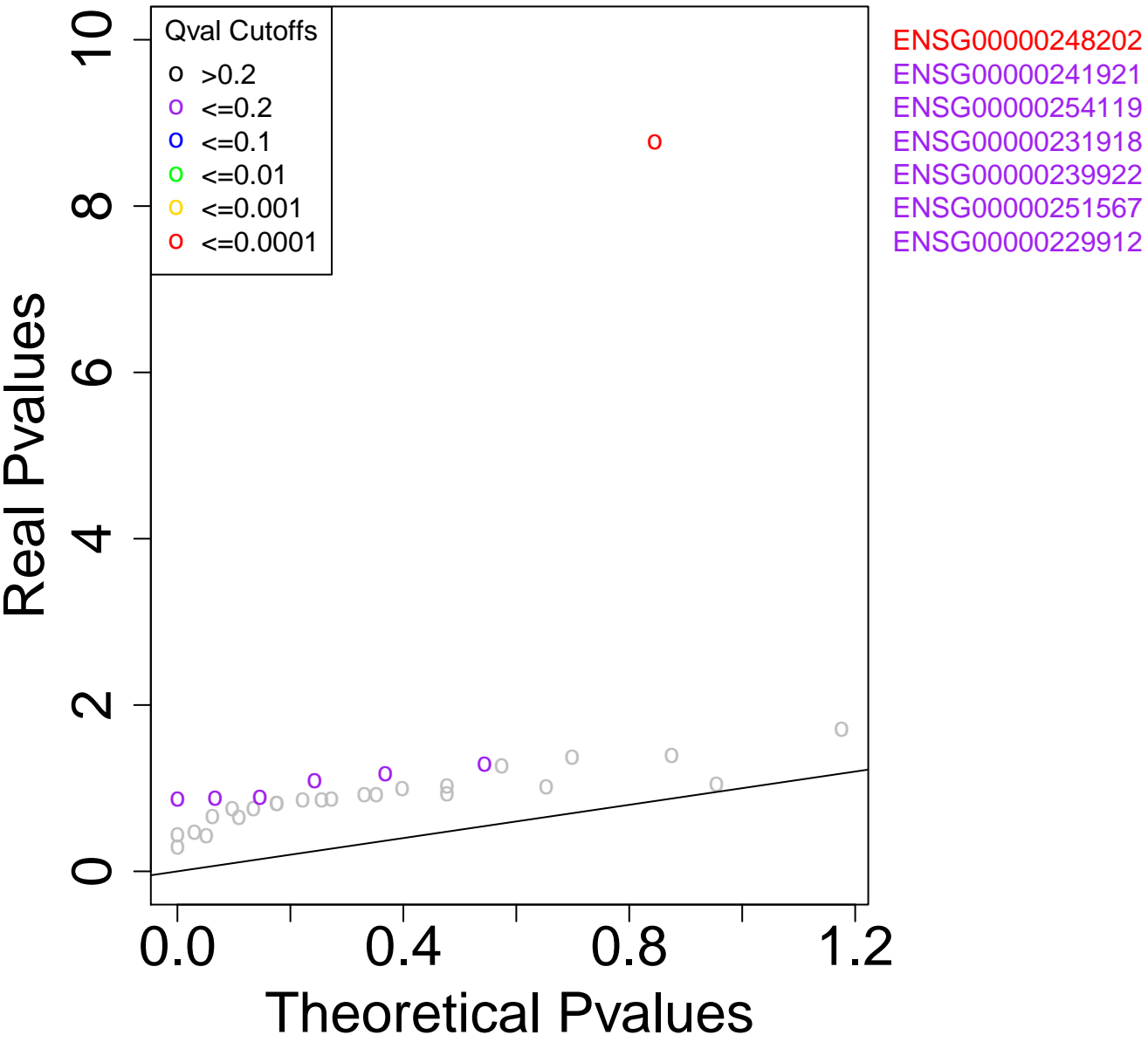
LncRNA Pancreas – 123 genes



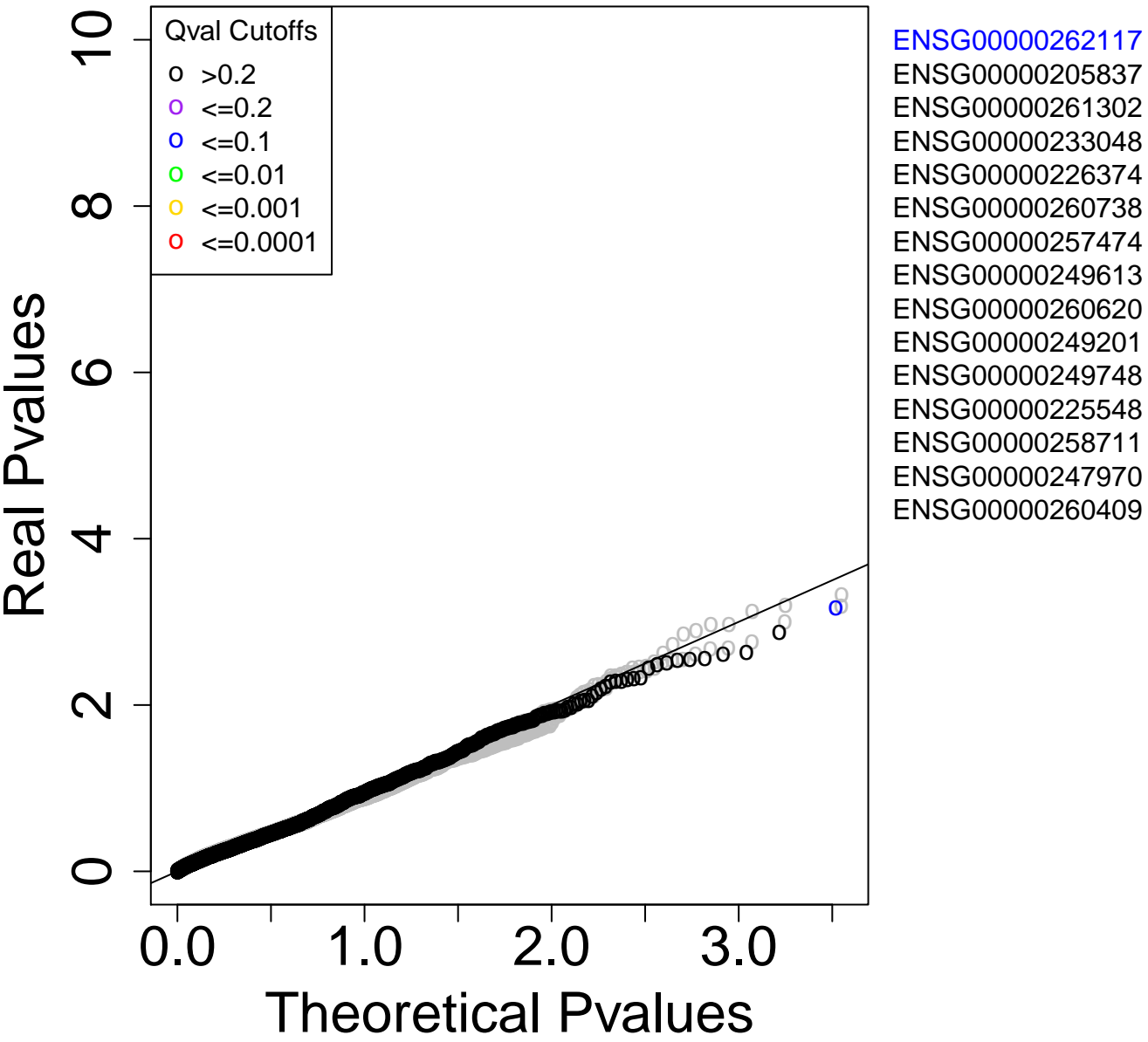
ncRNA Pilocytic_astrocytoma – 3 genes



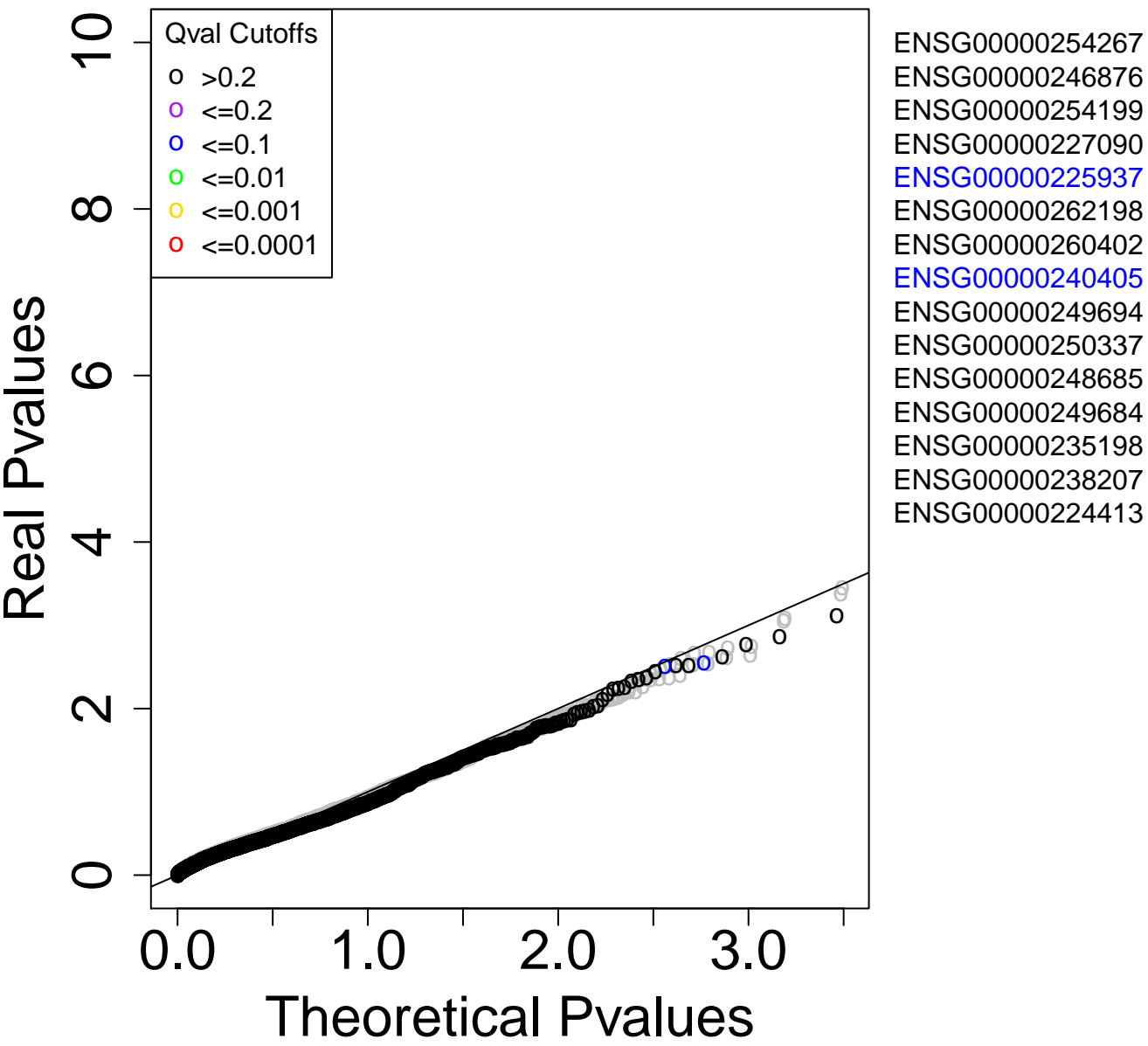
LncRNA PRAD – 7 genes



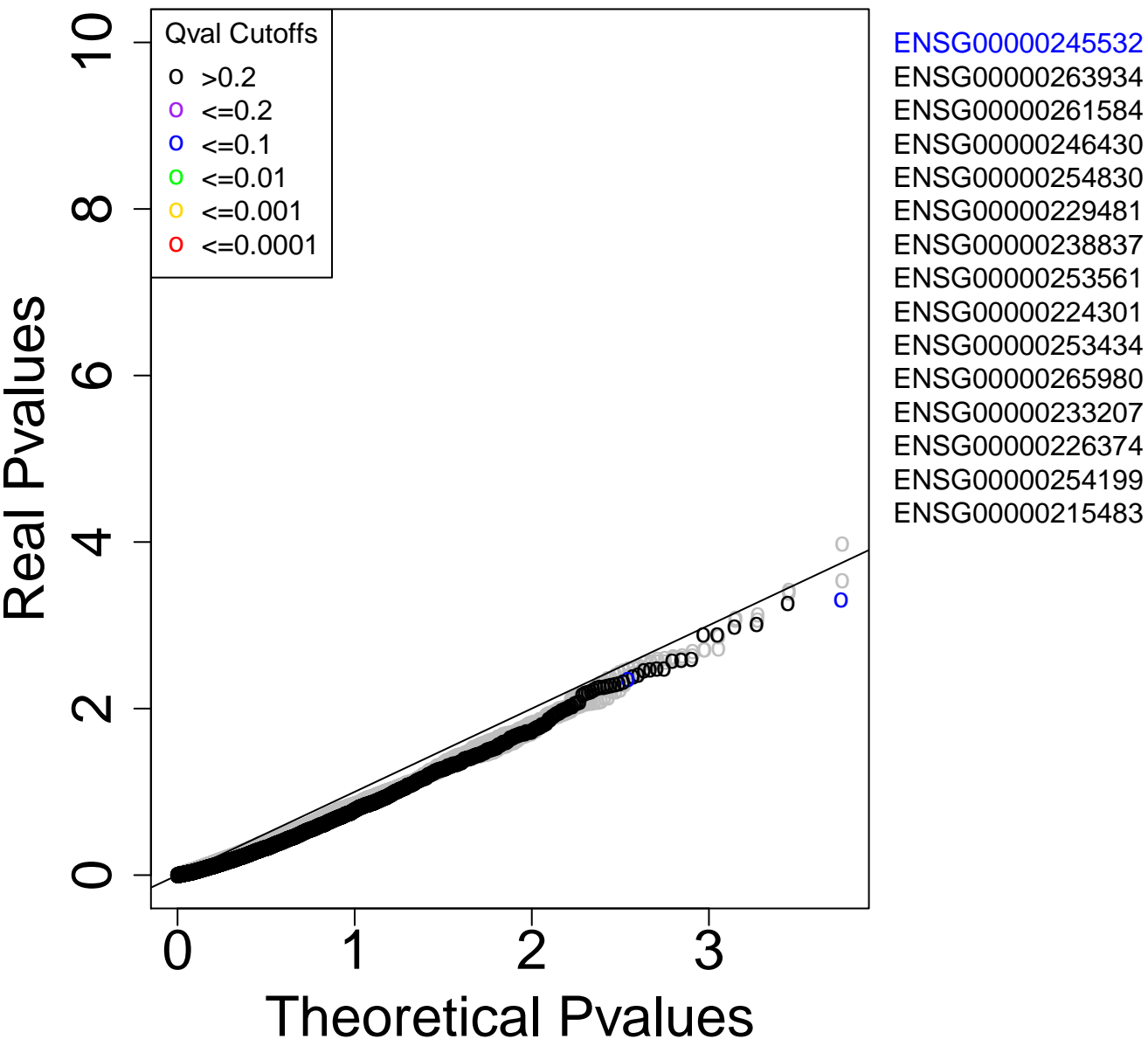
LncRNA SKCM – 3300 genes



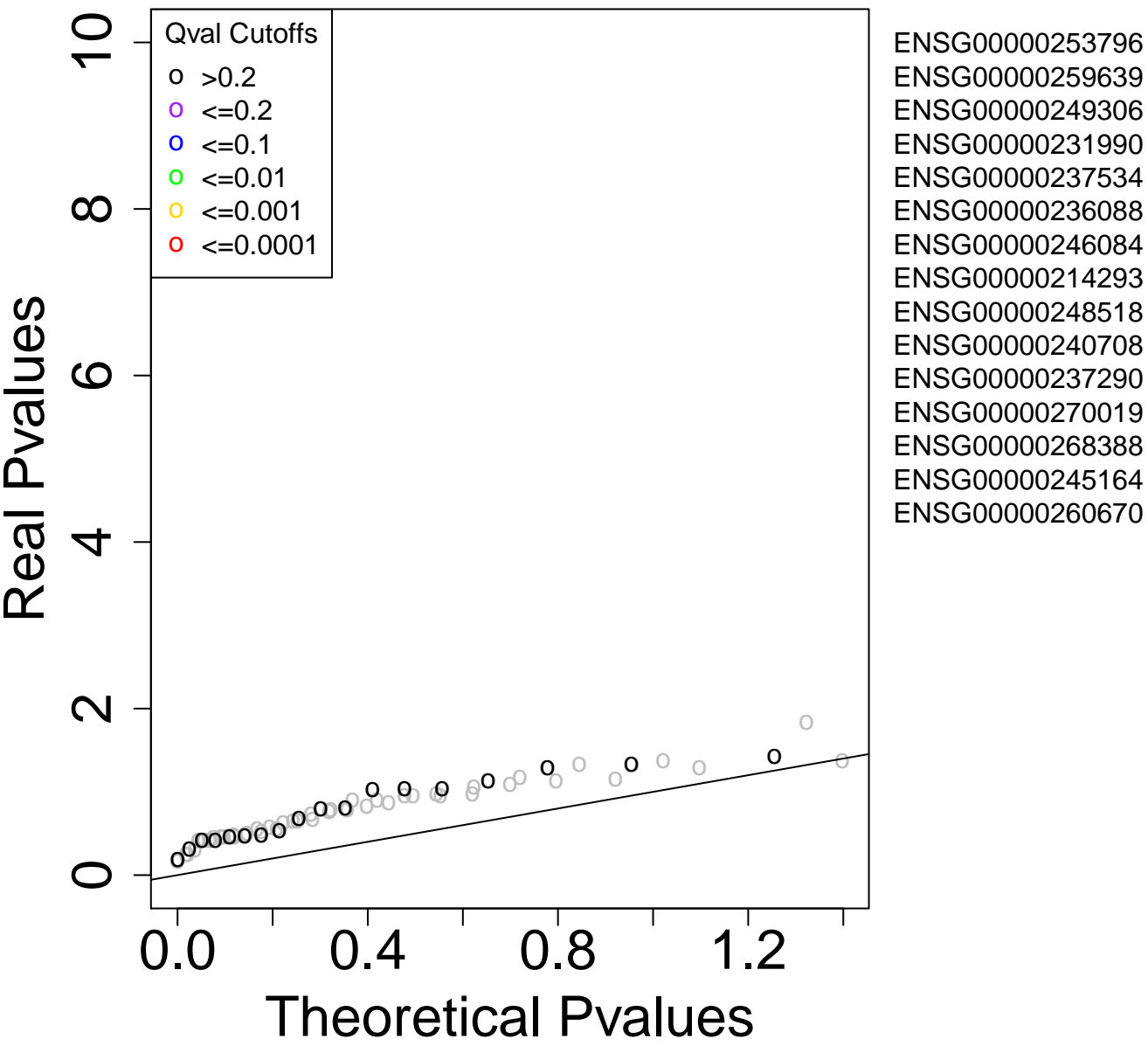
LncRNA Stad – 2912 genes



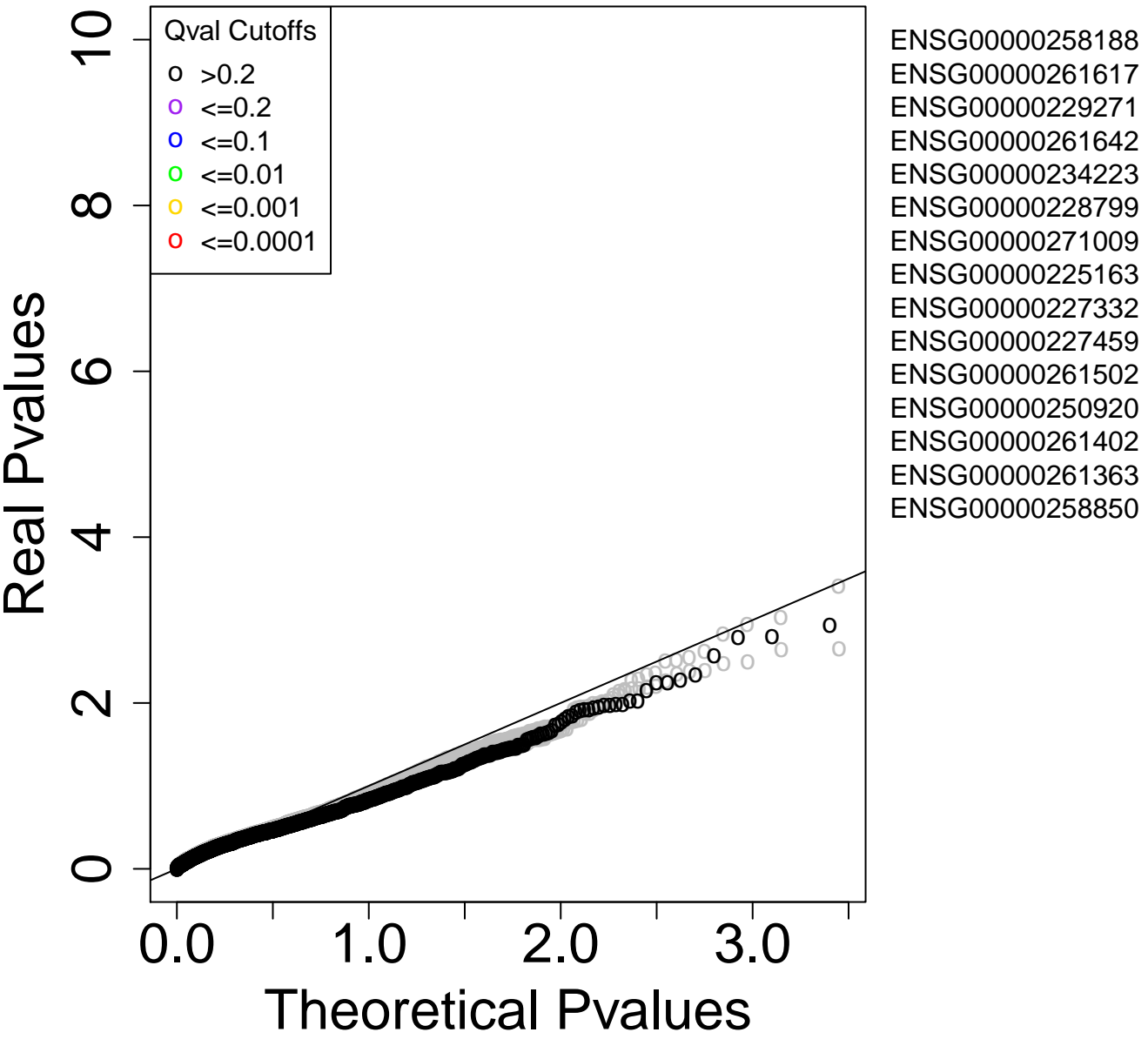
LncRNA Superpancancer – 5583 genes



LncRNA THCA – 18 genes



LncRNA UCEC – 2521 genes

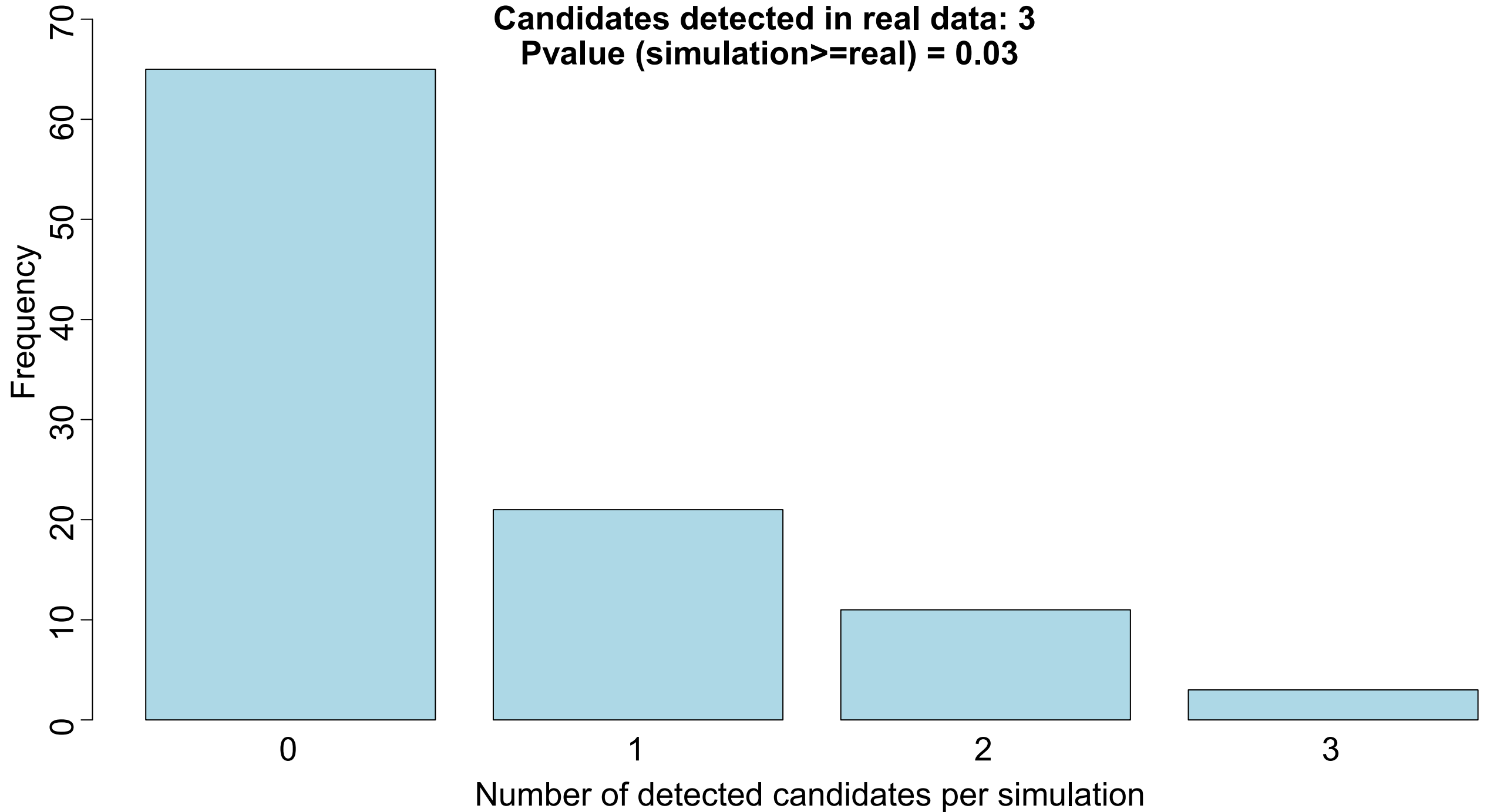


Supplementary Figure S4

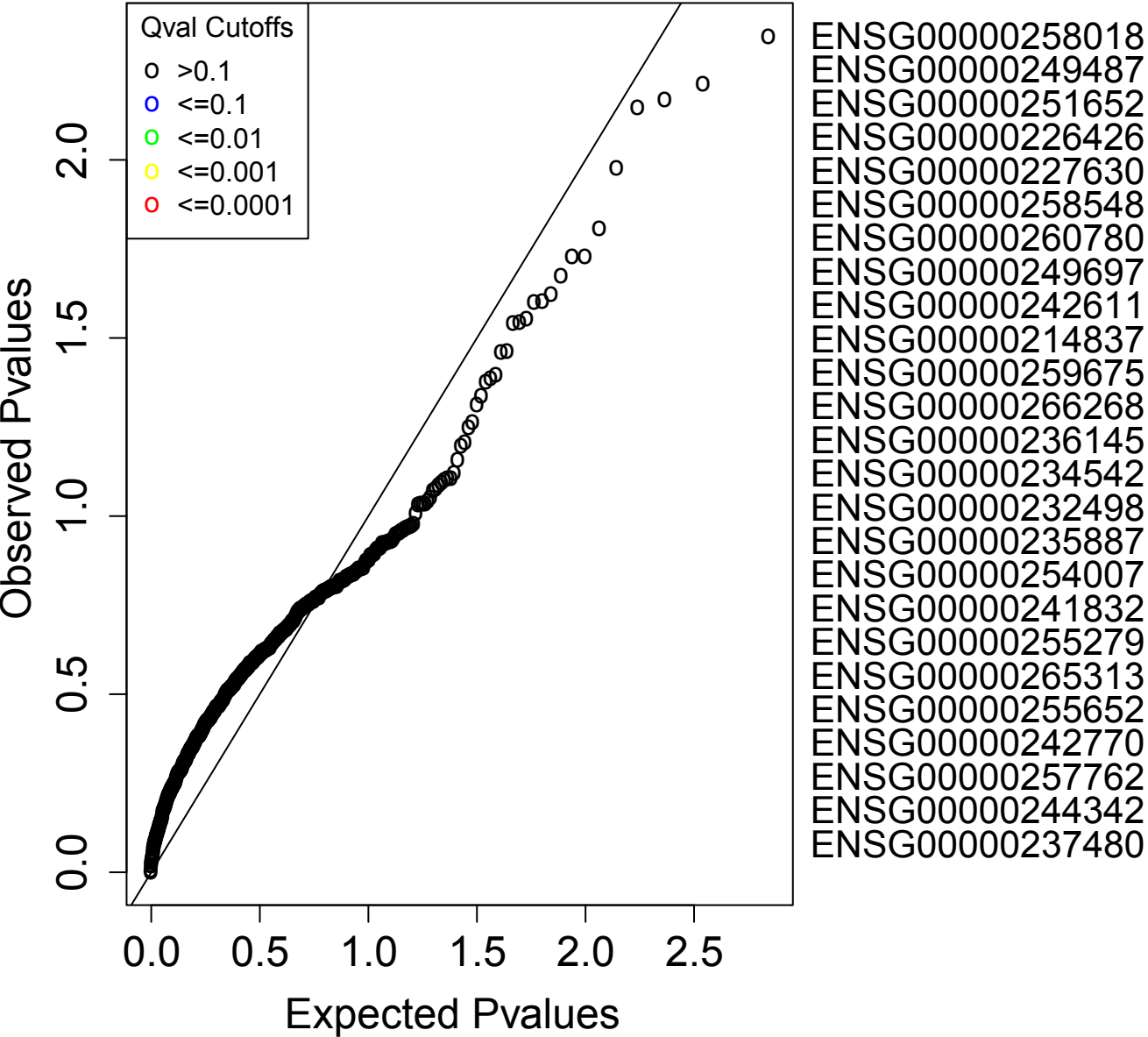
LncRNAs detected in Breast after 100 simulations

