

Final Consensus Gene Set

AugustusCGP-AugustusTMR-transMap consensus

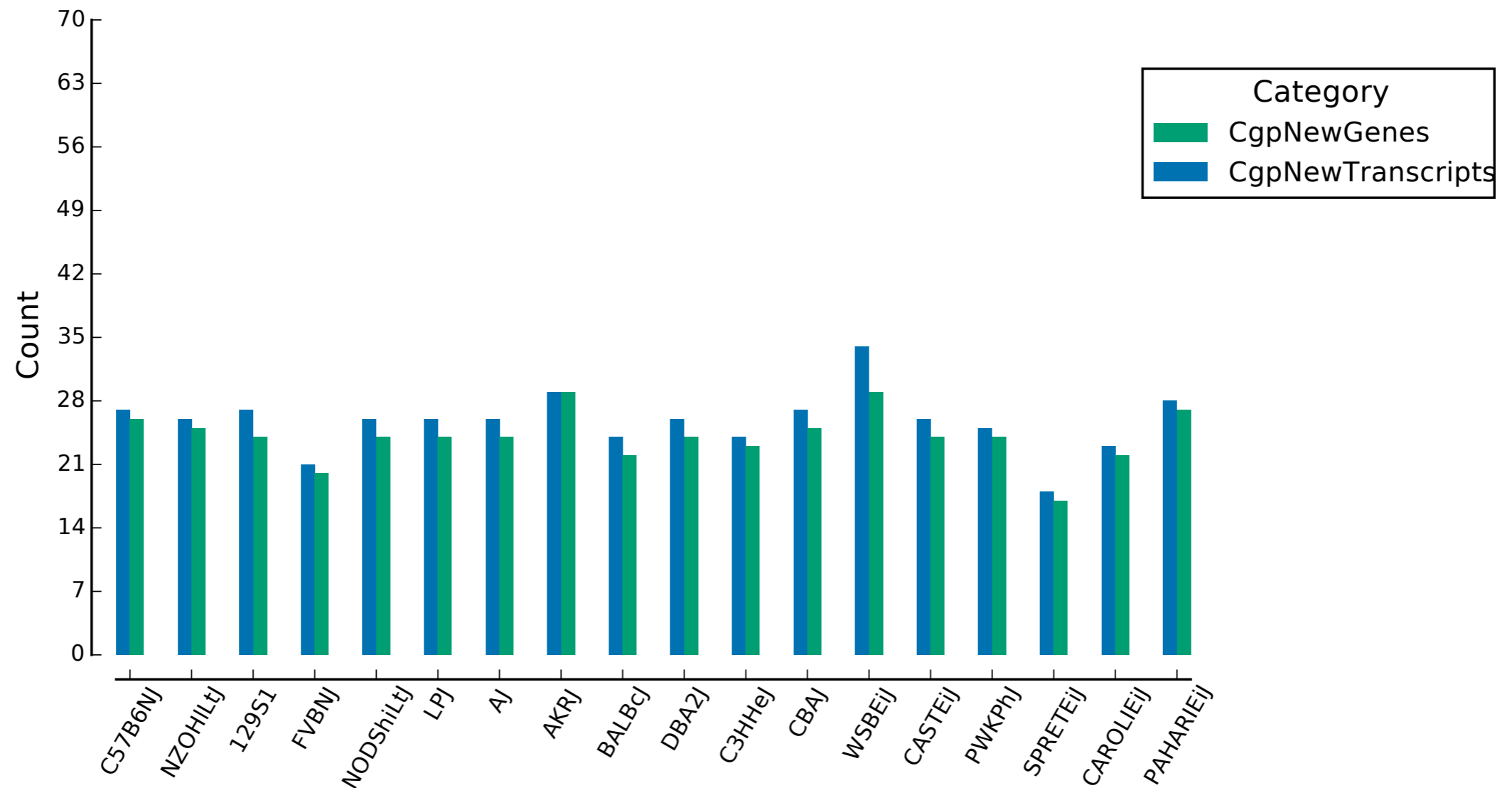
- Takes as input both the AugustusCGP transcripts and the AugustusTMR/transMap consensus
- Each AugustusCGP transcript is mapped to one or more GENCODE gene IDs by interval intersections
- Each AugustusCGP transcript is BLAT'd against every GENCODE transcript that are associated with each GENCODE gene that it overlaps

AugustusCGP-AugustusTMR-transMap consensus (2)

- For each AugustusCGP transcript:
 - If this transcript does not overlap a gene in the TMR consensus, include it as a new transcript.
 - If this transcript has a higher % identity with the same or higher % coverage to a reference transcript, replace that transcript in the TMR consensus.
 - UTR information will be lost.
 - If this transcript introduces splice junctions that are supported by RNAseq that are not present in the TMR consensus, include it as a new isoform.
 - If this transcript maps to more than one gene, make sure that the intergenic splice junctions are supported by RNAseq.

