

Supplementary Methods

Re-annotating PAS in regions with multiple genes

To avoid PAS being counted twice for being on different transcripts, the gene ID column was removed and duplicate lines in the table were removed. PAS were then split into two groups: PAS lying on overlapping genes and PAS that belonged to a single gene. PAS in overlapping genes were split apart by gene and PAS non-protein-coding genes were removed. Remaining PAS had their locations re-annotated according to ENSEMBL (1) release 96 using the same script that was used in the PolyA Site 2.0 workflow. To get the correct annotation, each PAS was annotated separately. The annotation script also required a GTF file that was restricted to the gene where the site is located. The --ds-range option was set to 1000 which specifies the number of nucleotides downstream of a gene's 3'UTR that gets designated as the "downstream region" (DS). For each group of PAS found in overlapping genes, if one or more PAS are in a terminal exon, all overlapping PAS outside of the terminal exon were removed. Next, PAS were removed if they were located outside of the 23 main chromosomes and if PAS were labeled by the annotation script as being in Antisense Exons (AE), Antisense Introns (AI), Intergenic (IG) and Downstream regions.

Selection of predominant hexamers

When there was only one candidate hexamer for a predominant PAS, that was selected as the predominant hexamer. When there were multiple candidate hexamers, candidate hexamers were ranked by strength in order of canonical-AATAAA, secondary-ATTAAA, and all other candidate hexamers. We also took into account the location of hexamers, which are preferred to be located 21 nucleotides upstream from the PAS (2-4). We considered the proximity of

hexamers from the -21 position from the representative site when selecting for predominant hexamers. If the strongest candidate hexamer was closest to the -21 position, that candidate hexamer was selected as the predominant hexamer. If the closest hexamer to the -21 position was not the strongest, we compared the distance of that stronger hexamer with a distribution of the distances at which similar hexamers occur as the lone hexamer upstream of pPAS. If that stronger hexamer was within 2 SD of this distribution, the stronger hexamer was chosen. In cases where two equally strong hexamers were equidistant to the -21 position relative to representative site, both were selected as predominant hexamers. The location of the pPAS and hexamers were then converted to hg19 using Crossmap v 0.5.4 (5).

Constraint analysis quality control and selection of trimer controls

To only include high quality hexamer variants in this analysis, we filtered for genomic regions with ≥ 20 read depth in 90% of gnomAD (6) samples. Hexamers that overlapped either Repeatmasker (7-9), segmental duplications (8, 10, 11), or regions with low mapability (12-15) were removed, SNVs in the remaining regions were retrieved. Hexamer SNVs with $> 50\%$ gnomAD population allele frequency were excluded from the analysis.

To define control sequences in the 3' UTR, 3' UTR intervals were downloaded (ENSEMBL BioMart v96). The 3' UTR intervals belonging to canonical gene transcripts were filtered for ENSEMBL transcript IDs belonging to canonical transcripts defined in the GENCODE v32 'known canonical' (16, 17) UCSC Table Browser file (8). Canonical 3' UTR regions immediately upstream of each predominant hexamer were retrieved using Bedtools Suite (18). From the defined canonical 3' UTR regions, the coordinates of each unique trimer found within the 18 unique hexamers were extracted by a custom python script. The set of unique trimers was derived from the 5' and 3' halves of all unique hexamers. Trimers in problematic regions were

filtered in the same manner as described in the filtering of the predominant hexamers. Trimers overlapping PAS and hexamers known to PolyASite 2.0 were removed. SNVs in trimer sequences were extracted from gnomAD and trimers containing variants with >50% allele frequency were excluded. Finally, from the remaining trimers, we selected the two most 5' trimers from the 3' UTR that matched the sequence of the corresponding downstream hexamers. In cases where two predominant hexamers existed, one set of trimers were chosen if the sequence of the two predominant hexamers were the same and two sets of trimers were selected if the sequence of the predominant hexamers differed.

Standard RNA sequencing

For standard RNA sequencing, RNA was depleted of globin mRNAs using the GLOBINclear kit (ThermoFisher Scientific, Waltham, MA). Libraries were then prepared using the TruSeq Stranded mRNA kit (Illumina). Dual-indexed barcode adapters (Illumina) were ligated to each library to mitigate index hopping. Equal amounts of each library were pooled together and pooled libraries were sequenced at NISC on an S4 flow cell of an Illumina NovaSeq 6000 with v1.0 reagents (Illumina). Sequences were derived using Real-Time-Analysis software version 3.4.4 (Illumina).

Paired-end reads were trimmed with Trimmomatic v.0.36 (22) using the sequences of the Truseq3 paired-end adapters provided with the program. The remaining reads were then aligned to hg19 with STAR v.2.7.3a. The options used were --OutSAMtype BAM Unsorted SortedByCoordinate --quantMode TranscriptomeSAM GeneCounts. This created two BAM files; one aligned to the transcriptome and another aligned to the reference genome. The BAM file that was aligned to the transcriptome had duplicate reads marked with Picard v.2.22.2 (23). The

BAM file was converted for RSEM (24) using the convert-sam-for-rsem tool in RSEM. Finally, gene expression was quantified using RSEM v.1.3.2 with the `--paired-end` option.

Preparation of 3' end sequencing and identification of PAS

Total RNA was isolated from whole blood using the PAXgene Blood RNA system (Qiagen, Gathersburg, MD) from 76 participants in the ClinSeq® cohort. For 3' end sequencing, libraries were prepared using the Quantseq_REV protocol (Lexogen, Greenland, NH) with the inclusion of the RS-Globin Block kit (Lexogen) to remove globin mRNAs. Each library was ligated with unique dual-indexed adapters to prevent index hopping. Equal concentrations of each library were pooled together and pooled libraries were sequenced at the NIH Intramural Sequencing Center (NISC) on an SP flow cell on a NovaSeq 6000 with v1.0 reagents (Illumina, La Jolla, CA). Sequences were derived using Real-Time-Analysis software v3.4.4 (Illumina). Replicate FASTQ files generated for each sample were merged.

To generate the BAM files for APA on the UCSC browser, QuantSeq_Rev generated FASTQ files were aligned to hg19 using STAR v.2.7.3a (19). Using PolyASite 2.0, we defined a set of PAS and hexamers that were used in the samples. PolyASite 2.0 uses reads from the 3' end cleavage site to identify PAS and candidate hexamers (4). FASTQ files were mapped to hg19 using workflow tools. Snakemake v.5.5.1 (20) was employed using container software run with Singularity v.3.6.4 (21). The specified 3' adapter was derived from the sequencing protocol: 'ACACTCTTTCCCTACACGACGCTCTTCCGATCTTTTTTTTTTTTTTTTTTTTTT'. For the required input table, all reads were set to `'paired'=False` and `'reverse_compl' = True`.

The identified PAS were annotated with gene symbols according to the hg19 gene model. To annotate PAS downstream of known transcripts, all transcript coordinates from the hg19 GTF file were extended on the 3' end by 1,000 bases using bedtools. Re-annotating PAS in regions

with overlapping genes was done using the same methods mentioned for the PolyASite 2.0 database (See reannotation section). PAS occurring in antisense regions were also removed.

Identification of variants in predominant hexamers in RNAseq data

SNVs occurring in our set of predominant AATAAA and ATTAAA hexamers were retrieved from exome data using bcftools view. Variants with >50% internal frequency were removed as it suggested they were the major allele in the population, and variants in genes that were not expressed in the 3' end RNA sequencing data were removed. CisASE v.1.0.2 was used for quality control and to generate final read counts (25). Fasta encoding options -F and -f were both set to 33 and -M 3 was used. Mpileups were generated using SAMtools mpileup v.1.9 (26) (with options -q, -Q=0, -C=50) on both exome, and 3' end RNA sequencing BAM files. Samples that did not produce any output were run with cisASE a second time using 3' end RNA sequencing mpileups only. If CisASE data showed there were no ALT reads present or if variants overlapped indels, we manually checked BAM files at the variant location to ensure the correct genotypes were called.

Selection of control 3'UTR variants upstream of predominant hexamers

To compare the effect on mRNA processing of SNVs found in predominant hexamers against SNVs found in non-hexamer 3' UTR regions, we randomly sampled 60 genes that contained a single SNV in the 3' UTR region upstream of the predominant hexamer in an individual from our data set. The originally created canonical 3' UTR bed file (See Constraint analysis of hexamer variants) was lifted over to hg19. The 64 hexamers that contained SNVs in the samples were removed and the remaining 15,144 predominant hexamers were intersected with the canonical UTRs. The 3' UTR regions upstream of this limited set of predominant hexamers were extracted using bedtools subtract and bedtools slop as described in the section 'Constraint

analysis of hexamer variants'. SNVs and indels were retrieved from the 3' UTRs and variants occurring in regions upstream of predominant hexamers were flagged. UTR regions were removed if they belonged to genes not expressed in the PAS derived from the PolyASite 2.0 workflow. The 3' UTRs with a sample containing a single upstream SNV and void of any indels were considered for further analysis. Sixty 3' UTR variants belonging to 60 genes were chosen at random for analysis. Wig files of these 60 UTR regions were created from both RNA Seq and QuantSeq REV BAM files as mentioned in the 'Viewing Alternative Polyadenylation in UCSC Genome Browser' section.

Processing files to view on the UCSC Browser

To view alignments of 3' QuantSeq and RNA Seq files, WIG files were from RNA Seq and 3' QuantSeq BAM files. All 152 BAM files were filtered for the 64 hexamer variant genes and 60 control variant genes using the gene intervals from the hg19 gtf file. All gene intervals had 1.5 kb added at the 3' end using bedtools command: `bedtools slop -l 0 -r 1500 -s -g <genome file> -i <bedfile of genes>`. Wig files were generated from filtered BAM files using `bam2wig v.3.0.1` with the `-q` option set to 0. Files were viewed in the UCSC Genome Browser (27).

REFERENCES

- 1 Yates, A.D., Achuthan, P., Akanni, W., Allen, J., Allen, J., Alvarez-Jarreta, J., Amode, M.R., Armean, I.M., Azov, A.G., Bennett, R. *et al.* (2020) Ensembl 2020. *Nucleic Acids Res*, **48**, D682-d688.
- 2 Gruber, A.J., Schmidt, R., Gruber, A.R., Martin, G., Ghosh, S., Belmadani, M., Keller, W. and Zavolan, M. (2016) A comprehensive analysis of 3' end sequencing data sets reveals novel polyadenylation signals and the repressive role of heterogeneous ribonucleoprotein C on cleavage and polyadenylation. *Genome Res*, **26**, 1145-1159.

- 3 Beaudoin, E., Freier, S., Wyatt, J.R., Claverie, J.M. and Gautheret, D. (2000) Patterns of variant polyadenylation signal usage in human genes. *Genome Res*, **10**, 1001-1010.
- 4 Herrmann, C.J., Schmidt, R., Kanitz, A., Artimo, P., Gruber, A.J. and Zavolan, M. (2020) PolyASite 2.0: a consolidated atlas of polyadenylation sites from 3' end sequencing. *Nucleic Acids Res*, **48**, D174-d179.
- 5 Zhao, H., Sun, Z., Wang, J., Huang, H., Kocher, J.P. and Wang, L. (2014) CrossMap: a versatile tool for coordinate conversion between genome assemblies. *Bioinformatics*, **30**, 1006-1007.
- 6 Karczewski, K.J., Francioli, L.C., Tiao, G., Cummings, B.B., Alföldi, J., Wang, Q., Collins, R.L., Laricchia, K.M., Ganna, A., Birnbaum, D.P. *et al.* (2020) The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature*, **581**, 434-443.
- 7 Jurka, J. (2000) Repbase update: a database and an electronic journal of repetitive elements. *Trends Genet*, **16**, 418-420.
- 8 Karolchik, D., Hinrichs, A.S., Furey, T.S., Roskin, K.M., Sugnet, C.W., Haussler, D. and Kent, W.J. (2004) The UCSC Table Browser data retrieval tool. *Nucleic Acids Res*, **32**, D493-496.
- 9 Smit, A., Hubley, R & Green, P. (1996-2010), in press.
- 10 Bailey, J.A., Gu, Z., Clark, R.A., Reinert, K., Samonte, R.V., Schwartz, S., Adams, M.D., Myers, E.W., Li, P.W. and Eichler, E.E. (2002) Recent segmental duplications in the human genome. *Science*, **297**, 1003-1007.
- 11 Bailey, J.A., Yavor, A.M., Massa, H.F., Trask, B.J. and Eichler, E.E. (2001) Segmental duplications: organization and impact within the current human genome project assembly. *Genome Res*, **11**, 1005-1017.
- 12 Zook, J.M., Catoe, D., McDaniel, J., Vang, L., Spies, N., Sidow, A., Weng, Z., Liu, Y., Mason, C.E., Alexander, N. *et al.* (2016) Extensive sequencing of seven human genomes to characterize benchmark reference materials. *Sci Data*, **3**, 160025.
- 13 Kent, W.J., Baertsch, R., Hinrichs, A., Miller, W. and Haussler, D. (2003) Evolution's cauldron: duplication, deletion, and rearrangement in the mouse and human genomes. *Proc Natl Acad Sci U S A*, **100**, 11484-11489.
- 14 Chiaromonte, F., Yap, V.B. and Miller, W. (2002) Scoring pairwise genomic sequence alignments. *Pac Symp Biocomput*, in press., 115-126.
- 15 Schwartz, S., Kent, W.J., Smit, A., Zhang, Z., Baertsch, R., Hardison, R.C., Haussler, D. and Miller, W. (2003) Human-mouse alignments with BLASTZ. *Genome Res*, **13**, 103-107.
- 16 Harrow, J., Denoeud, F., Frankish, A., Reymond, A., Chen, C.K., Chrast, J., Lagarde, J., Gilbert, J.G., Storey, R., Swarbreck, D. *et al.* (2006) GENCODE: producing a reference annotation for ENCODE. *Genome Biol*, **7 Suppl 1**, S4.1-9.
- 17 Harrow, J., Frankish, A., Gonzalez, J.M., Tapanari, E., Diekhans, M., Kokocinski, F., Aken, B.L., Barrell, D., Zadissa, A., Searle, S. *et al.* (2012) GENCODE: the reference human genome annotation for The ENCODE Project. *Genome Res*, **22**, 1760-1774.
- 18 Quinlan, A.R. and Hall, I.M. (2010) BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*, **26**, 841-842.
- 19 Dobin, A., Davis, C.A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M. and Gingeras, T.R. (2013) STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*, **29**, 15-21.

- 20 Mölder, F., Jablonski, K.P., Letcher, B., Hall, M.B., Tomkins-Tinch, C.H., Sochat, V., Forster, J., Lee, S., Twardziok, S.O., Kanitz, A., Wilm, A., Holtgrewe, M., Rahmann, S., Nahnsen, S., Köster, J. (2021) Sustainable data analysis with Snakemake. *F1000Res*, **10**.
- 21 Kurtzer, G.M., Sochat, V. and Bauer, M.W. (2017) Singularity: Scientific containers for mobility of compute. *PLoS One*, **12**, e0177459.
- 22 Bolger, A.M., Lohse, M. and Usadel, B. (2014) Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics*, **30**, 2114-2120.
- 23 , in press.
- 24 Li, B. and Dewey, C.N. (2011) RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC Bioinformatics*, **12**, 323.
- 25 Liu, Z., Gui, T., Wang, Z., Li, H., Fu, Y., Dong, X. and Li, Y. (2016) cisASE: a likelihood-based method for detecting putative cis-regulated allele-specific expression in RNA sequencing data. *Bioinformatics*, **32**, 3291-3297.
- 26 Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G. and Durbin, R. (2009) The Sequence Alignment/Map format and SAMtools. *Bioinformatics*, **25**, 2078-2079.
- 27 Kent, W.J., Sugnet, C.W., Furey, T.S., Roskin, K.M., Pringle, T.H., Zahler, A.M. and Haussler, D. (2002) The human genome browser at UCSC. *Genome Res*, **12**, 996-1006.

Supplementary Figures

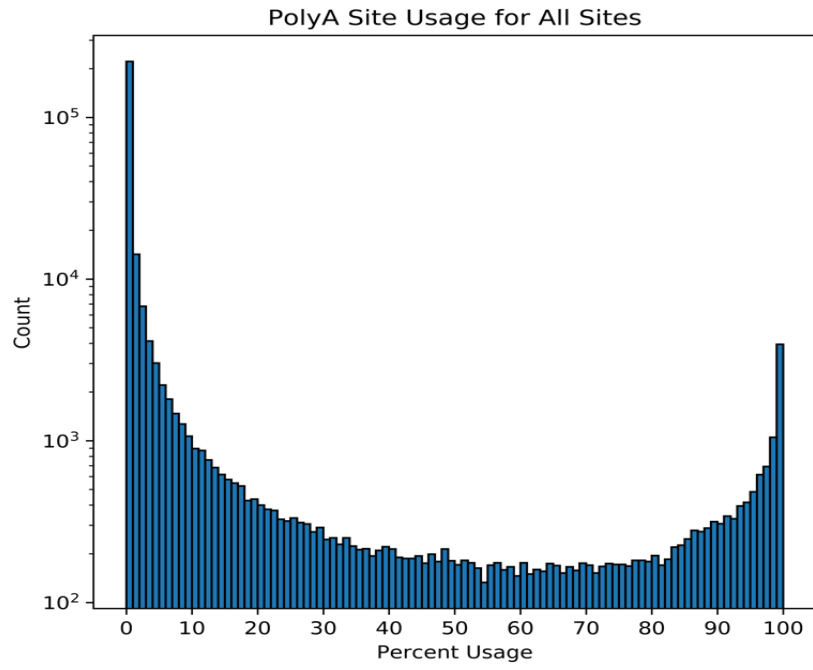


Figure S1. Histogram of usage for all protein-coding PAS

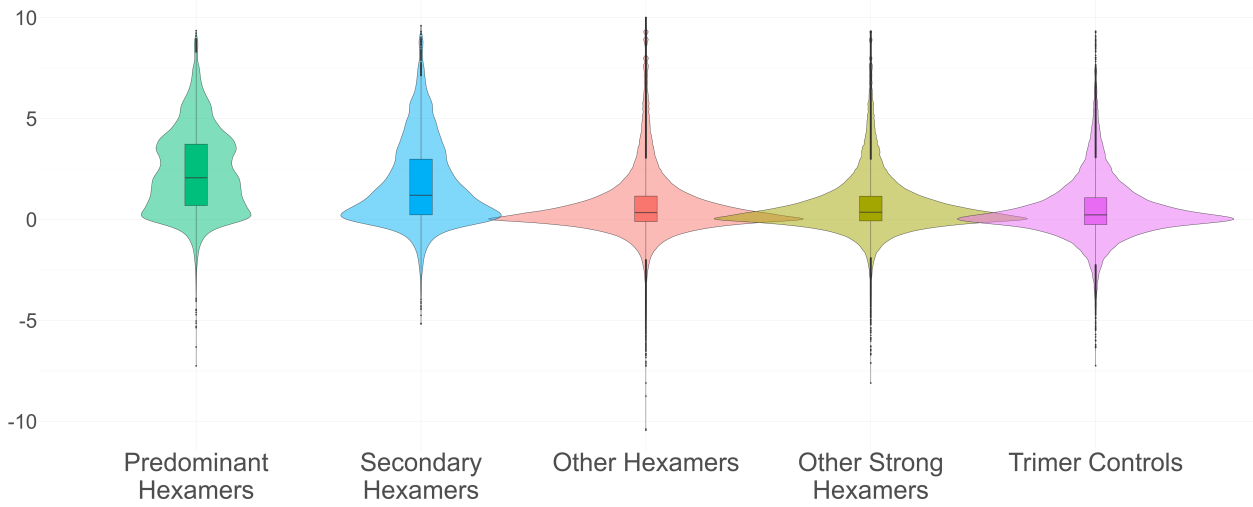


Figure S2: Comparison of PhyloP scores, between predominant hexamers, secondary hexamers, other hexamers, other strong hexamers, and trimer controls.

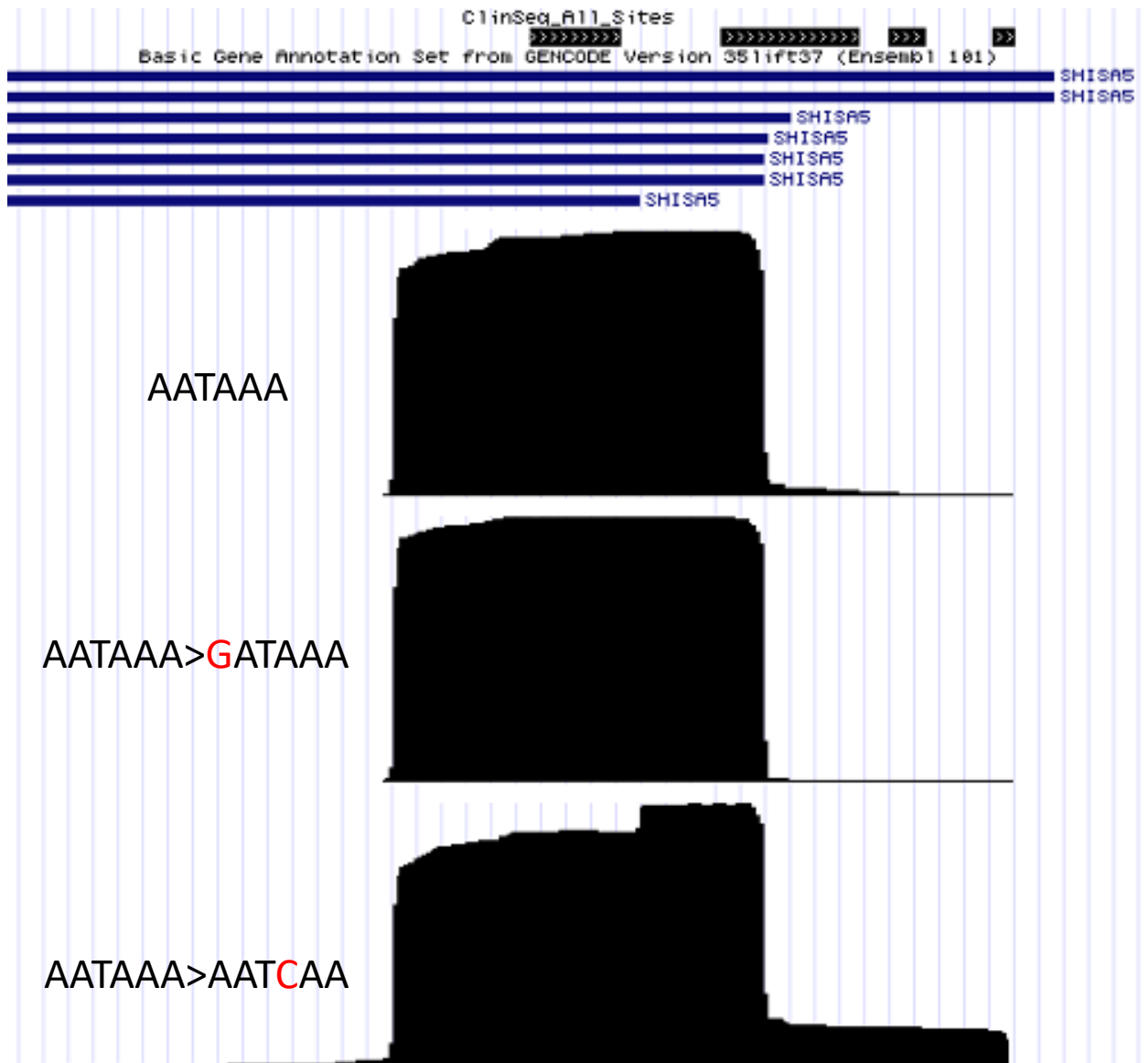


Figure S3. Comparing the effect of different SNVs in the same PAS of gene SHISA5.

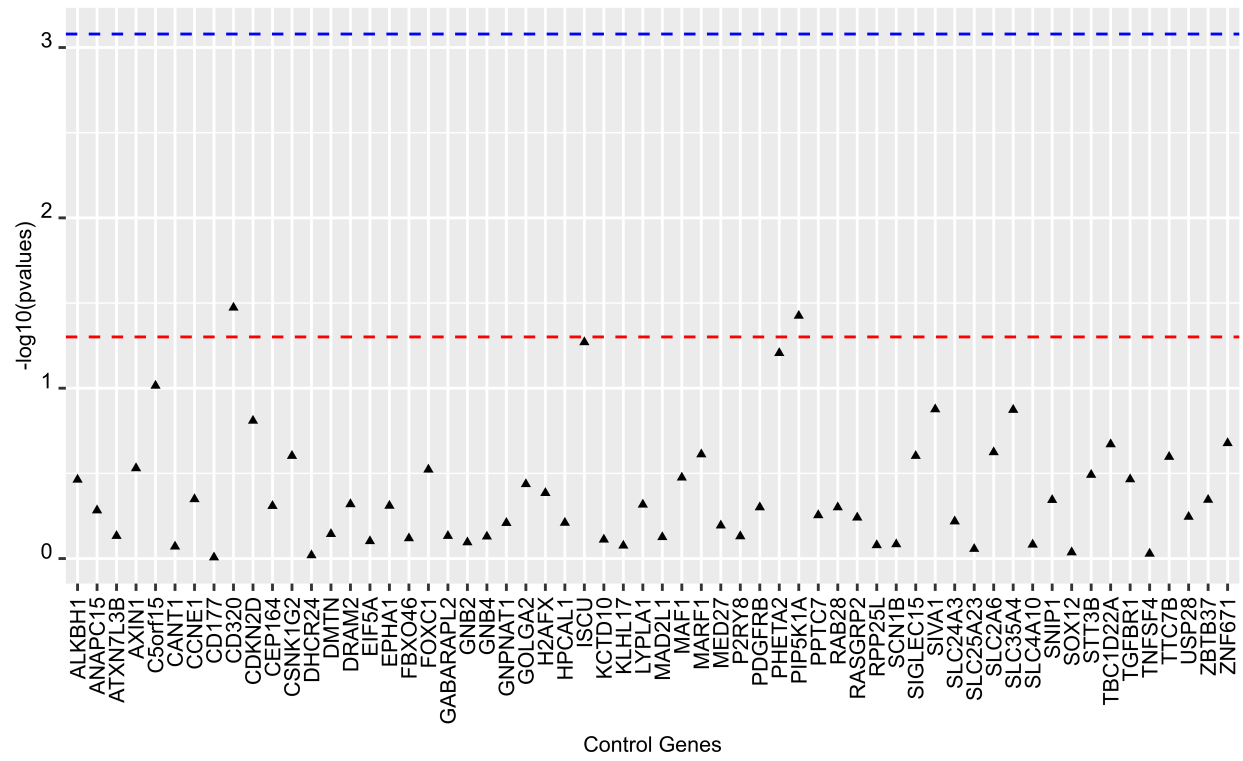


Figure S4. ClinSeq gene expression plot for controls

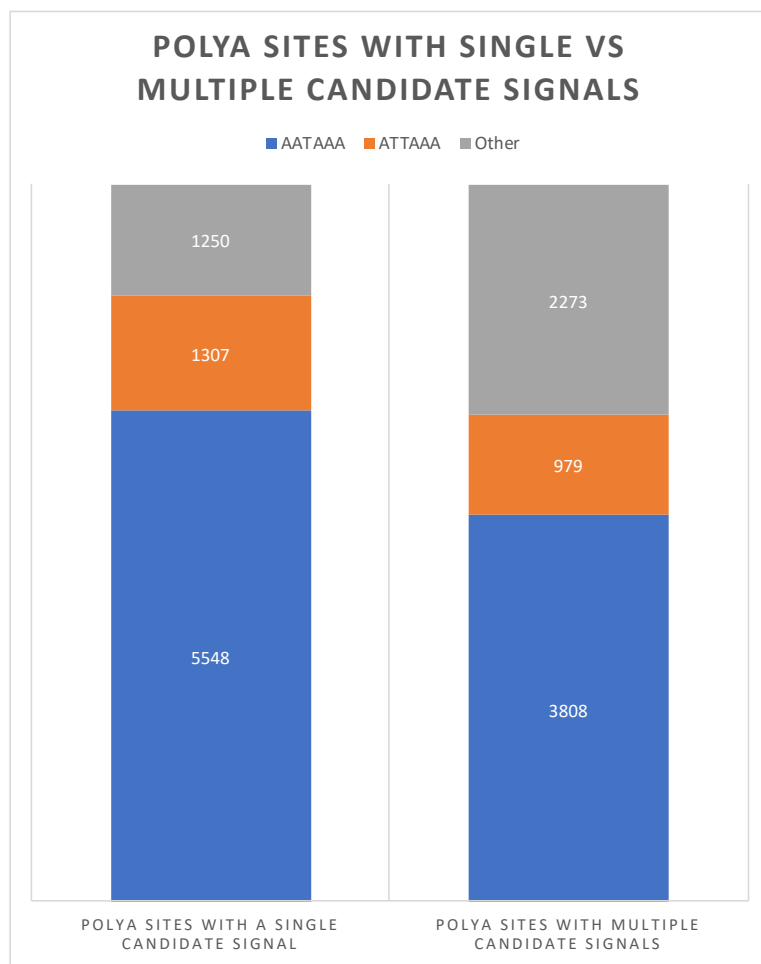


Figure S5. Predominant hexamers sequence in pPAS with single vs multiple candidate signals

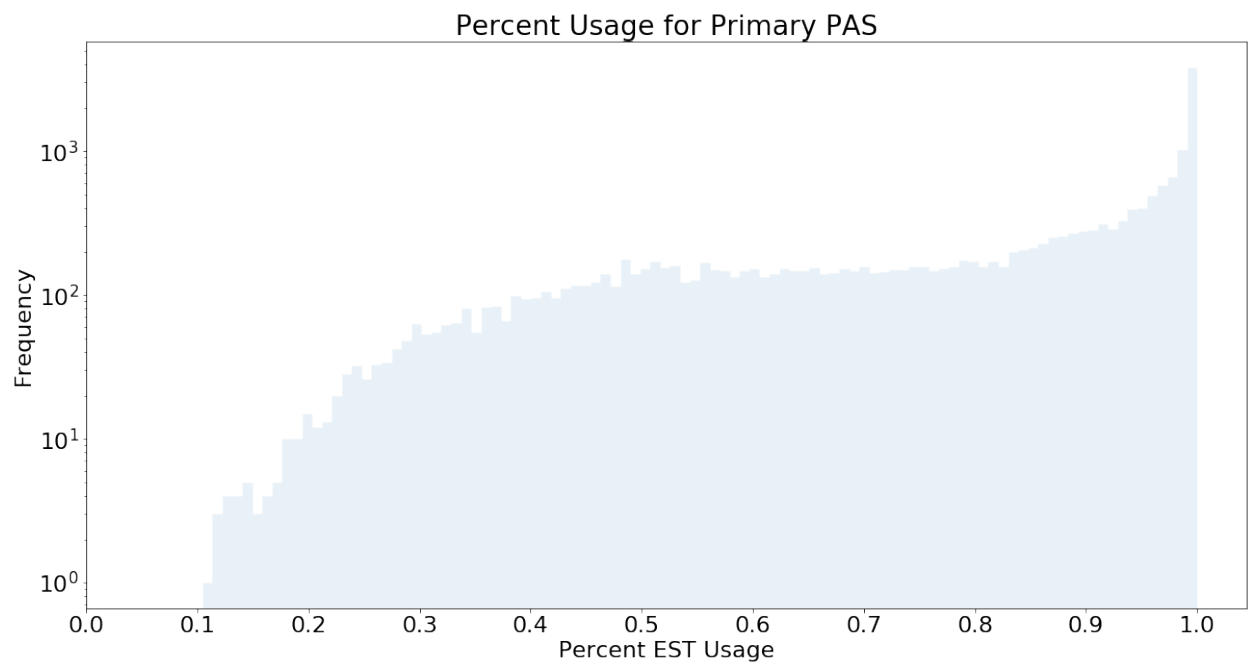


Figure S6. Percent usage of primary sites across all genes in PolyA Site 2.0

Primary and Secondary Site Usage for Genes Without a Predominant PolyA Site

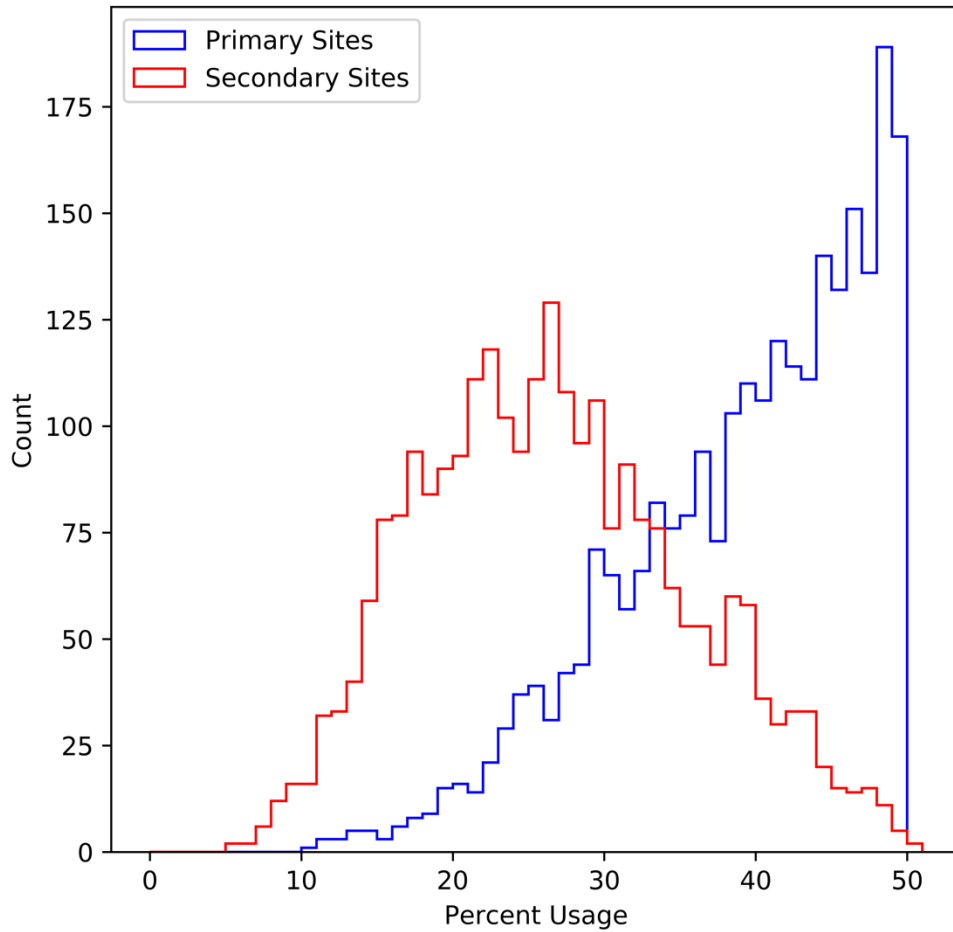


Figure S7. The usage of the primary sites and secondary sites for protein coding genes with no pPAS. Primary PAS is defined as the highest usage site and secondary PAS is defined as the second highest usage PAS.

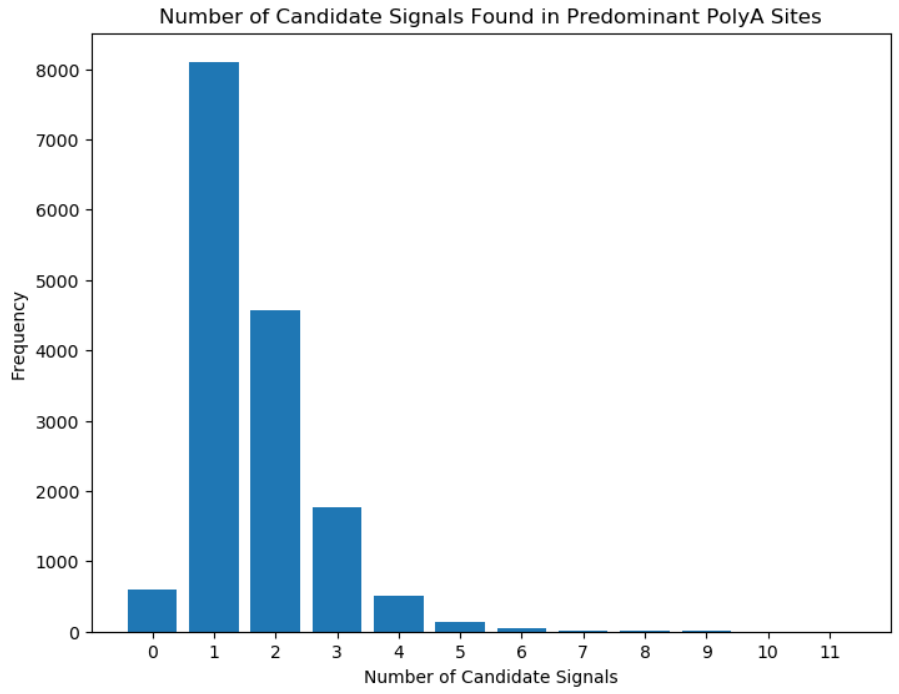


Figure S8. Histogram of number of candidate polyA signals for pPAS

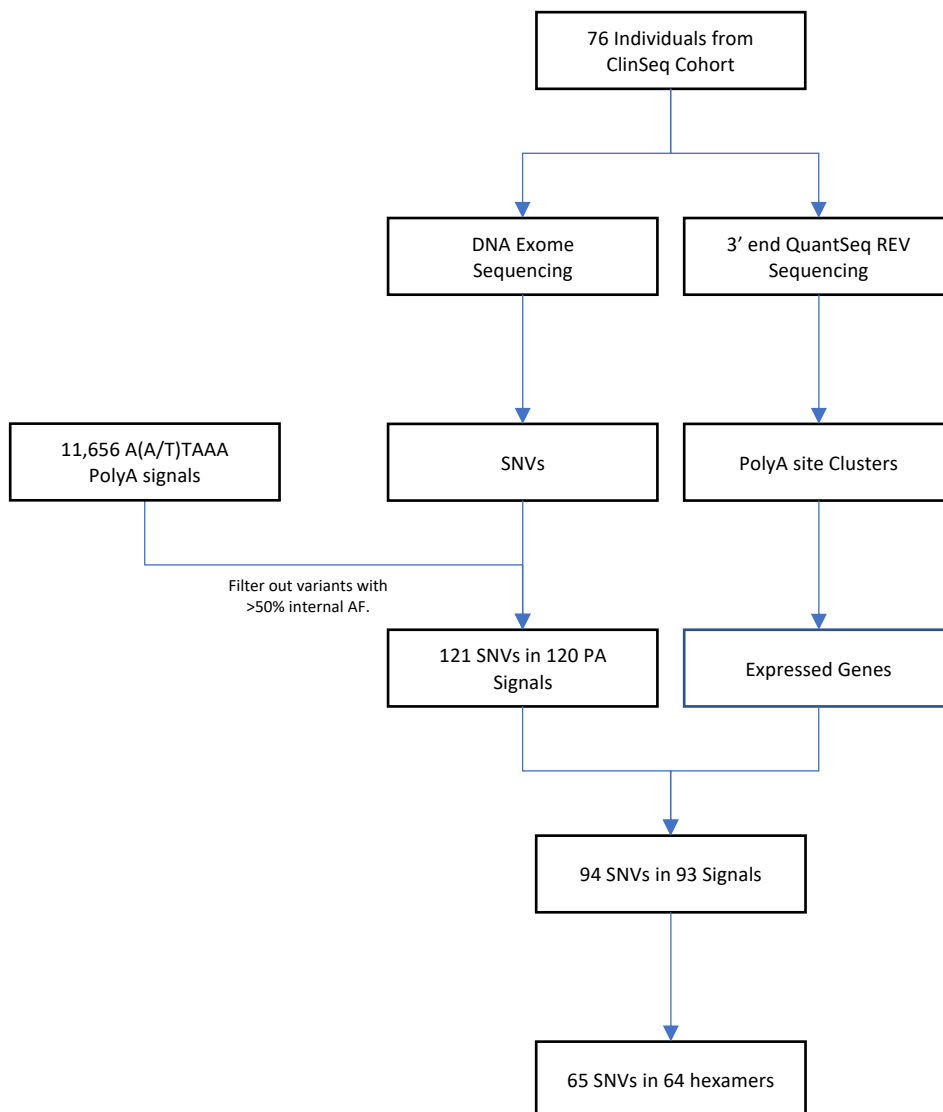


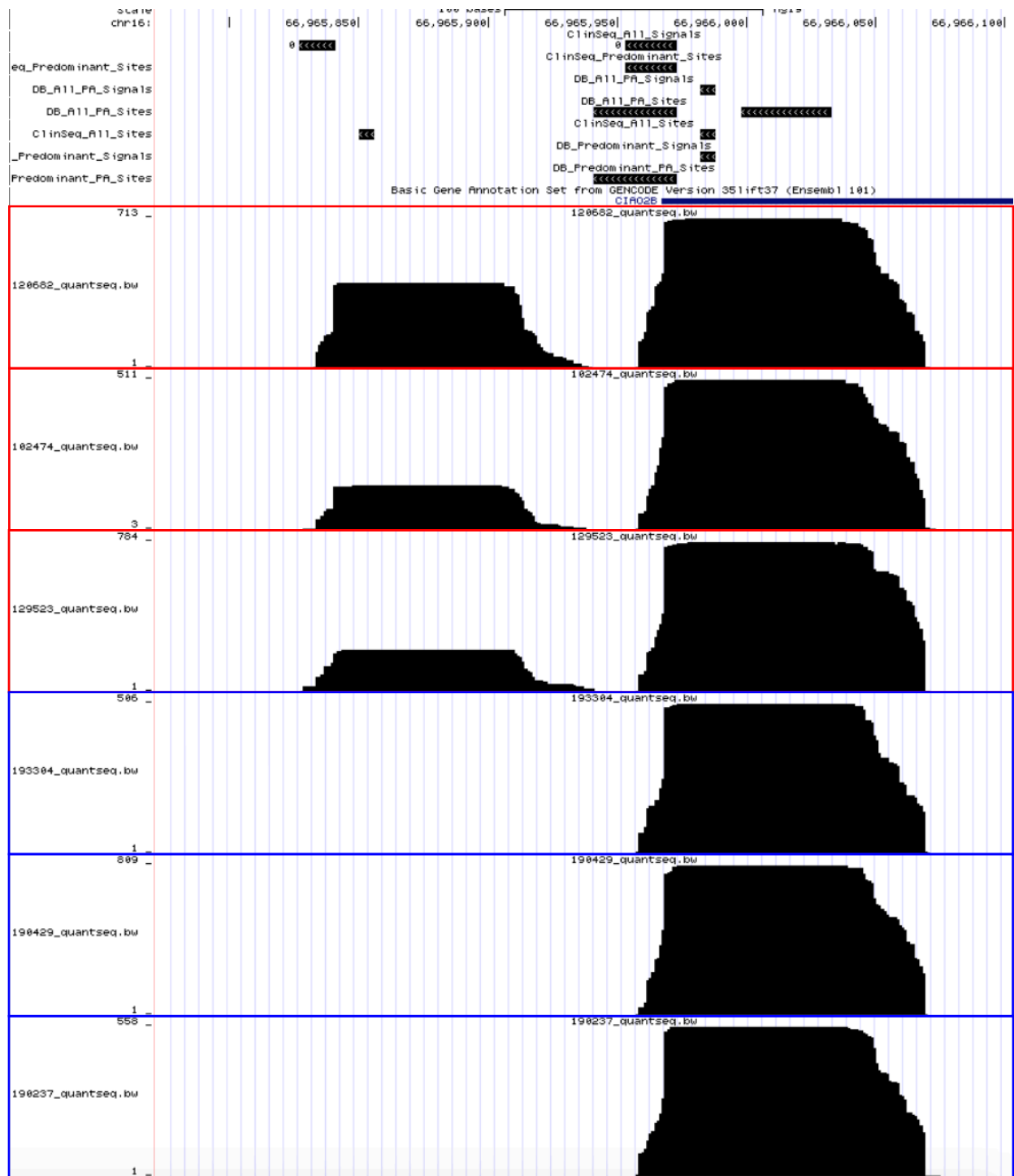
Figure S9. Workflow of assessing polyA signal variants in ClinSeq cohort

