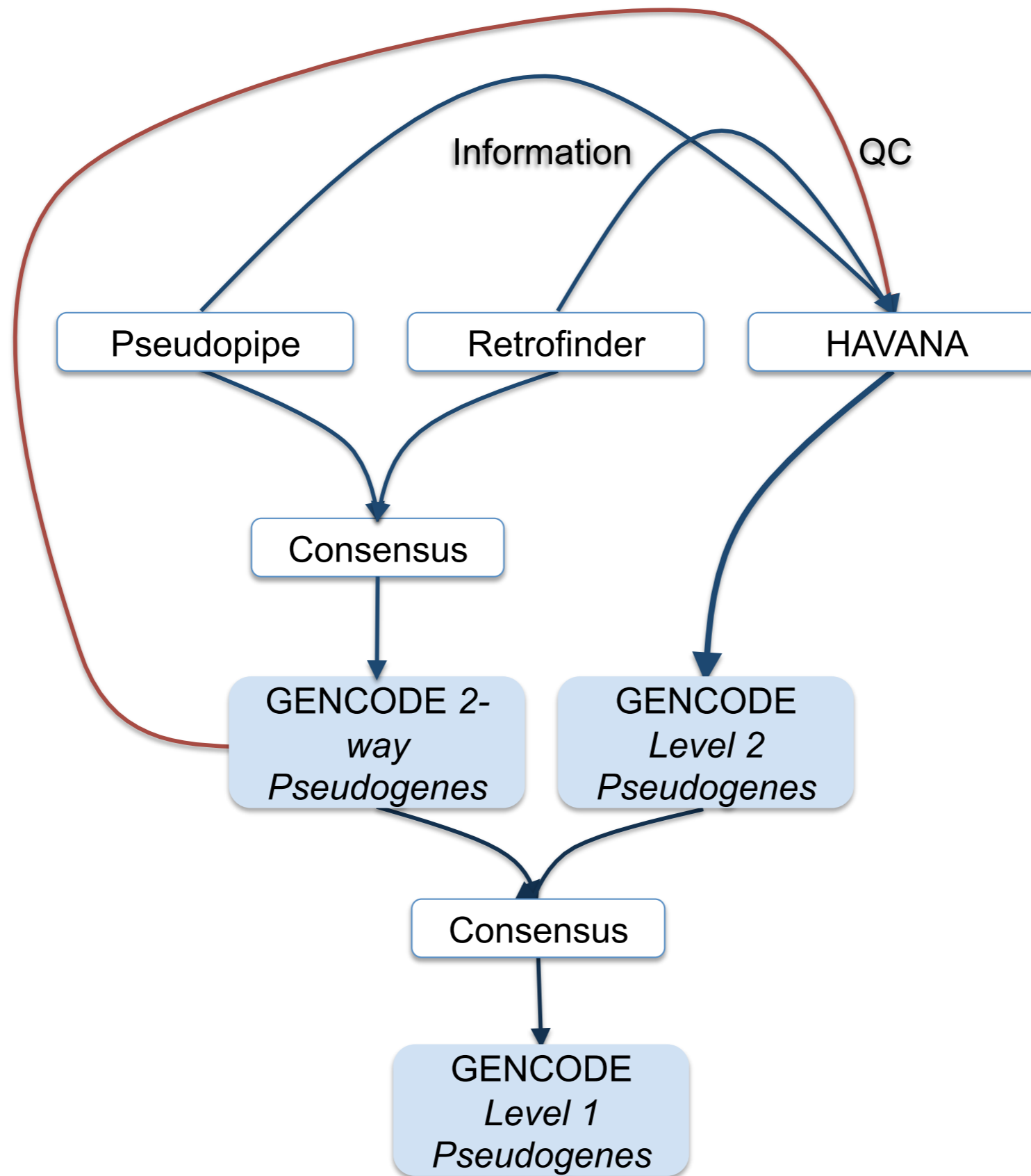# Annotation Pipeline



Figure 1. Pipeline for pseudogene annotation
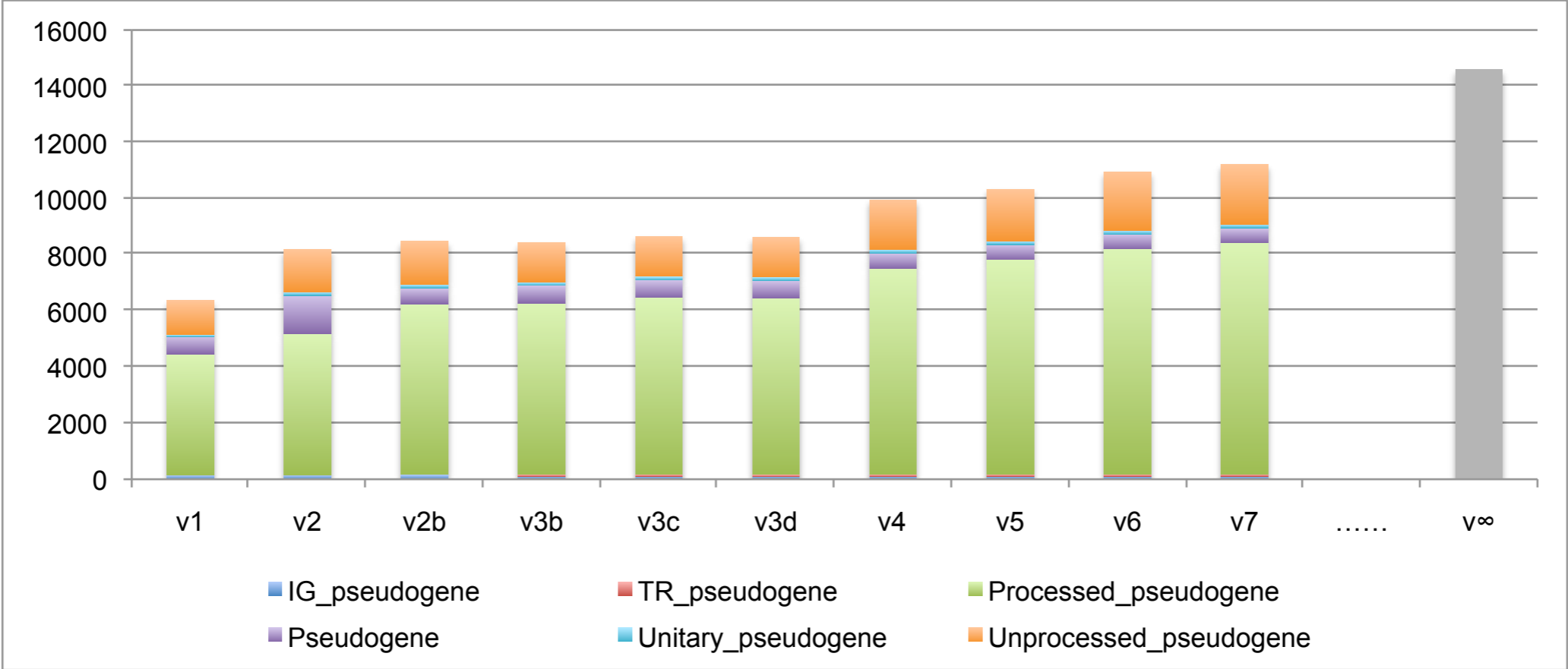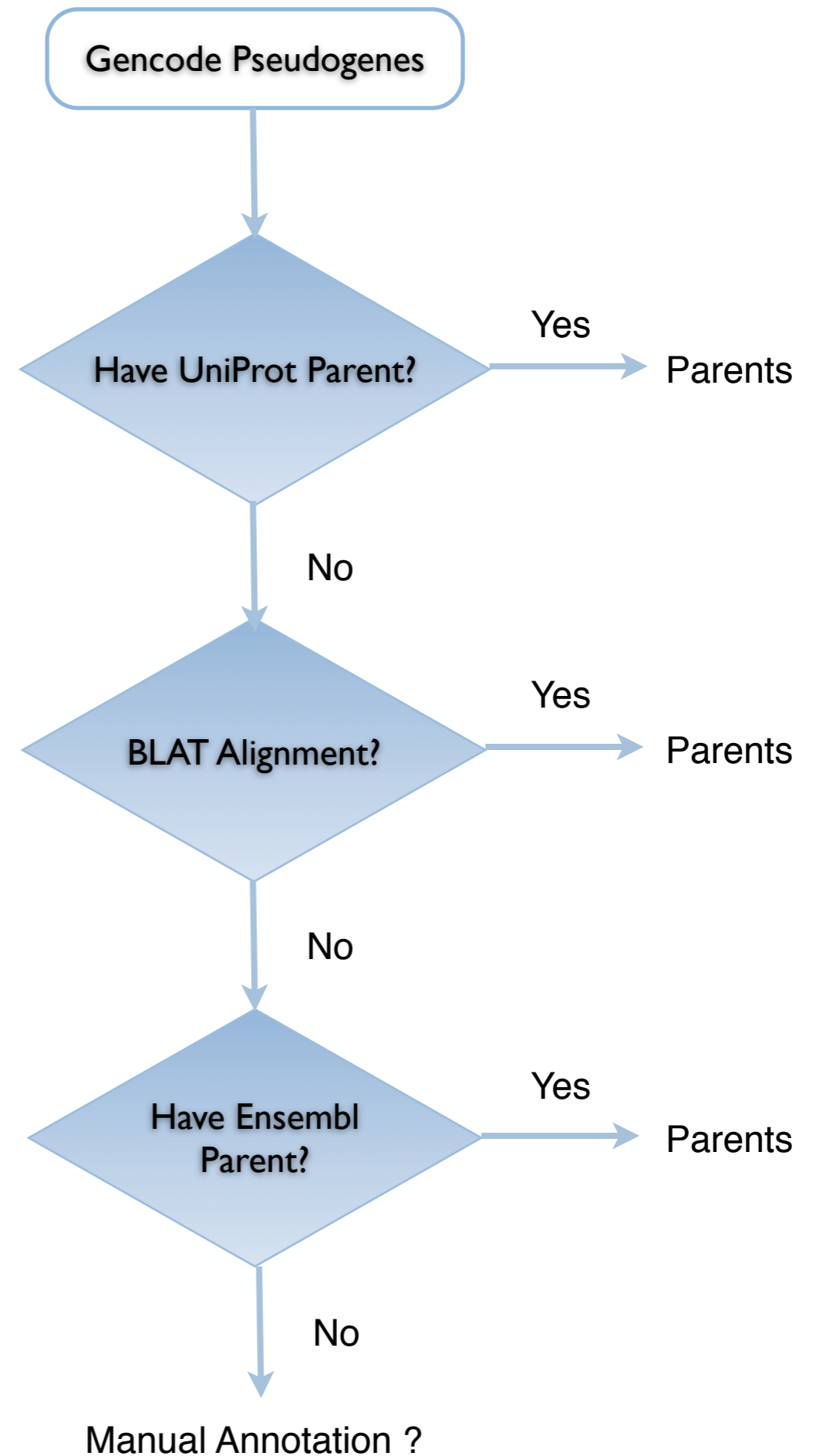
I

# GENCODE Pseudogene Annotation Growth



Figure 2. Growth of pseudogene annotation. Number of pseudogene in the whole genome was extrapolated.

# GENCODE Pseudogene Parents

Figure 3. Identification of GENCODE pseudogene parents.

- UniProt proteins used by HAVANA annotation;

  Parent gene symbols were parsed and mapped to Ensembl 62