New genes

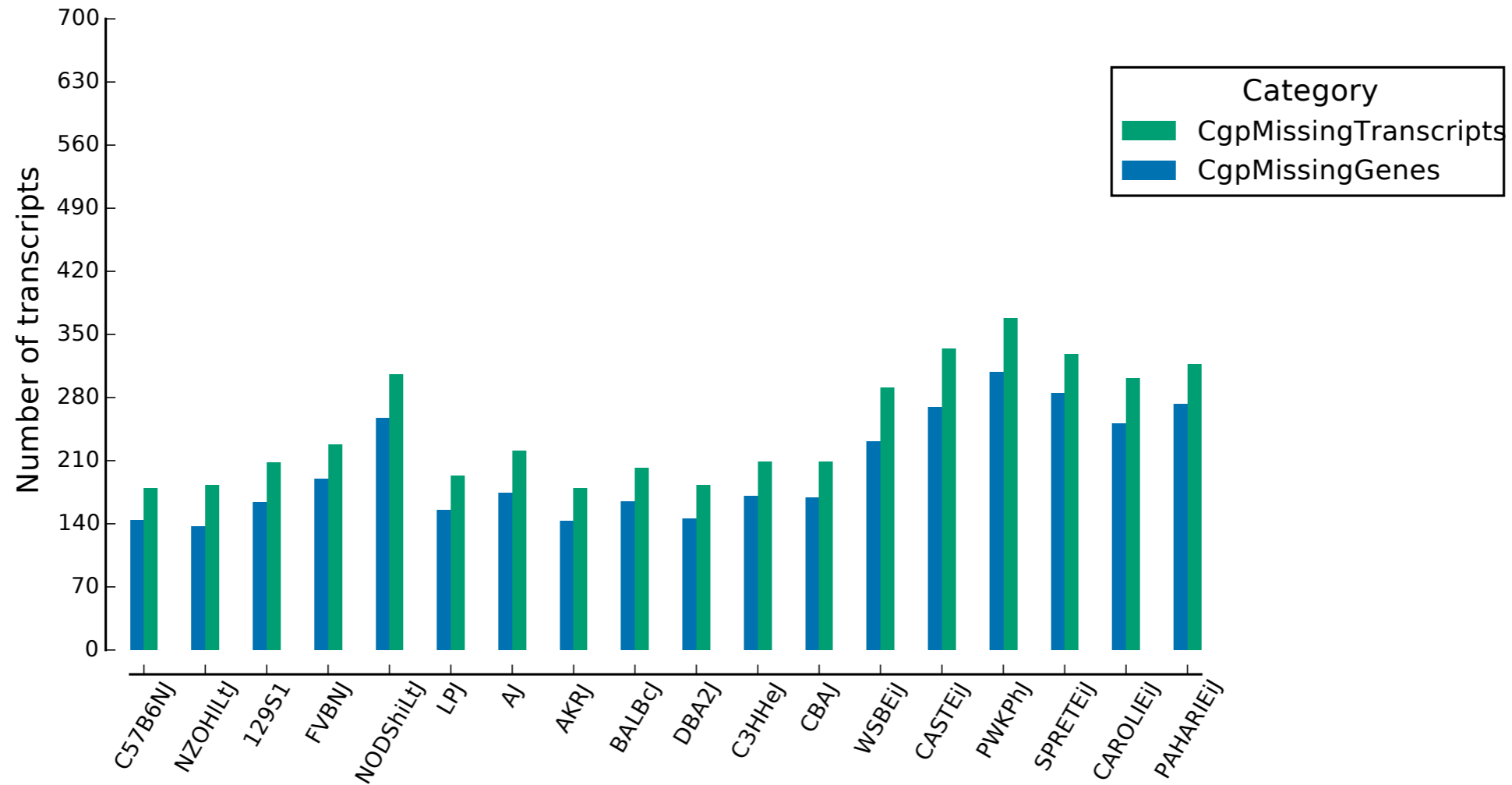
Breakdown of the number of new genes/transcripts introduced by Comparative Augustus to the consensus gene set derived from the annotation set GencodeCompVM7