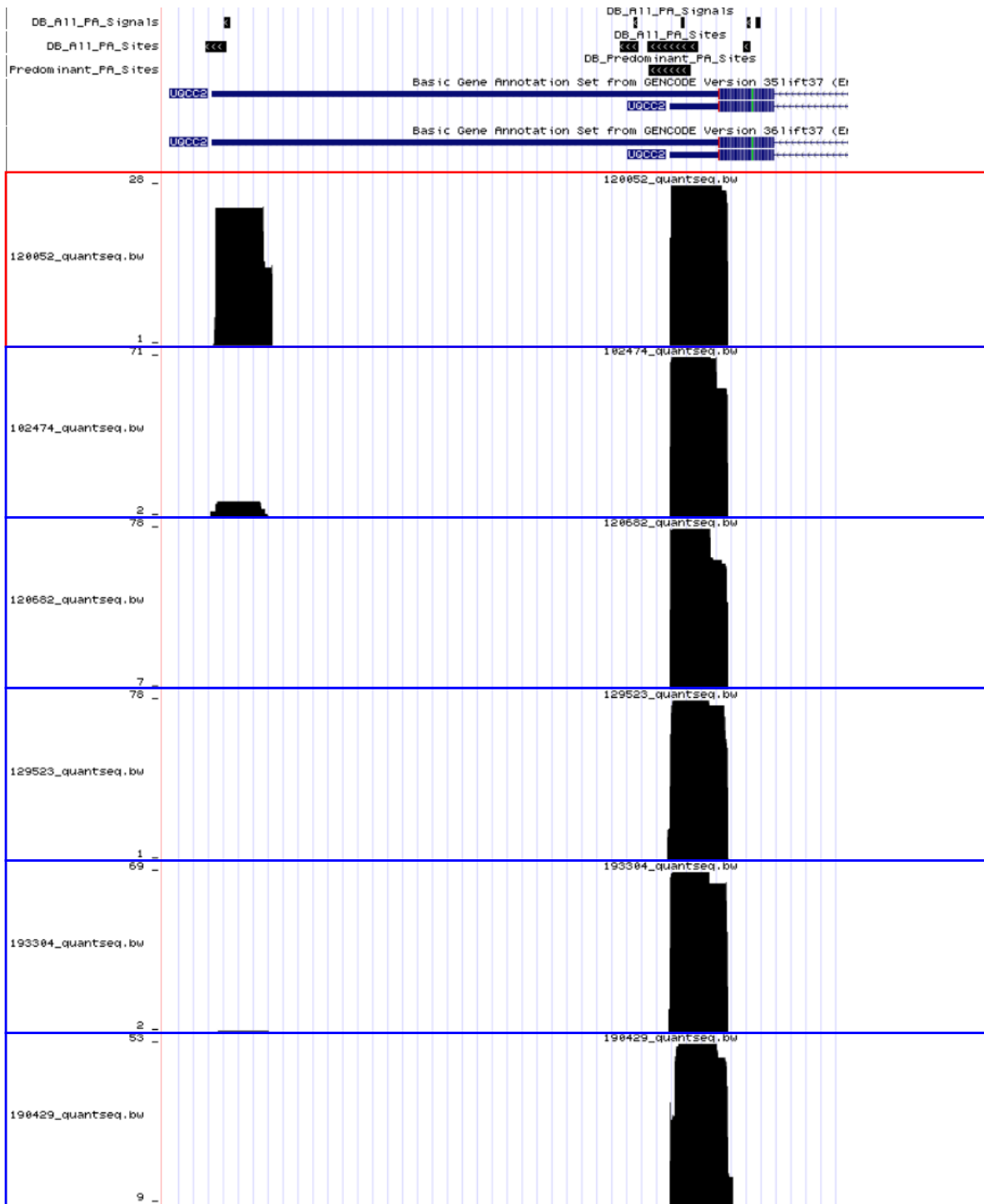


Figure S10. ClinSeq predominant PAS usage by genotypes

Figure S11. ClinSeq Screenshots of all case genes in UCSC Browser.

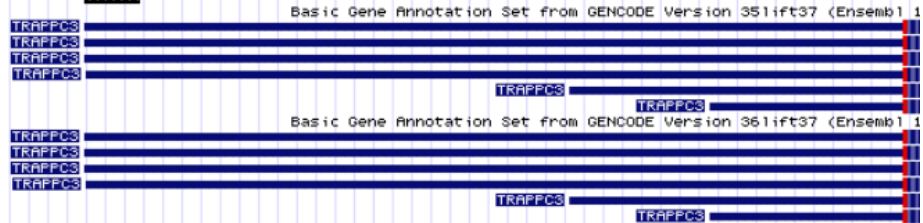
Screenshots of aligned reads from 3' QuantSeq are shown for all 65 variants found in Hexamers. RNA Seq alignments are shown for only for 3 genes where RNA extension is found.





DB_A11_PA_Signals
edominant_PA_Signals
DB_A11_PA_Sites
Predominant_PA_Sites

DB_A11_PA_Signals
DB_Predominant_PA_Signals
DB_A11_PA_Sites
DB_Predominant_PA_Sites



750

122579_quantseq.bw

122579_quantseq.bw

5
697

130689_quantseq.bw

130689_quantseq.bw

1
505

133170_quantseq.bw

133170_quantseq.bw

1
422

116964_quantseq.bw

116964_quantseq.bw

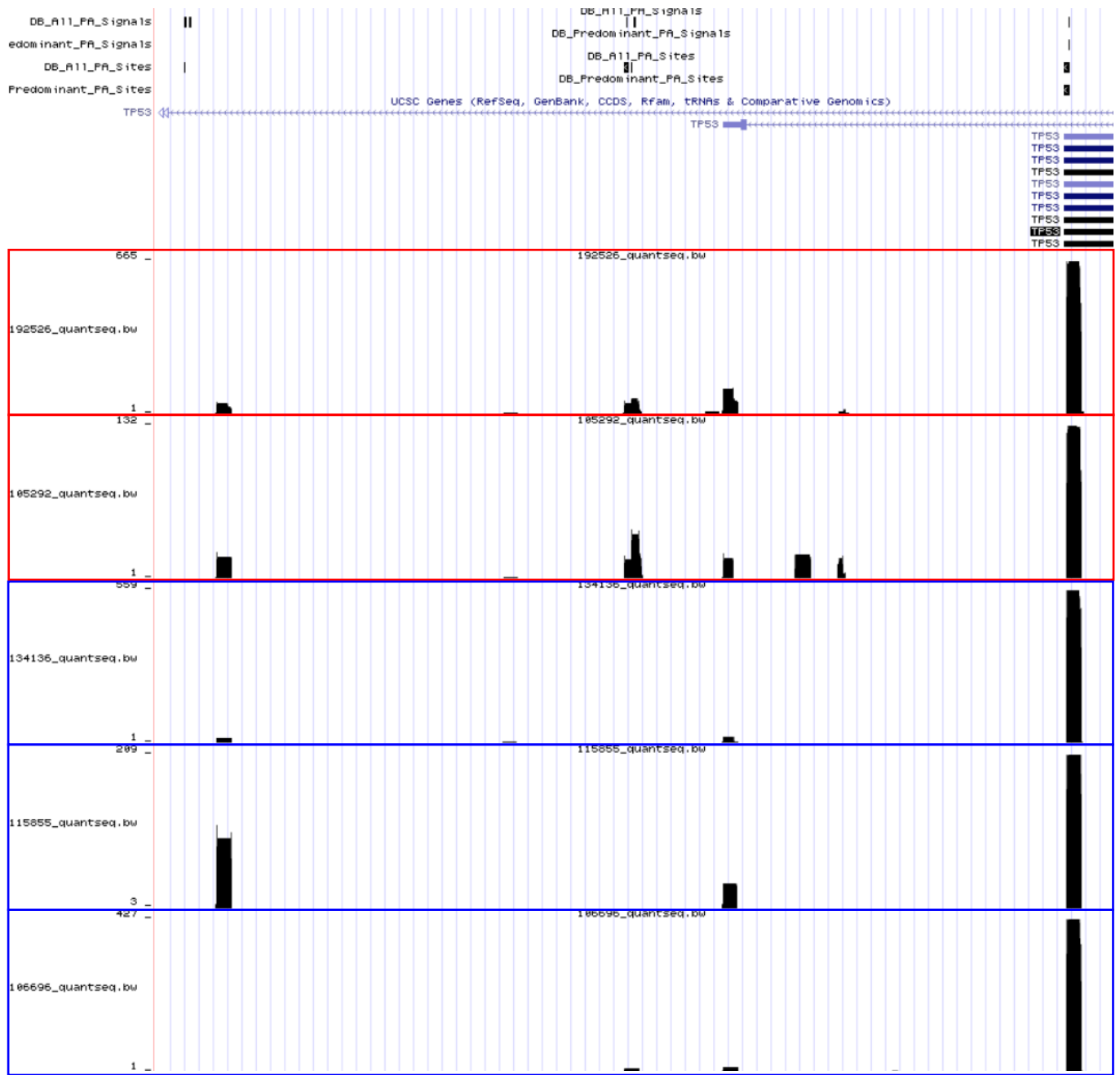
3
983

173071_quantseq.bw

173071_quantseq.bw

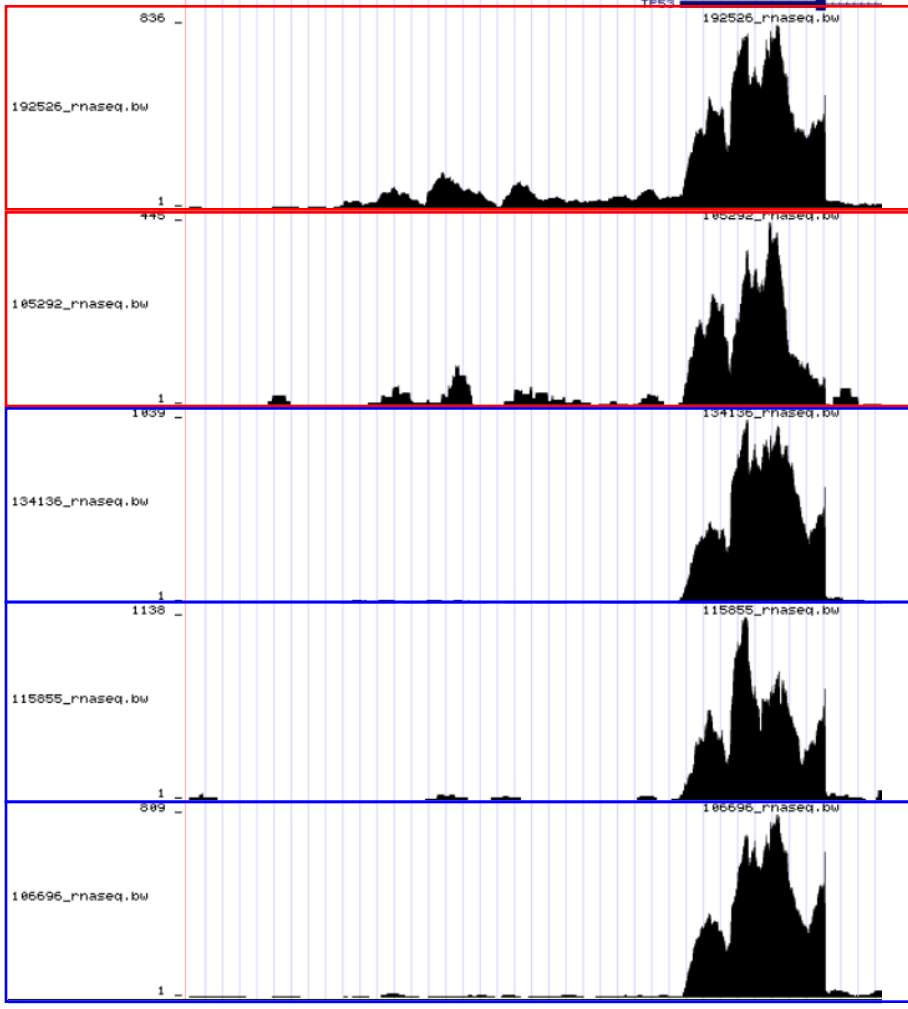
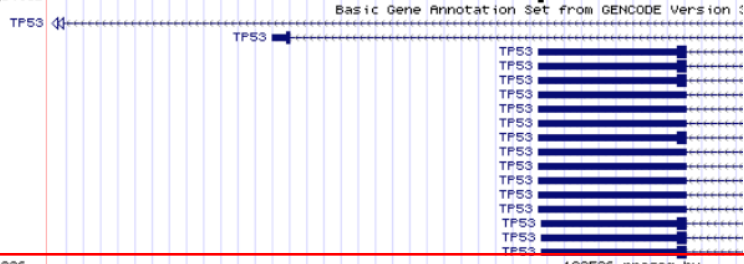
1

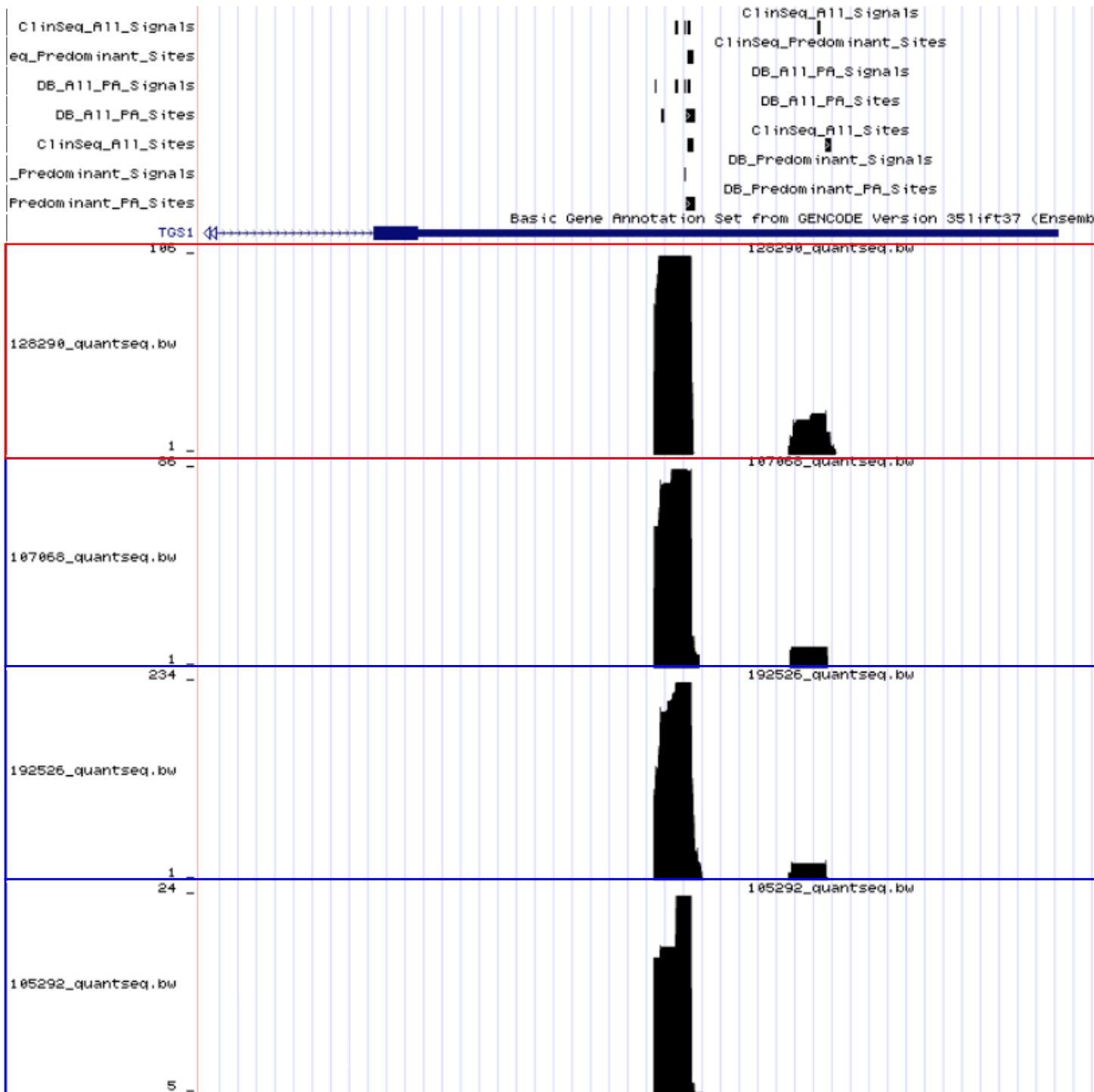


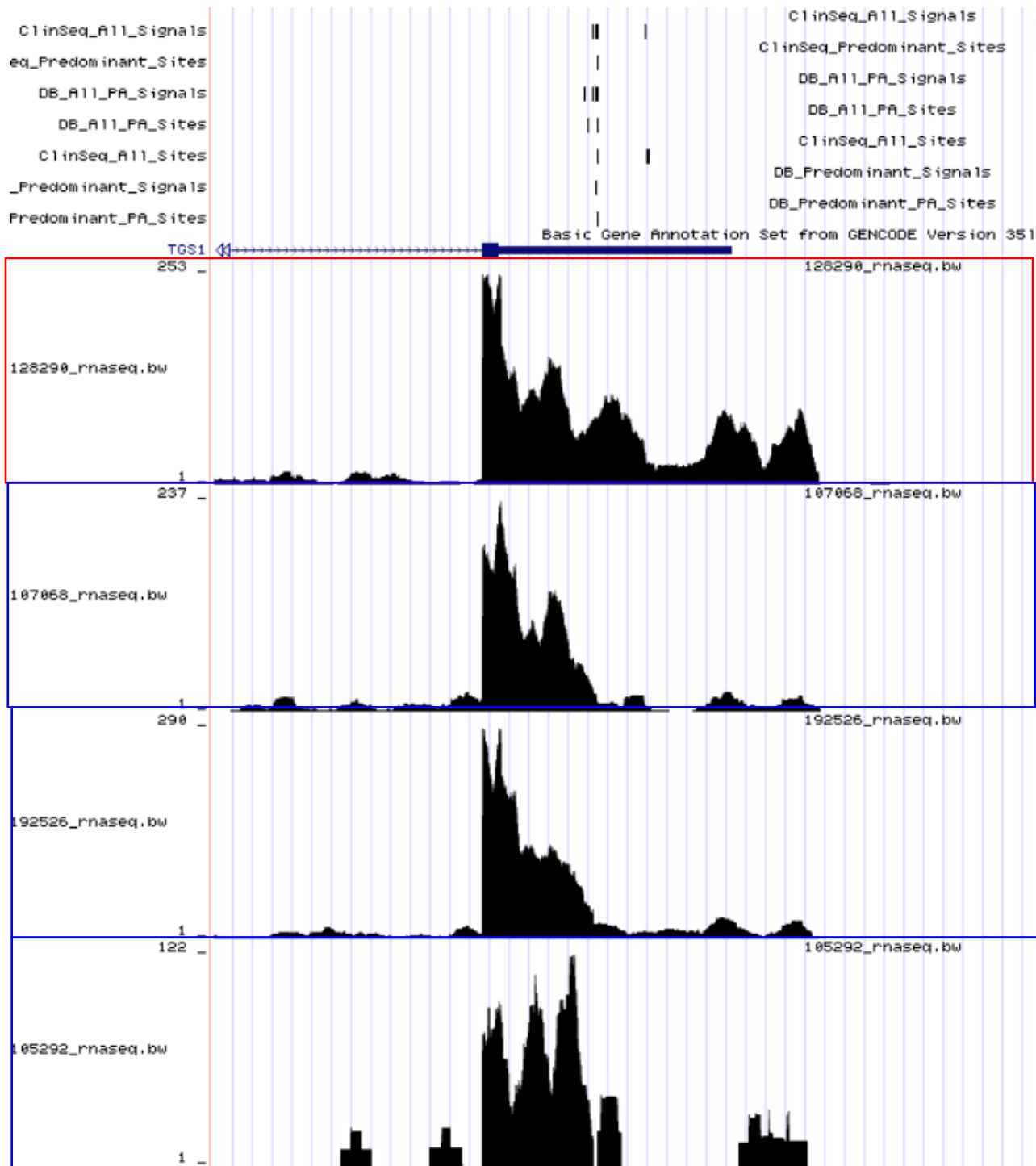


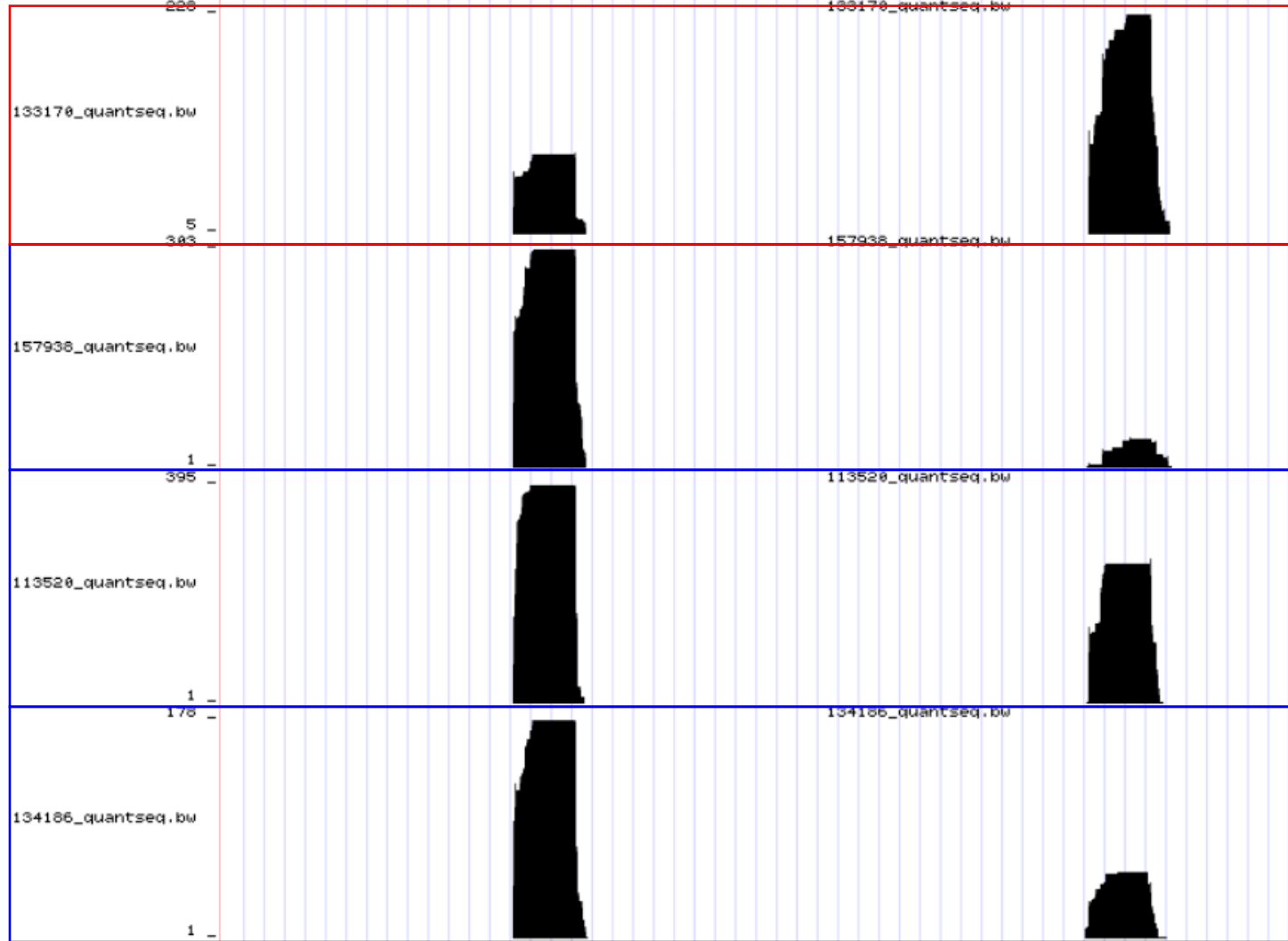
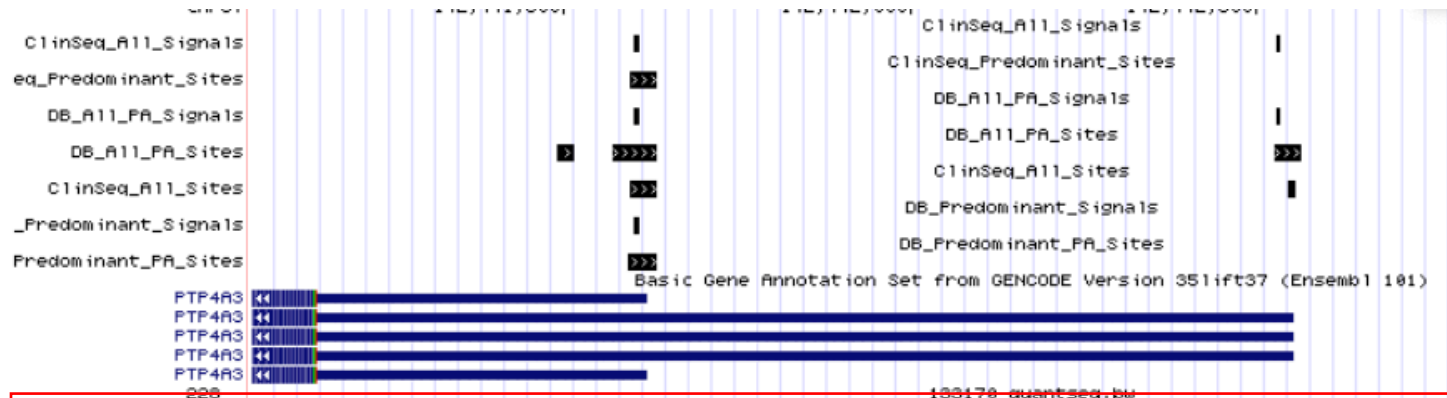
ClinSeq_A11_Signals
 eq_Predominant_Sites
 DB_A11_PA_Signals
 DB_A11_PA_Sites
 ClinSeq_A11_Sites
 _Predominant_Signals
 Predominant_PA_Sites

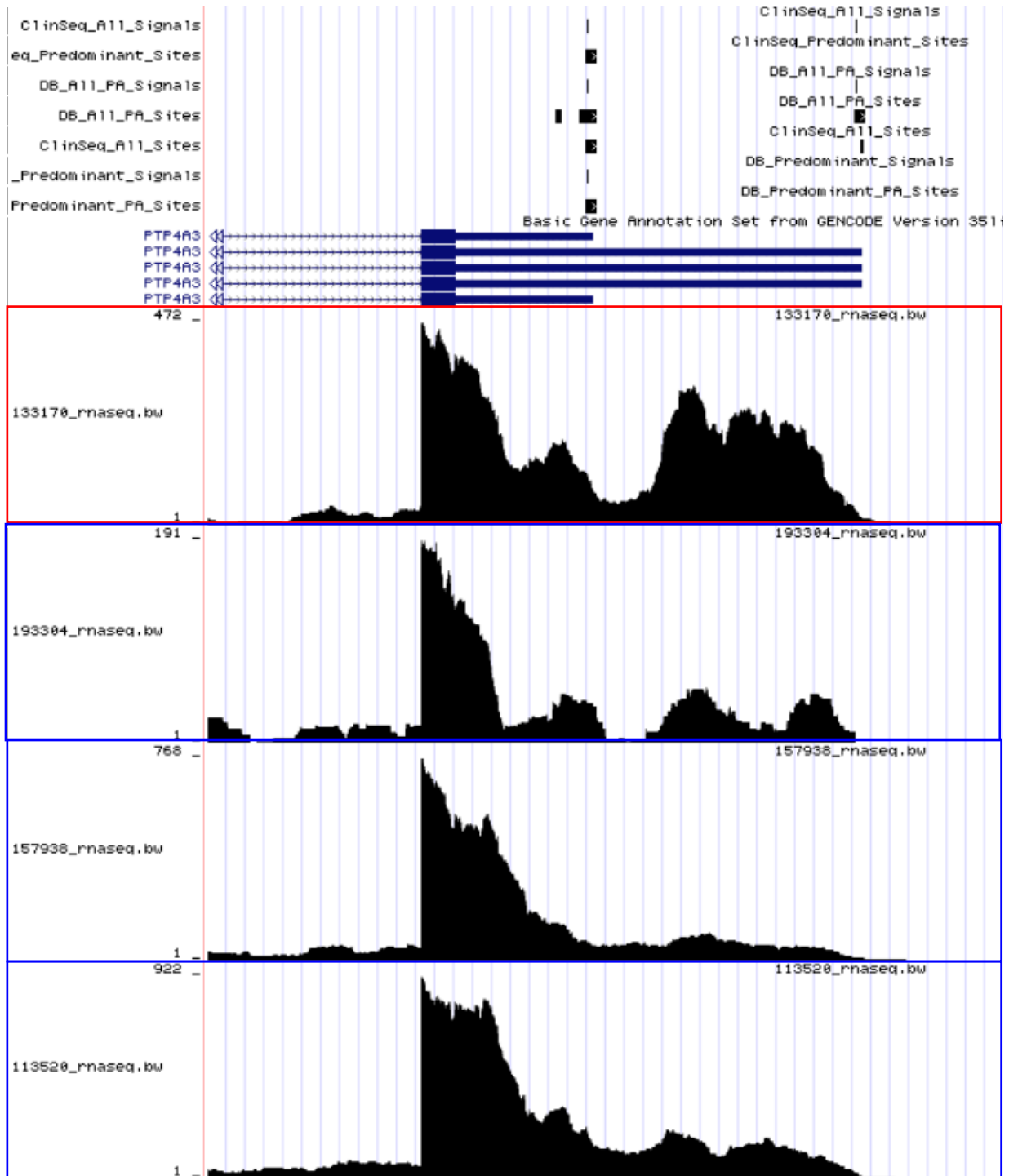
ClinSeq_A11_Signals
 ClinSeq_Predominant_Sites
 DB_A11_PA_Signals
 DB_A11_PA_Sites
 ClinSeq_A11_Sites
 DB_Predominant_Signals
 DB_Predominant_PA_Sites







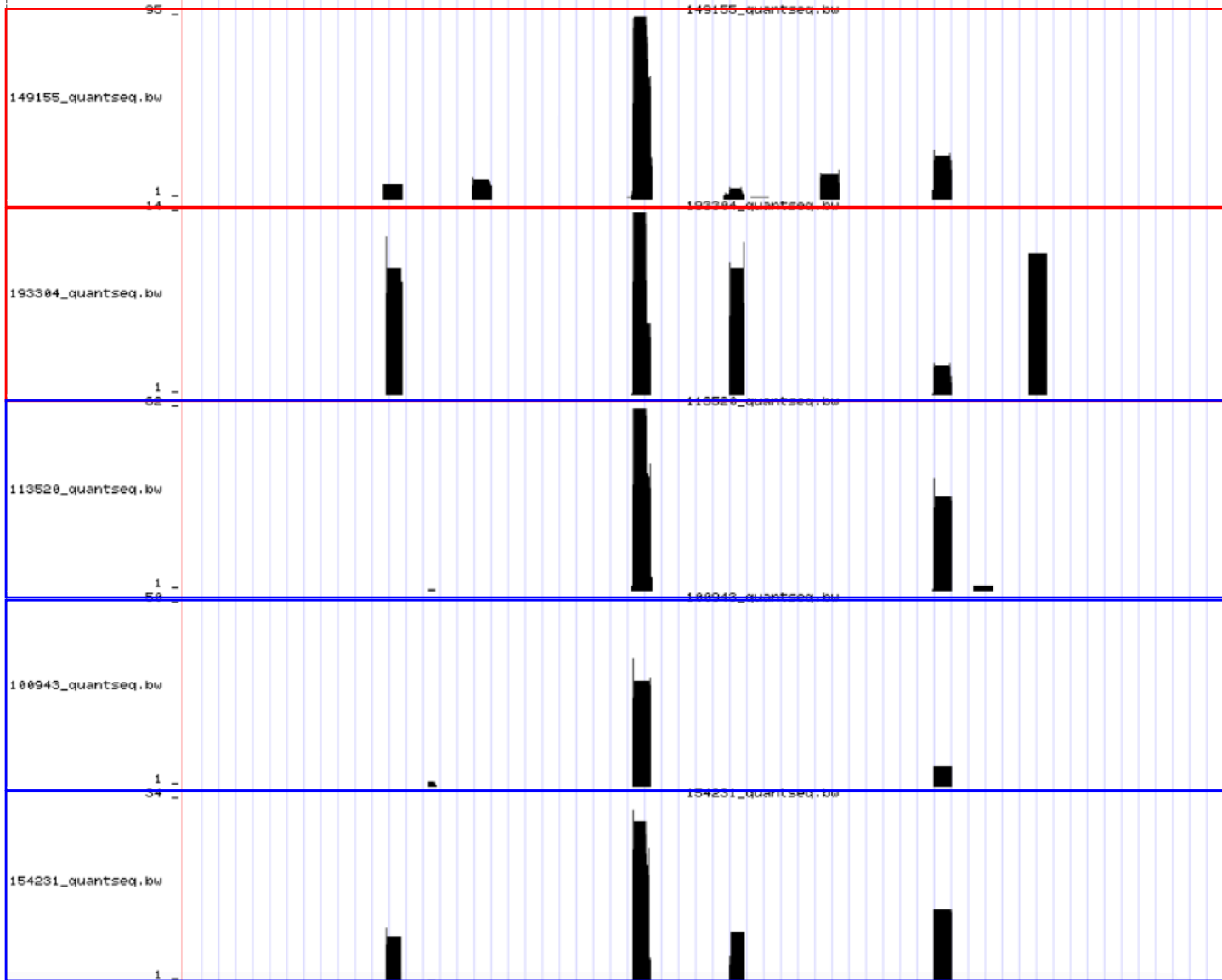


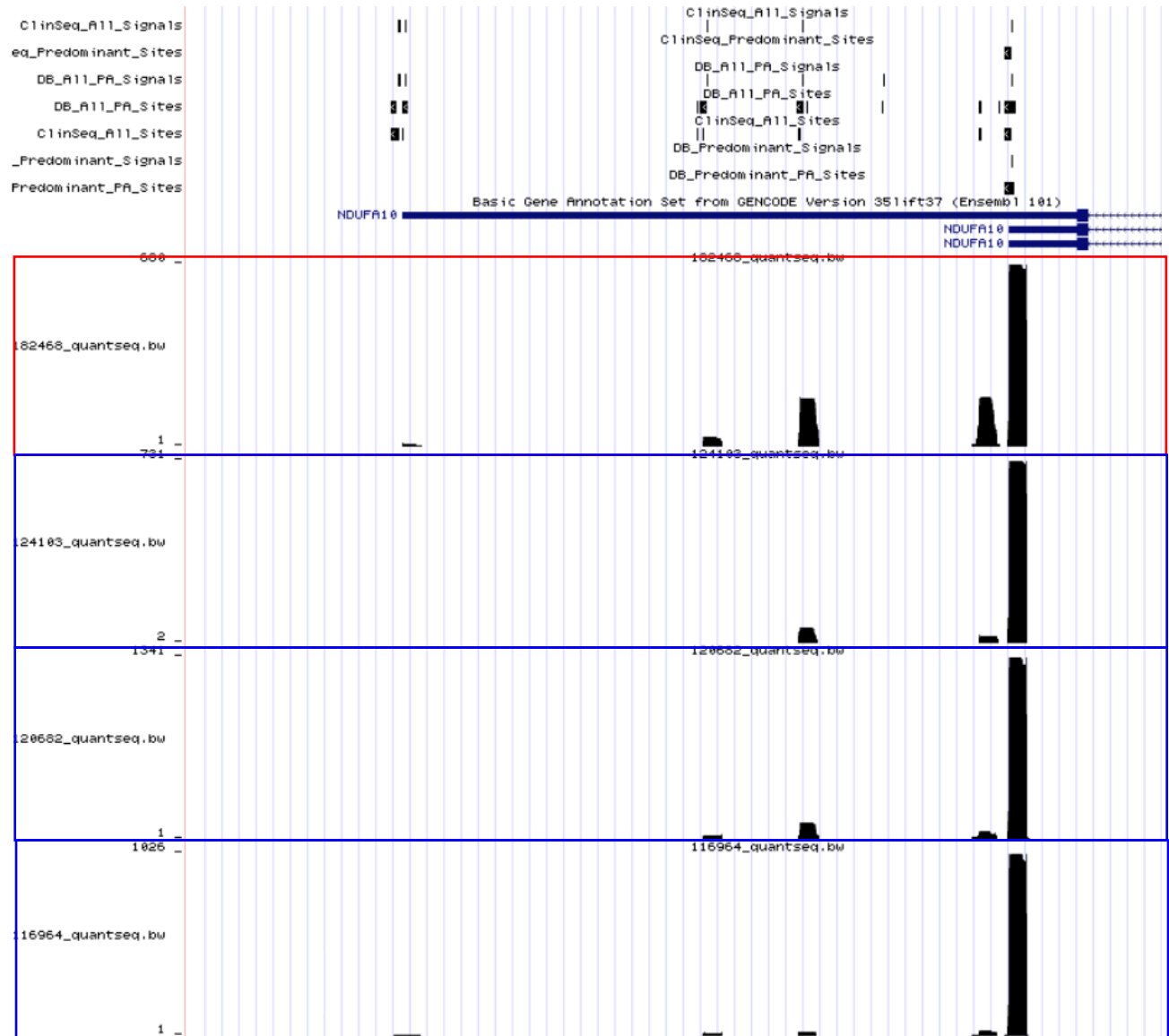


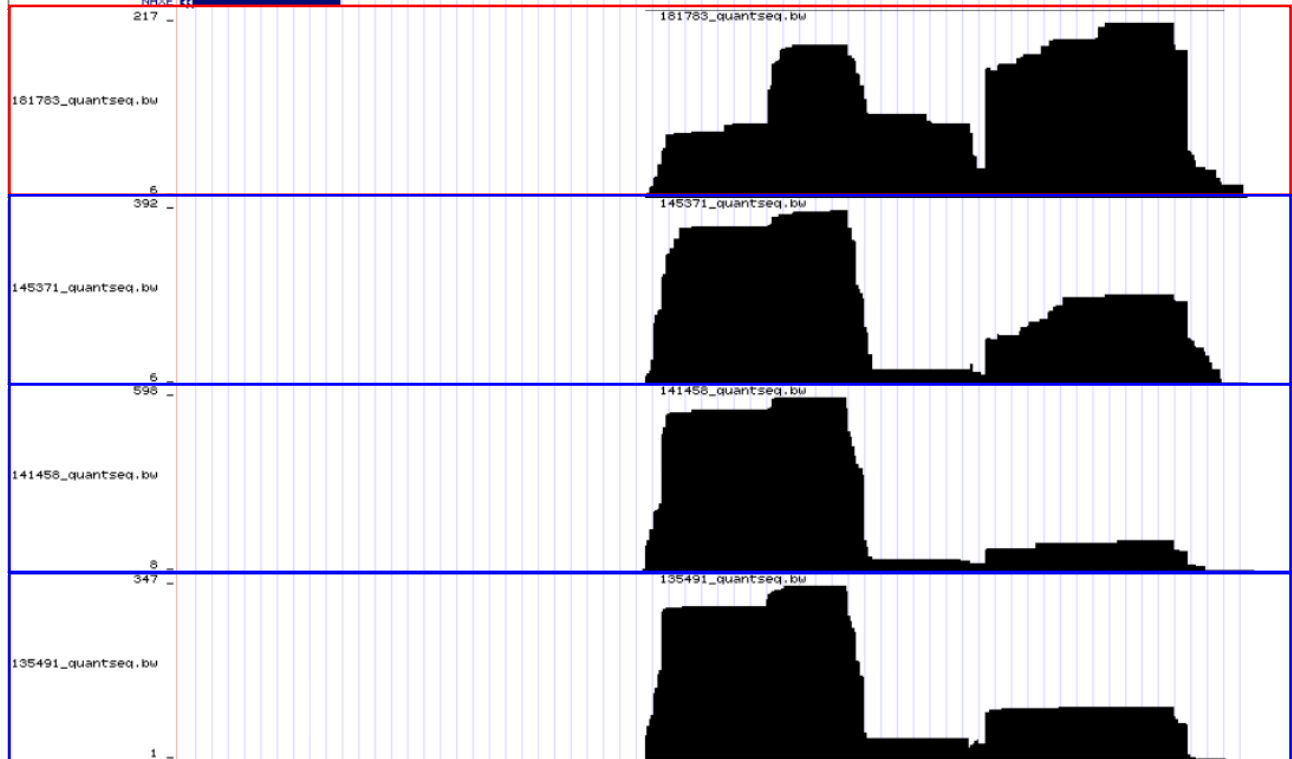
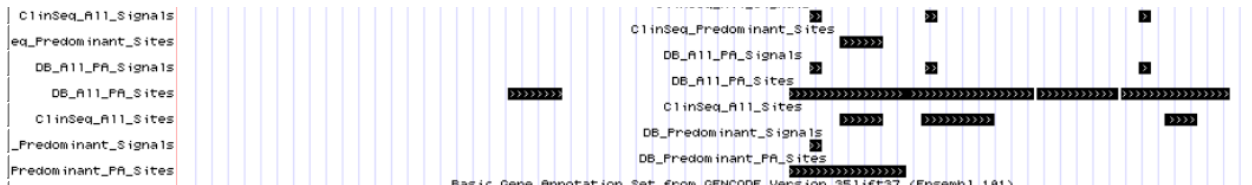
ClinSeq_A11_Signals
 eq_Predominant_Sites
 DB_A11_FA_Signals
 DB_A11_FA_Sites
 ClinSeq_A11_Sites
 _Predominant_Signals
 Predominant_FA_Sites

