

Cis-regulatory G-quadruplex Motifs are Preferentially Associated with Splice Sites in the Protein-Coding Human Genome

! " # \$! " %

April 2020

Abstract

Expression of mammalian genes involves regulated RNA splicing. Most human genes undergo alternative splicing during gene expression. As a result, the human protein-coding genome provides a rich variety of proteins with complex and diverse functions. It is estimated that up to one-fifth of human diseases are associated with altered splicing.

Our study studies the role of cis-regulatory motifs, such as G-quadruplex forming G-rich sequences (GRGs) in RNA processing. We focus on computationally identifying GRG distribution patterns near splice sites in the human protein-coding genome and investigate their role in regulated splicing. Our dataset consists of 1,000 genes, 1,000 exons, 1,000 introns, and 1,000 unique splice sites based on the hg38 human genome assembly extracted from the human Ensembl database. We have developed scripts in Python and Perl, based on our previously established GRG Mapper program, to map GRG motifs.

Our analysis discovered a preferential association of GRG motifs with splice sites in exons