- Alignment of pseudogenes to genome

  BLAT pseudogene exons against genome, and identify one-to-one matching regions which overlap coding sequences

- Ensemble peptides used by PseudoPipe

  Find the PseudoPipe output consensus to GENCODE pseudogene, and get PseudoPipe parent information.

Gencode Pseudogenes

Have UniProt Parent? — Yes → Parents

No

BLAT Alignment? — Yes → Parents

No

Have Ensembl Parent? — Yes → Parents
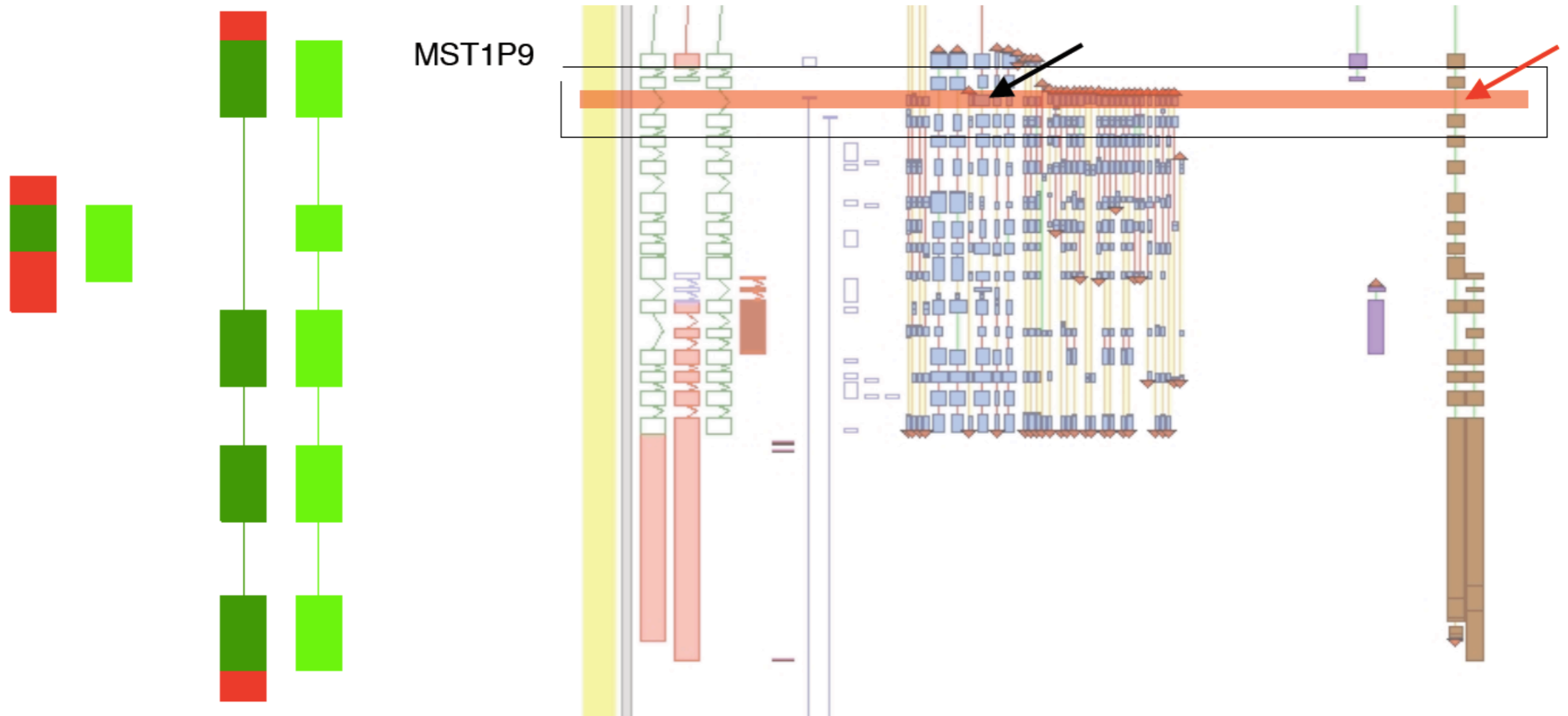
No

Manual Annotation ?