Candidates detected in real data: 3

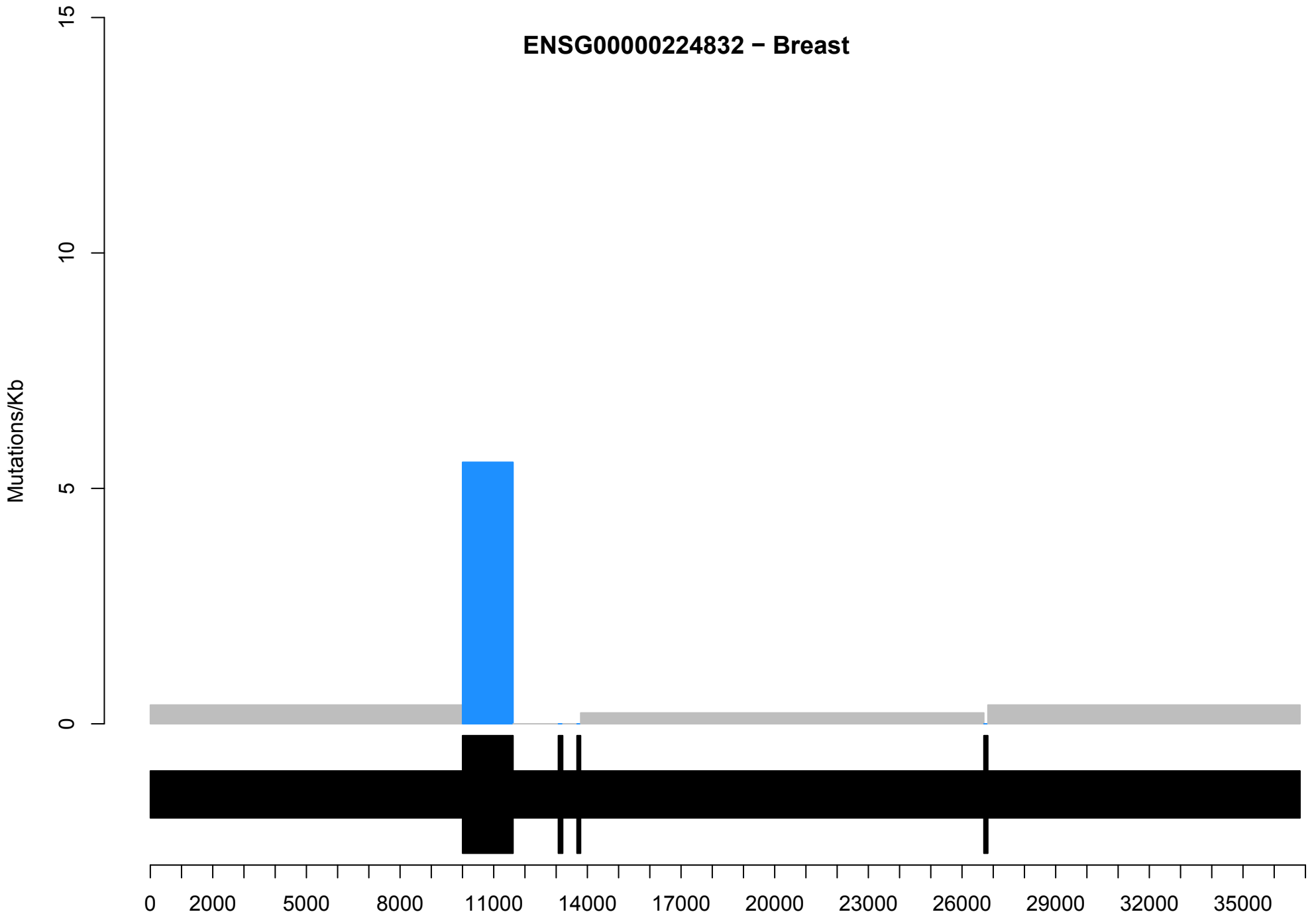
Pvalue (simulation \geq real) = 0.03



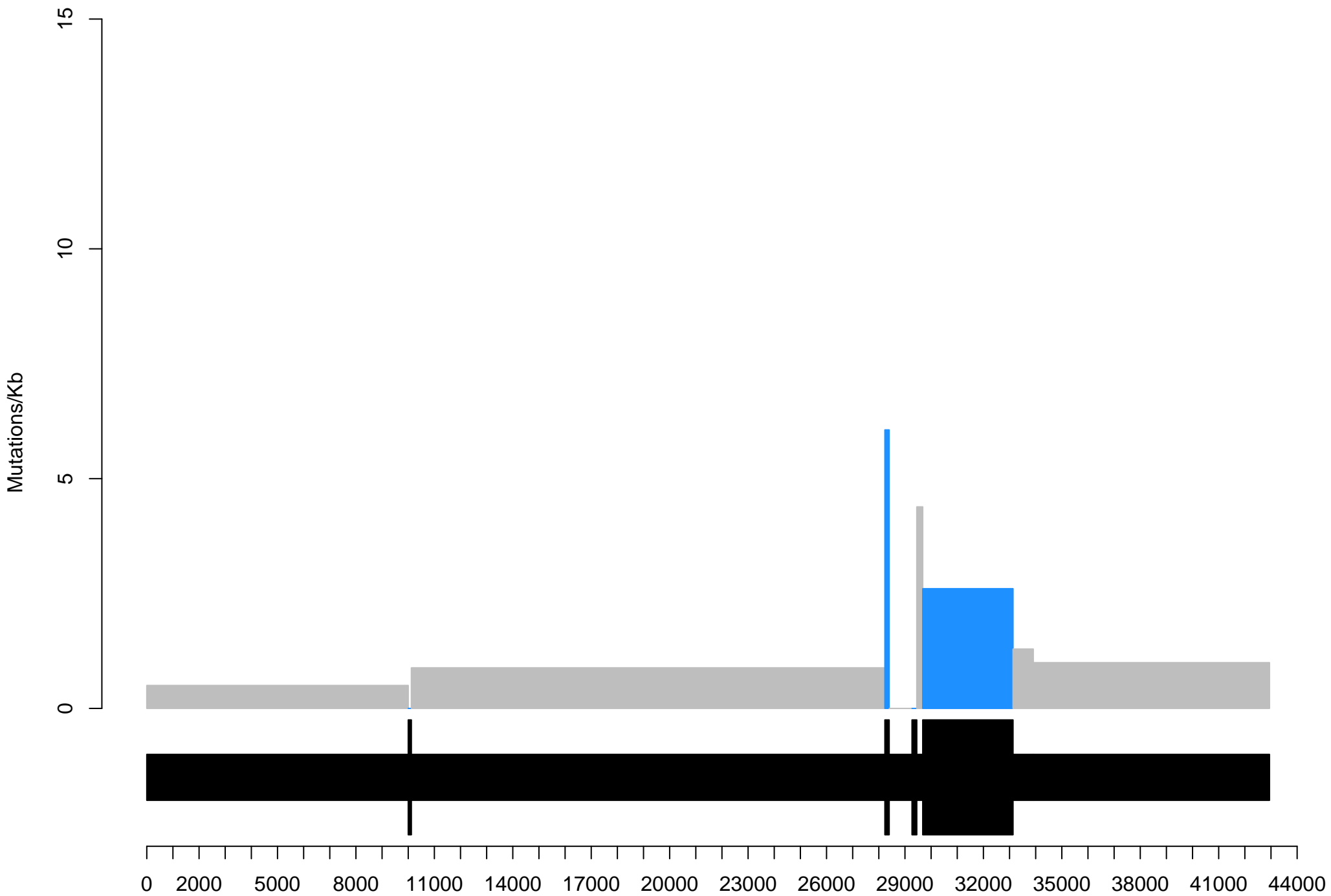
Breast QQplot – 694 genes



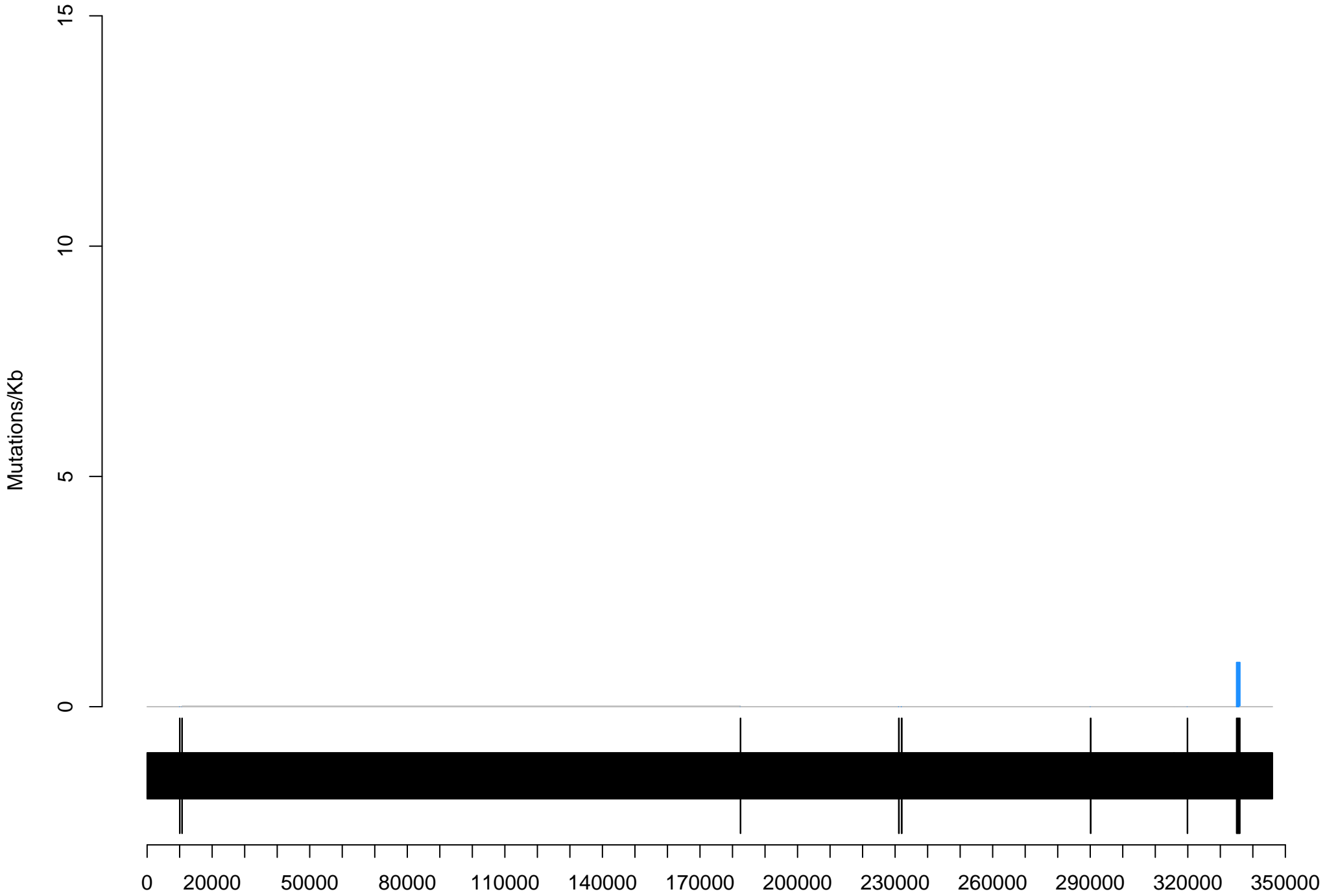
Supplementary Figure S7



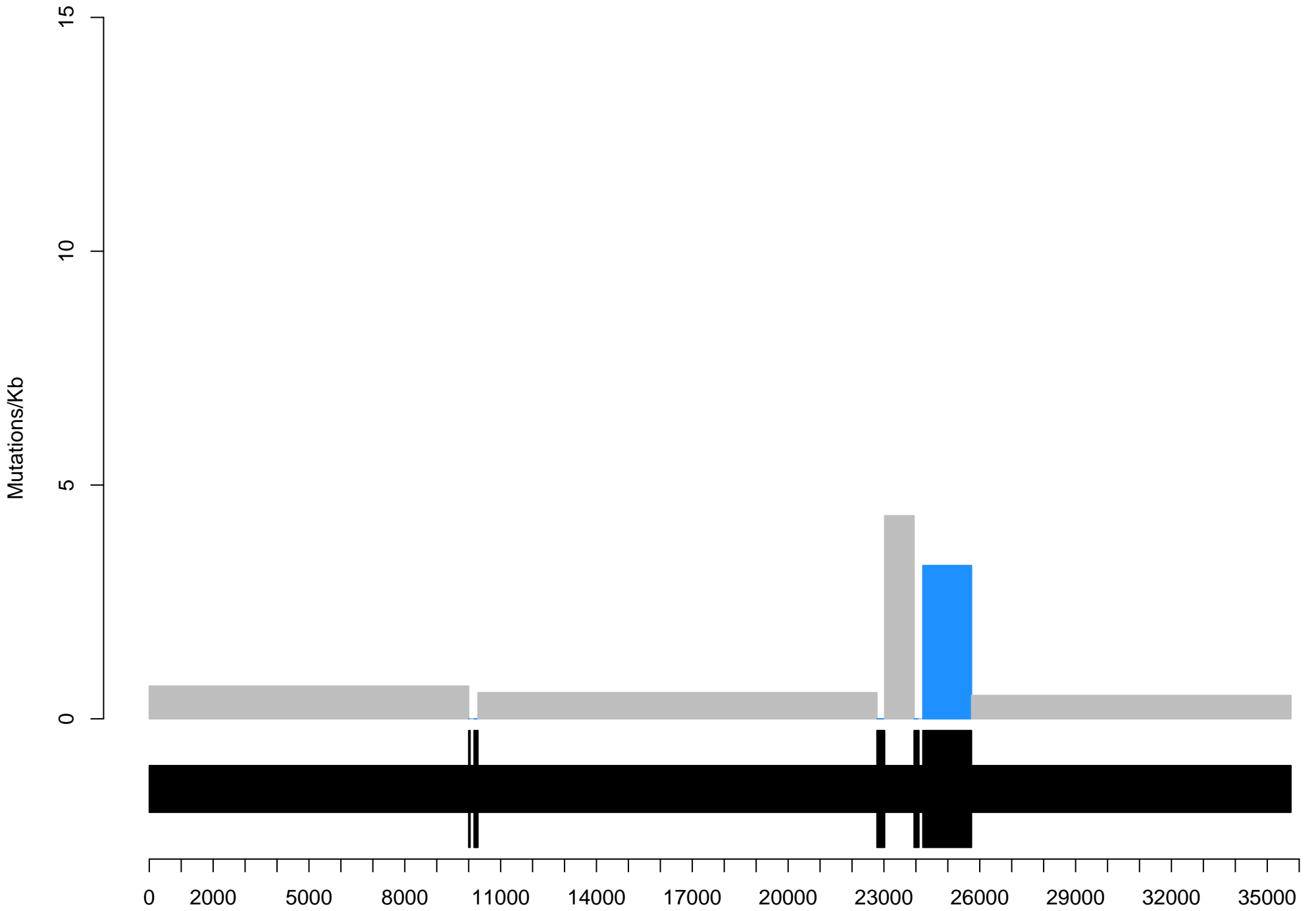
ENSG00000225937 – Stad



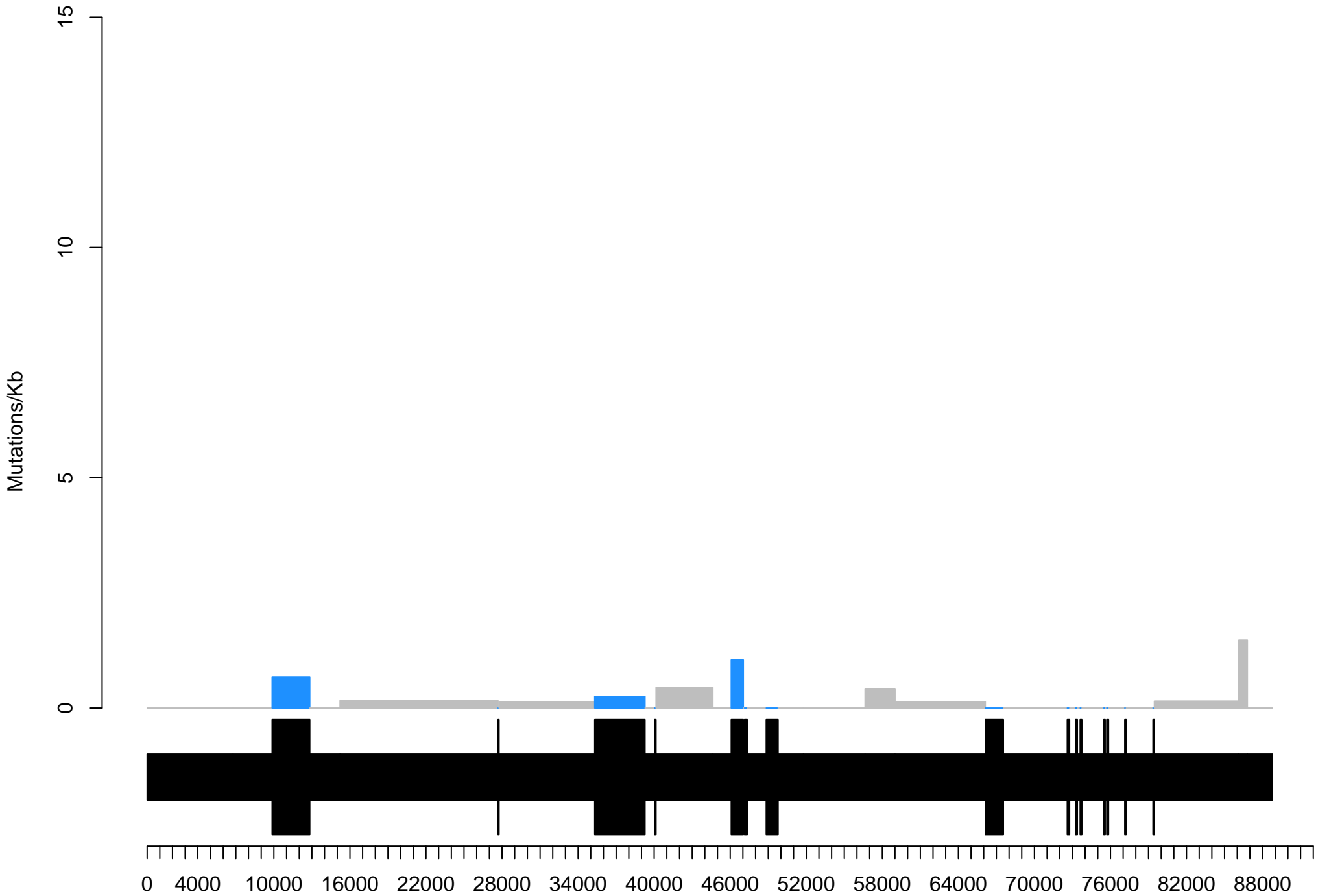
ENSG00000234323 – KICH



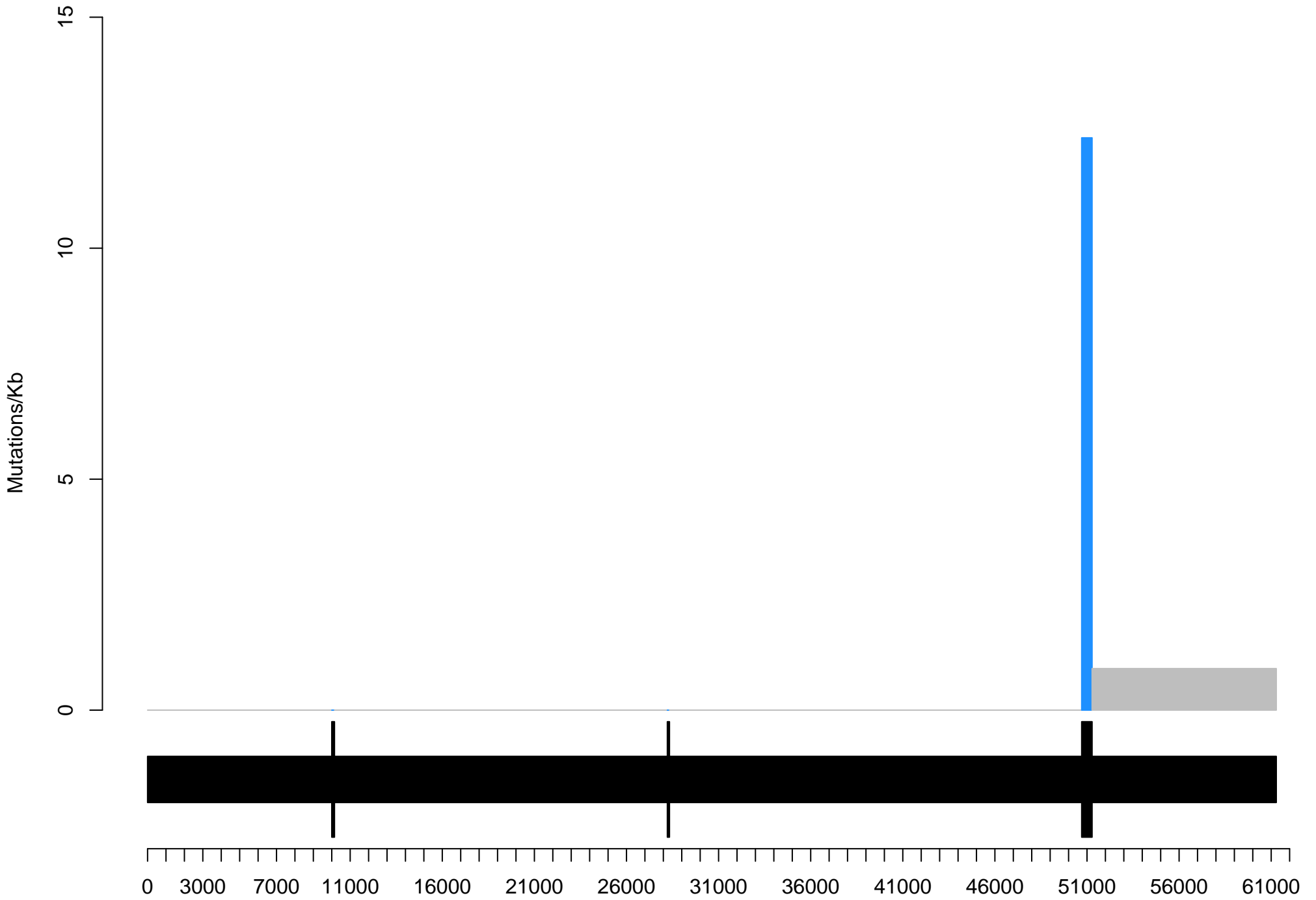
ENSG00000240405 – Stad



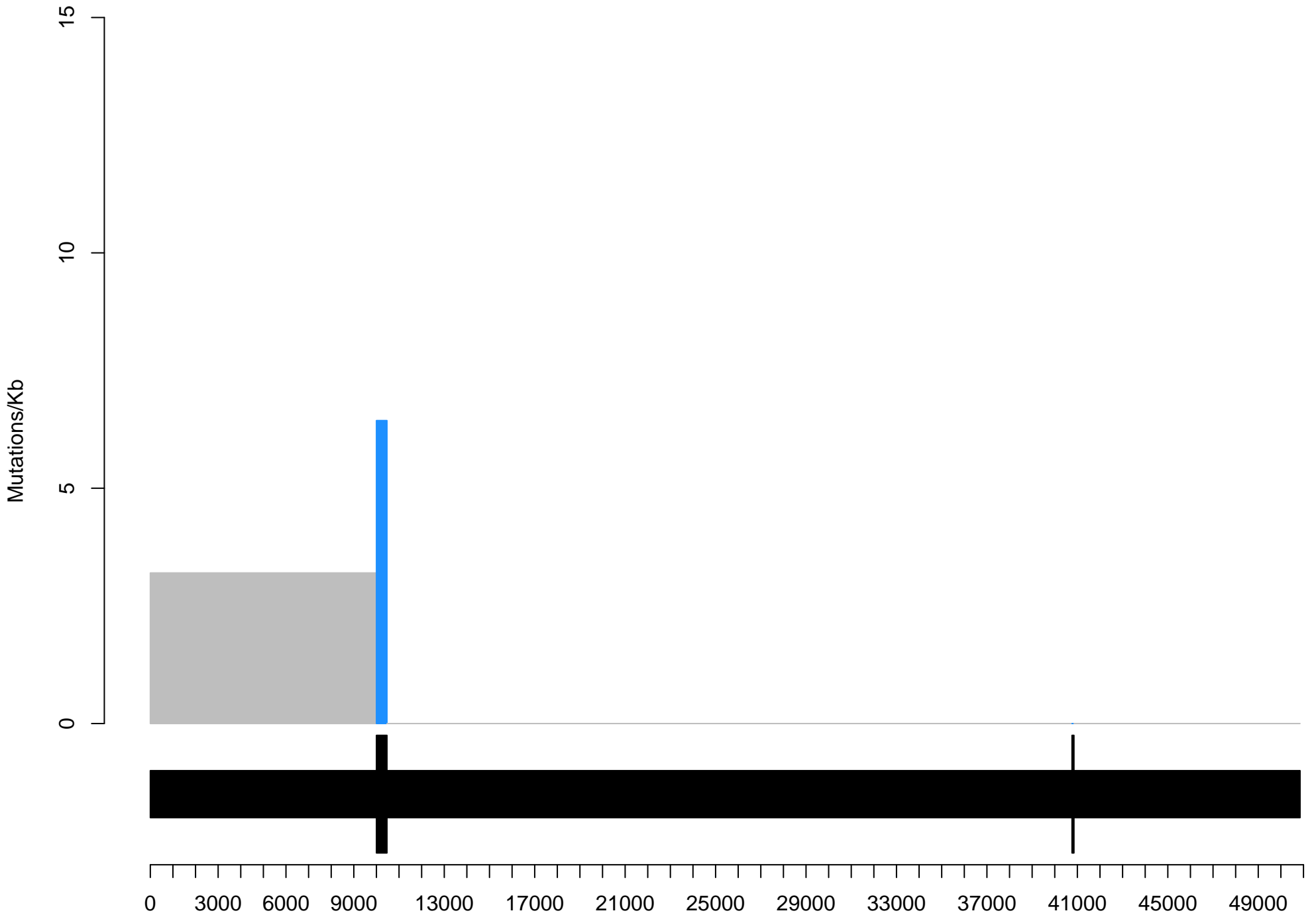
ENSG00000244306 - GBM



