



Global Analysis of Transcription Start Sites and Transcription Units in Bacterial Genomes



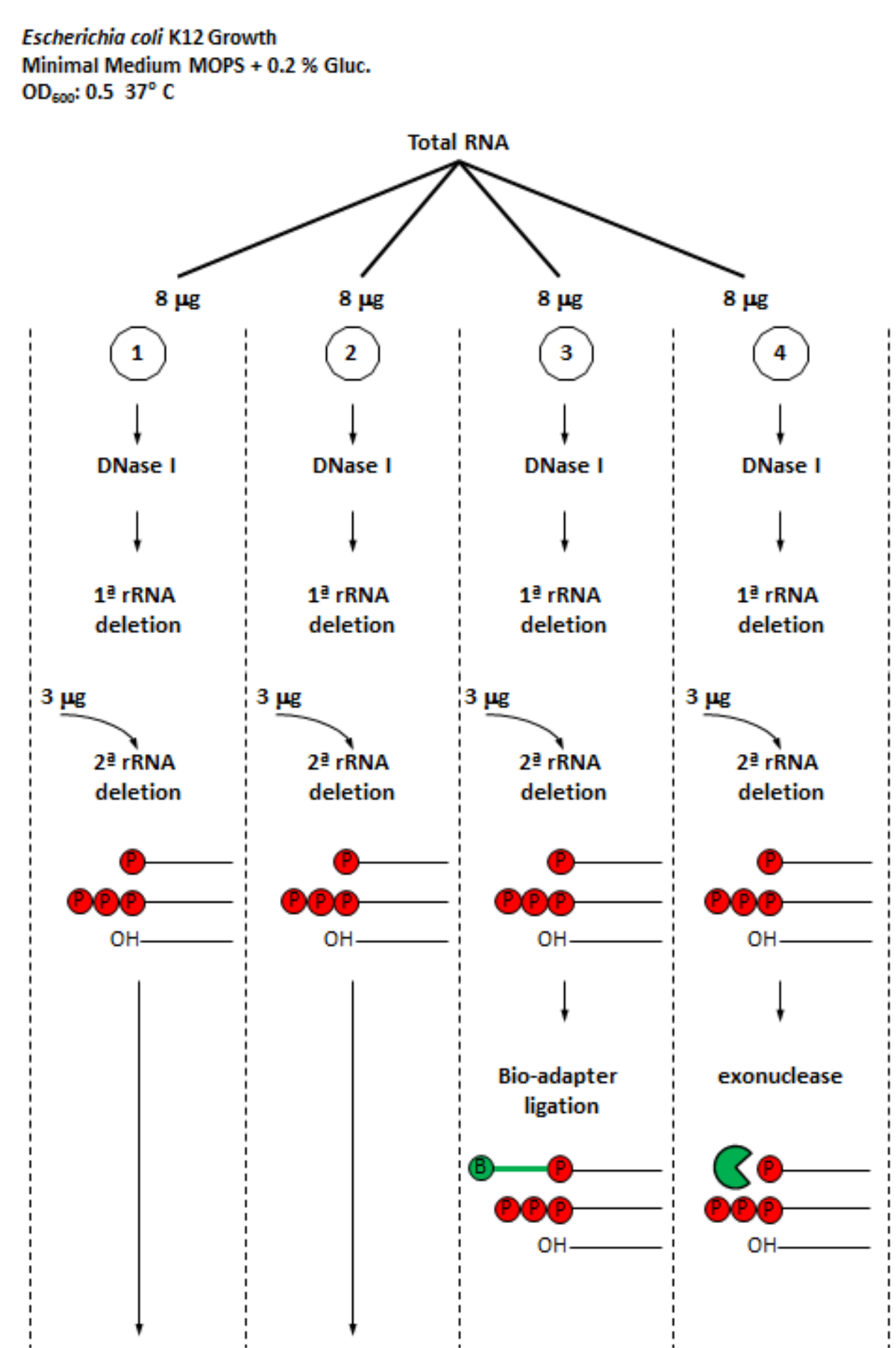
Collado-Torres L.^{1,5}, Reyes-Quiroz A.², Cuellar-Partida G.², Moreno-Mayar V.², Vargas-Chávez C.^{1,5}, Taboada B.³, Vega-Alvarado L.³, Jiménez-Jacinto V.⁴, Mendoza-Vargas A.⁵, Grande R.⁴, Olvera L.⁵, Olvera M.⁵, Juárez K.⁵, Collado-Vides J.⁶, Morett E.⁵