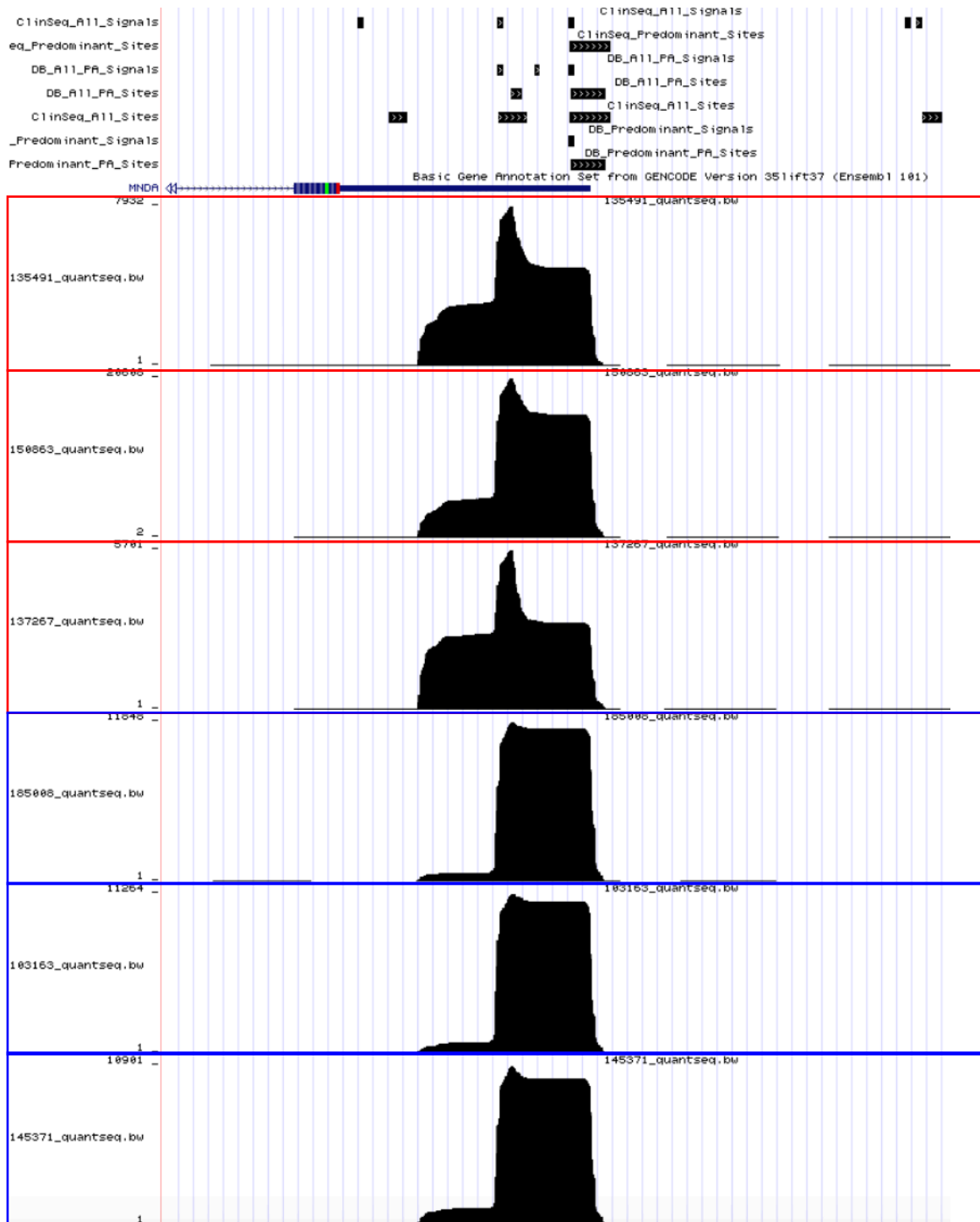
ClinSeq_A11_Signals
 ClinSeq_Predominant_Sites
 DB_A11_FA_Signals
 DB_A11_FA_Sites
 ClinSeq_A11_Sites
 DB_Predominant_Signals
 DB_Predominant_FA_Sites

Basic Gene Annotation Set from GENCODE Version 35 lift37 (Ensembl 101)
 NUP155







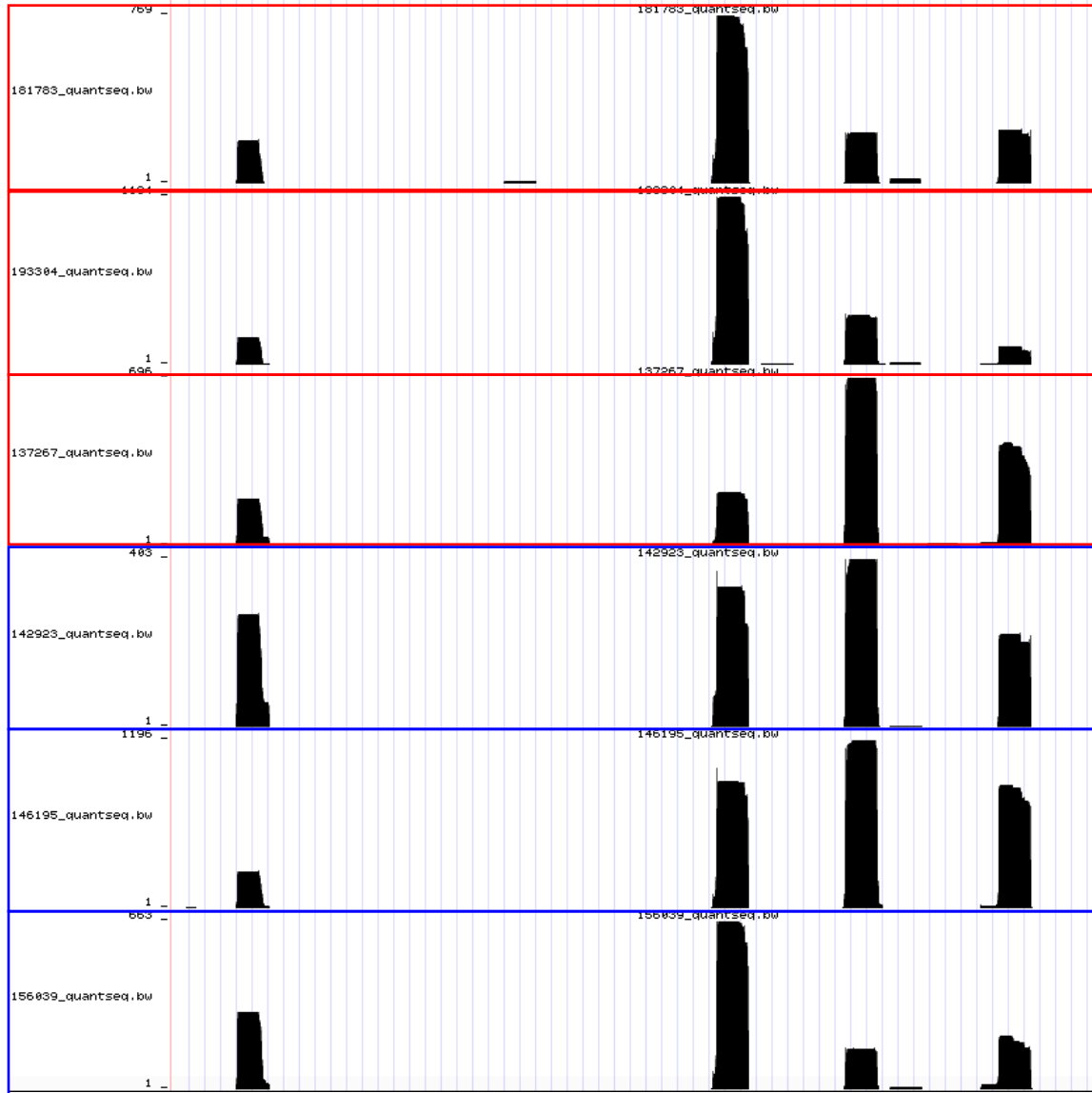


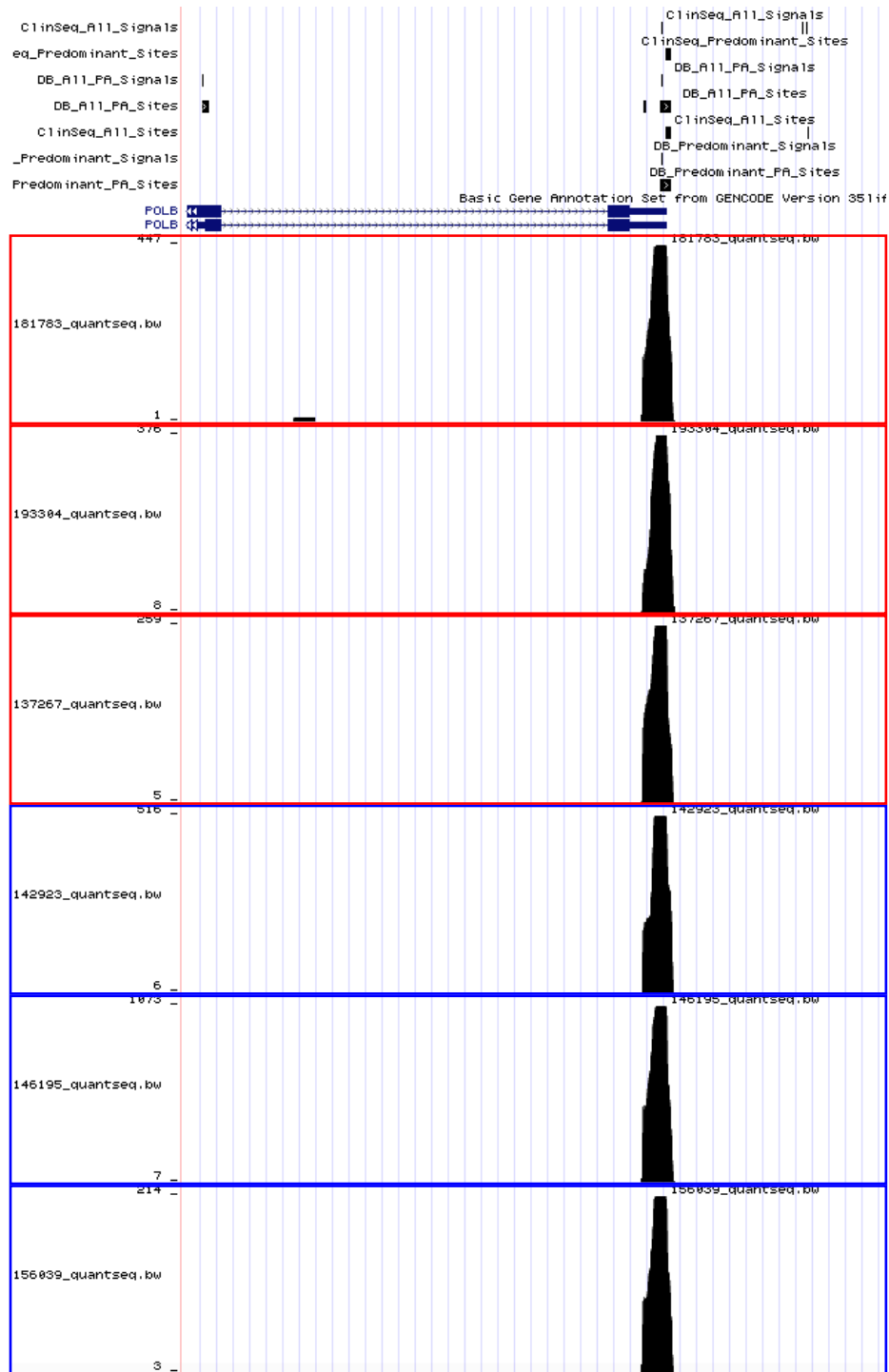
ClinSeq_A11_Signals
ed_Predominant_Sites
DB_A11_PA_Signals
DB_A11_PA_Sites
ClinSeq_A11_Sites
_Predominant_Signals
Predominant_PA_Sites

Basic Gene Annotation Set from GENCODE Version 35 (11/13/17) (Ensembl 101)

ClinSeq_A11_Signals
ClinSeq_Predominant_Sites
DB_A11_PA_Signals
DB_A11_PA_Sites
ClinSeq_A11_Sites
DB_Predominant_Signals
DB_Predominant_PA_Sites

MS4A6A
MS4A6A
MS4A6A
MS4A6A
MS4A6A
MS4A6A

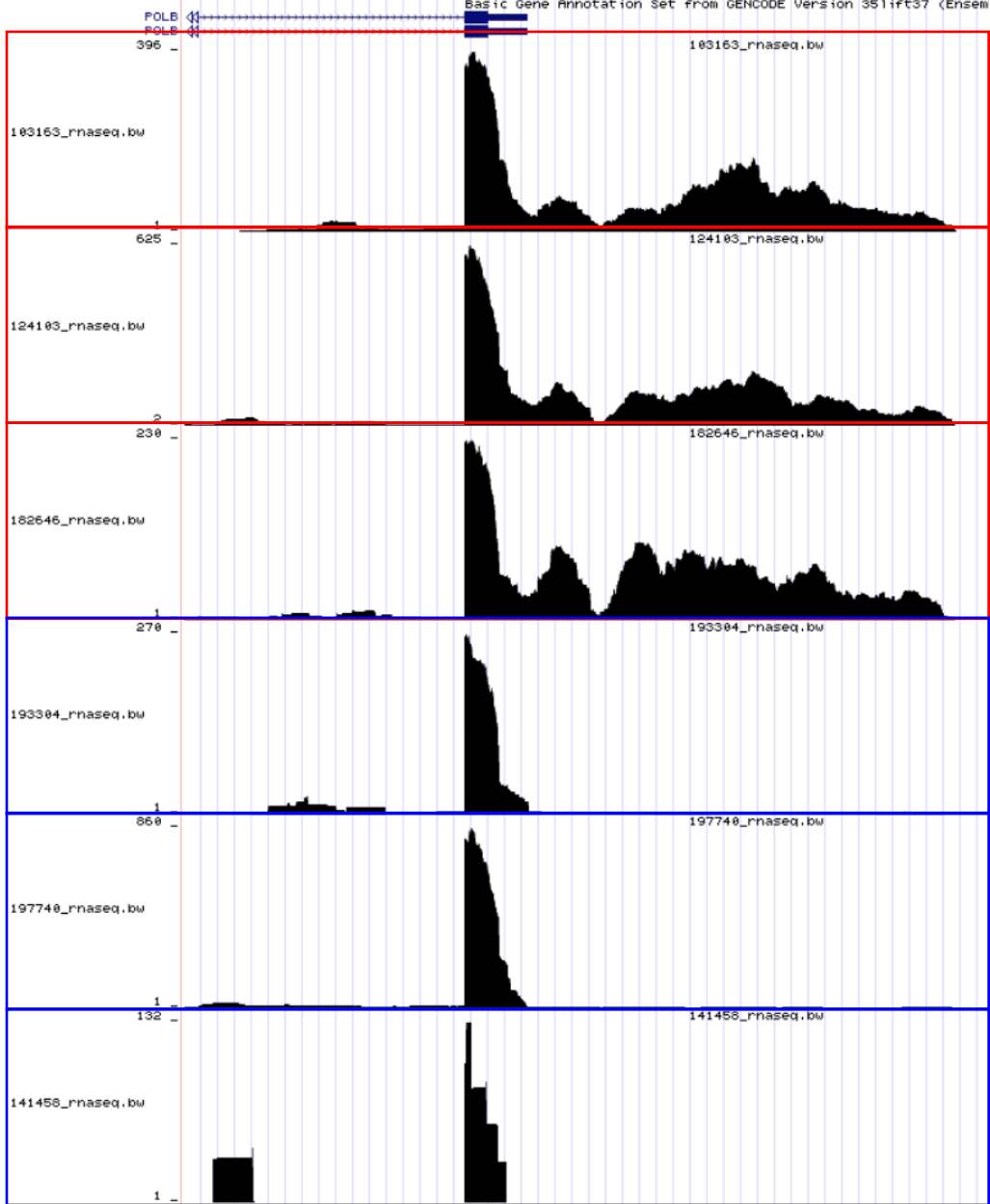


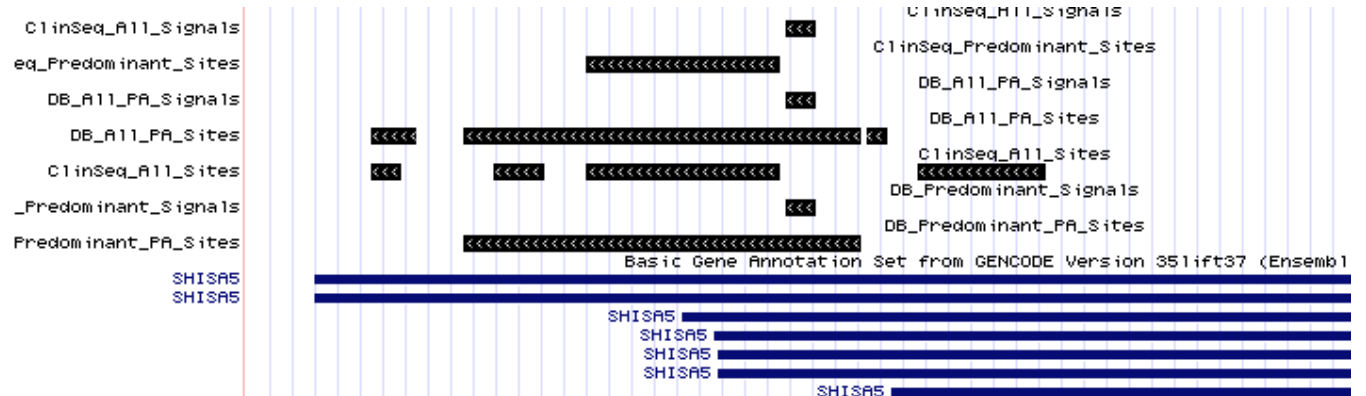


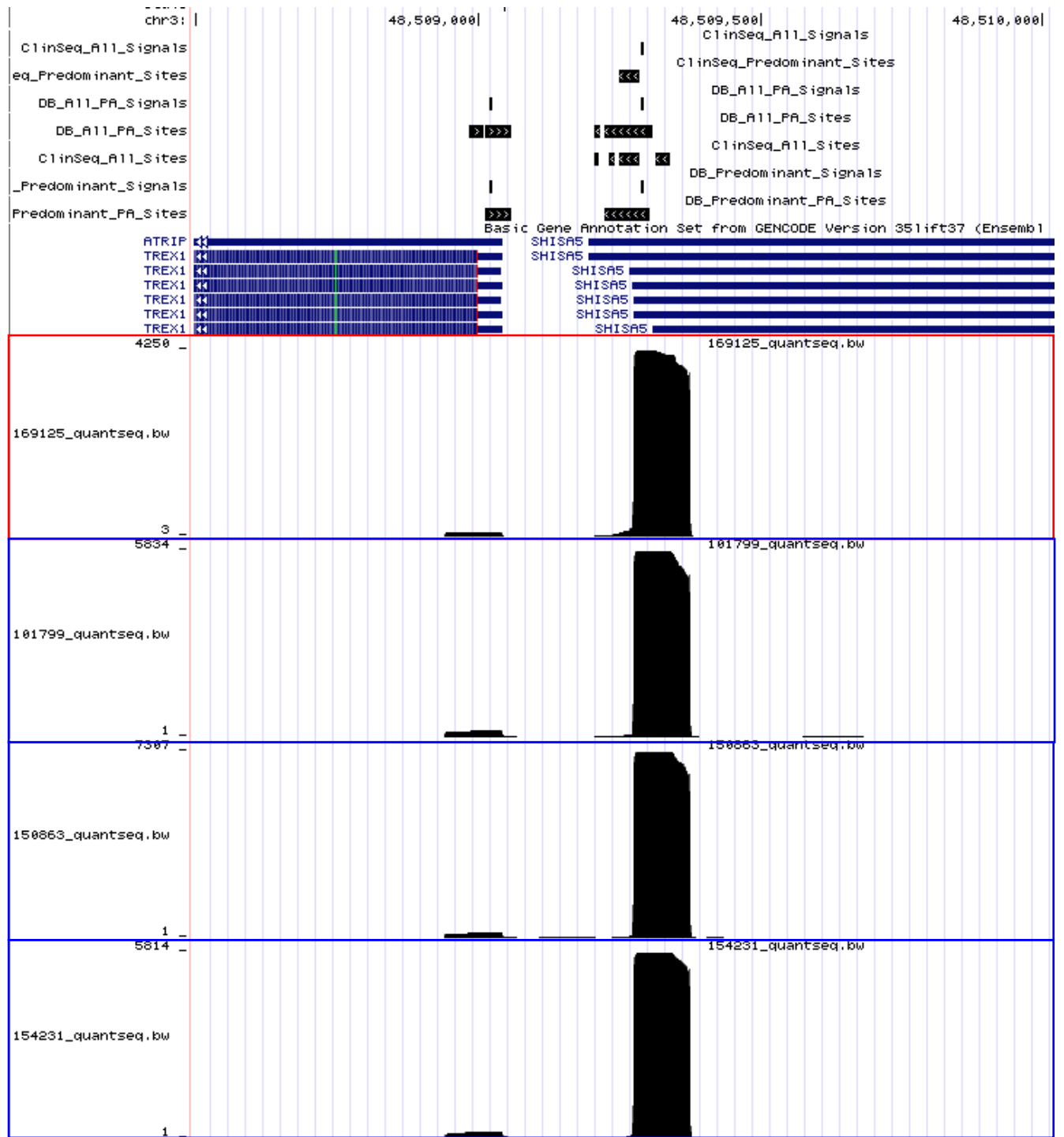
ClinSeq_A11_Signals
 ed_Predominant_Sites
 DB_A11_PA_Signals
 DB_A11_PA_Sites
 ClinSeq_A11_Sites
 _Predominant_Signals
 Predominant_PA_Sites

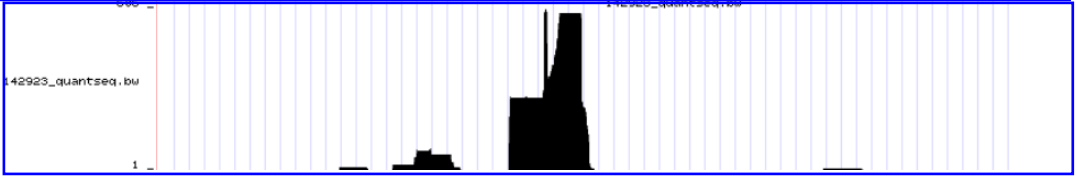
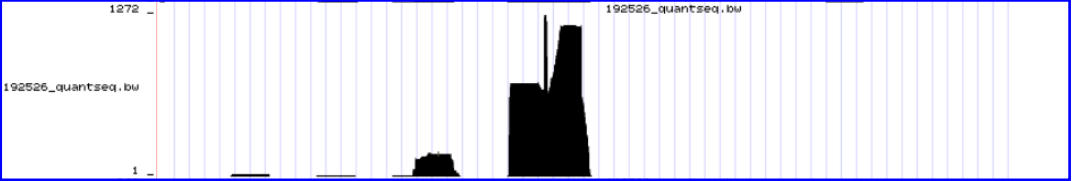
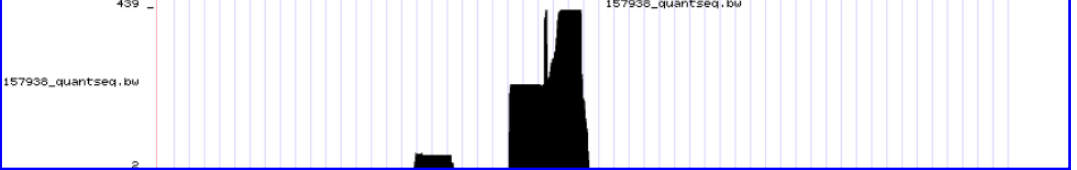
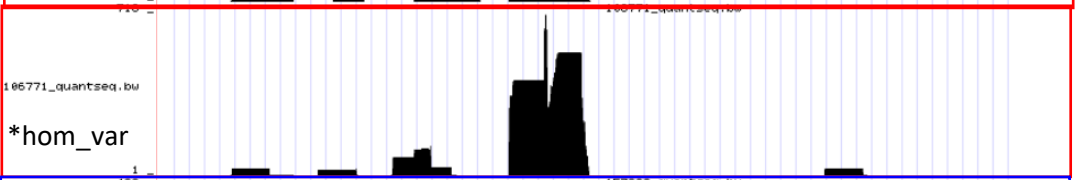
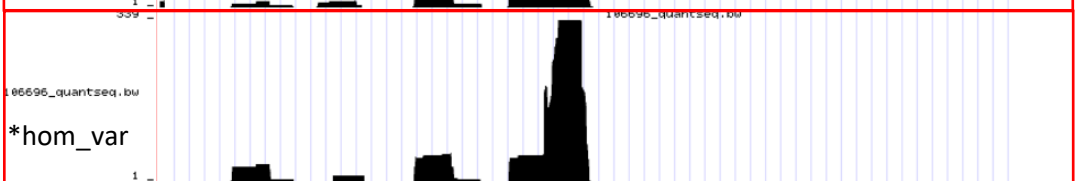
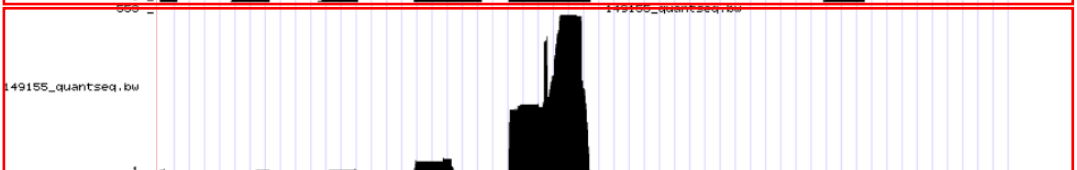
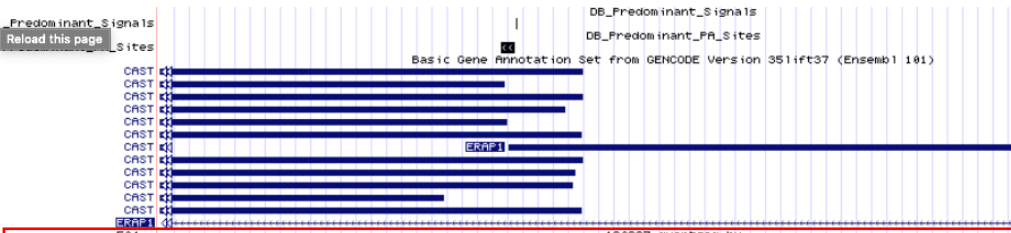
Basic Gene Annotation Set from GENCODE Version 35lift37 (Ensembl)

ClinSeq_A11_Signals
 ClinSeq_Predominant_Sites
 DB_A11_PA_Signals
 DB_A11_PA_Sites
 ClinSeq_A11_Sites
 DB_Predominant_Signals
 DB_Predominant_PA_Sites



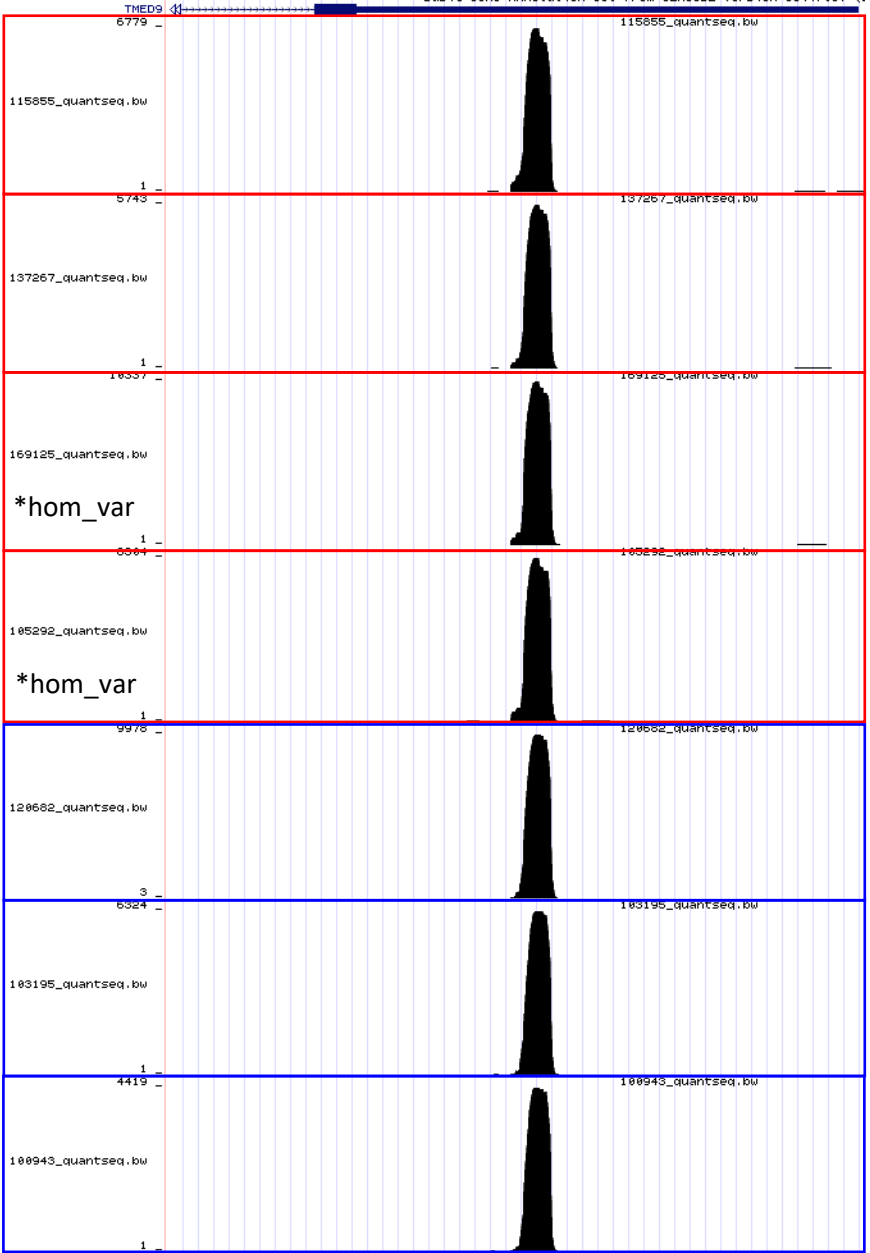


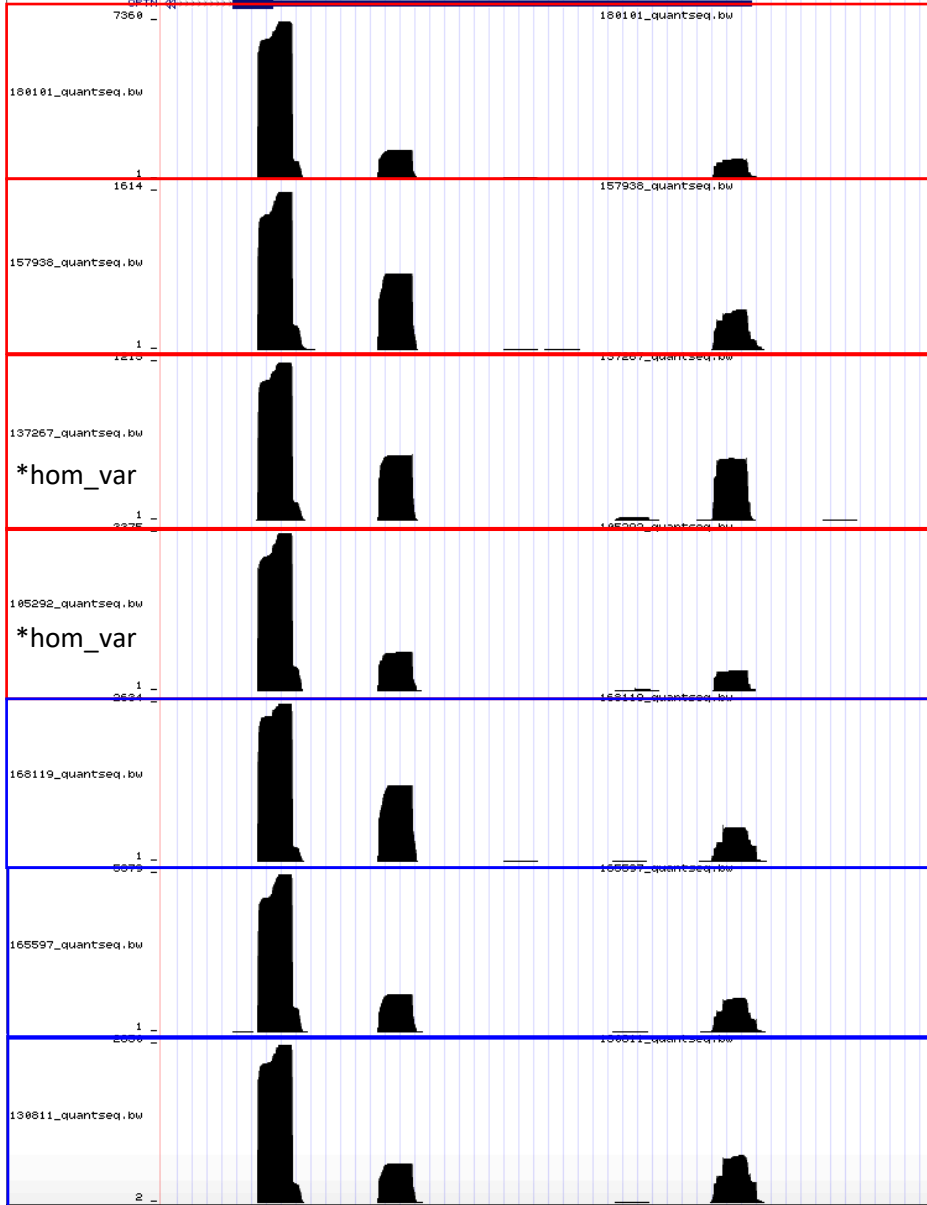
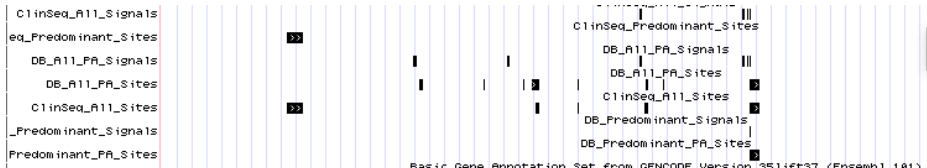