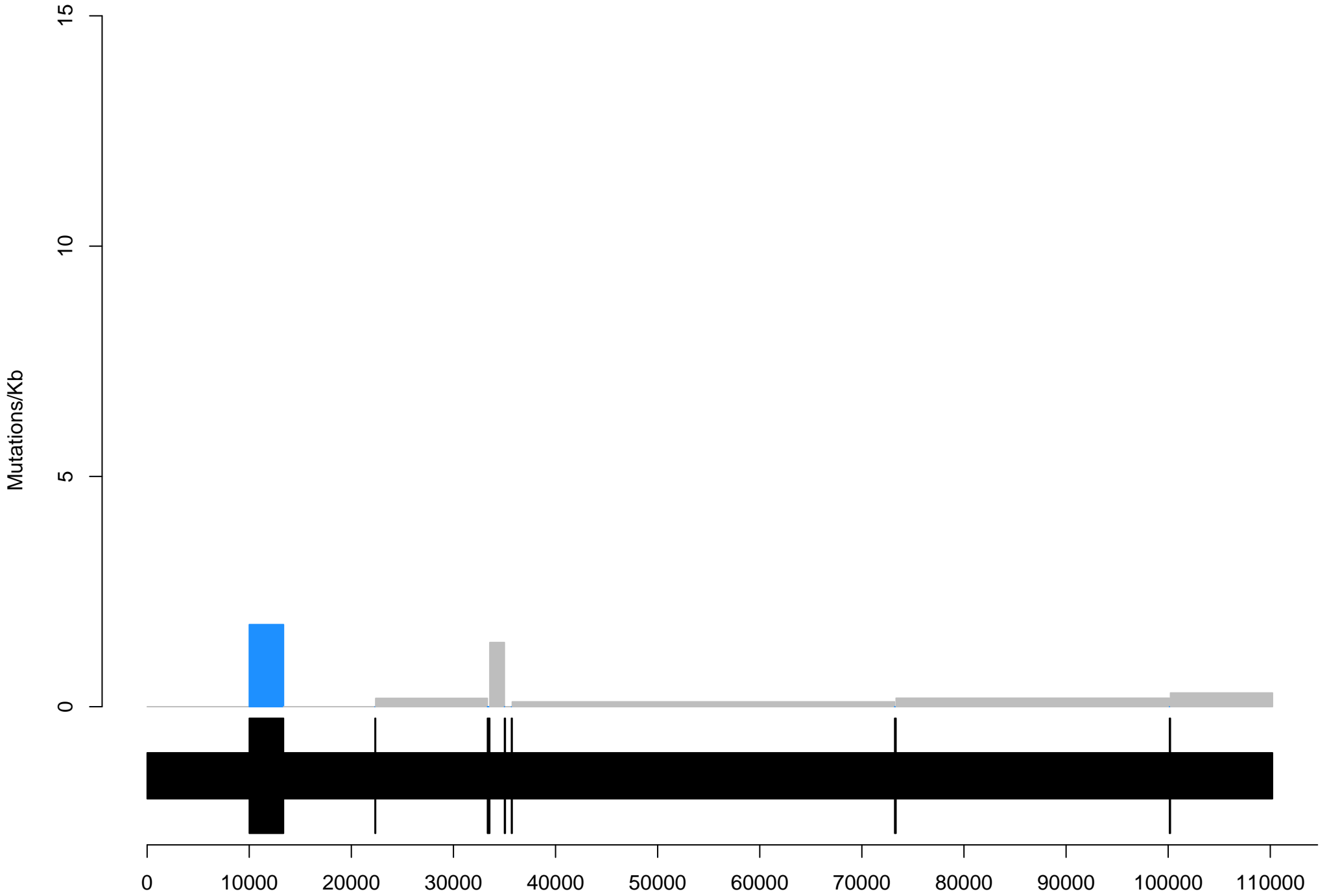
ENSG00000248202 – PRAD



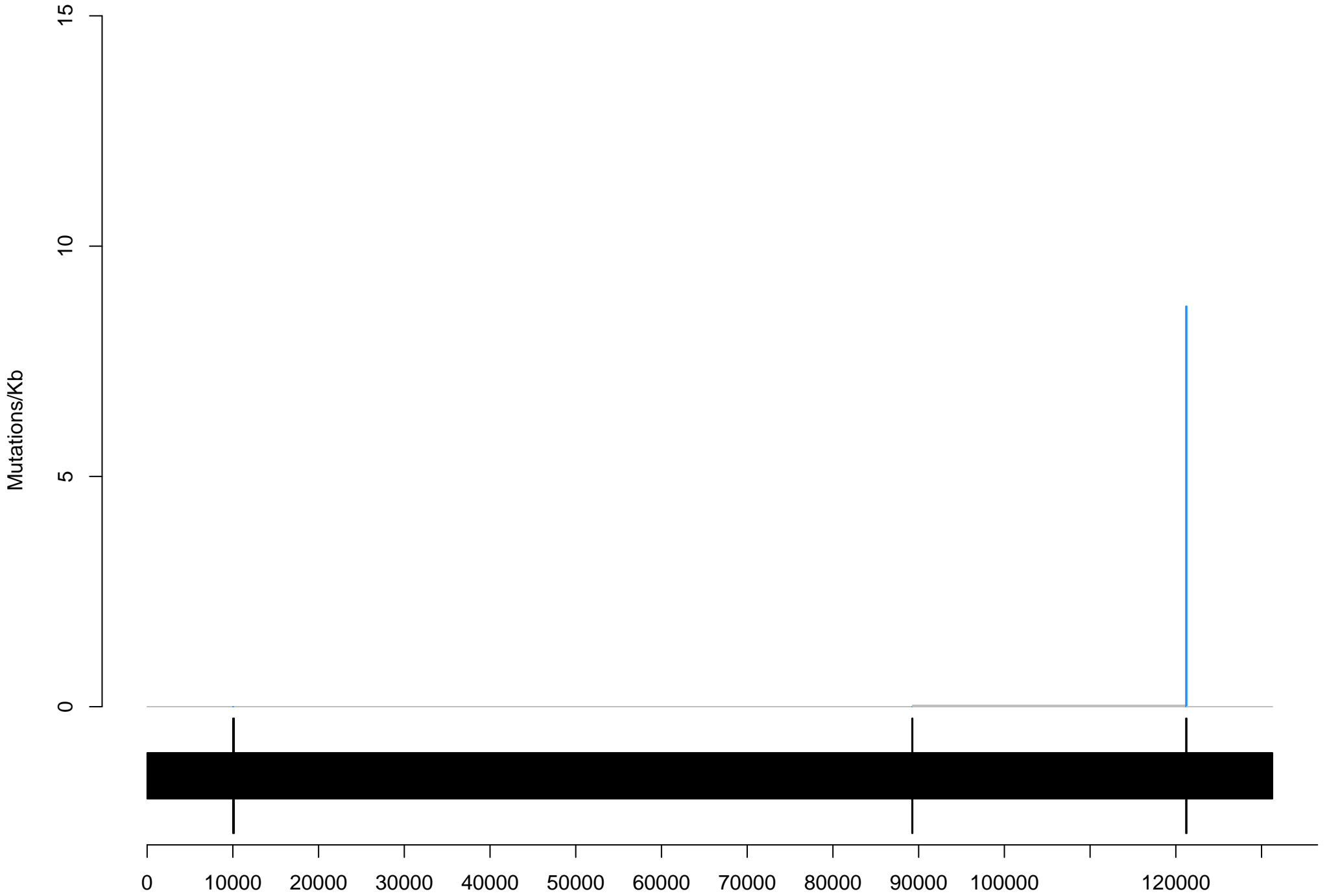
ENSG00000249734 - KICH



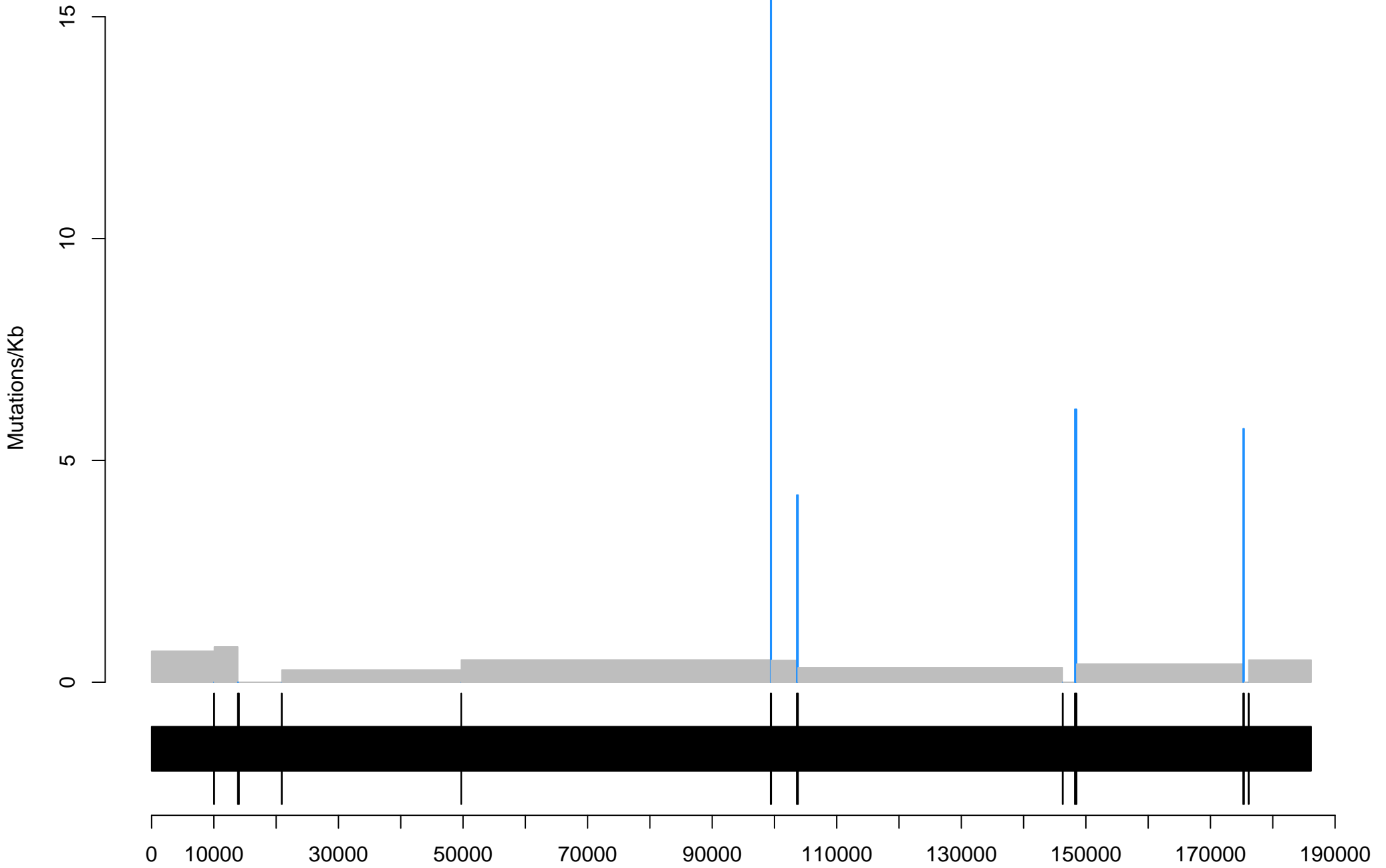
ENSG00000250125 – Breast



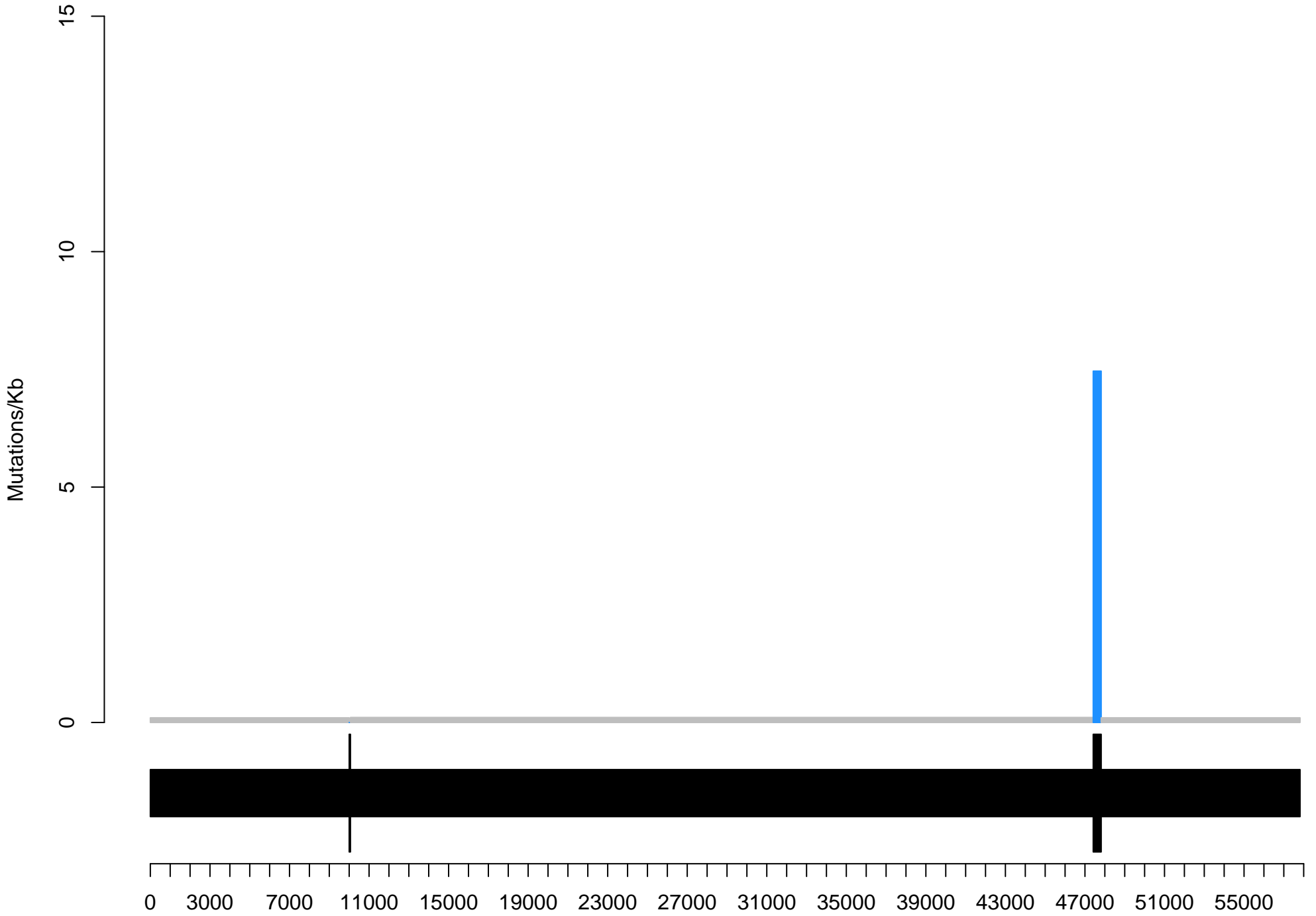
ENSG00000250488 – KICH



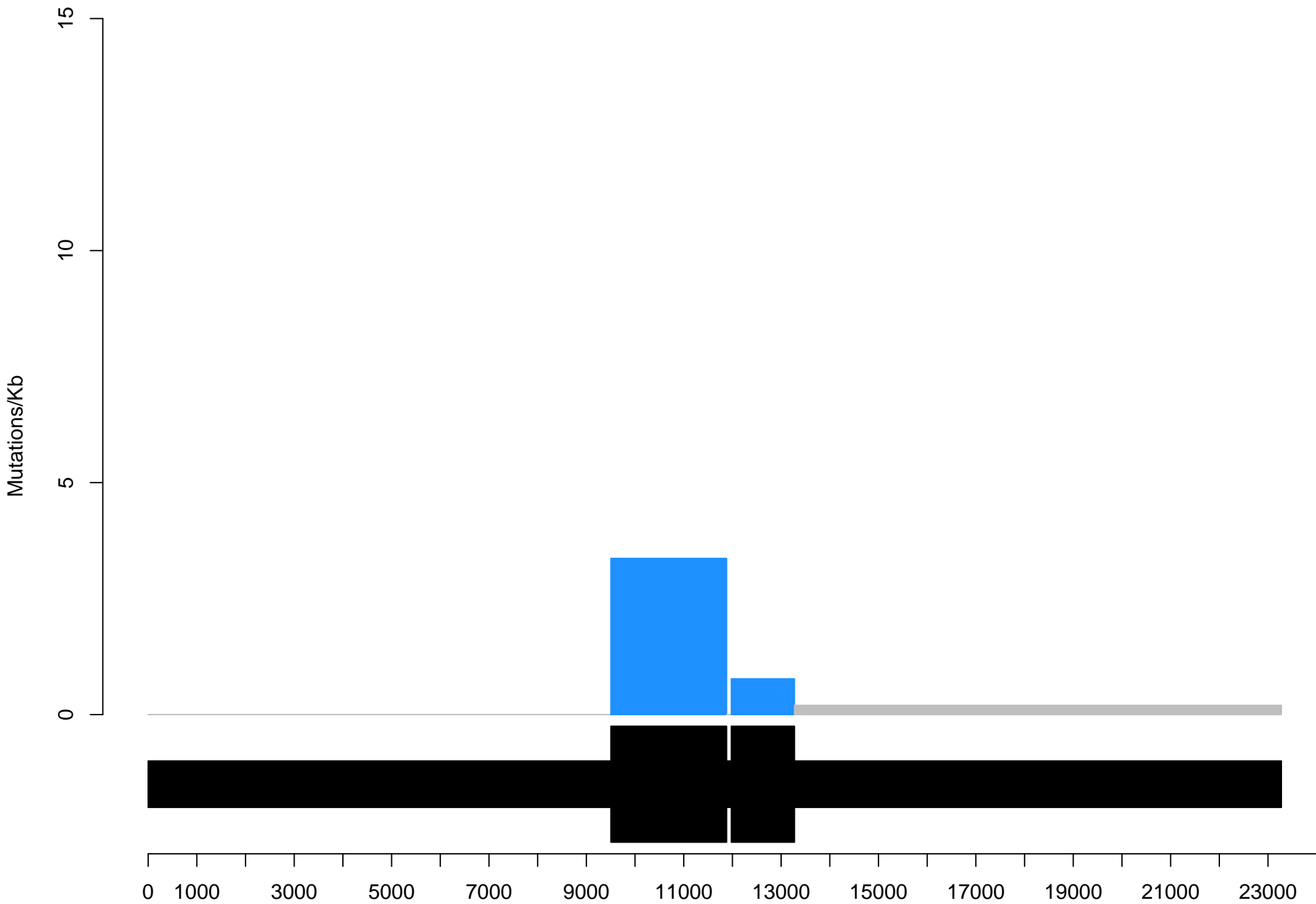
ENSG00000253434 – Breast



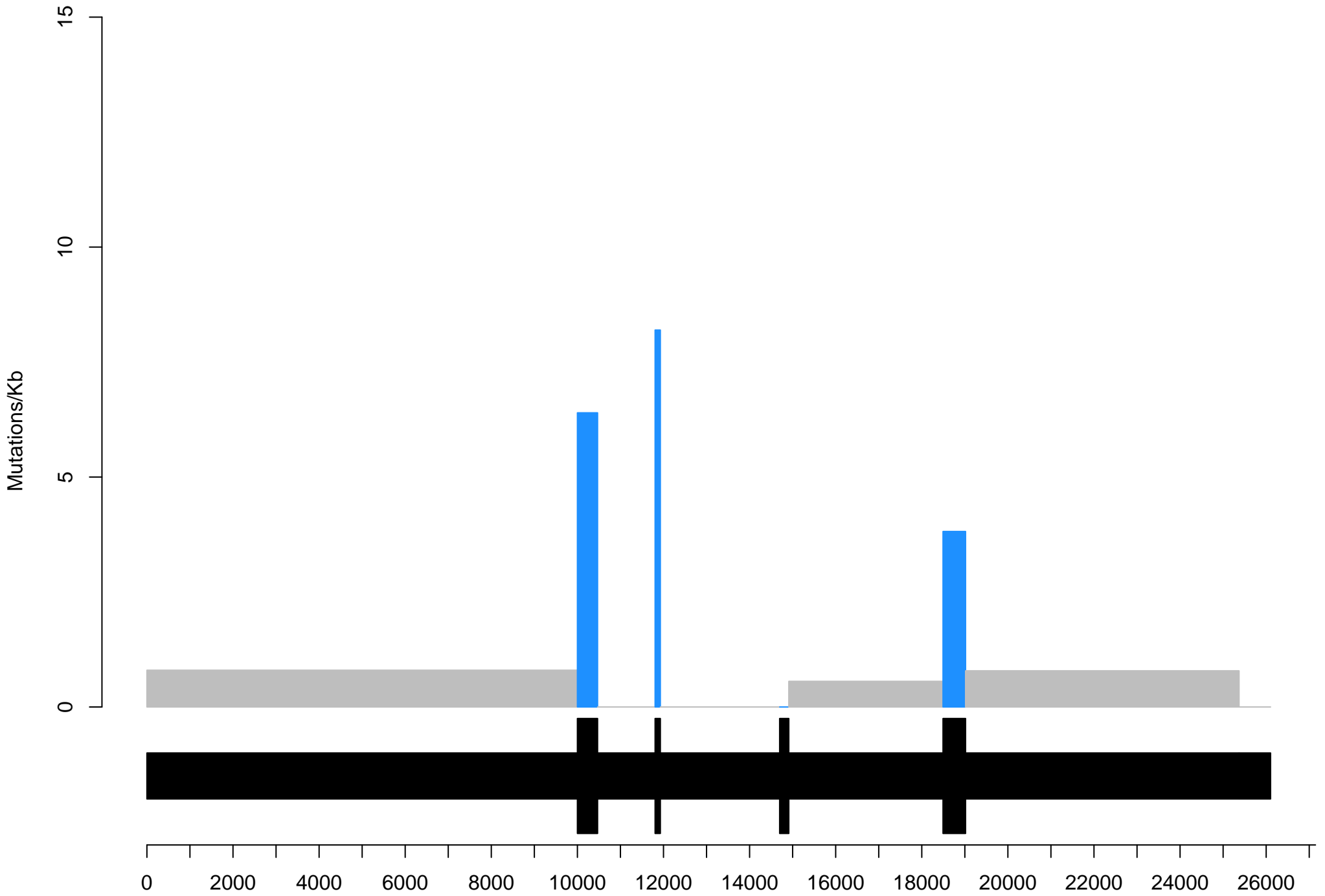
ENSG00000254689 – HNSC



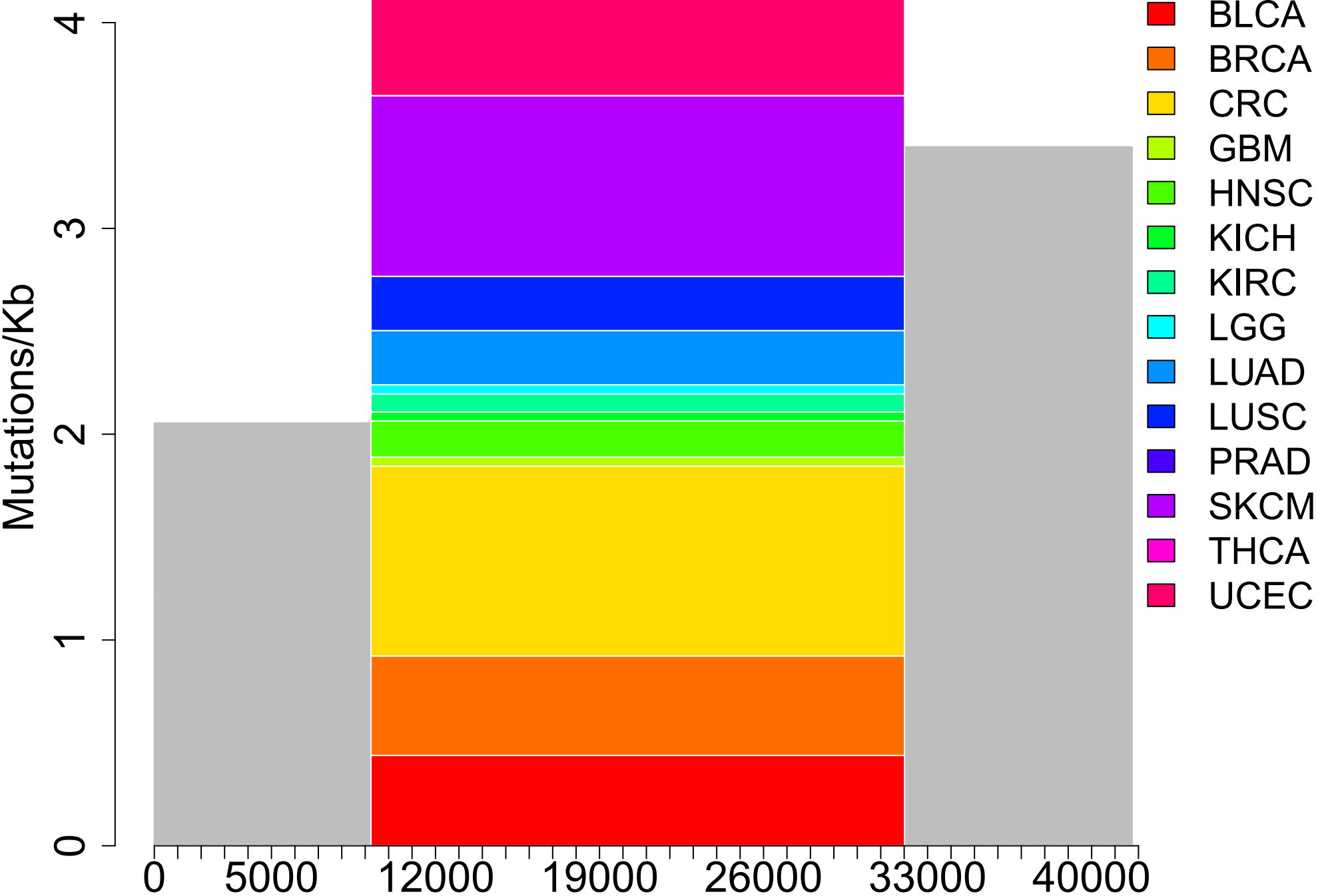
ENSG00000261623 – BRCA



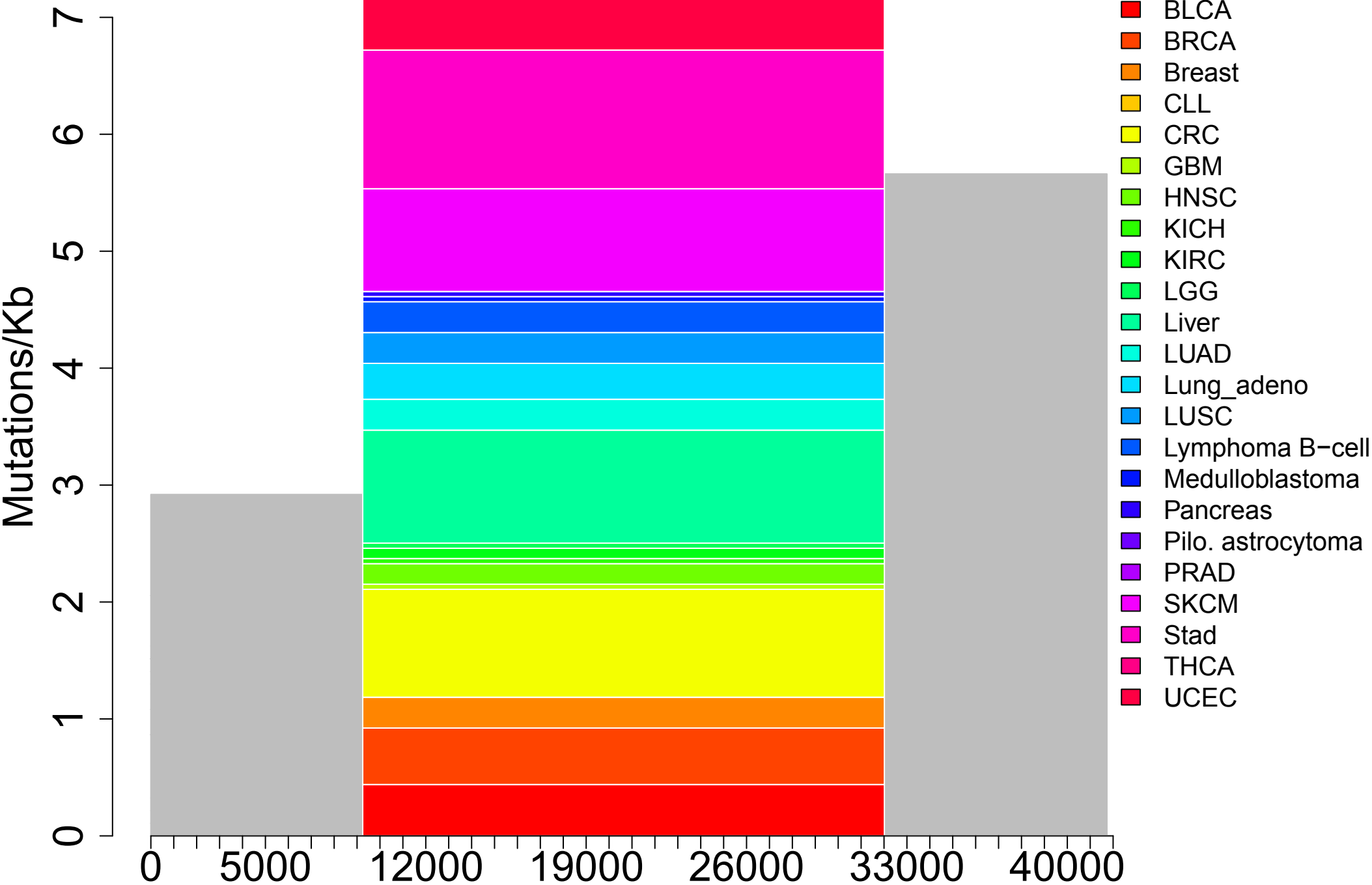
ENSG00000262117 – SKCM