¹Winter Genomics, ²LCG-UNAM, ³CADET-UNAM, ⁴UUSMD-UNAM, ⁵IBT-UNAM, ⁶CCG-UNAM

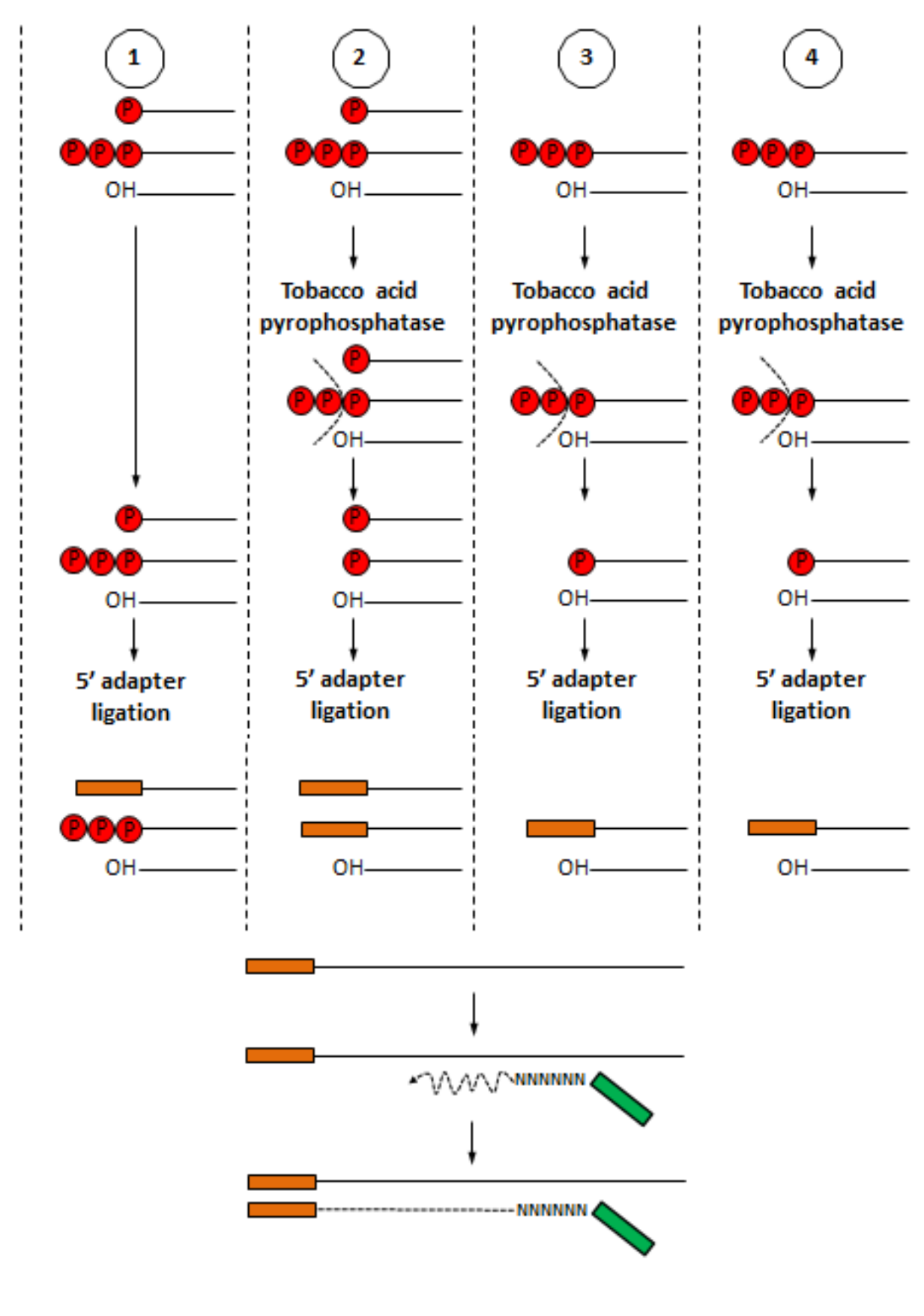
Summary

With high throughput sequencing it is possible to identify at the genomic scale the transcription start sites (TSSs) and transcription units (TUs) in bacterial genomes. Due to the biological and data complexity, these analyses are challenging and require the development of custom algorithms. Critical steps involve 1) maximizing the number of reads that can be used without introducing false alignments, 2) removing biological noise: random transcription and degradation products, 3) identifying the TSSs, 4) visualizing global TSSs patterns, and 5) identifying TUs. Condensing the analyses tools into a Bioconductor package will guarantee the reproducibility of the work. The methods have been developed with data from *Escherichia coli* and *Geobacter sulfurreducens*.