New genes do not overlap any TMR consensus genes

Missing genes rescued

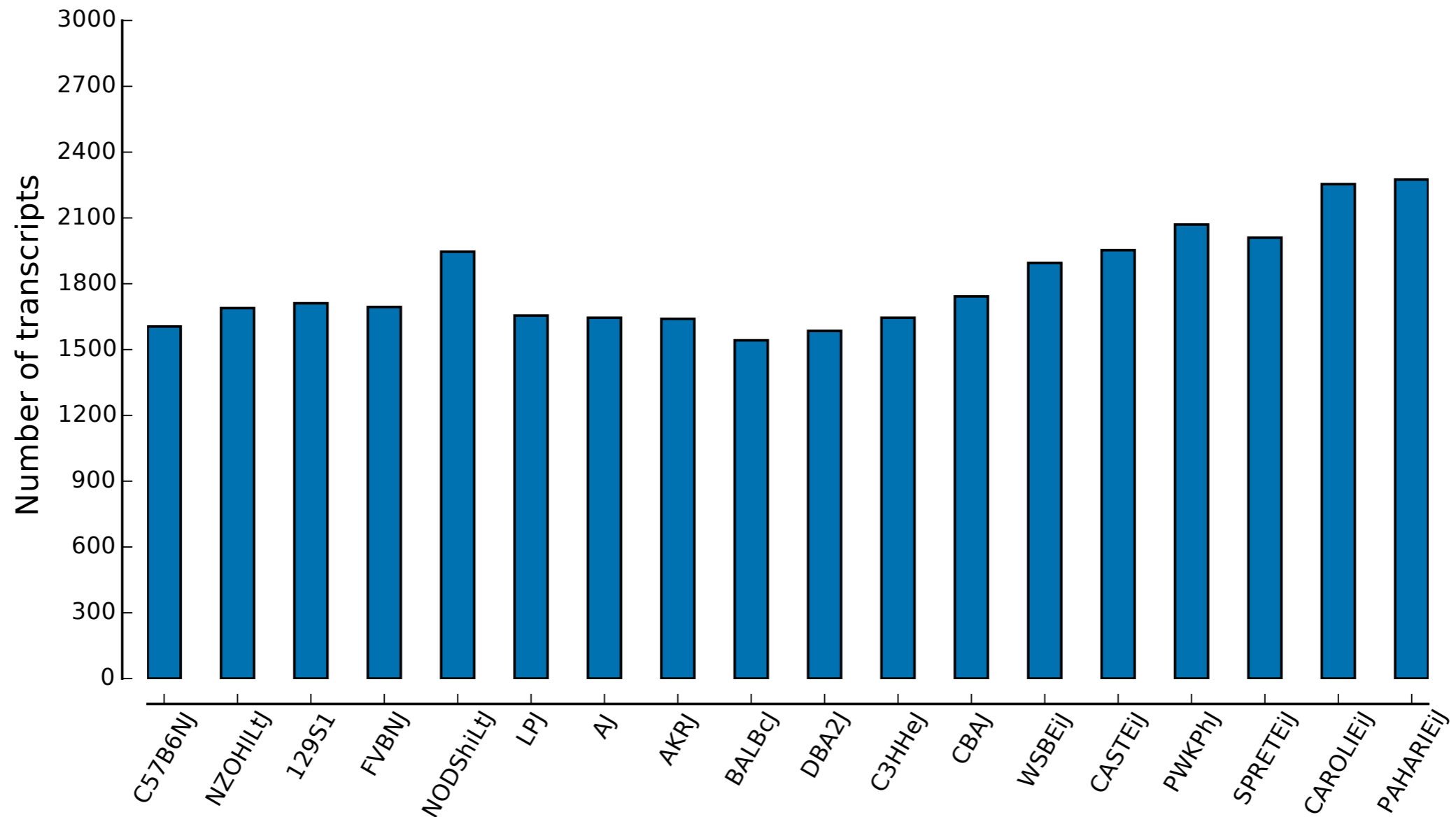
Breakdown of the number of missing genes/transcripts rescued by Comparative Augustus to the consensus gene set derived from the annotation set GencodeCompVM7



Missing genes are genes which got filtered out of TMR consensus due to coverage/identity