# Resurrected Unprocessed Pseudogene



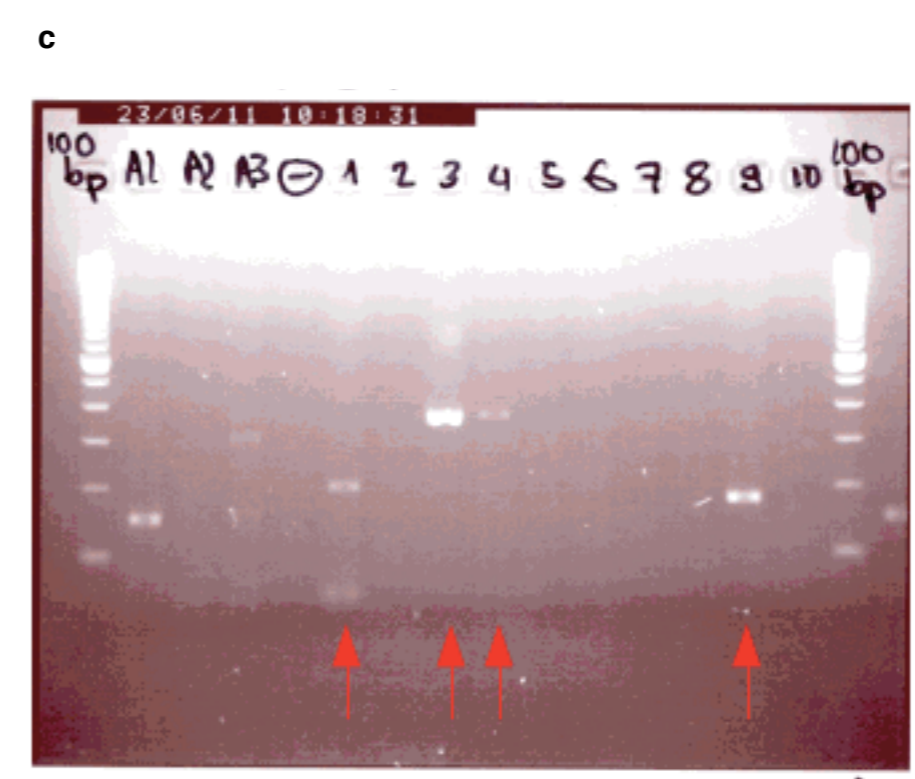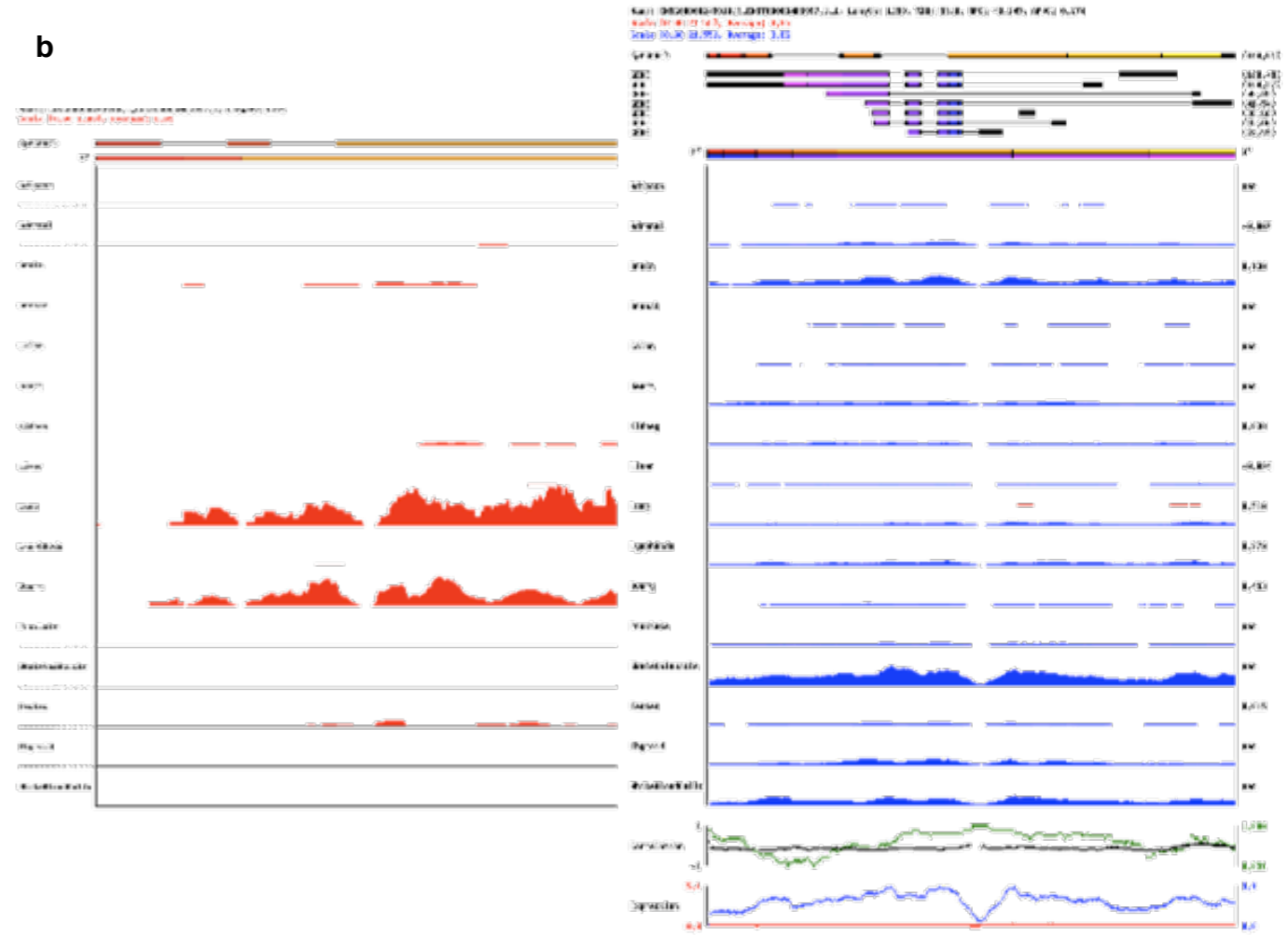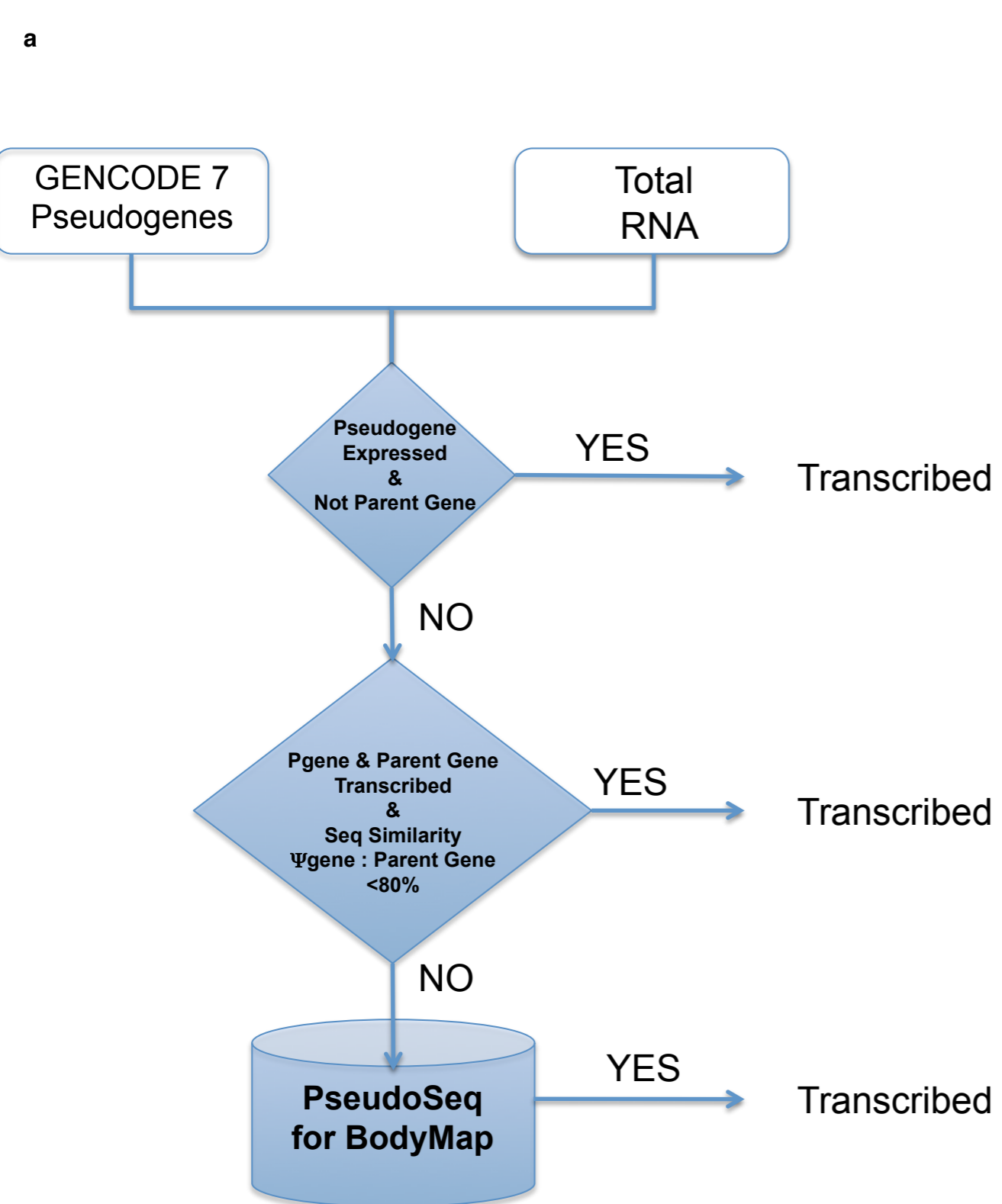Figure 4. Coding locus resurrected from pseudogene