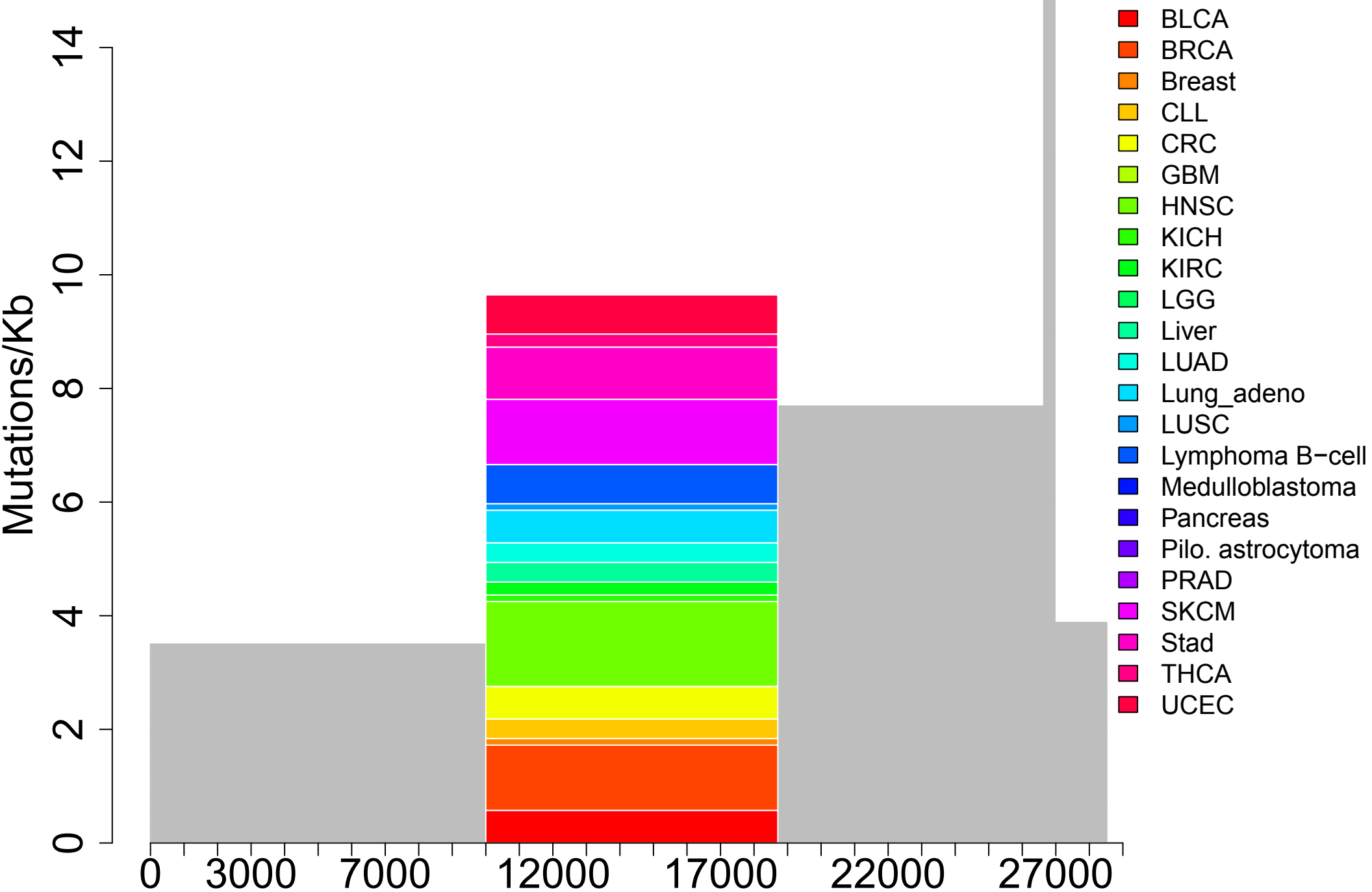
ENSG00000245532 - Pancancer TCGA



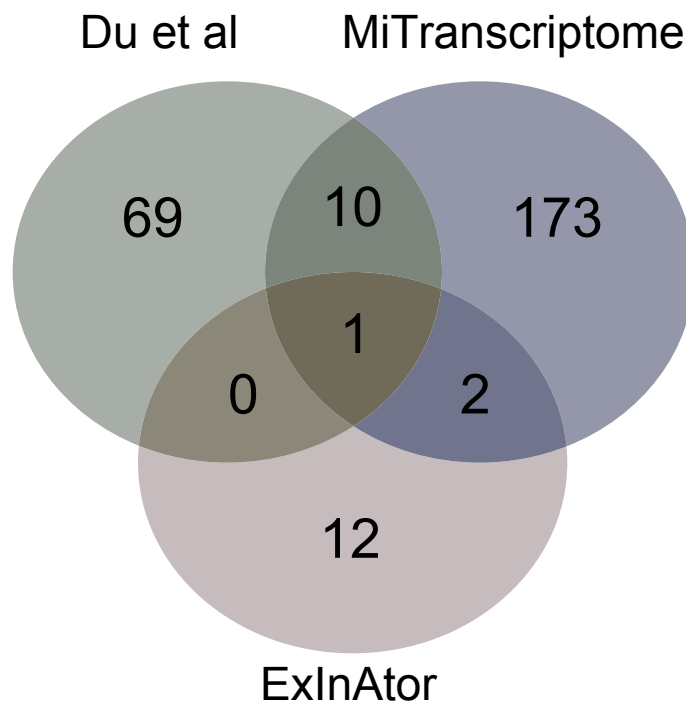
ENSG00000245532 - Superpancancer



ENSG00000251562 - Superpancancer

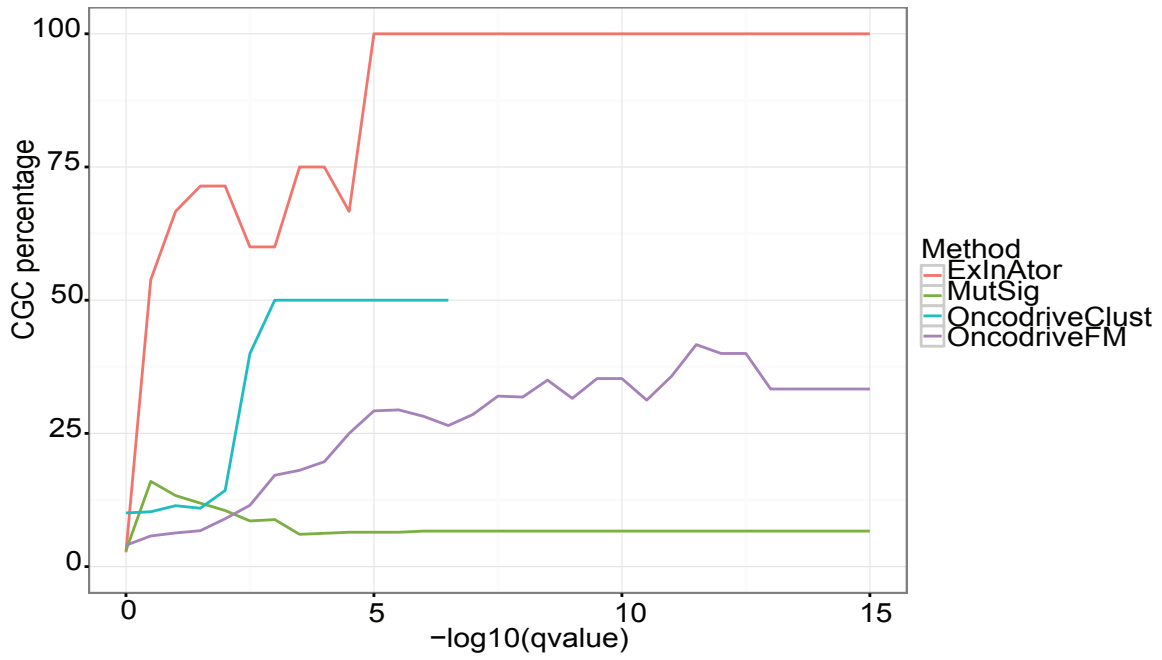


Supplementary Figure S8

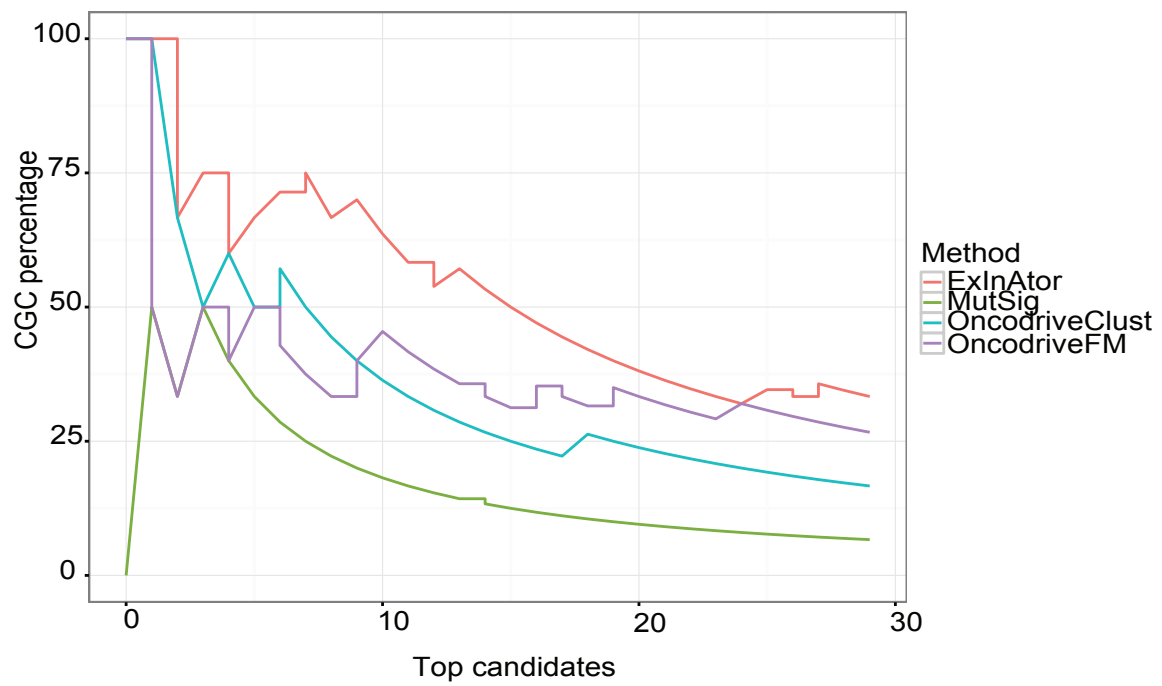


Supplementary Figure S9

A

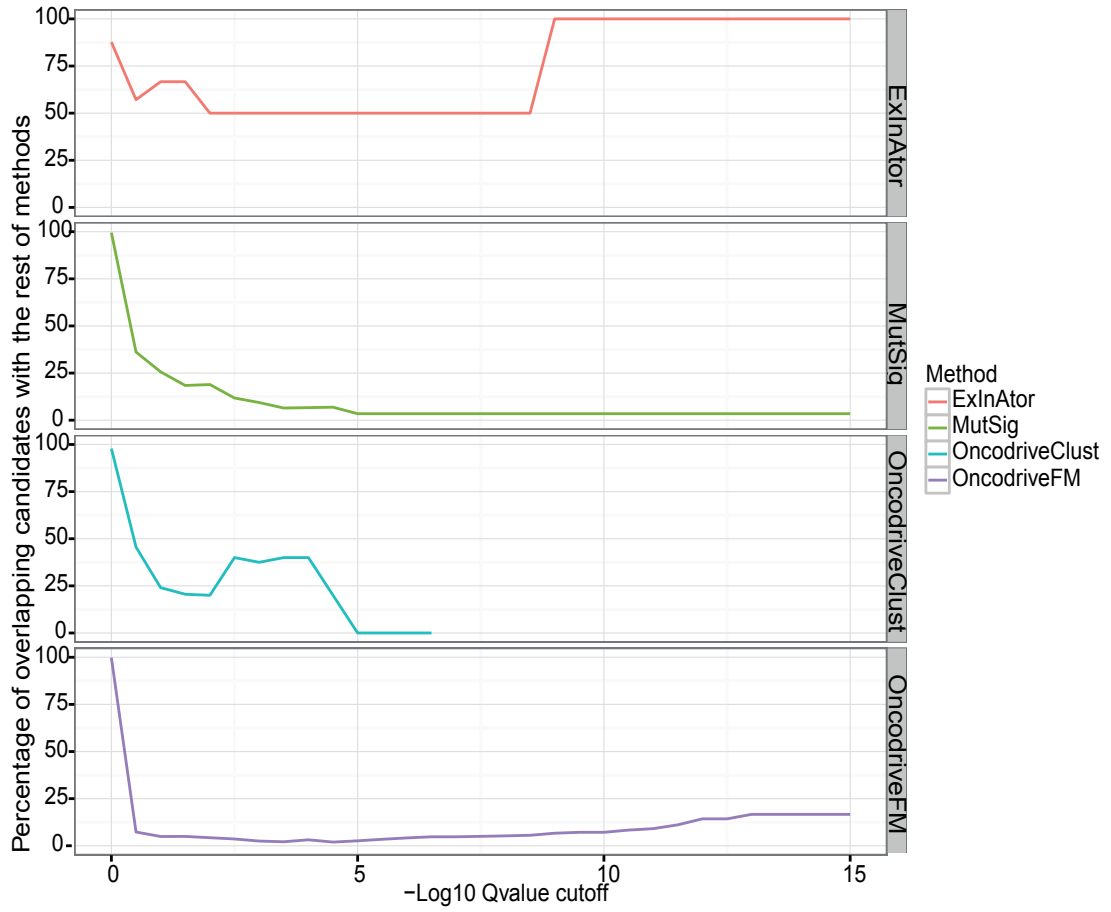


B

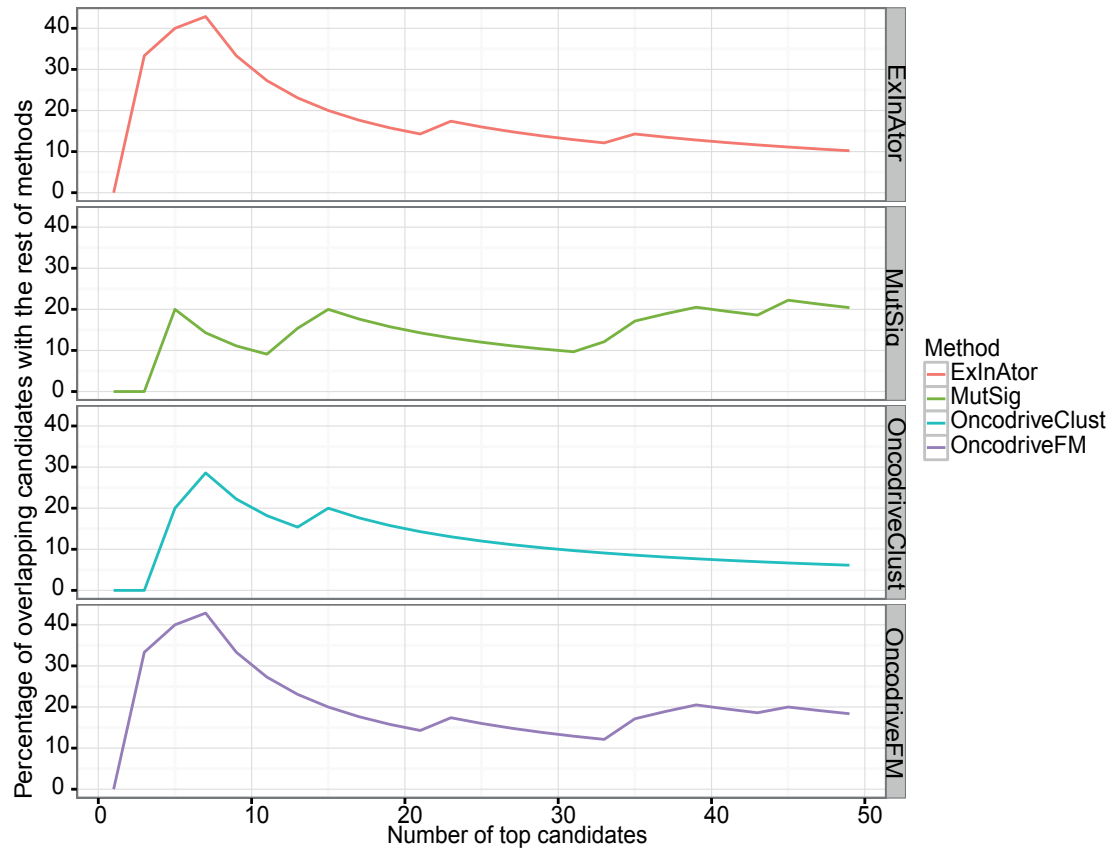


Supplementary Figure S10

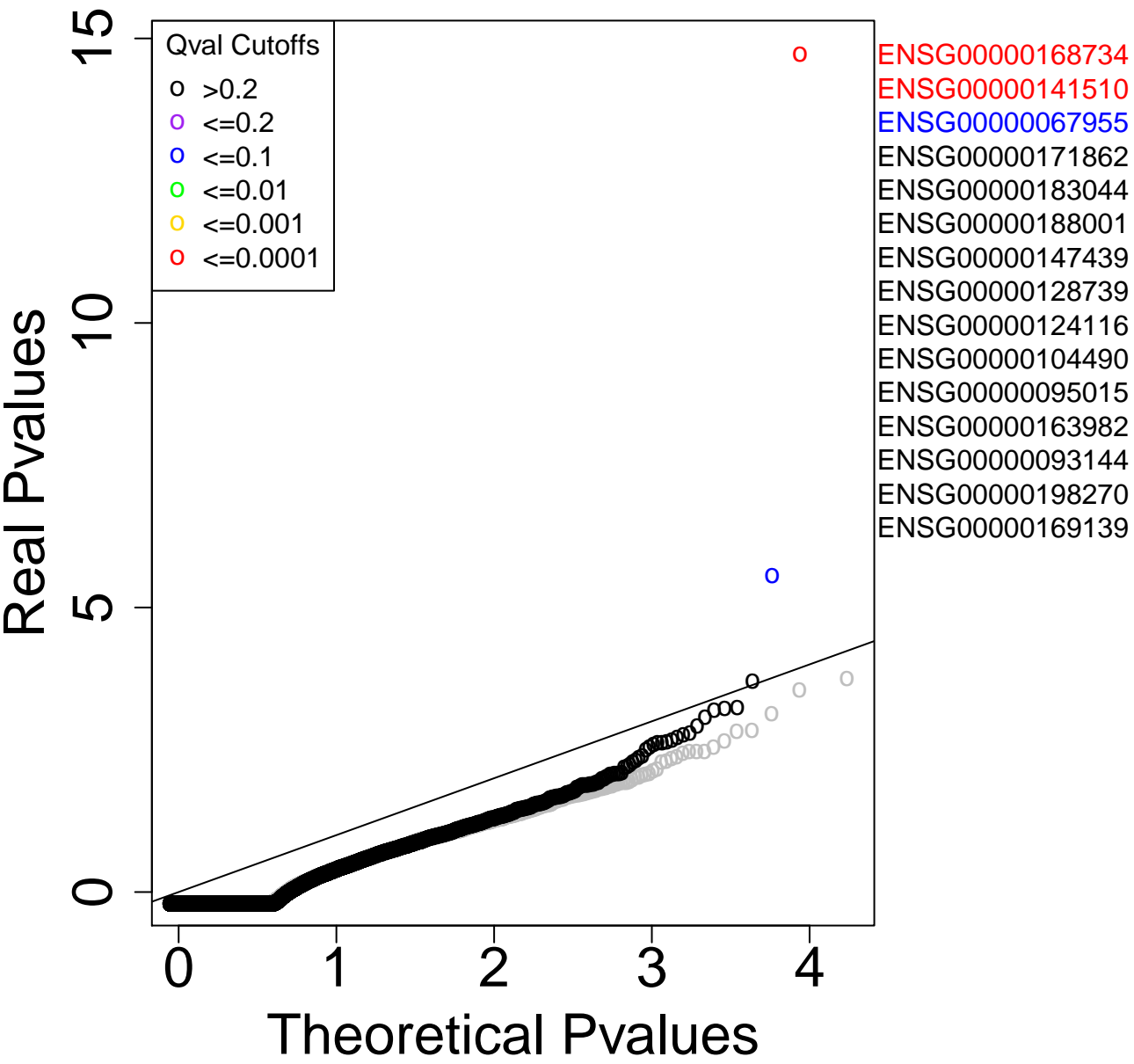
A



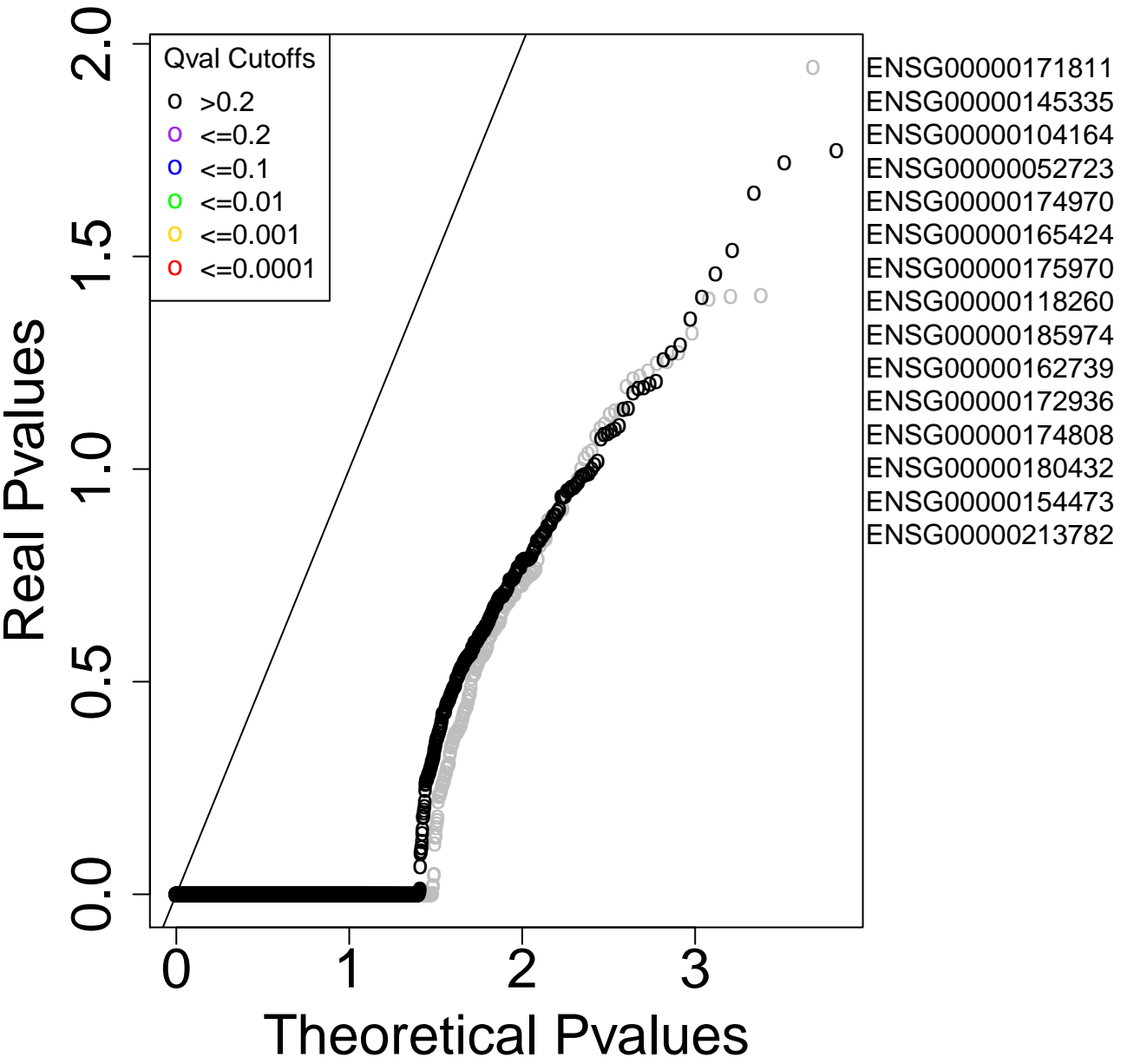
B