New Isoforms

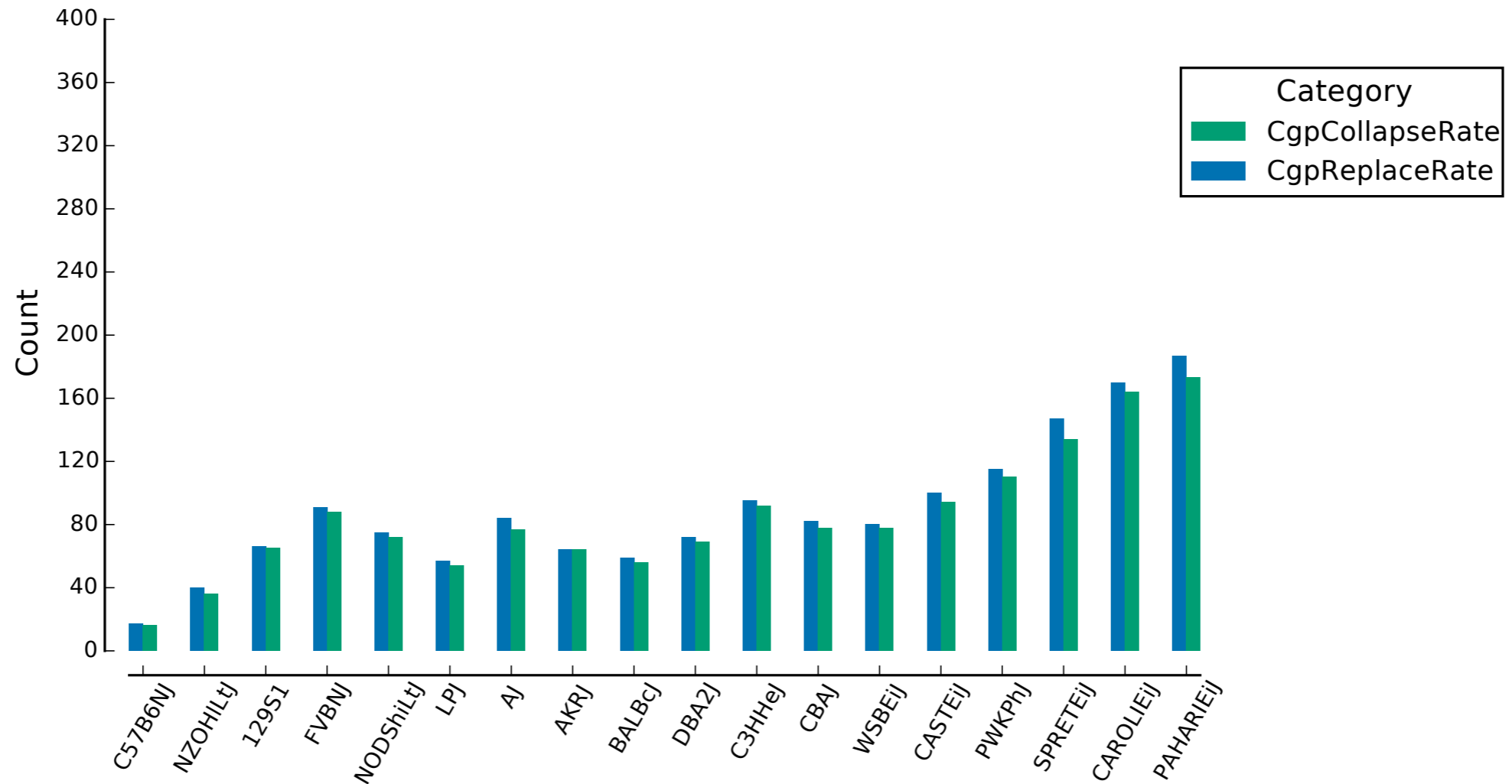
Breakdown of the number of new isoforms added by Comparative Augustus to the consensus gene set derived from the annotation set GencodeCompVM7



New isoforms are CGP transcripts who have splice junctions supported by RNAseq not present in the TMR consensus