Figure 5. Transcription of pseudogene. a. Flowchart to pick transcribed pseudogene from total RNA RNA-seq data and BodyMao RNAseq data. b. User Interface of PseudoSeq for transcription from a pseudogene and its parent side by side. c. RT-PCR validation of transcribed pseudogenes.

5

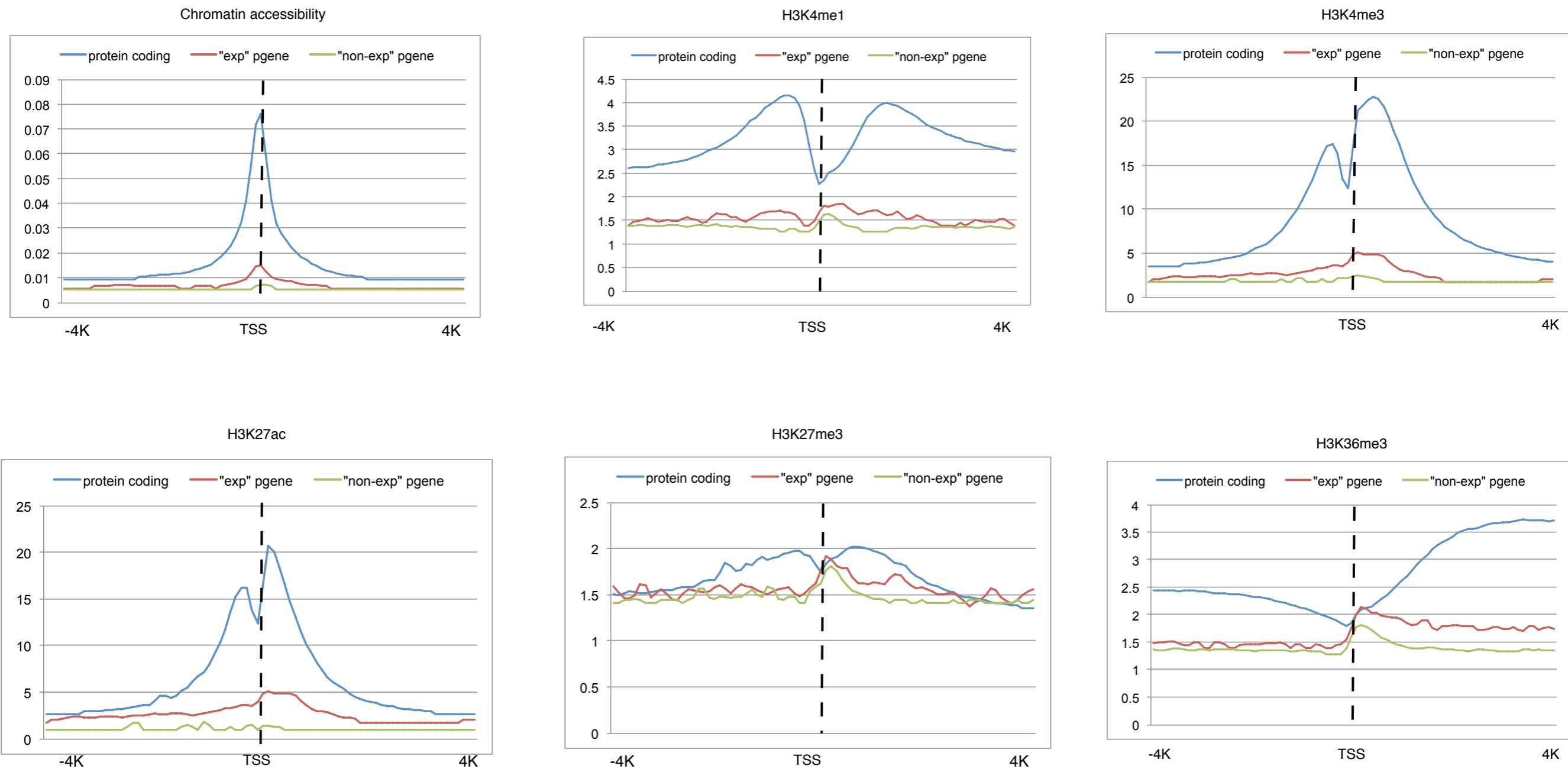# Average chromatin signatures



Figure 6. DNaseI hypersensitivity and histone marks of pseudogenes. Data from Chip-Seq results of GM12878.
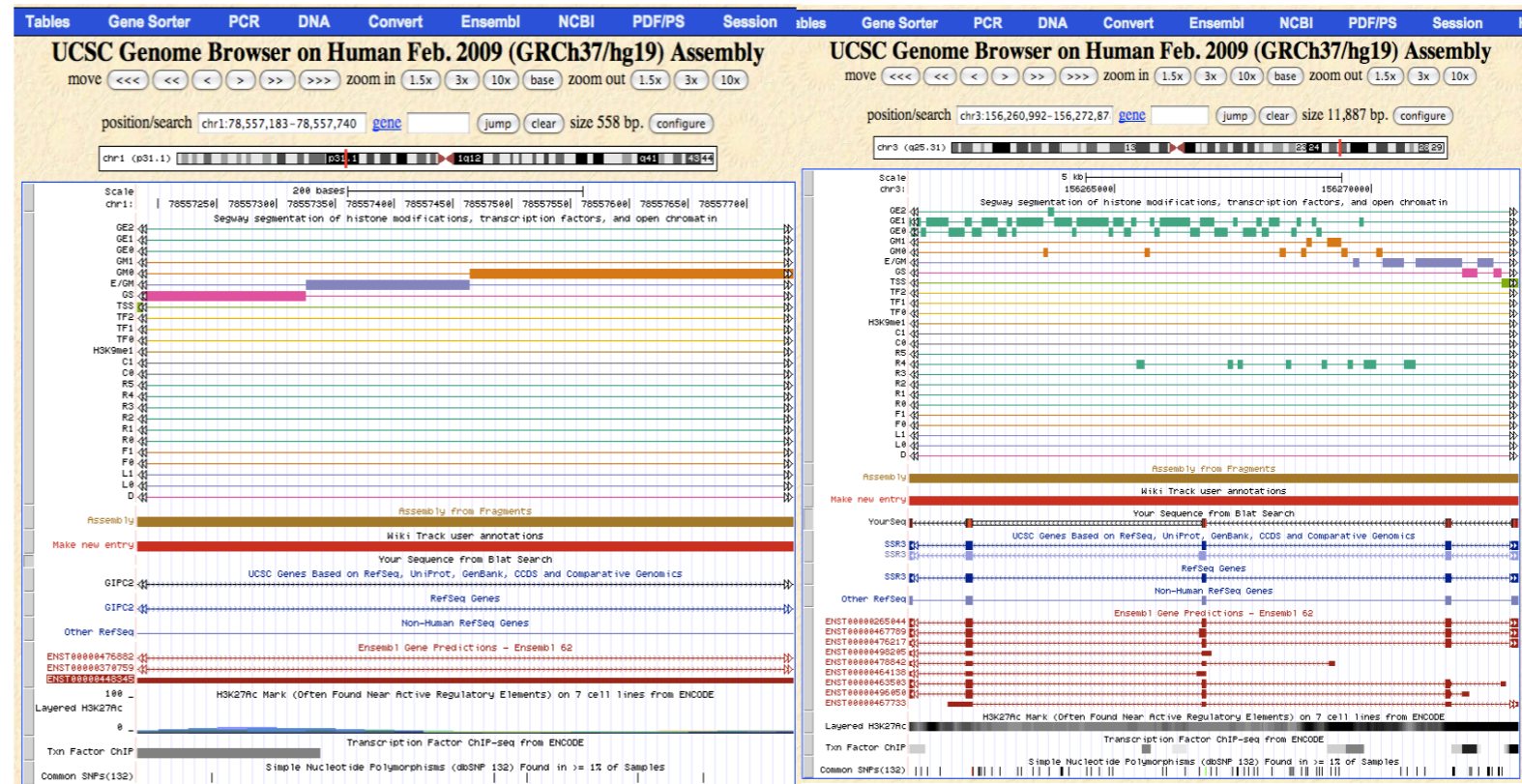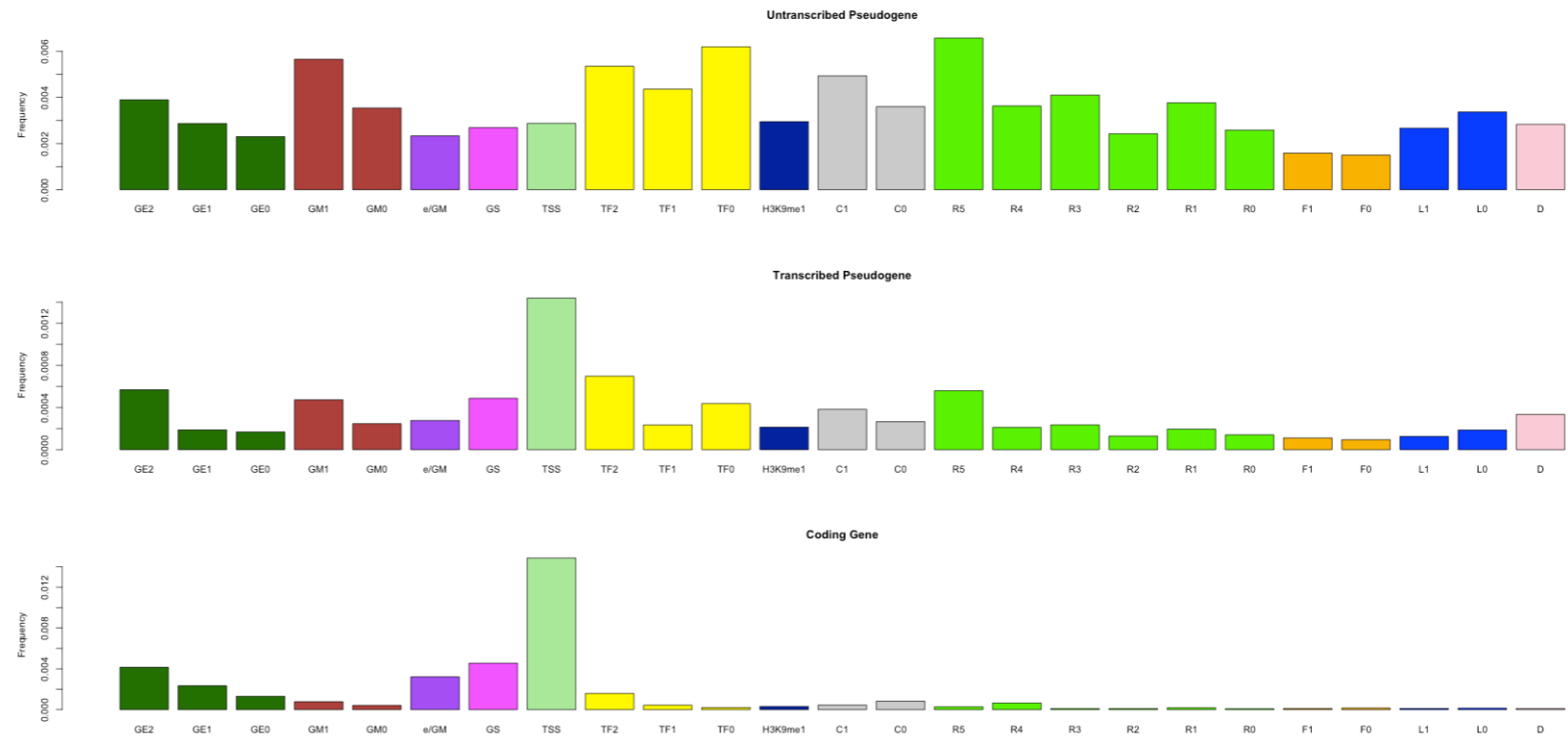
# Segmentation