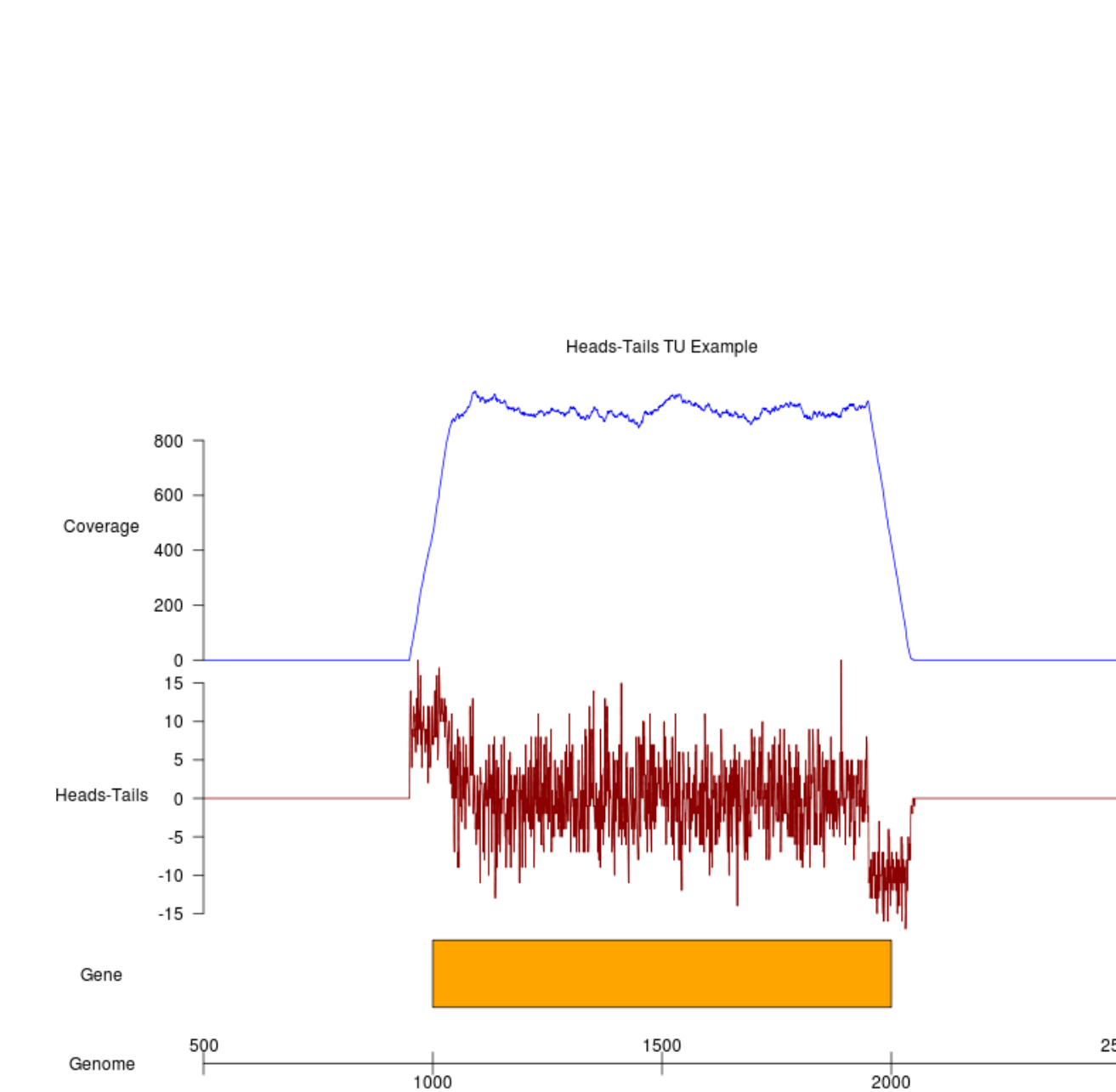
Methods



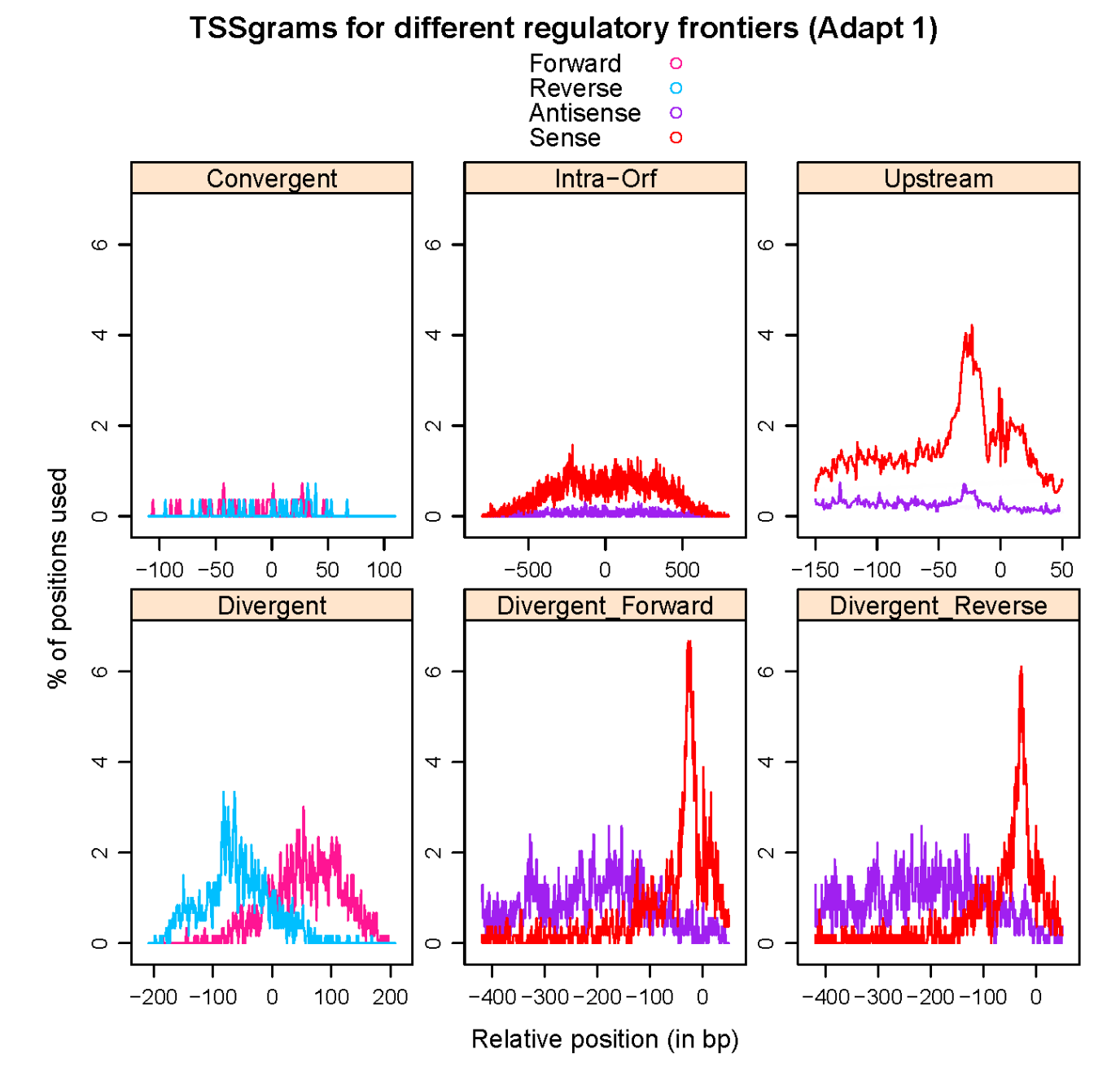
TSSs experimental methods