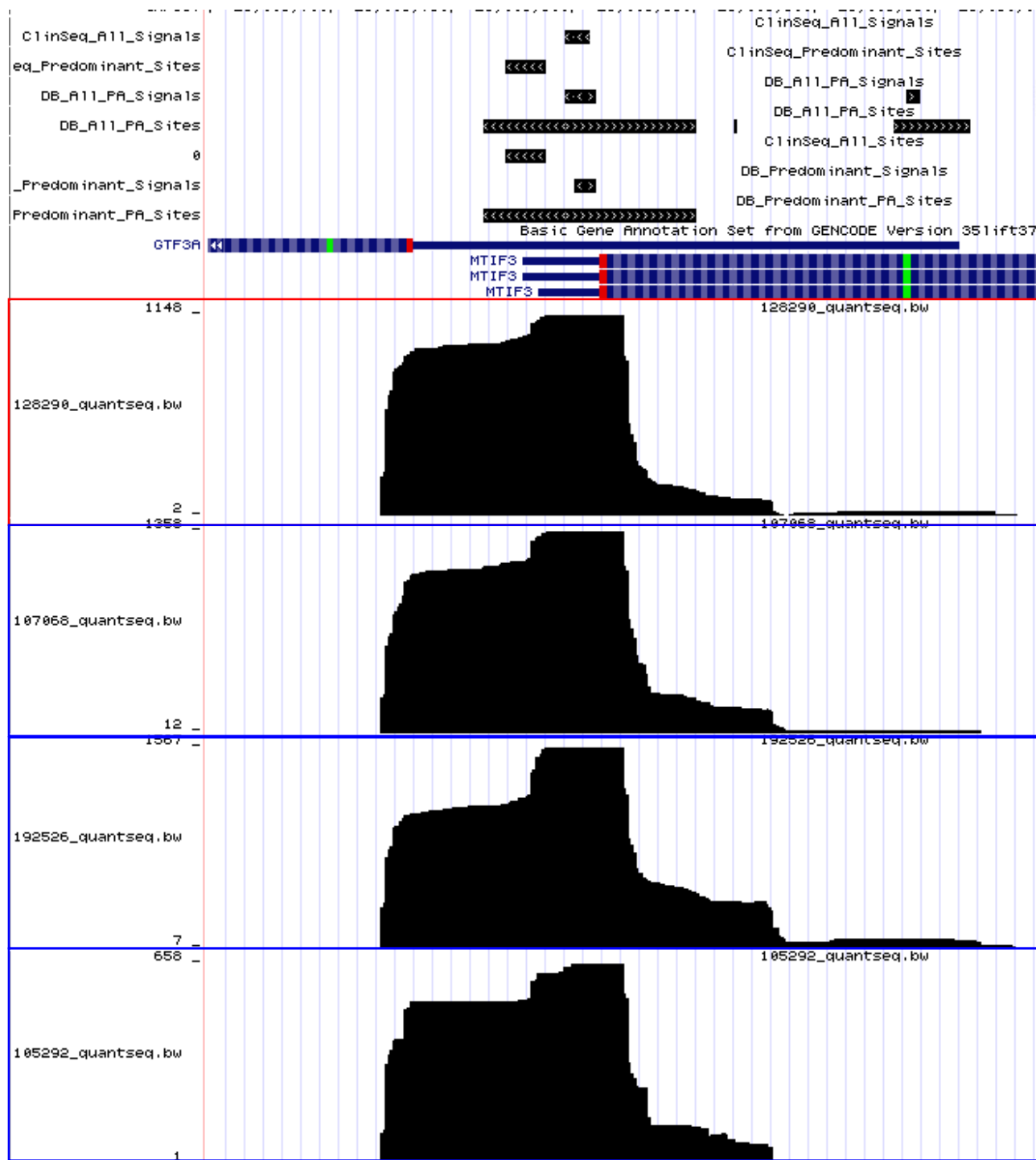


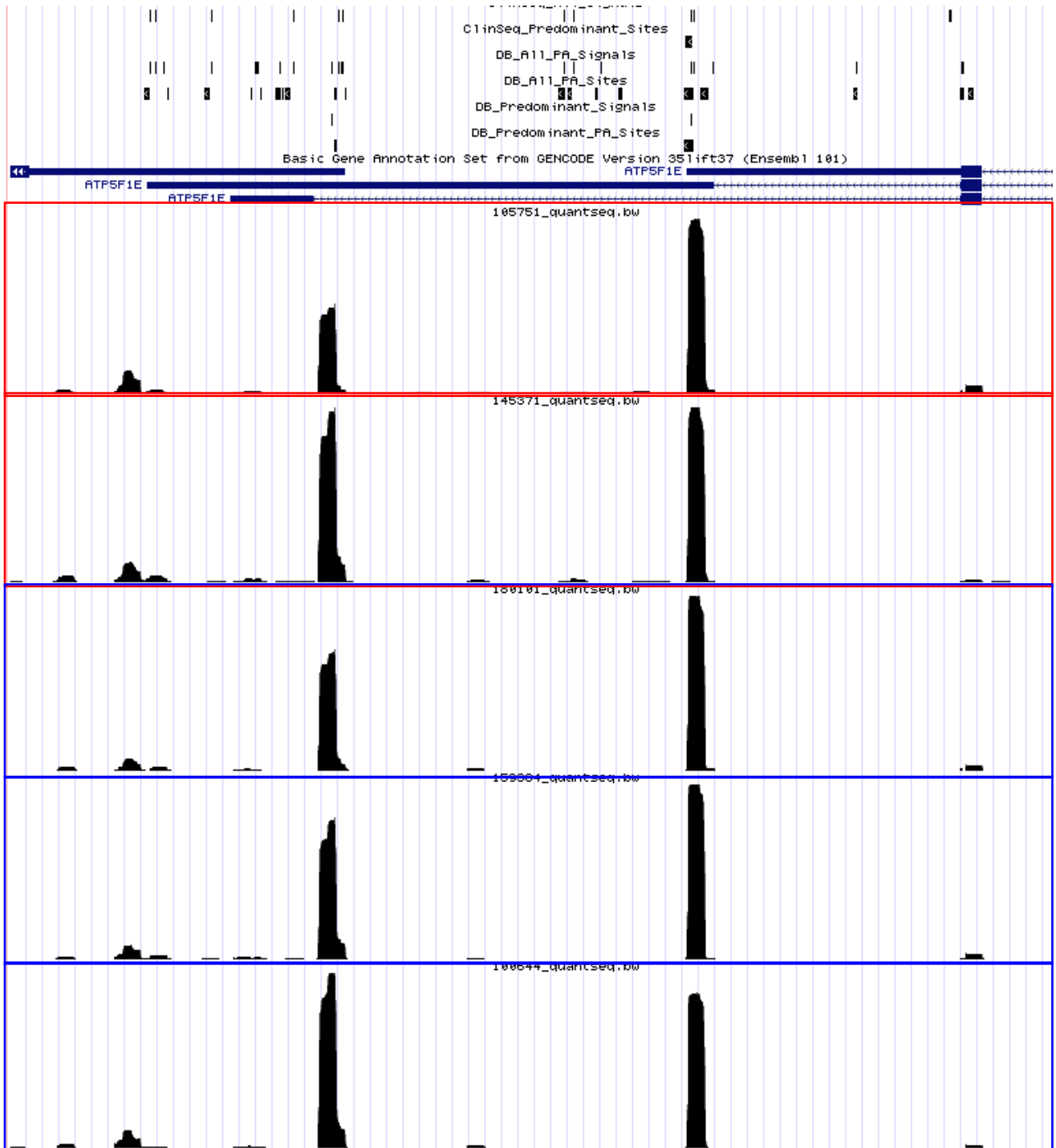
C1inSeq_A11_Signals
 eq_Predominant_Sites
 DB_A11_PA_Signals
 DB_A11_PA_Sites
 C1inSeq_A11_Sites
 _Predominant_Signals
 Predominant_PA_Sites

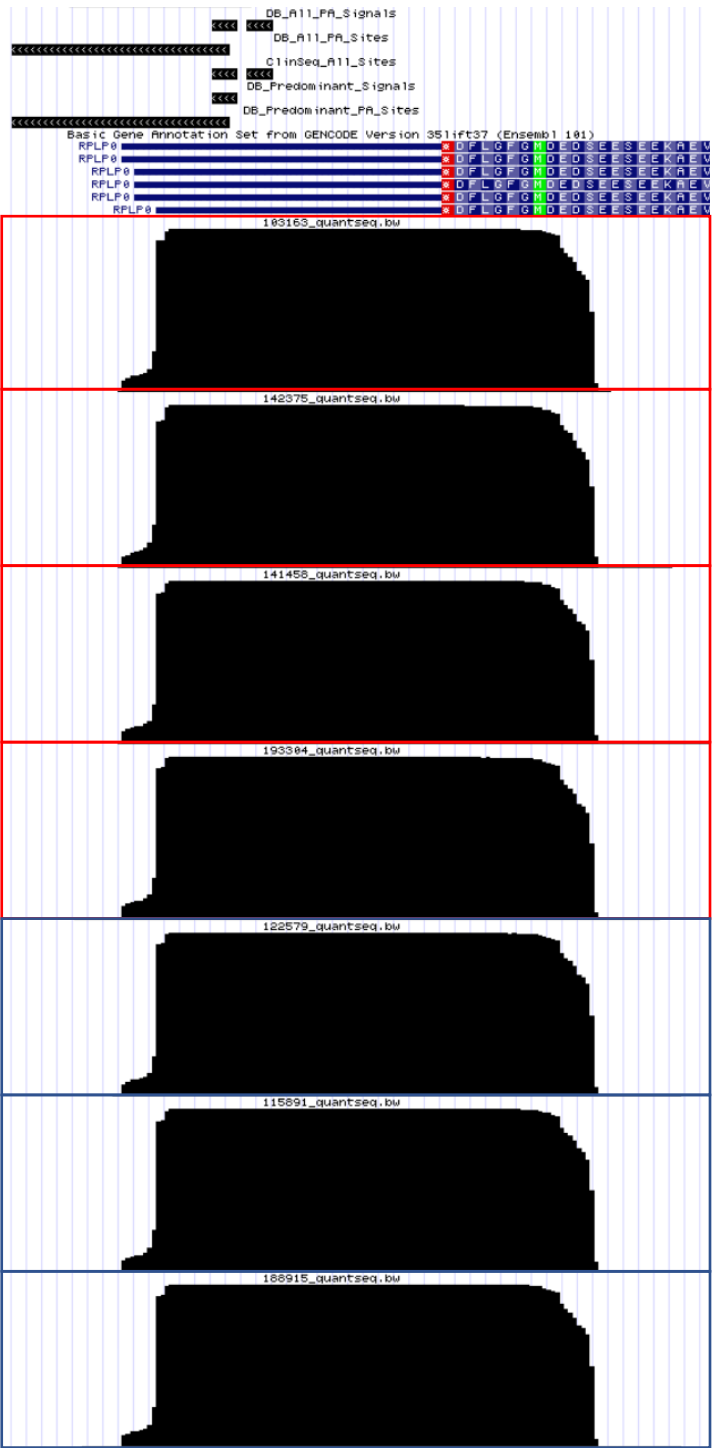
Basic Gene Annotation Set from GENCODE Version 35 lift37 (1





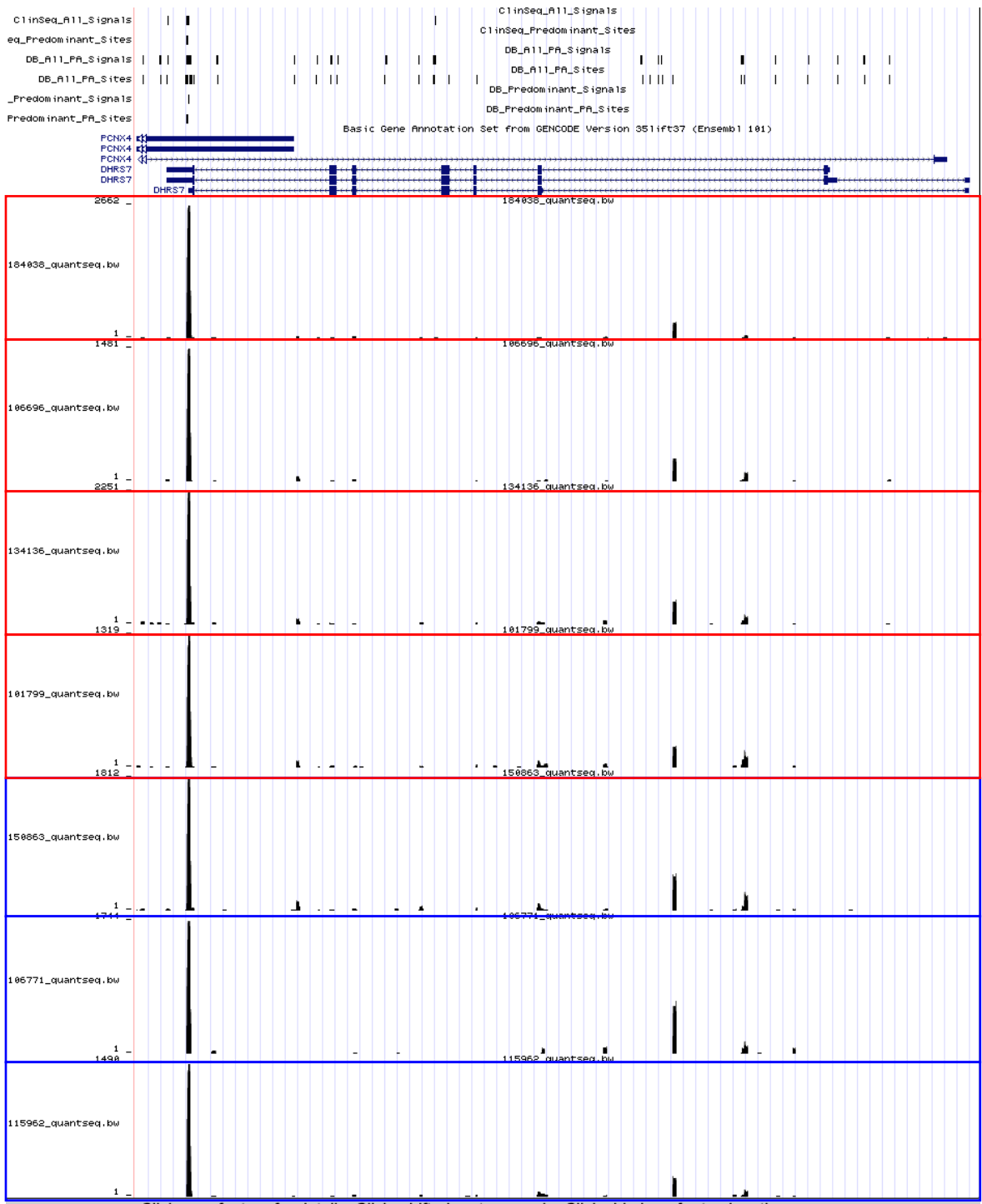




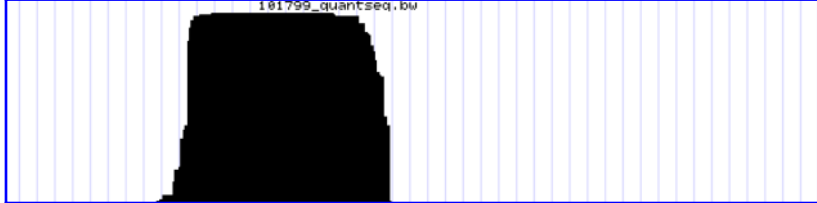
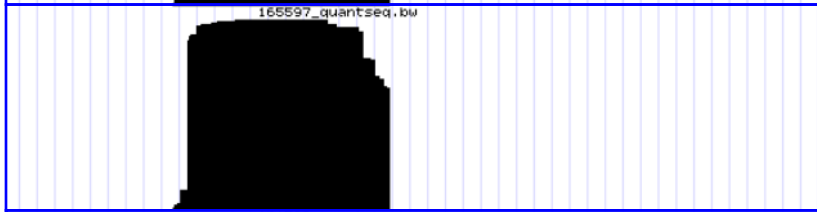
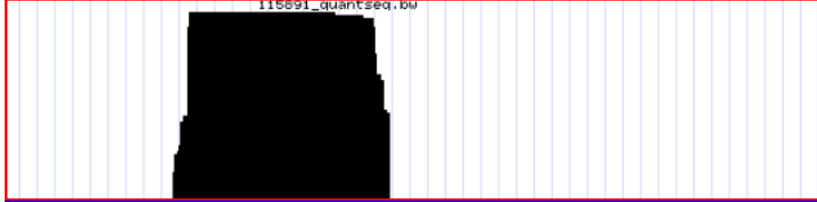
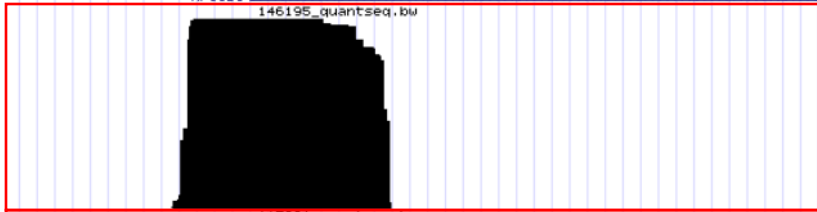


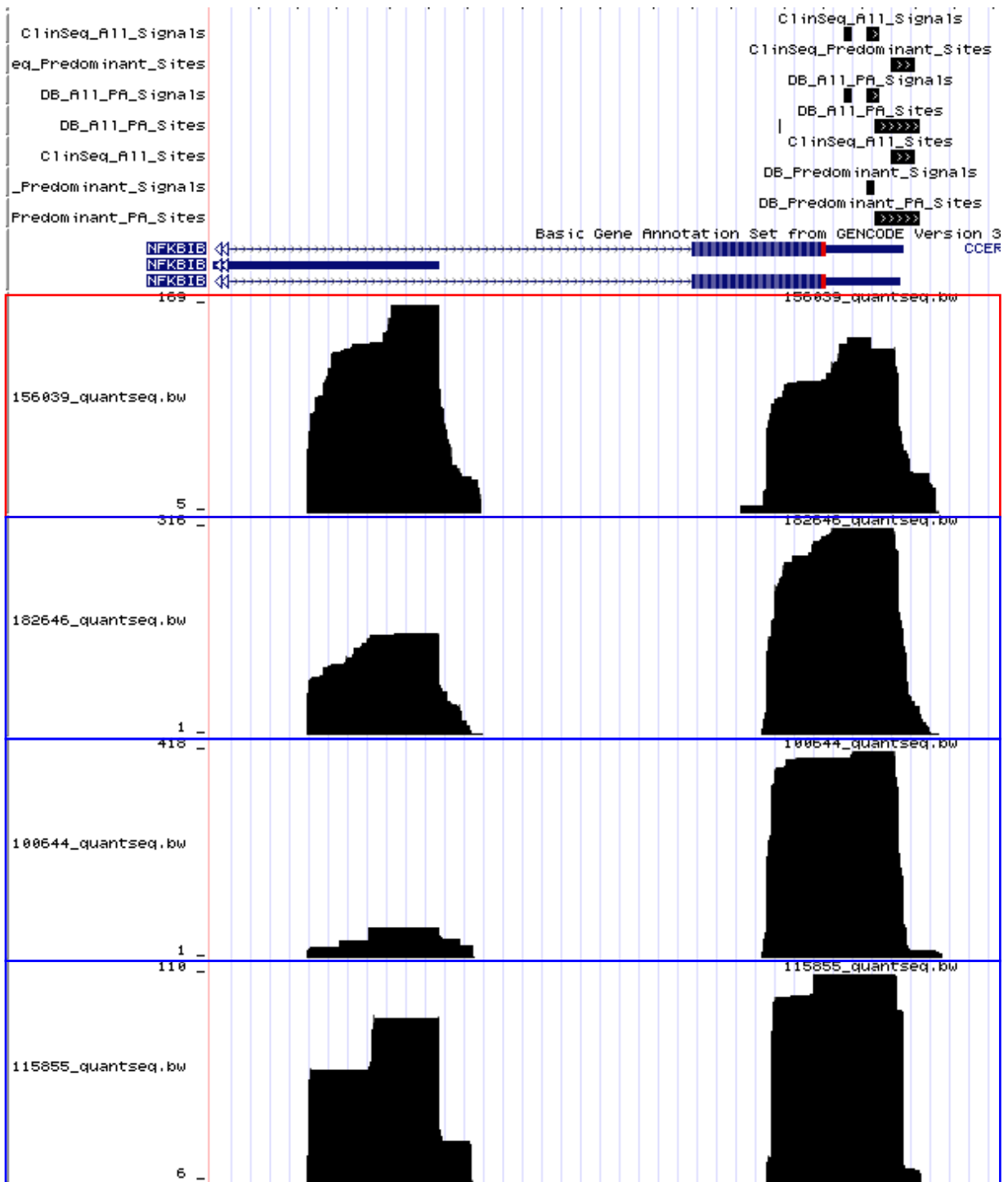
*hom_var

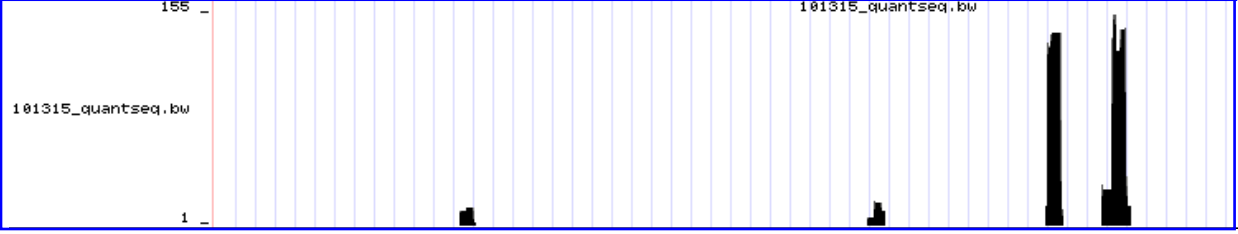
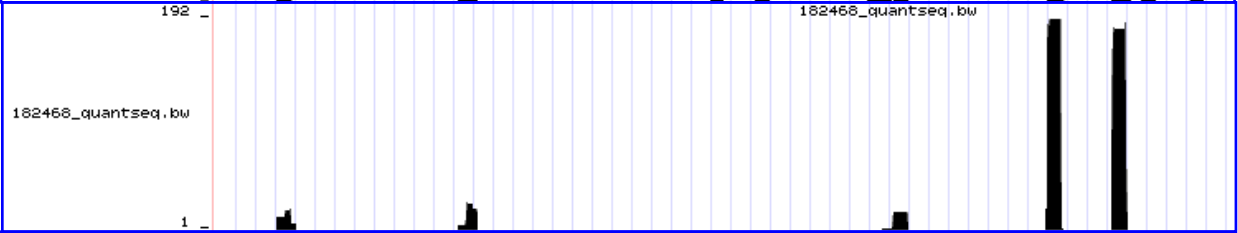
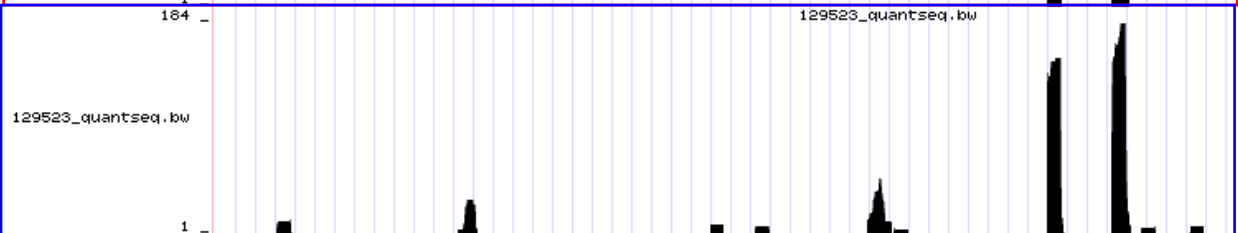
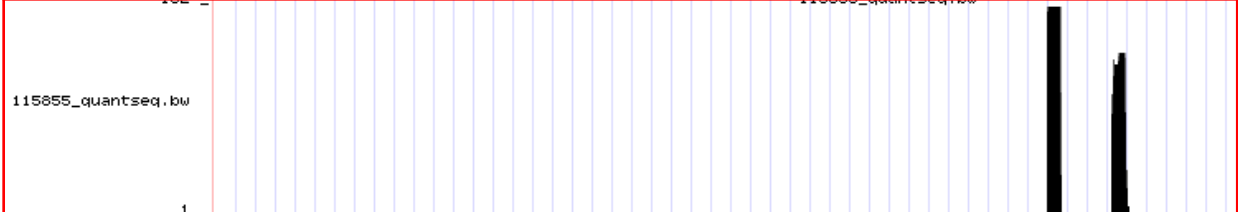
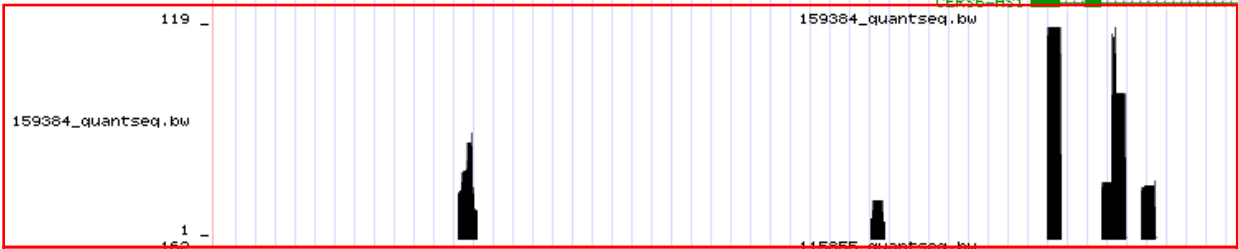
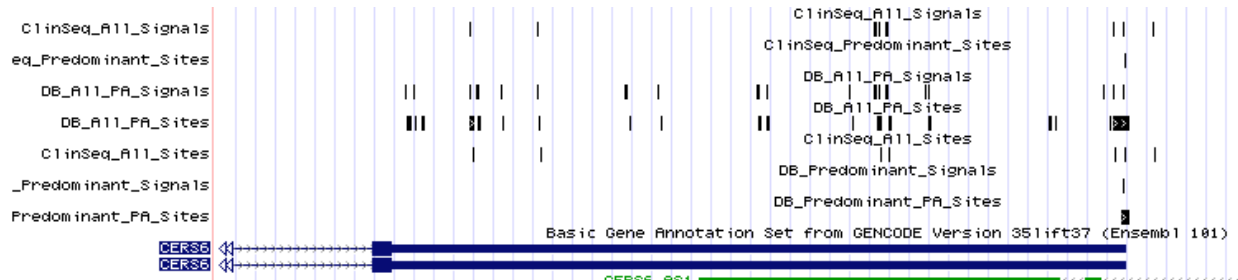
*hom_var



DB_Predominant_Signals
DB_Predominant_PA_Sites
Basic Gene Annotation Set from GENCODE Version 35 lift37 (Ensembl 101)
RPUSD1
RPUSD1
RPUSD1



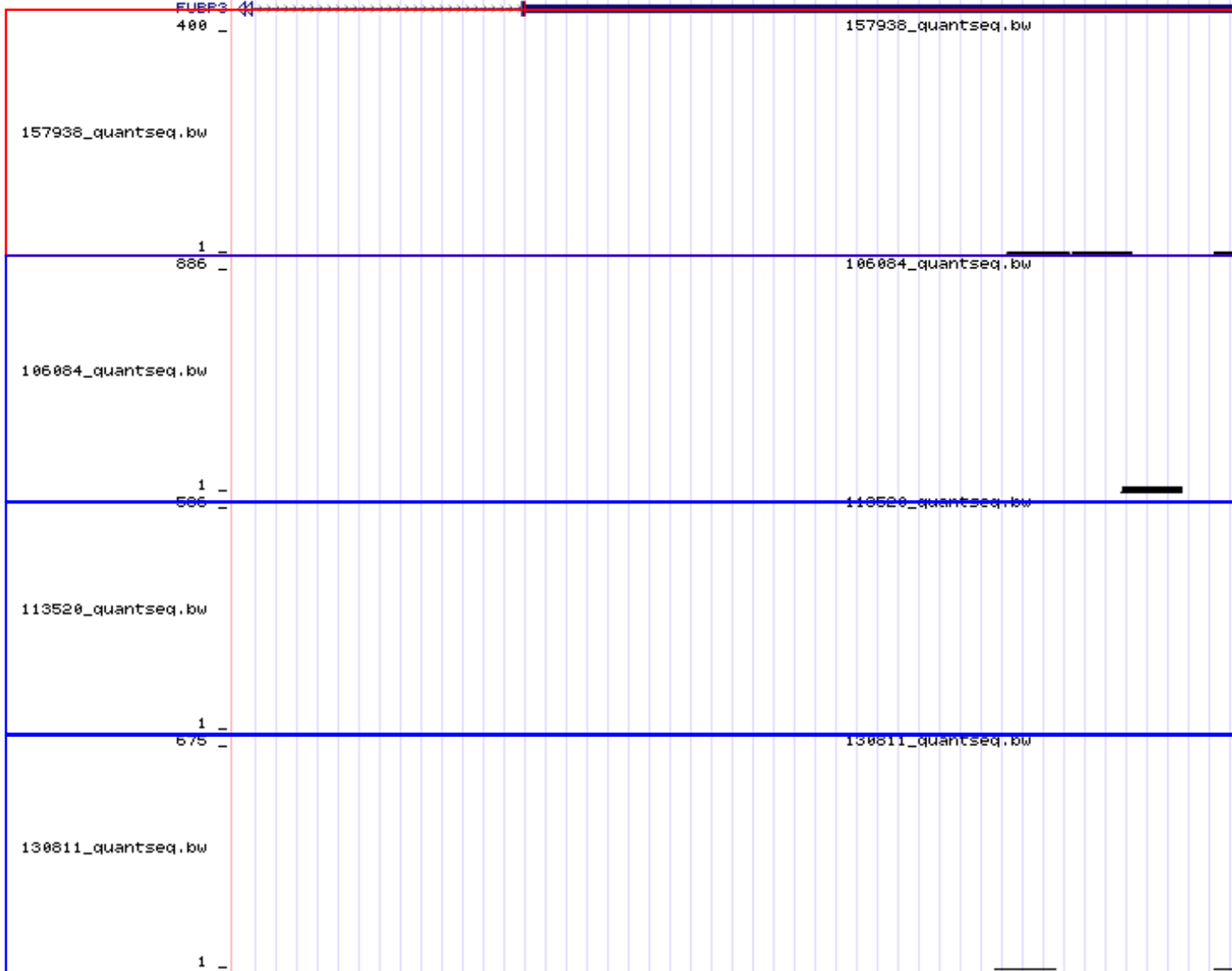


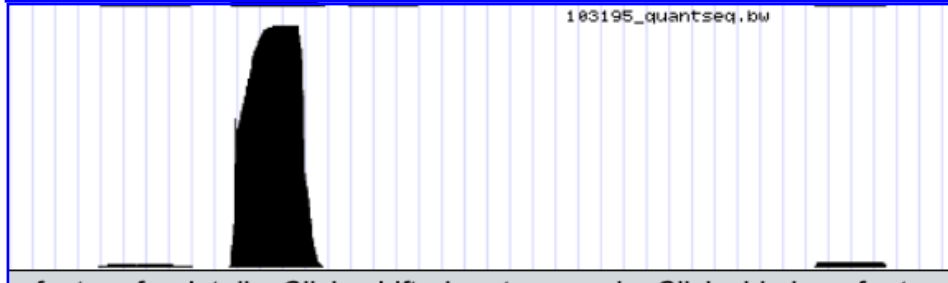
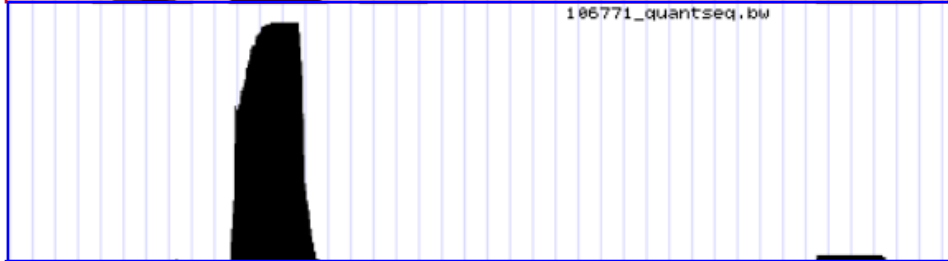
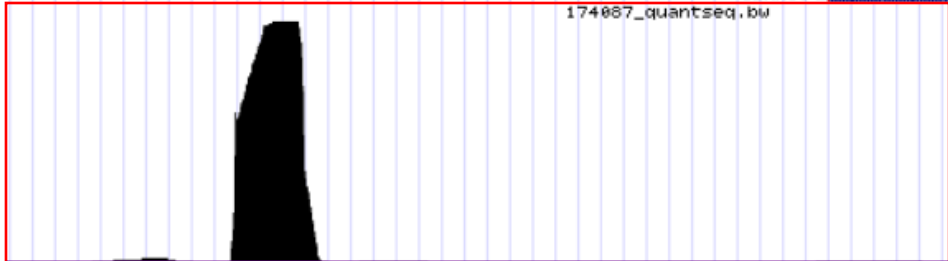


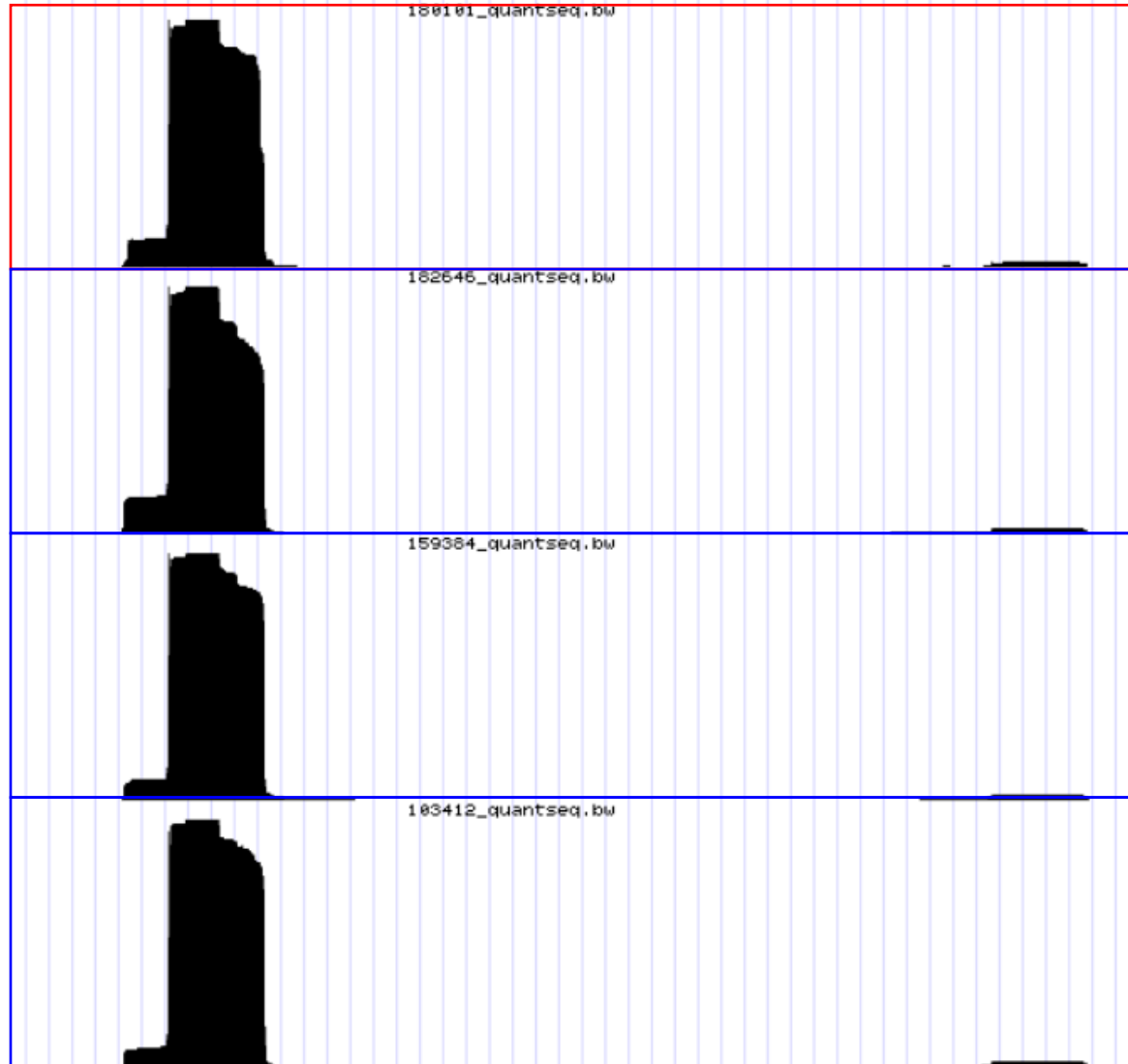
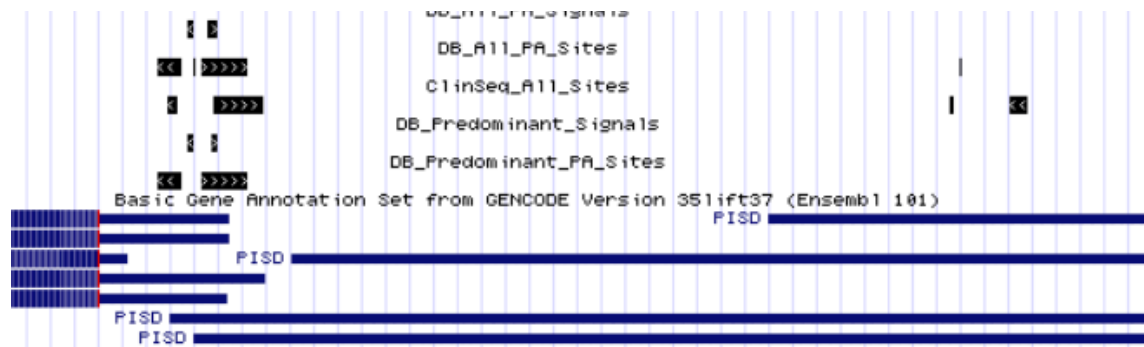
ClinSeq_A11_Signals
eq_Predominant_Sites
DB_A11_PA_Signals
DB_A11_PA_Sites
ClinSeq_A11_Sites
_Predominant_Signals
Predominant_PA_Sites

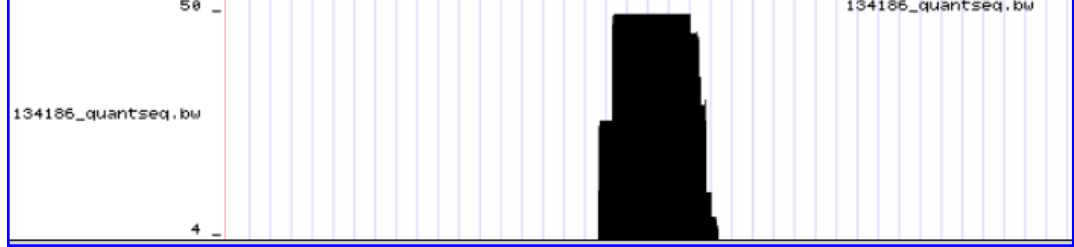
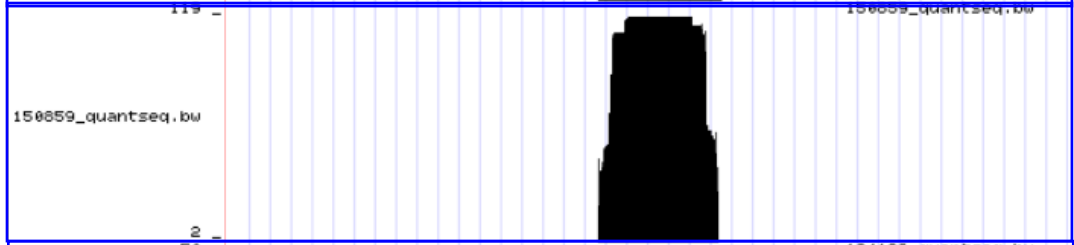
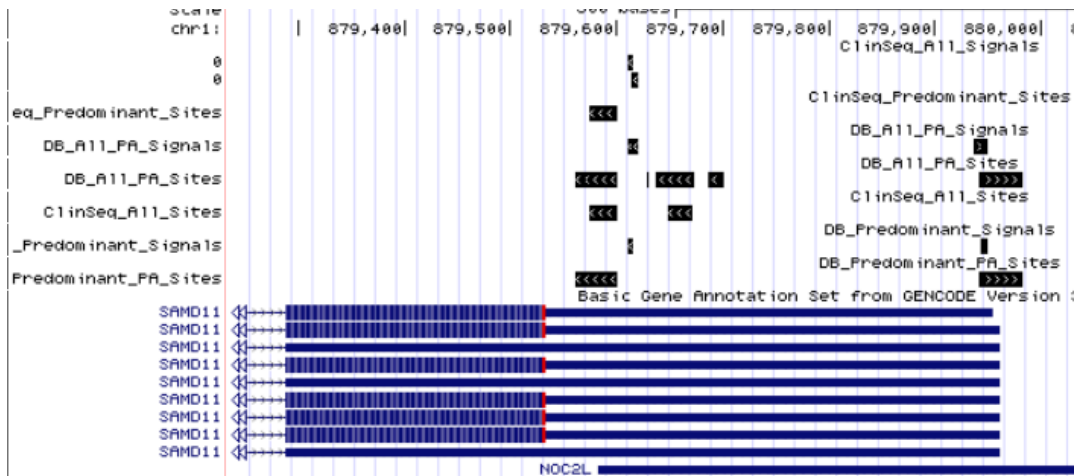
ClinSeq_A11_Signals
ClinSeq_Predominant_Sites
DB_A11_PA_Signals
DB_A11_PA_Sites
ClinSeq_A11_Sites
DB_Predominant_Signals
DB_Predominant_PA_Sites

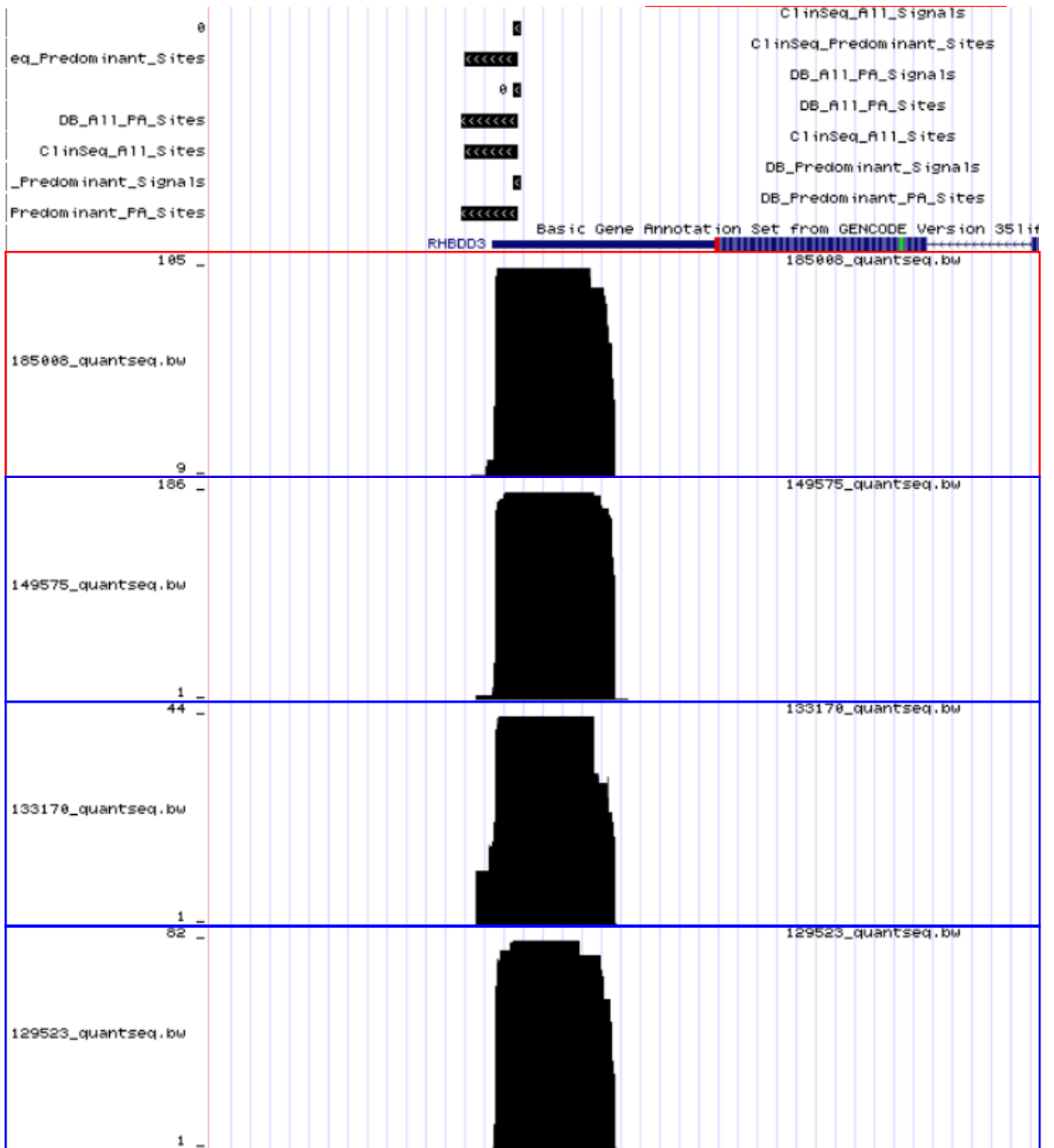
Basic Gene Annotation Set from GENCODE Version 35lift37 (Ensembl)

