

**Functional analysis of repeat regions in the eukaryotic  
genomes**

**Functional analysis of repetitive DNA derived from transposable  
elements in the human genome**

Lu Zeng

A thesis submitted for the degree of Master of Philosophy

Discipline of Genetics

School of Molecular and Biomedical Science

The University of Adelaide

July 2013

## Table of Contents

<b>Abstract</b> .....	<b>II</b>
<b>Declaration</b> .....	<b>III</b>
<b>Acknowledgements</b> .....	<b>IV</b>
<b>STATEMENT OF AUTHORSHIP</b> .....	<b>V</b>
<b>Chapter 1</b> .....	<b>1</b>
<b>Introduction</b> .....	<b>2</b>
<b>1.1 Background</b> .....	<b>2</b>
<b>1.2 Research questions</b> .....	<b>3</b>
<b>1.3 Aims and objectives</b> .....	<b>3</b>
<b>1.4 Significance</b> .....	<b>4</b>
1.4.1 Definition and classification of TEs .....	4
1.4.2 Functions of TEs.....	5
1.4.3 Association between RNAs and TEs .....	7
1.4.4 Conclusions .....	8
<b>2 Methods</b> .....	<b>8</b>
<b>2.1 Theoretical framework and methods</b> .....	<b>8</b>
2.1.1 The pipeline for the identification and distribution of functional repetitive elements from human genome.....	8
<b>2.2 Functional analysis of human/bovine repeats</b> .....	<b>11</b>
2.2.1 Repetitive element expression in different human tissues.....	12
<b>2.3 Relationship between lincRNAs and TEs</b> .....	<b>12</b>
<b>3 Results</b> .....	<b>13</b>
<b>3.1 The distribution of chromatin state associated transposable elements (CSTEs) from six different cell lines</b> .....	<b>13</b>
<b>3.2 The proportions of different repeat classes in active chromatin from six distinct cell lines</b> .....	<b>15</b>
<b>3.3 Repeat sequence distribution in the human genome</b> .....	<b>15</b>
<b>3.4 Functional representation of repeat consensus sequence in the human genome</b> .....	<b>16</b>
<b>3.5 The effect of Alu, L1 and LTR on gene expression in 6 human tissues</b> .....	<b>17</b>
<b>3.6 Are specific repeat sequences present in lincRNAs?</b> .....	<b>18</b>
<b>Discussion</b> .....	<b>20</b>
<b>Future Directions</b> .....	<b>23</b>
<b>Abbreviations List:</b> .....	<b>24</b>
<b>Figures and table legends:</b> .....	<b>24</b>
<b>Supplementary Materials</b> .....	<b>47</b>
<b>Reference</b> .....	<b>85</b>

# Abstract

Nearly half of the human genome is made up of transposable elements (TEs). With the rapid progress of sequencing technologies, we are now much better able to systematically analyze these TEs. We have used multiple types of omics data, including the genomic sequences, epigenetic data and transcriptomic data, to investigate the potential functions of TEs across the entire human genome. Comparative analysis revealed that a large proportion of potentially functional transposable elements were located in introns, and they were mainly associated with gene repression. Functional classification from GO enrichment showed that different functions were enriched in protein coding regions containing TEs compared to non-protein coding regions. For example, protein coding genes with Alus in non-coding regions are enriched with respect to intracellular membrane-bounded organelles, while protein coding genes with Alus in coding regions are more enriched at intracellular non-membrane-bounded organelles. Significantly, transcriptome data showed that the genes with TEs had lower expression levels compared with genes without TEs, revealing a novel aspect of the impact of TEs on the human genome. In addition, genome wide analysis of repeats with regulatory elements showed that MIR and L2 repeats were more probable to be active regulators while L1 repeats were less probable to be regulators. In conclusion, the role of TEs is significant across the genome. Repeats reduce or repress the expression of related gene, either through the proximal promoter, 5'UTR or 3'UTR or perhaps as components of lincRNA exons.

# Declaration

This work contains no material which has been accepted for the award of any other degree or diploma in any university or other tertiary institution to Lu Zeng, and to the best of my knowledge and belief, contains no material previously published or written by another person, except where due reference has been made in the text.

I give consent to this copy of my thesis, when deposited in the University Library, being made available for loan and photocopying, subject to the provisions of the Copyright Act 1968.

The author acknowledges that copyright holder(s) of these works contained within this thesis (as listed below) resides with the copyright holder(s) of those works.

I also give permission for the digital version of my thesis to be made available on the web, via the University's digital research repository, the Library catalogue, and also through web search engines, unless permission has been granted by the University to restrict access for a period of time.

Signed..... Date.....

# Acknowledgements

I would like to express my sincere gratitude to the following people: Dr. Dave Adelson, gave me the opportunity and support to join the Master by Research program in the University of Adelaide. I am so lucky to have such a great supervisor. Dave, the knowledge that I learned from you is not just valuable for my Master, but will support me for my whole academic career.

Dr. Chaochun Wei, my supervisor in SHANGHAI JIAOTONG University, introduced me to a completely new field of bioinformatics. Over the year I stayed in Adelaide, he was always in contact with Dave and me to co-advise my research. Furthermore, he gave me this opportunity to do the Master in Adelaide.

Dan Kortschak, my co-supervisor in the University of Adelaide, helped me do research and always provide valuable suggestions.

Joy Raison, helped me generate nice figures. Sim Lim Lin, helped me sort out the tough problems I encountered in my research, and kept encouraging me. Zhipeng Qu, made both my life and research easier and fluent. Reuben Buckley, helped me practice English and providing me valuable advices. All other members of the Adelson lab, past and present, made it such a supportive and enjoyable place to work.

I would like to also specially thank SHANGHAI JIAOTONG University and The University of Adelaide, who provided me this opportunity to come here for my research.

# STATEMENT OF AUTHORSHIP

## Functional analysis of repeat regions in the human genome

Submitted, July 2013

**Lu Zeng** (Candidate)

Designed and performed experiments, analyzed results and wrote the manuscript.

I hereby certify that the statement of contribution is accurate

Signed..... Date.....

**David L. Adelson & Chaochun Wei**

Supervised development of work and assisted in analyzing results and writing the manuscript.

I hereby certify that the statement of contribution is accurate and I give permission for inclusion of the paper in the thesis.

Signed..... Date.....

# **Functional analysis of repeat regions in the human genome**

**Lu Zeng**

**School of Molecular and Biomedical Science**

**The University of Adelaide**

**Adelaide, SA**

**Australia**

**School of Life Sciences and Biotechnology**

**SHANGHAI JIAOTONG University**

**Shanghai**

**P.R.China**

# **Chapter 1**

## **Functional analysis of repeat regions in the human genome**

Lu Zeng and David L. Adelson and Chaochun Wei

School of Molecular and Biomedical Science, The University of Adelaide,  
Adelaide, SA, Australia

School of Life Sciences and Biotechnology, SHANGHAI JIAOTONG  
University,  
P.R.China



# Introduction

## 1.1 Background

Eukaryotic genomes contain vast amounts of repetitive DNAs derived from TEs that contribute significantly to biological activity and genome evolution. Furthermore, TEs are mutagens; they may damage their host cells through various mechanisms [1]. For example, a transposon or a retrotransposon that inserts itself into a functional gene may disrupt or alter it, disrupting gene function. Similarly, a DNA transposon that excises from a genome may result in a deletion that cannot be repaired. Due to the presence of multiple copies of repetitive elements, such as Alu sequences, precise chromosomal pairing during meiosis may be deleted, causing unequal crossovers and deletion or insertion of genetic materials. Through these mutagenic mechanisms, repeats are known to cause a variety of human genetic disorders [2].

A number of recent studies have shown that TEs can influence host genes by providing novel promoters, splice sites or post-transcriptional modification to re-wire different developmental regulatory and transcriptional networks [3-5]. TEs tend to regulate gene expression through several mechanisms [5-7]. For example, the expression levels of protein coding genes containing repeats are significantly associated to the number of repeats in rodent genomes [6]. Moreover, TEs have been shown to influence gene expression through non-coding RNAs, resulting in the reduction or silencing of gene expression [8]. Past studies have also found that TEs have contributed to nearly half of the active regulatory elements to the human genome [9], such as altering gene promoters, creating alternative promoters and enhancers to regulate gene activity [10-12]. According to previous research, 60% of TEs in both human and mouse were located in intronic regions and all TE families in human and mouse can exonize [13], supporting the view that TEs may create new genes and exons by promoting the formation of

novel or alternative transcripts [14, 15]. The association between repeats and RNAs has also been investigated, some findings showed that tRNA can use TinT events to drive the formation of novel SINE [16]. Telomeric repeats may be transcribed as telomeric RNAs or telomeric repeat-containing RNAs [14, 15] and the insertion of TEs may also drive the evolution of lincRNAs and alter their biological functions [17].

## **1.2 Research questions**

In order to uncover the hidden information of TEs, my research will focus on these following questions:

- 1) What is the distribution of TEs in functional regions in the human genome? Functional regions here include protein coding genes, ncRNAs and regulatory elements like TFBS, promoters and enhancers.
- 2) What is the association between repetitive elements and functional elements in the human genome?
- 3) Do repetitive elements impact on gene expression?

## **1.3 Aims and objectives**

- 1) To analyze the distribution of repeat-associated functional elements in human genes.
- 2) To classify and analyze repeats in different cell lines
- 3) To build the consensus sequences for some important classes of repeats and use these sequences to identify full-length repeats and their function and distribution in the human genome.
- 4) To analyze the relationship between repeats and lincRNAs.
- 5) To conduct analysis of the expression level of specific repeats in the human genome.

## 1.4 Significance

### *1.4.1 Definition and classification of TEs*

TEs are DNA sequences that can change their position within the genome, potentially giving rise to mutations or altering genome size and structure [18]. These characteristics of TEs can affect biological activities and thus may contribute to genome evolution. Moreover, TEs are able to insert at new locations without having a sequence relationship with the target locus. TEs make up about 50 percent of the human genome.

Transposons fall two major classes: RNA (retrotransposons) and DNA (DNA transposons), according to whether their replication is via RNA or DNA [19] intermediates. DNA transposons use a cut-and-paste transposition mechanism instead of involving the RNA intermediate.

Retrotransposons include two classes of elements, autonomous and non-autonomous. Autonomous transposons contain open reading frames (ORFs), which encode proteins essential for transposition and are thus able to autonomously transpose. Non-autonomous transposons do not encode these functions and so rely on replication machinery provided by autonomous transposons.

Retrotransposons can also be separated into two groups with respect to different characteristics: Long terminal repeat (LTR), and Non-LTR. LTR retrotransposons have transcription control sequences and open reading frames encoding retrotranspositional activities [19]. They range in size from ~100bp to over 5kb. About 8% of the human genome and approximately 10% of the mouse genome are composed of LTR transposons [20].

Non-LTR retrotransposons include two sub-types; long interspersed elements (LINEs) and short interspersed elements (SINEs) respectively, both of which are widespread in eukaryotic genomes. Furthermore, LINEs are autonomous retrotransposons, while SINEs are non-autonomous retrotransposons.

LINEs [21] are genetic elements that contribute significantly to eukaryotic genomes, they are transcribed into RNA using an RNA polymerase II promoter. LINEs account for 17% of the human genome.

SINEs [21] are short DNA sequences, usually less than 500 bases long [22] originally transcribed by RNA polymerase III into tRNA, 5s ribosomal RNA and other small nuclear RNAs. The most common SINEs are Alu sequences, which account for 10.6% of the human genome.

According to previous research, SINEs and LINEs have similar nucleotide sequences at the 3' end, and SINEs usually dependent on LINE RT/EN function for transposition [21]. This finding was the starting point for the concept of LINE machinery involved in the retrotransposition of SINEs [23, 24]. Moreover, from the most recent human genome sequence, I found that there are 1500000 SINEs and 850,000 LINEs that account for 34% of the human genome in total. 70% of SINEs are Alu elements.

#### *1.4.2 Functions of TEs*

Retrotransposons can impact on human genome structure, which can dramatically affect genome evolution.

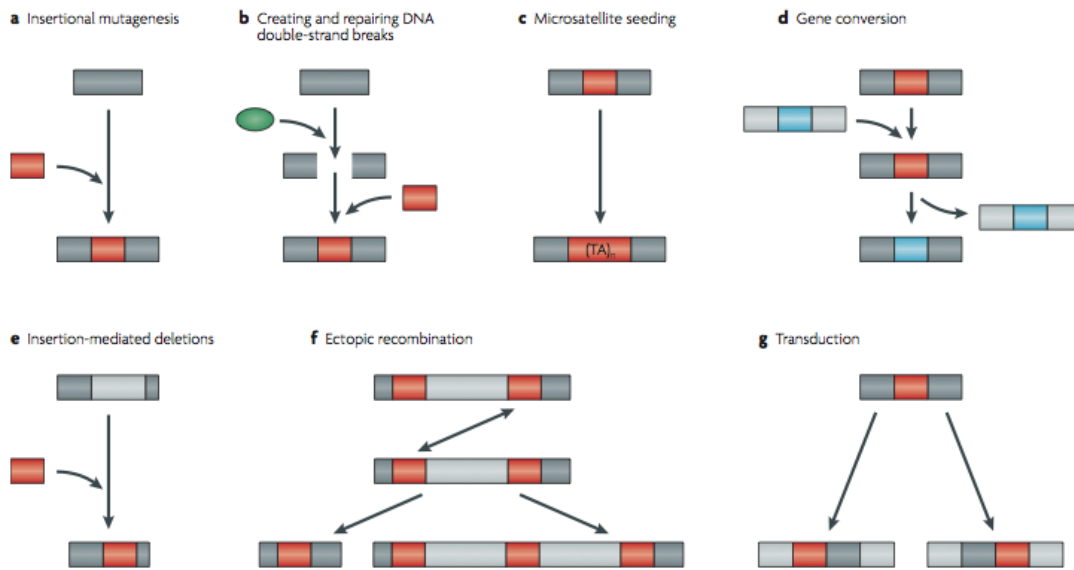


Figure 1: Richard et al [25] showed that how retrotransposons have an impact on human genome structure through 7 mechanisms.

Retrotransposons can affect the human gene expression.

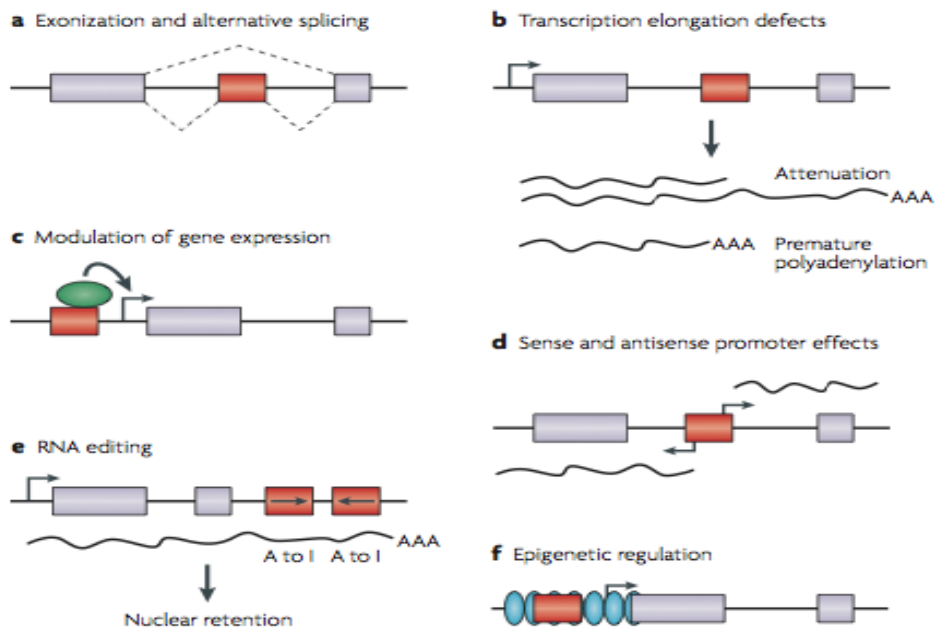


Figure 2: Retrotransposons impact on human gene expression from Richard et al [25]

#### *1.4.2.1 SINEs/Alu elements are primate-specific repeats and influence gene expression*

Alu insertion is ongoing in modern human genomes, including somatic insertion events, generating genetic diversity and causing disease through insertional mutagenesis as well as causing copy number variation. Many Alu elements affect polyadenylation [26, 27], splicing [28-30], and double-stranded RNA-specific adenosine deaminase (ADAR) editing [31, 32].

#### *1.4.2.2 LINE/L1 insertions have a high frequency of retrotransposition*

L1 elements can cause human disease by inserting into human genes. After transcription to RNA, they can be reverse-transcribed into cDNA and integrated into other genomic locations. LINE/L1 has two open reading frames, ORF1 encodes a nucleic acid binding protein [33, 34], and ORF2 encodes a protein with endonuclease activity [35], reverse transcriptase activity [36] and a C-terminal cysteine-rich motif [37]. The 5'UTR of LINEs contains an internal promoter sequence, while the 3'UTR has a polyadenylation signal and a poly-A tail.