Transcription Units

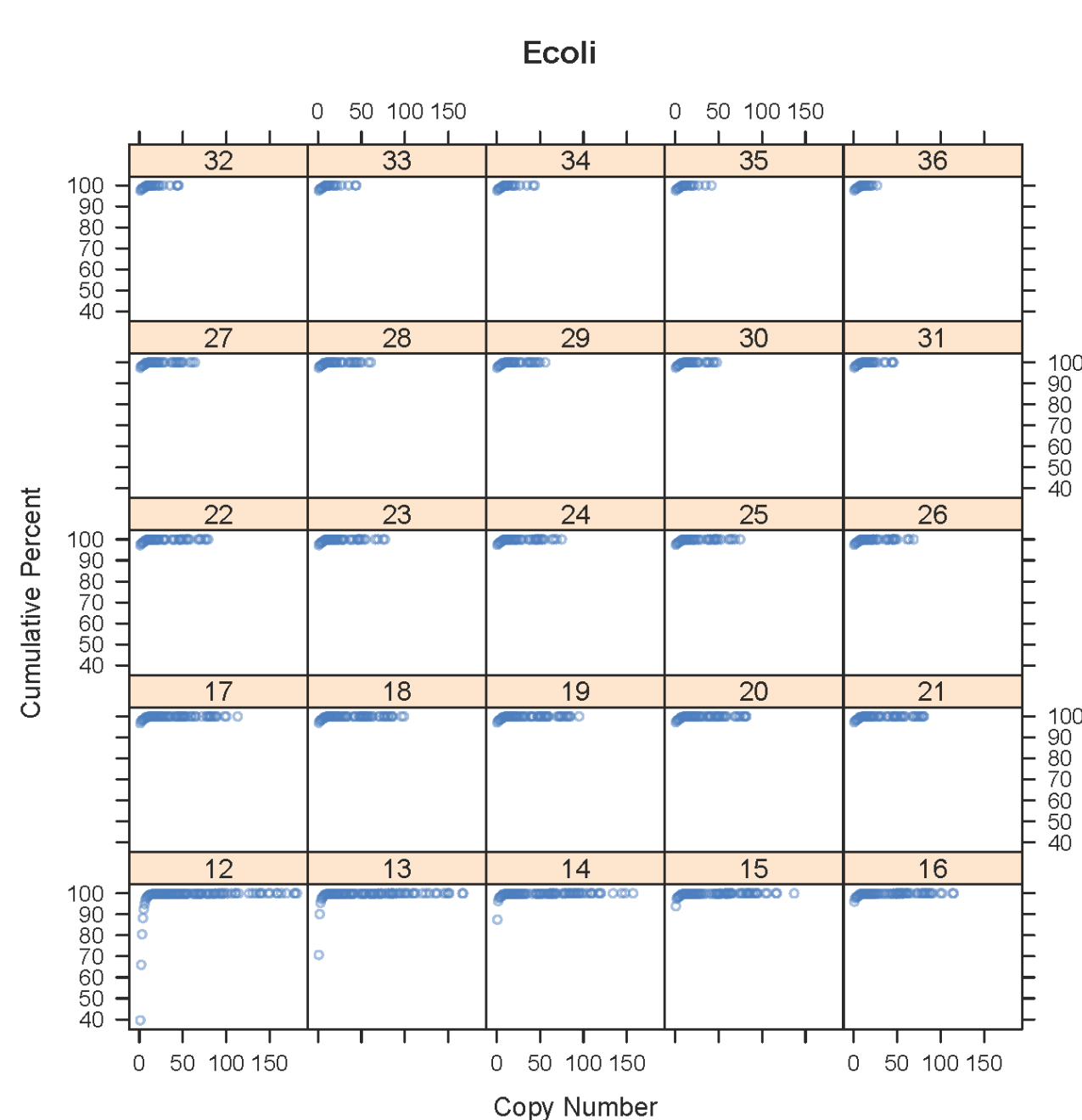


Heads minus tail method

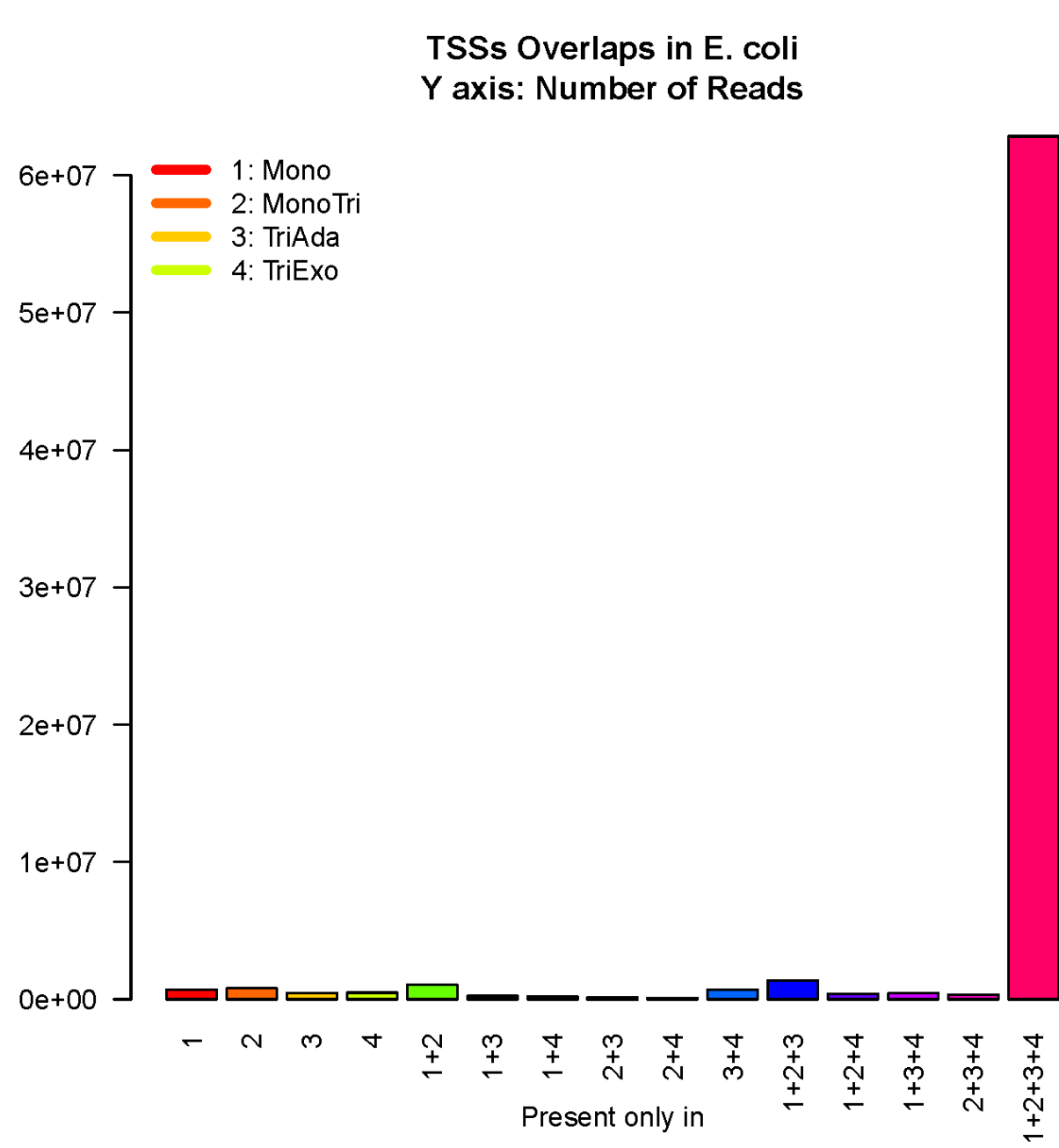