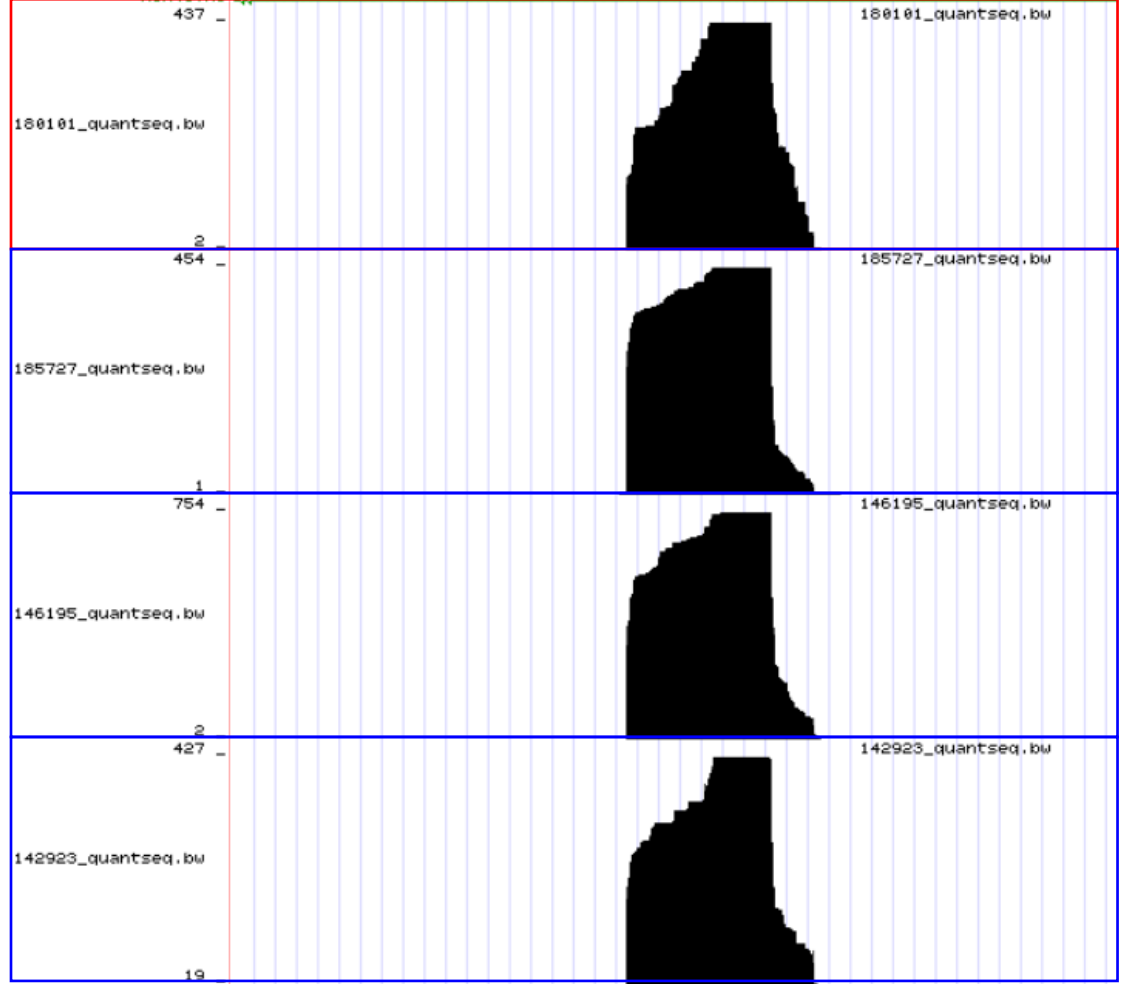
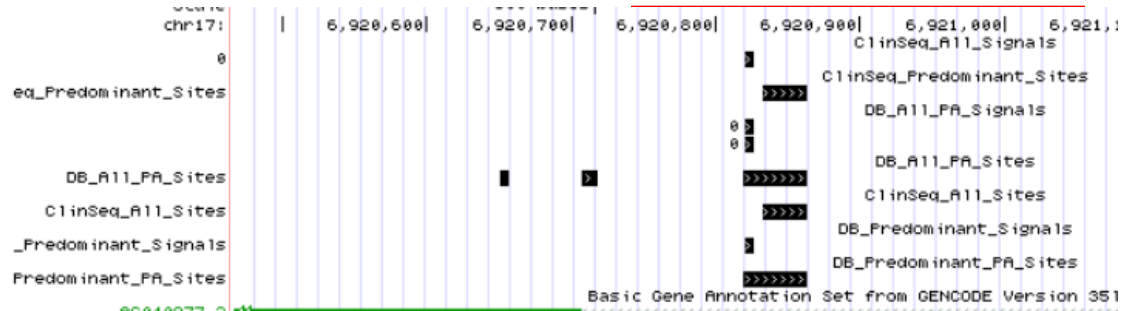


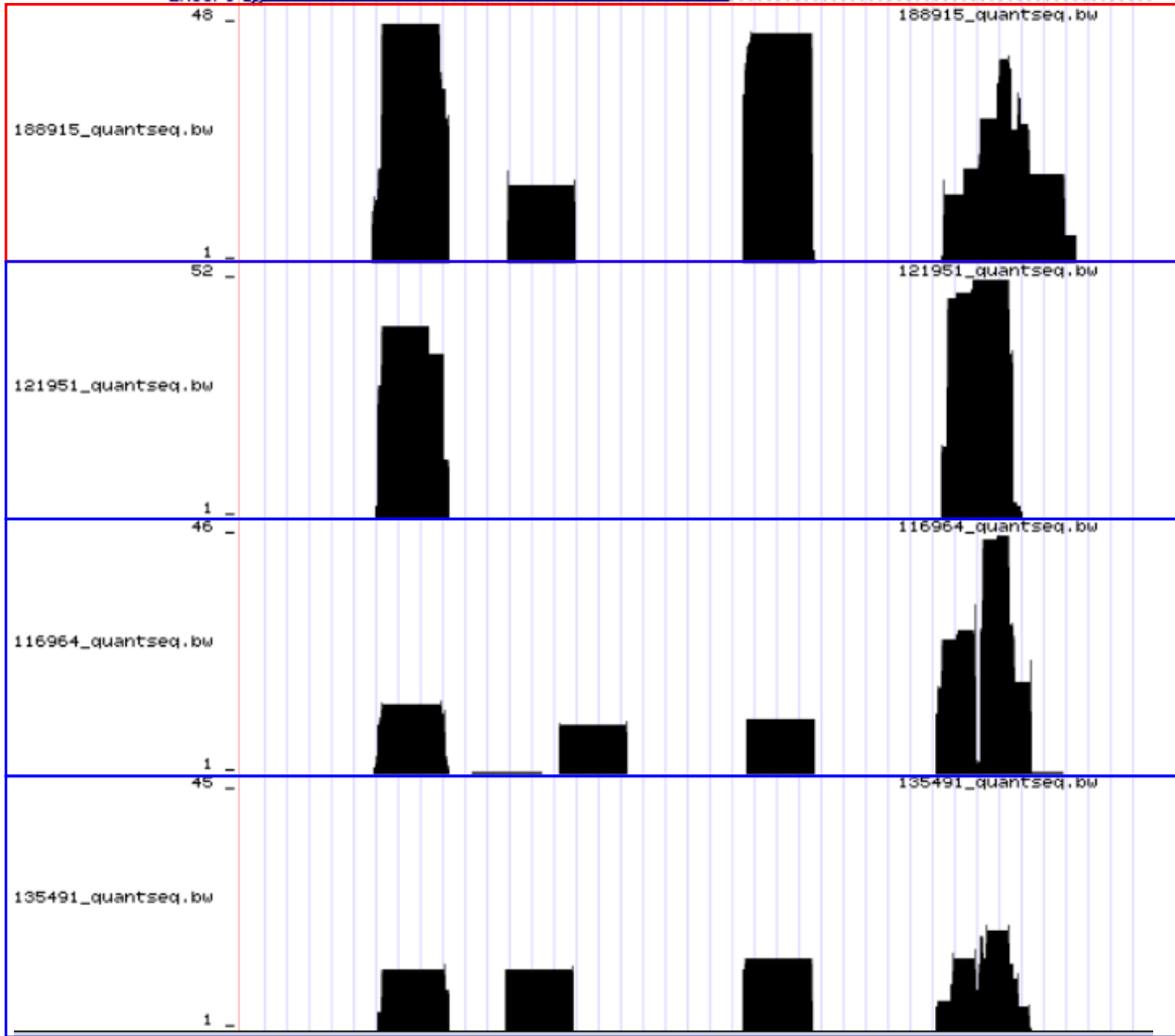
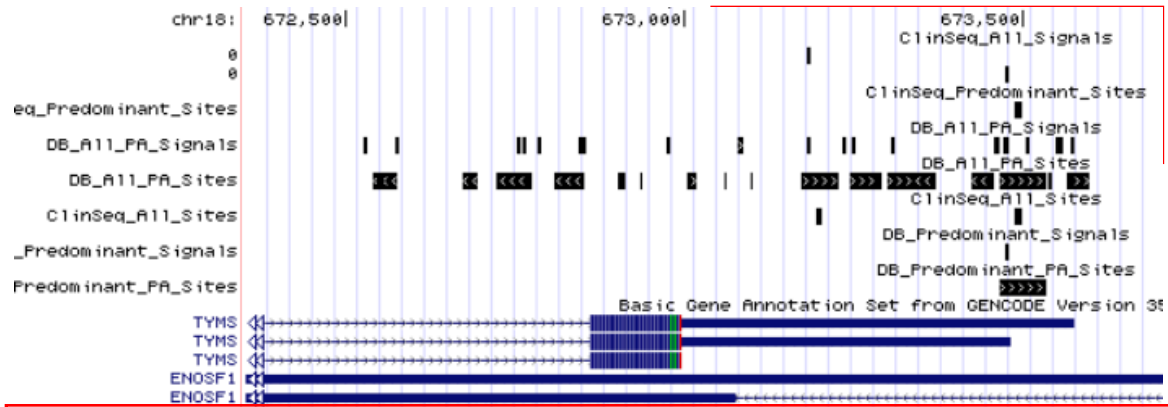


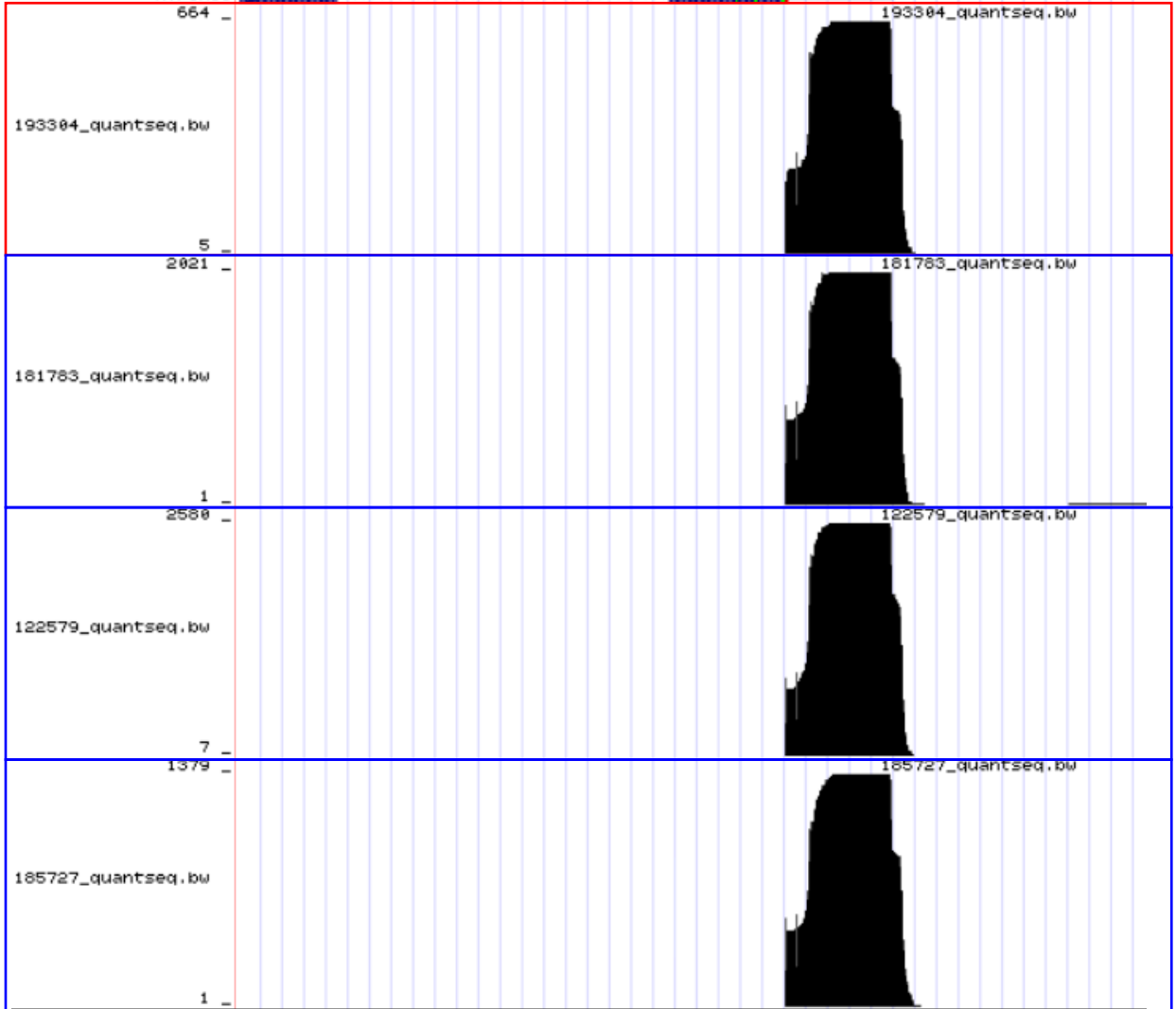
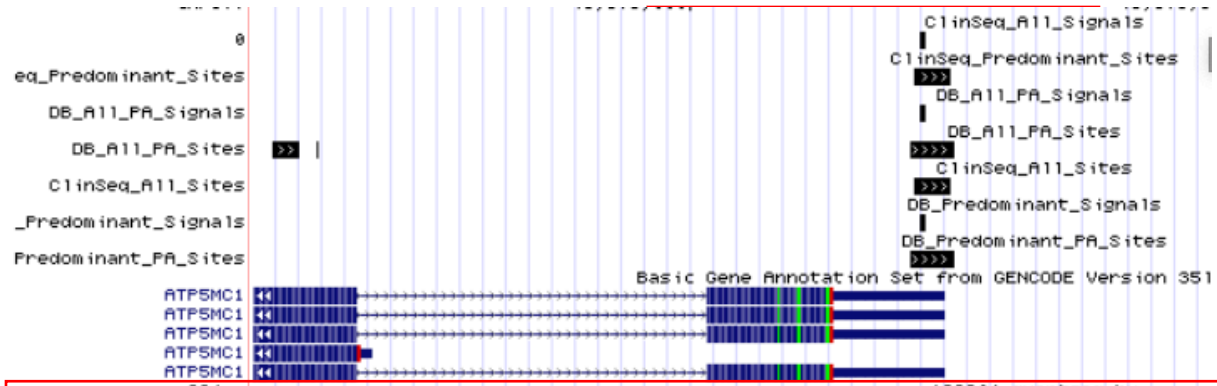


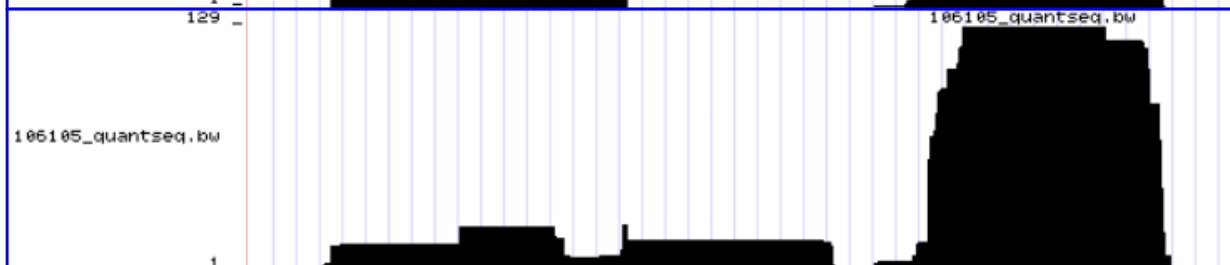
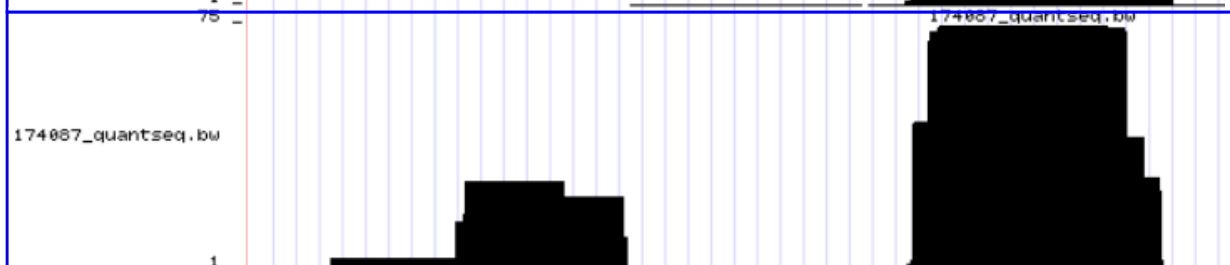
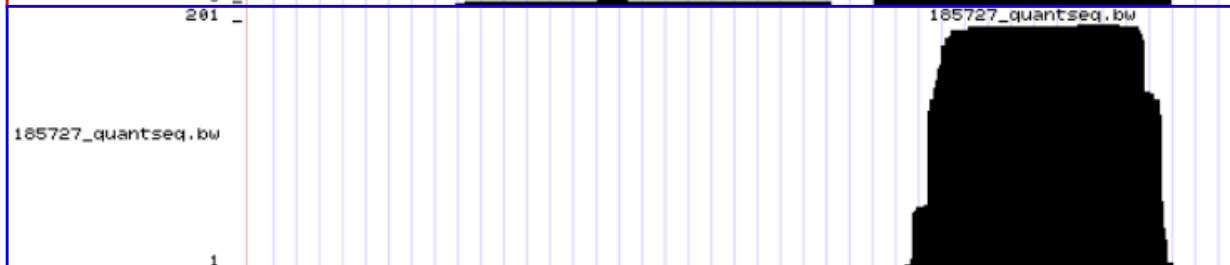
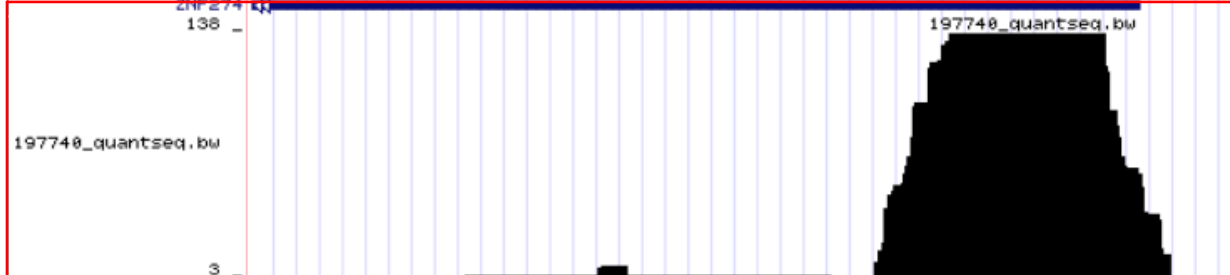
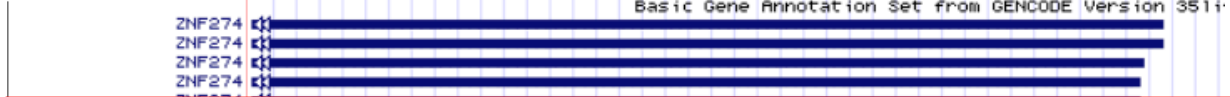
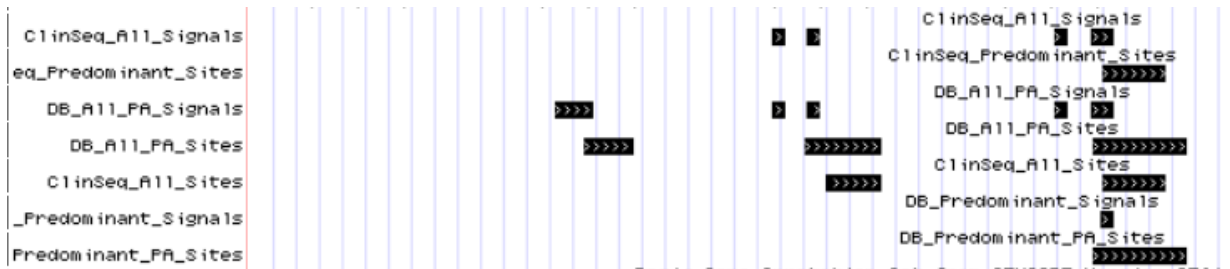


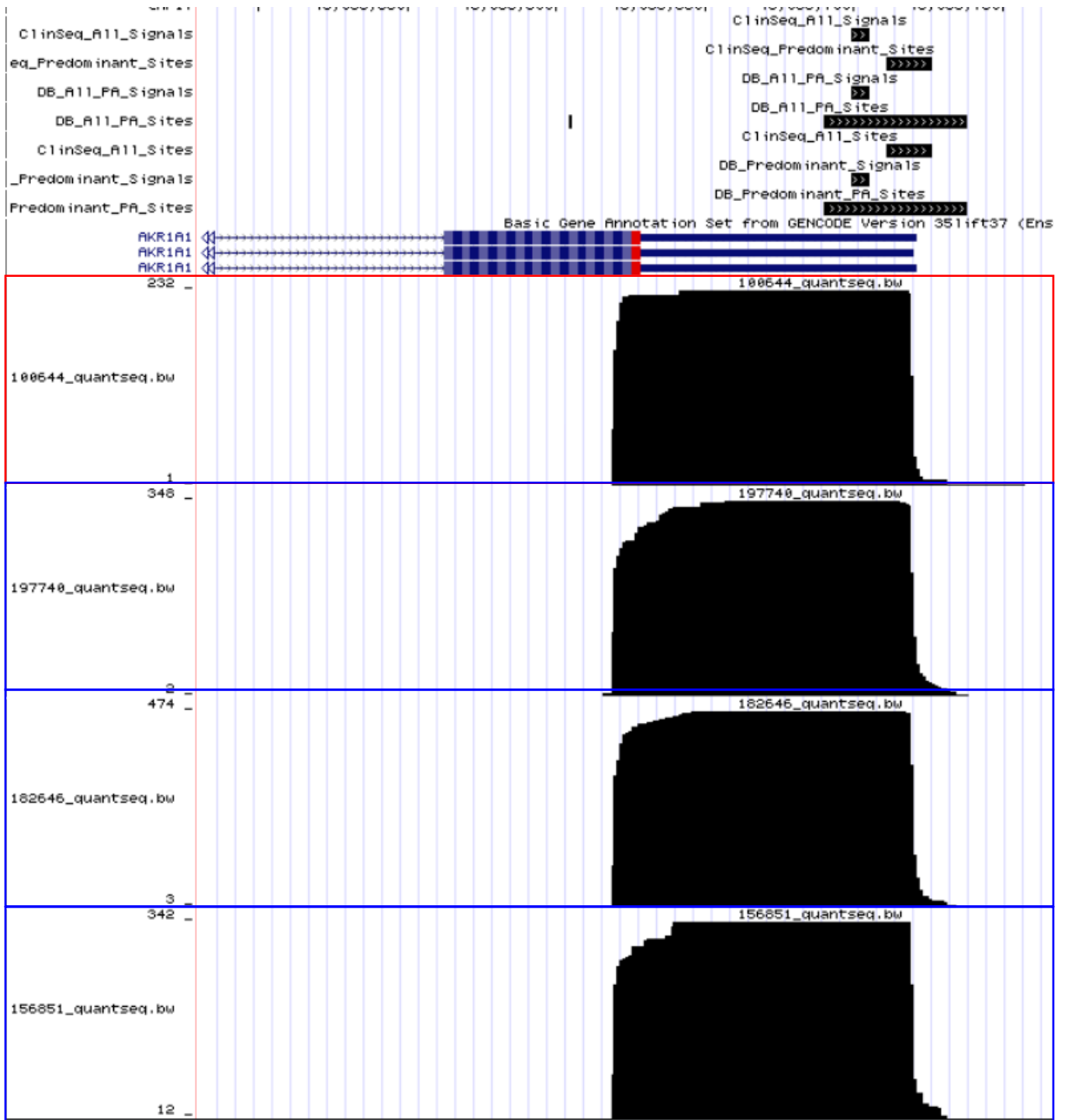












ClinSeq_A11_Signals
ClinSeq_Predominant_Sites
DB_A11_PA_Signals
DB_A11_PA_Sites
DB_Predominant_Signals
DB_Predominant_PA_Sites

Basic Gene Annotation Set from GENCODE Version 35 lift37 (Ensembl 101)

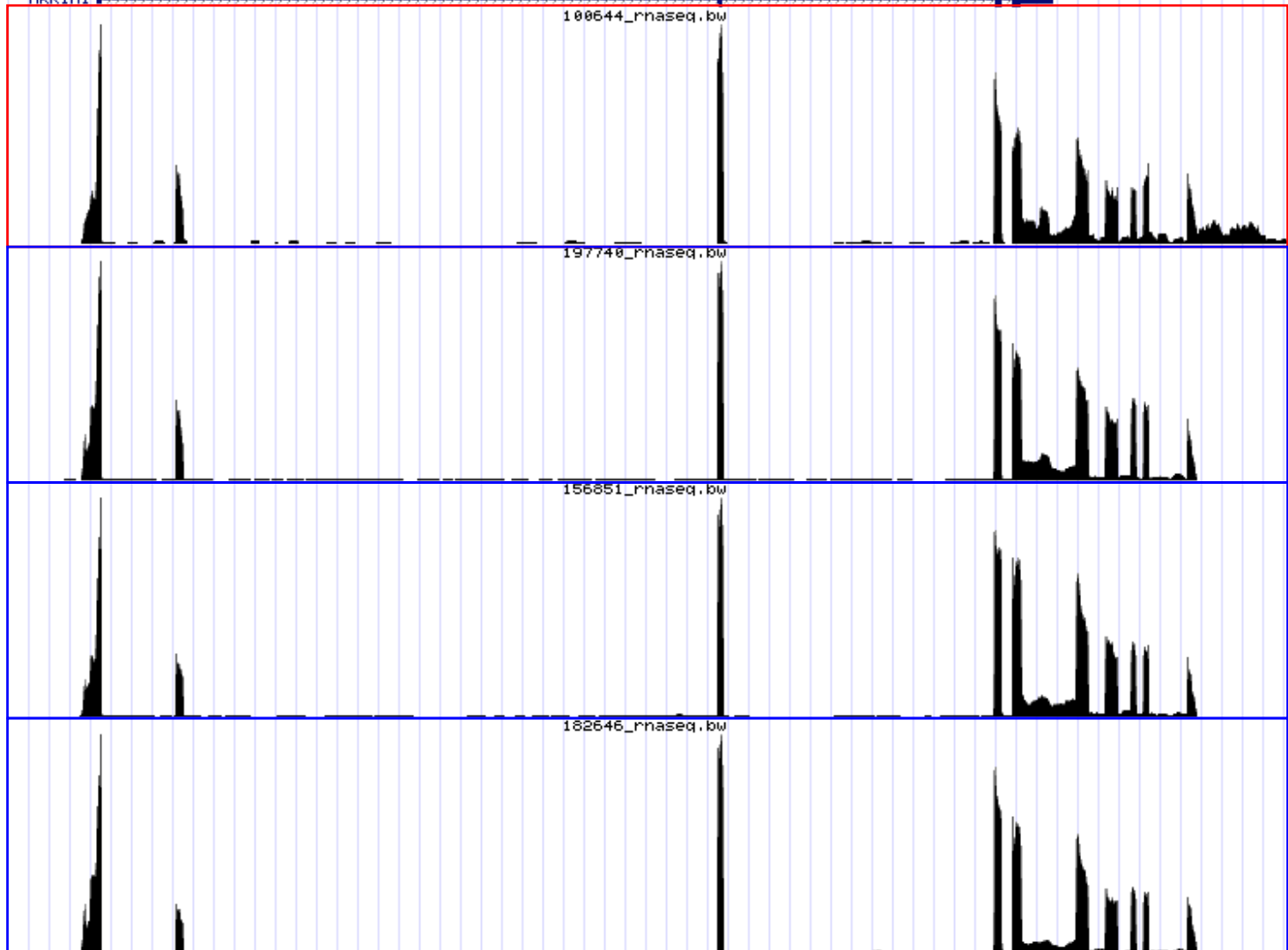


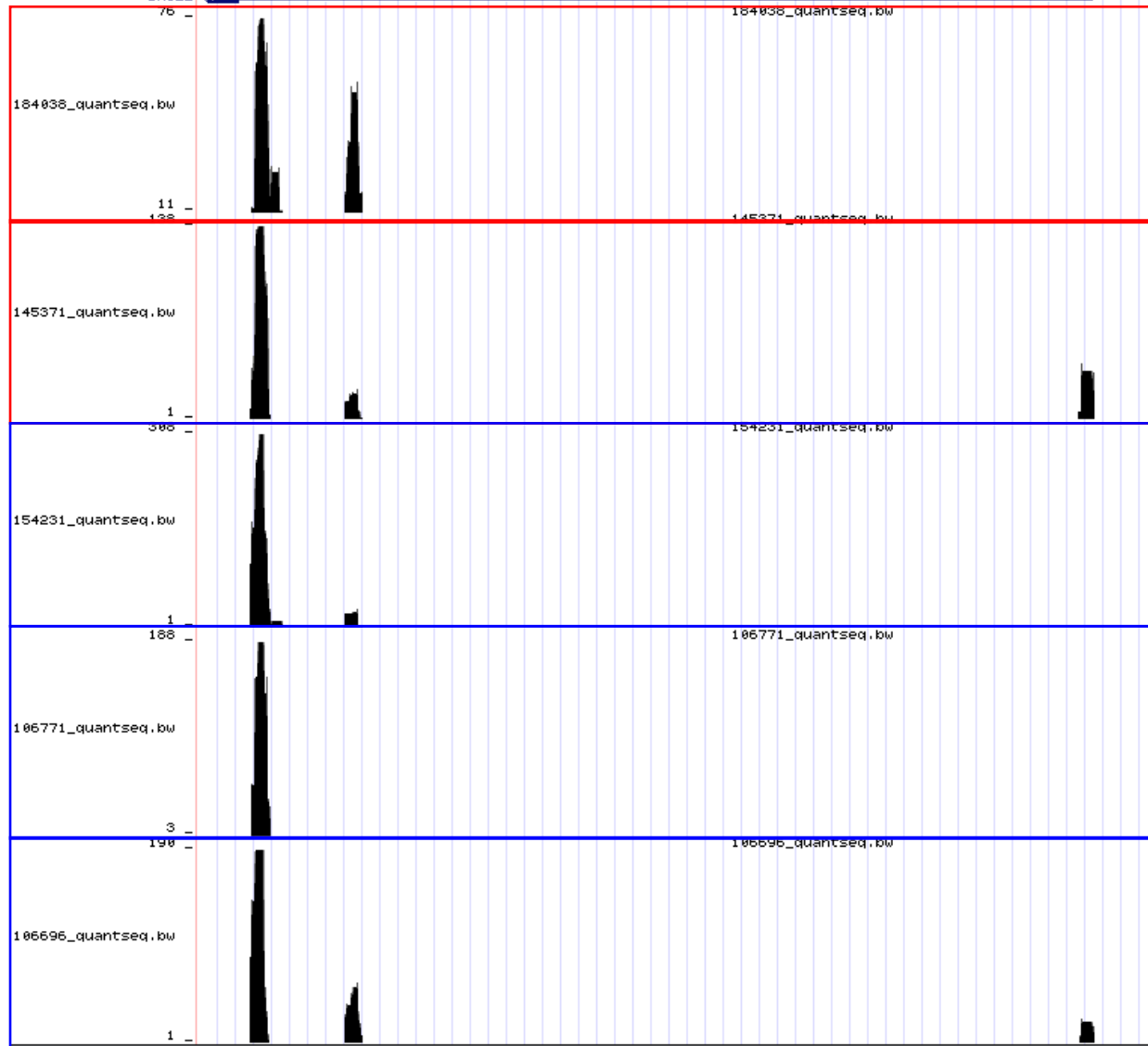
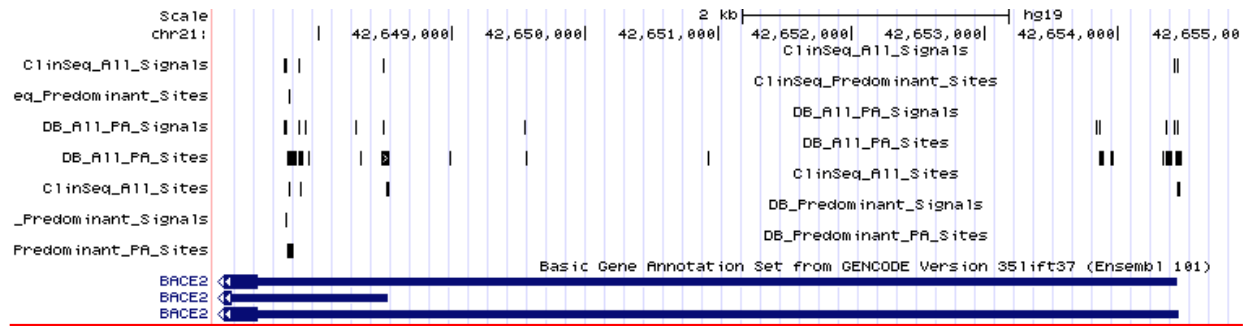
188544_rnaseq.bw

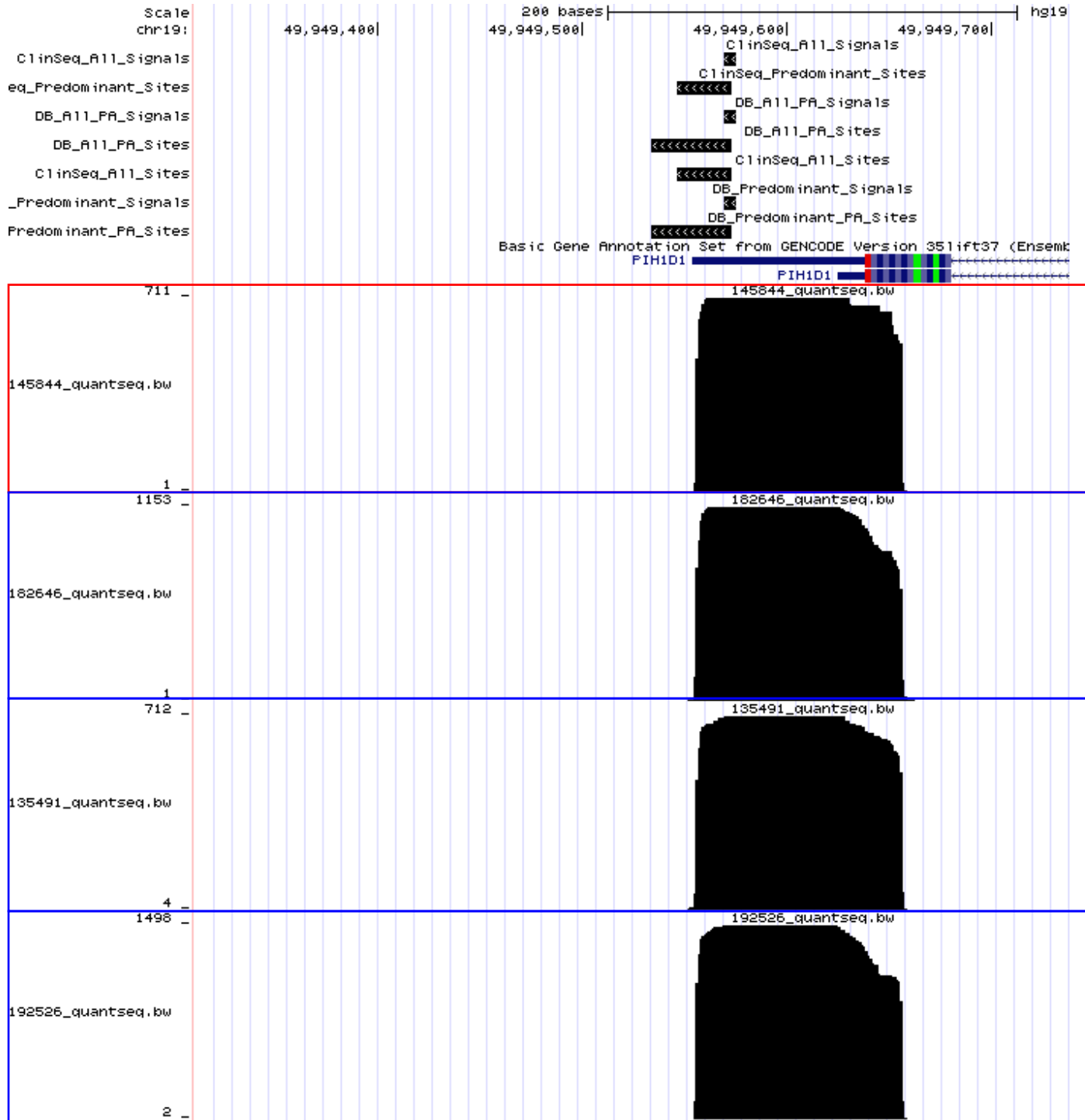
197748_rnaseq.bw

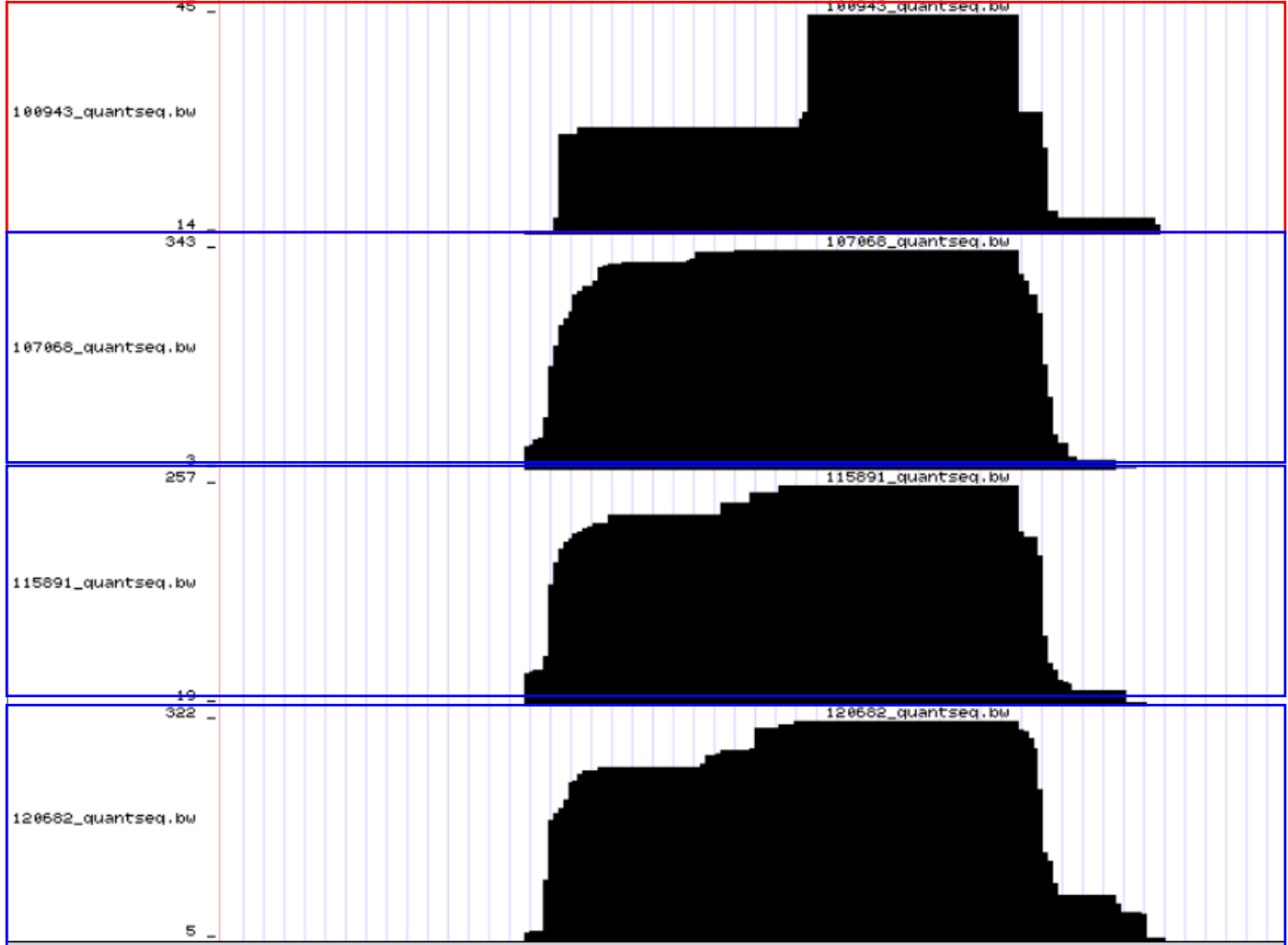
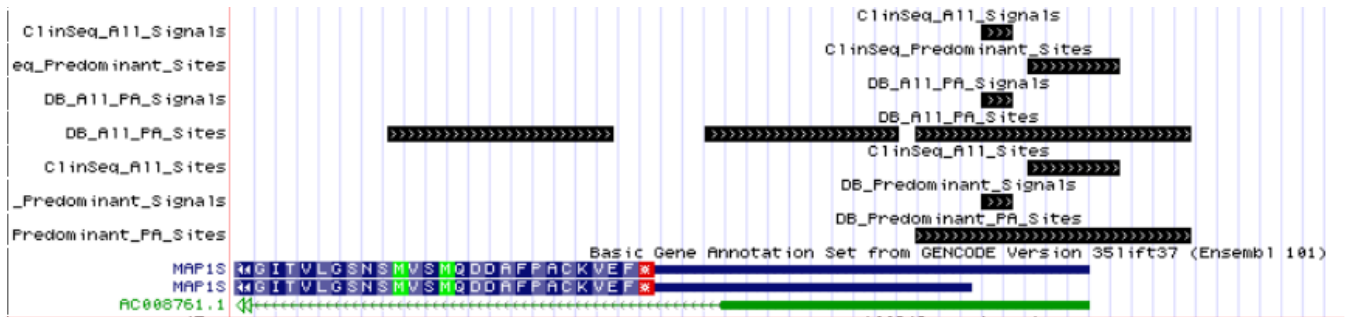
156851_rnaseq.bw