#### *1.4.3 Association between RNAs and TEs*

SINEs are derived from RNA [38], for example, Alu elements come from the ubiquitous 7SL RNA [39]. A functional sequence within Alu RNA transcripts has revealed a modular structure analogous to the organization of domains in protein transcription factors. According to recent studies, telomeric repeats may be transcribed as telomeric RNAs or telomeric repeat-containing RNAs [14, 15]. Furthermore, ncRNAs can modulate the function of transcription factors, and as far as I know, retrotransposons also contain transcription factor binding sites, which can combine with transcription factors to alter gene expression.

#### *1.4.4 Conclusions*

The recent explosion of retrotransposon studies has brought about a great improvement in understanding of TEs. It is clear that gene-regulatory networks are complicated, but this is not just the realm of genes and proteins, but also repeats. It appears that genome structure, especially for complex organisms, is very complicated as well. Genomes possess a high proportion of repeat regions, which may represent a hidden level of gene regulation. TEs impact the transcriptome through both transcriptional and post-transcriptional [5], as well as some disease-related mechanisms. Although evidence suggests that TEs are highly expressed from different regions of genomes, and possess a wide range of functionality in gene regulation, these discoveries still constitute just a glimpse of the hidden repeats.

In general, most studies of TEs are constrained to several model organisms, such as human, mouse and cow. There are few, and sometimes no studies focusing on other well-known organisms, such as chicken, pigs and so on.

## **2 Methods**

### **2.1 Theoretical framework and methods**

The main framework of this project was to build a pipeline to analyze the distribution, function and expression of repeats from both human genes and bovine genes.

#### *2.1.1 The pipeline for the identification and distribution of functional repetitive elements from human genome*

The identification and classification of TEs from the human and bovine genome was conducted by developing a pipeline based on free software, Perl, R and The UCSC Genome

Browser (University of California, Santa Cruz) database. Perl is a programming language that can be used for a large variety of tasks. One of the most powerful functions of Perl is for extracting information from a text file and printing out a report or for converting a text file into another form. This feature makes Perl popular in bioinformatics. In this project, Perl was used as a glue language to conduct result parsing and program linking. R is a free software programming language and a software environment for statistical computing and graphics. The R language is widely used among statisticians and data miners for developing statistical software [40, 41]. In this project, R was used to build graphs in order to illustrate the distribution, classification and function of TEs. The UCSC Genome Browser is an on-line genome browser [42, 43] that offers access to genome sequence data from a variety of vertebrate and invertebrate species and major model organisms. In this project, I have used this service to retrieve the repeat data and RefGene annotation data for my experiment.

#### *2.1.1.1 Identification and distribution of functional TEs*

In order to study the distribution of repetitive elements in human genes as well as the classes of various classes of TEs that exists in the human genome, I collected the datasets that applied to my experiment. First, NCBI's human Reference Gene Collection (RefSeq hg19) [44] and the associated annotation table were downloaded from the UCSC genome browser [42, 45]. In order to analyze the function of repeat regions, I have downloaded the regulatory elements data of nine human cell lines from UCSC. These regulatory element annotations, including active promoters, weak promoters, strong enhancers, weak enhancers, insulators and polycomb repressed regions, which were derived from different chromatin states that have been marked by histone methylation, acetylation as well as histone variant H2AZ, PolIII, and CTCF [45]. I also chosen six human cell types from those nine cell lines that are useful for studying human disease, they are GM12878, HepG2, HMEC, HUVEC, K562, NHLF, Table2



shows the resource and information of these human tissues. I also have downloaded these datasets from group regulation track Broad ChromHMM and divided them into six parts according to their functional roles: active-promoter, weak-promoter, strong-enhancer, weak-enhancer, insulator and polycomb-repressed regions. Then, I retrieved the most recent human RefGene (hg19) from UCSC [46], separating it into different sections according to the human genome regions, which include 5'UTR, start codon, CDS exon, CDS intron, and stop codon and 3' UTR. Next, BED intersection was applied to get the overlap between RepeatMasker and Human regulatory elements, and then rerun overlap between the union data and Human RefGene respectively. From this operation I determined the distribution of repeat-associated regulatory elements with respect to human gene sections. I normalized my results with respect to the number of base pairs in each gene region.

#### *2.1.1.2 Analyze the classes of regulatory repetitive elements*

In this part, I have used BED intersection to get the overlap between regulatory elements in human cell lines and human RepeatMasker annotation. As I have already described, I have acquired the various classes of regulatory repetitive elements in different cell lines with respect to various regulators. According to the results I obtained in this part, in the next step, I built six consensus sequences in order to study specific TEs.

#### *2.1.1.3 Building consensus sequences*

I have identified thousands of short fragments of repetitive elements. Using consensus sequences from RepeatMasker database I can identify the 5'/3'UTRs only, often annotating repeats as having 5'/3'ends from different repeats. Thus, in order to study the impact of full-length specific retrotransposons on human gene structure and function; I have built complete consensus sequences of specific TEs.

Multiple sequence alignments of full-length sequences were performed using MUSCLE software with default parameters [47], these alignments were then used to run fasttree [48] to generate full-length repetitive elements respectively and used the tool archaeopteryx to generates the classification (See Supplementary Material S1) [49]. Next, I reran the multiple alignment tool MUSCLE to get the alignment sequences between different classes of each TEs, in the last step, I used Gblocks [50] to pile up these alignment results to acquire complete consensus sequences (See Supplementary S2).

Next, I used BED intersection to obtain the distribution of these consensus sequences in different human genes' sections; I only kept the first intron when I encountered alternative splicing. I normalized my data with respect to the length of relevant repetitive elements.

## **2.2 Functional analysis of human/bovine repeats**

To demonstrate the functional significance of repetitive elements, I used DAVID (The Database for Annotation, Visualization and Integrated Discovery) to perform the GO (Gene Ontology) classification, which represents gene product properties. First, I extracted the gene-IDs from the results that were overlapped with repeat consensus sequences in the human genome; I then submitted these gene-IDs to the DAVID Gene Functional Classification Tool [51]. From the results I chose the third level of GO terms to acquire the over-represented function terms of genes that contained repetitive elements. According to the GO term hierarchy; the third level of GO terms contains annotation categories for my analysis. Then I visualized the functional over-representation of genes overlapped with those six specific repeat consensus sequences in the human genome. The thresholds for over-represented GO terms were set as gene count >5 and p-value (EASE score) <0.05. The web

server REVIGO was used to reduce redundancy and visualize the overrepresented GO terms based on semantic similarity.

### *2.2.1 Repetitive element expression in different human tissues*

In order to analyze the expression level of TEs, I have taken advantage of the Human RNA-seqdata Illumina bodyMap2 transcriptome (<http://www.ebi.ac.uk/ena/data/view/ERP000546>) datasets. These RNA-seq data represent 16 different human tissues. Then, the normalized expression for each transcript dataset is based on my colleague (Z. Qu, unpublished) previous work in our lab. In order to conduct my work, I obtained the proximal promoters that belong to the gene upstream 1000bp. First, I obtained the overlap between repeat and human genes; I then intersected this data with proximal promoter dataset, 3'UTR and 5'UTR of the human genome. Next, I divided the data into repeats in proximal promoter and repeats not in proximal promoter. In the last step, I applied the intersection between the data I have already obtained and the transcriptome datasets, and obtained the expressed genes' interacting repeats in 16 human tissues. I found liver, kidney, testes, brain, skeletal muscle and adipose had relatively higher expression levels in my datasets. Therefore, I extracted these six tissues to analyze the expression level with respect to those three specific repetitive elements. To make the results more clear and significant in my graphs, I used log transformation to better separate the visualize differences in my expression data.

### **2.3 Relationship between lincRNAs and TEs**

Considering the previous study, we have known that TEs are a source of endogenous small RNAs in animals and plants, and they are considered as functionally significant from gene-regulating small RNAs [52]. Furthermore, according to past research, it was found that there are many Alu elements inserted into RNAs. Thus, I planned to study the association between lincRNAs and different classes of TEs.

I have acquired a dataset that contains 8196 putative human lincRNAs (Long intergenic non-coding RNA) [53], then I extracted the sequences of these lincRNAs from whole human sequences, and used BLASTN [54] to align these lincRNAs sequences against the previously constructed repeat consensus sequences I had built. In the next step, I extracted the aligned sequences from the BLAT result in order to analyse their coverage characteristic within lincRNAs.

## 3 Results

### 3.1 The distribution of chromatin state associated transposable elements (CSTEs) from six different cell lines

I used the pipeline I have built to study the distribution of transposable elements (TEs) in the human and cow genomes (Figure 3). The large numbers of TEs overlapping with genome intervals from the human and cow RefGene datasets showed that TEs was enriched in different gene sections.

I found that TEs have a similar distribution between the human and cow genome; in total, they are mainly located on 5' UTR intron, CDS intron, 3'UTR intron and intergenic regions. Comparing the TE locations between these two species, I found there were more TEs in each section of the human genome. Intergenic regions in particular have a higher number of TEs in the human compared to cow. The exception is for 5'UTR introns, where there are more TEs in the cow genome (Figure 3).

When I examined to chromatin state functional transposable elements (CSTEs), which represent the TEs that overlap with chromatin state segmentation [45] for each of six human

cell types, I chose six states that have predicted functional elements. Considering the CSTE distribution in the human genome, I divided my results into two parts, the first part is based on CSTE association with regulators in these functional cell lines, the second part is the relationship between CSTEs and different types of those chromatin human cell lines (Figure 4A, B), the human cell lines based on the tissues karyotype (Table 2). My aim was to look at how these two datasets are distributed in the human genome. When we are looked at the CSTEs that overlap with regulatory region distributions in the human genome, I found that the TEs that overlap with active promoters are the most highly represented in the gene start region. However, CSTEs that overlap with weak enhancers and polycomb-repressed regions are enriched in other regions of the human genome (Figure 4A). Among them, I discovered that most of CSTEs' that overlap with regulatory elements are located within the intergenic regions, followed by 5' UTR introns and CDS introns, which may indicate that TEs affect gene transcription through alteration of intron length. Moreover, the CSTEs that may down regulate the gene transcription were always located within the intergenic region, while regulators that can up regulate gene expression were usually located in the 5'UTRs and CDS introns.

In the analysis of the distribution of histone modification that overlapped with TEs in the six human cell lines (Figure 4B), I found that blood cancer cells had CSTEs that were highly enriched in every gene section, while in other cell types, CSTEs were less enriched in various gene features, except for intergenic regions, all of the six cell lines that contained CSTEs were highly enriched in this part, especially blood cancer, lung normal and blood vessel normal. To the contrary, breast normal had CSTEs were always located on the lowest enrichment in those gene sections, which may indicates that the chromatin state breast normal is less active when it contains TEs.

### **3.2 The proportions of different repeat classes in active chromatin from six distinct cell lines**

In this section I set out to study the distribution of specific classes of TEs that overlap with active chromatin from functional regions of the human genome (Figure 5A, B).

Active chromatin from blood cancer cells contained the highest percentage of repeats from virtually all repeat classes, but active chromatin from liver cancer cells ranked either in the mid-range or towards the bottom of repeat percentage (Figure 5B). Over-represented repeats in active chromatin include SINE Mir, LINE L2 and DNA transposons, which are normally represent as much lower percentages of the genome as a whole. LINE L1 elements were present at a surprisingly low level in active chromatin from all cell lines.

When I looked at the overlap of functional regulators (extracted from active chromatin regions) with TEs, I found a similar pattern in each repeat type, with repressed region containing the highest proportion of repeat sequence and active promoters the smallest proportion of repeat sequence (Figure 5A). I also found that MIR and L2 elements overlapped more with regulators compared to other TEs, while L1 had the lowest coverage percentage for all regulators. Because MIR and L2 are molecular fossils, while L1 is a currently active repeat type, this may indicate that ancient TEs may have been exapted during the course of evolution. From this graph I also can conclude that while the single biggest association of repeat coverage is with repressor activity, there is also a strong association with up-regulators of gene expression.

### **3.3 Repeat sequence distribution in the human genome**

I have generated six repeat consensus sequences that represent the principal classes of TEs in the human genome; they are MIR, Alu, L1, LTR, SVA and ERV. These classes of retrotransposons cannot only affect the human genome structure, but also gene transcription activity. I have determined the distribution of these six specific TE classes with respect to different gene features (Figure 6). I found that a large proportion of these TEs are located in 5'UTR, cds-intron and intergenic regions, with L1 and Alu the most common repeats represented.

### **3.4 Functional representation of repeat consensus sequence in the human genome**

I have mapped six primary repeat consensus sequences to the human genome, including gene models. I have used Gene Ontology (GO) annotation of the genes that overlapped with repeats to assess the functional contribution of repeats. By comparing the P-value as my standard index to acquire the definite functions of these repeat consensus sequences.

GO category results fall into three categories: cellular component, molecular function and biological process, with respect to their gene product properties (Table 1). I discovered that most of TEs' are associated with specific functional products that execute biological processes. Figure.7A illustrates that L1 and MIR are significant associated with genes that have many functions in the human genome. L1 in particular is associated with genes that have many functions in biological processes, such as establishment of localization in cell, cell-cell adhesion and cellular component morphogenesis. Moreover, gene associated with MIR participates in cell development processes, protein transport and neuron projection development. ERV associated genes are significantly over-represented with respect to cell-cell adhesion. Furthermore, Mir and L1 have similar clustering with respect to their functions if the genes they are associated with in terms of biological processes.

The second most annotated gene product property is cellular components (Figure 7B). Genes containing L1 and MIR have a number of associated functions. For example, MIR containing genes act as a part of plasma membrane. L1 containing genes are apparently enriched in the functions extrinsic to membrane, dendritic spine, dendritic shaft, extracellular matrix and protein serine/threonine. ERV containing genes have several functions as well, for example, endomembrane system and constitution of the membrane.

The third gene product property is molecular functions (Figure 7C). L1 containing genes functions include: purine nucleotide binding, substrate-specific transmembrane transporter activity and calmodulin binding. Genes containing ERV have functions including cell adhesion molecular binding and GTPase activator activity and enzyme binding. MIR containing genes have functions such as passive transmembrane transporter activity, while genes containing LTR participate in cell adhesion molecular binding.

### **3.5 The effect of Alu, L1 and LTR on gene expression in 6 human tissues**

I analyzed the effect of different TEs on human gene expression by comparing the expression of genes with TEs in either the proximal promoter region, the 5'UTR or the 3'UTR to genes without TEs in those regions. I selected six human tissue transcriptome sets (kidney, liver, brain, testes, skeletal muscle and adipose tissue) from the Illumina BodyMap2 dataset for this analysis.

In order to carry out the analysis, I compared the genomic intervals for TEs to the genomic intervals for gene models to identify genes with and without TEs in the proximal promoter, 5'UTR and 3'UTR. I then plotted log normalized gene expression levels (determined by



TopHat/Cufflinks) of genes from the different TE categories (Figure 8, 9 and 10) to see if overall gene expression levels were correlated with the presence or absence of TEs.

I found that Alus, L1s and LTRs appeared to be associated with lower levels of gene expression in all tissues if they were present in the proximal promoter regions. The exception to this was for testes expressed (Figure 8), where the presence of L1 in the proximal promoter made no difference to gene expression and where the effects of other TEs were also less obvious. I did not see changes in gene expression associated with repeats in the 5'UTR, but in the 3'UTR the presence of LTR was associated with decreased gene expression in all six tissues. These 2 suggest that TEs in proximal promoters and 3'UTR could possibly act as repressors of gene expression.

### **3.6 Are specific repeat sequences present in lincRNAs?**

In order to determine if TEs contribute specific sequences to Long Intergenic Non-Coding RNAs (lincRNAs), I aligned lincRNAs [53] and my TE sequences with BLASTN [54]. I found that many lincRNA sequences contained almost full length Alu, but that the alignments showed a greater contribution of left and right monomer sequences from Alu elements (Figure 11A). In particular, the left monomer sequence peak (Figure 11A) might correspond with the internal Alu promoter region. This result suggests a potential regulatory role for Alu sequences in lincRNAs expression.

Mir elements contributed the sequences between position 50 and 200 of their consensus sequence to lincRNAs (Figure 11B). When I extracted this subsequence and scanned it for regulatory motifs (<http://asp.ii.uib.no:8090/cgi-bin/CONSITE/consite>), I found this 150bp

subsequence contained low complexity sequence with motifs for Myf and Snail binding snail, CREB and bZIP910.

L1 elements are generally about 7kbp long and L1 alignments to lincRNAs showed a preponderance of sequences from the 3'end of L1, consistent with known 5'truncation of inserted L1 sequences (Figure 11C).

LTR sequences were generally found at low levels in lincRNAs except for a sharp 17bp peak in the middle of the LTR sequence (Figure 11D). When this sequence was extracted from the LTR it was found to be a low complexity tetranucleotide repeat sequence. This 17bp sequence has been annotated as containing regulatory motifs (See Supplementary Material S3).

ERV sequences were also found a generally low level in lincRNAs, but also showed some sharp spikes (Figure 11E). I extracted the subsequences for the two tallest spikes, 21bp at about 3kbp and 17bp at about 7kbp within the ERV consensus. The 21bp subsequence contains a binding site for the Pbx transcription factor and the 17bp subsequence contains low complexity sequence with motifs for Myf and Snail binding (See Supplementary Material S3).

Initial alignment of SVA sequences to lincRNAs showed that the Alu domain of SVA was probably aligning to the Alu regions of lincRNAs (Figure 11F). The VNTR domain of SVA showed large numbers of hits in lincRNAs and the 100bp subsequence from this region contained binding motifs for 8 transcription factors (See Supplementary Material S3).