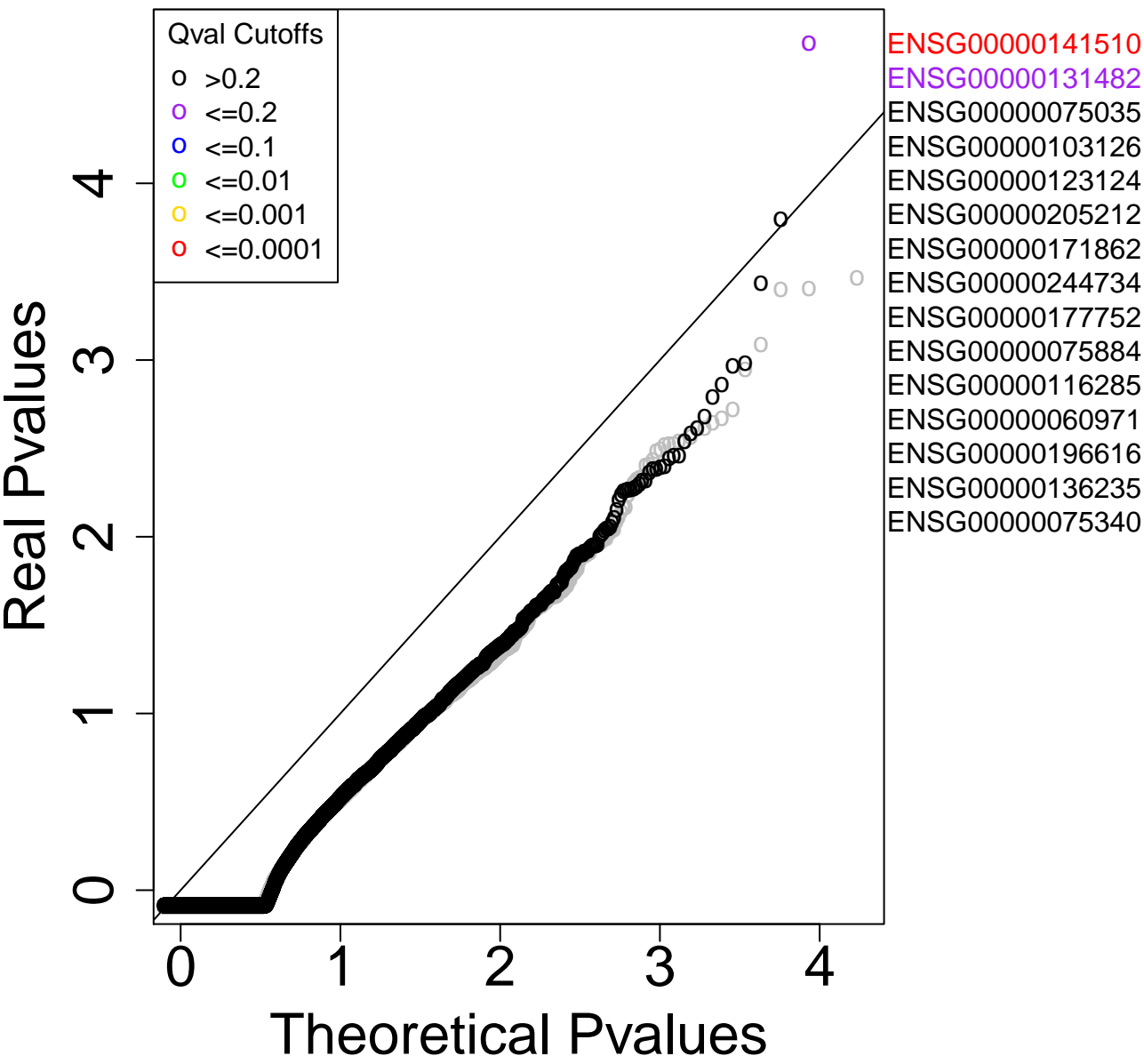
MutSig Breast - 34659 genes



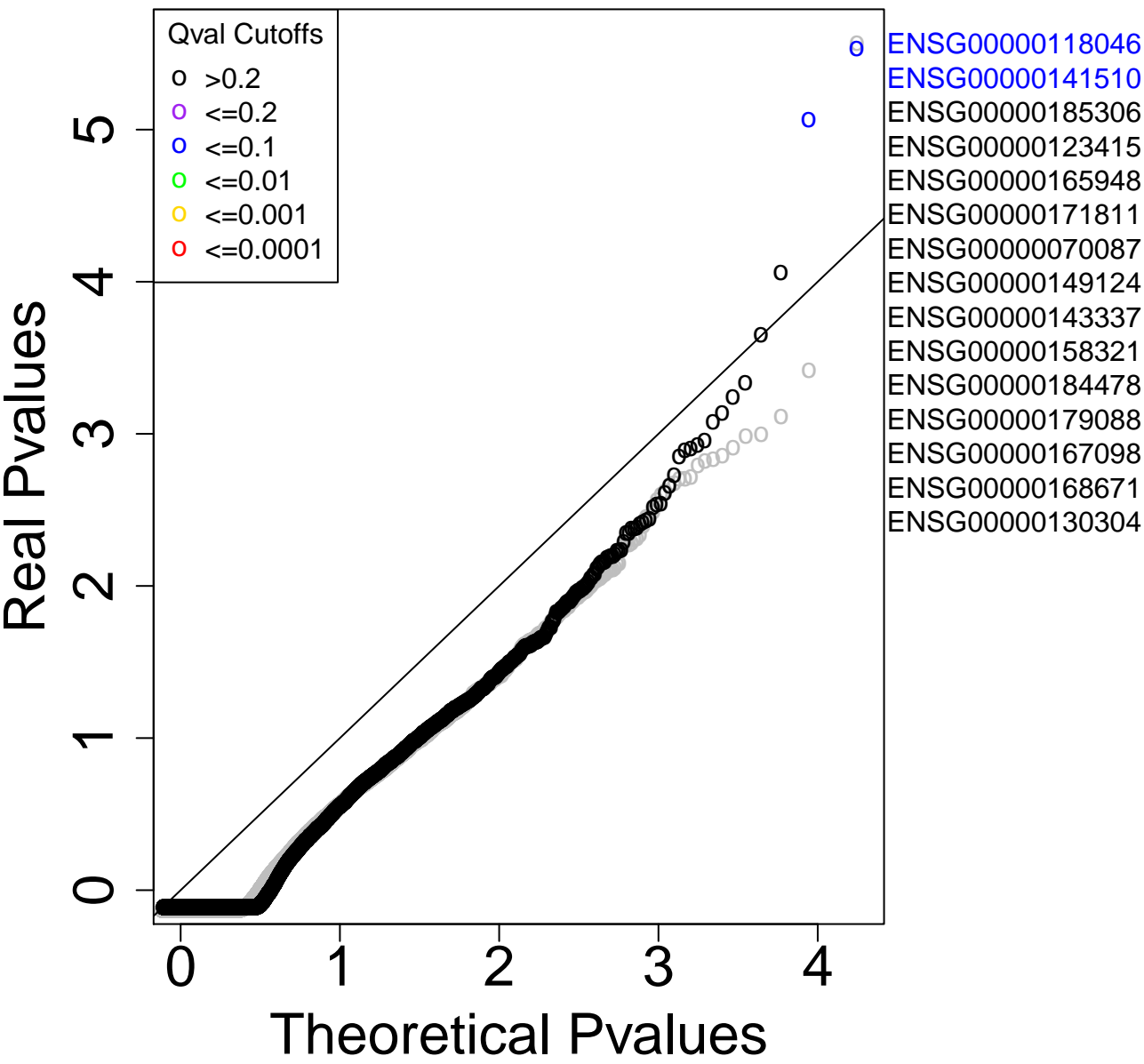
MutSig CLL - 11375 genes



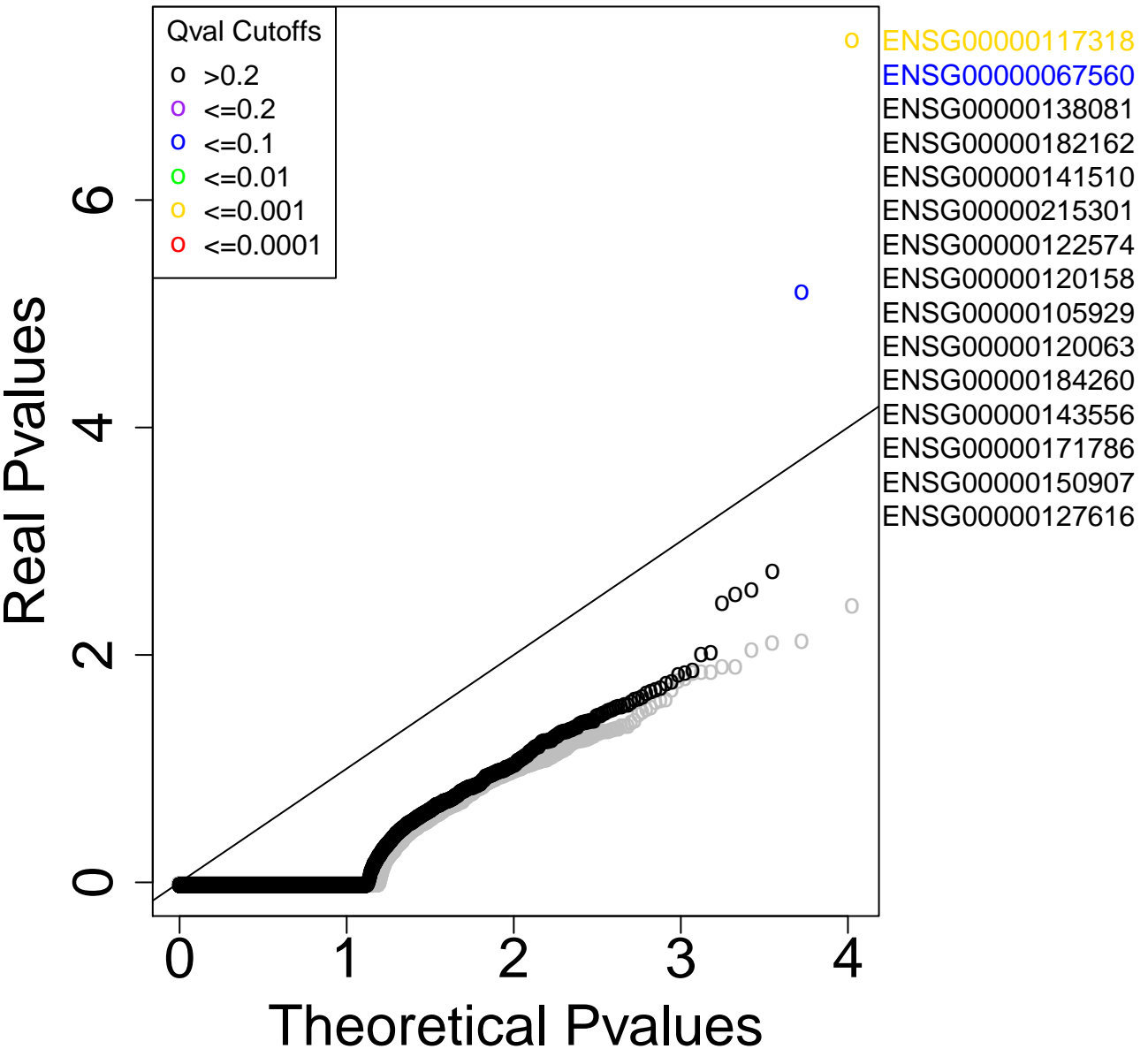
MutSig Liver – 34248 genes



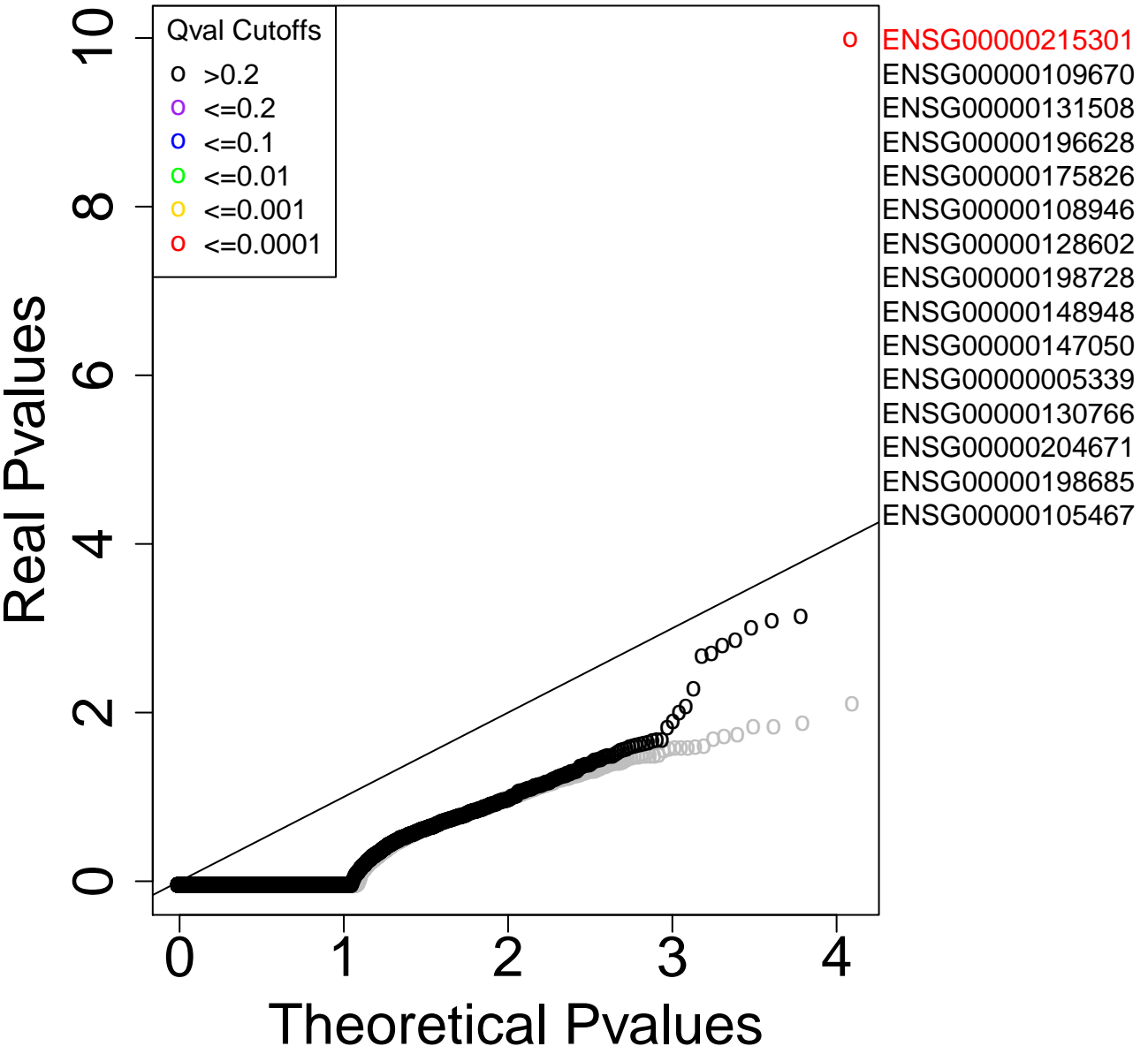
MutSig Lung adeno- 35211 genes



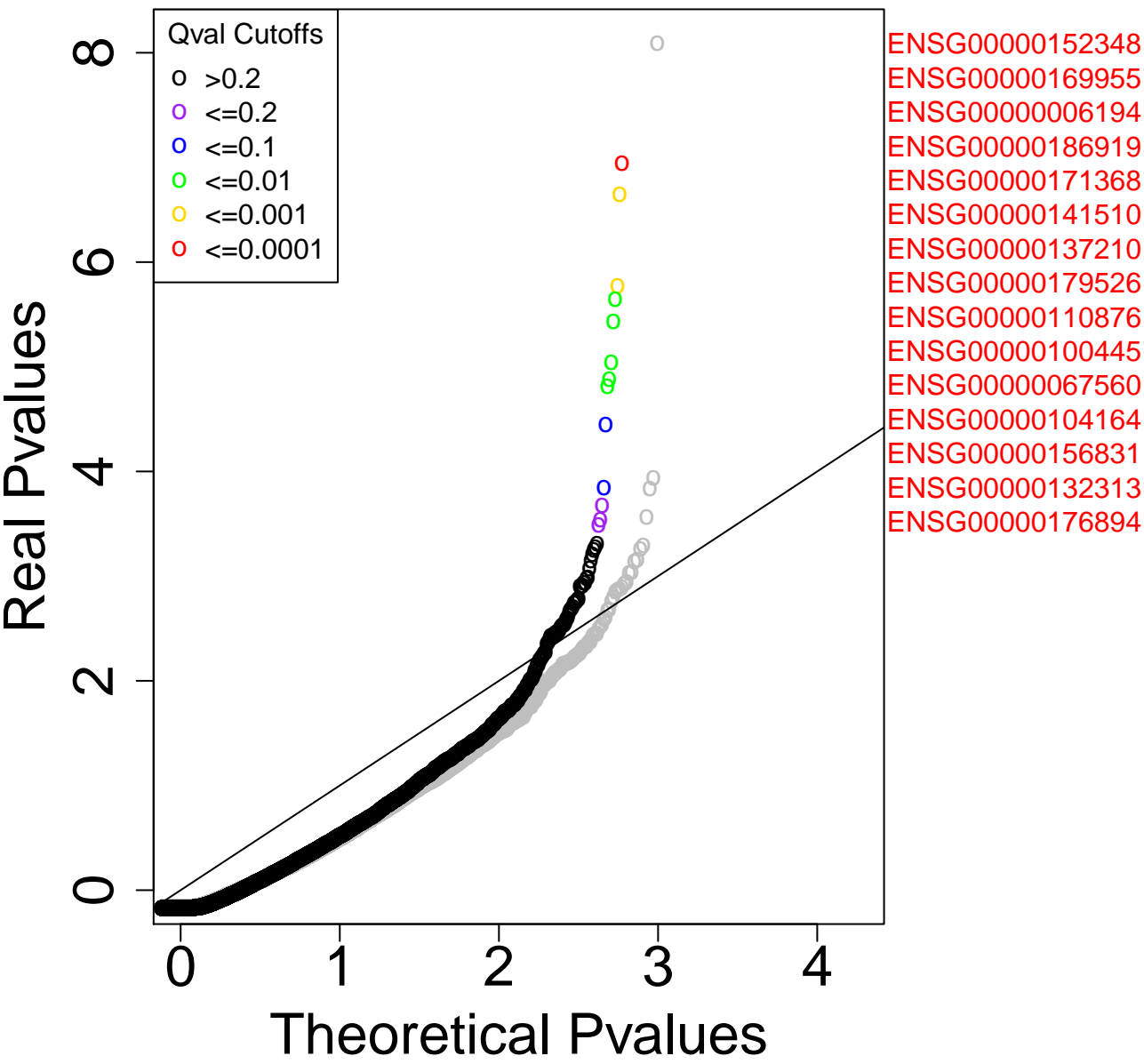
MutSig Lymphoma B-cell – 21158 genes



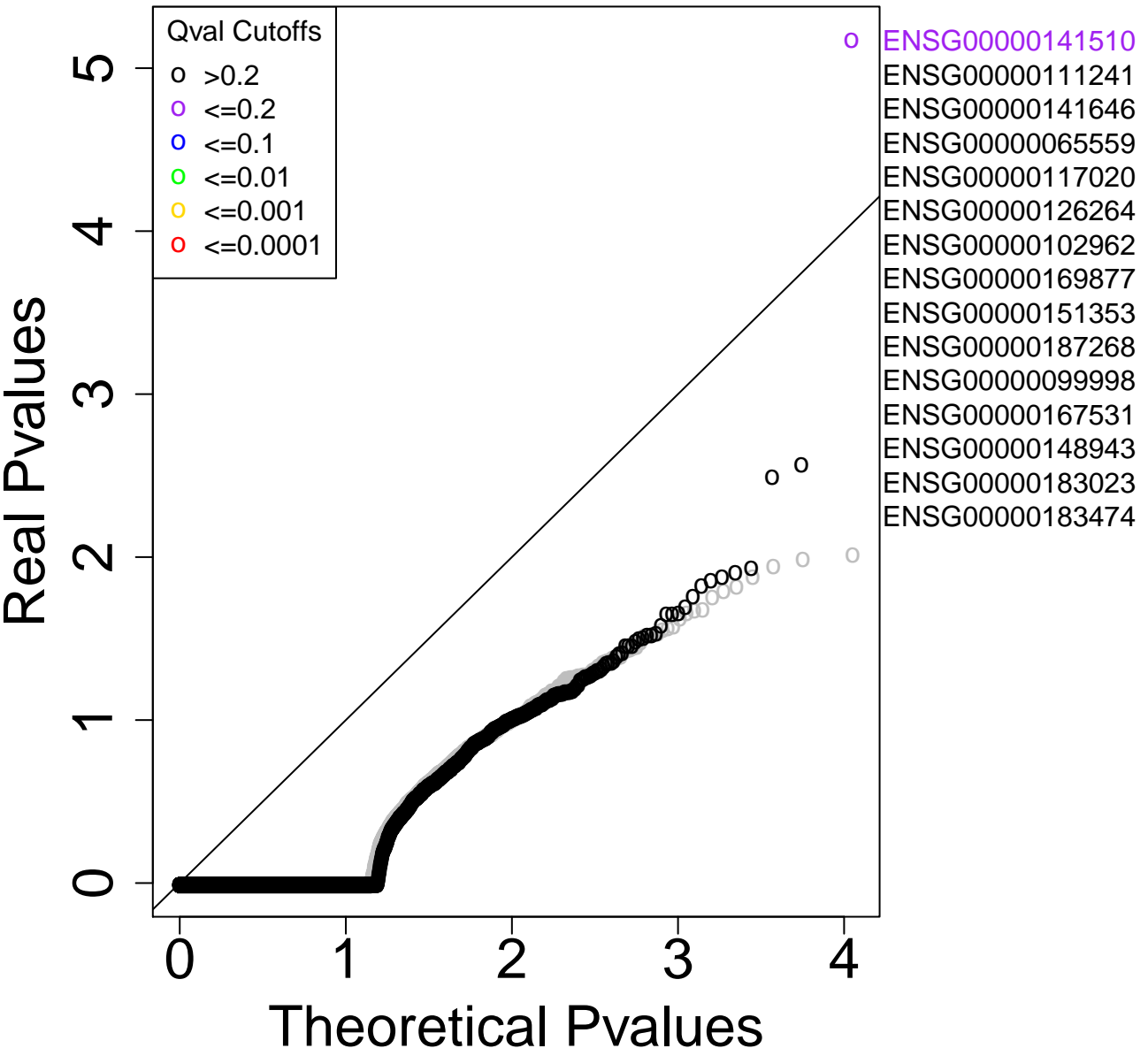
MutSig Medulloblastoma – 24518 genes



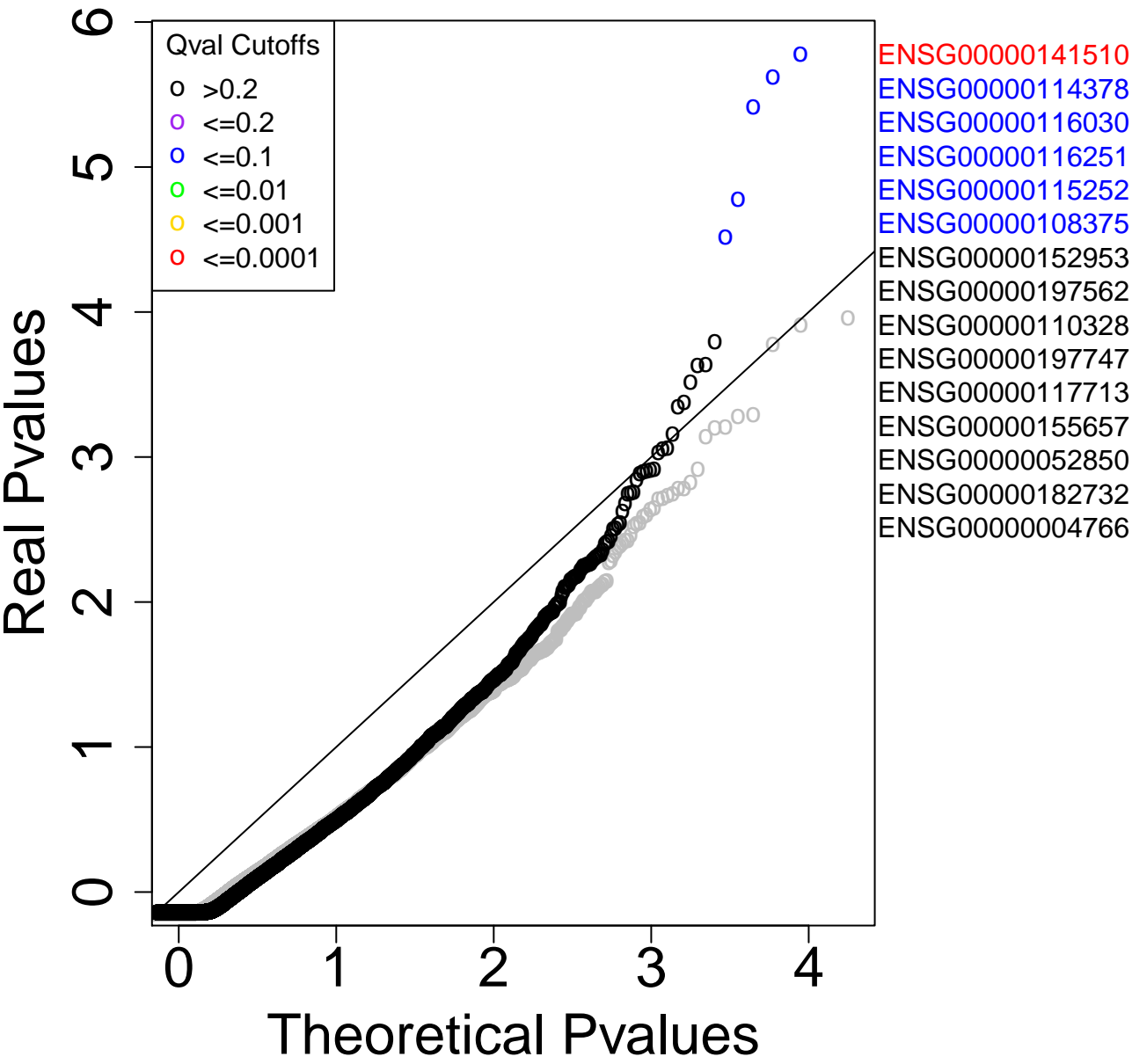
MutSig Pancancer – 35544 genes