Figure 7. Segway segmentation of pseudogene.
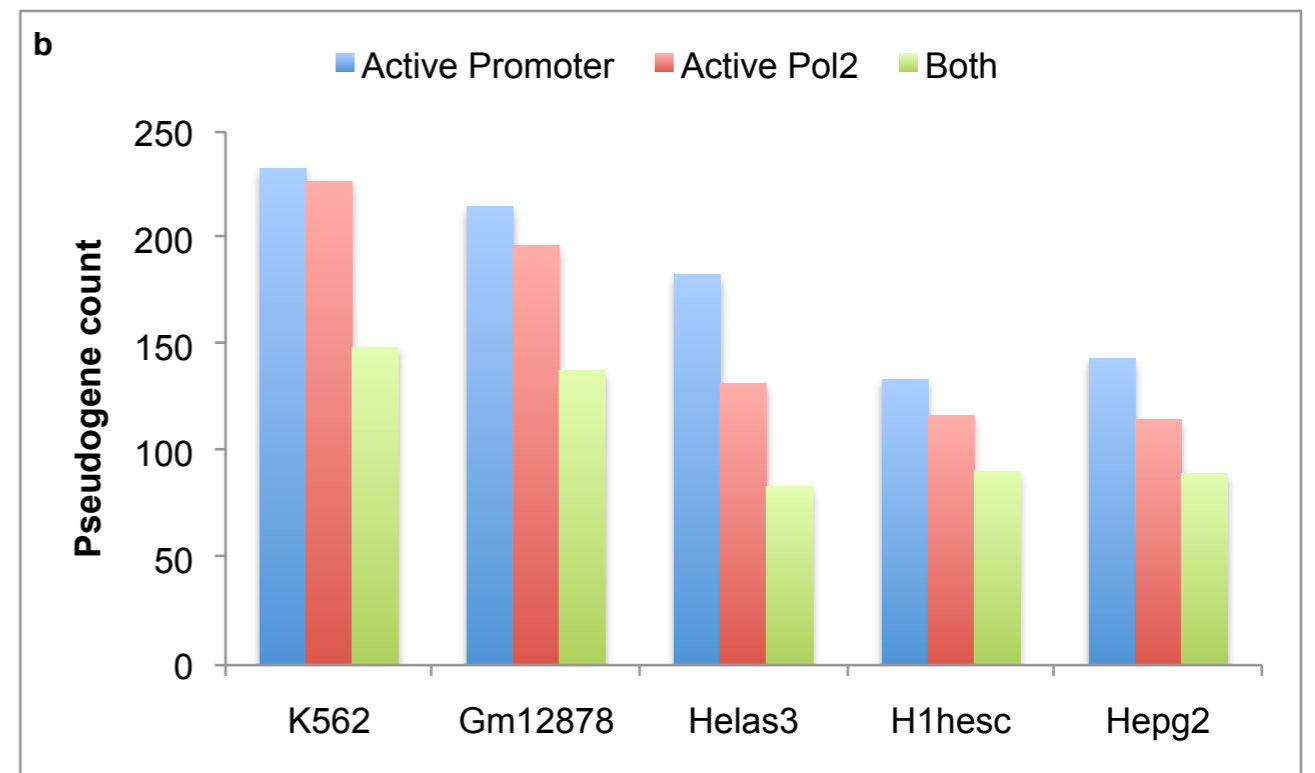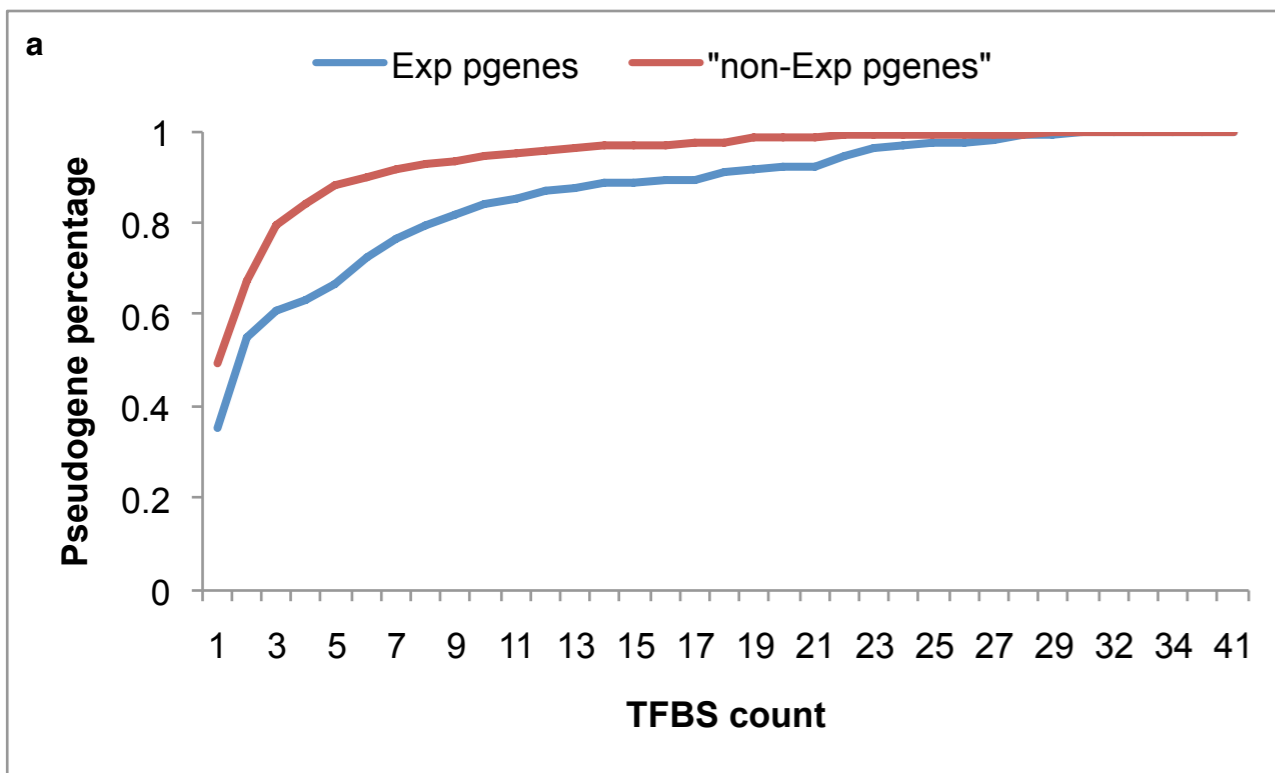
# TFBS of Pseudogenes