The fact that alignment peaks in lincRNAs contained transcription factor binding motifs contributed from 4 TEs is suggestive of a role for these motifs in lincRNAs regulation or function.

## Discussion

In this work, I have analyzed the function, distribution and expression of various classes of TEs in the human genome. I have found that TE distribution is similar in human and cow that TEs appear to be able to alter gene expression, based on their distribution in functional regions and their correlation with gene expression. Furthermore, I found that genes containing TEs were over-represented in terms of their GO annotations in all three categories. Finally, I found that TEs are present in ncRNAs, specifically lincRNAs, and that some repeat classes appeared to contribute specific subsequences that contain functions to the lincRNAs.

TEs had similar distributions in human and cow (Figure 3), with a greater proportion of sequence originating from TEs in the 5' and 3'UTRs compared to coding exons or start and stop codons. This is not surprising considering the potential adverse effect of TE insertion in a protein coding sequence, but it is also relevant with respect to the known regulatory functions within the UTRs [55, 56]. The repeat content for 5'UTR introns was comparable to remaining introns, but this may be significant in the context of transcriptional repression, where genes with shorter 5' UTR introns were expressed at higher levels [57, 58]. It would appear that in some genes the exaptation of repeats into UTRs is neutral or functional (see below).

The distribution of TEs in genomic regions epigenetically modified to regulate transcription (active chromatin) was consistent with a potential role as regulators of gene expression. From my results (Figure 4A, B), I found that some functional regions of active chromatin contained

higher levels of TE, specifically polycomb-repressed regions and weak enhancers. This was most striking for the polycomb-repressed regions, which was consistent with epigenetic silencing of TEs. This was also consistent with the higher level of repeat content seen in 5'UTR, which are also known to regulate gene [59]. Furthermore, it has been shown that TEs in 3'UTR are associated with lower transcript abundance [60]. This suggests that exaptation of repeat into regulatory regions is most often associated with repression of gene expression. In fact, older repeat, such as fossil L2 and MIR, are more prevalent in active chromatin than recently inserted repeats, such as L1 or Alu, supporting the argument for exaptation.

Different repeat classes were clearly present at different levels in active chromatin or specific regulatory regions, such as polycomb-repressed (Figure 5), implying an association with gene expression. My analysis of gene expression as a function of regulatory region repeats content indicated that there was a general association between repeat content and lower gene expression (Figures 8, 9 and 10). These results are consistent with previous reports showing TEs such as Alus can be exapted as regulators of alternative splicing, as transcription factor binding sites and as sense and antisense promoters [32, 61]. Furthermore, the presence of Alu pairs in opposite orientations in gene transcripts can lead to adenosine to inosine editing, resulting in suppression of expression through nuclear retention of edited RNA transcripts [25]. My results are consistent with TEs mainly down regulating gene expression if they were present in the proximal promoter or 3'UTR region. This association was weakest or non-existent for L1 elements (Figures 8, 9 and 10). Furthermore, L1 elements were less prevalent in regulatory regions or active chromatin, compared to other repeat classes (Figure 5). This makes sense because most L1 elements in the human genome are 5' truncated [60], lacking promoters and other regulatory sequences. LTR repeat were similarly distributed, compared to non-LTR repeats, and associated with repression of gene expression. This is in contrast to

previous work that implicates LTRs as alternative promoters [62]. However, an LTR in the first intron of the equine TRPM gene suppresses gene expression by acting as an alternative poly-A site (Bellone et al, unpublished), and the insertion of LTRs in introns has been associated with premature termination of transcription [63].

In addition to regulating gene expression, TEs may also be associated with specific functional characteristics of expressed protein coding genes. When I examined the functional annotation of repeat containing genes, I found that some functions were over-represented (Figure 7). This might occur because exaptation of TEs into coding sequences is associated with particular functions or that particular functional classes of genes are co-regulated and that repeats are involved in this co-regulation.

Non-coding genes also included exapted TEs or functionally important TE subsequences. Specifically I found that lincRNAs, which are known to regulate gene expression through epigenetic mechanisms and competition for transcription factors [64-66], contained TE sequences. Not all TEs contributed sequences equally, most of the Alu consensus was found in lincRNAs, but for L1s, 3' truncated sequences were dominant and for MIRs, the central 150bp was most often found in lincRNAs. This was in contrast to LTR, SVA and ERV TEs, which contributed much more defined, shorter sequences that were shown to contain transcription factor binding motifs. This result is consistent with a role for lincRNAs as competitors for transcription factors that bind to promoter regions.

In conclusion I can say that TEs are most strongly associated with repression of gene expression, either through the proximal promoter or 3'UTR or by contributing sequences to lincRNAs. These results extend previous published work in a novel and more comprehensive

fashion. Finally, I speculate that TEs may influence gene expression by co-regulating functionally similar genes, and I believe this to be a novel hypothesis.

## Future Directions

I have analyzed the repeats distribution, association with different gene sections and functional elements as well as expression of various classes of repeats in the human genome. However, my research target was limited to Homo sapiens. Thus, applying these methods and pipelines to other species is my next work, which may include chimpanzee, cow and elephant. Furthermore, in my analysis of transposable elements influences on human gene expression, I only focused on Alu, L1 and LTR in six human tissues. In next step, I may explore other classes of TEs impact on human gene transcript level, and also study the association between different types of TEs and tissues, to find out whether the influence of TEs in genes has any tissue specific. In this paper, I have studied the association between repeats and various regulatory elements, which were only based on active chromatin states. Therefore, I may transfer to the association between repeats and heterochromatin in our next work, giving us better understanding of the association between TEs and epigenetic regulation.

According to recent studies, piRNAs from RNA-protein complexes through interactions with piwi proteins, and these piRNA complexes have been linked to both epigenetic and post-transcriptional gene silencing of retrotransposons and other genetic elements in germ line [67]. This finding can help me to further explore the mechanism about how retrotransposons interact with piwi interacting RNA to impact the gene expression.

# Abbreviations List:

<b>TEs:</b>	transposable elements
<b>ORFs:</b>	open reading frames
<b>LTR:</b>	long terminal repeats
<b>LINES:</b>	long interspersed elements
<b>SINEs:</b>	short interspersed elements
<b>ADAR:</b>	double-stranded RNA-specific adenosine deaminase
<b>GO:</b>	Gene ontology
<b>lincRNAs:</b>	long intergenic non-coding RNA
<b>CSTE:</b>	chromatin states transposable elements

# Figures and table legends:

## Figure 3. Gene feature distributions of repetitive elements in human and cow

The y-axis has been normalized by the equation

$$percentage = \frac{repeats\ in\ specific\ genome\ section\ (bp)}{genome\ specific\ section\ (bp)} \times 100$$

## Figure 4. Gene feature distributions of biological activity of repetitive elements

The biological activity of repetitive elements was measured as the percentage of bases that overlapped repetitive elements and genomic areas of a particular chromatin state. Areas of genome that have been named biologically active are those that do not overlap heterochromatic areas. Shown are the specific biological activities of repetitive elements, averaged over six different cell lines (A), and the total biological activity of repetitive elements for each cell line (B). Percentages in A and B were normalized by the proportions of

genome occupying a particular gene feature. The method of normalization is same with the figure 3.

### **Figure 5. Biological activity of different repeat classes**

The biological activity of a repeat class was measured as the percentage of bases that overlapped with particular repeat class and genomic areas of a particular chromatin state. Areas of genome that have been named biologically active are those that do not overlap heterochromatic areas. Shown are the specific activities of repeat classes, averaged over six different cell lines (A), and the total activity of repeat classes for each cell line (B). Percentages in A and B have been normalized by the proportions of genome occupied by a particular repeat class. Normalization equation:

$$Percentage = \frac{Specific\ chromatin\ state\ TEs\ (bp)}{Specific\ TEs\ (bp)} \times 100$$

### **Figure 6. Gene feature distributions of repeat classes**

Percentages of genome that overlap a certain repeat class and gene feature. Percentages are normalized by the proportion of genome occupied by a particular gene feature. The method of normalization is same with the figure 3.

### **Figure 7. Enrichment of GO terms of various repeat classes in the human genome**

Enrichment is based on a log (p-value) transformation, which is shown over three GO domains, biological process (A), cellular compartment (B), and molecular function (C).

### **Figure 8. Expression level of three repeats in kidney (A), liver (B), brain (C), testes (D), skeletal muscle (E) and adipose (F) in the part of proximal promoter of human genes. I**

have used log transformation to acquire these violin plots. The white dot in the middle



represents the medium value of each sub graphs; the width of each subplot shows the gene number that enriched with respect to the y-axis.

**Figure 9. Expression level of three repeats in kidney (A), liver (B), brain (C), testes (D), skeletal muscle (E) and adipose (F) in 5'UTR region of human genes.** I have used log transformation to obtain these violin plots. The white dot in the middle represents the medium value of each sub graphs; the width of each subplot shows the gene number that enriched with respect to the y-axis.

**Figure 10. Expression level of three repeats in kidney (A), liver (B), brain (C), testes (D), skeletal muscle (E) and adipose (F) in 3'UTR section of human genes.** I have used log transformation to get these violin plots. The white dot in the middle represents the medium value of each sub graphs; the width of each subplot shows the gene number that enriched with respect to the y-axis.

**Figure 11. Coverage plot showing the length distribution of aligned sequence between lincRNA and different types of repetitive elements.** The plot shows Alu, Mir, LTR, SVA, L1 and ERV these six repeat consensus sequences association with lincRNAs. X-axis shows the whole length of each relevant repeat consensus sequence, y-axis signifies the number of the alignment length that covered in that region.

Figure 3:

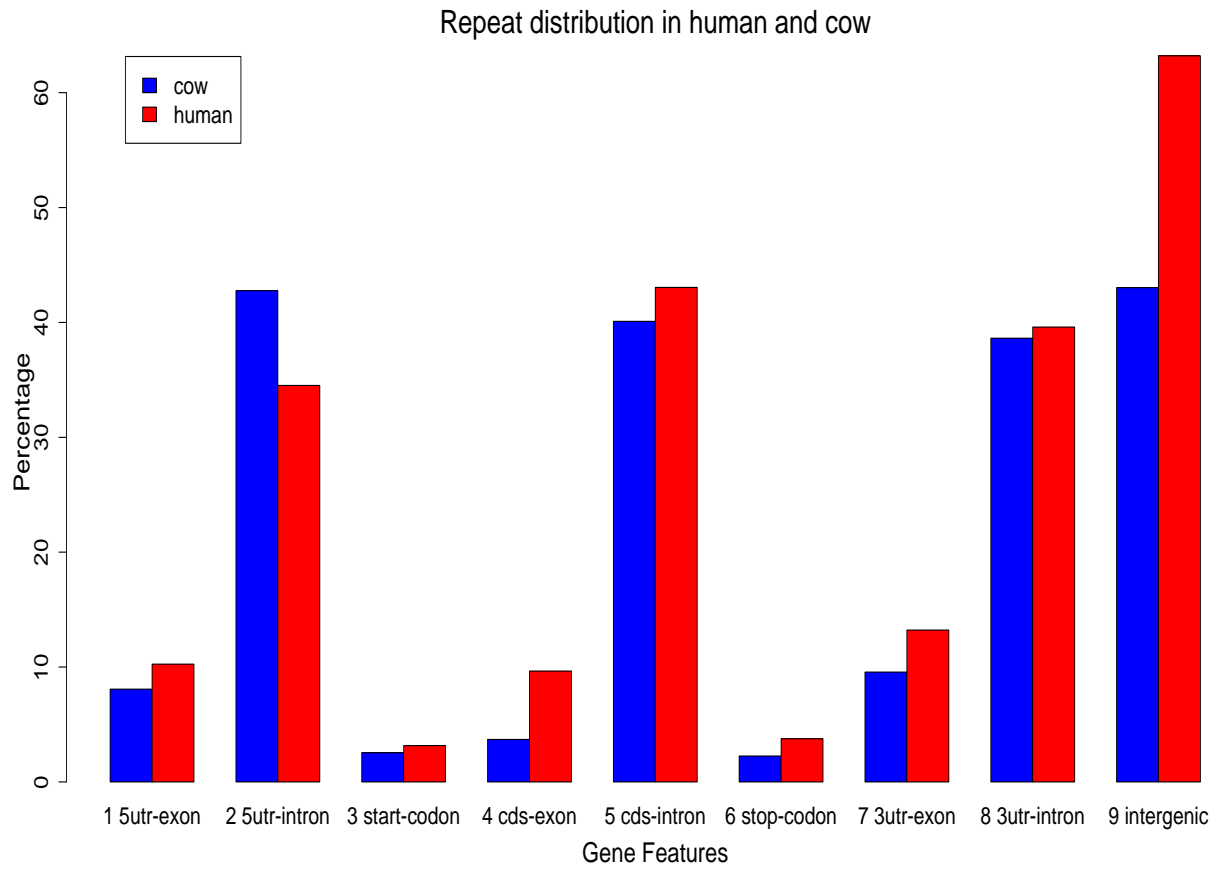


Figure 4:

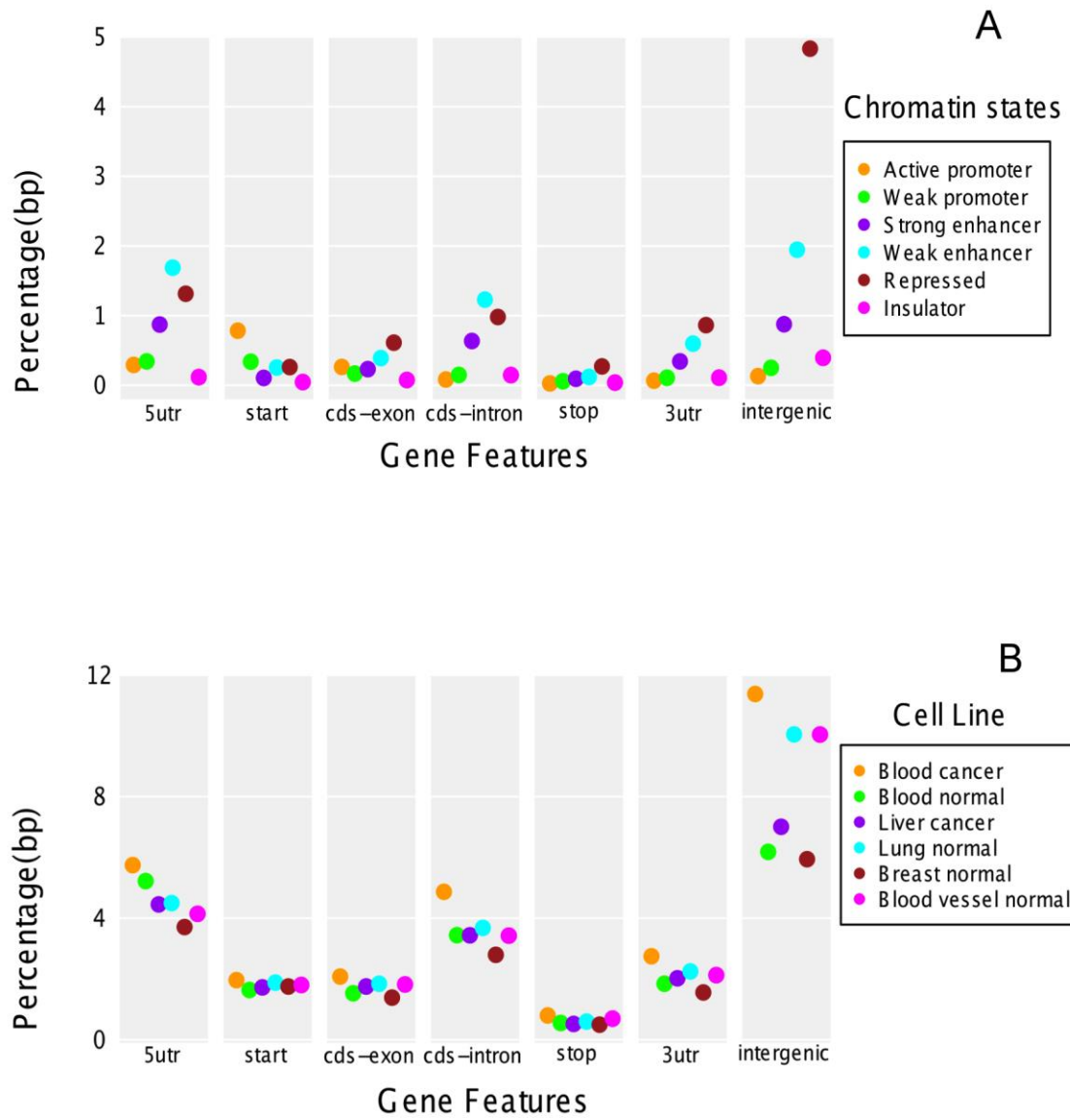


Figure 5:

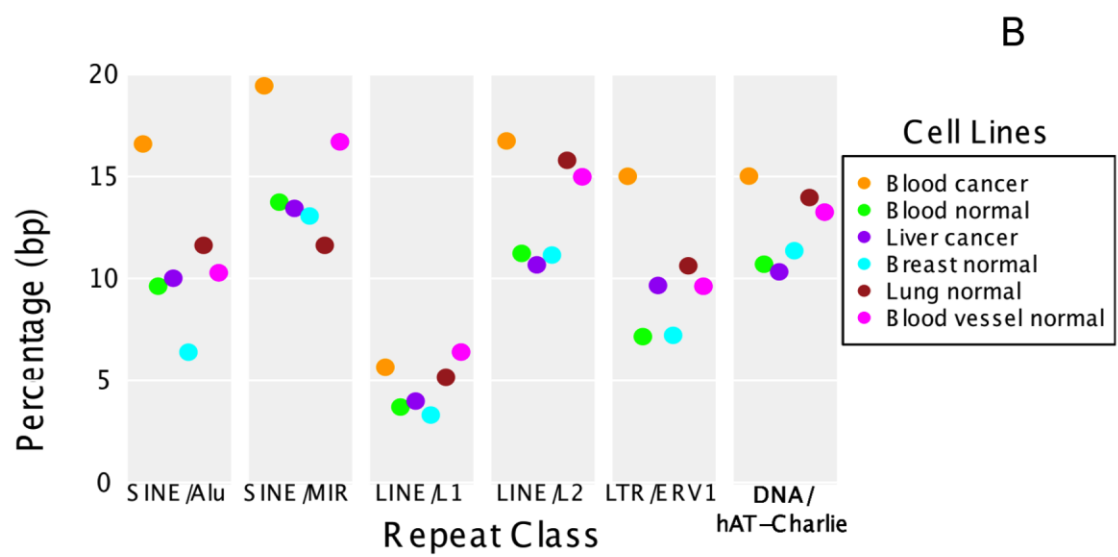
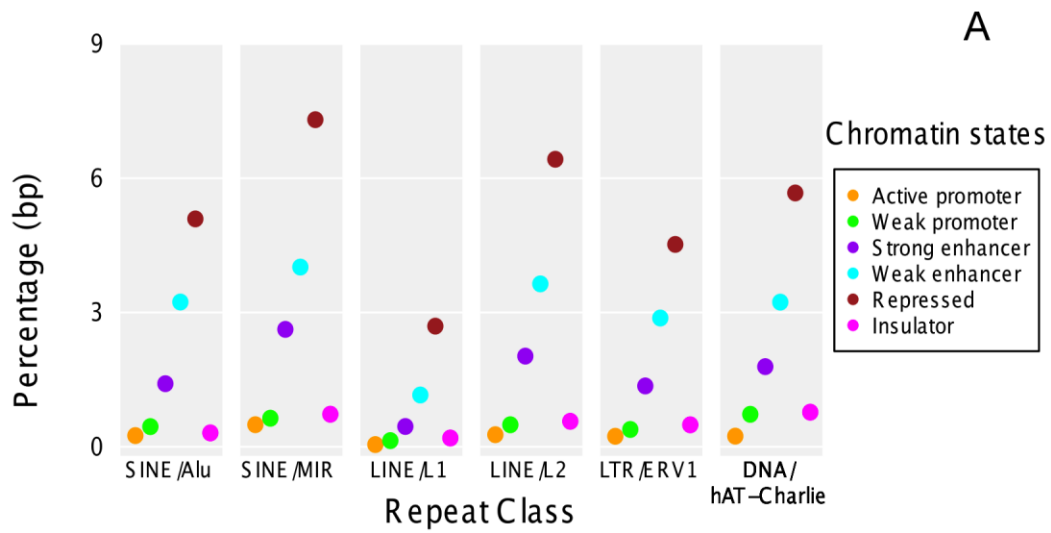


Figure 6:

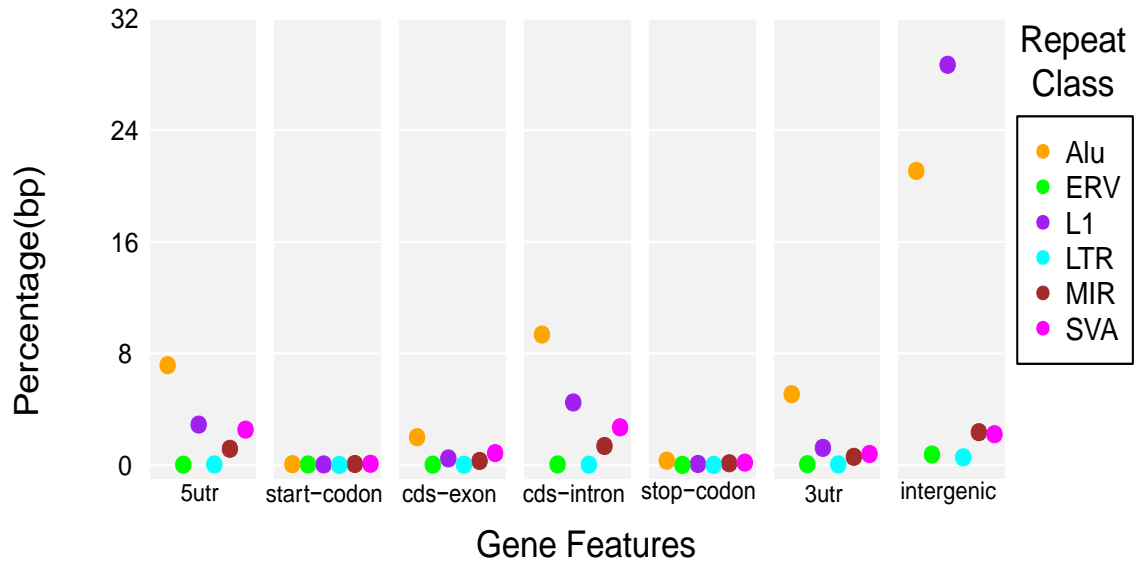
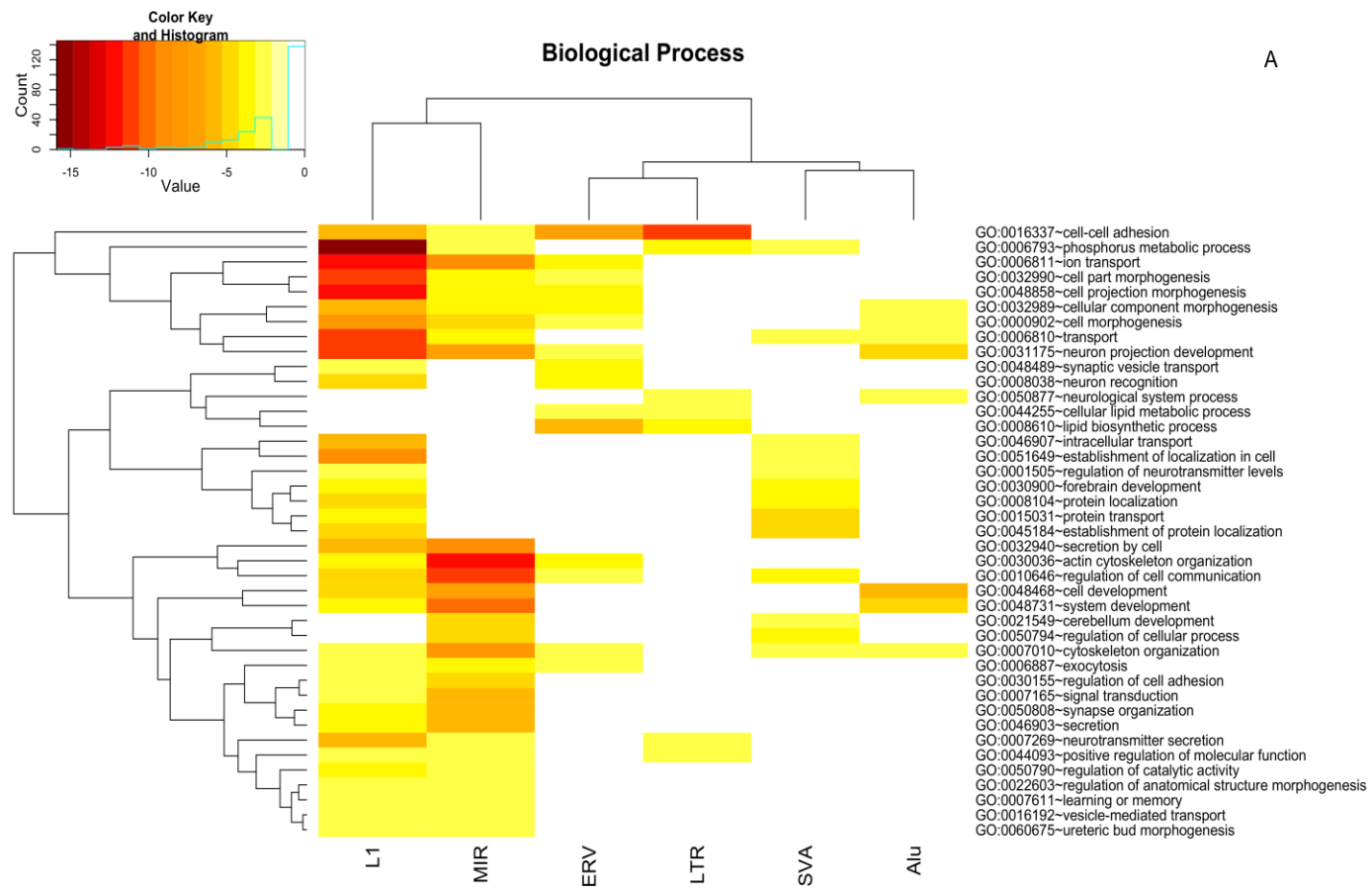
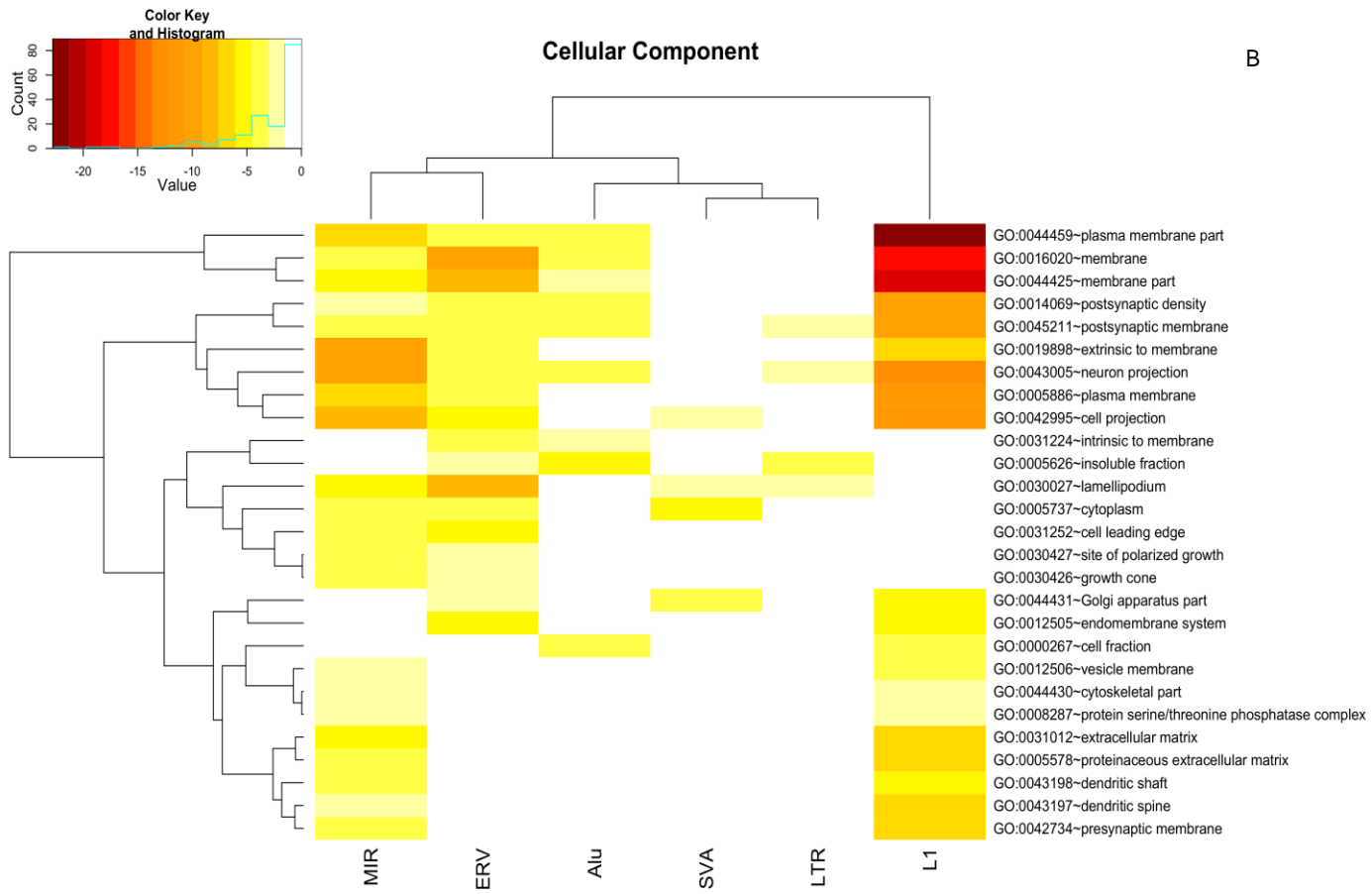
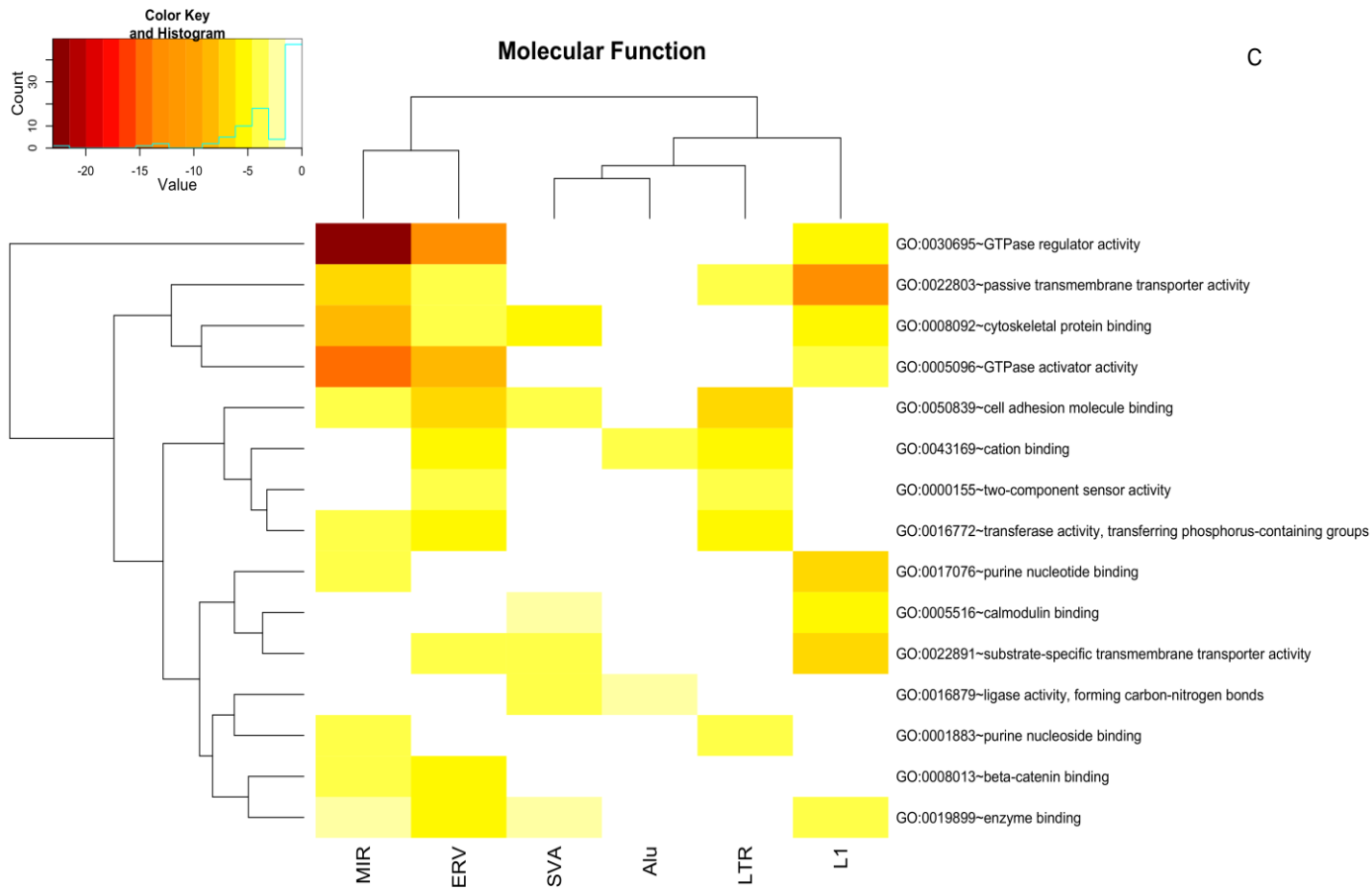


Figure 7:





B



C



Figure 8:

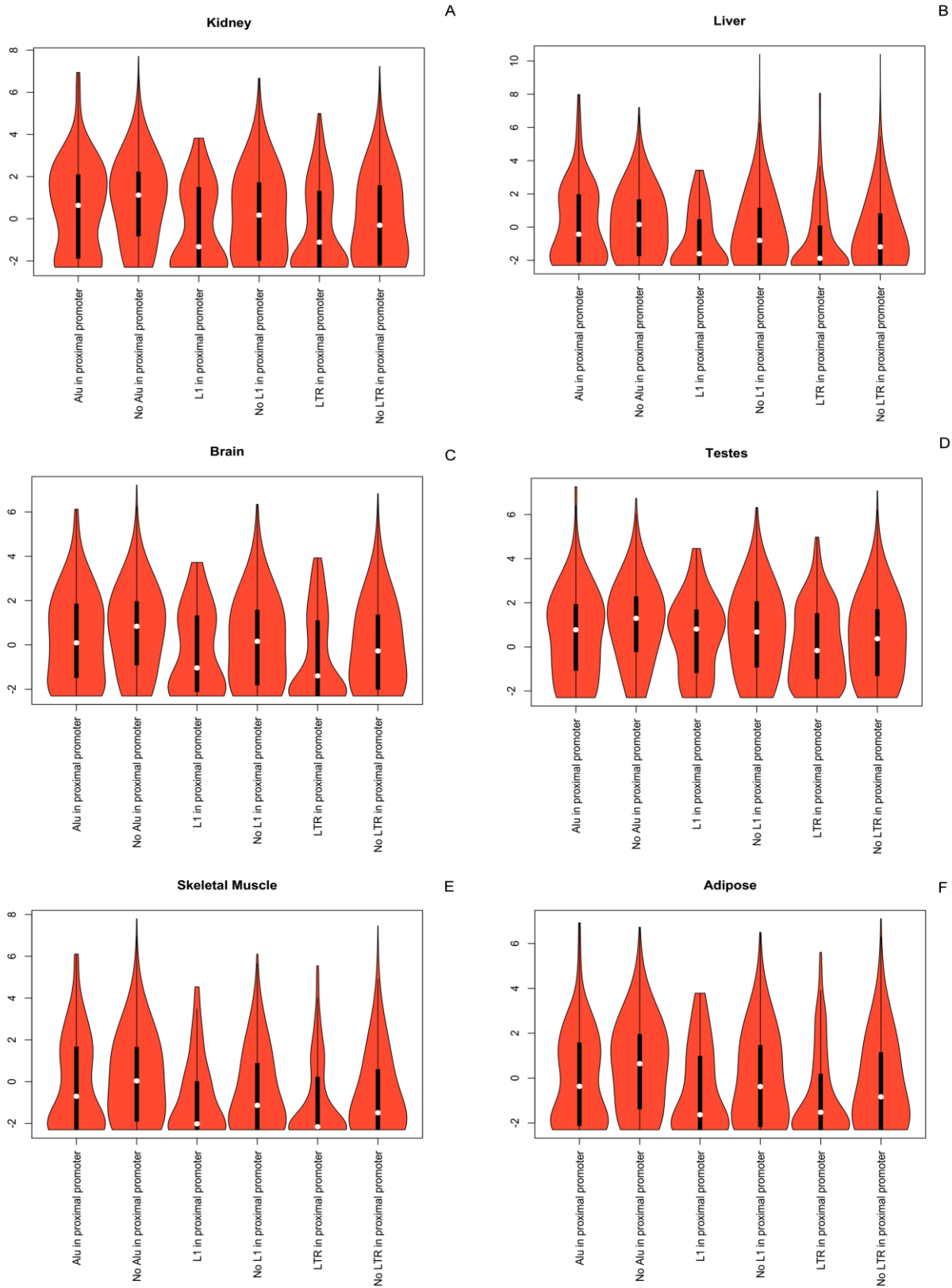


Figure 9:

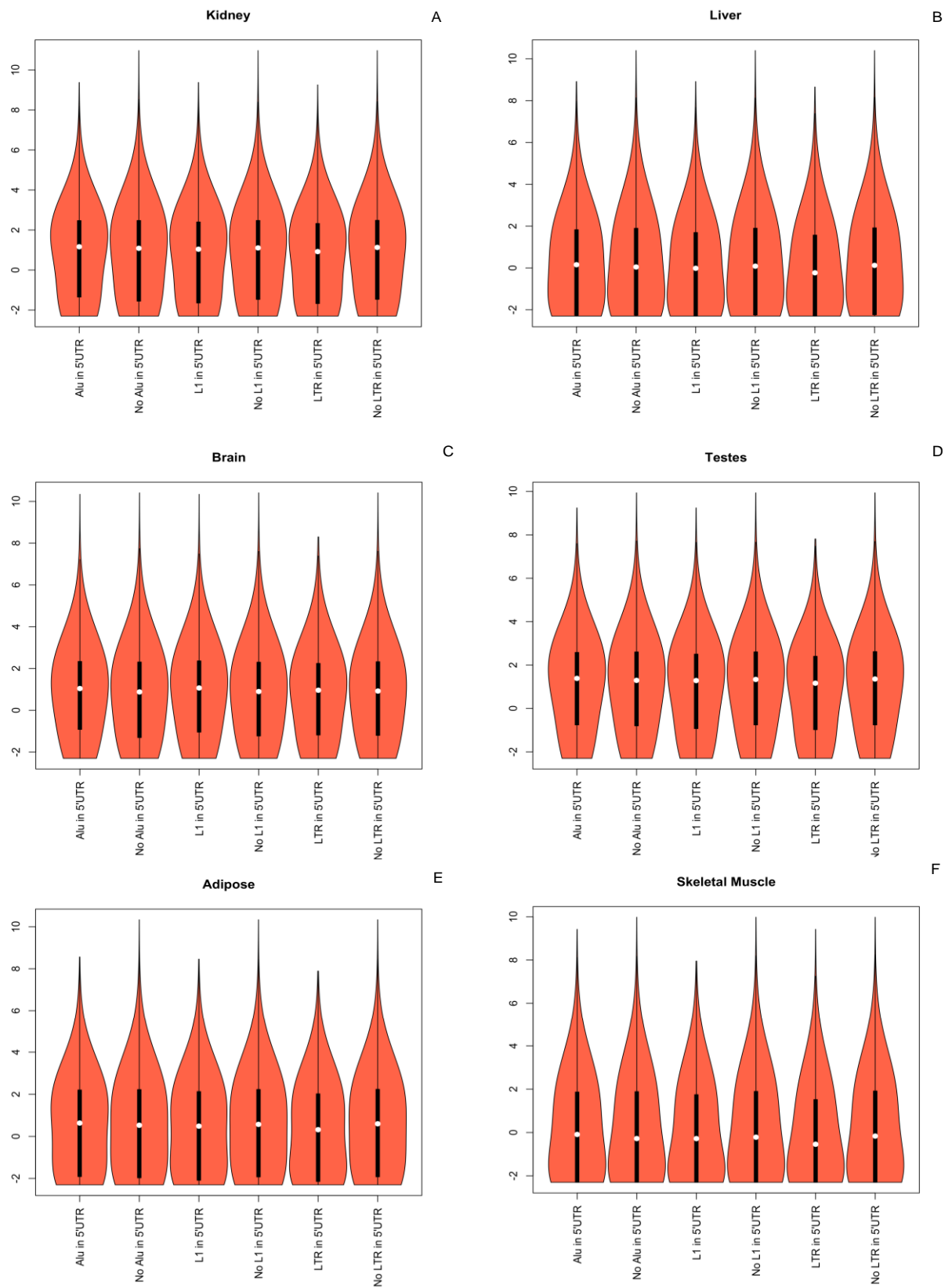


Figure 10:

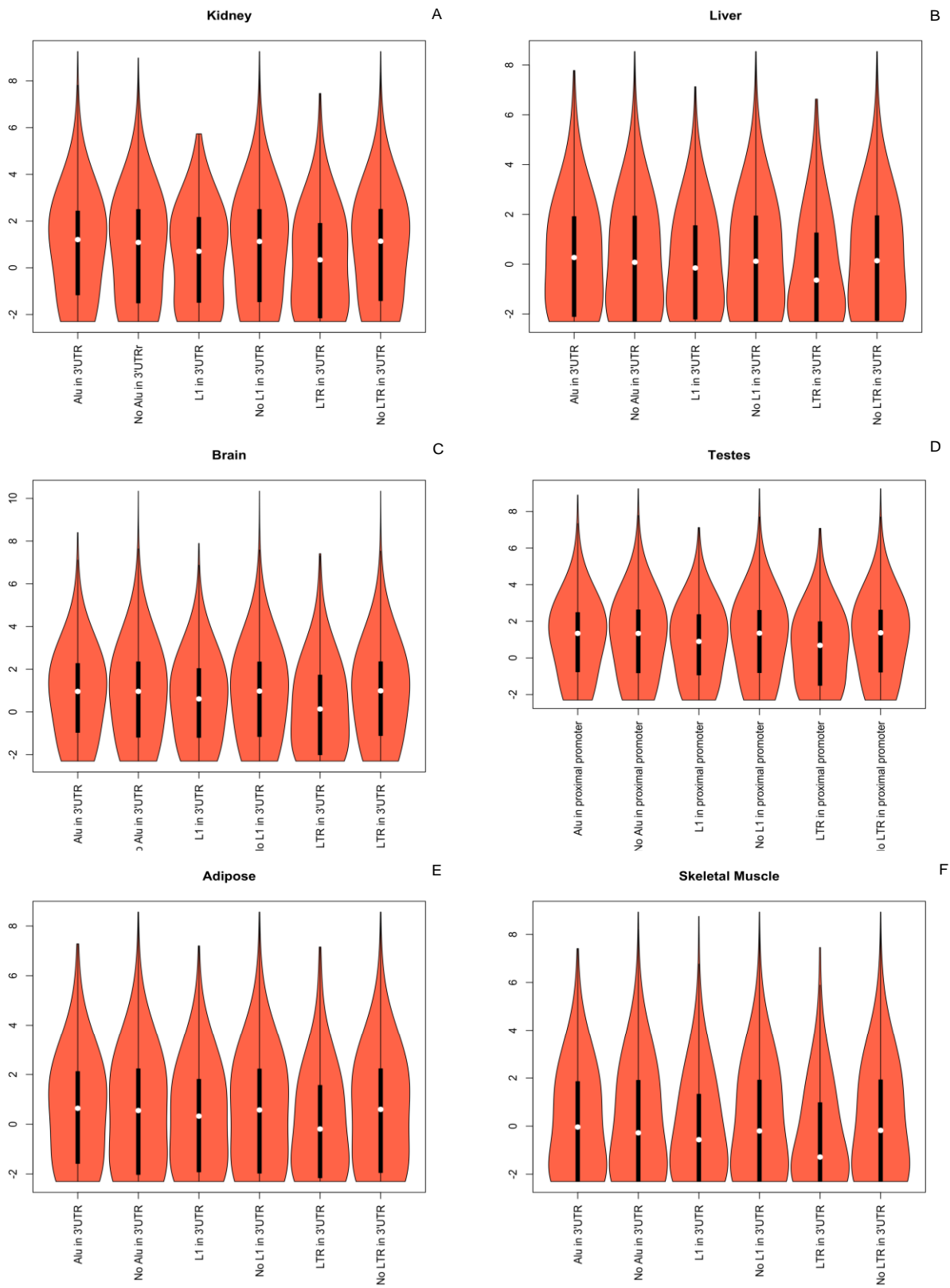
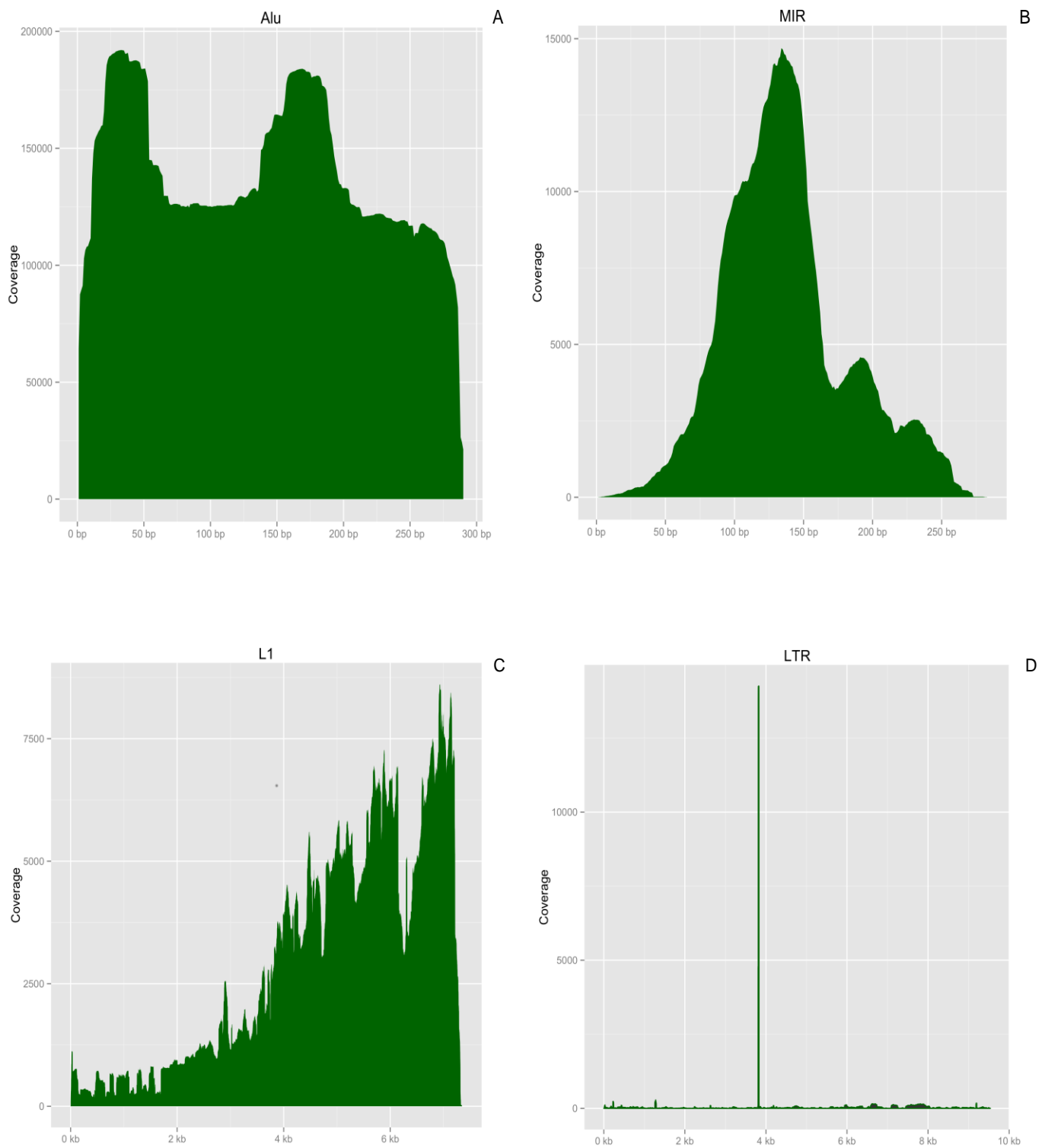


Figure 11:



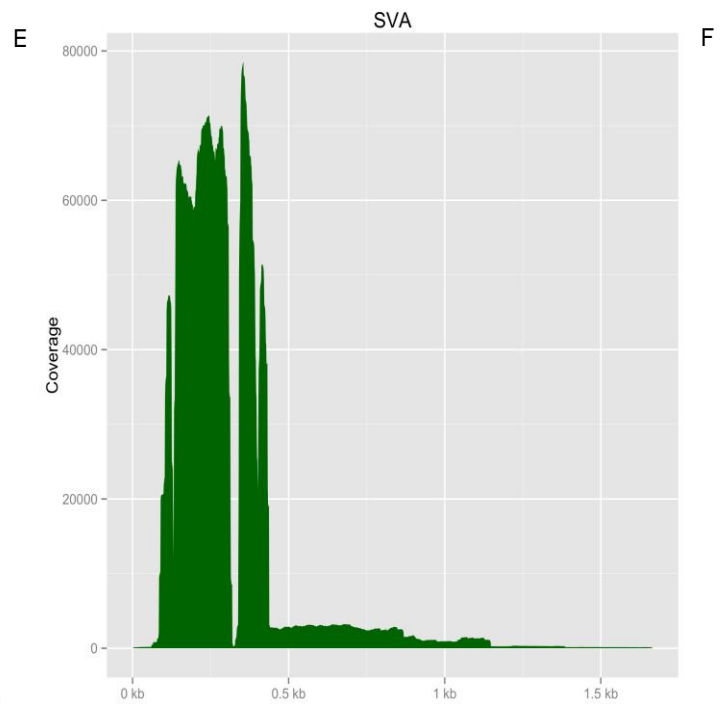
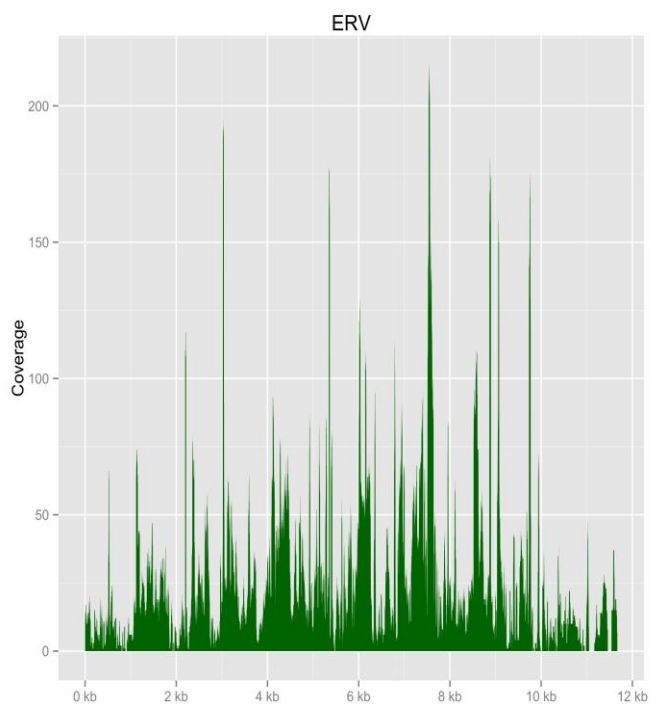