TSSgram sample start

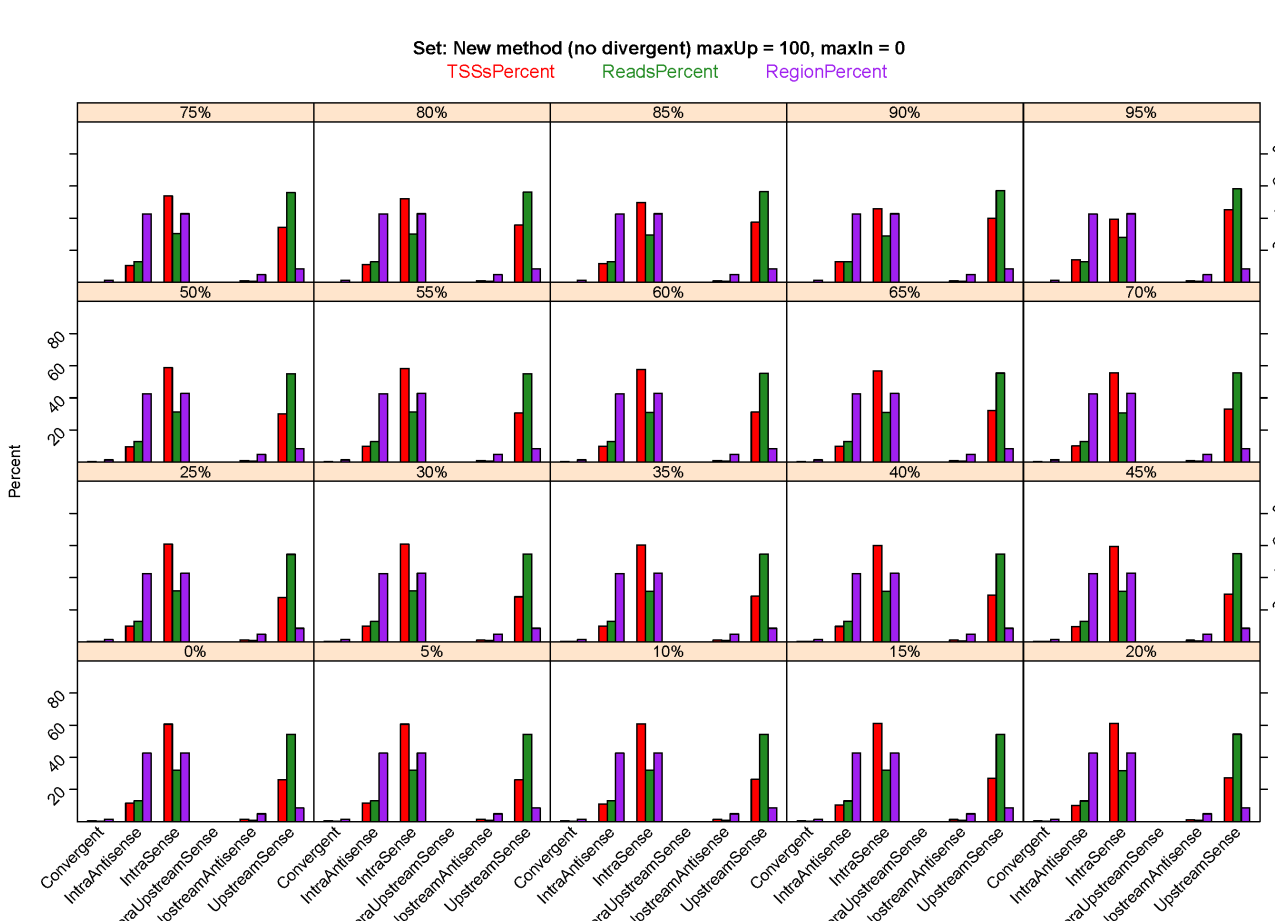
Transcription Start Sites



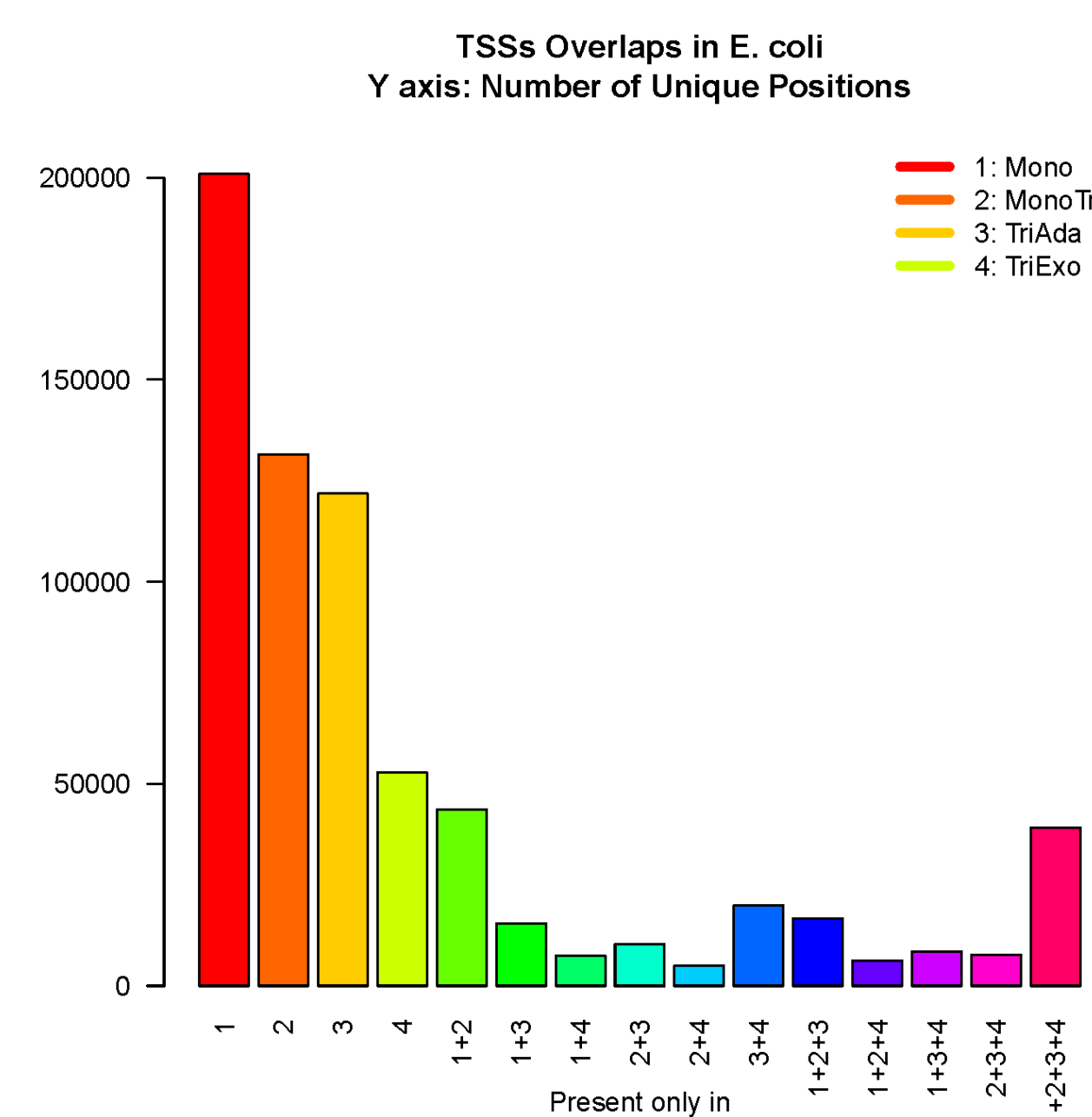