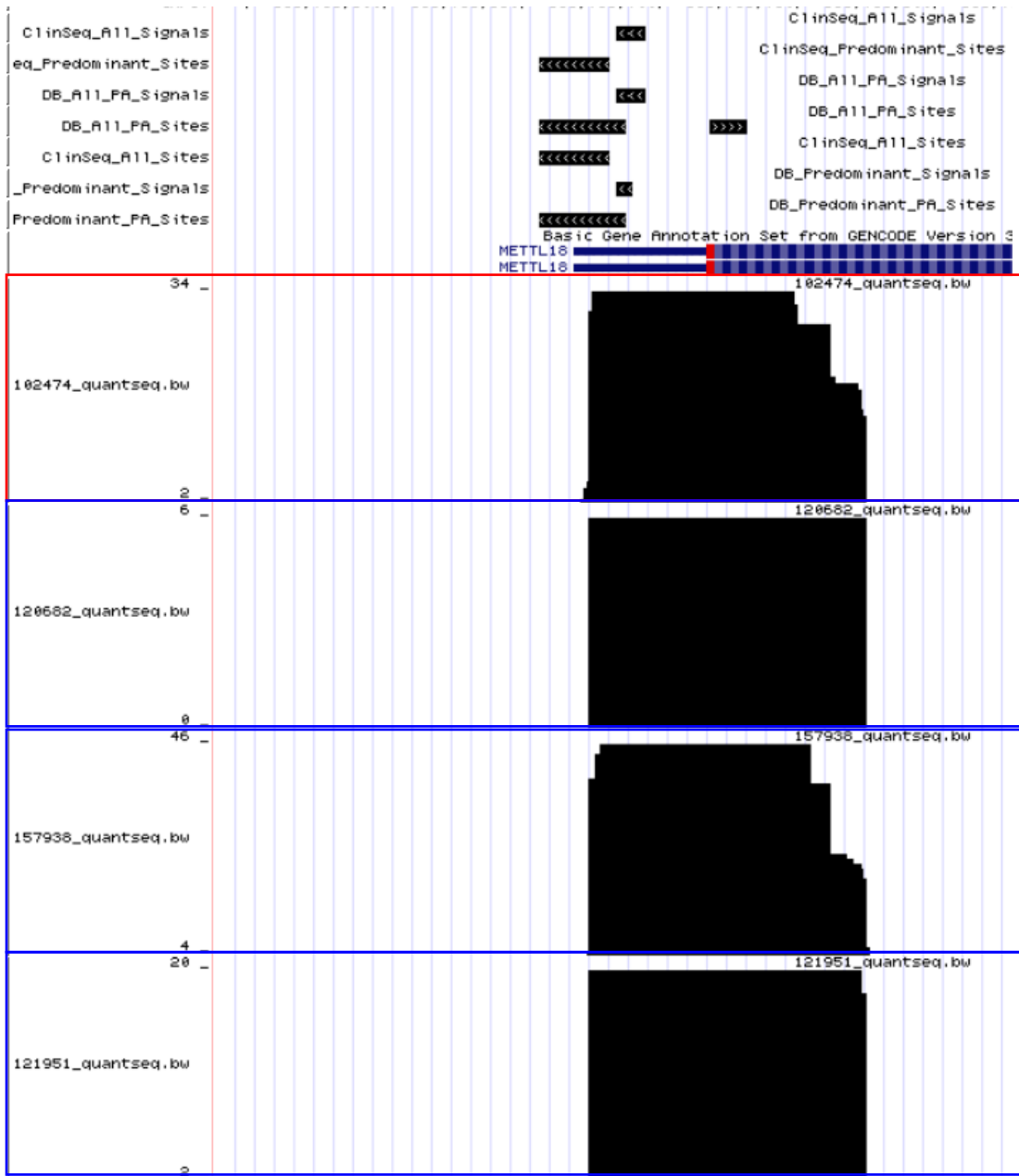
182546_rnaseq.bw









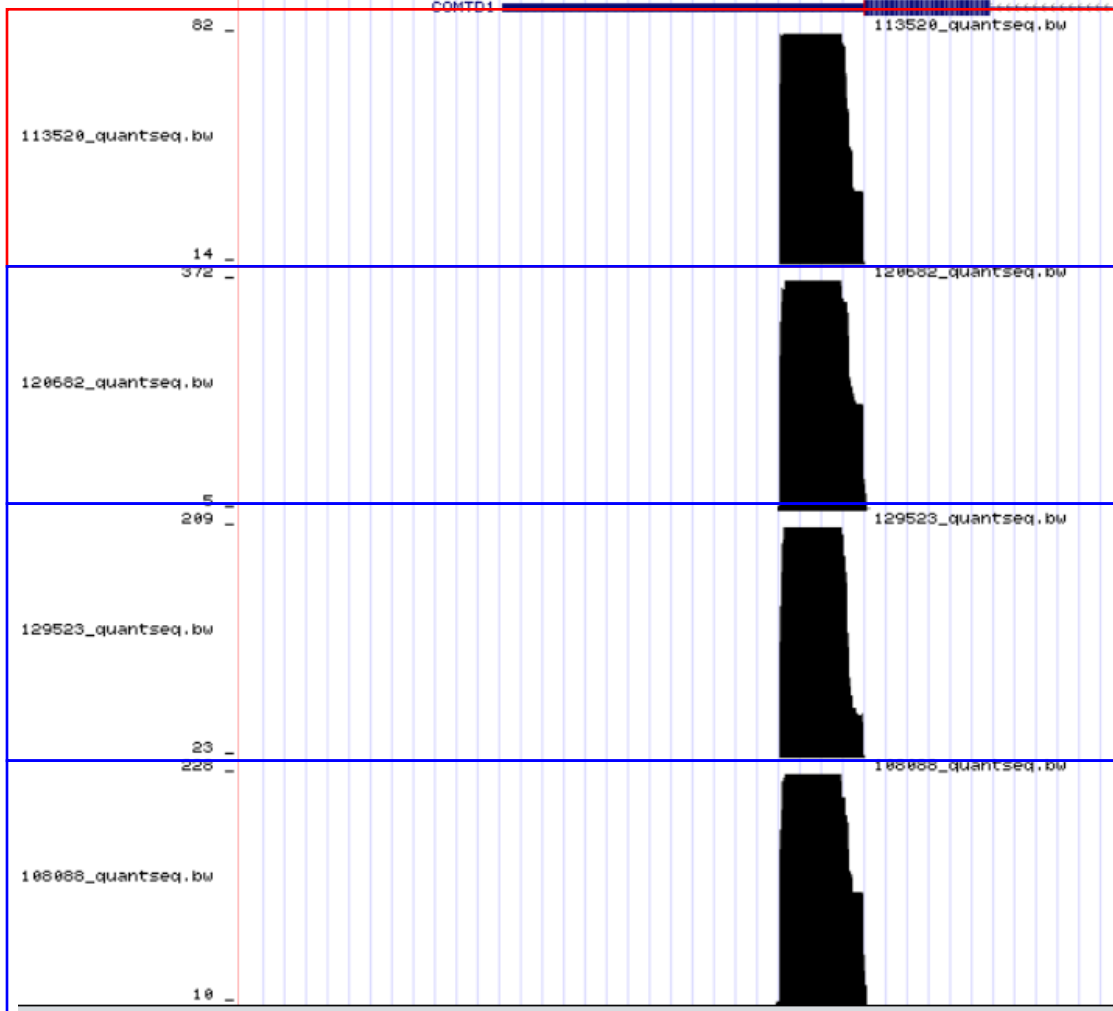


chr10: 76,993,500 76,994,000

ClinSeq_A11_Signals
 ed_Predominant_Sites
 DB_A11_PA_Signals
 DB_A11_PA_Sites
 ClinSeq_A11_Sites
 _Predominant_Signals
 Predominant_PA_Sites

ClinSeq_A11_Signals
 ClinSeq_Predominant_Sites
 DB_A11_PA_Signals
 DB_A11_PA_Sites
 ClinSeq_A11_Sites
 DB_Predominant_Signals
 DB_Predominant_PA_Sites

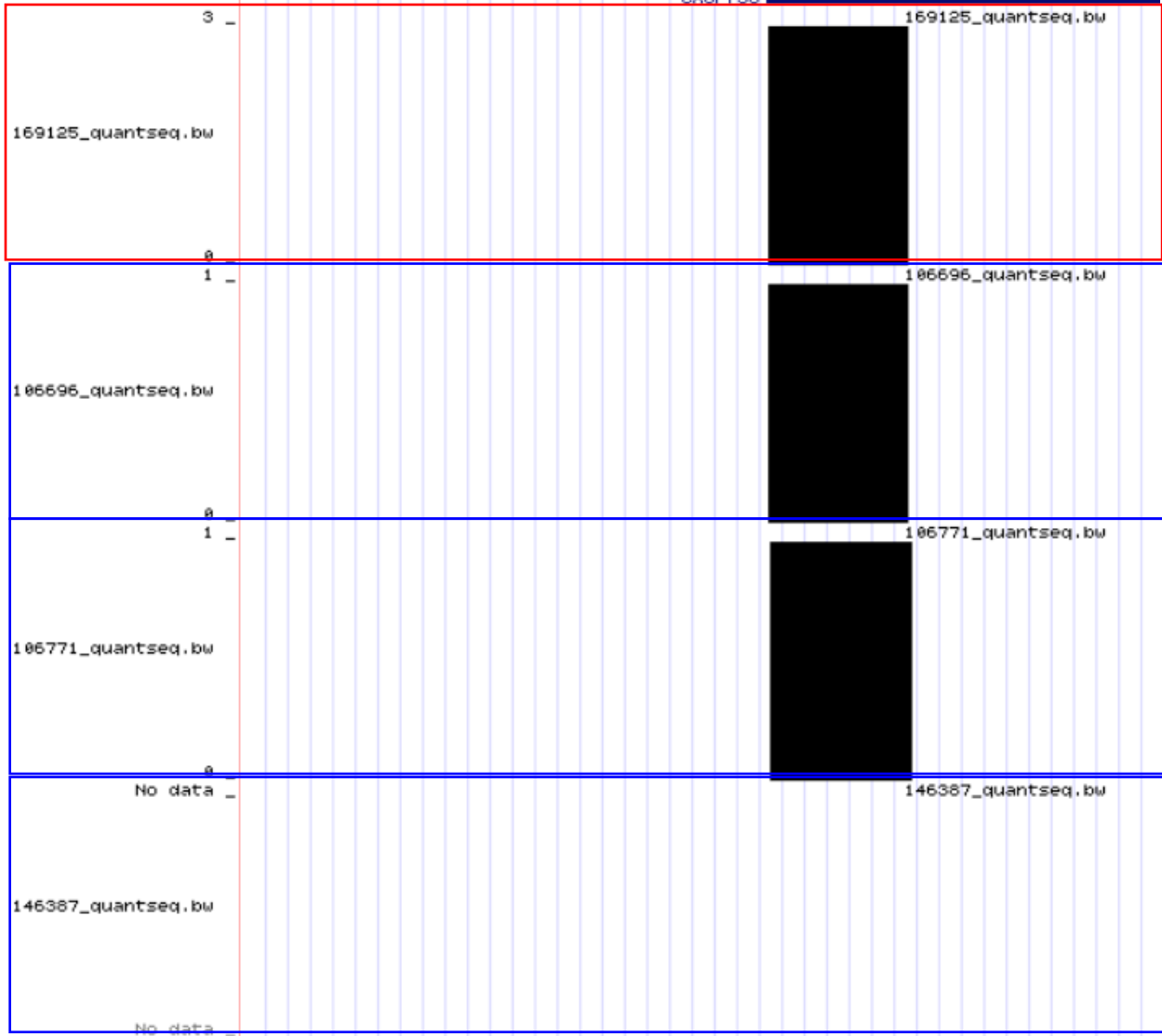
Basic Gene Annotation Set from GENCODE Version 35

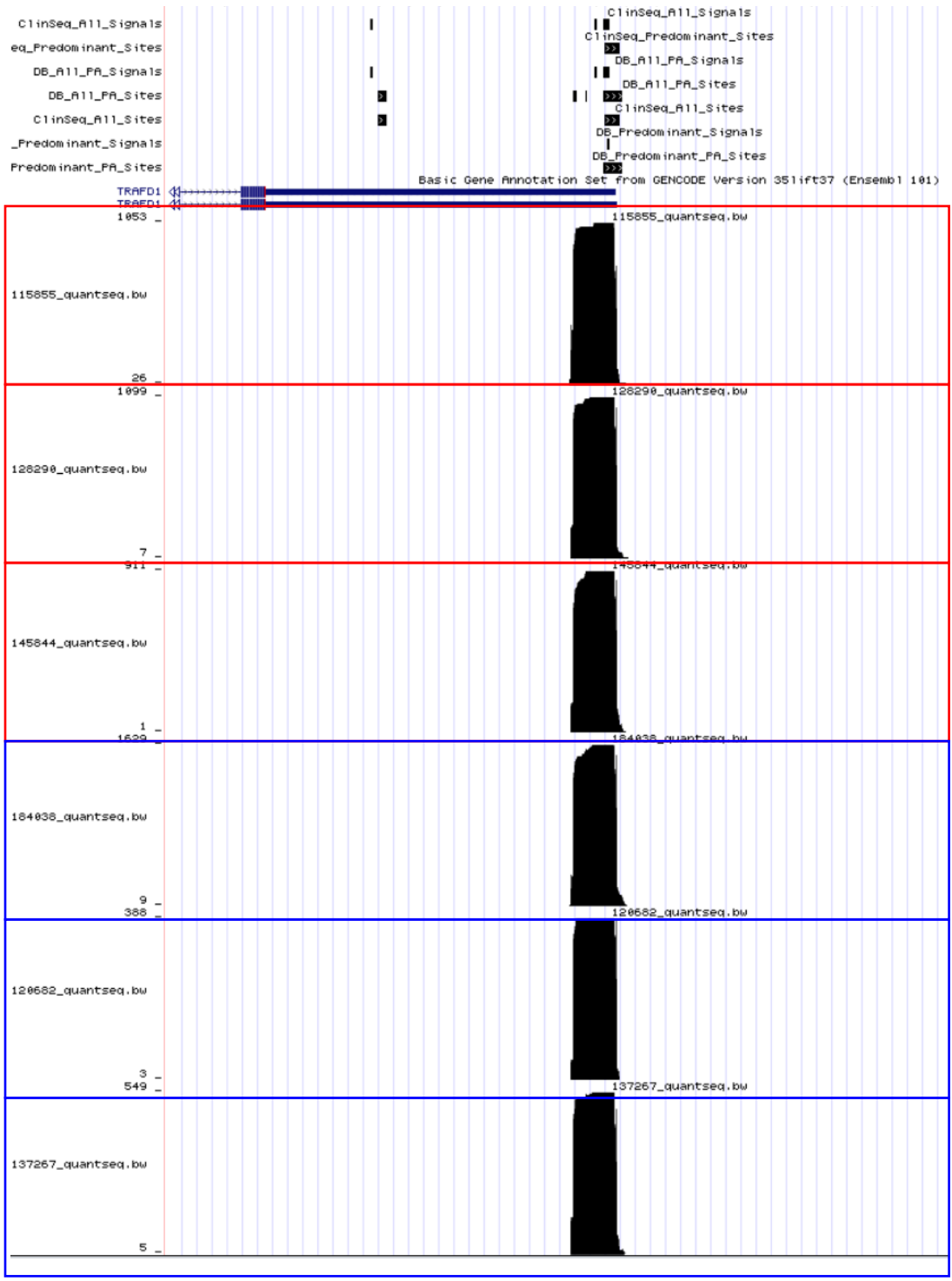


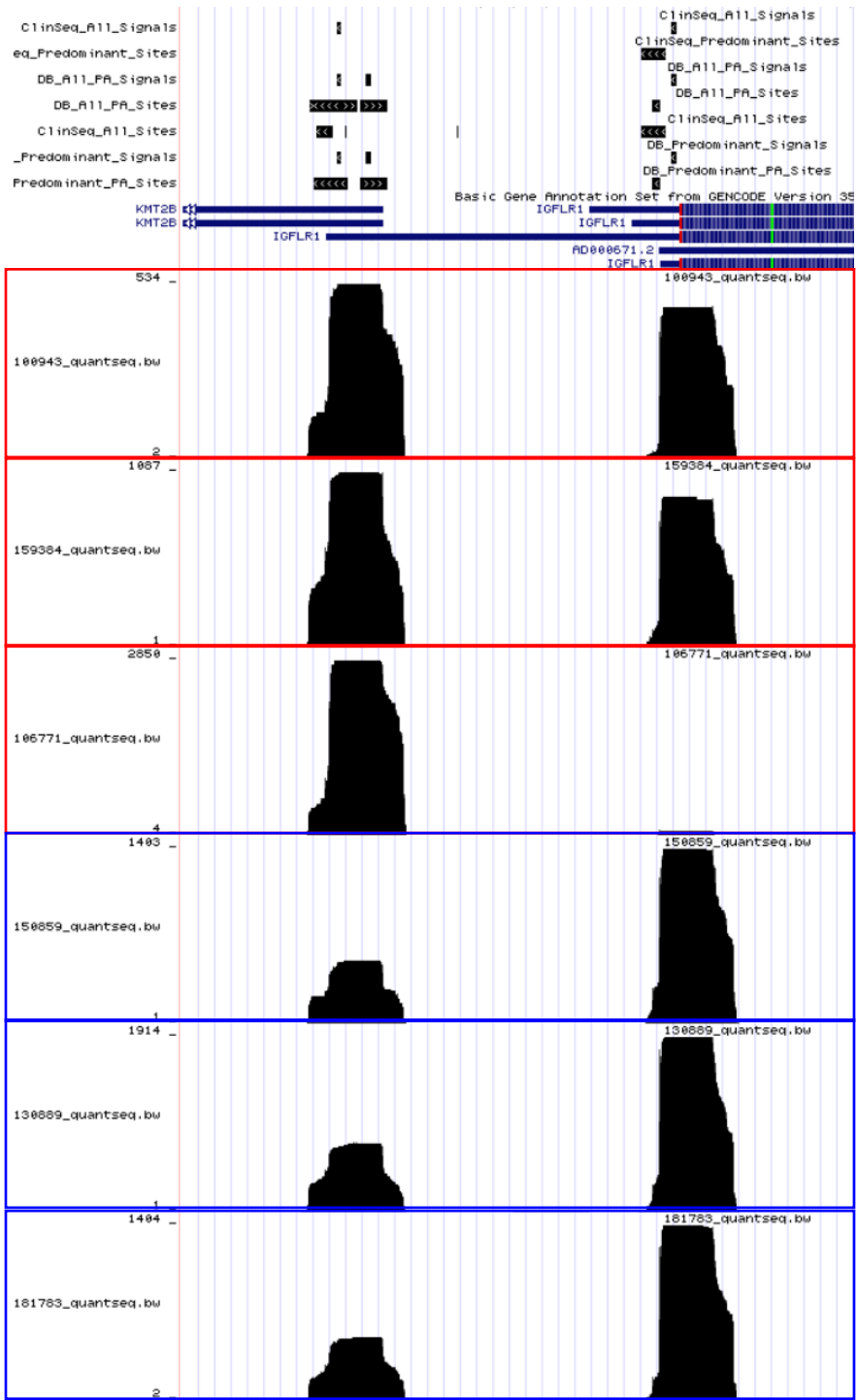
ClinSeq_A11_Signals
eq_Predominant_Sites
DB_A11_PA_Signals
DB_A11_PA_Sites
ClinSeq_A11_Sites
_Predominant_Signals
Predominant_PA_Sites

ClinSeq_A11_Signals
ClinSeq_Predominant_Sites
DB_A11_PA_Signals
DB_A11_PA_Sites
ClinSeq_A11_Sites
DB_Predominant_Signals
DB_Predominant_PA_Sites

Basic Gene Annotation Set from GENCODE Version 31
CXorf38
CXorf38



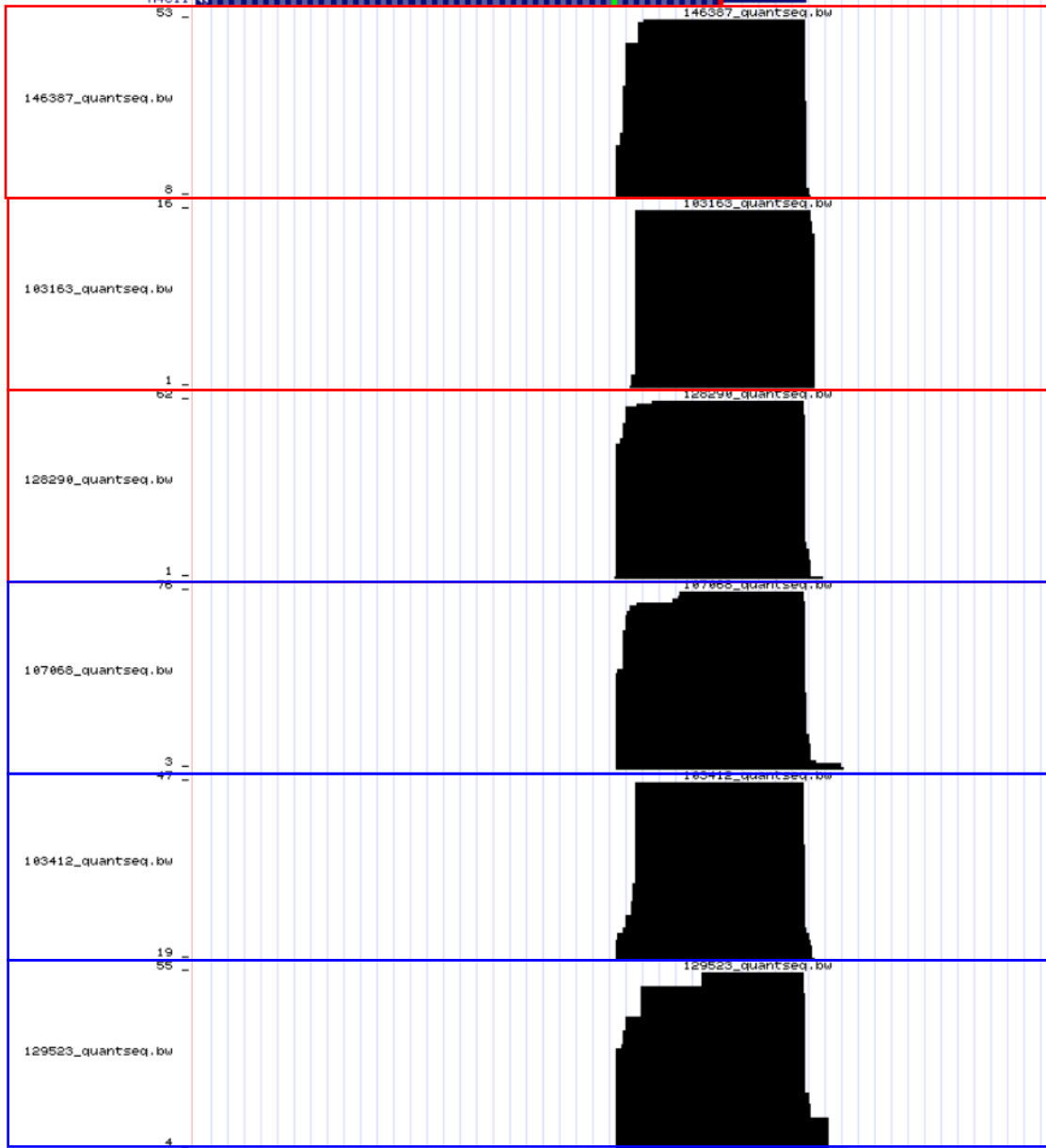


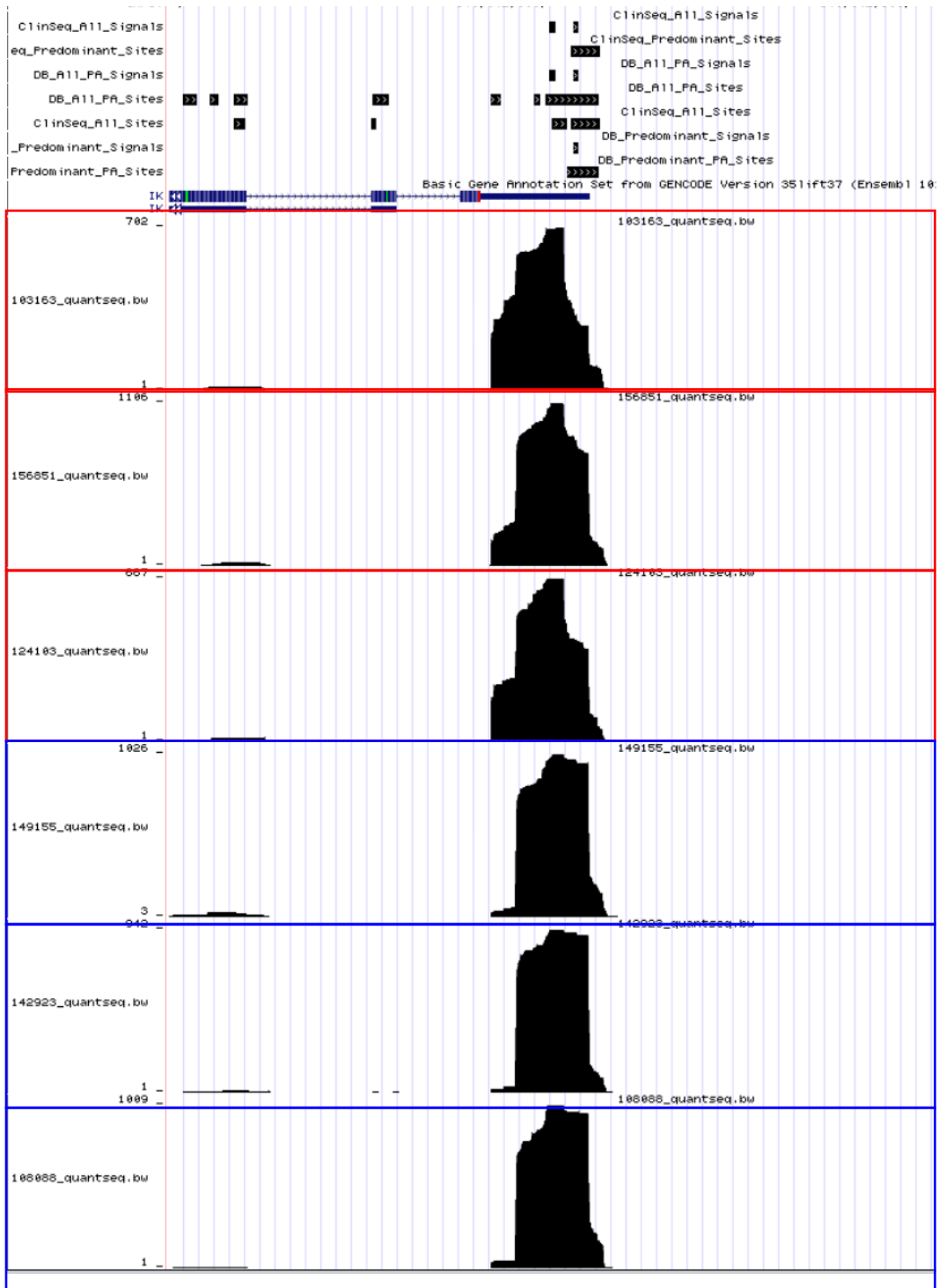


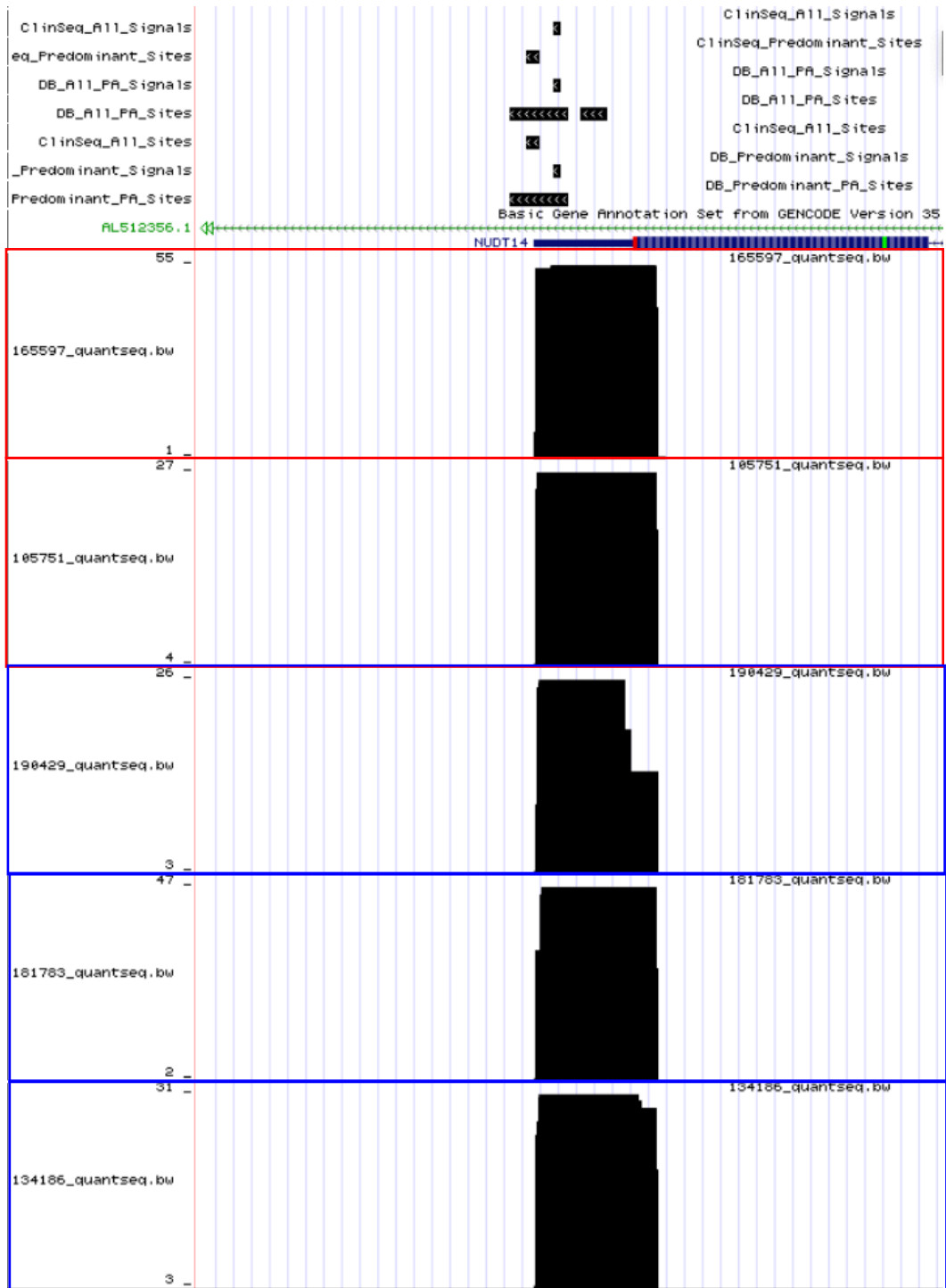
ClinSeq_A11_Signals
eq_Predominant_Sites
DB_A11_PA_Signals
DB_A11_PA_Sites
ClinSeq_A11_Sites
_Predominant_Signals
Predominant_PA_Sites

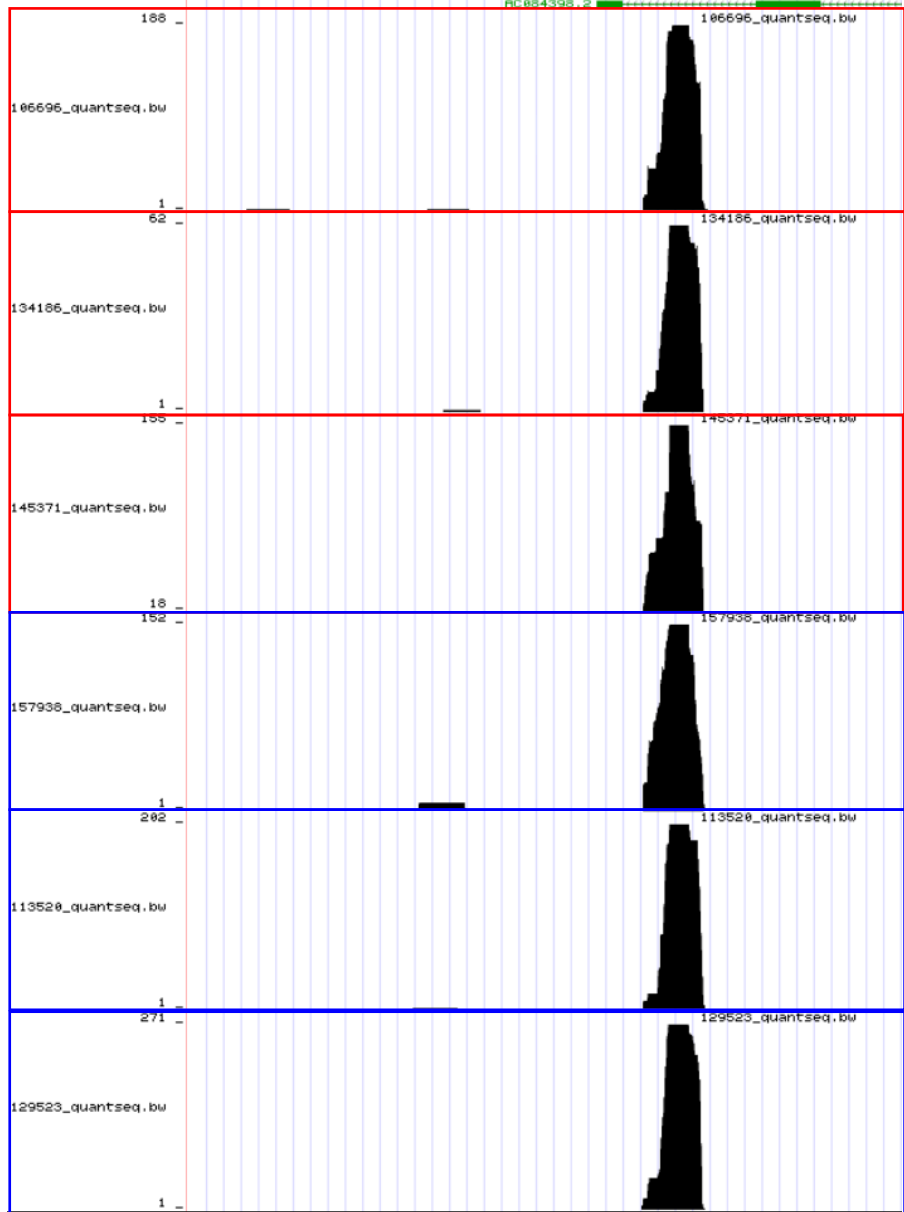
ClinSeq_A11_Signals
ClinSeq_Predominant_Sites
DB_A11_PA_Signals
DB_A11_PA_Sites
ClinSeq_A11_Sites
DB_Predominant_Signals
DB_Predominant_PA_Sites

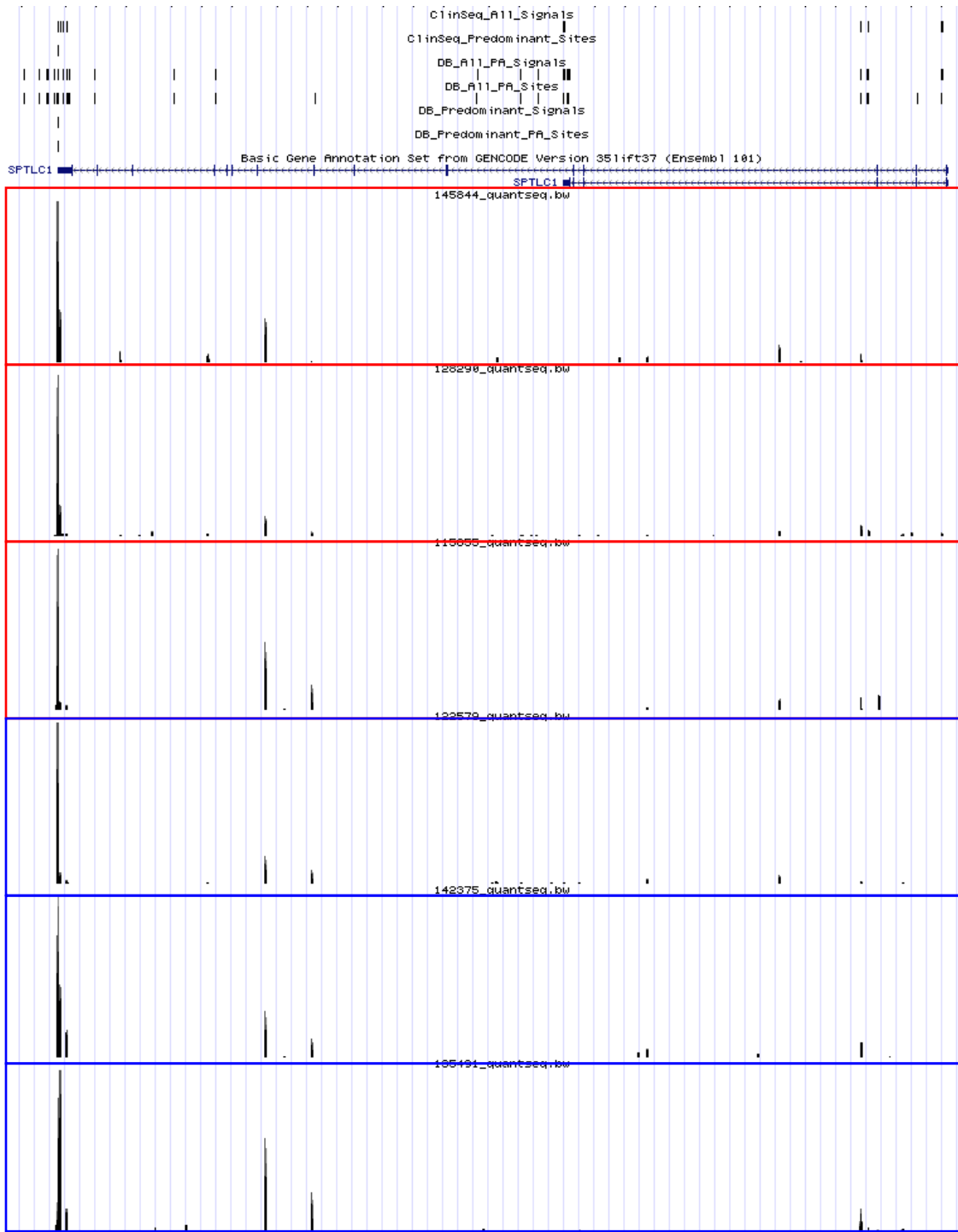
Basic Gene Annotation Set from GENCODE Version 35 lift37 (Ensembl 101)