Figure 8. TFBS in the upstream of pseudogenes. a. Distribution of pseudogenes with different number of TFBS in their upstream sequences. Data is for K562. b. Number of pseudogenes with active promoter, active pol2 binding sites or both in different cell lines.
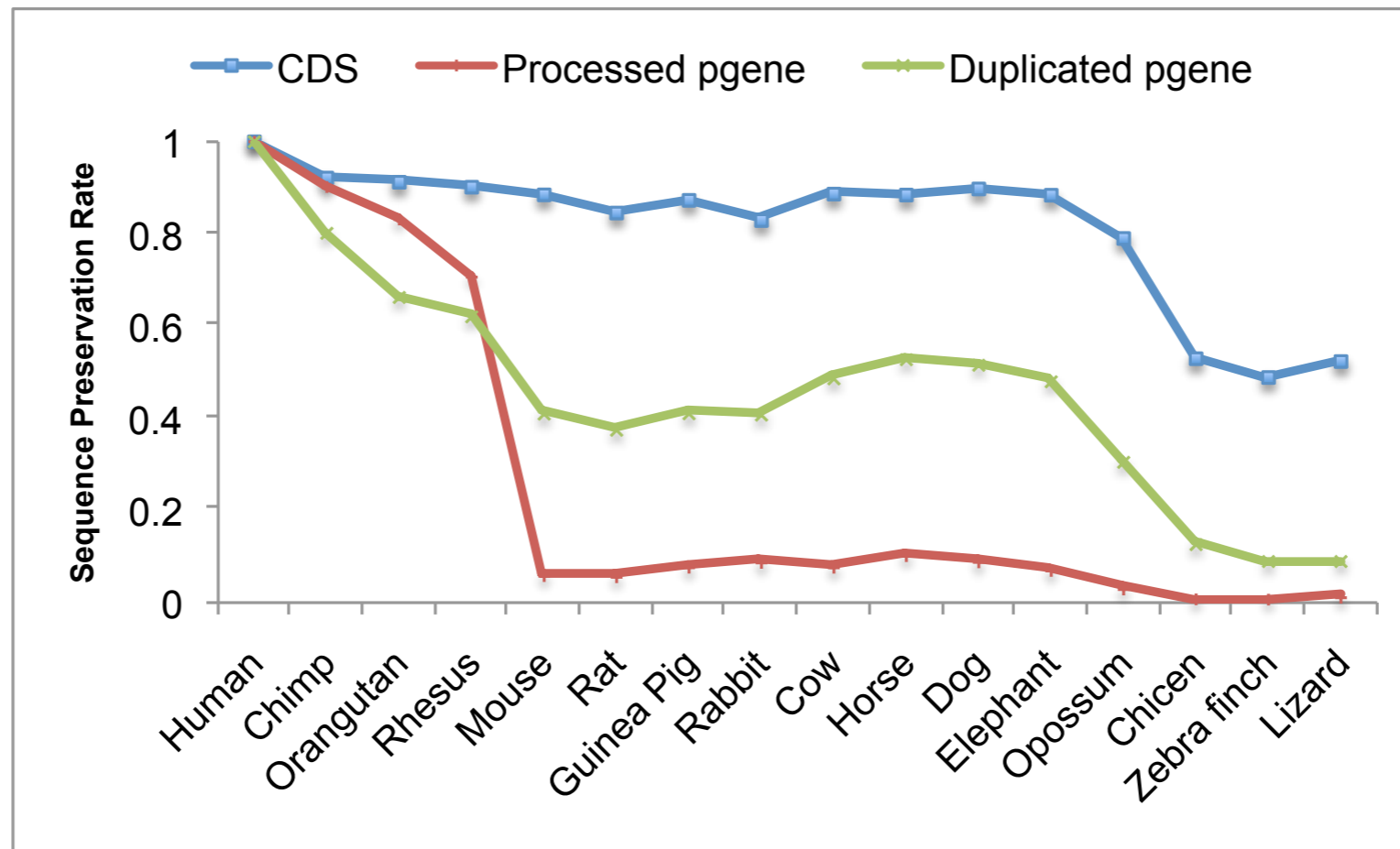
# Sequence Preservation of GENCODE Pseudogenes