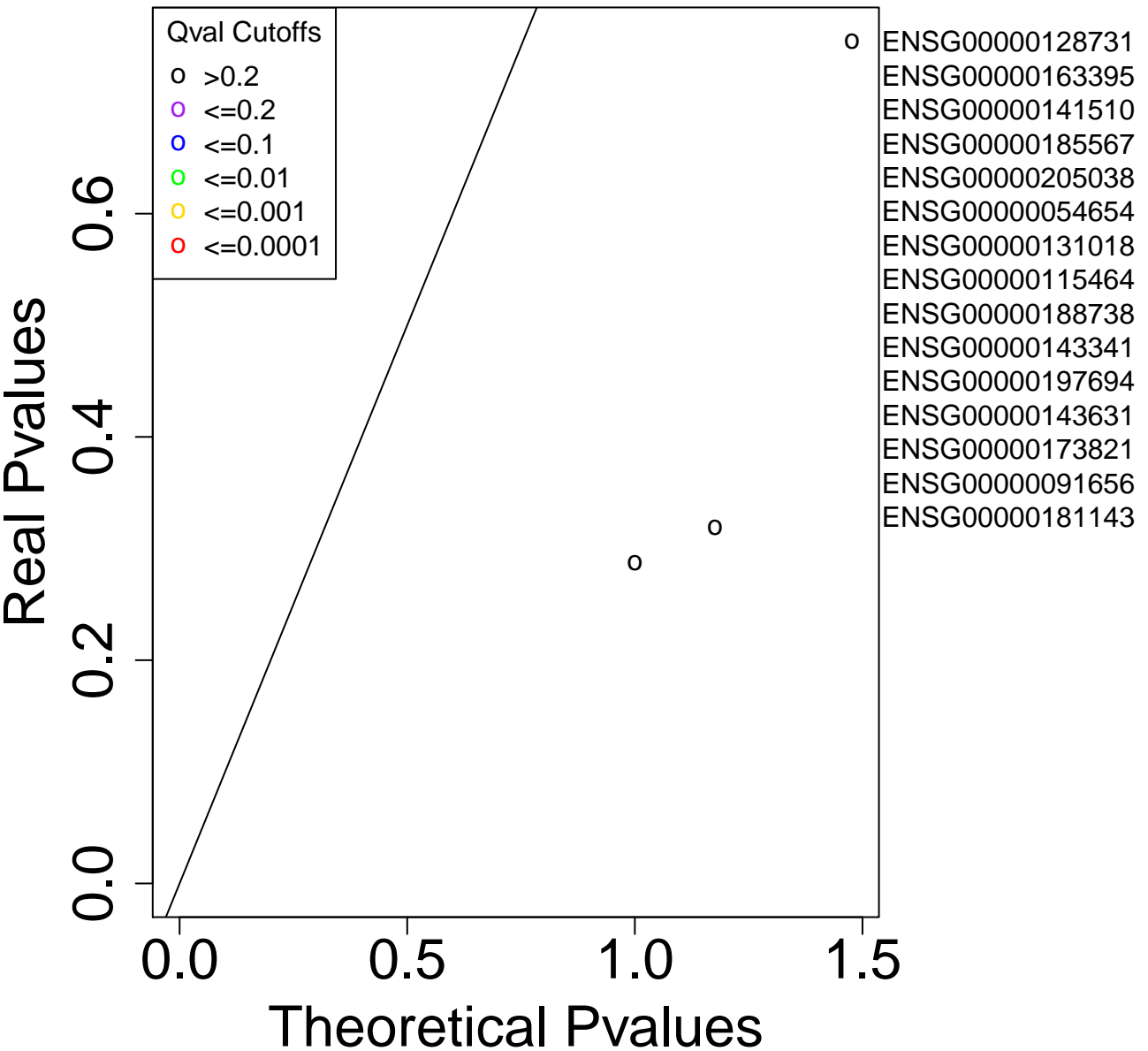
MutSig Pancreas – 22319 genes



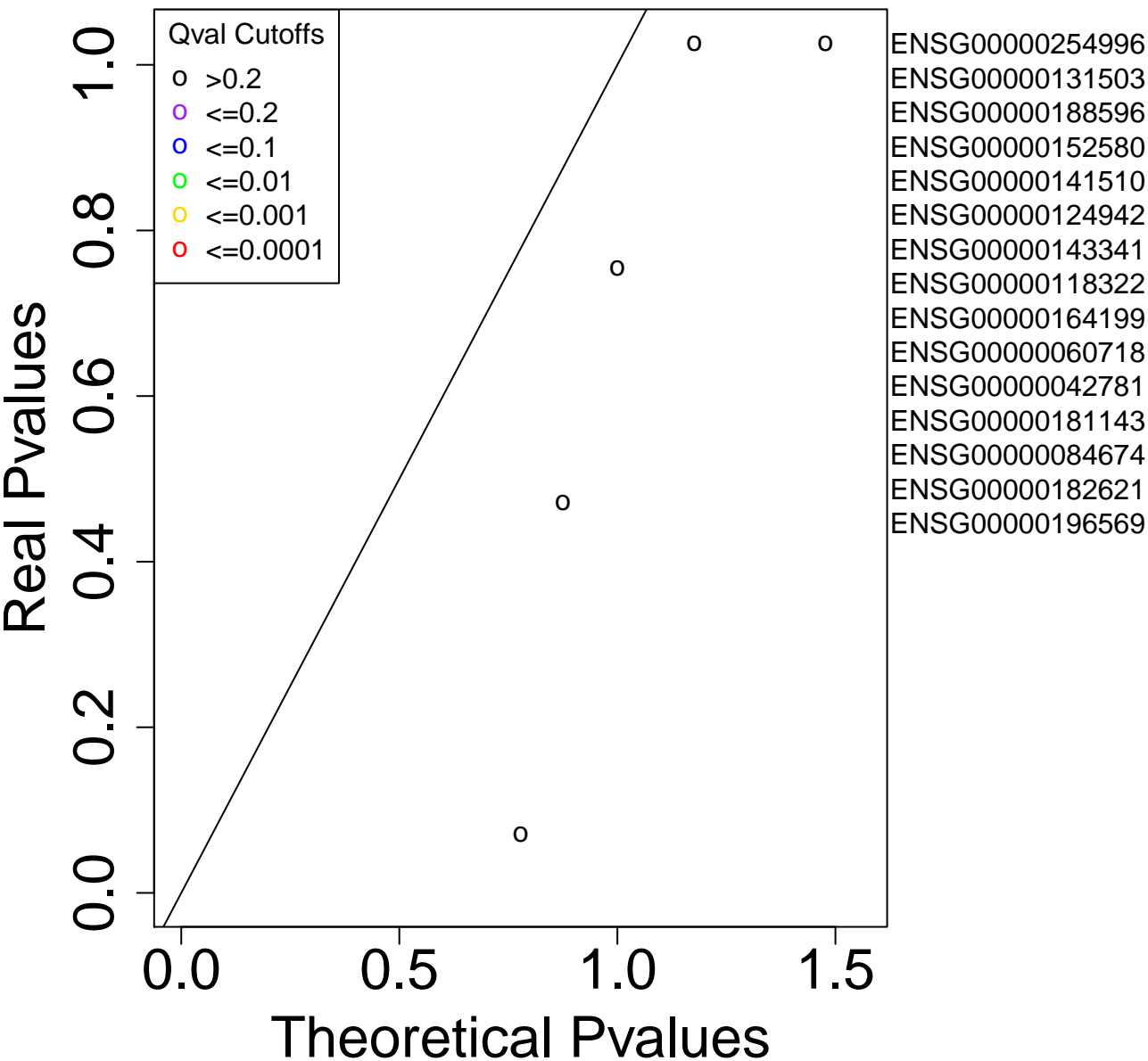
MutSig Stad – 35513 genes



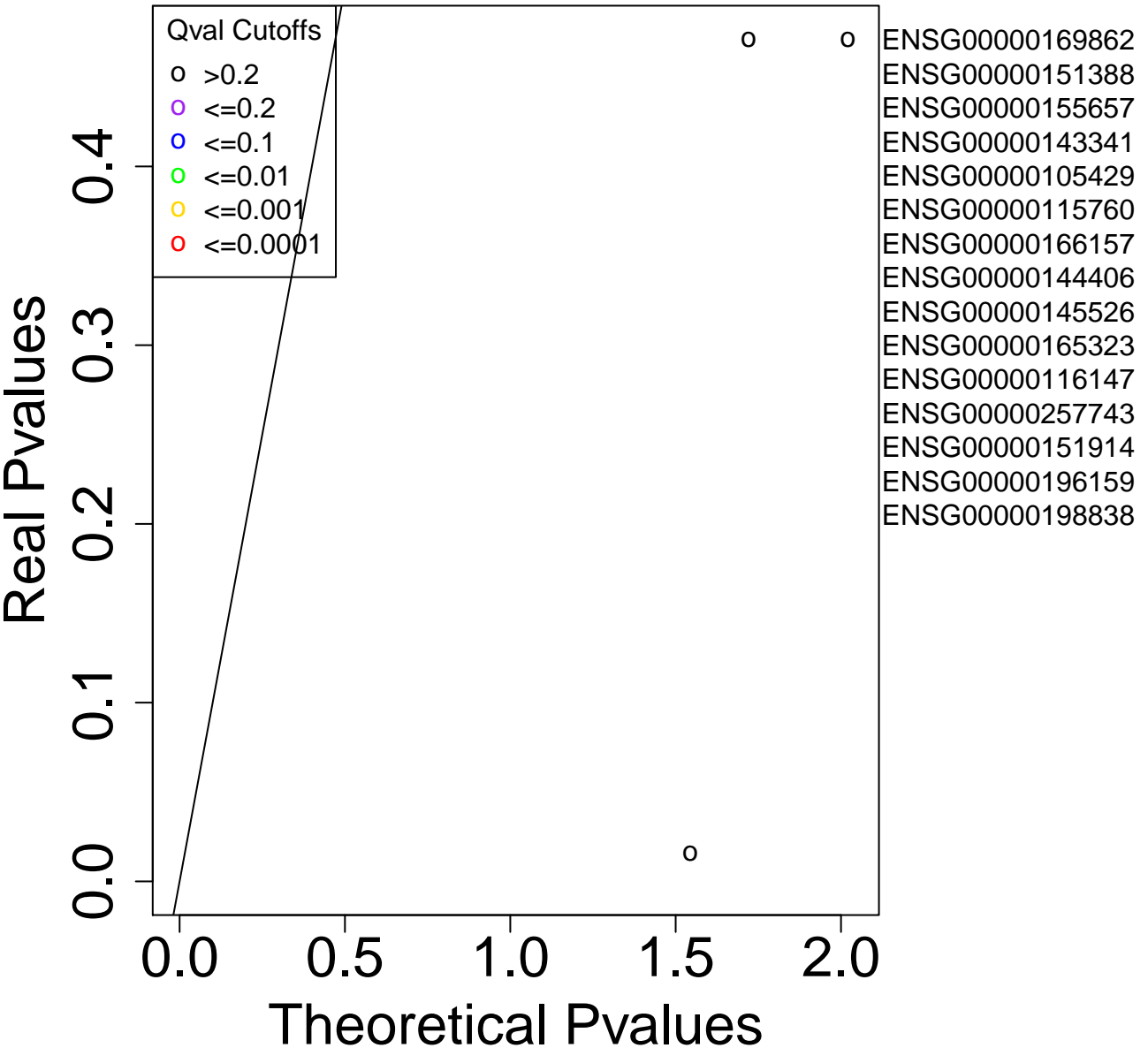
OncodriveClust Breast – 47 genes



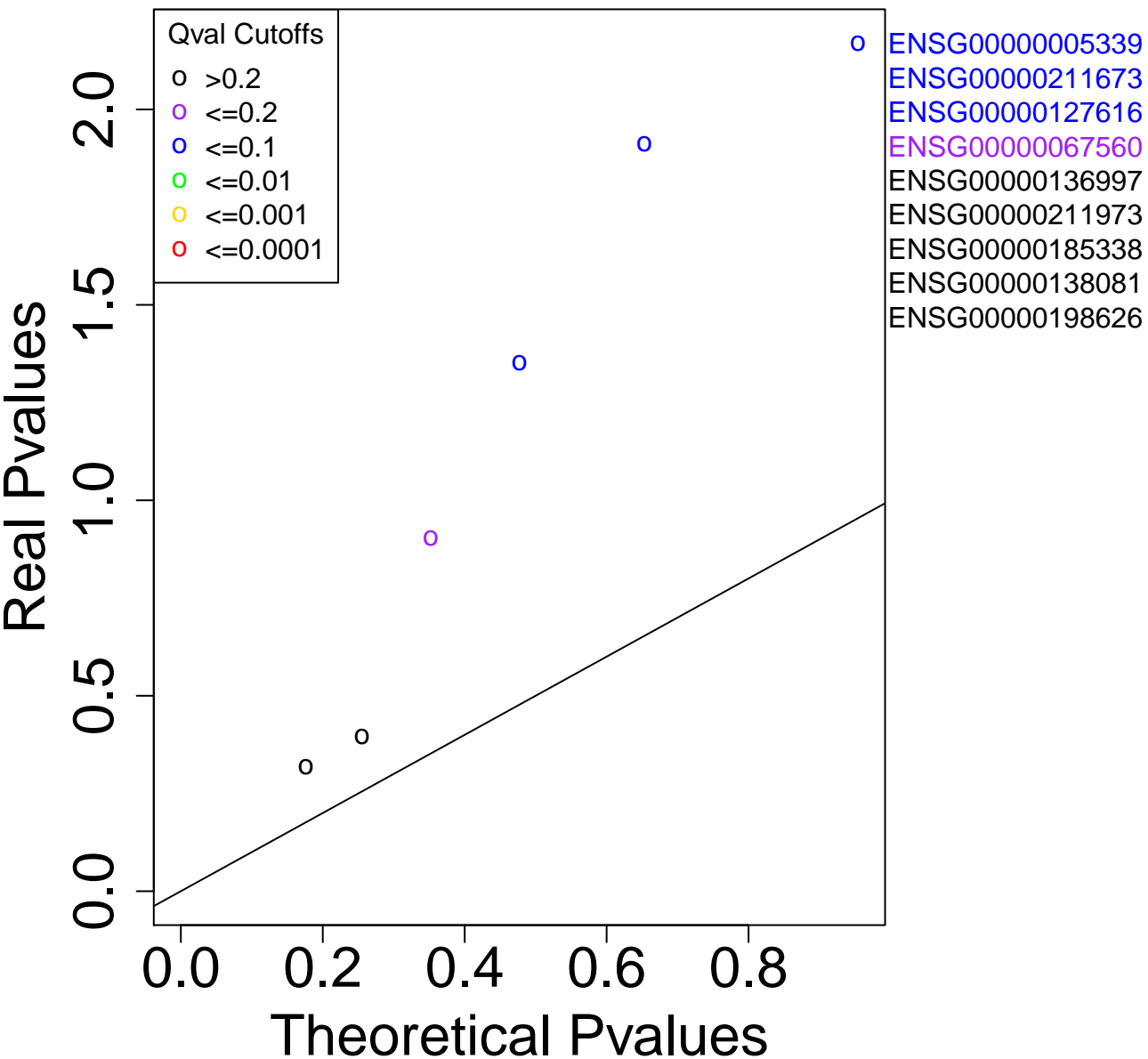
OncodriveClust Liver – 66 genes



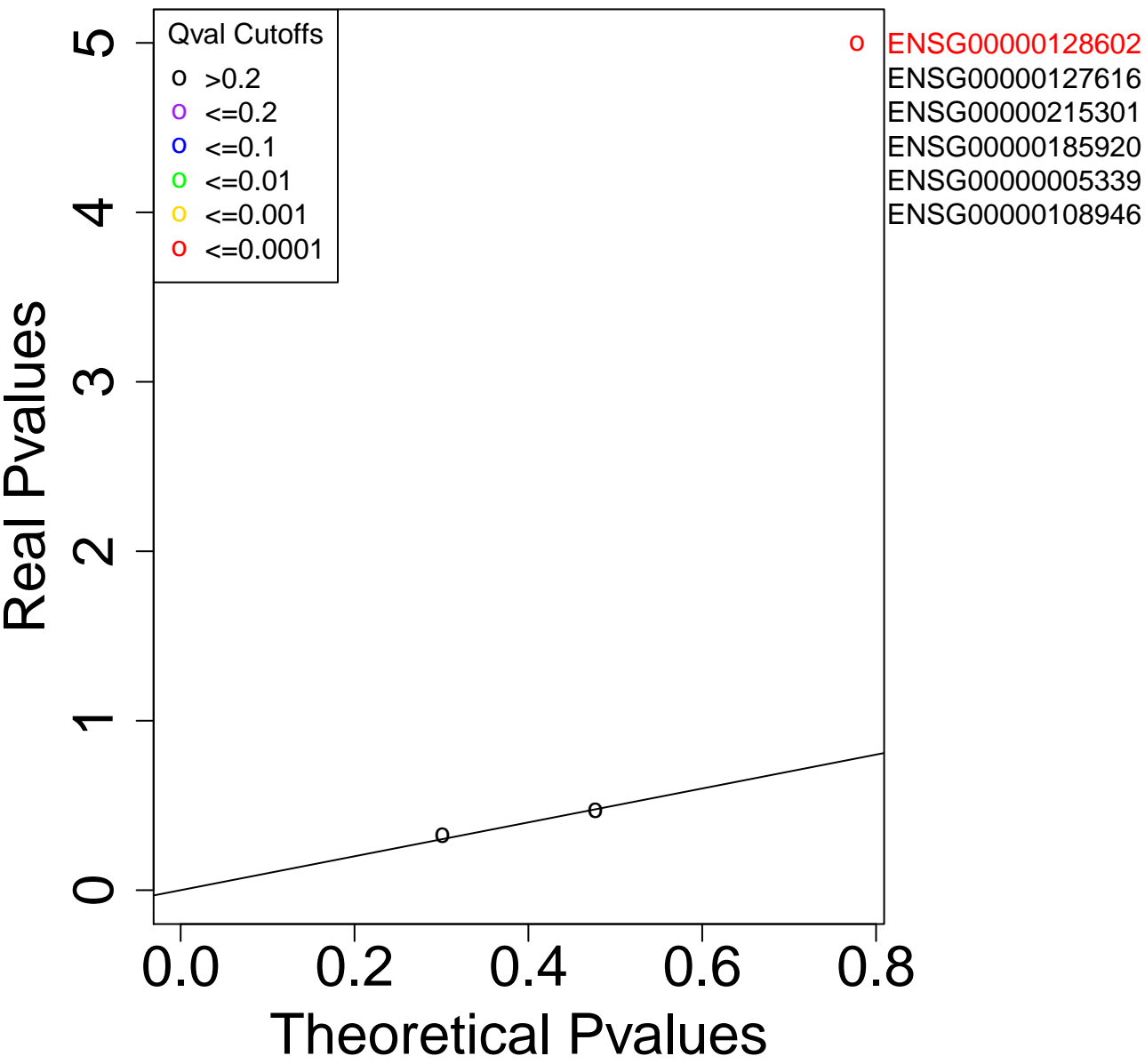
OncodriveClust Lung adeno – 213 genes



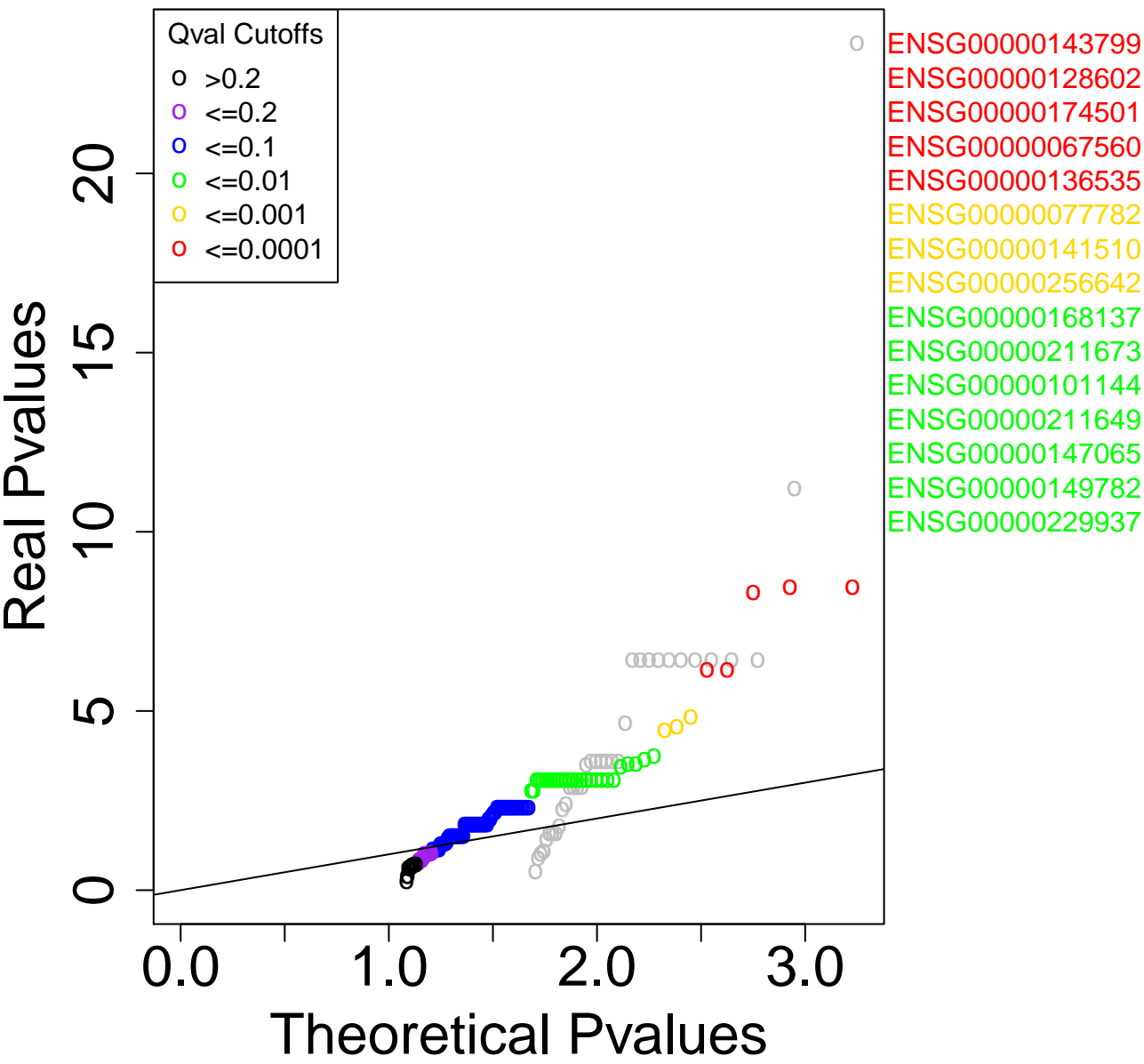
OncodriveClust Lymphoma B-cell – 11 genes



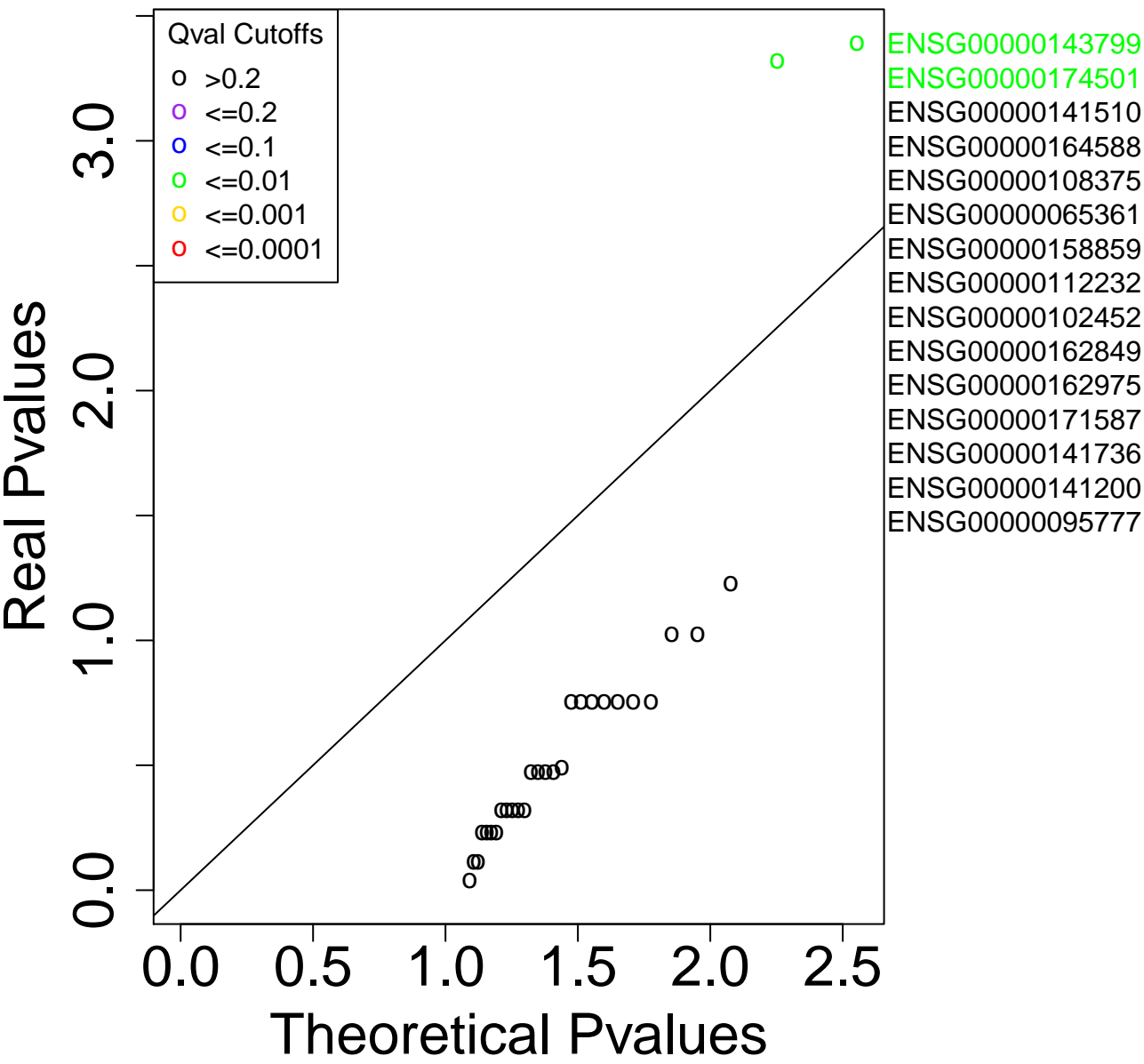
OncodriveClust Medulloblastoma – 7 genes