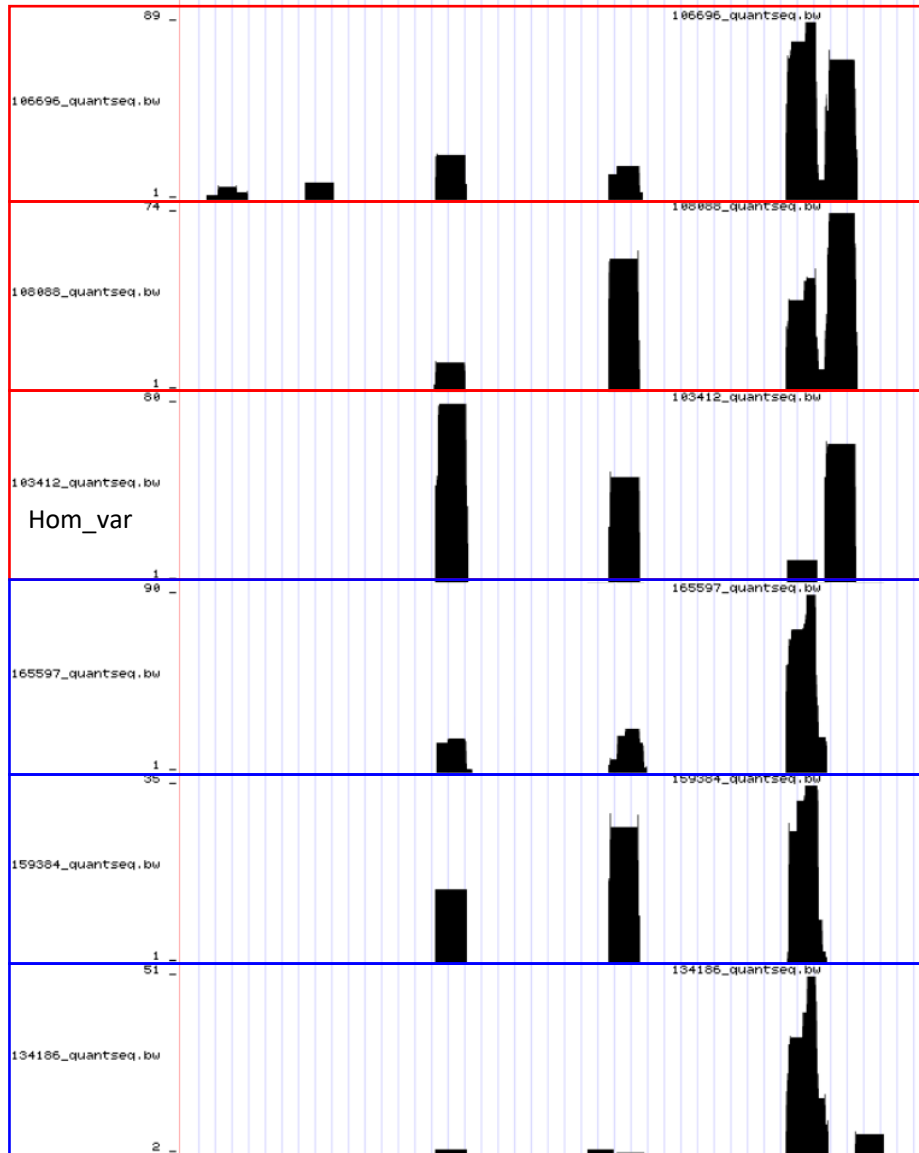
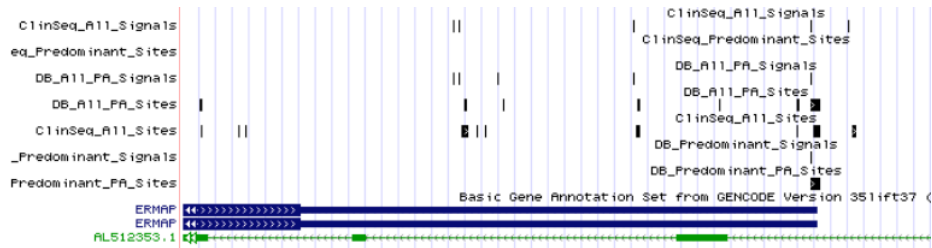