Figure 9. Preservation of human coding sequences, processed pseudogenes and duplicated pseudogenes. Orthologous to human genomic regions from different species were studied. The sequence preservation rate was calculated as percentage of sequences aligned to human sequence from each species. The calculation was based on MultiZ sequence alignment.
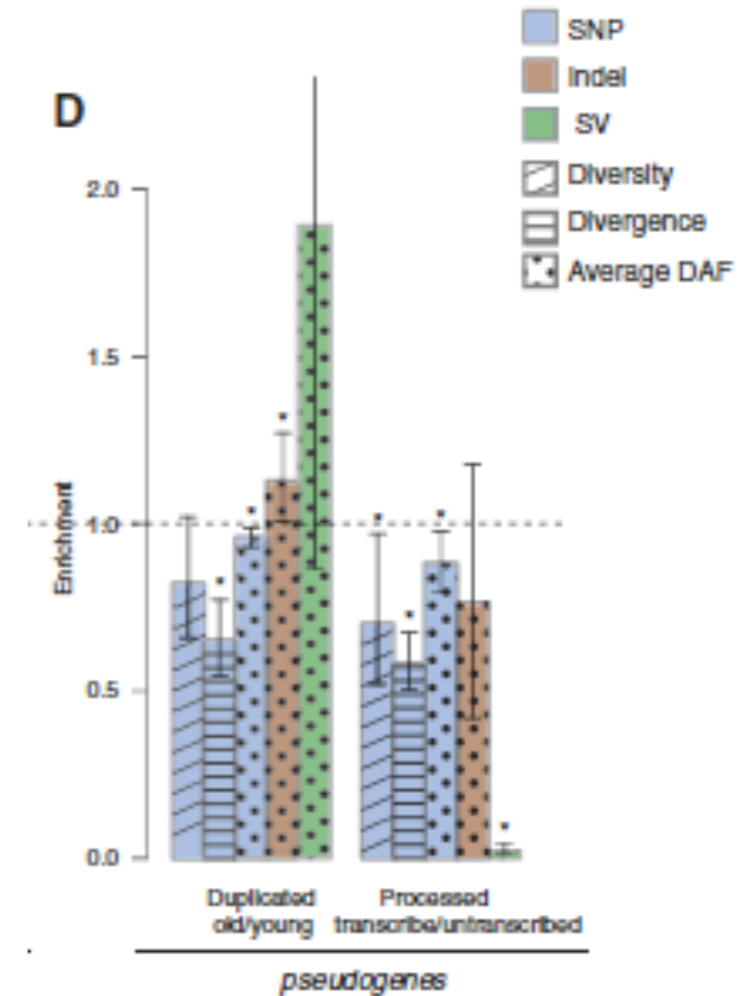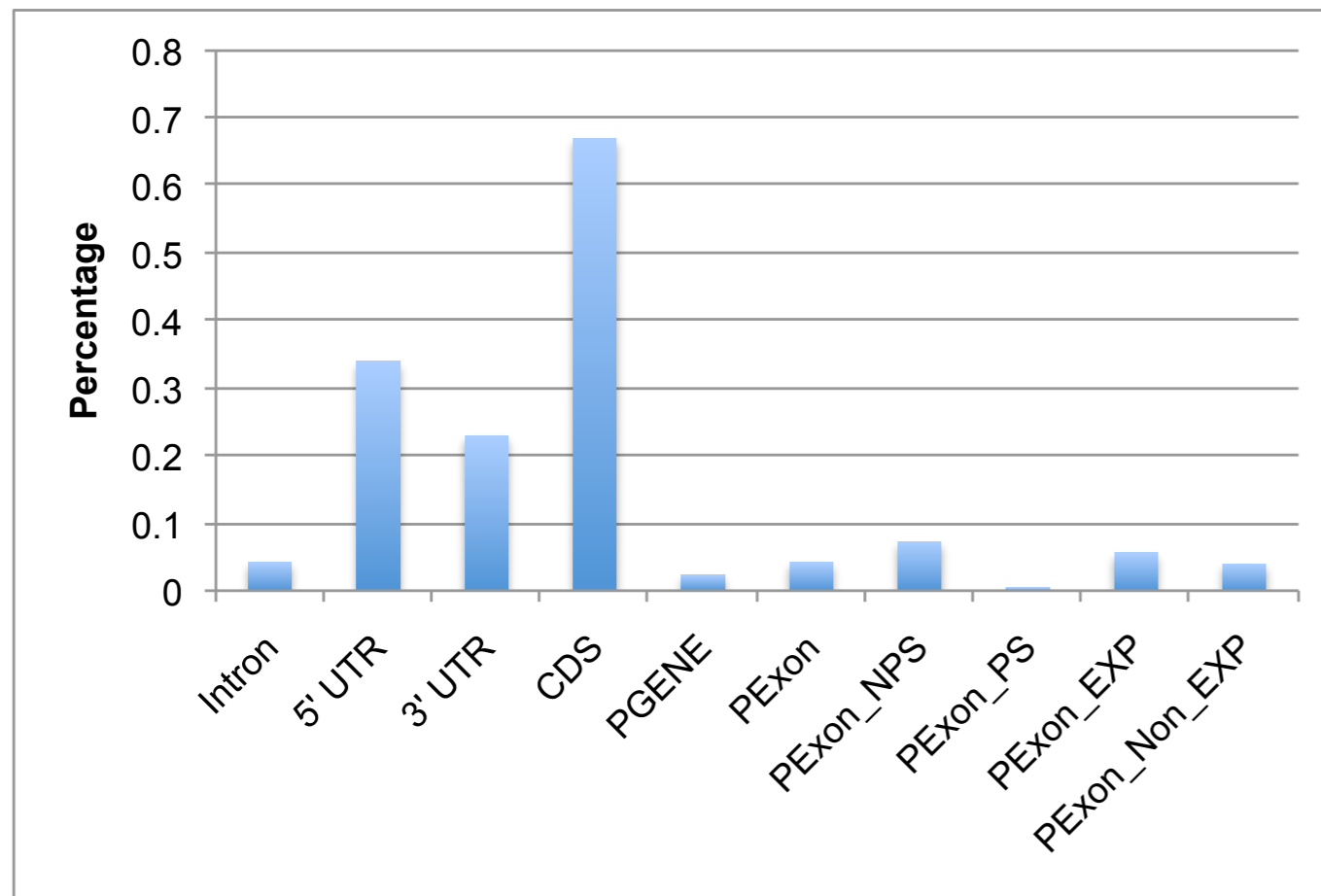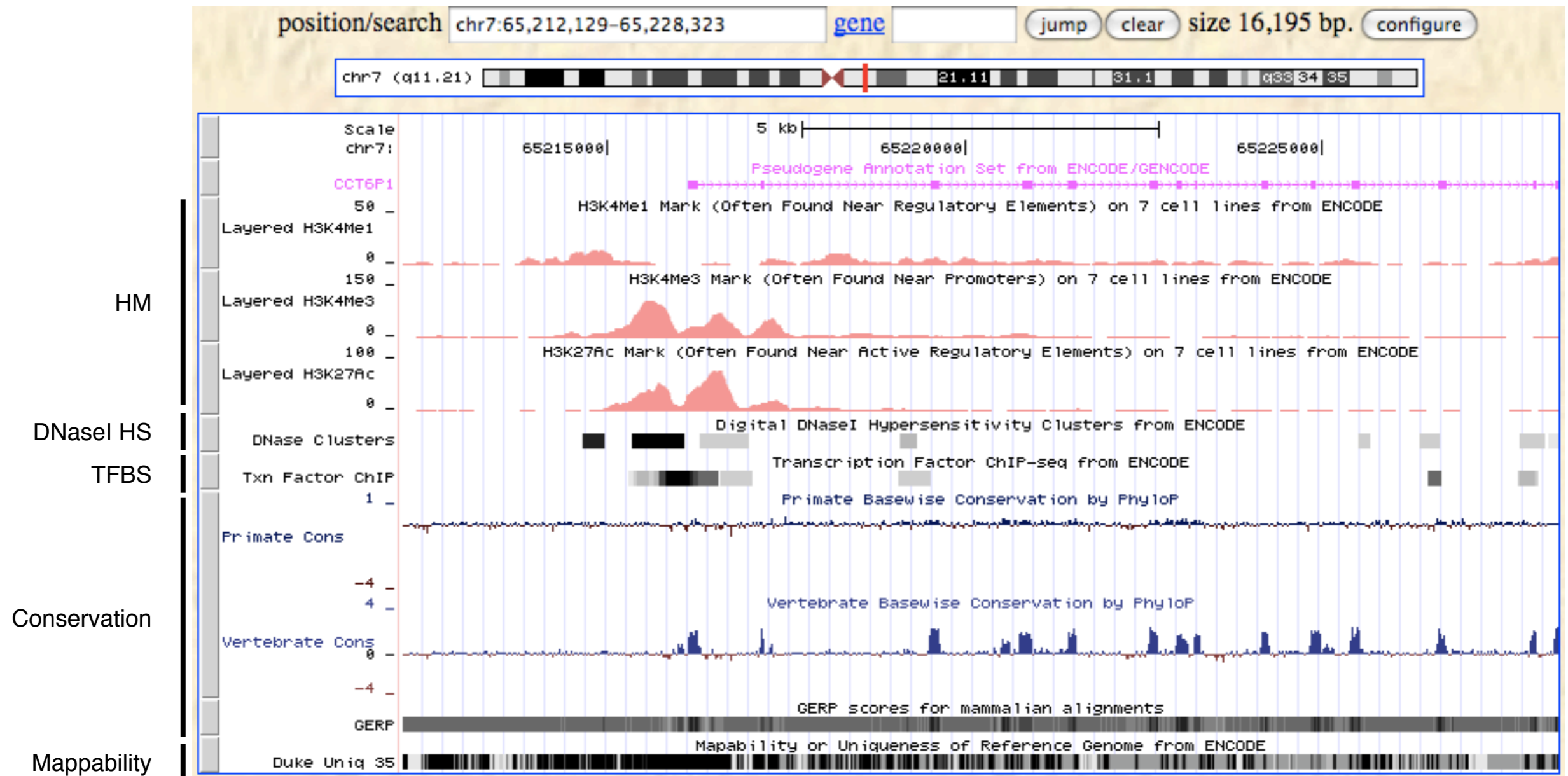
# Pseudogene Conservation



Figure 10. Overlap between GERP constraint elements and human genomic regions. Pseudgenes were divided to different classes, as nonprocessed pseudogene, processed pseudogene, transcribed pseudogenes and non-transcribed pseudogenes.

# Case 1: transcribed pseudogene with active chromatin states

# Case 2: transcribed pseudogene with inactive chromatin states