

# *De novo* assembly of *Trypanosoma cruzi* strain CL Brener transcriptome

Eddie Luidy Imada<sup>1</sup>, Mainá Bitar<sup>2</sup>, Máira Ribeiro Rodrigues<sup>1</sup>, Daniela Ferreira Chame<sup>1</sup>, Helaine Grazielle dos Santos Vieira<sup>3</sup>, Michele Araújo Pereira<sup>1</sup>, André Martins Reis<sup>1</sup>, Dominik Kaczorowski<sup>3</sup>, Willian Santos Prado<sup>1</sup>, Andréa M. Macedo<sup>1</sup>, Carlos R. Machado<sup>1</sup>, Martin A. Smith<sup>3</sup>, Glória R. Franco<sup>1</sup>

<sup>1</sup>Departamento de Bioquímica e Imunologia, UFMG, Belo Horizonte – Brazil, <sup>2</sup>QIMR Berghofer Medical Research Institute, Brisbane – Australia, <sup>3</sup>Garvan Institute of Medical Research, Sydney – Australia

Chagas disease is a neglected tropical disease caused by *Trypanosoma cruzi* that is estimated to affect at least 6 million people worldwide. The first genome draft of the hybrid clone CL Brener was published in 2005 as several scaffolds and contigs, which were latter partially assembled into 82 chromosomes in mid 2009. Despite the great improvement of the *T. cruzi* genome assembly it still presented many gaps and unplaced contigs. Another caveat of the current genome is that its current annotation are almost exclusively based on automatic ORF detection, which might skew transcriptome analysis due to the lack of unannotated genomic elements. To address these problems we have assembled a *de novo* transcriptome of *T. cruzi* CL Brener using ~70 million paired end Illumina Hiseq 2500 reads with Trinity. The assembly resulted in 57,084 transcripts that were clustered using CD-HIT-EST into 24,844 non-redundant clusters. The clusters representatives (herein referred simply as transcripts) were annotated using a 2-steps methodology by first searching for homologues with BLASTN at nucleotide level against current predicted CL Brener transcripts and then searching at protein level with BLASTX against the TriTrypDB for those failing to align at nucleotide level. As a result, 93% of the transcripts were assigned to 14883 annotations, with 6038 genes being completely covered and 8,665 genes assembled by at least half of its predict length. Transcripts expression evaluation with Kallisto and GO annotation showed that the most expressed genes were related to basic metabolism as expected since steady epimastigote cultures were used in this work. Furthermore, using the assembled transcriptome we were able to close over 400 gaps in the current genome, improve/create UTR annotations for 9206 genes and correct some misplaced contigs in the assembled chromosomes. Our approach shows how *de novo* transcriptomes can be leveraged not only to functional studies, but to improve genome assemblies and annotations as well.

Financial Support: CAPES, FAPEMIG