Table 1. GO terms with high enrichment in the human genome within six different classes of TEs

Category	Term	P-value	Repeats
Biological Process	GO:0048468~cell development	0.002433068	Alu
Biological Process	GO:0031175~neuron projection development	0.010950207	Alu
Biological Process	GO:0048731~system development	0.014218732	Alu
Biological Process	GO:0000902~cell morphogenesis	0.045430207	Alu
Biological Process	GO:0016337~cell-cell adhesion	6.76E-04	ERV
Biological Process	GO:0008610~lipid biosynthetic process	0.004836686	ERV
Biological Process	GO:0007565~female pregnancy	0.007958811	ERV
Biological Process	GO:0010942~positive regulation of cell death	0.021215202	ERV
Biological Process	GO:0030036~actin cytoskeleton organization	0.022252018	ERV
Biological Process	GO:0006629~lipid metabolic process	0.029784268	ERV
Biological Process	GO:0006811~ion transport	0.030293183	ERV
Biological Process	GO:0048489~synaptic vesicle transport	0.033240211	ERV
Biological Process	GO:0048858~cell projection morphogenesis	0.041187851	ERV
Biological Process	GO:0032989~cellular component morphogenesis	0.04121194	ERV
Biological Process	GO:0008038~neuron recognition	0.0416136	ERV
Biological Process	GO:0000902~cell morphogenesis	0.0448662	ERV
Biological Process	GO:0006793~phosphorus metabolic process	1.26E-07	L1
Biological Process	GO:0006811~ion transport	5.44E-06	L1
Biological Process	GO:0048858~cell projection morphogenesis	6.02E-06	L1
Biological Process	GO:0032990~cell part morphogenesis	1.54E-05	L1
Biological Process	GO:0031175~neuron projection development	1.54E-05	L1
Biological Process	GO:0006810~transport	1.83E-05	L1
Biological Process	GO:0051649~establishment of localization in cell	1.13E-04	L1

Biological Process	GO:0000902~cell morphogenesis	2.82E-04	L1
Biological Process	GO:0007411~axon guidance	7.57E-04	L1
Biological Process	GO:0046907~intracellular transport	0.001803868	L1
Biological Process	GO:0032835~glomerulus development	0.002073621	L1
Biological Process	GO:0032989~cellular component morphogenesis	0.002253662	L1
Biological Process	GO:0007267~cell-cell signaling	0.002650669	L1
Biological Process	GO:0016337~cell-cell adhesion	0.003030857	L1
Biological Process	GO:0055085~transmembrane transport	0.00371976	L1
Biological Process	GO:0007269~neurotransmitter secretion	0.003912961	L1
Biological Process	GO:0032940~secretion by cell	0.003993018	L1
Biological Process	GO:0048468~cell development	0.007052635	L1
Biological Process	GO:0008038~neuron recognition	0.007550272	L1
Biological Process	GO:0008104~protein localization	0.00963923	L1
Biological Process	GO:0045184~establishment of protein localization	0.010047077	L1
Biological Process	GO:0030031~cell projection assembly	0.011954353	L1
Biological Process	GO:0010646~regulation of cell communication	0.013264599	L1
Biological Process	GO:0048731~system development	0.015254729	L1
Biological Process	GO:0046903~secretion	0.016725195	L1
Molecular Function	GO:0043169~cation binding	0.030378079	Alu
Molecular Function	GO:0005102~receptor binding	0.035545206	Alu
Molecular Function	GO:0030695~GTPase regulator activity	3.76E-06	ERV
Molecular Function	GO:0005096~GTPase activator activity	4.18E-04	ERV
Molecular Function	GO:0050839~cell adhesion molecule binding	9.41E-04	ERV
Molecular Function	GO:0019899~enzyme binding	0.002702865	ERV
Molecular Function	GO:0043169~cation binding	0.005550116	ERV
Molecular Function	GO:0008013~beta-catenin binding	0.00576055	ERV
Molecular Function	GO:0016772~transferase activity, transferring phosphorus-containing	0.007591362	ERV

	groups		
Molecular Function	GO:0000155~two-component sensor activity	0.012793986	ERV
Molecular Function	GO:0008092~cytoskeletal protein binding	0.015298604	ERV
Molecular Function	GO:0022803~passive transmembrane transporter activity	0.021237808	ERV
Molecular Function	GO:0008656~caspase activator activity	0.022962103	ERV
Molecular Function	GO:0022891~substrate-specific transmembrane transporter activity	0.035814672	ERV
Molecular Function	GO:0022803~passive transmembrane transporter activity	2.97E-06	L1
Molecular Function	GO:0043169~cation binding	2.99E-06	L1
Molecular Function	GO:0016772~transferase activity, transferring phosphorus-containing groups	1.20E-05	L1
Molecular Function	GO:0001883~purine nucleoside binding	1.37E-04	L1
Molecular Function	GO:0032553~ribonucleotide binding	7.33E-04	L1
Molecular Function	GO:0017076~purine nucleotide binding	0.001280726	L1
Molecular Function	GO:0022891~substrate-specific transmembrane transporter activity	0.002093544	L1
Molecular Function	GO:0005516~calmodulin binding	0.003082402	L1
Molecular Function	GO:0030695~GTPase regulator activity	0.005323101	L1
Molecular Function	GO:0008092~cytoskeletal protein binding	0.007427365	L1
Molecular Function	GO:0005096~GTPase activator activity	0.011907	L1
Molecular Function	GO:0019899~enzyme binding	0.013789565	L1
Molecular Function	GO:0005543~phospholipid binding	0.016279339	L1
Molecular Function	GO:0019887~protein kinase regulator activity	0.030812232	L1
Molecular Function	GO:0050839~cell adhesion molecule binding	0.037913538	L1
Molecular Function	GO:0008093~cytoskeletal adaptor activity	0.044567764	L1
Molecular Function	GO:0050839~cell adhesion molecule binding	6.11E-04	LTR
Molecular Function	GO:0016772~transferase activity, transferring phosphorus-containing groups	0.00807089	LTR
Molecular Function	GO:0043169~cation binding	0.009748826	LTR
Molecular Function	GO:0005496~steroid binding	0.010019163	LTR



Molecular Function	GO:0000155~two-component sensor activity	0.012243649	LTR
Molecular Function	GO:0022803~passive transmembrane transporter activity	0.020931795	LTR
Molecular Function	GO:0001883~purine nucleoside binding	0.046141345	LTR
Molecular Function	GO:0030695~GTPase regulator activity	9.81E-11	MIR
Molecular Function	GO:0005096~GTPase activator activity	5.53E-07	MIR
Molecular Function	GO:0008092~cytoskeletal protein binding	3.32E-04	MIR
Molecular Function	GO:0022803~passive transmembrane transporter activity	6.24E-04	MIR
Molecular Function	GO:0043169~cation binding	0.00144037	MIR
Molecular Function	GO:0022891~substrate-specific transmembrane transporter activity	0.008004068	MIR
Molecular Function	GO:0005516~calmodulin binding	0.008041464	MIR
Molecular Function	GO:0046983~protein dimerization activity	0.015808567	MIR
Molecular Function	GO:0001883~purine nucleoside binding	0.023661931	MIR
Molecular Function	GO:0019992~diacylglycerol binding	0.026389336	MIR
Molecular Function	GO:0032553~ribonucleotide binding	0.028150117	MIR
Molecular Function	GO:0016772~transferase activity, transferring phosphorus-containing groups	0.035853993	MIR
Molecular Function	GO:0008013~beta-catenin binding	0.037247081	MIR
Molecular Function	GO:0050839~cell adhesion molecule binding	0.037247081	MIR
Molecular Function	GO:0017076~purine nucleotide binding	0.042483365	MIR
Molecular Function	GO:0008092~cytoskeletal protein binding	0.002789855	SVA
Molecular Function	GO:0050839~cell adhesion molecule binding	0.011650127	SVA
Molecular Function	GO:0030695~GTPase regulator activity	0.012213165	SVA
Molecular Function	GO:0005543~phospholipid binding	0.022869129	SVA
Molecular Function	GO:0016879~ligase activity, forming carbon-nitrogen bonds	0.031038513	SVA
Cellular Component	GO:0005626~insoluble fraction	0.002754261	Alu
Cellular Component	GO:0014069~postsynaptic density	0.01455023	Alu
Cellular Component	GO:0045211~postsynaptic membrane	0.016366668	Alu

Cellular Component	GO:0000267~cell fraction	0.023283704	Alu
Cellular Component	GO:0043005~neuron projection	0.037089783	Alu
Cellular Component	GO:0016020~membrane	0.038961333	Alu
Cellular Component	GO:0044459~plasma membrane part	0.039979154	Alu
Cellular Component	GO:0016020~membrane	1.03E-04	ERV
Cellular Component	GO:0030027~lamellipodium	2.77E-04	ERV
Cellular Component	GO:0044425~membrane part	4.90E-04	ERV
Cellular Component	GO:0012505~endomembrane system	0.003112417	ERV
Cellular Component	GO:0031252~cell leading edge	0.005978395	ERV
Cellular Component	GO:0042995~cell projection	0.008750212	ERV
Cellular Component	GO:0043005~neuron projection	0.010567625	ERV
Cellular Component	GO:0005886~plasma membrane	0.015943166	ERV
Cellular Component	GO:0014069~postsynaptic density	0.017268836	ERV
Cellular Component	GO:0031224~intrinsic to membrane	0.018157998	ERV
Cellular Component	GO:0044459~plasma membrane part	0.019160607	ERV
Cellular Component	GO:0031965~nuclear membrane	0.019851813	ERV
Cellular Component	GO:0031300~intrinsic to organelle membrane	0.022330112	ERV
Cellular Component	GO:0045211~postsynaptic membrane	0.031613444	ERV
Cellular Component	GO:0019898~extrinsic to membrane	0.034741239	ERV
Cellular Component	GO:0005737~cytoplasm	0.039586133	ERV
Cellular Component	GO:0044459~plasma membrane part	1.28E-10	L1
Cellular Component	GO:0044425~membrane part	6.73E-09	L1
Cellular Component	GO:0016020~membrane	3.08E-08	L1
Cellular Component	GO:0043005~neuron projection	1.73E-06	L1
Cellular Component	GO:0005886~plasma membrane	6.19E-06	L1
Cellular Component	GO:0042995~cell projection	1.89E-05	L1
Cellular Component	GO:0045211~postsynaptic membrane	5.98E-05	L1

Cellular Component	GO:0014069~postsynaptic density	7.58E-05	L1
Cellular Component	GO:0044463~cell projection part	8.04E-05	L1
Cellular Component	GO:0005578~proteinaceous extracellular matrix	0.00100319	L1
Cellular Component	GO:0019898~extrinsic to membrane	0.001395374	L1
Cellular Component	GO:0042734~presynaptic membrane	0.001400923	L1
Cellular Component	GO:0043197~dendritic spine	0.001680735	L1
Cellular Component	GO:0031012~extracellular matrix	0.001904228	L1
Cellular Component	GO:0005875~microtubule associated complex	0.003288719	L1
Cellular Component	GO:0044431~Golgi apparatus part	0.003519217	L1
Cellular Component	GO:0012505~endomembrane system	0.003851358	L1
Cellular Component	GO:0043198~dendritic shaft	0.008184011	L1
Cellular Component	GO:0000267~cell fraction	0.0254723	L1
Cellular Component	GO:0009986~cell surface	0.032262217	L1
Cellular Component	GO:0012506~vesicle membrane	0.040162796	L1
Cellular Component	GO:0044441~cilium part	0.04790736	L1
Cellular Component	GO:0044420~extracellular matrix part	0.050716684	L1
Cellular Component	GO:0009897~external side of plasma membrane	0.052729331	L1
Cellular Component	GO:0005581~collagen	0.057803073	L1
Cellular Component	GO:0035085~cilium axoneme	0.061178965	L1
Cellular Component	GO:0044430~cytoskeletal part	0.086692602	L1
Cellular Component	GO:0008287~protein serine/threonine phosphatase complex	0.091900824	L1
Cellular Component	GO:0005604~basement membrane	0.09873182	L1
Cellular Component	GO:0005626~insoluble fraction	0.047082603	LTR
Cellular Component	GO:0030027~lamellipodium	0.063143676	LTR
Cellular Component	GO:0043005~neuron projection	0.079788863	LTR
Cellular Component	GO:0005929~cilium	0.088069784	LTR
Cellular Component	GO:0045211~postsynaptic membrane	0.099886093	LTR

Cellular Component	GO:0019898~extrinsic to membrane	6.70E-05	MIR
Cellular Component	GO:0043005~neuron projection	8.70E-05	MIR
Cellular Component	GO:0042995~cell projection	1.68E-04	MIR
Cellular Component	GO:0044459~plasma membrane part	6.25E-04	MIR
Cellular Component	GO:0005886~plasma membrane	0.001138769	MIR
Cellular Component	GO:0034702~ion channel complex	0.001420407	MIR
Cellular Component	GO:0030027~lamellipodium	0.00671882	MIR
Cellular Component	GO:0044425~membrane part	0.009441945	MIR
Cellular Component	GO:0031012~extracellular matrix	0.00961077	MIR
Cellular Component	GO:0005578~proteinaceous extracellular matrix	0.010895064	MIR
Cellular Component	GO:0016020~membrane	0.015034447	MIR
Cellular Component	GO:0031252~cell leading edge	0.018580787	MIR
Cellular Component	GO:0016459~myosin complex	0.020015626	MIR
Cellular Component	GO:0005737~cytoplasm	0.02689117	MIR
Cellular Component	GO:0042734~presynaptic membrane	0.031856375	MIR
Cellular Component	GO:0030426~growth cone	0.040758303	MIR
Cellular Component	GO:0030427~site of polarized growth	0.043152743	MIR
Cellular Component	GO:0043198~dendritic shaft	0.043472116	MIR
Cellular Component	GO:0045211~postsynaptic membrane	0.044889397	MIR
Cellular Component	GO:0002142~stereocilia ankle link complex	0.047579822	MIR
Cellular Component	GO:0002139~stereocilia coupling link	0.047579822	MIR
Cellular Component	GO:0002141~stereocilia ankle link	0.047579822	MIR
Cellular Component	GO:0032420~stereocilium	0.049378586	MIR
Cellular Component	GO:0005622~intracellular	2.59E-04	SVA
Cellular Component	GO:0044424~intracellular part	0.001178181	SVA
Cellular Component	GO:0043229~intracellular organelle	0.001696152	SVA
Cellular Component	GO:0043231~intracellular membrane-bounded organelle	0.003537599	SVA

Cellular Component	GO:0005737~cytoplasm	0.009380962	SVA
Cellular Component	GO:0000776~kinetochore	0.033366997	SVA
Cellular Component	GO:0044431~Golgi apparatus part	0.043688535	SVA
Cellular Component	GO:0044444~cytoplasmic part	0.044648231	SVA

Table 2. Cell lines source and kinds that used to analyze the function of retrotransposons in the human genome

Cell	Tier	Description	Lineage	Tissue	Karyotype	Sex	Documents	Vendor ID	Term ID	Label
GM12878	1	B-lymphocyte,lymphoblastoid,international HapMap Project-CEPH/Utah-European Caucasian,Epstein-Barr Virus	mesoderm	blood	normal	F	ENCODE	Coriell GM12878	BTO: 0002062	GM 12878
HepG2	2	Hepatocellular carcinoma	endoderm	liver	cancer	M	ENCODE	ATCC HB-8065	BTO: 0000599	HepG2
HMEC	3	Mammary epithelial cells	ectoderm	breast	normal	U	Bernstein Crawford stam	Lonza CC-2551	BTO: 0002178	HMEC
HUVEC	2	Umbilical vein endothelial cells	mesoderm	blood vessel	normal	U	Encode	Lonza CC-2517	BTO: 0001949	HUVEC
NHLF	3	Lung fibroblasts	endoderm	lung	normal	U	Bernstein stam	Lonza CC-2512	BTO: 0000161	NHLF
K562	1	leukemia, "the continuous cell line K-562 was established by Lozzio and Lozzio from the pleural effusion of a 53-year-old female with chronic myelogenous leukemia in terminal blast crises." - ATCC	mesoderm	blood	cancer	F	ENCODE	ATCC CCL-243	BTO: 0000664	K562