Choosing the minimum read length



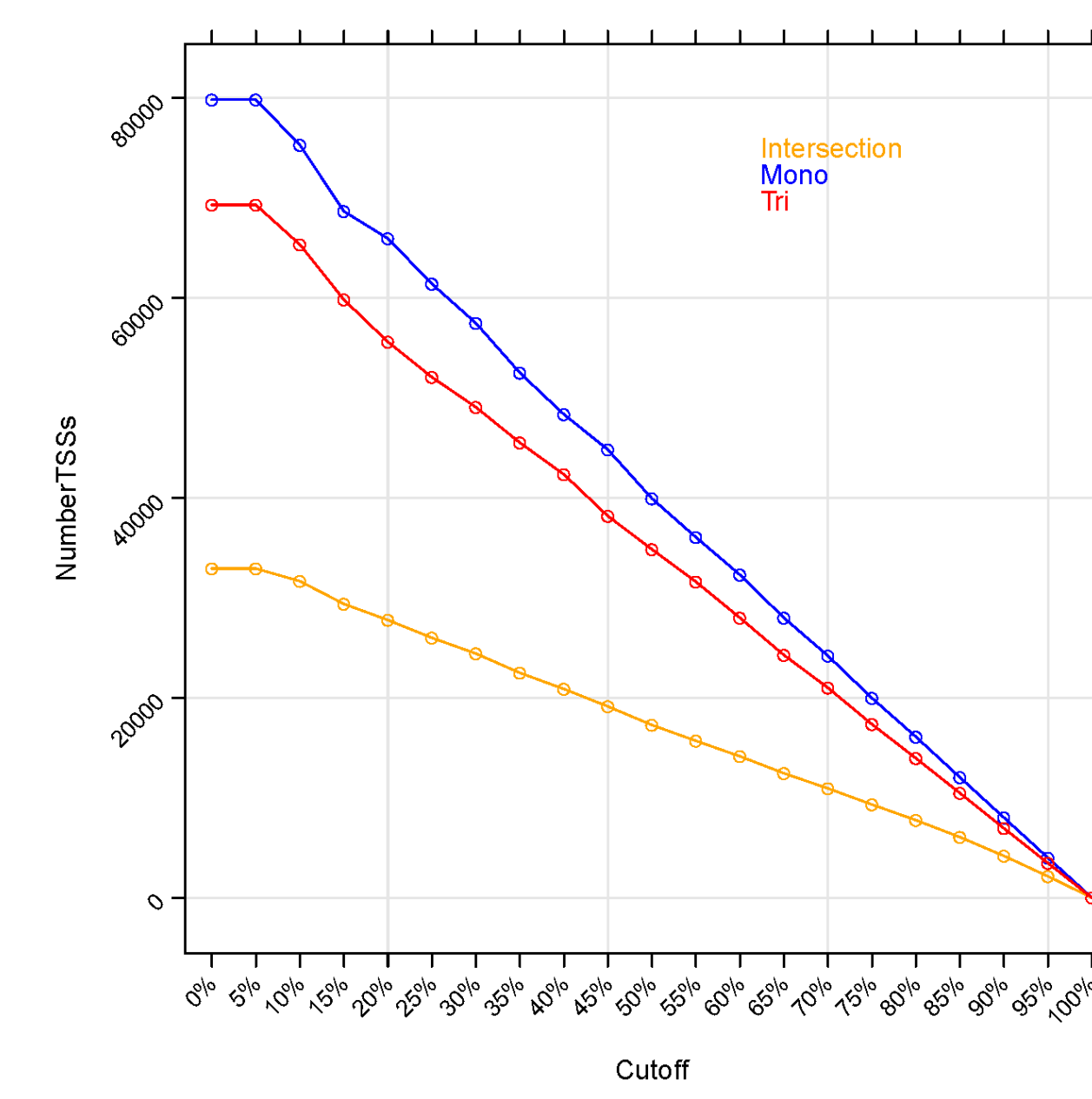
Read coverage (frequency) of the putative TSSs