Replace rate

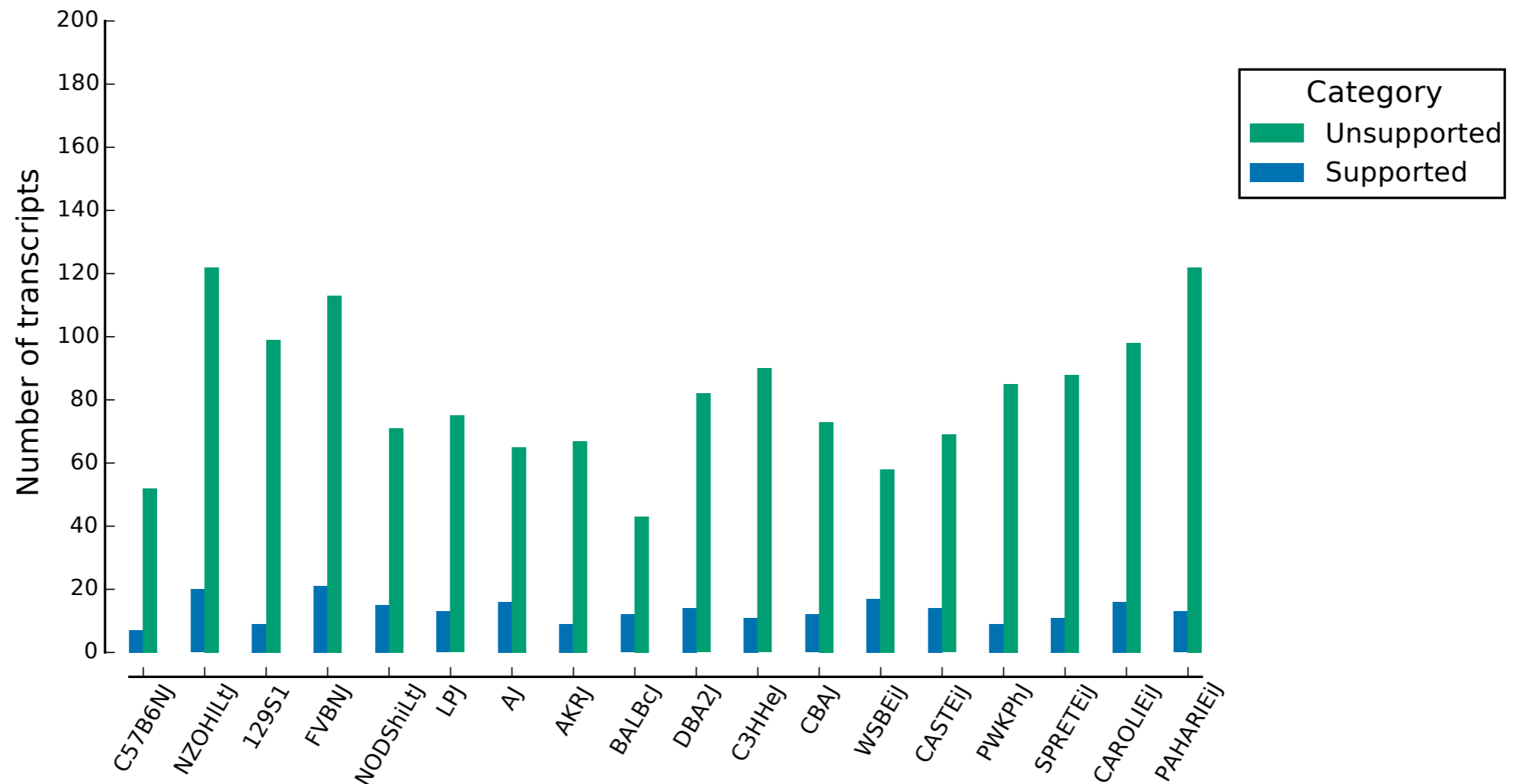
breakdown of the number of transMap/augustusTMR consensus transcripts replaced by augustusCGP from the consensus gene set derived from the annotation set GencodeCompVM7



How many transcripts more closely matched the reference?
After replacement, how many now-identical transcripts will be collapsed?