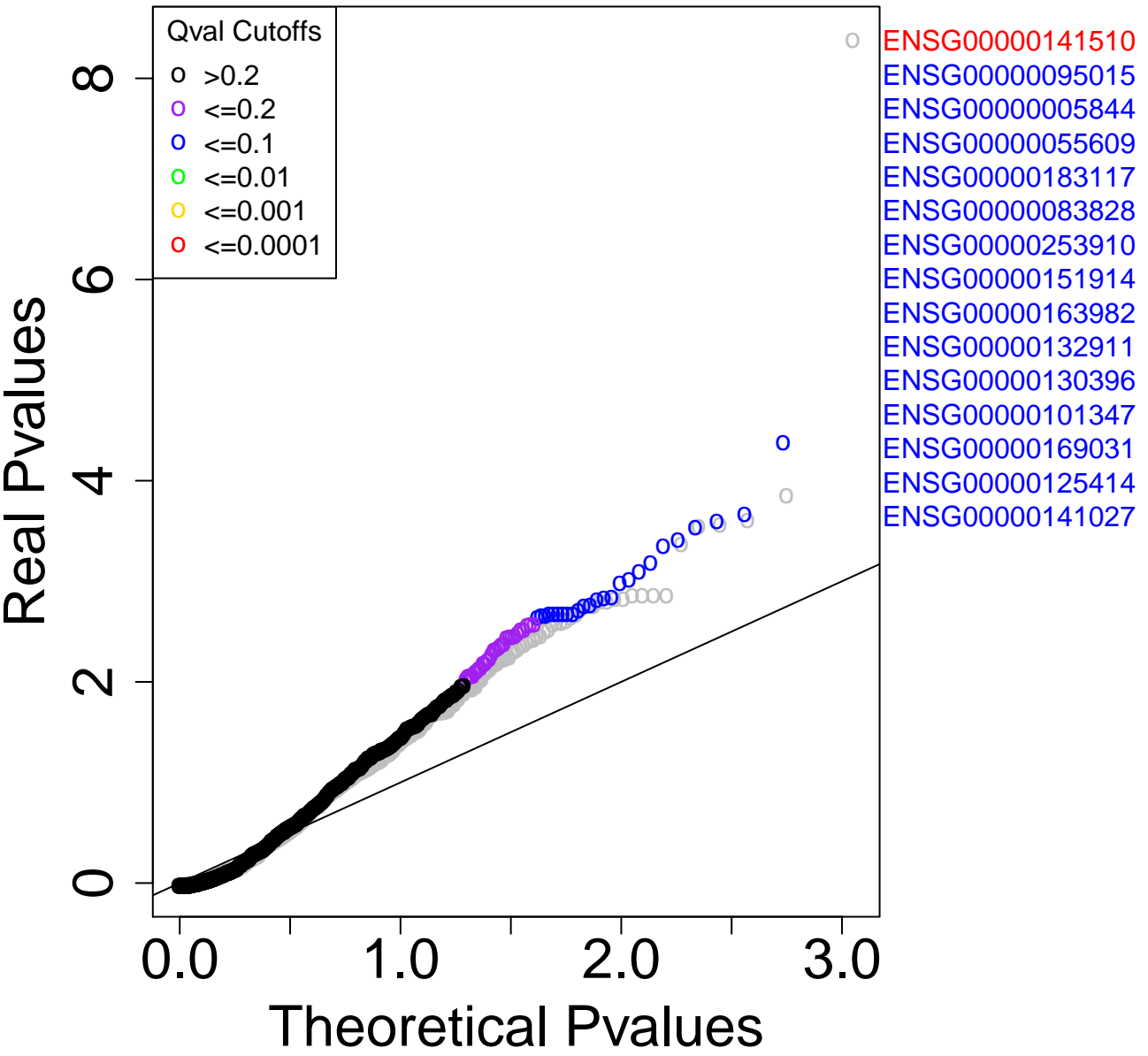
OncodriveClust Pancancer – 3464 genes



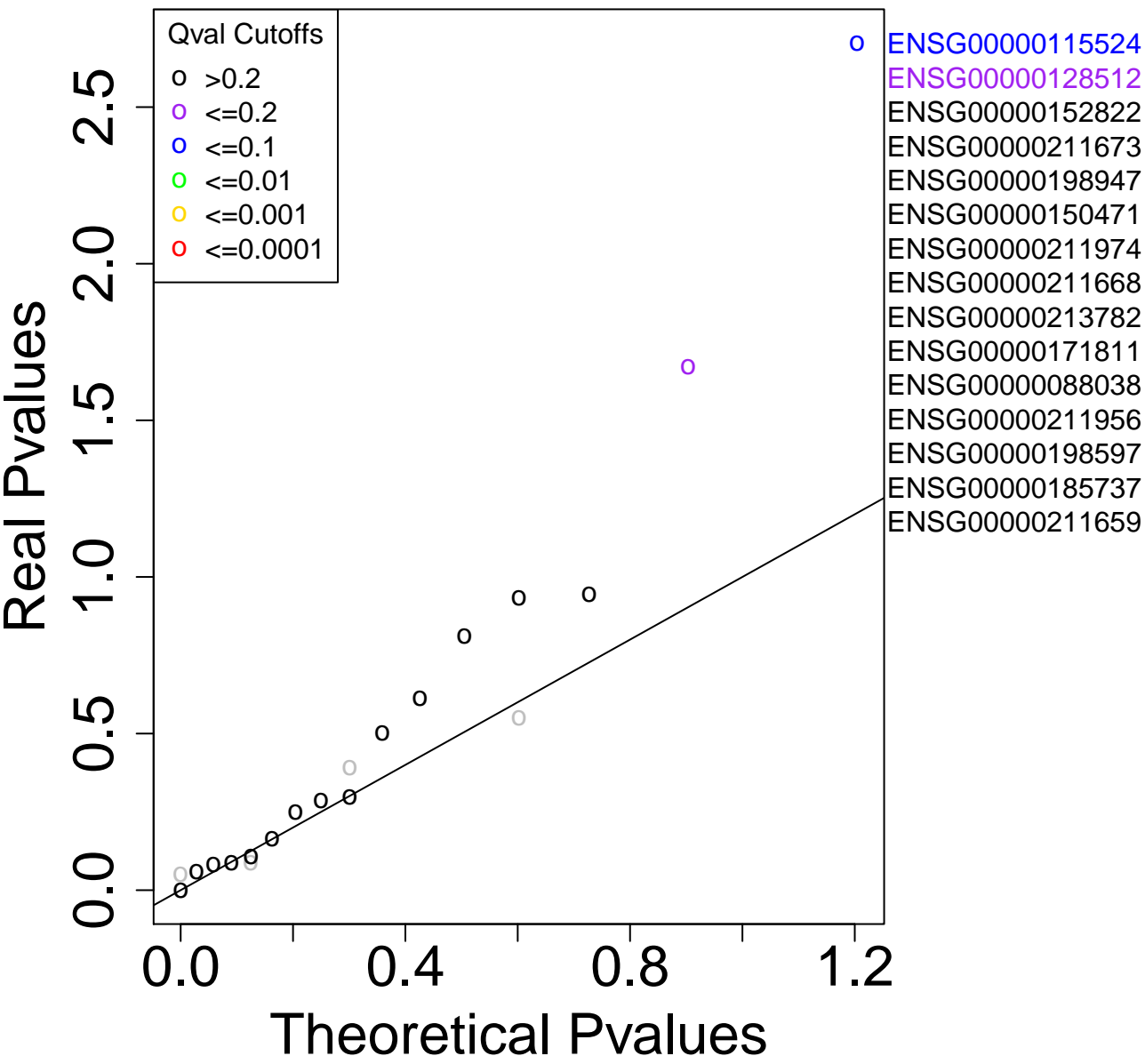
OncodriveClust Stad – 616 genes



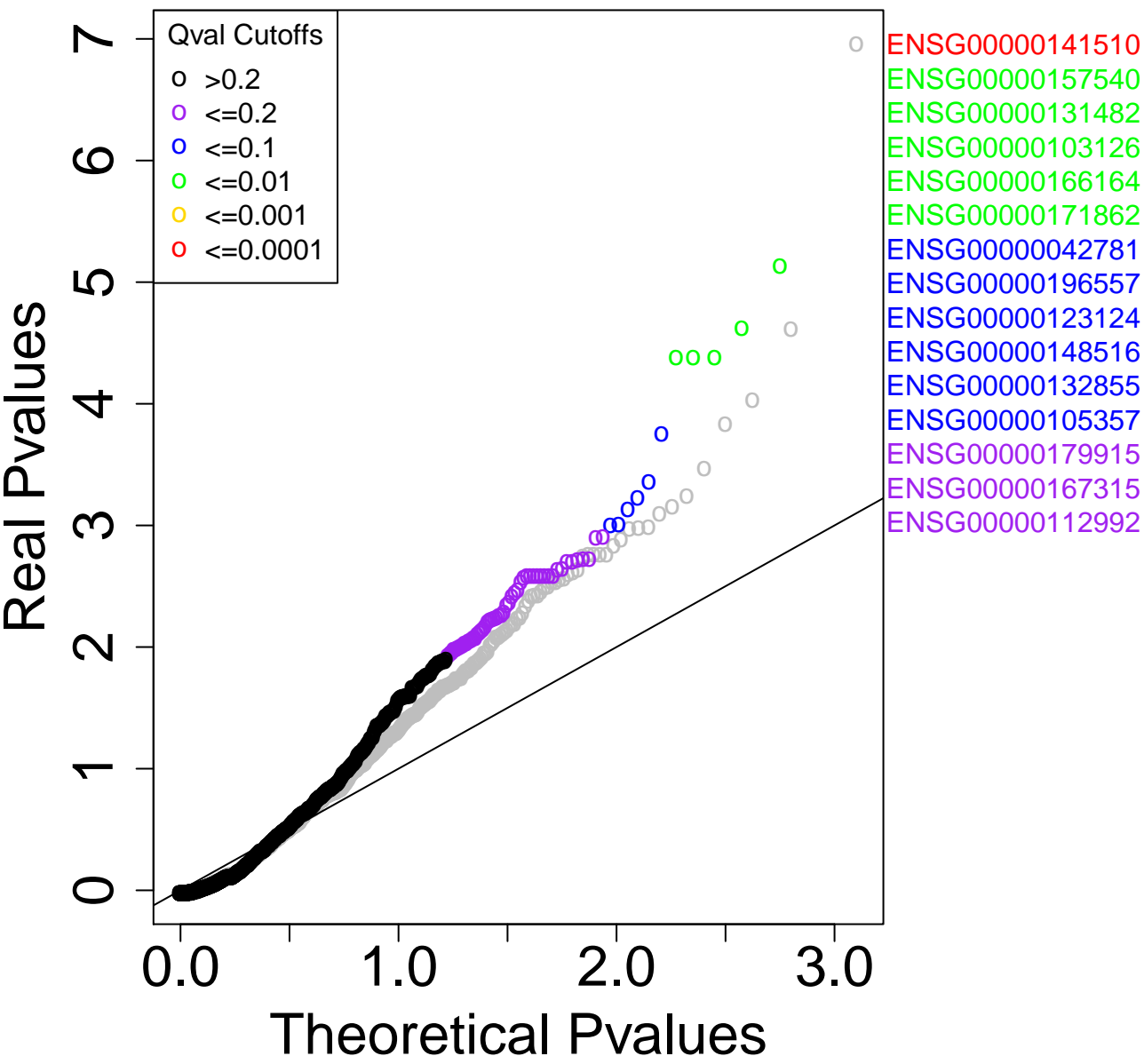
OncodriveFM Breast – 2201 genes



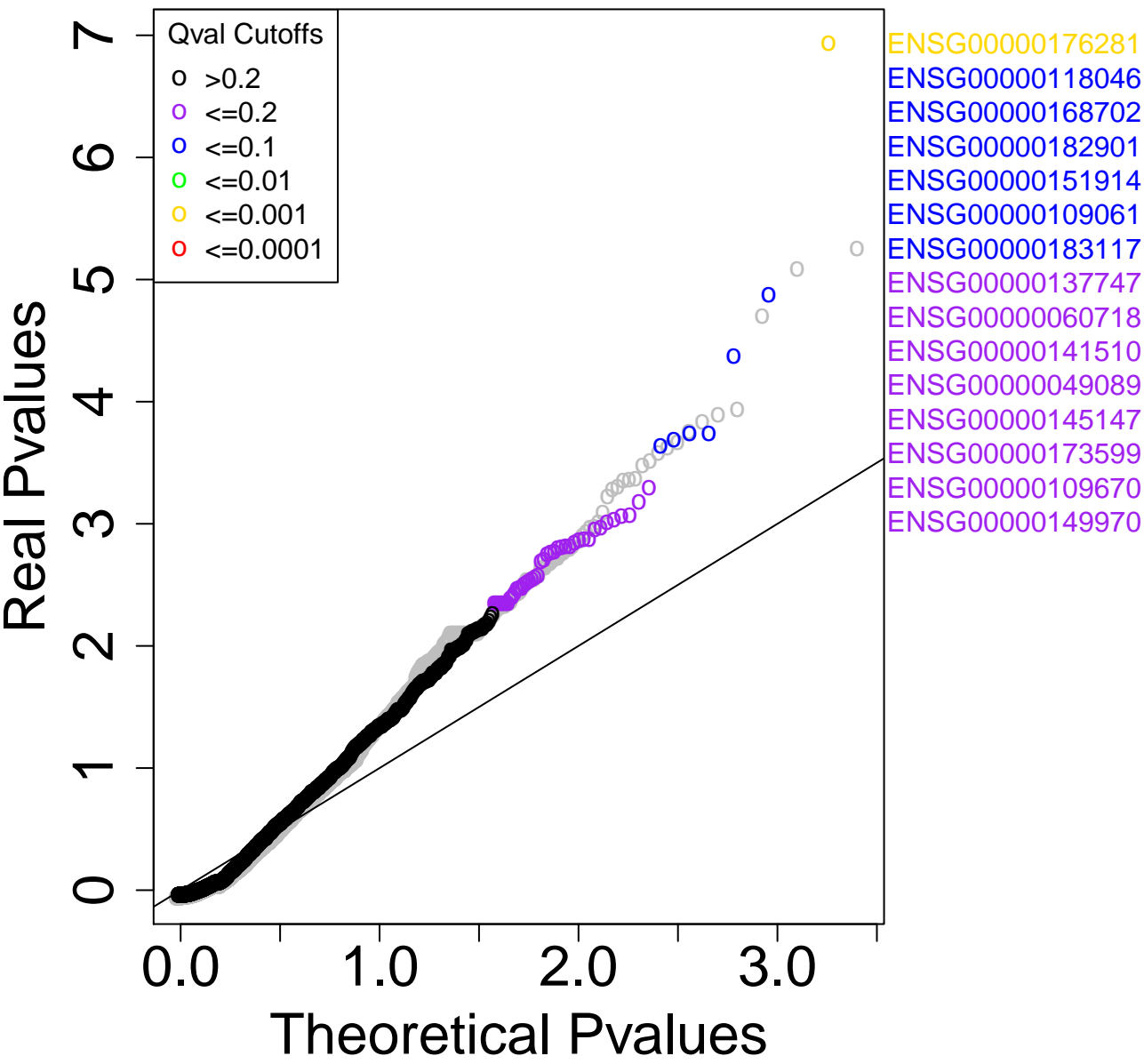
OncodriveFM CLL – 20 genes



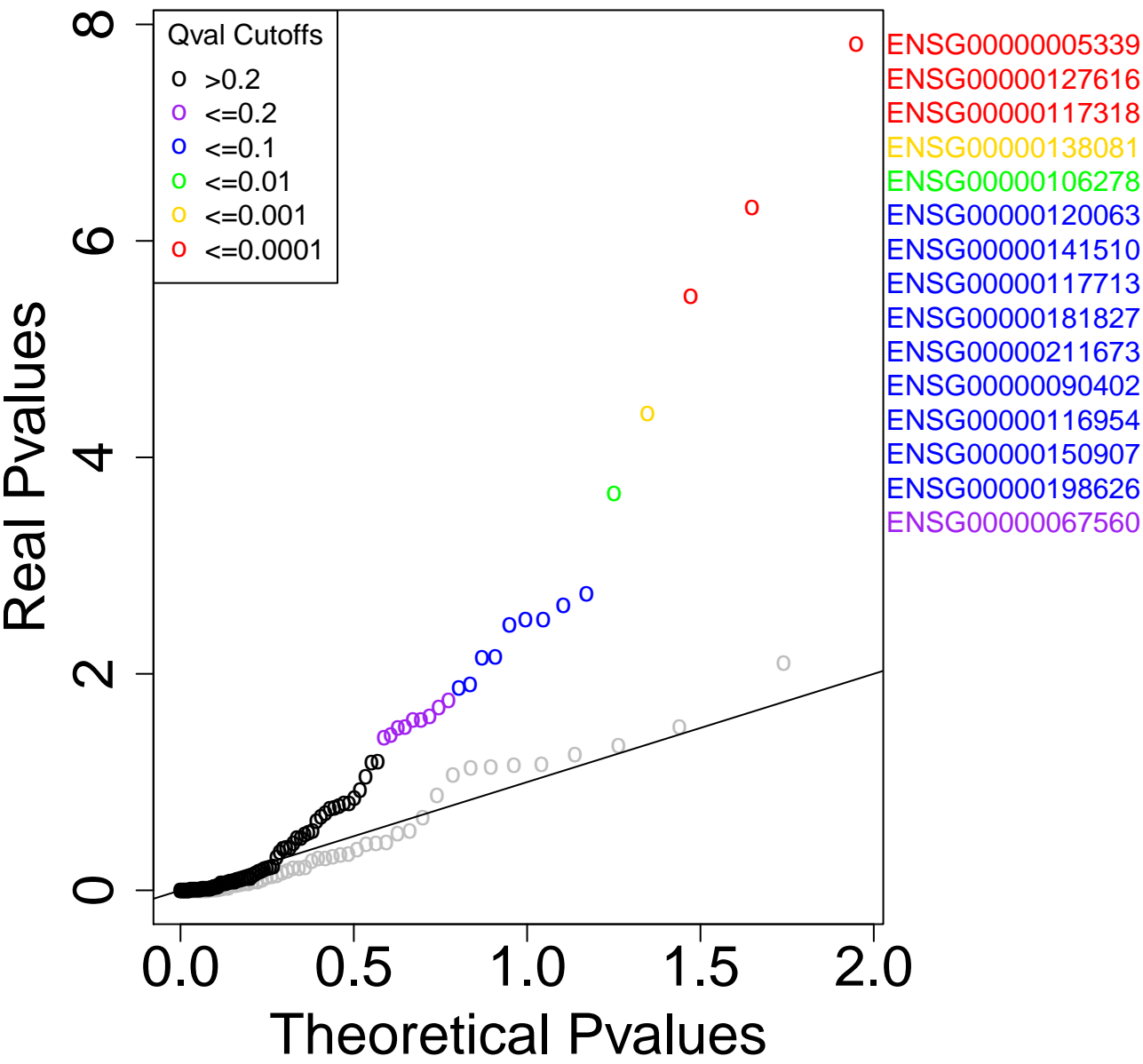
OncodriveFM Liver – 2383 genes



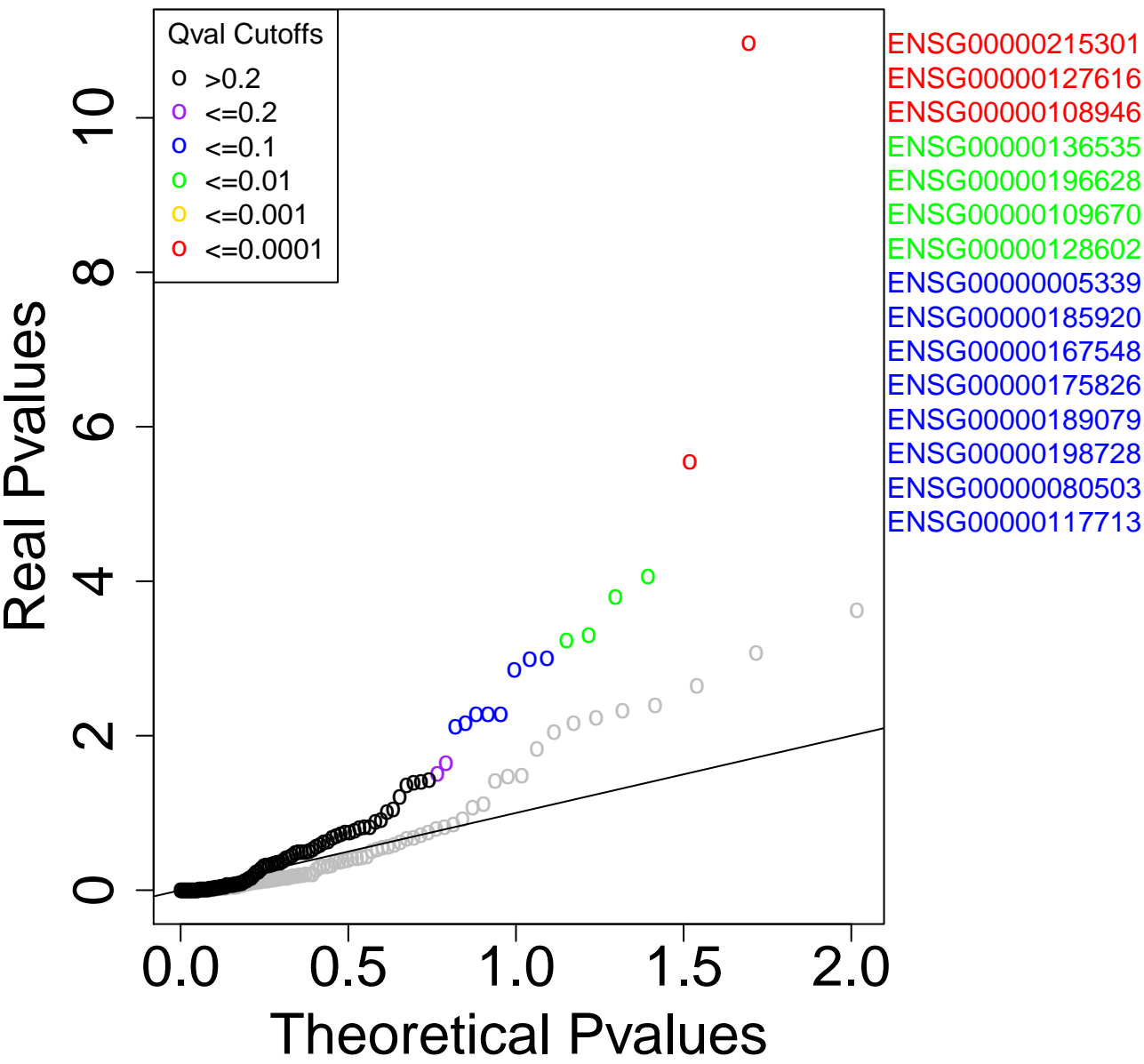
OncodriveFM Lung adeno – 4319 genes



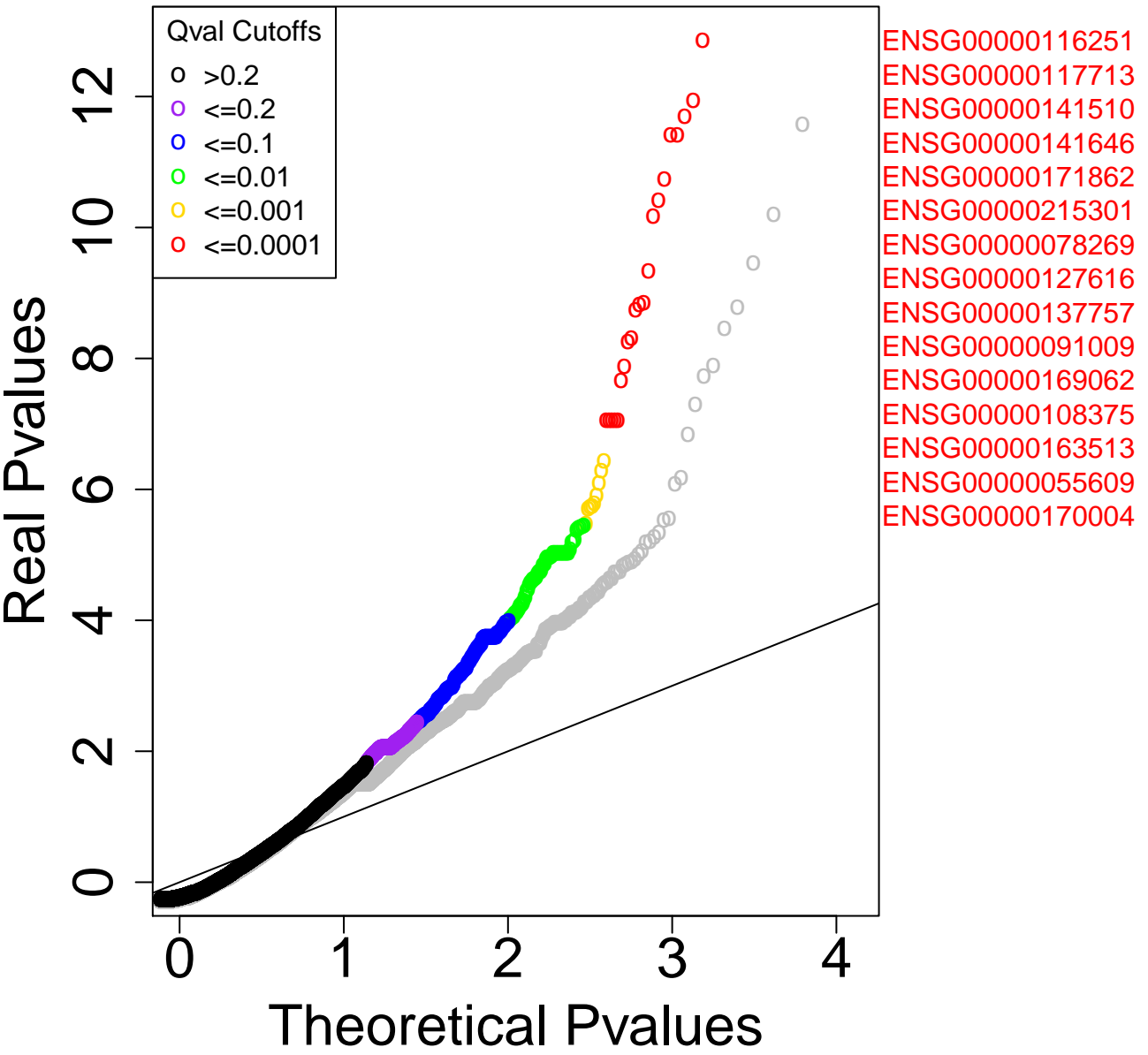
OncodriveFM Lymphoma B-cell – 144 genes



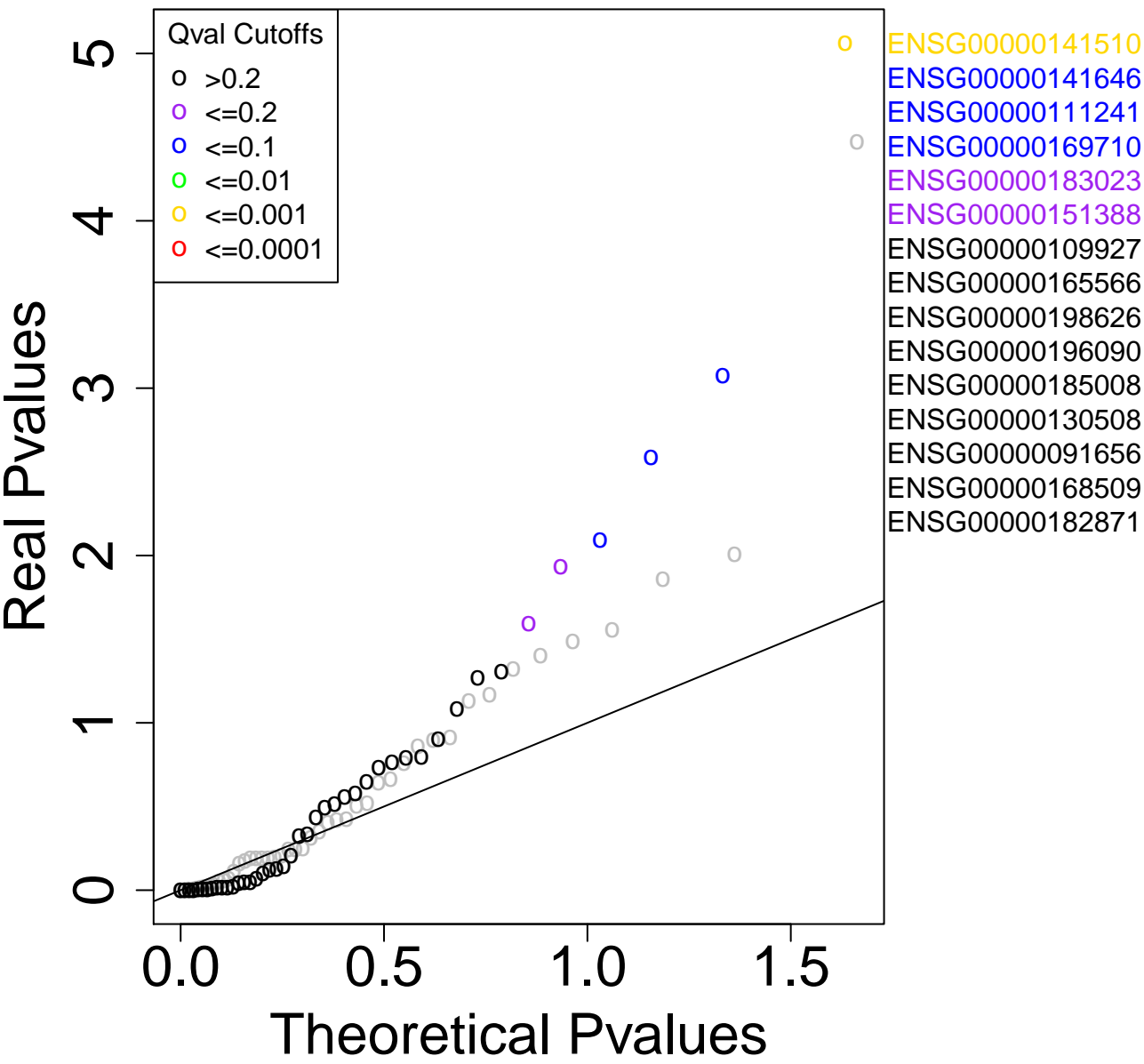
OncodriveFM Medulloblastoma – 203 genes



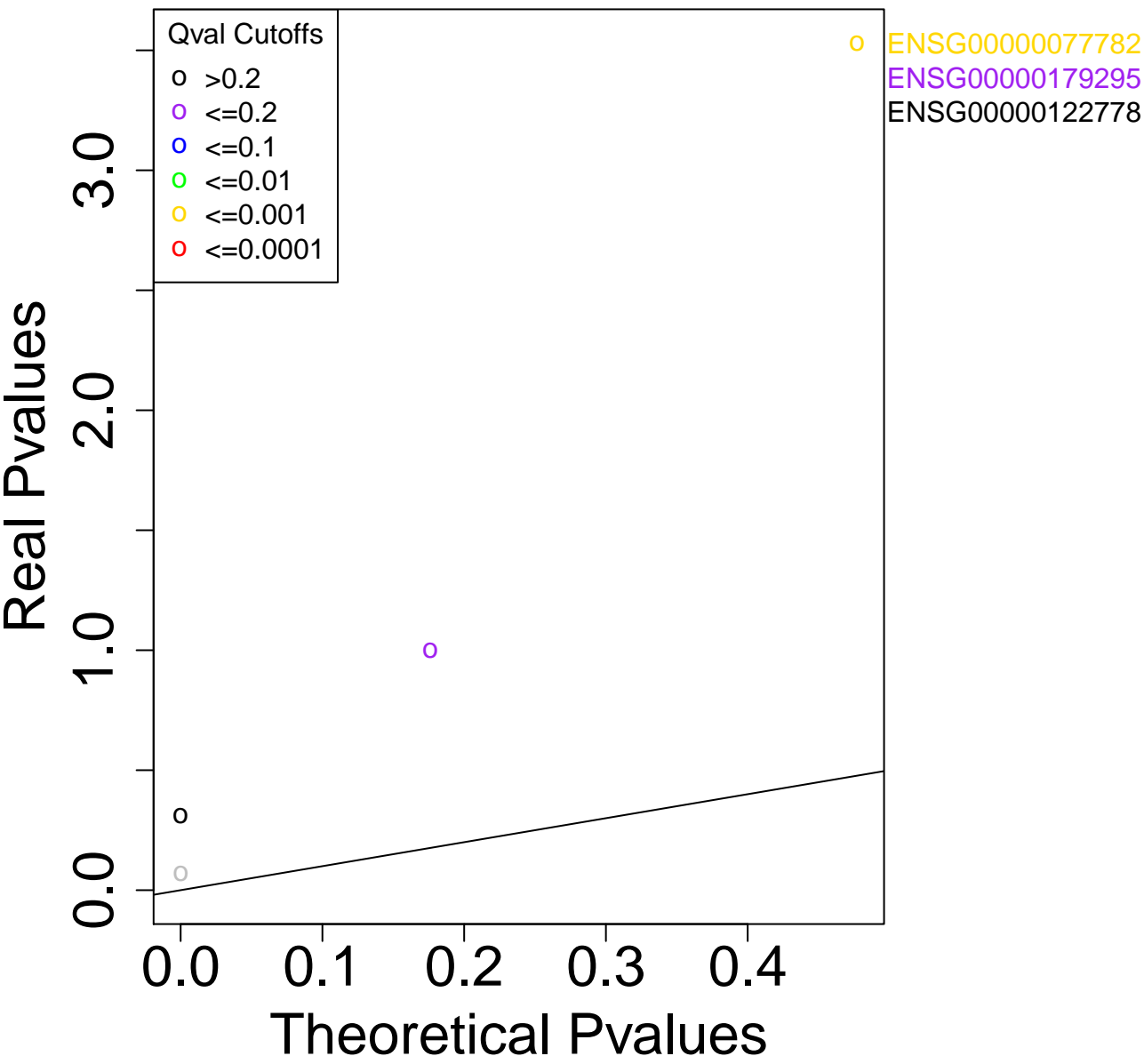
OncodriveFM Pancancer – 23195 genes



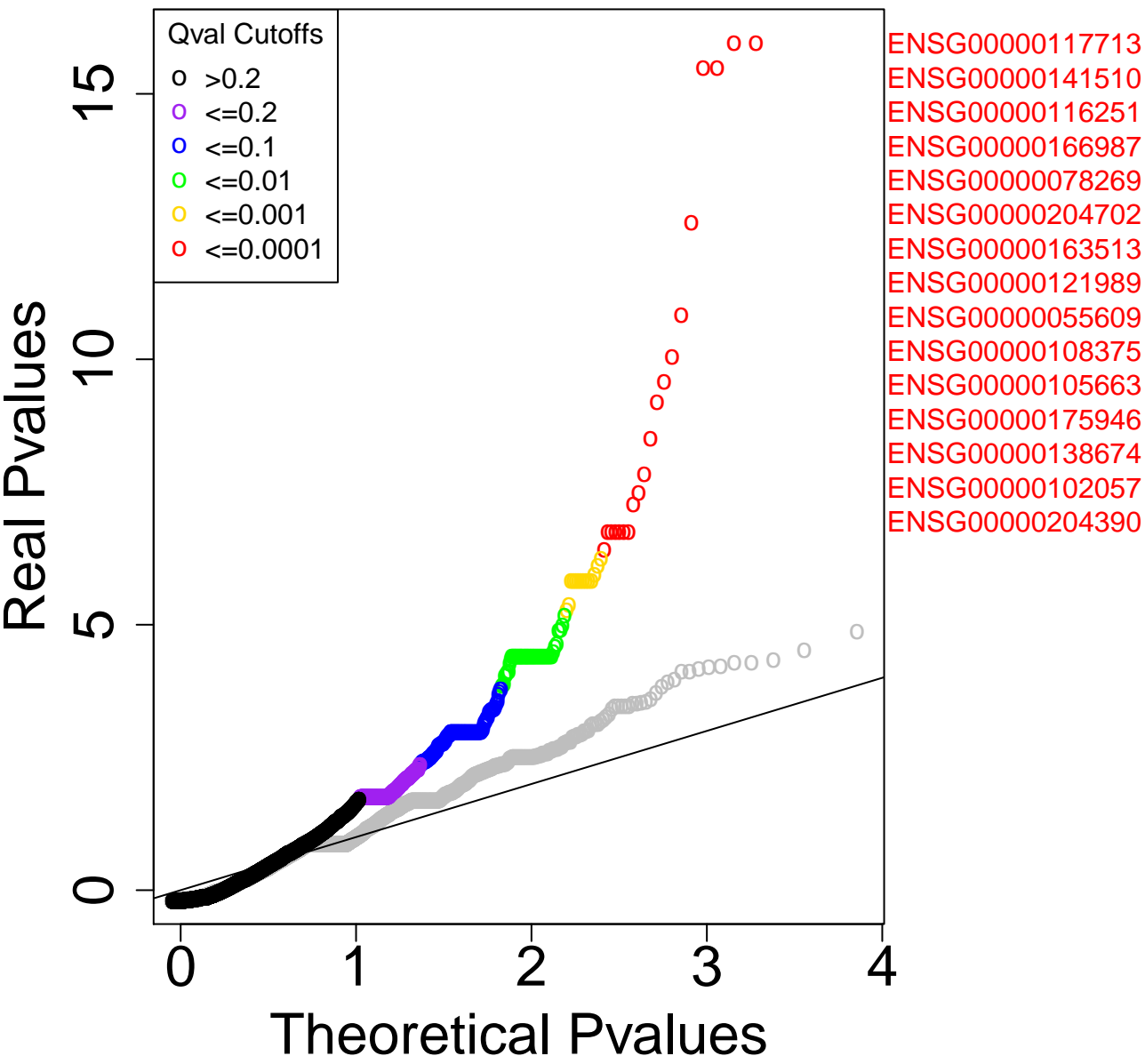
OncodriveFM Pancreas – 89 genes



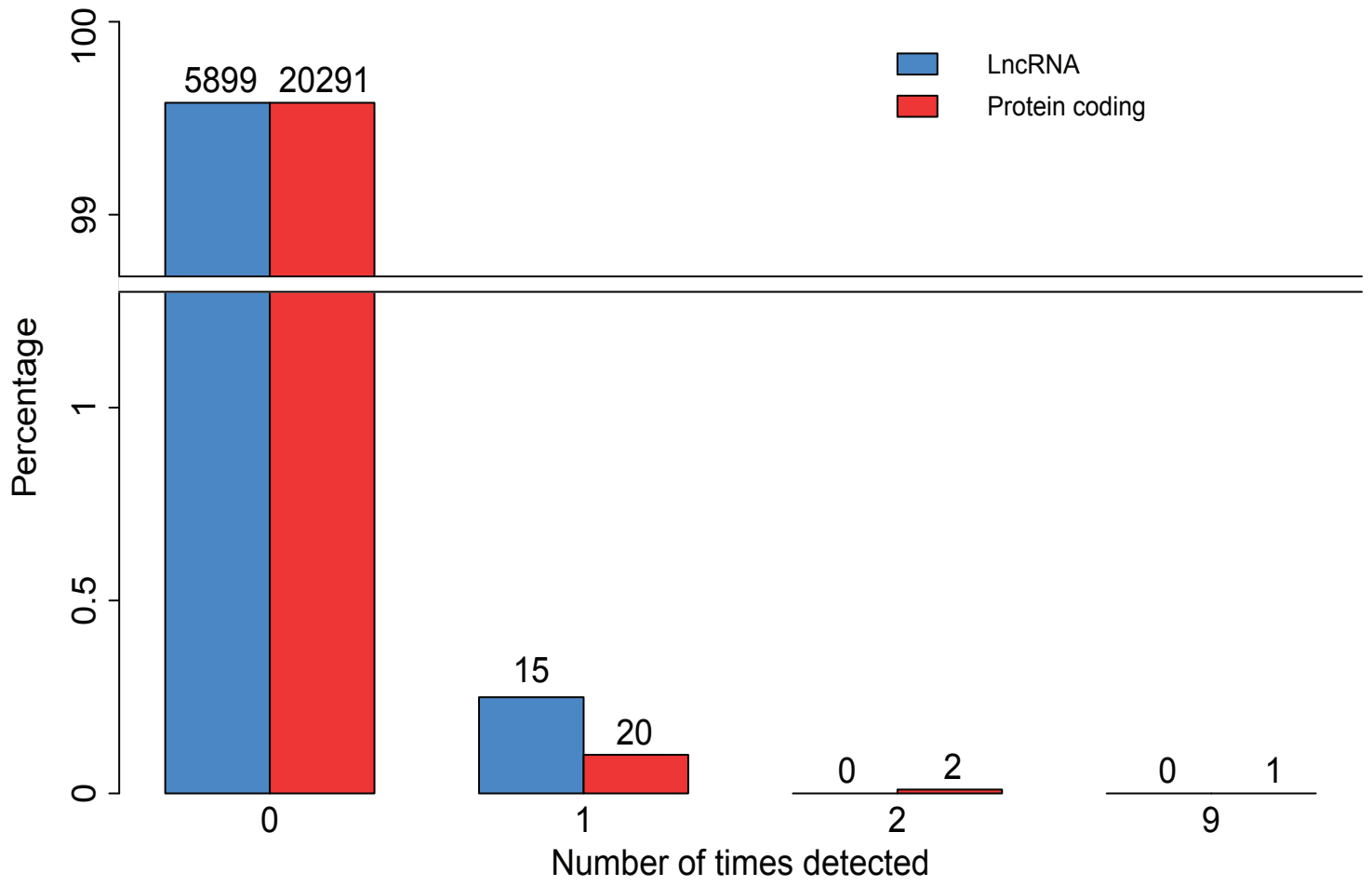
OncodriveFM Pilocytic astrocytoma – 4 genes