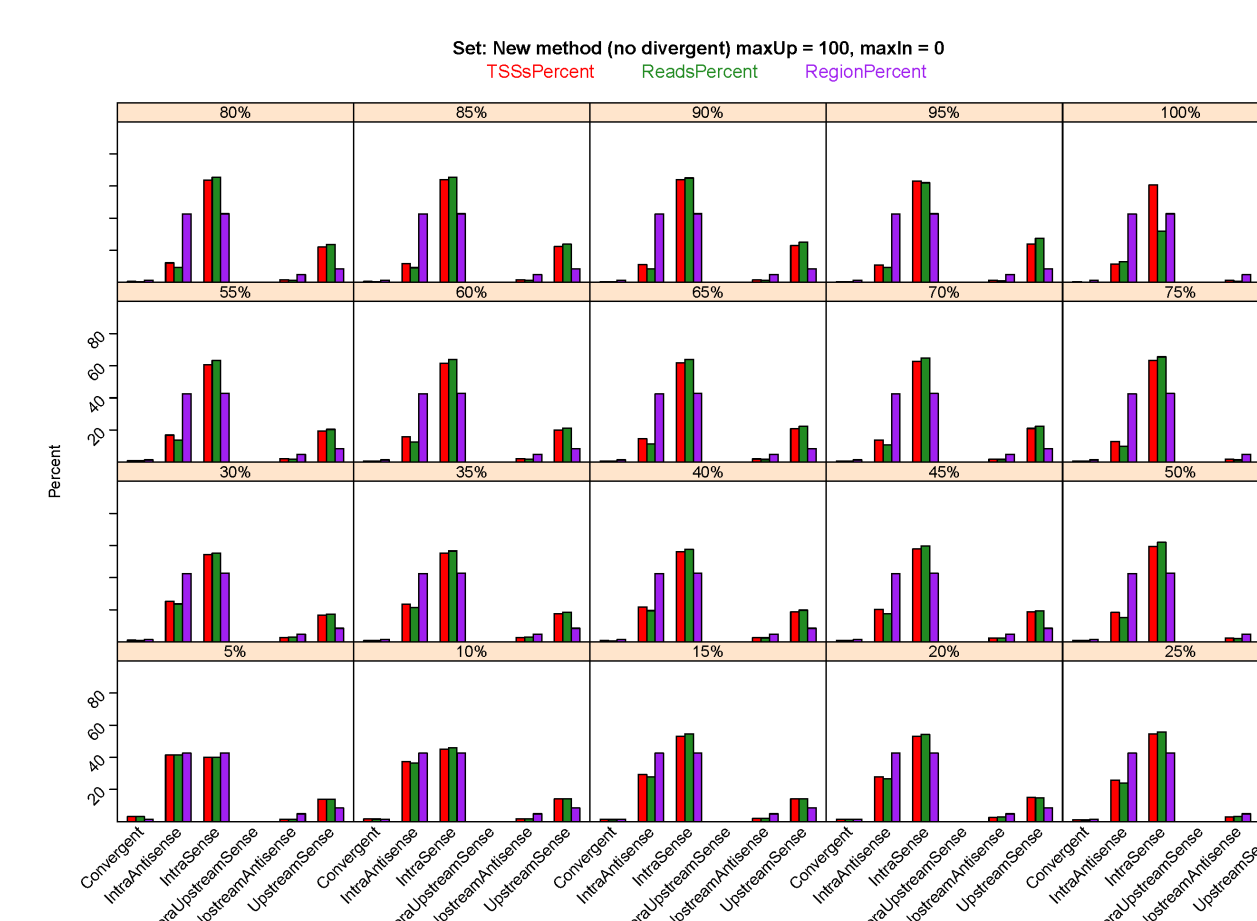
TSSs in genomic regions be removing low frequency TSSs