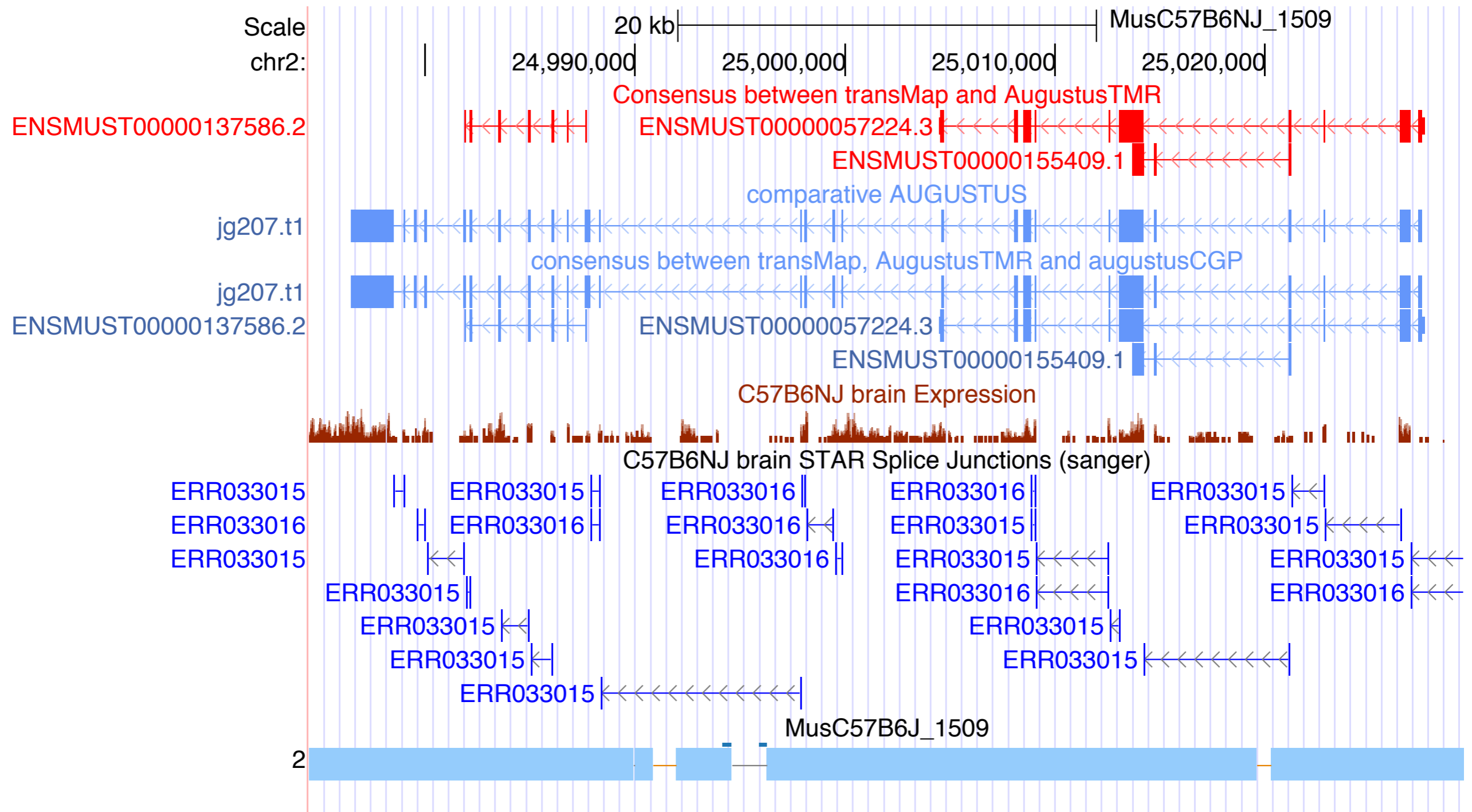
Joined genes

How many CGP consensus transcripts join TMR transcripts in a supported fashion to the consensus gene set derived from the annotation set GencodeCompVM7