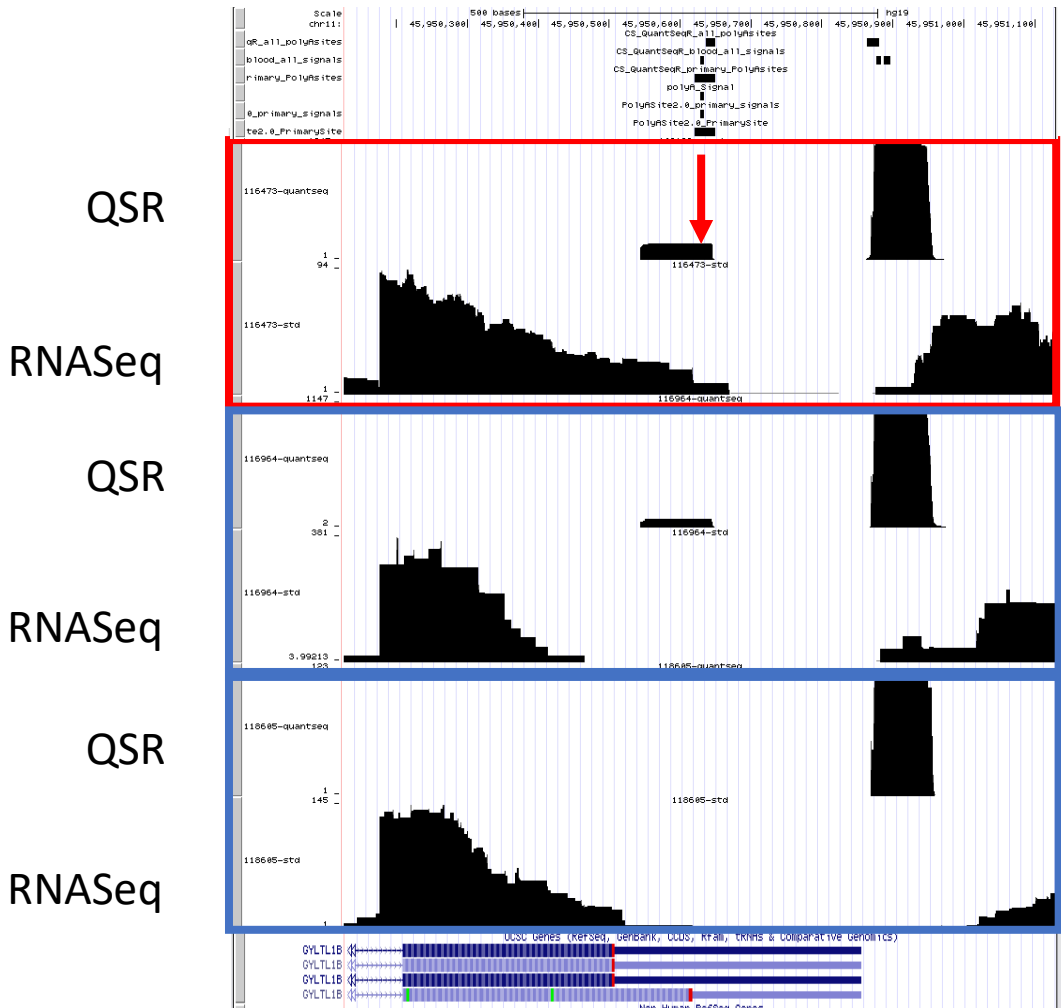




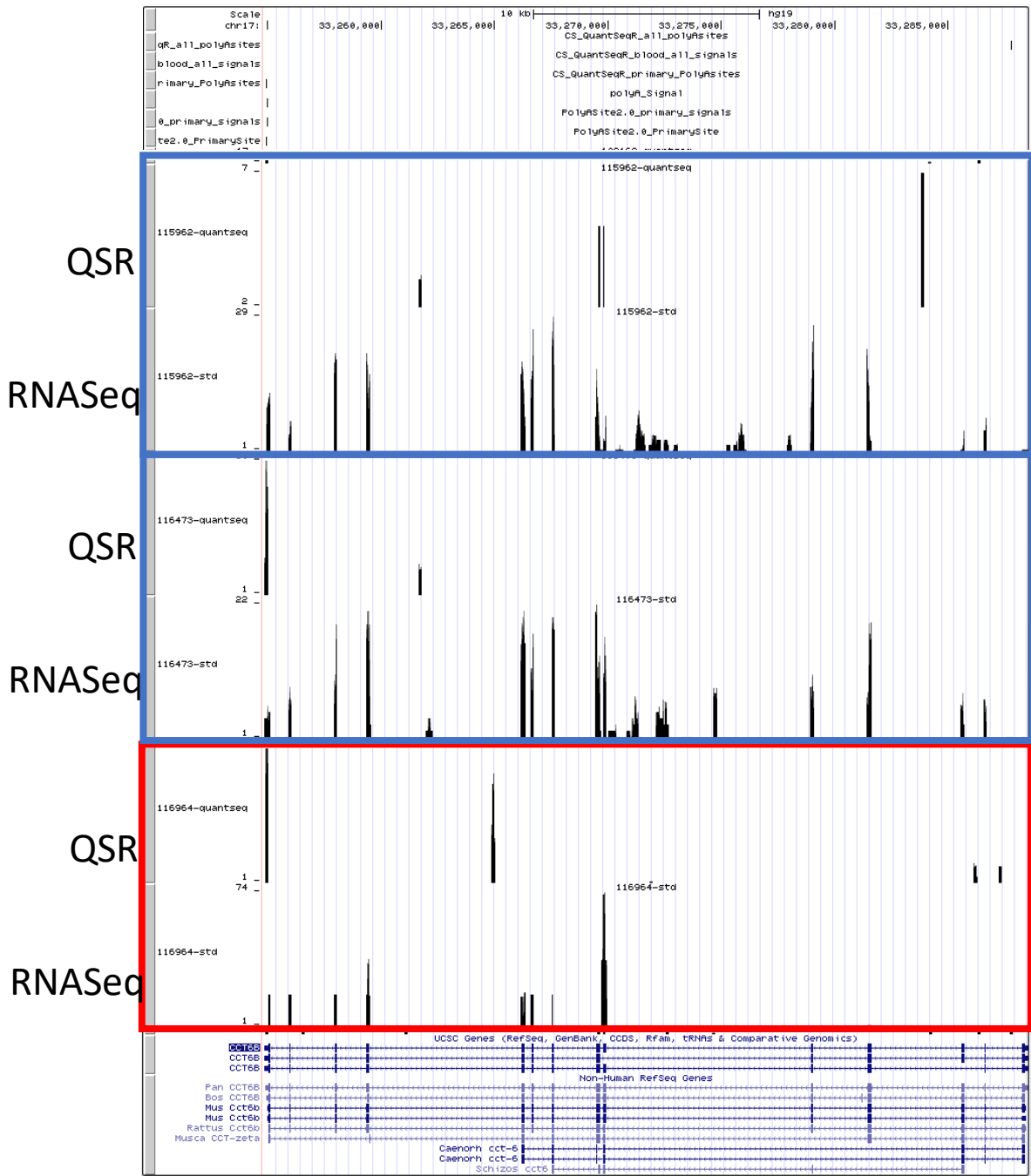




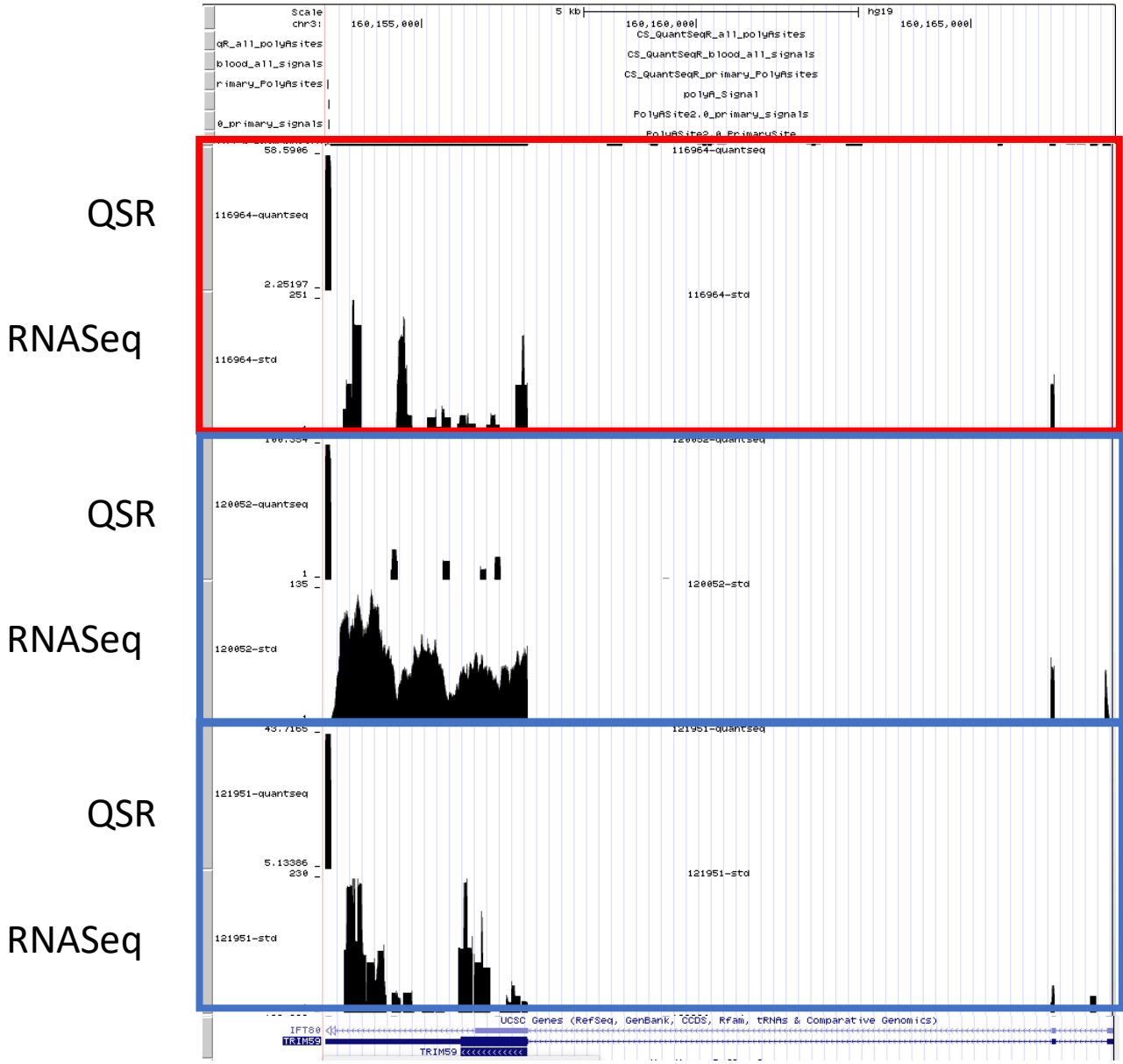




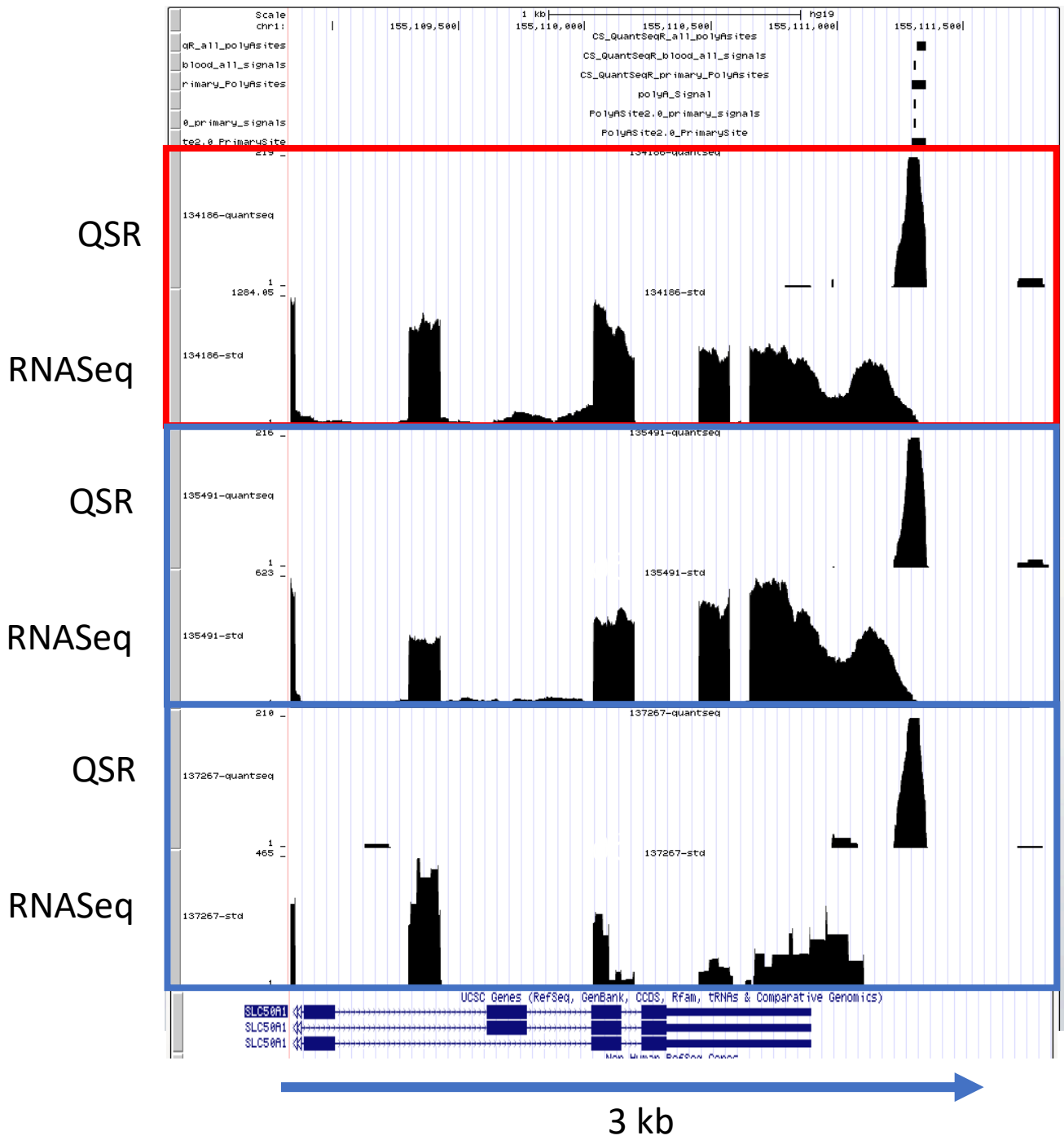
1 kb

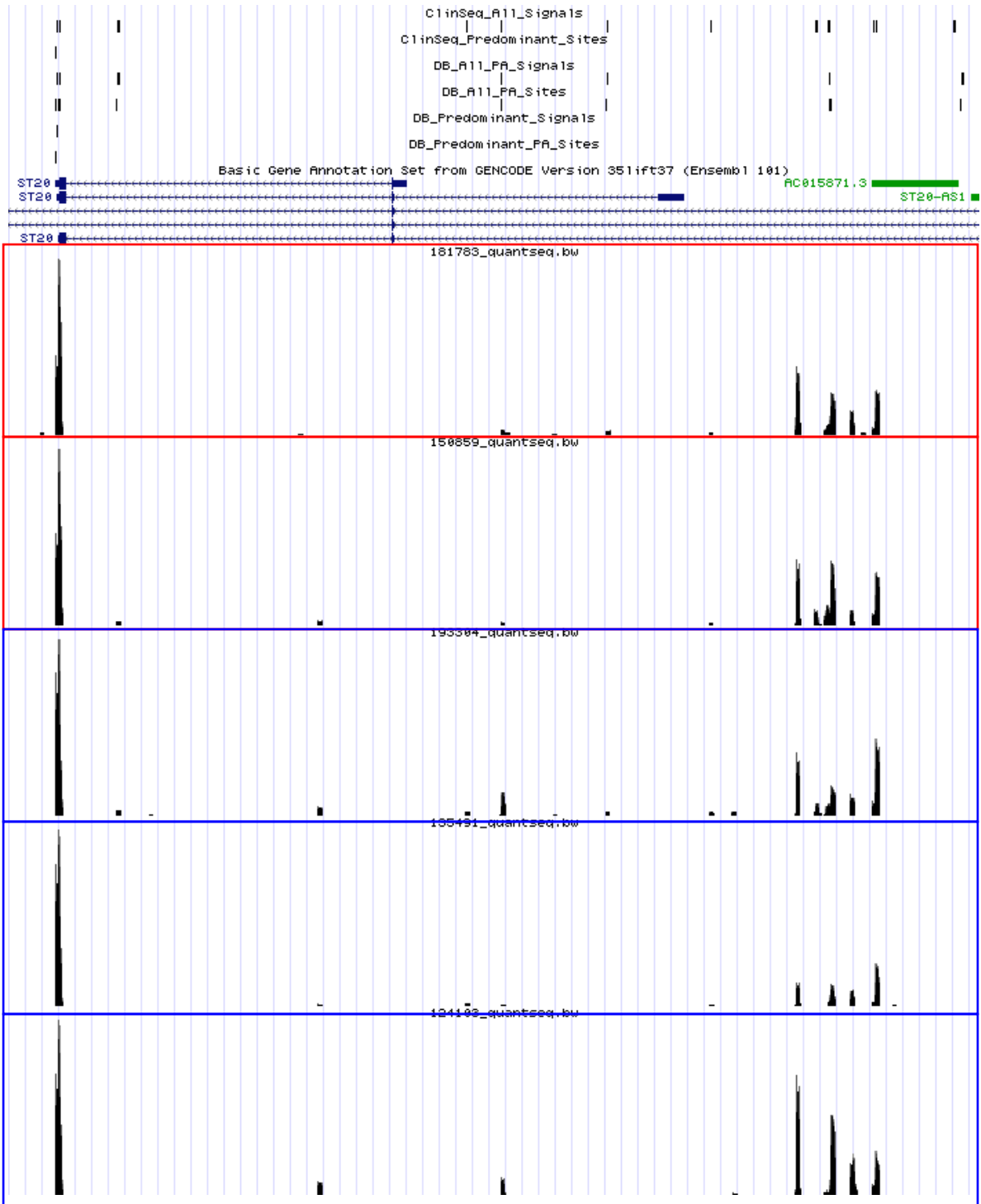


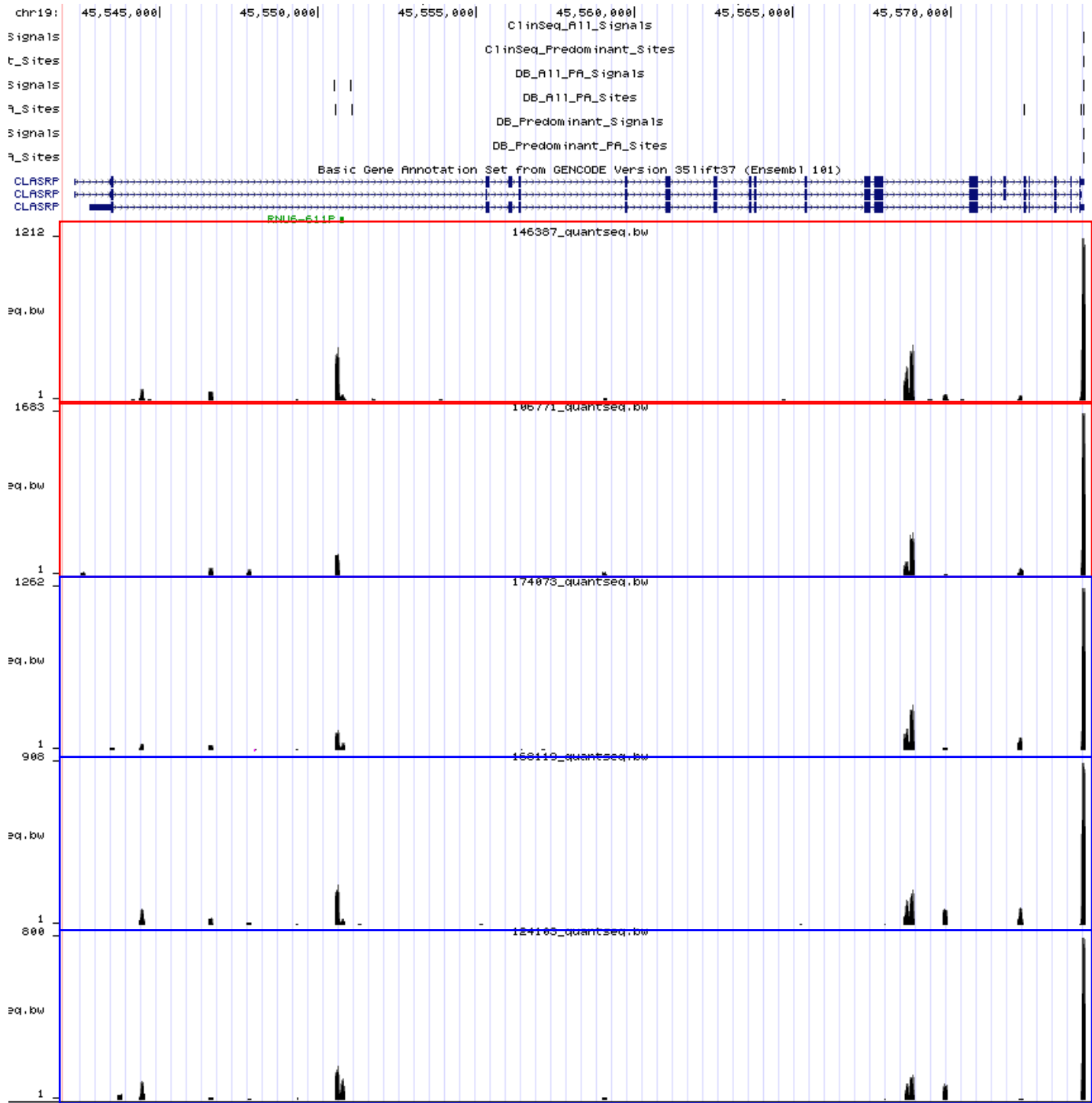
← 33.5 kb

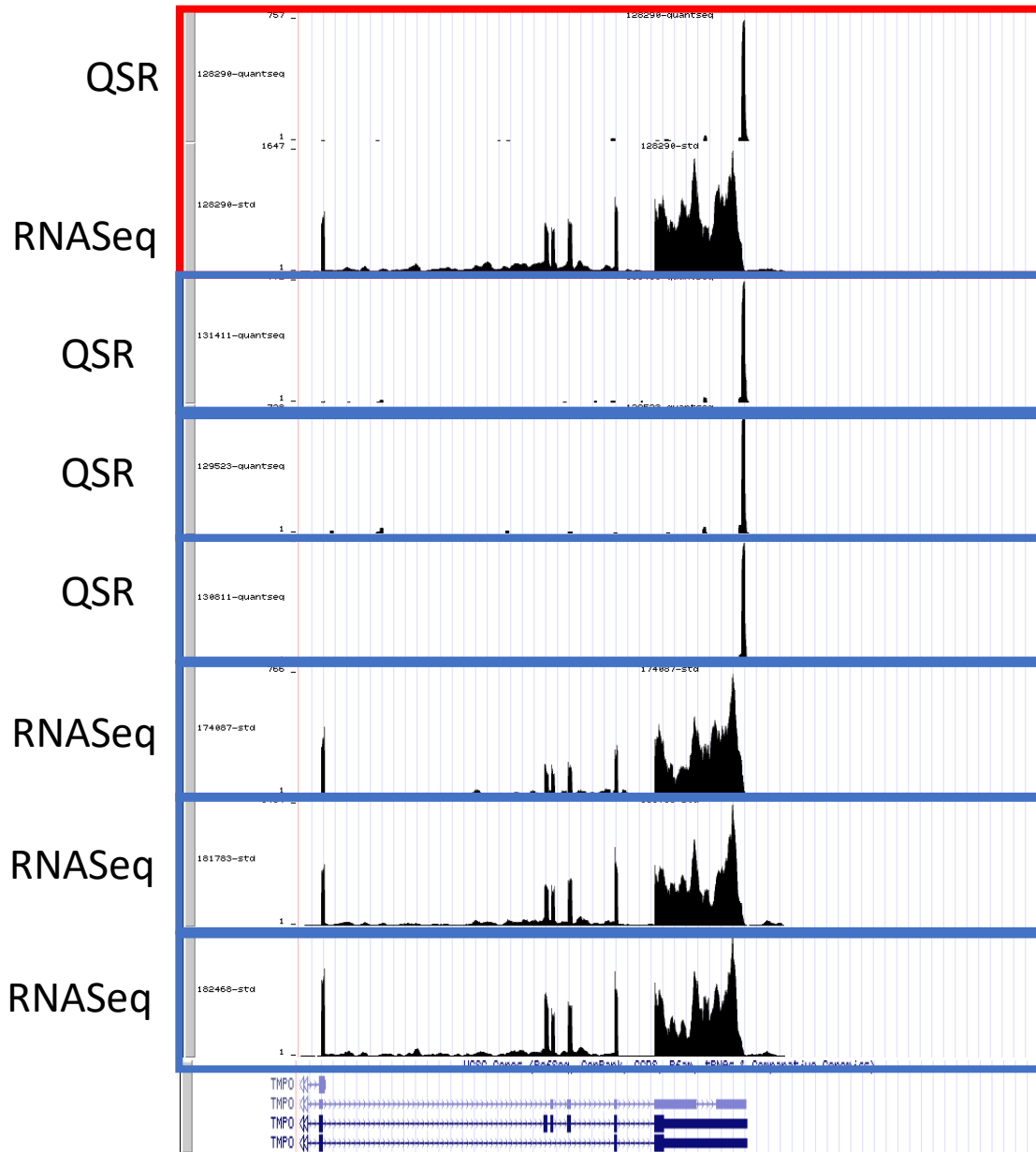


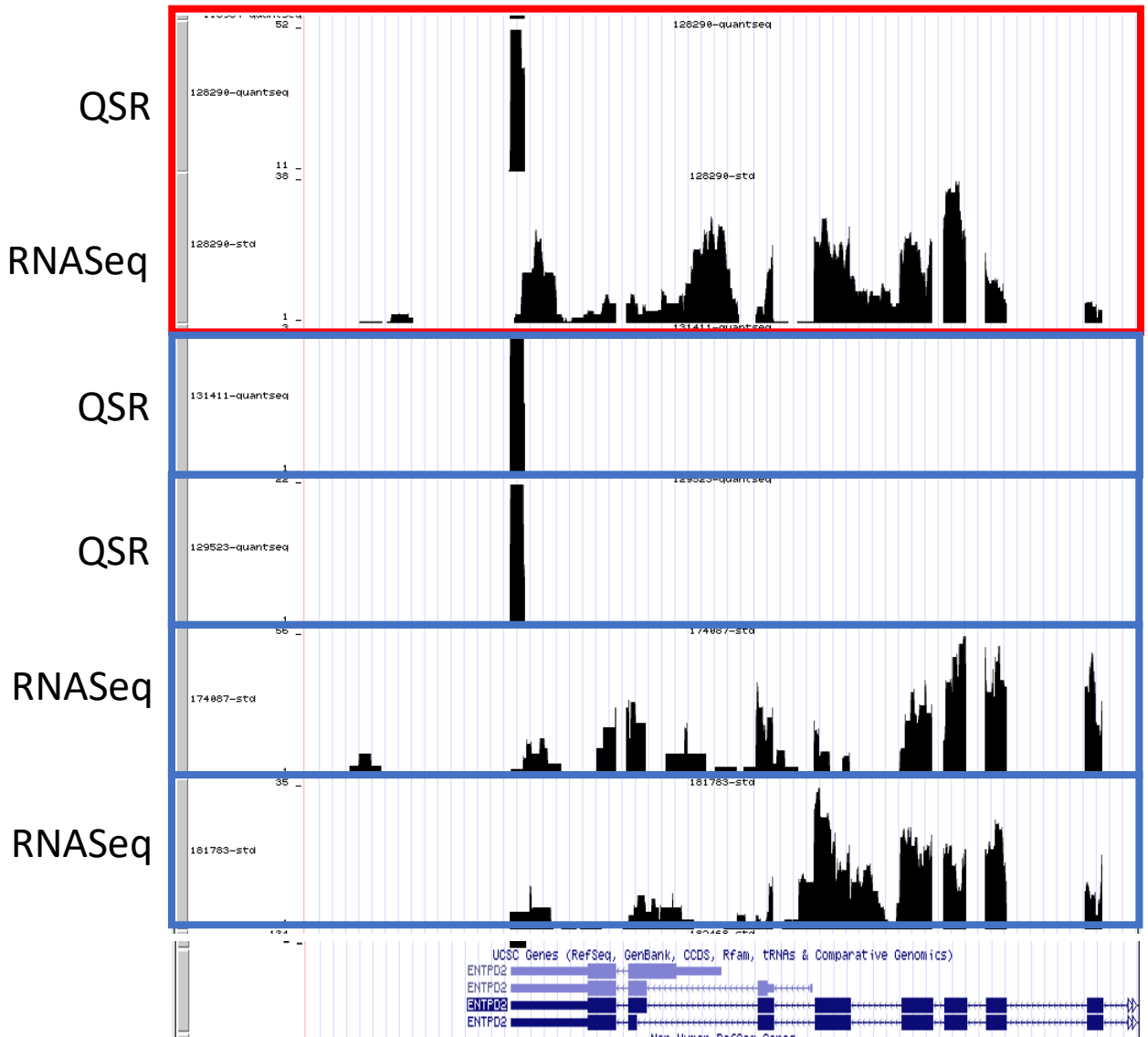
← 14.3 kb

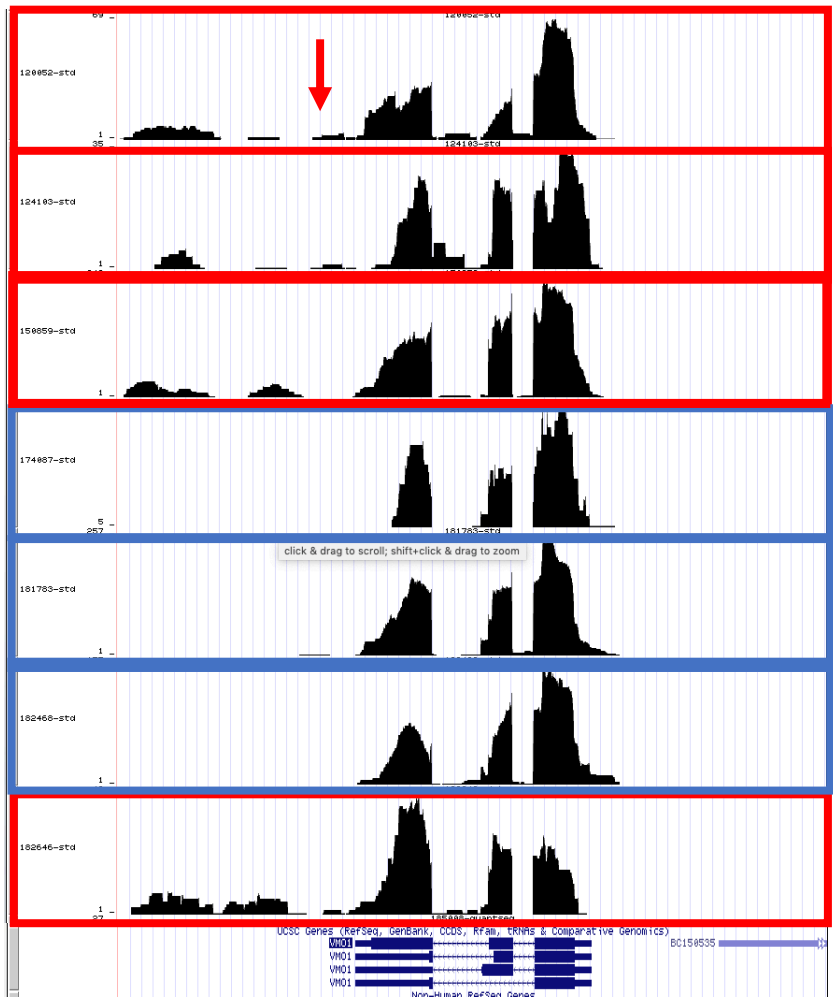


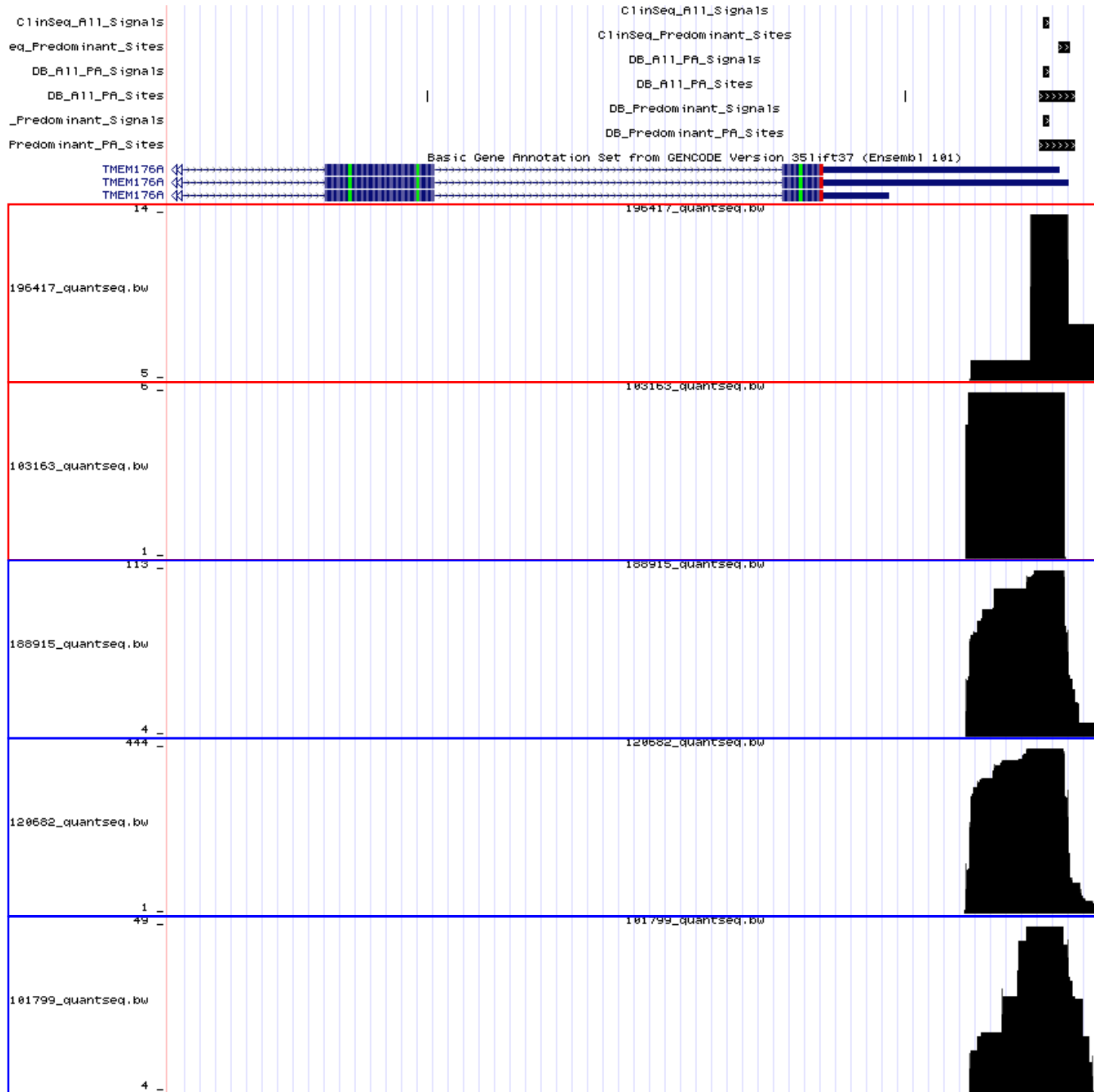


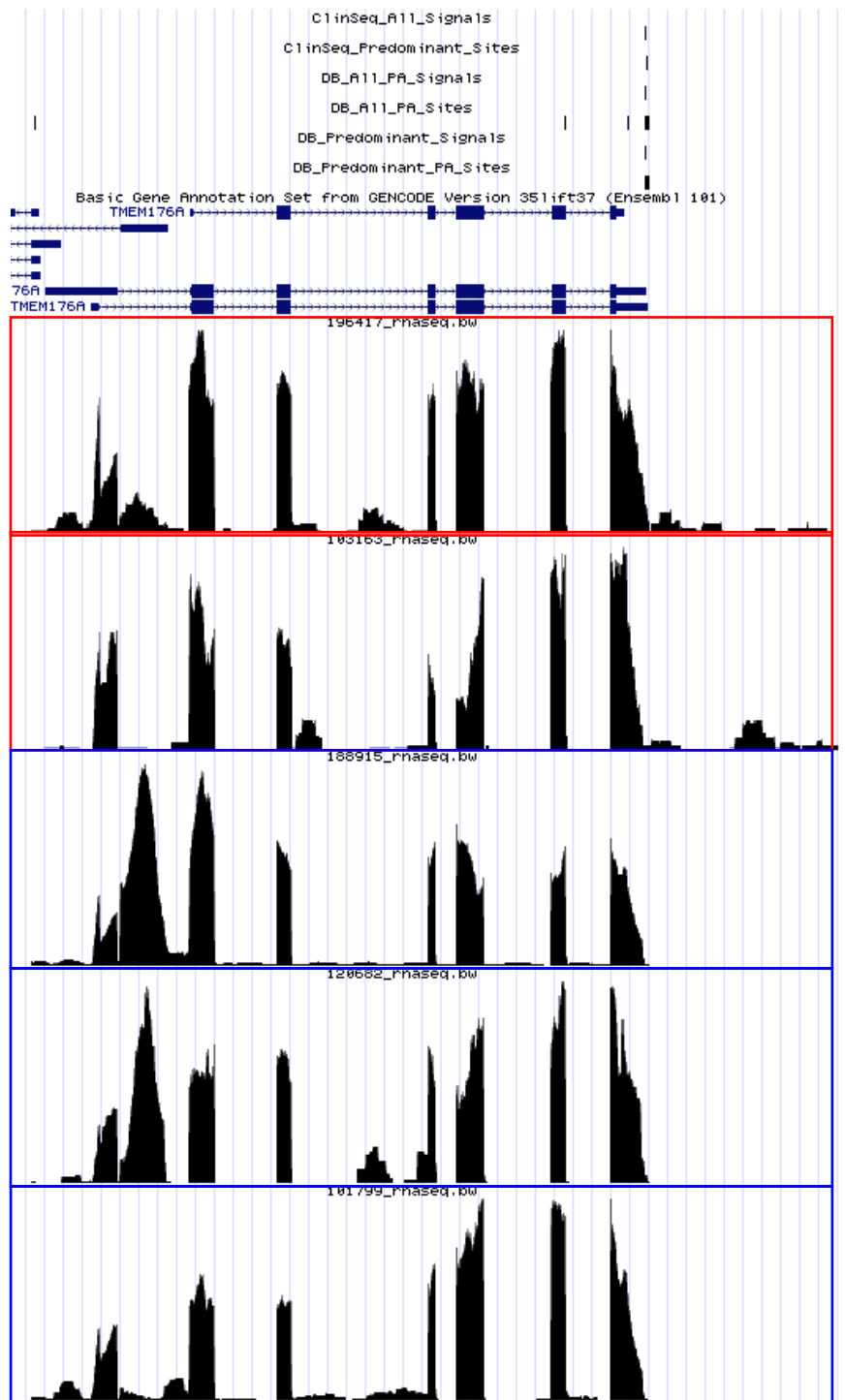


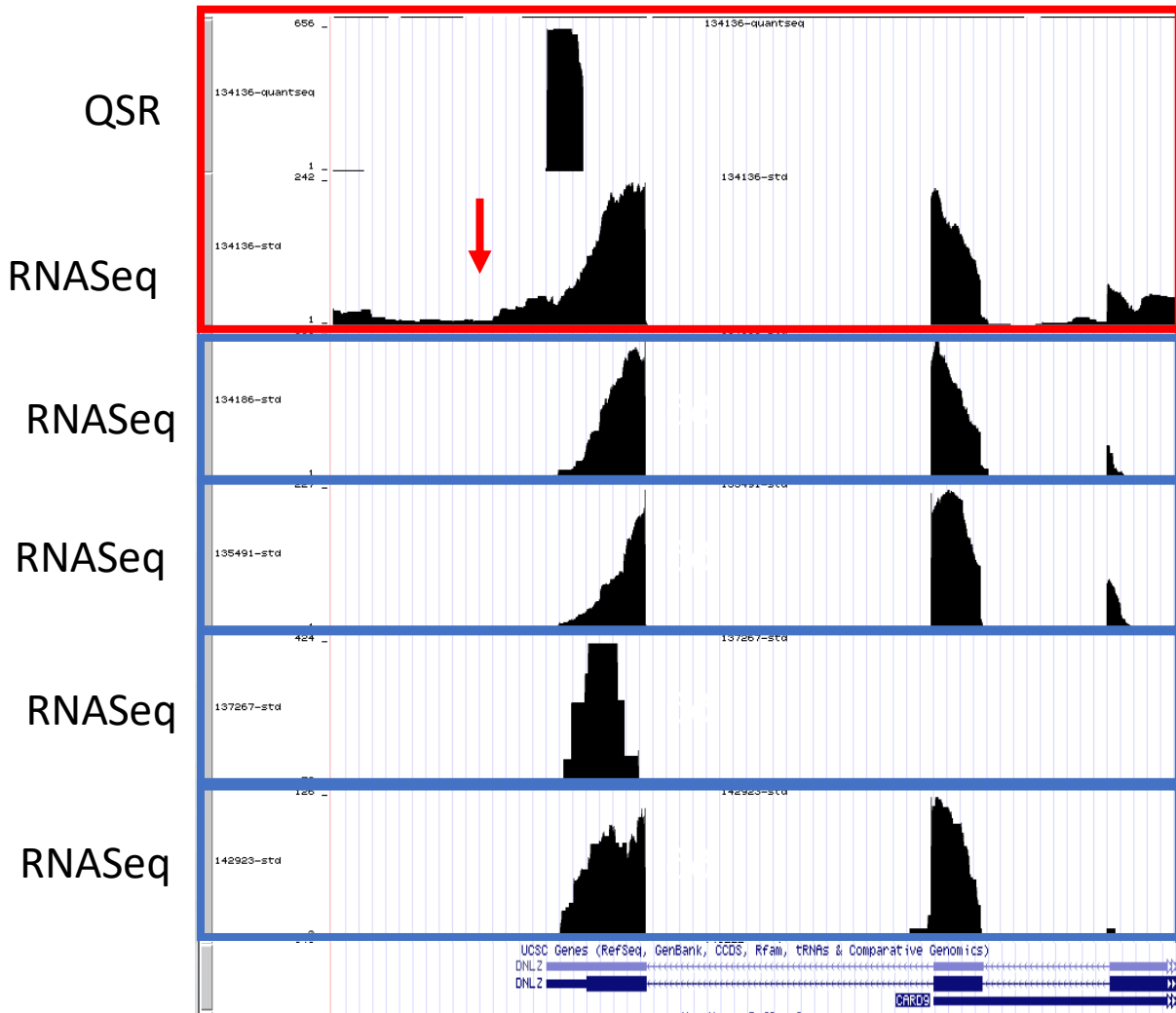


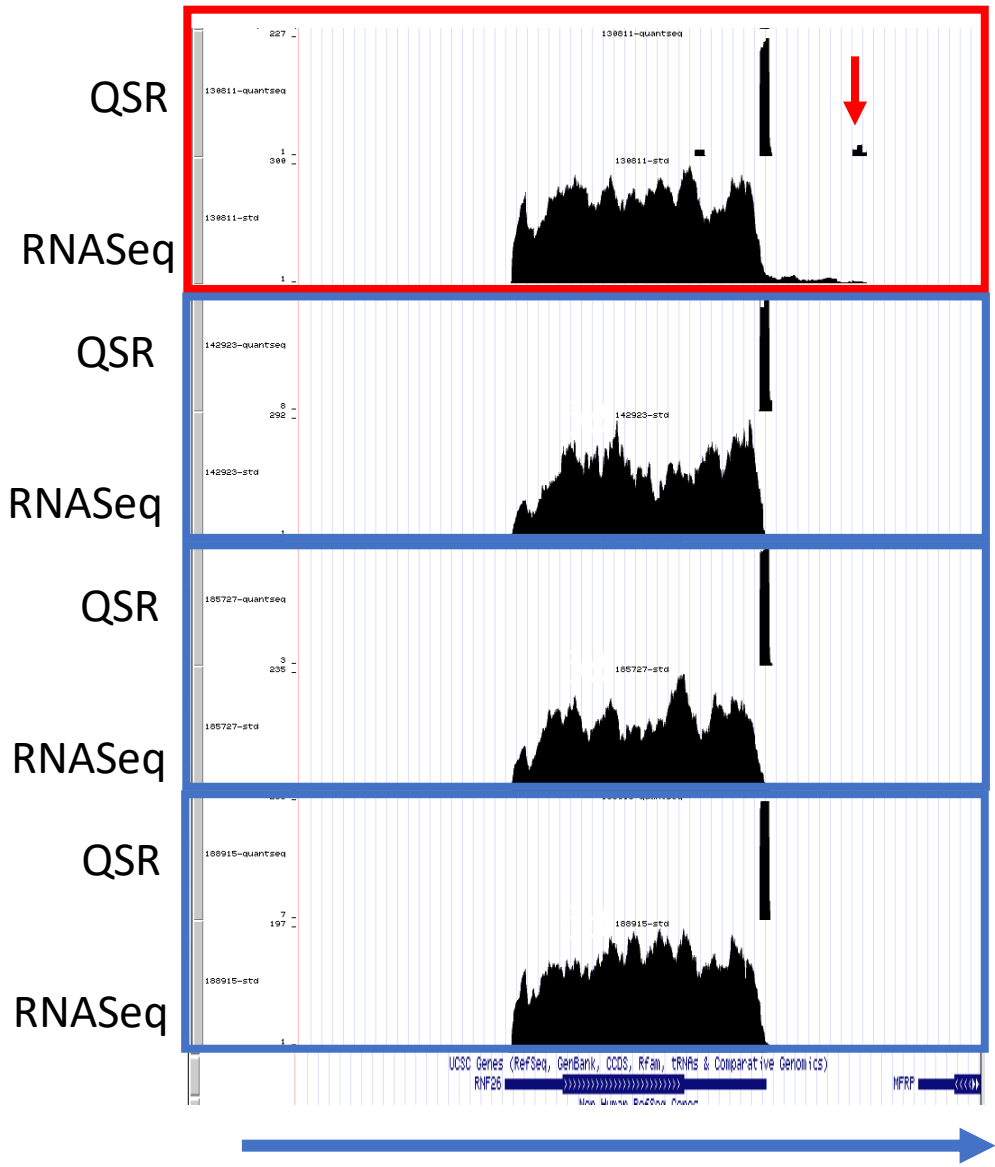


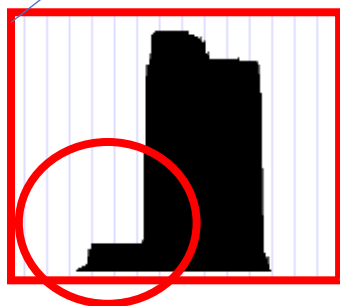
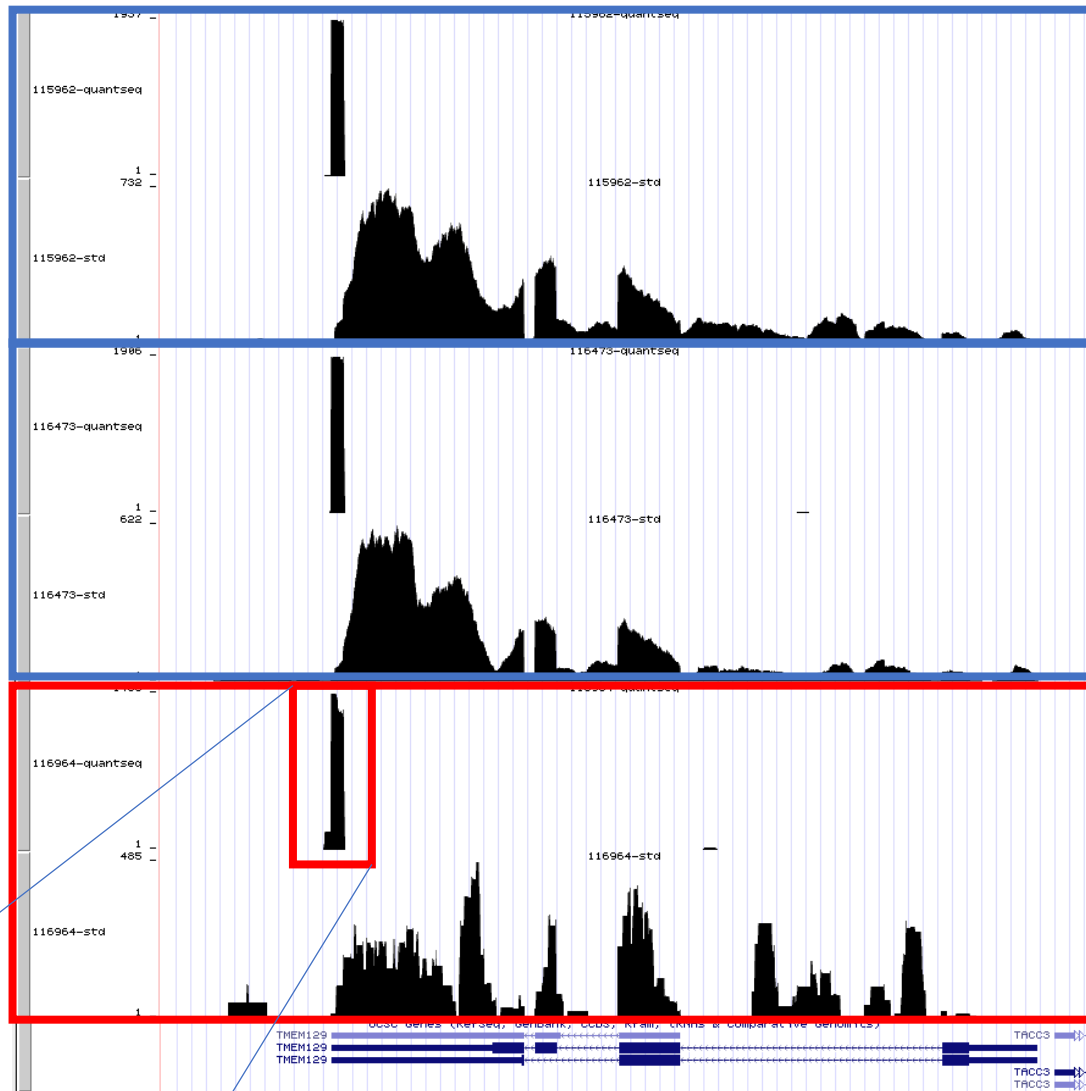




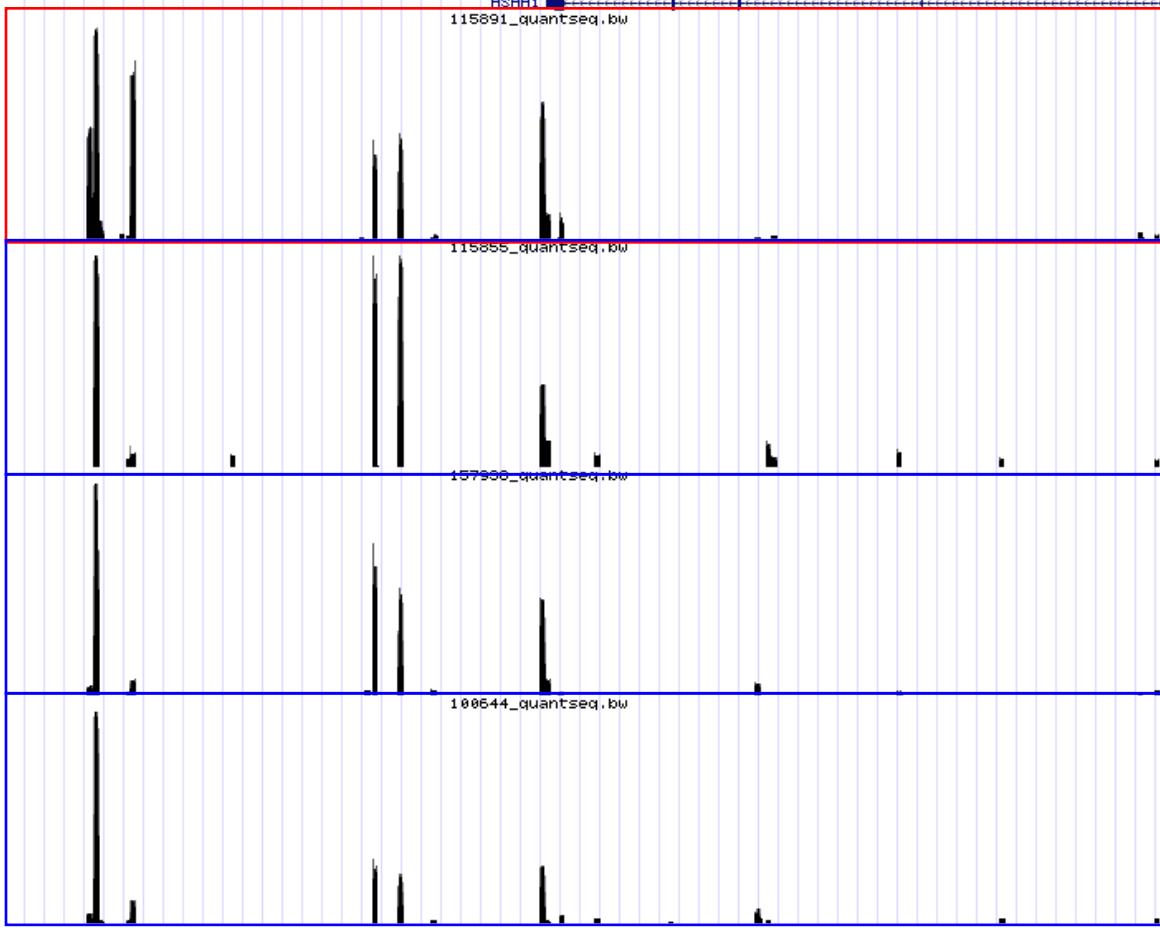
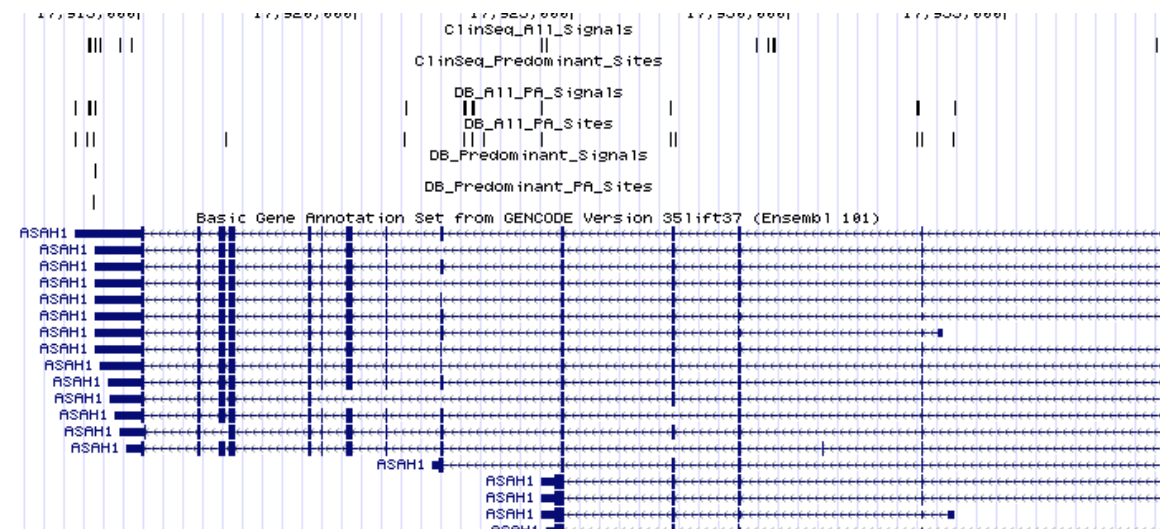


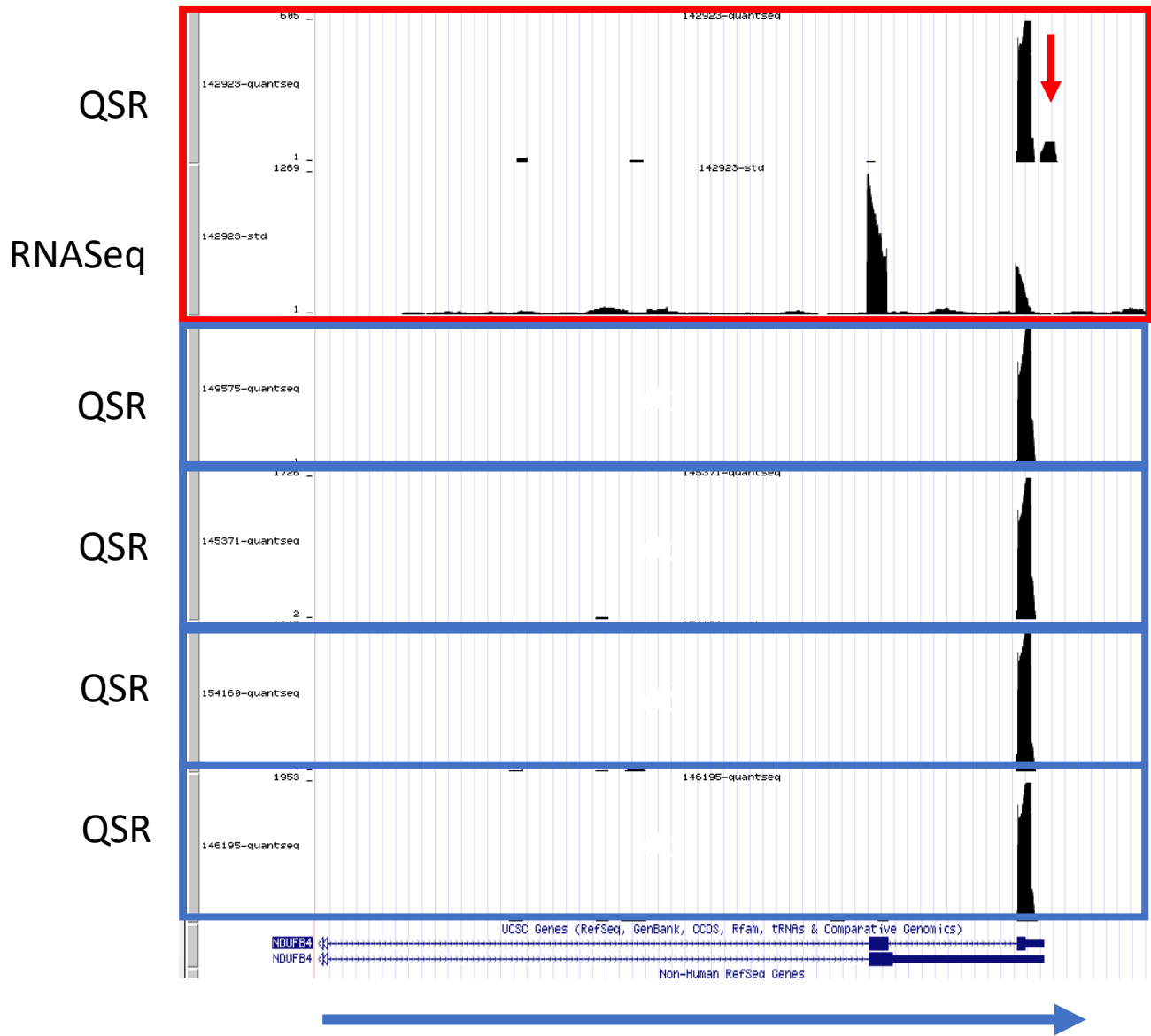






5.5 kb





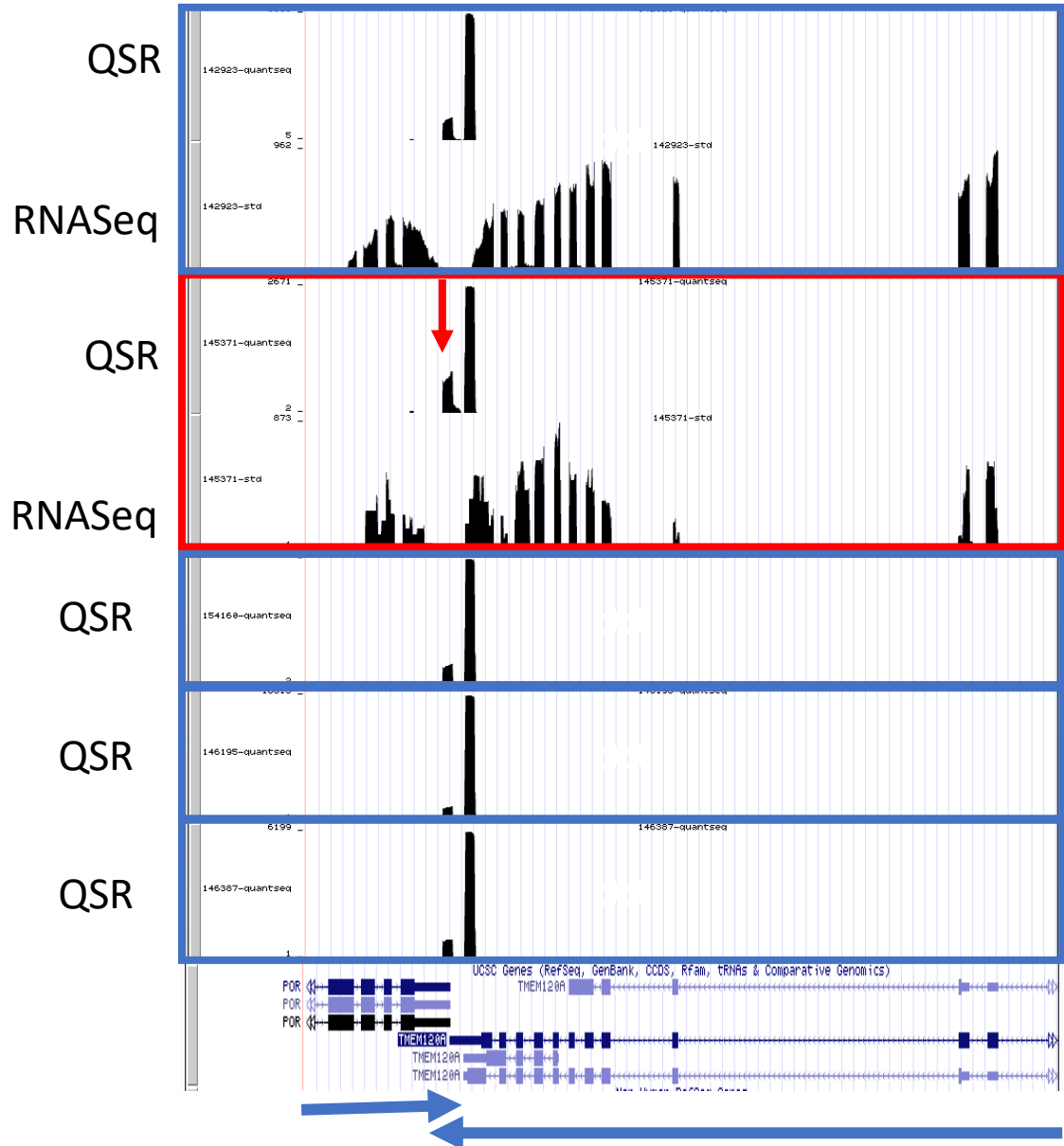


Table S1. The 18 recognized hexamers in PolyA Site 2.0

AATAAA
ATTAAA
TATAAA
AGTAAA
AATACA
CATAAA
AATATA
GATAAA
AATGAA
AATAAT
AAGAAA
ACTAAA
AATAGA
ATTACA
AACAAA
ATTATA
AACAAG
AATAAG

Table S2: Explanation of modified ACMG criteria for polyadenylation hexamer variants.

VERY STRONG EVIDENCE OF PATHOGENICITY

PVS1	Null variant (nonsense, frameshift, canonical +/- 1 or 2 splice sites, initiation codon, single or multi-exon deletion) in a gene where loss of function (LOF) is a known mechanism of disease. Poly A: PVS1 is not applicable, non-coding.
PS2/PM6_ Very Strong	<i>De novo</i> in a patient with the disease and no family history. Counts BOTH proven and unproven <i>de novo</i> cases. Note: Confirmation of paternity only is insufficient. Egg donation, surrogate motherhood, errors in embryo transfer, etc. can contribute to non-maternity. Poly A: PS2/PM6 follow SVI recommendation for <i>de novo</i> criteria. ^a Each proven <i>de novo</i> case gets 2 points, each unproven <i>de novo</i> case gets 1 point, PS2/PM6_Very Strong applied if ≥8 points.

STRONG EVIDENCE OF PATHOGENICITY

PS1	Same amino acid change as a previously established pathogenic variant regardless of nucleotide change. Poly A: PS1 is not applicable, non-coding.
------------	---

PS2/PM6_ Strong	<p><i>De novo</i> in a patient with the disease and no family history. Counts BOTH proven and unproven <i>de novo</i> cases. Note: Confirmation of paternity only is insufficient. Egg donation, surrogate motherhood, errors in embryo transfer, etc. can contribute to non-maternity. Poly A: PS2/PM6 follow SVI recommendation for <i>de novo</i> criteria.^a Each proven <i>de novo</i> case gets 2 points, each unproven <i>de novo</i> case gets 1 point, PS2/PM6 Strong applied if 4-7 points.</p>
PS3	<p>Well-established <i>in vitro</i> or <i>ex vivo</i> functional studies or knock-in mouse studies supportive of a damaging effect on the gene or gene product. Poly A: See PS3_Moderate and PS3_Supporting.</p>
PS4	<p>The prevalence of the variant in affected individuals is significantly increased compared to the prevalence in controls. Poly A: If ClinGen VCEP guidelines are available for the gene in question they should be used. If not the following general guidelines can be considered.</p> <ul style="list-style-type: none"> • PS4_Strong requires ≥ 7 unrelated cases. Popmax MAF in gnomAD ≤ 0.00006. • For variants with popmax MAF in gnomAD > 0.00006, and below the BA1 cutoff, MedCalcs online calculator can be used to calculate the OR using cases from the literature, these numbers are based on an approximation of 3,000 cases (6,000 alleles) reported in the literature and allele counts from gnomAD (MedCalc; https://www.medcalc.net/statisticaltests/odds_ratio.php). An OR of ≥ 18.7 is required for PS4_Strong. • Proband with multiple variants in the gene in question classified as VUS, likely pathogenic or pathogenic are not considered.
PP1_Strong	<p>Co-segregation with disease in multiple affected family members. Poly A: ≥ 7 meioses, only consider phenotype positive/variant positive individuals. To use PP1, no phenotype positive/variant negative individuals can be identified in a pedigree.^b</p>

MODERATE EVIDENCE OF PATHOGENICITY

PM1	<p>Located in a mutational hot spot and/or critical and well-established functional domain (e.g., active site of an enzyme) without benign variation. Poly A: See PM1_Supporting</p>
PM2	<p>Absent from controls (or at extremely low frequency if recessive) in Exome Sequencing Project, 1000 Genomes or ExAC. Poly A: PM2 is not used alone. However, if PS4 is utilized according to a ClinGen VCEP guideline, PM2 should also be considered using the specified guidelines.</p>
PM3	<p>For recessive disorders, detected in trans with a pathogenic variant Note: This requires testing of parents (or offspring) to determine phase. Poly A: PM3 is used as described by the ClinGen SVI.^a</p>
PM4	<p>Protein length changes due to in-frame deletions/insertions in a non-repeat region or stop-loss variants. Poly A: PM4 is not applicable. If a change in protein length or amount is experimentally demonstrated this should be considered for PS3.</p>
PM5	<p>Missense change at an amino acid residue where a different missense change determined to be pathogenic has been seen before. Poly A: See PM5_Supporting.</p>
PS2/PM6_ Moderate	<p><i>De novo</i> in a patient with the disease and no family history. Counts BOTH proven and unproven <i>de novo</i> cases. Note: Confirmation of paternity only is insufficient. Egg donation, surrogate motherhood, errors in embryo transfer, etc. can contribute to non-maternity. Poly A: PS2/PM6 follow SVI recommendation for <i>de novo</i> criteria.^a Each proven <i>de novo</i> case gets 2 points, each unproven <i>de novo</i> case gets 1 point, PS2/PM6 Moderate applied for 2-3 points.</p>

PS3_ Moderate	Well-established <i>in vitro</i> or <i>ex vivo</i> functional studies or knock-in mouse studies supportive of a damaging effect on the gene or gene product. Poly A: <i>Ex vivo</i> assays showing decreased RNA or protein levels, cell lines from three unrelated individuals who harbor variant tested.
PS4_ Moderate	The prevalence of the variant in affected individuals is significantly increased compared to the prevalence in controls. Poly A: If ClinGen VCEP guidelines are available for the gene in question they should be used. If not the following general guidelines can be considered. <ul style="list-style-type: none"> • PS4_Moderate requires 2-6 unrelated cases. Popmax MAF in gnomAD ≤ 0.00006. • For variants with popmax MAF in gnomAD > 0.00006, and below the BA1 cutoff, MedCalcs online calculator can be used to calculate the OR using case points from the literature, these numbers are based on an approximation of 3,000 cases (6,000 alleles) reported in the literature and allele counts from gnomAD (MedCalc; https://www.medcalc.net/statisticaltests/odds_ratio.php). An OR of ≥ 4.33 is required for PS4_Moderate. • Probands with multiple variants in the gene in question classified as VUS, likely pathogenic or pathogenic are not considered.
PP1_ Moderate	Co-segregation with disease in multiple affected family members. Poly A: 5-6 meioses, only consider phenotype positive/variant positive individuals. ^b In order to use PP1 no phenotype positive/variant negative individuals can be identified in a pedigree.

SUPPORTING EVIDENCE OF PATHOGENICITY

PP1	Co-segregation with disease in multiple affected family members. Poly A: 3-4 meioses, only consider phenotype positive/variant positive individuals. ^b In order to use PP1 no phenotype positive/variant negative individuals can be identified in a pedigree.
PP2	Missense variant in a gene that has a low rate of benign missense variation and where missense variants are a common mechanism of disease. Poly A: PP2 is not applicable.
PP3	Multiple lines of computational evidence support a deleterious effect on the gene or gene product (conservation, evolutionary, splicing impact, etc.). Poly A: CADD score > 10 .
PP4	Patient's phenotype or family history is highly specific for a disease with a single genetic etiology. Poly A: PP4 is used as described in Richards et al. (2015).
PP5	Reputable source recently reports variant as pathogenic but the evidence is not available to the laboratory to perform an independent evaluation. Poly A: PP5 has been dropped from the ACMG framework for variant assessment.
PS2/PM6_ Supporting	<i>De novo</i> in a patient with the disease and no family history. Counts BOTH proven and unproven <i>de novo</i> cases. Note: Confirmation of paternity only is insufficient. Egg donation, surrogate motherhood, errors in embryo transfer, etc. can contribute to non-maternity. Poly A: PS2/PM6 follow SVI recommendation for <i>de novo</i> criteria. ^a Each proven <i>de novo</i> case gets 2 points, each unproven <i>de novo</i> case gets 1 point, PS2/PM6 Supporting applied for 1 point.
PS3_ Supporting	Well-established <i>in vitro</i> or <i>ex vivo</i> functional studies or knock-in mouse studies supportive of a damaging effect on the gene or gene product. Poly A: Two independent <i>ex vivo</i> studies (tissues from unrelated individuals) all showing decreased RNA or protein levels.
PS4_ Supporting	The prevalence of the variant in affected individuals is significantly increased compared to the prevalence in controls. Poly A: If ClinGen VCEP guidelines are available for the gene in question they should be used. If not the following general guidelines can be considered. <ul style="list-style-type: none"> • PS4_Supporting requires one case. Popmax MAF in gnomAD ≤ 0.00006. • For variants with popmax MAF in gnomAD > 0.00006, and below BA1 cutoff of 0.0038, MedCalcs online calculator can be used to calculate the OR using case points from the literature, these numbers are based on an approximation of 3,000 cases (6,000 alleles) reported

	<p>in the literature and allele counts from gnomAD (MedCalc; https://www.medcalc.net/statisticaltests/odds_ratio.php). An OR of ≥ 2.08 is required for PS4_Supporting.</p> <ul style="list-style-type: none"> • Probands with multiple variants in the gene in question classified as VUS, likely pathogenic or pathogenic are not considered.
PM1_ Supporting	<p>Located in a mutational hot spot and/or critical and well-established functional domain (e.g., active site of an enzyme) without benign variation. Poly A: Presence of a variant in a predominant polyadenylation hexamer.</p>
PM3_ Supporting	<p>For recessive disorders, variants in the homozygous state Poly A: PM3_Supporting is used as described by the ClinGen SVI.^a</p>
PM5_ Supporting	<p>Missense change at an amino acid residue where a different missense change determined to be likely pathogenic has been seen before. Poly A: PM5 has been modified to allow for use for non-missense variation. <ul style="list-style-type: none"> • Variation in a hexamer where a previously pathogenic variant has been identified can be considered for PM5. The new hexamer must be predicted to be less functional than the original hexamer according to the table from Sheets et al. </p>

STAND ALONE EVIDENCE OF BENIGN IMPACT

BA1	<p>Allele frequency is >0.05 in any general continental population dataset of at least 2,000 observed alleles and found in a gene without a gene- or variant-specific BA1 modification. Poly A: BA1 is used as described in Richards et al. (2015). If specifications have been provided by an expert panel BA1 should be determined as set by the expert panel for the gene.</p>
------------	---

STRONG EVIDENCE OF BENIGN IMPACT

BS1	<p>Allele frequency is greater than expected for disorder. Poly A: BS1 is used as described in Richards et al. (2015). If specifications have been provided by an expert panel BS1 should be determined as set by the expert panel for the gene.</p>
BS2	<p>Observed in a healthy adult individual for a recessive (homozygous), dominant (heterozygous), or X-linked (hemizygous) disorder with full penetrance expected at an early age. Poly A: BS2 is used as described in Richards et al (2015).</p>
BS3	<p>Well-established <i>in vitro</i> or <i>ex vivo</i> functional studies or knock-in mouse studies show no damaging effect on protein function. Poly A: BS3 downgraded to BS3_Moderate/BS3_Supporting.</p>
BS4	<p>Lack of segregation in affected members of a family Poly A: BS4 is used as described in Richards et al. (2015).</p>

MODERATE EVIDENCE FOR BENIGN IMPACT

BS3_ Moderate	<p>Well- established <i>in vitro</i> or <i>ex vivo</i> functional studies or knock-in mouse studies show no damaging effect on protein function. Poly A: Three or more independent <i>ex vivo</i> studies all showing NO significant decrease in protein or RNA level. Cells from unrelated individuals who harbor the variant. <ul style="list-style-type: none"> • <i>Ex vivo</i> studies using patient derived samples need to be interpreted with the understanding that unidentified variants may be present. </p>
----------------------	---

SUPPORTING EVIDENCE FOR BENIGN IMPACT

BP1	Missense variant in a gene for which primarily truncating variants are known to cause disease. Poly A: BP1 is not applicable.
BP2	Observed in trans with a pathogenic variant for a fully penetrant dominant gene/disorder; or observed in cis with a pathogenic variant in any inheritance pattern. Poly A: BP2 is applicable for variants shown to be in cis with a known pathogenic variant.
BP3	In-frame deletions/insertions in a repetitive region without a known function Poly A: BP3 is not applicable.
BP4	Multiple lines of computational evidence suggest no impact on gene or gene product (conservation, evolutionary, splicing impact, etc.). Poly A: A CADD score <5 can be used in support of benign status.
BP5	Variant found in a case with an alternate molecular basis for disease. Poly A: BP5 is applicable as described in Richards et al. (2015).
BP6	Reputable source recently reports variant as benign, but the evidence is not available to the laboratory to perform an independent evaluation. Poly A: BP6 has been dropped from the ACMG framework for variant assessment.
BP7	A synonymous (silent) variant for which splicing prediction algorithms predict no impact to the splice consensus sequence nor the creation of a new splice site AND the nucleotide is not highly conserved. Poly A: BP7 is not applicable.
BS3_ Supporting	Well- established <i>in vitro</i> or <i>ex vivo</i> functional studies or knock-in mouse studies show no damaging effect on protein function. Poly A: One or two independent <i>ex vivo</i> studies all showing NO significant decrease in RNA or protein level. Cell lines from unrelated individuals who harbor the variant. <ul style="list-style-type: none"> <i>Ex vivo</i> studies using patient derived samples need to be interpreted with the understanding that unidentified variants may be present.
Key: ^a https://clinicalgenome.org/working-groups/sequence-variant-interpretation/ ; ^b Kelly et al (2018). For criteria included in the original ACMG/AMP framework but not utilized in these specified criteria the row is shown in gray.	