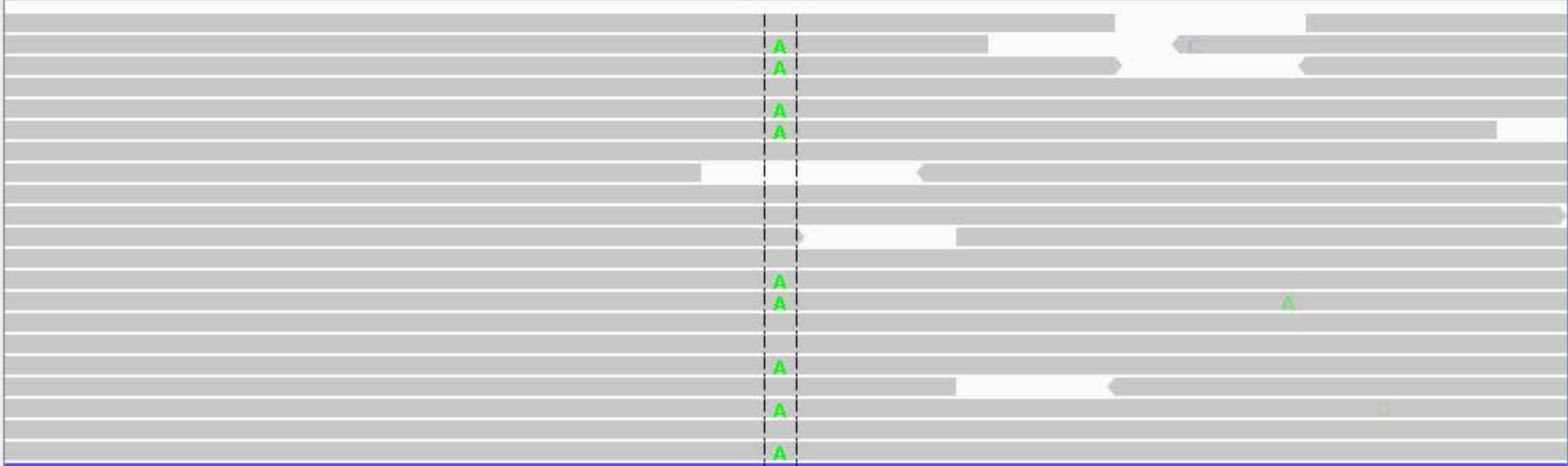
OncodriveFM Stad – 12899 genes



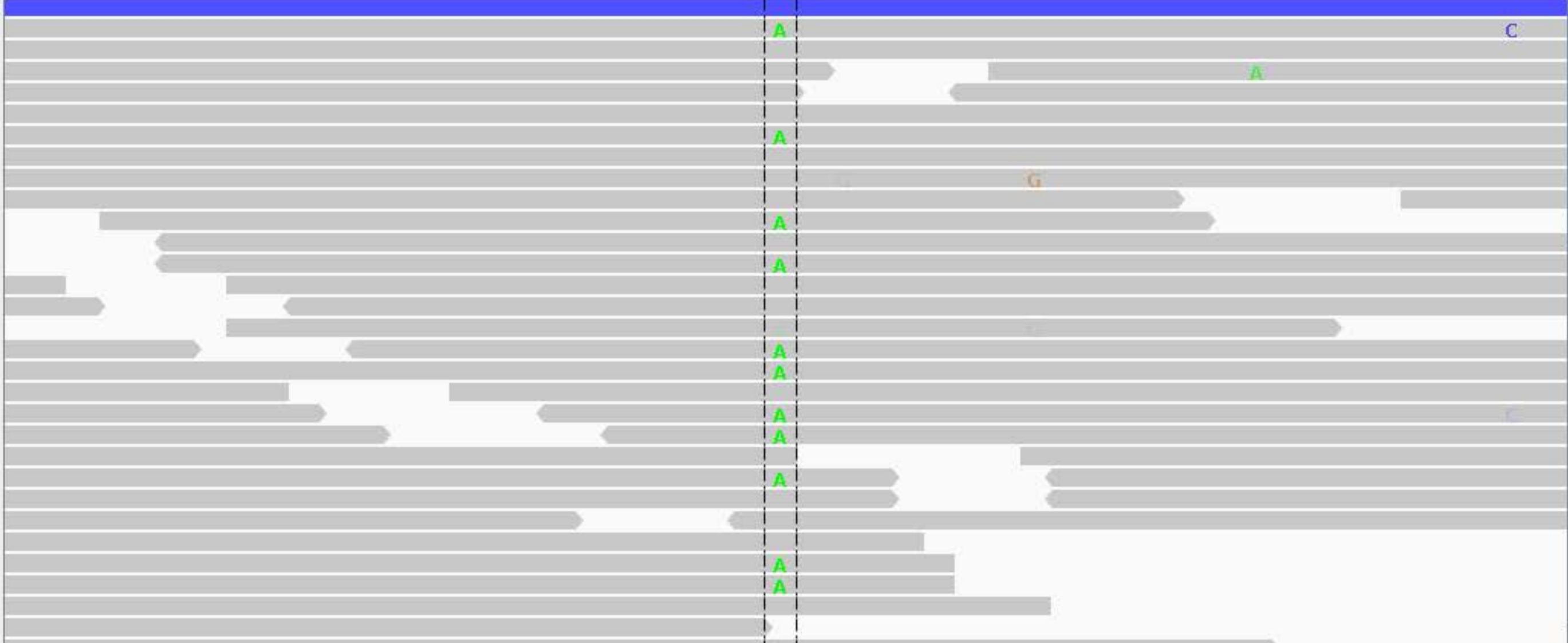
Supplementary Figure S12



12a5ba81-91e2-40a4-8fa6-b7f6
0463_9_109262730.sam Covera



12a5ba81-91e2-40a4-8fa6-b7f6
0463_9_109262730.sam

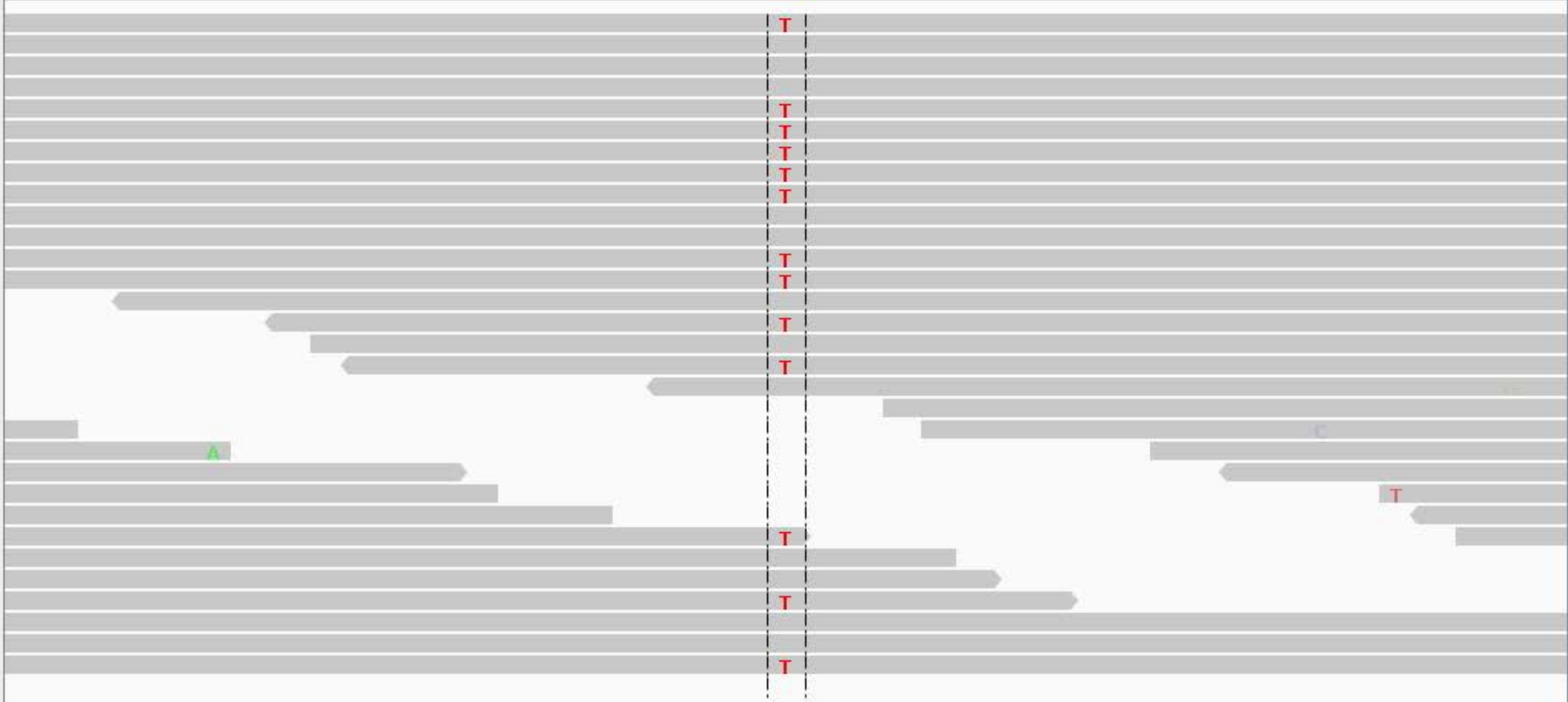


146c97a5-69dd-4863-be0b-81b
957e_9_109366290.sam Coverd

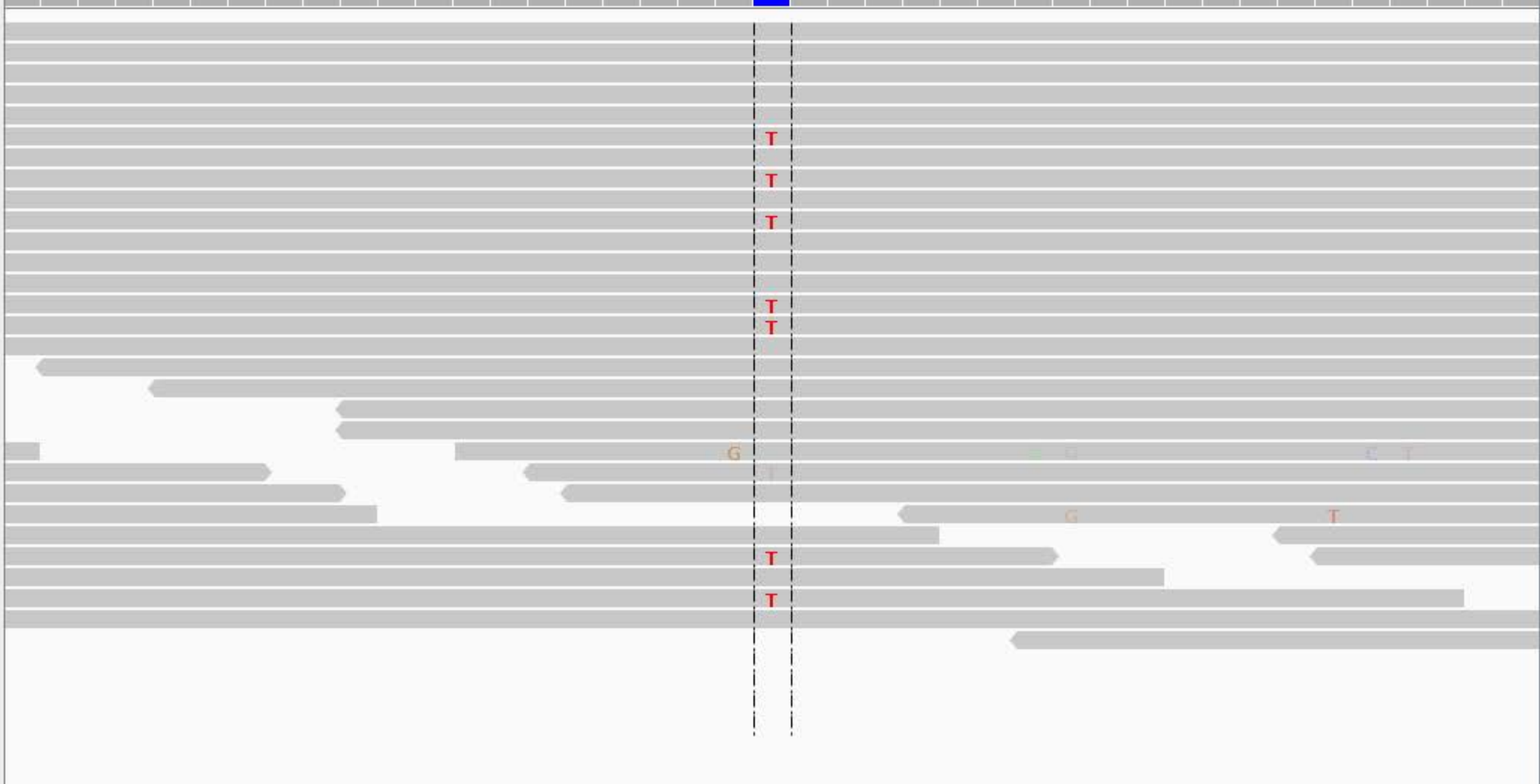
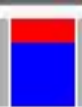
[0 - 31]



146c97a5-69dd-4863-be0b-81b
957e_9_109366290.sam







T
T
T

T
T

G
T

T
T

G

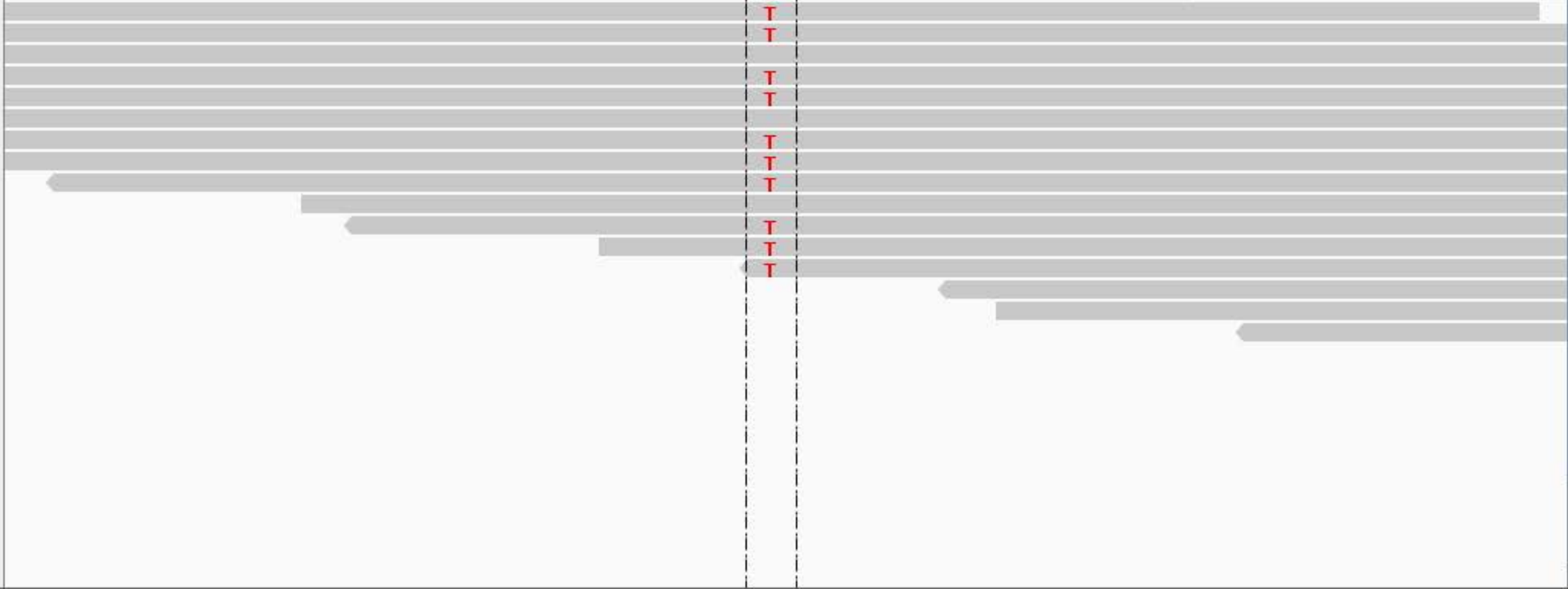
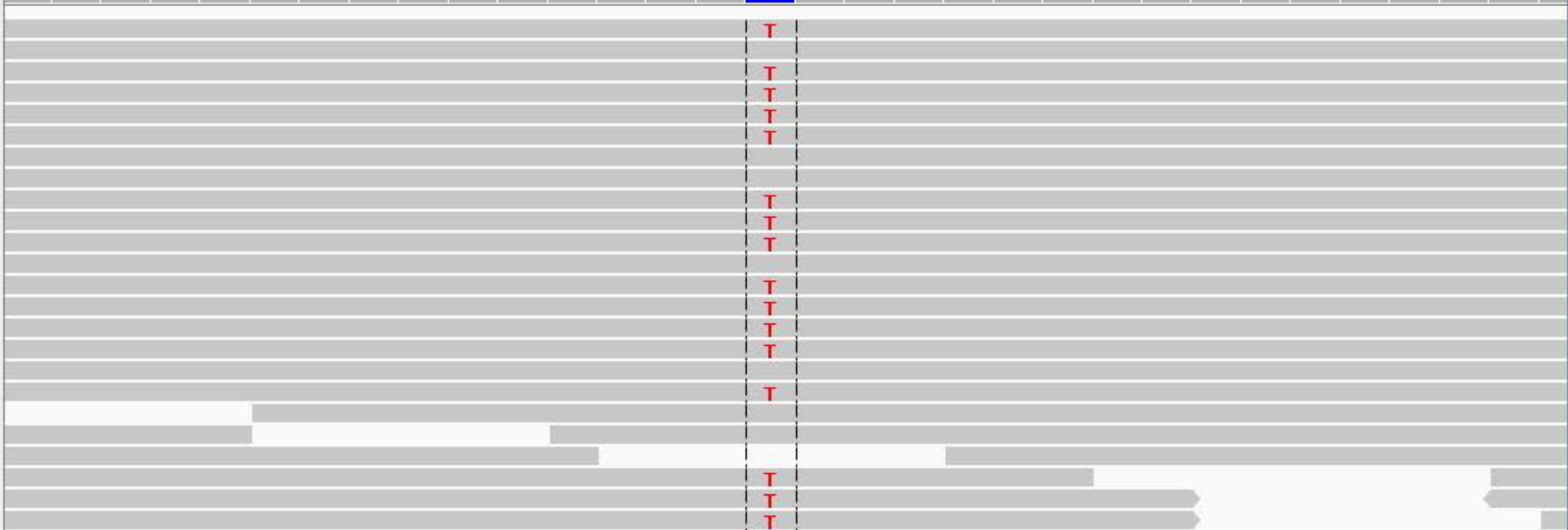
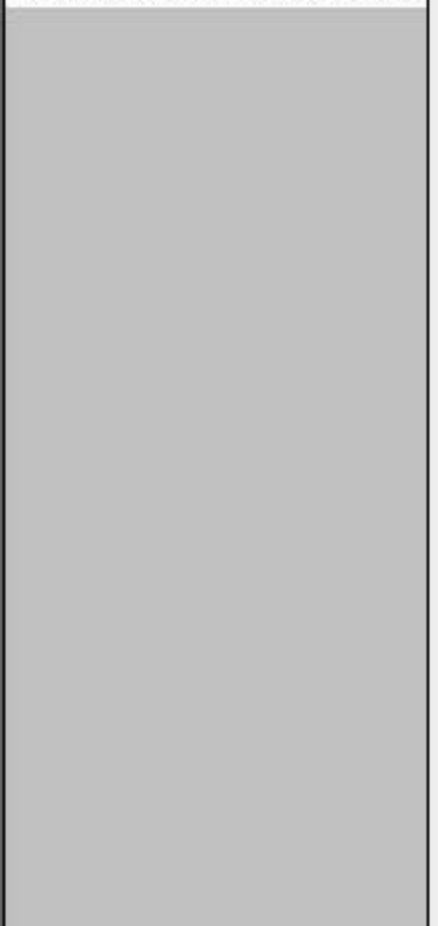
G G

G T

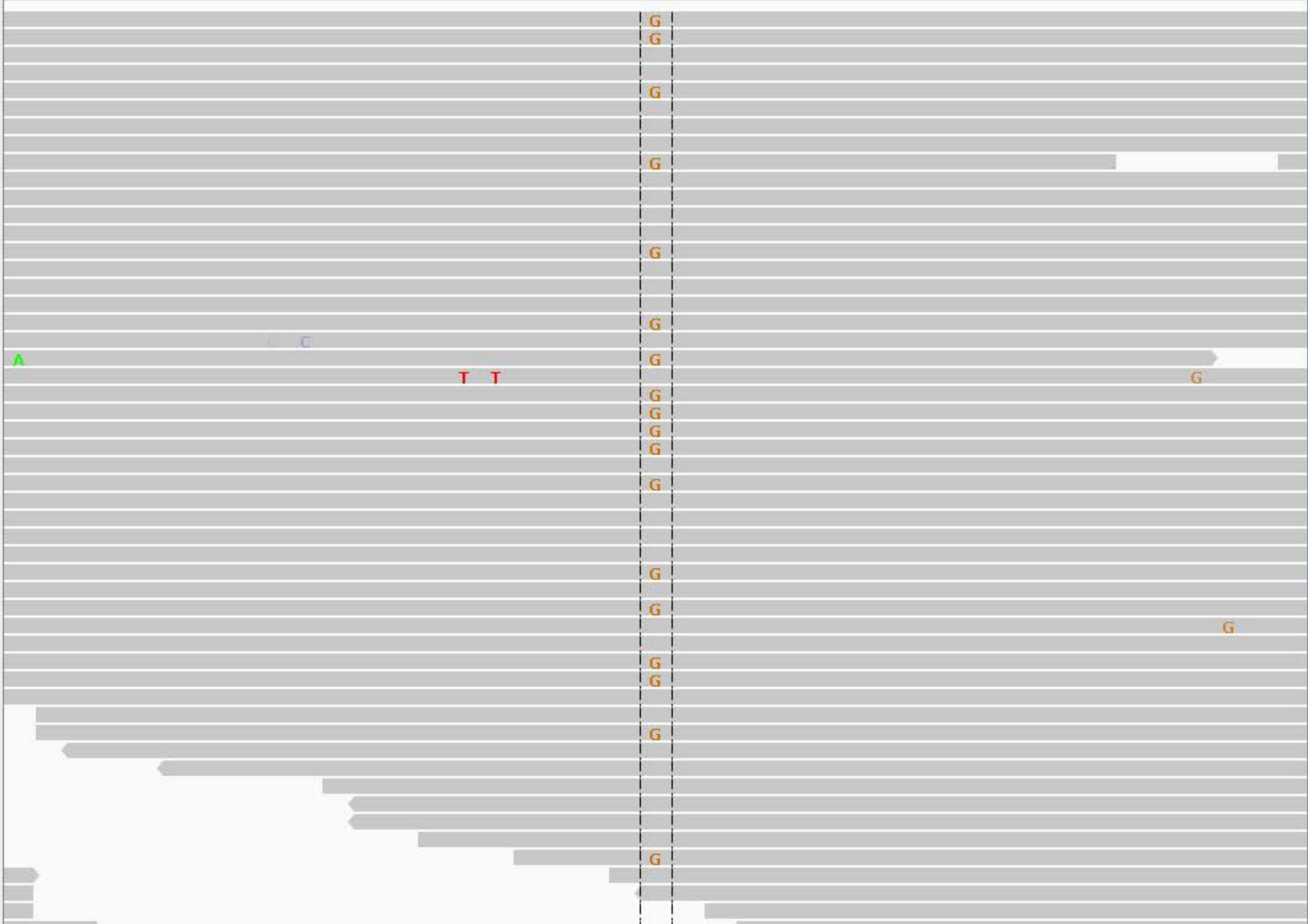
G

T

T
T

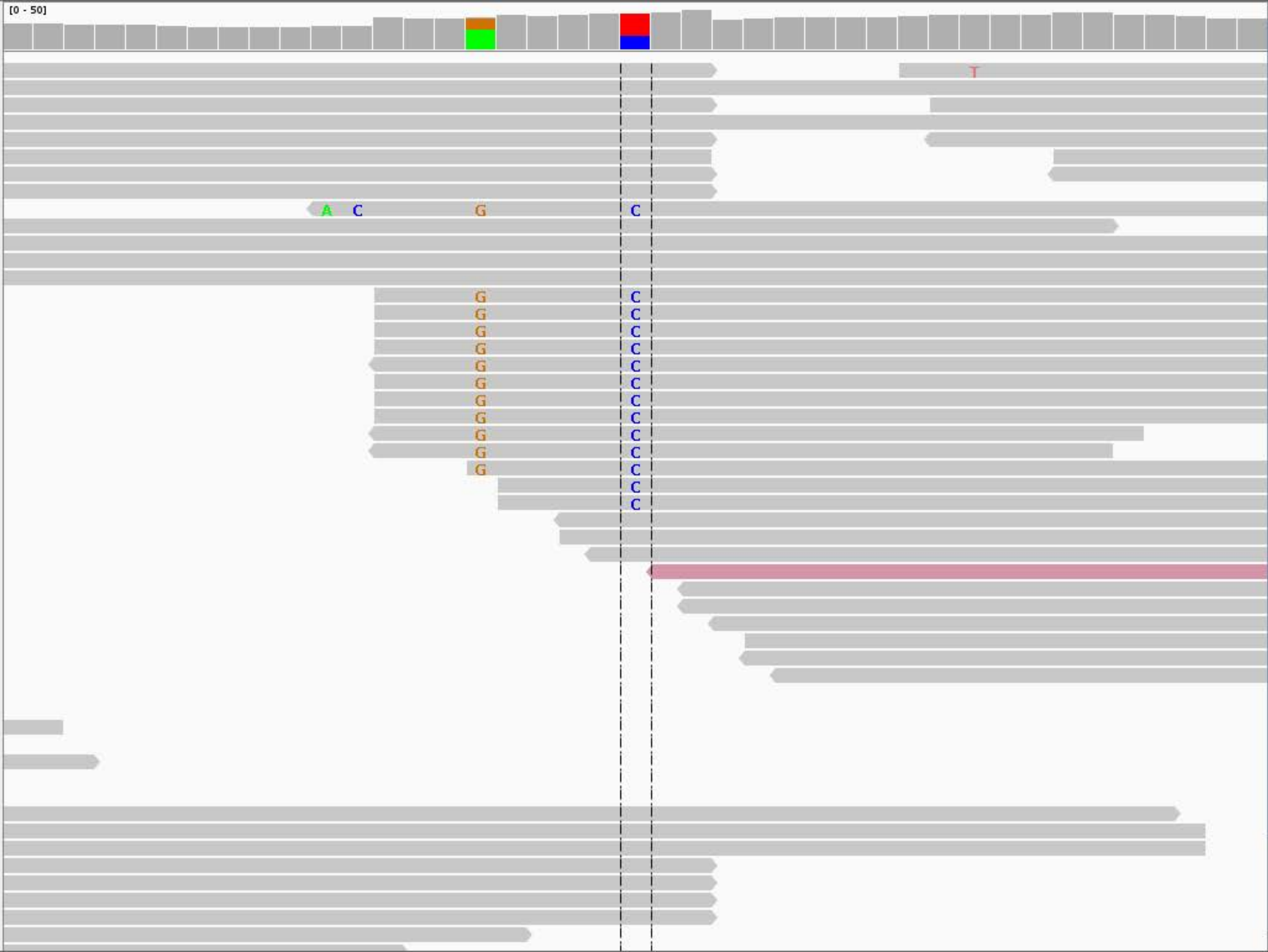






d5115b50-b8bd-4ac4-9c94-4e3
f9b1_21_25677409.sam Coverage

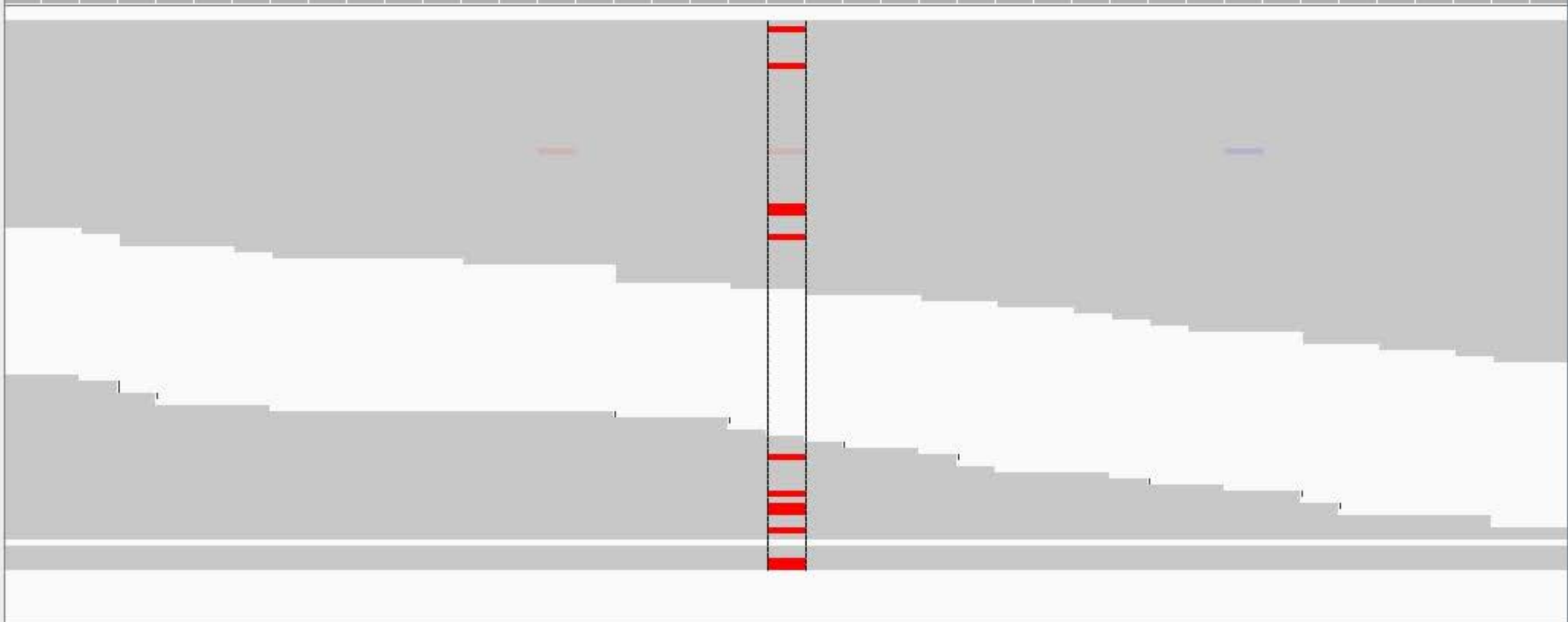
[0 - 50]



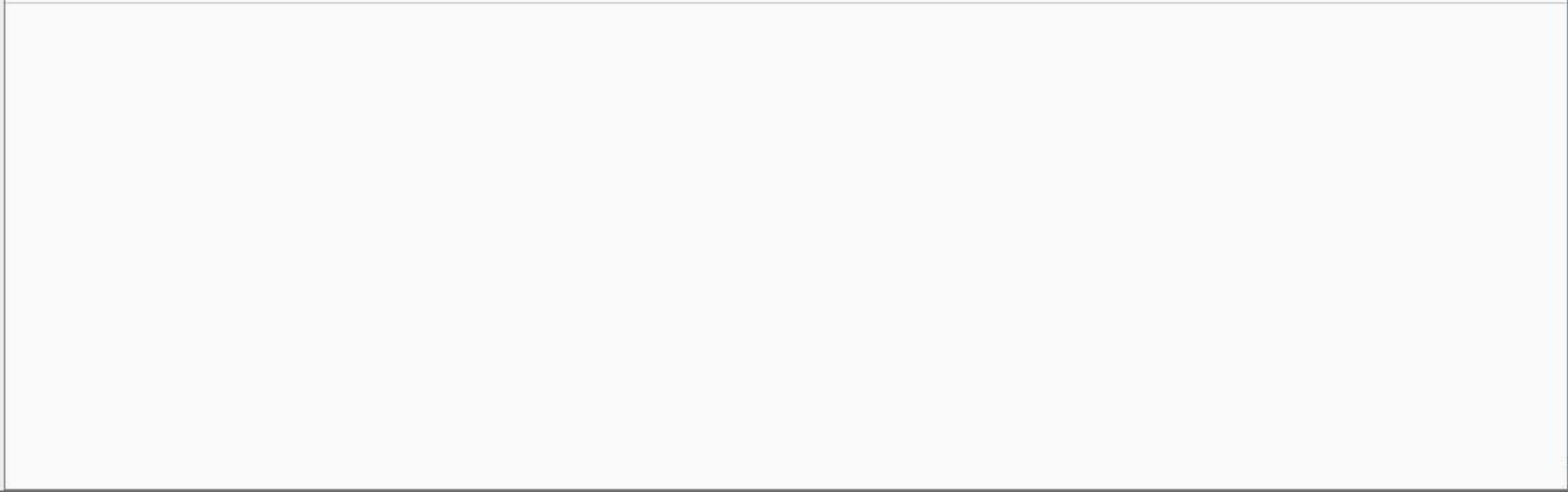
d5115b50-b8bd-4ac4-9c94-4e3
f9b1_21_25677409.sam

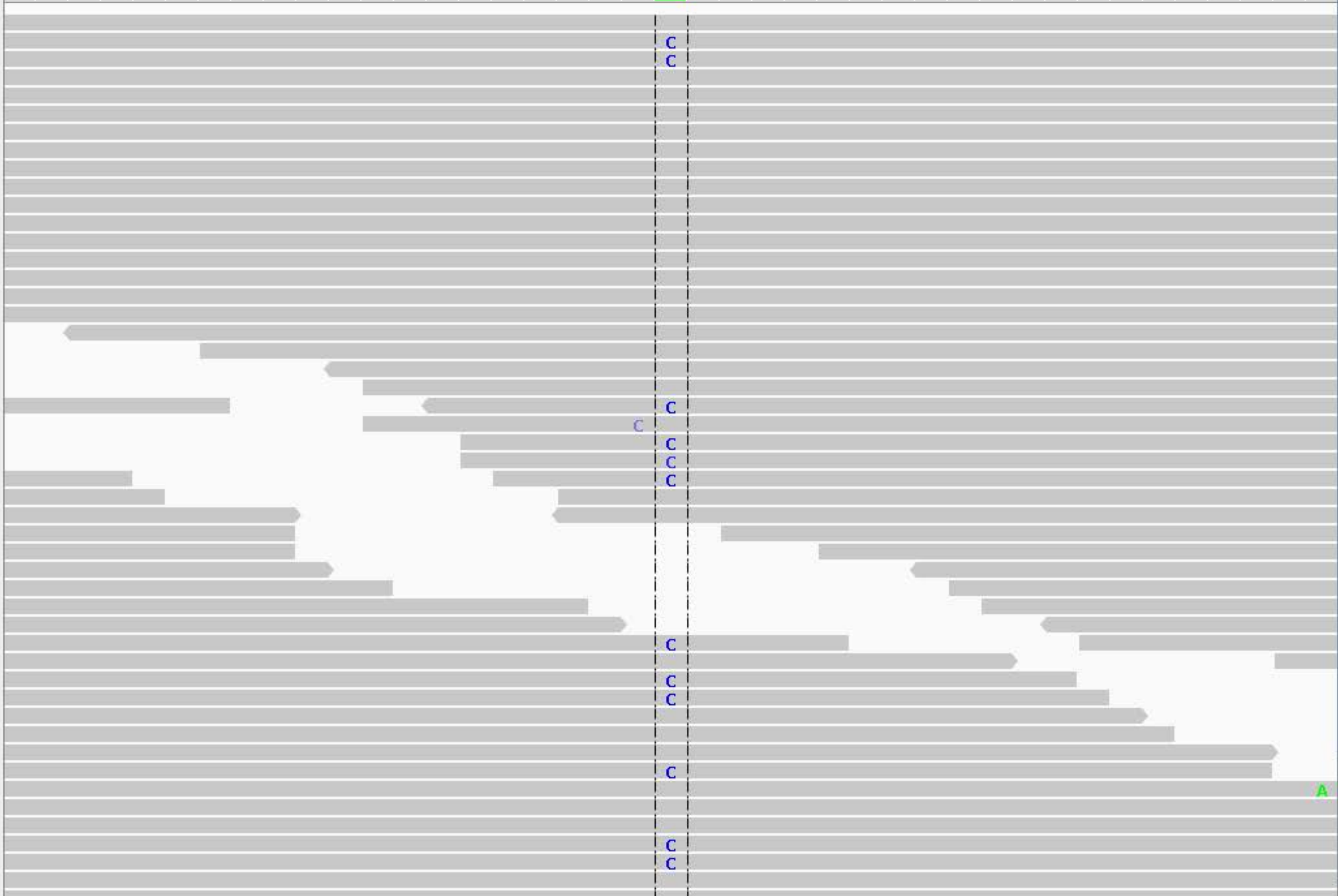
e72c727c-557e-424d-b045-356:
f9c6_9_109366501.sam Covera

[0 - 86]



e72c727c-557e-424d-b045-356:
f9c6_9_109366501.sam





Supplementary Figure S14