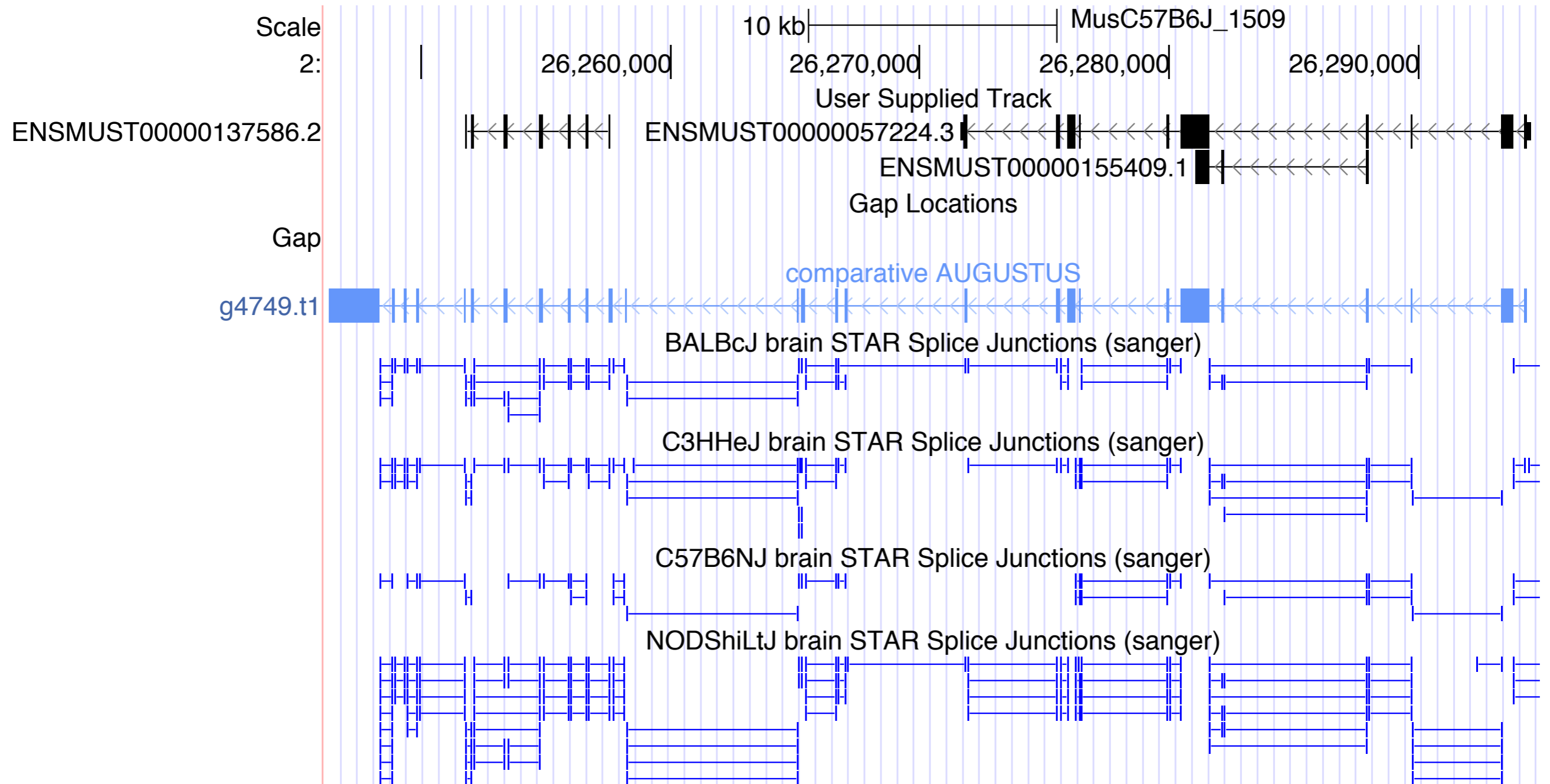


How many transcripts that join together two genes with junctions supported by RNAseq? Only supported joins make it to the consensus set.

Joined gene in reference



Joined gene in reference (2)



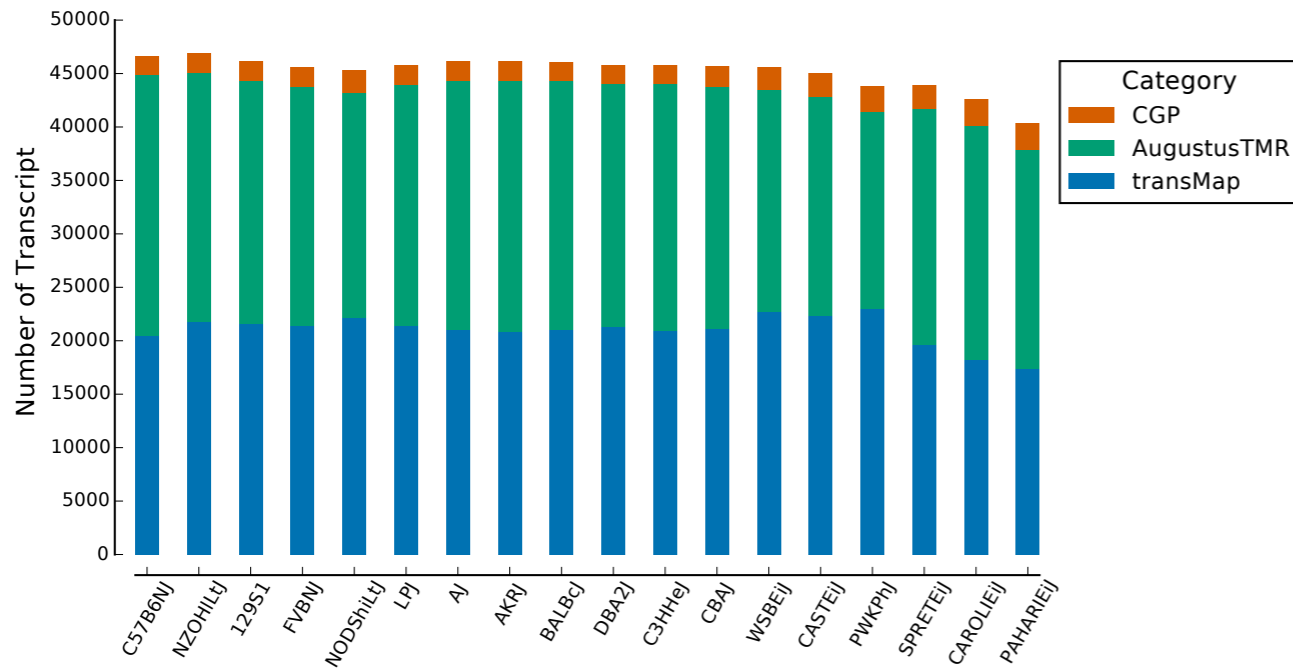
Mixed RNAseq splice junction support across strains

