

# Creating Regular Expressions as mRNA Motifs with GP to Predict Human Exon Splitting

W. B. Langdon, J. Rowsell, A. P. Harrison

Crest, Department of Computer Science, King's College, London, WC2R 2LS, UK  
Departments of Mathematical and Biological Sciences, University of Essex, UK

Wi11iam.Langdon@kcl.ac.uk, Mark.Harman@kcl.ac.uk, Yue.Jia@kcl.ac.uk

## ABSTRACT

RNAnet [3] <http://bioinformatics.essex.ac.uk/users/wlangdon/rnanet/> allows the user to calculate correlations of gene expression, both between genes and between components within genes. We investigate all of Ensembl and find all the Homo Sapiens exons for which there are sufficient robust Affymetrix HG-U133 Plus 2 GeneChip probes. Calculating correlation between mRNA probe measurements for the same exon shows many exons whose components are consistently up regulated and down regulated. However we identify other Ensembl exons where sub-regions within them are self consistent but these transcript blocks are not well correlated with other blocks in the same exon. We suggest many current Ensembl exon definitions are incomplete.

Secondly, having identified exon with substructure we use machine learning to try and identify patterns in the DNA sequence lying between blocks of high correlation which might yield biological or technological explanations. A Backus-Naur form (BNF) context-free grammar constrains strongly typed genetic programming (STGP) to evolve biological motifs in the form of regular expressions (RE) (e.g. TCTTT) which classify gene exons with potential alternative mRNA expression from those without. We show biological patterns can be data mined by a GP written in `gawk` and using `grep` from NCBI's GEO database. The automatically produced DNA motifs suggest that alternative polyadenylation is not responsible. (Full version in TR-09-02 [7].)

Blocky exons can be found in <http://bioinformatics.essex.ac.uk/users/wlangdon/tr-09-02.tar.gz>

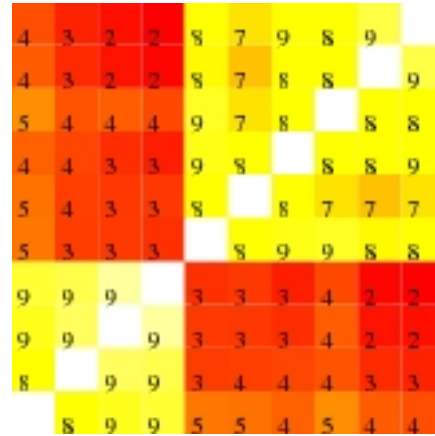
**Categories and Subject Descriptors:** J.3 [Life and Medical Sciences]: Biology and Genetics

**General Terms:** Experimentation

## 1. INTRODUCTION

Particularly in Man, there are multiple ways messenger RNA is transcribed from its gene. E.g. alternative splicing may lead to the removal of exons or to exons being repeated. Usually it is assumed that exons are indivisible. Usually measurements of mRNA taken at various positions in an exon are highly correlated. However Figure 1 shows this is not true for an exon in Aspartyl-tRNA synthetase. Indeed we have found several hundred other Homo sapiens Ensembl exons with pronounced structure in their correlation heatmaps.

Copyright is held by the author/owner(s).  
GECCO '09, July 8–12, 2009, Montreal, Canada  
ACM 978-1-60558-325-9/09/07.



**Figure 1: Correlation ( $\times 10$ ) between different locations in an Ensembl exon across 2757 tissue samples. The first four locations (lower left) clearly fall into a different block than the others (top right).**

Having identified these DNA sequences, we used strongly typed grammar based [5] GP [1, 6, 9] to evolve biologically meaningful motifs to explain them using only the mRNA sequences. A motif gets high fitness if it matches many blocky exons (i.e. positive examples) but fails to match many exons with high correlations but without blocks (negative examples).

We excluded suspect GeneChip data [8], suspect probe sequences [4] and probes which map to more than one exon [11] from the training data. The regular expression grammar in [4] was used, except  $\sim$  and  $\$$  were omitted and the maximum Kleene closures was 7 instead of 9.

## 2. FITNESS OF STGP RE MOTIFS

In each generation, GP generates a unix command file which contains an `grep -c 'RE'` command for each individual in the population [2]. (*RE* is the individual's regular expression.) The command is run on a file holding the 100 sequences lying between two blocks. (Collars of 50 additional bases mean the last 50 bases of the first well correlated block and the first 50 bases of the trailing well correlated block are also included.)

`grep -c` counts the number of probes which match the evolved motif (*RE*). The same command is also run on a file holding the 100 sequences of exactly the same length taken from exons which do not contain well separated blocks.