## Supplementary Material S2:

### Alu consensus sequences:

#### >1\_cons

GGCCGGGCGCGGTGGCTCACGCCTGTAATCCCAGCACTTTGGGAGGCCGAGGCGGGCGGA  
TCACCTGAGGTCAGGAGTTTCGAGACCAGCCTGGCCAACATGGTGAAACCCCGTCTCTACT  
AAAAATACAAAAAATTAGCCGGGCGTGGTGGCGCGCGCCTGTAATCCCAGCTACTCGGGA  
GGCTGAGGCAGGAGAATCGCTTGAACCCGGGAGGCGGAGGTTGCAGTGAGCCGAGATCGC  
GCCACTGCACTCCAGCCTGGGCGACAAGAGCGAGACTCCGTCTCAA

#### >2\_cons

GGCCGGGCGCGGTGGCTCACGCCTGTAATCCCAGCACTTTGGGAGGCCGAGGCGGGCGGA  
TCACGAGGTCAGGAGATCGAGACCATCCTGGCTAACACGGTGAAACCCCGTCTCTACTAA  
AAATACAAAAAATTAGCCGGGCGTGGTGGCGGGCGCCTGTAGTCCCAGCTACTCGGGAG  
GCTGAGGCAGGAGAATGGCGTGAACCCGGGAGGCGGAGCTTGCAGTGAGCCGAGATTGCG  
CCACTGCAGTCCGAGTCCAGCCTGGGCGACAGAGCGAGACTCCGTCTCA

#### >3\_cons

GGCCGGGCGCGGTGGCTCACGCCTGTAATCCCAGCACTTTGGGAGGCCGAGGCGGGCGGA  
TCACGAGGTCAGGAGATCGAGACCATCCTGGCTAACACGGTGAAACCCCGTCTCTACTAA  
AAAATACAAAAAATTAGCCGGGCGTGGTGGCGGGCGCCTGTAGTCCCAGCTACTCGGAG  
AGGCTGAGGCAGGAGAATGGCGTGAACCCGGGAGGCGGAGCTTGCAGTGAGCCGAGATC  
GCGCCACTGCACTCCAGCCTGGGCGACAGAGCGAGACTCCGTCTCA

### MIR consensus sequences:

#### >1\_cons

CAGAGGGACAGCATAGCACAGTGGTAAAGAGCACGGACTCTGGAGCCAGACAGACCTGGG  
TTCGAATCCCGGCTCTGCCACTTACTAGCTGTGTGACCTTGGGCAAGTCACTTAACCTCT  
CTGAGCCTCAGTTTCCTCATCTGTAATAATGGGGATAATAATAGTACCTACCTCACAGGT  
TGTGTGAGGATTAATGAGATAATACATGTAAAGCGCTTAGAACAGTGCCTGGCACACA  
GTAAGCGCTCAATAAATGGTAGCTCTATTATT

#### >2\_cons

TTCTCGAAGCAGTATGGTACAGTGGAAAGAACAACCTGGACTAGGAGTCAGGAAGACCTGG  
GTTTCGAGTCCTAGCTCTGCCACTAAGCTGTGTGACCTTGGGCAAGTCACTTAACCTC  
TCTGAGCCTCAGTTTCCTCATCTGTAATAATGAGGATAATAATACCTGCCCTGCCTACCT  
CACAGGGTTGTTGTGAGGATCAAACGAGATAATCTATGTGAAAGCGCCCTGCAAACCTTA  
AAATGCTATACAAATGTAAGGGGATACTATGATTCTAAAAAAA

### L1 consensus sequences:

#### >1\_cons

GGGGGAGGAGCCAAGATGGCCGAATAGAAACAGCACCGGTCTACAGCTCCCAGCGTGAG  
CGACCACAGAAGACGGGTGATTCCTGCATCGCCAACCTGAGGTACCAGGTTTCATCTACTA  
GGAAGTGCCAGACAGCGGGCGCACCCACAGACCCTCTGAAGGAAGCGGACTGCTCCTGCA  
GGACCCGGGAGACACCCCAAATACTGTGAGTGCCCAAACCTGCGGAAGTGGGAAAGGGAGA  
TCCTCCGCTCCCGAACACACACCCCACTGGGGAAACTGAAGGTCTAGTTTTCGGGGAGAA  
GTTTCCGACCTTACCTGGAGCTGAGTCAATTTAGAGAGCCGAGCGAAATACAGGGGTAGA  
GGAAGCAGCGAGGAAAGGCCCTGGGAGCTCGCTGGGTCCCAAGCAGGCCATTCTGCCT  
GGCACCACAGGATCCTTCGGGAGGGCGGACAGAGGAGCGAGCGCACCGAGCGCAAGCCG  
AAGCAGGGCGAGGCAGCGCCTCACCTGAGAAGCGCAAGGGGTTCAGGGAATCCCTTTCCC  
AGTCAAAGAAAGGGGGGACGGACGCCACCTGCAAATCGGGTCACTCCCGCCCTAACAAAT

GCCTTTTCCGACCCACTTAAGAAACGGCGCACCACGAGAATATACCCACAGGGGCCTA  
GGTCCCAACCCTGGAGCCGCGCAGATTCTCAACAGCCTCTCAGCTGGAATCTGCTTAAG  
CCTGCCGAGCTCCTGGCTCGGAGGGTCTACGCACACGGACTCTCGCTGATTGCTAGCAC  
AGCAGTCTGAGATCAAAGTCAAGGGGCGAGCCAGCACTGGGACTCATAACTGCCTAA  
CACACTAAGCTCCCGCAACGAGGCTGGGCGAGGGGCGCCCGCAATGCCAGGCTTCCC  
TAGGTAAACAAAGCAGCCGGGAAGCTCGAACGGGGTGGACCCACCACAGCTCAAGCACA  
CCTCCATGCACCTGTAGGCTCCACCTCTGGCCGACGCGCACAGACAAAACAAAAACACAGC  
AGTAACCTCCGCAGACTTAAGTGTCCCTGTCCGACAGCTTCGAAAAGAGCAGTCGTTCTC  
CCAGCACGCAGCTGCAGATCTGACTGGGCTGAGCCCCTAGAGGGAGGGGTGGCCGCAGT  
CTCTGCGGACCAGCAGACTTAGCCTTCTCCTGGTAGTTCTGAGGAATCCGGGCAGCCC  
AGATGAGTGGGTTTCCCCCAGCGAAGCACACCCCTGAACGAGCAGACTGCCACCTCAA  
GTGGGTAAATGACCCCTGACCCCGAGCAGCCTAACTGGGAGGCACCCCGCAGAGGGGC  
GGACTGACACCCACACGGCAGCGATCCTACTGGCATCAGGTTGGTACACCTCGAAGACA  
AAAATACCAGAAGAACGACCAGGCAACAACTCTGCCGTTCTCCAATATCCACCACTGAC  
ACCACCAACGCGGAAGAGAACCAGAAAAACAGGGTCTGAAGTGAACCCCGCAGAAACTC  
CAACAGACCTGCAGCAGAGGGTCTGACTATTAGAAGGAAAATAACAAACAGAAAGGAC  
ATCCACACCCAAAACCCATATAAACATCACTGCAGCTCGGCTCACAGGAAGCCACATCCA  
TAGGAAAAGGGGGAGAGTACTACATCAAGGGAACACCCCGTGGGAAAAAAAAACCCAAAC  
AACAGCCAGCAGCATCAAAGACCAAACTAGATAAAACCACAAAGATGAGAAAAAACAG  
AAAAGAAAACTGGAACTCTAAAAAACAGAGCGCTCTCCTCCTCAAAGGAACGCAGC  
TCCTCACCAGCAACGGAACAAAACCTGGACGGAGAATGACTTTGACGAGTTGACAGAAGAA  
GGCTTCAGAAGATGAGTAATAAAAACTACTCTGAGCTACGGGAGGAAATTCAAACCAA  
GGCAAAGAAGTTAAAACTTTGAAAAAAATTAGAGGAATGGATAACTAAGGGAATAACC  
AATACAGAGAAGAACTTAAAGGACCTGATGGAGCTGAAAAACAAAGCACGAGAACTACGT  
GAAGAATGCAGAAGCCTCAGTAGCCGAAGCGATCAACTGGAAGAAAGGATATCAGAGATG  
GAAGATCAAATCAATGAAATAAAACAAGAAGAGAAGATTAGAGAAAAAGAATAAAAAAGA  
AATGAACAAAGCCTCCAAGAAATATGGGACTATGTAAAAAGACCAAATCTACGACTGATT  
GGTGTACCTGAAAGAGACGGGGAGAATGGAACCAAGTTGGAAAACACACTGCAGGATATT  
ATCCAGGAGAACTTCCCAACCTAGCAAGACAGGCCAACATTCAAATTCAGGAAATACAG  
AGAACACCACAAAGATACTCCTCGAGAAGATCAACTCCAAGACACATAATTGTCAGATTC  
ACCAAAGTTGAAATGAAGGAAAAATCTTAAGGGCAGCCAGAGAGAAAGGCCGGGTAACC  
TACAAAGGAAAGCCCATCAGACTAACAGCAGATCTCTCAGCAGAAACCCTACAAGCCAGA  
AGAGAGTGGGGGCCAATATTCAACATTCTCAAAGAAAAGAATTTTCAACCCAGAATTTCA  
TATCCAGCCAACTAAGCTTCATAAGTGAAGGAGAAATAAAATACTTTACAGACAAGCAA  
ATGCTGAGAGATTTTGTACCACCAGGCCCTGCCCTAAAAGAGCTCCTGAAAGAAGCGCTA  
AACATGGAAAGGAACAACCGGTACCAGCCACTGCAAAAAACAGCCAAAATATAAAGACCA  
TCGAGACTAGGAAGAACTGCATCAACTAATGAGCAAAAATAACCAGCTAACATCATAATG  
ACAGGATCAAATTCACACATAACAATATTAACCGAATAGTACCTCACATCTCAATACTAA  
AATTAATGTAAATGGACTAAATGCTCCAATTAAGACACAGACTGGCAAAGTGGATAA  
AAAGTCAACAACCAACAGTGTACTGTACTCAGGAGACCCACCTCACATATAAAGACACAC  
ATAAGCTCAAAAATACTGAAATATAAAGGATGGAGAAAGATCTAGCCAAGCAAATGGAA  
ACCCAAAAAAGCAGGAGTTGCAATCCTAGTCTCAGACAAAACAGACTTTAAACCAACA  
AAGATCAAAAAAGACAAAGAAGGCCATTACATAATGATAAAGGATCAATTCAACAAGAA  
GAGCTAACTATCCTAAATATATATGCACCCAATACAGGAGCACCCAAATTCATAAAGCAA  
ATACTGAGAGACCTACAAAGAGACTTAGACACCCACACAATAATAGTGGGAGACTTTAAC  
ACCCACTGTCAACATTAGACAGATCAACGAGACAGAAAGTCAACAAGGATACACAGGAA  
TTGAACTCAGCTCTGCACCAAGTGGACCTAATAGACATACTACAGAACTCTCCACCCCAA  
ATCAACAGAATATACATTCTACTACCAGCACACCACACATATTCCAAAATAGACCACAT  
AATTGGAAACAAAACCTCGCCTCAGCAAATGTAAAAGAACAGAAATCATAACAAACAATCT  
CTCAGACCACAGTGCAATAAACTAGAACTCAGGAATAAAGAAATTCCTCAAACCCGCA  
CAACTACATGGAAATTGAACAACCTGCTCCTGAATGACTACTGGGTAGAGGGGCCCTCTC  
TGCTCCACGCCAGGCAGATCTCCAGGCATCTGGAGCACCCACTCTCCTGAATAACGAAA  
TCAAGGCGCCCCACCCTTCCCGTGCAGAGAACTTGAAATTAAGAAGTTCTTTGAAACCAA  
CGAGAACAAAGACACAACATAACAAAAATCTCTGGGACACAGCTAAAGCAGTGTGTAGAG  
GGAAATTTATAGCACTAAATGCCACATGAAAAAGCAGGAAAGATCTCAAATCGACACCC  
TAACATCACAACATAAAGAACTAGAAAAACAAGAGCAAACAAATCCAAAAGCTAGCAGAA  
GACAAGAAATAACTAAAATCAGAGCAGAACTGAAGGAAATAGAGACAAAAAAACACTAC

AAAAAATCAATGAAACCAAGAGCTGGTTTTTTGAAAAGATCAACAAAATTGATAGACCAC  
TAGCAAGACTAACAAAGAAAAAAGAGAGAAAAATCAAATAGACTAAATAAAAAAAGATA  
AAGGAGATATCACCACAGATCCCACAGAAATACAACTACCATCAGAGAATACTATAAAC  
ACCTCTATGCAAATAAACTAGAAAATCTAGAAGAAATGGATAAATTCCTCGACACATACA  
CCCTCCCAAGACTAAACCAGGAAGAAGTAGAATCTCTGAATAGACCAATAACAAGCTCTG  
AAATTGAGGCAGTAATTAATAGCTTACCAACCAAAAAAAGTCCAGGACCAGATGGATTCA  
CAGCCGAATTCTACCAGAGGTACAAAGAGGAACTGGTACCATTCTTCTGAAACTATTCC  
AAAAAATAGAAAAAGAGGGAATCCTCCCTAACTCATTTTTATGAAGCCAGCATCACCTGA  
TACCAAACCAGGCAAAAGACACAACAAAAAAGAAAACTTCAGACCAATATCCCTGATG  
AACATAGATGCAAAAATCAGCTTTTTAGTGCCTCACTCAATAATAAATACTGGCAAACCA  
AATCCAACAGCACATCAAAAAGATTATCCACCATGATCAAGTGGGCTTAATCCAAGGGAT  
GCAAGGCTGGTTAACATACGAGAAGAATATGCAAATCAAGAAAAGAAATTCATAAAATA  
ATACAAGAAACAAAGTCAAAGACAAAAACCATATGATTATATCAATAGATGCAGAACTGA  
AAGCATTAGACAAAATTCAACAACAATTCATGATAAAAACTCTCAATAAAATAGGTATAG  
AAGGGACATATCTAAAAATAATAAGAGCTATTTATGACAAACCTACAGAAAACATCAAAC  
TGAAAGGGCAAAAAATGAAAGCATTCCCAATAAAAAACAGGAACAAGACAAGGATGCCCTC  
TCTCACCCTCCTATTCAACATAGTATTGGAAGTACTGGCCAGGGCAATAAGGCAAAAAGA  
AAGAAATAATGGGTATTCAAGAAGGAAAAGAAAGTCAAATTGTCCTGTTTGCAGATG  
ACATGATTGTATATCTAGAAAACCCCATCGTCTCAGCCCAAAAACTCCTTAAGCTGATAA  
GCAACTTCAGCAAAGTCTCAGGATACAAAATCAATGTACAAAAATCACAAGCATTCTAT  
ACACCAACAACAGACAAGCAGAGAGCCAAATCATGAATGAACTCCCATTCAACAATAGCTA  
CAAAGAGAATAAAATACCTAGGAATACAACCTTACAAGGGATGTGAAGGACCTCTTCAAGG  
AGAACTACAAACCACTGCTCAAGGAAATAAAAGAGGACACAAACAAATGGAAAAACATTC  
CATGCTCATGGATAGGAAGAATCAATATCGTGAAAATGGCCATACTGCCCAAAGTAATTT  
ACAGATTCAATGCAATTCCCATCAAACCTACCAATGACATTCTTACAGAATTAGAAAAAA  
CTACTTTAAAATTTCATATGGAACCAAAAAAGAGCCCGCATAGCCAAGACAATCCTAAGCA  
AAAAGAACAAGCTGGAGGCATCACACTACCTGACTTCAAACCTATACTACAAGGCTACAG  
TAACCAAAACAGCATGGTACTGGTACCAAAACAGACATATAGACCAATGGAACAGAACAG  
AGACCTCAGAAATAAAACCACATATCTACAACCATCTGATCTTCGACAAACCTGACAAAA  
ACAAGCAATGGGGAAAGGATACCCTATTCAATAAATGGTGCTGGGAAAACCTGGCTAGCCA  
TATGTAGAAAACCTGAAACTGGATCCCTTCTTACACCTTATACAAAAATCAACTCAAGAT  
GGATTAAGACTTAAATGTAAGACCTAAAACCATAAAAAACCCTAGAAGAAAACCTAGGCA  
ATACCATTGAGGACATAGGCATGGGCAAGGACTTCATGACCAAAACACCAAAAGCAATTG  
CAACAAAAGCCAAAATTGACAAATGGGATCTAATTAACCTAAAGAGCTTCTGCACAGCAA  
AAGAACTATCATCAGAGTGAACAGGCAACCTACAGAATGGGAGAAAATTTTCGCAATCT  
ACACATCTGACAAAGGACTAATATCCAGATCTTGGTTTTCTAAATACCAATCTCCATTAGC  
AGACAAGGAACTTACAGAACAACCTCCTTGGAGAAAATGGCTGATTCCAGGACTCTTACAAG  
AAAAAGACATGAACCAACAATAAATGAACCTGGACAACCCCATCGTGCCACAAGCTAAAAAA  
GTGGGCAAAGGACATGAACACCAAGACACTGTCTGATACACCAGTCAAACCTATAGAAAAA  
AAAAAAAACCTTCAAAGATAAAGAAGACATAAGCCAACCTGAACATGCGGCCAACAAAC  
AATATGAAAAAATGCTTTAAAAAGCAACATCACTAAGAAGATCATCAGACTTCTCAACAG  
CGACAATGCCAATAGTTTGAAATCAAAACATCCAAGGGGGCACAATAGAAGTGAGAAAAA  
GTACCATCTTCAAATTTCTCACGGAAAAAACTAACCAACCAGTCATTCTATACCCAGCC  
AAACTATCCTTCAAGAATGAAGGAGAAAGAATGGCTATTATTAGACAAGAAAAAAGTCAA  
TAAAAAAAATTTACAAGCCTTTCCAGATGAACCCTCTGTACAAAAGAGCCCCCTCTGAA  
GAAACAAACATGGATACGCGAGGATGATACCTGCTACCACAAAAACACACTTAAAGTACAT  
AGCCCACAGACCCTATAAAGCTGCTCCACCAAAAAAACAACAGAGAAATAGTCCTAGGAA  
AAGAACACAGGATTGAAAGGAACACACTTTTACAAAAGAAAGAAGGAATGAAGACCAAAA  
CTGTTGGTGGTCTAAATGCCCACTTAAAAGACACAGAGTGGCAAATTGGATAAAAAAAA  
AAAAAACAAGACCCATCCATCTGCTGTCTTCAAGAGACCCATCTCGAATGTAATTAGAAT  
TAGTGGATCAACCATTGTGGAGAAGCAGGATATTTGCCTAGACAAAAAGCAACCACCCAA  
AAAGAAAAAACAAGTGTGGCGATAACTAATCTAGACAAAAATATAACTGACATTTCT  
CCCTCAAAGAACTAGAACATTAAGAAGCTAGAAACCTGTACCATAAAAGTGGAAATTGAC  
CCAGAGAGGCCAAAAAAATCCCATTAACCTGGGTATATACAAAACCAACATCCGAGCACC  
CAGAGGGAGGAATATAAATGGAGGAATCATTCTAATTCCAAGAAAAGACTTAGACAGCCA  
ACAATAGACAGCCTGTCAATATAAAATAAGACACATTGCACACGTATGGCCTTTAAACA  
AGTAAACAAGCCACTTGCACAATTGGCACAGAAATAGAGTATAGAAAAGTCACAATAGCA