Joined genes problems

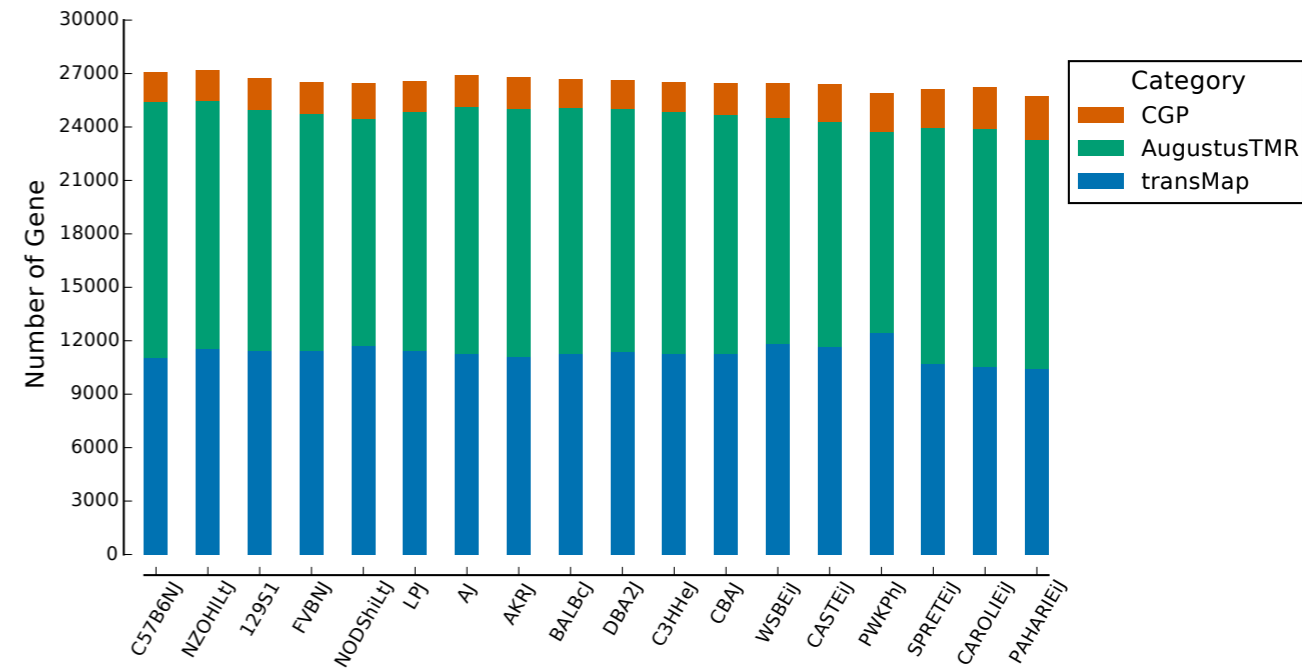
- Most 'unsupported' are in fact artifacts of alignment and are not real candidates for join-genes at all
- Need to better filter intervals before comparing to transMap results

Provenance of final consensus sets

Breakdown of the origins of the final consensus Transcript set to the consensus gene set derived from the annotation set GencodeCompVM7



Breakdown of the origins of the final consensus Gene set to the consensus gene set derived from the annotation set GencodeCompVM7



Consensus Gene Set Conclusions

- CGP introduces thousands of new isoforms, rescues missing genes.
- CGP shows examples of joined genes that can help improve the GENCODE reference
- In my opinion, we are hitting a point of diminishing returns for refinement. Consensus gene set is ready, and on the browser!

