## TITLE: DECREASING THE EFFECTS OF GC BIAS ON DE NOVO GENOME ASSEMBLY

AUTHORS : PINHEIRO, K.C.; MAUES, D.; SILVA, A.; RAMOS, R. T. J.

**INSTITUTIONS:** 1-INSTITUTE OF BIOLOGICAL SCIENCES, LABORATORY OF BIOLOGICAL ENGINEERING, FEDERAL UNIVERSITY OF PARA (R. AUGUSTO CORRÊA, 1 - GUAMÁ, BELÉM - PA, 66075-110)-BRAZIL

## **ABSTRACT:**

The emergence of high throughput sequencing (HTS) platforms increased the amount of data making feasible to obtaining complete genomes. Despite the advantages and the throughput produced by these platforms, the different genomic coverage in the regions of the genome can be related to GC content. This GC bias may affect genomic analyzes and the genomic/transcriptomic analysis based on de novo and reference approach. In addition, the ways to evaluate the GC bias should be fit to data with different profiles of the relationship GC vs coverage, such as linear and non-linear. We used the data of 5 species: Pseudomonas fluorescens, Shewanella amazonensis, Escherichia coli, Staphylococcus aureus and Mycobacterium tuberculosis sequenced by Illumina platform. This paper proposes, the use of a new metric to measure GC bias regardless of pattern observed between GC content and coverage, then, based on a degree of GC Bias, use a new approach to identify and reduce gaps from such regions by run a sliding window in a reference genome and separate the reads by GC Content range, creating groups, to submit them to a tools to discover the best k for a k-mer de novo assembly. Our analysis showed that the new metric called "Median - IQR" (difference between Median and Interguartile range - IQR) is able measure different bias patterns (monotonic and non-monotonic) and also could capture different degrees of bias in different samples. Also, the new method to assembly data with high degree of GC bias could represent regions that would not be identified in data assembled with default parameters.

Keywords: GC Bias, Median, Interquartile range, assembly

Development Agency: CNPq, CAPES











Contig alignment viewer. Contigs aligned to "Mycobacterium tuberculosis H37Rv"