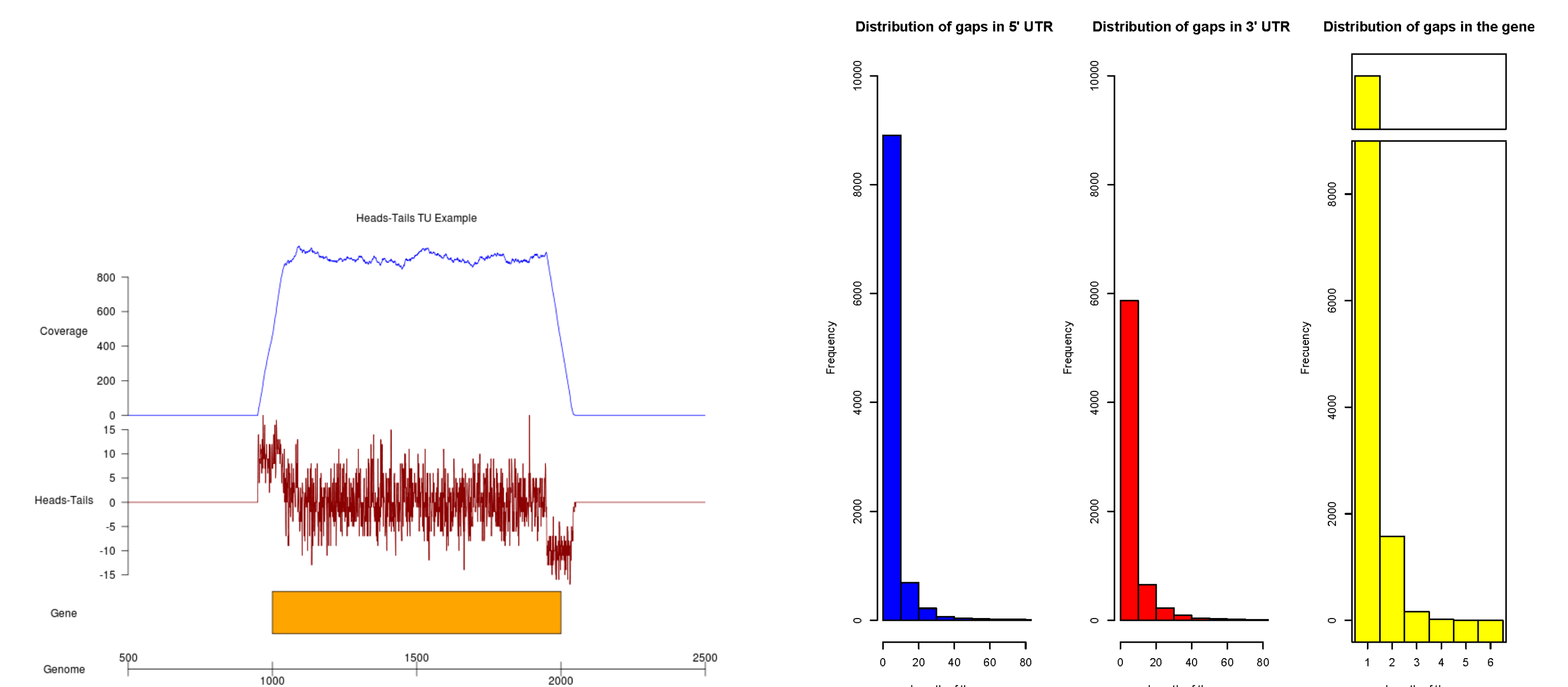
Number of putative TSSs positions



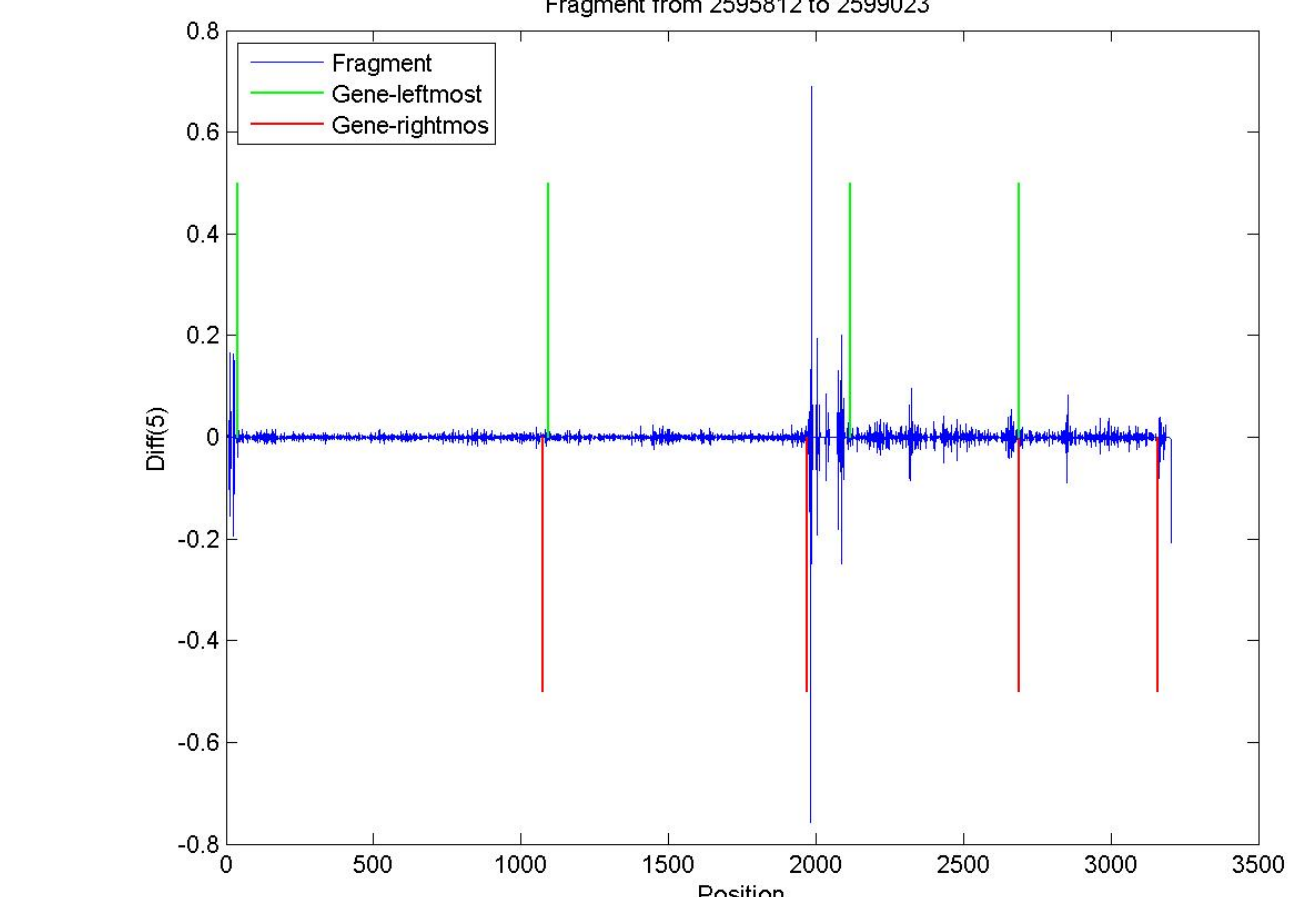
Number of TSSs related to the stringency



TSSs in genomic regions be removing high frequency TSSs



Left position gaps method



Differentiation method

Conclusions

These preliminary analyses will lead to the improvement of the accuracy of promoter prediction, operon structure and regulatory networks and support a new understanding from a genome perspective of the complex regulatory network that governs transcription and regulation in bacterial genomes such as *E. coli* and *G. sulfurreducens*.

References

Mendoza-Vargas, A. *et al.* Genome-Wide Identification of Transcription Start Sites, Promoters and Transcription Factor Binding Sites in *E. coli*. PLoS ONE (2009).
Gama-Castro, S. *et al.* RegulonDB version 7.0: transcriptional regulation of Escherichia coli K-12 integrated within genetic sensory response units (Gensor Units). Nucleic Acids Res (2010).
Collado-Torres, L. *et al.* manuscript in preparation.
BacterialTranscription Bioconductor package in preparation. More on the EMBL 2010 Bioconductor Developer Meeting.
We acknowledge support from NIGMS-NH grant R01 GM071962-05 and from CONACYT México (83686 G.I.).