TAATACAAGACATAAACCAATGGAACGCAAAAAACAACTCACCCAAATGGTATTCATAC  
CAAAGAACCATAACCTCAACCCATCAATAAAACCACATATATACAACCTAACCAATATTTG  
ATAGATACTCTAAGATTGACCACATGCTTGGCCATAAAGCAAGTCTCAATGGATAAAGAA  
AAAAAATAACGTGGTACAATTCTATACATGCTTCCATATGCTCAGTGCTCATCGGAATAA  
AACTATGCAGCCATAAAAAAGAGATCTCTCAAACCACACAAGATAAGACAGATGAGATC  
ATGTCCTAACCAAAAACATGAGAGAAAAAATAAACAATATGCAGGGACTAAAAAACAT  
TTGAAATGGATGAAACACTTAAATGTAAAGCTGGAAACTGAACATAAAAAACCTAGAAGA  
TATTCTAGAACAGCAAACCTAACGCAGGACATTGGCATAAGGCAAAGATTATTTGGATAAAA  
CAGGAAAACCAATAACAACCGCCTGGAGATTAGATAATGTTATTGAATCAATGCTCACTT  
ATAAGTATGTTGGGAGATGTAAAATTAACAATGAGAGAATCCCCATCTTTTAAAGAAAC  
ACACAAAAACCTTGGACACAGGAAGGGGAAAGCAACACACACTGGGGACTTACTGTCCGA  
GGGTGGAGGGAAAAAATAAATAAAATGGGAGGAGGGTGAGGATAGCATTAGGATAAAA  
AAATACCTAATGGGAATTACAAAGCTTAATACCTGGGTATAAAAAGAACTCAGATGGGTT  
AATGGGAATACAGCAAACCACCATGGCACATAGTATACATATGTAACAAACCTGCACGTT  
CTGCAAGAAGATATACATGTACCCACCAAACTTTTTATAATATAAAAAATAAAATTTAAA  
AAAAAAGAAAAGAAAATAAAAA

>2\_cons

AATGTATTAATAAAATTAATAAAAAAAAAAAGAAAGAGAGTTCGGTTAGAGGGAGGCCGGC  
AATGGCGGCCGAGTAAGAGGCACCTGCAGCCCGCCTCTCCCACTAAAAGGAACACCAAAT  
TGAACAATAAAAAAAAACTAGAAAAATCTGCACCTATAAAACAACAATTTGAGAGAAGAT  
CACAGTACCTGGTTTTAACTTCATATCACTGAAAGAGGCACTGAAGAGGGTAGGAAAGAC  
AGTCTTGAATCGCCGACGCCACCCCTCCCCATCCCCGGCAGCGGCCGTGTGGCGCGGA  
GAGAGAATCTGTGCACTTGGGGGAGGGAGAGCGCAGCGATTGTGGGACTTTGCATTGGAA  
CTCAGTGCTGCTGTGCACAGCGGAAATTAACAAAGAGAAGGAACCAACCGCCACAAACG  
GAGAGAAAATTCAGGAAACACCAGAAACATGGAAAGAAACCGTCCCAGCGGTCGGAACC  
TGAGTTCCGGCAAGCCTCGCCACCGCGGGCTAAAGTGCTCTGGGGTCTAAATAAGAAAA  
ACAGCCAGCCAGCCAAGAGCGGCCGAGATCCAGGAGAGAGTCCAAAATCTGTGCATGCC  
CAAAAGCCGGGAAAGCAAAGAGGACGCGACCTAGTGAGACACCAGCCTGAGTTCCCCAG  
CGGTCCACCGCGCAACACACTTTCTCACACGTACCCCCAGCTCTGAGCCACGGGAGAG  
CCACCCAACCCGCACGCGACCTGAGACTAGCATAGGGAGCAGAGGACTTTGTCTGGAAAG  
CCGGGAAGCCGCACTGCTACAGAGAGATAGGGCACCGAGCAACTTCATGAGGCACCCAC  
GCACCGCCAAAACCCCAAGCAGAGCTTGACGCTCCAGCTTGAGCAGAGGGAACCCAGCT  
CCCACCATAGCACACAAGACTGCATCCTGCCCTGGGGCCCCCTTCTCCACCCAGCCCC  
ACTCATAGAACCAGGCACGACCCAAGCCCCGGCAGCCGAGAATCCTGGGGCCCTGATT  
GCTCTCCCCAGAACCAACCAAGAACGCGCCAGGGGCACAAGACCAGCTGGACCCAGCAG  
TGCAGCAGGGTCCCAGCACTCTAGCCACACAGTGTCCCCGCTCACAGGCGTGCTGGCT  
TCAAGAATGAACCATCCAGTTCACCACTGTGGTAACTATGGCACCACACGCATCCACCCC  
GACAAAAGGAGGGGGCACAAAAGGGAACCAAAAGCGGCTACACCTCCAGCCTGGGAGCT  
GCCTCTAAGGCAGGGCGTCACCCACAGAGAGGACCCCTGCAGCCCCAGCTGCAGGGCT  
GCCACACACCTGCACGCACCCTGAGGACAGGCTCTCCCTACCTACTGCTGCTGCCAGTGC  
AAAGTGTATGCTCCCCAGAGCCTGAGGACCACCTGCCTGGGGCTGCGCCACGGACAGCA  
ACCCTGCCCGACCCTGGCCCAACAACAGCGCACGGCCGTGAAGGCGGGGCAAAGGCC  
AGGACTCTAAAAGCTCCCCCCTGCCTACTAATGAGACTCGAGGCATAGAGTGGACAAC  
CCAAGAGCCTGACAGCACCCGCCCTGGGCCGCGGCCACTGGTAGCAAACTGCCCCCC  
CCCAGCAGCAGGGCAGCTGAGCAACGGCTGGCACCCCGAGGACAGACTGCCCTTCCCAC  
CACTGCTGCTGCCACCCAGACACTAAGGCCTTCTCCAGCGGCCTGGGGATCAGCCCGGC  
CCCACCAACAAAGGCCACCGAAACACCCTGAGAGCACCACCTGGAGGCCTGAGAACCG  
GACCCCCGGCCGACTCCAGCAACACCACTCCCTATCCGTGCCGGCTTGGGGCCAGAGC  
AAGCCTAACCTGCCCGATCGGATGGTCTTCTCTACCAACCCTGGTAGCAGAAGACAA  
AAGACATAATCTCTTGGGAGCTCCCCAGCCAGCCACCACCACTGGAACCTAAGCACTC  
ATCCCGGAGCCTGAGGGCAAGGCCACAAAACCCACCACCAAAAGCGGCAAGAAACCAC  
CACAACCTGGCCCCCTCTGAACGTGCCACCCAGGGCCCTGGGGACTGGCCCGCCAGCCC  
ACCACAGACAAAGAAGAGAGCAGTAAGGAACCACTCGGGCGCCAGAGGAGGGGCCACAA  
CTGCTAATGCCACTGCCACGCCATGCCTGCTACCCAGGGGCCGAGAACCTGCCACCC  
AGCAGGCCACTGCTACCACTGCTCCAACCCGAGAAAGCCACCGCAGGCCTAGGGACCC  
CGTTTGTGGGGCCAGCCAAAAGAAGAGAAAACCAGTGCCAGCCTAGGAATCCCGAGAG

CCCCAGCACAGACCAACCAAGCCCACCGCTGCCGCCCCAGACCAGTCCCCAAGAGCCAC  
AGCTGAGAGACCCATAGATGGTTCACATCTACAACCAAGGACCCTCACAGAGTCCACTTC  
ACTCCCCTGCTACCTCCACCGGCCAAGAAGTGGCATACTGGCTTCCAAATCCACAGAA  
AAACAACATCACAGAAATCCTACAACACACCCCCCAAATAACTACCAAAGCCTAGAACT  
CCACTGGGTGGCTAGAGCCACTGAGGAAATCACAGACACCACTGATACTGTTTAAAGACG  
AAAAAACCTAAAGGAACGGGGAGAGCACCACATCAAGGGAGCAACCCCCACAATGCAAG  
AAGGCAAAAACCAGATCCAGAGTCTCCTATCCTCCCTACAACATAGGTACACCTAAAGGA  
AAAATGTCCTCCCCTACGAAAGCAAATTCAAAAAATTGGAAGAAGCAACTGTTACACCAG  
ATGCACAGATTTCAATAAAAAGAAACAAAAAATTACAAGAAATACGGAGGAAAGAAACAG  
GAAAAAGGACTACTATGACCCCTCCAAAGAAGCAGAAGAAAATCAATCACCAGAAACAGA  
CCCCGAGGAGACACAGAATCGTTAGTATTTGAATTACCAGACAAAGAACTTTAAAAATA  
CTATTATAAAAGATGCTCAATGAAATAAAAGAGAAACATGGATAAAAAAAAAAAGGAAGA  
CAGAAAAACAATAAAGGATAAGATCAGAGAAATTCAACAAAGAGATAAGAAATTATAAAA  
AAGAACCAACAACAAGAAATTCTAAGGAACTGAAAAATAAAACAATAACTGAAATGAA  
AAAACATATACTGGAGGGGCTCAAAACAGCAGAATTGATACAATTGCAGAAGAAAAGA  
ATCAGTGAACCTGAAGACAGGTCAATAGAAATTATCCAATCTGAAGAACAGAAAGAAAA  
AAAATTTAAAAAATTAATAATGAACAGAGCCTCAGAGACATGTGGGACACCATCAAGAGAT  
CCAACATATACTACATGGACAGACAAAAACATATGAATTATTGGAGTCCAGAAAGGAG  
AGGAGAAGAGAGAGAAAGGGGCAGAAAAAATATTTGAAGAAATAATGGCTGAAACTTCC  
CAAATTTGACCCTGAAAAGACAATAATATTAATATACAGATAGATTCAAGAAGCTCAACG  
AACCCCAAGCACACCTGGGAGAATAAACCCAAAGAGATCCACACCAAGGCACATCATAGT  
CAAATGCTGAAAACCAAGACAAAGAAAAAAGAAATCTTAAAAGCAGCCAGAGAAAA  
AAAGACACATTACCTACAAGGGAAAAAACAACAATAAGACTAACAGCAGATTTCTCAGCA  
GAAACCTTACAGGCCAGAAGACAGTGGAACGACAATTCTTCAAAGTGCTGAAAGAAAA  
AAAAAAAAAAAAAAAAACAAAAAACACAGAATTCTACAATAAAAAAAAAAACTGTCAACCA  
AGTAAATTCTATATCGAAGCAGCAAAAAATAATCCTTCAAAAATGAAGGAGAAATAAAGAC  
ATTCTCAGATAACGACAAAAGCTGAGAGAATTCATCACCAGCAGACCTGCAAACCTACAA  
GAAATGCTAAAGGAAGAGTTCTTCAGGCTGAAACGAAAGGAACACAACAAAACAAAATCAA  
AACAAAAGAAAAAACAACACACACACAAAAAACAACAAAACAAAAGAAAAAACA  
GCAGAAAGACAAGAAAGGAATTATGAAAGAAATCAACAGAAAAAAGAAAAAAGAAAAA  
AACAGAAGAAAGCCTACGGCACTTAACACACCAAAAAAGAAATAAAAAAACAATATAG  
AAATAACAGAAGGAAAGACATAGAAAAAAGAAAAAGAAAAACATATTTAACAAAAAAGA  
ACAAAAACGTCCCAAGGAGCAAAACAGAAATGAACATCCAGATCAGGAAGCTCAAAGAT  
CCCCAATTAGATTCAACCCAAAAAGATCCTCTCTGAGGCACATTATAATCAAATGTCAA  
AAGTCAAAGACAAAGAGAGAATTCTAAAAGCTGCAAGAGAAAAGCATCAAGTCACATATA  
AGGGAATCCCCATTAGACTATCAGCAGATTTCTCAGCAGAACTTTGCAGGCCAGGAGAG  
AATGGGATGATATATTCAAAGTGCTGAAAGAAAAAATAAAAAAAAAAACTGTCA  
GCCAAGAATACTATATCCAGCAAAGCTATCCTTCAGAAATGAAGGAGAAATAAAGACCTT  
CCCAGACAAGCAAAGCTGAGGGAATTCATCACCCTAGAACTGGCCTTACAAGAAATGC  
TTAAGGGAGTGCTACAACCTGGAAAGGAAATGATACACAGATAGAAACATGAAACCACAGG  
AATCATGGAATAAAAGACTGGCAGTGATCACCAGAAACAGGTAATACATGGGTAAATAT  
AAAAAATAATAATACTGTAATTATGGTATGTAATTCTTTTTTTCTTTCTTAAATTCTC  
TAAAAGACAAATGAATGTTTAAAGATAAAATGTACAAAAACAAAAATAACTATAACAAGT  
GTCTGTAACCTCACATTTTGTCTTCTACATAATTTAAGAGACTAATGCATTTAAAAAAT  
TATTAGATAAGGATGTGGGGCATATAATATACAAAAATGTAGATATAATAAGACTGTAAA  
ATGTATGACAACAATAACAAAAATGGAAGTAGGGGGGGGAAATTAGAAAGGAGAGTT  
TTTTATAGTGTGTAAGGTTTTTACATTTTATACAACCTGGAAGTAAGATAATATCAATTC  
AAAGTAGACTGTGATAATTAGAGAAATATAGGGTTAAGGATGTATATTGTAATCCCTAGA  
GTAACCACTAAAAAATTATAAAAAAAGATATACAAAAAACAATTAAGAAAT  
AAAAACGTAACAATAAAAAAACATTCAAATAAACCAAAAAGAAAGCCAGAAAAAAGAAA  
AAAAGGAACAAAAACAGATAAGACAAAAGACAAACAGAAAAACAATAACAAAATGACAG  
AAATAAGTCCAAACATATCAATATAACATTAATGTGATTATGGATTAAAAAATGGCA  
AAAGCTGTCAGACTAGAGATTTAATATATAAAATCCAAATAAATAGTTAAAATGATAAG  
ACAGATAATACAAATATCAATAATAGGCTACATTAATGTAAATGGACTAACTCTCCAA  
TTAAAAGACAAAGAGATTGTCAGAATGGATAATTAAAAAAATAAAAAACAAGACCCAAC  
TATATGCTGTCTACAAGAGACACACTTCAAATATAAAGACACAAATAGATTGAAAGTAAA  
AGGATGGAAAAAGATATATCATGCAAACGAAACCAAAAAGAAAGCAGGAGTAGCTATACT

AATATCAGATAAAATAGACTCTTCAAACAAAAAATATAACAAGAGACAAAGAAGGTCAT  
TATATAATATAATGATAAAAGGGATAAAGGGGTCAATCCAACAAGAAGATATAACAATTA  
TAAACATATATGCATAAATATATATGCACCAAACAACAGAGCCCCAAAAATACATGAAGC  
AAAACTGACAGAAATGAAAGGAGAAATAGACAAATCAACAATAATAGTTGGAGACTTCA  
ACAACCCACTCTCAACAATGGATAGAACAACACTAGACAGAAAATAAGAAAAGAAACAAAA  
AACACAACAATACAACAAAACAATAAACCAACTAGACCTAACAGACATCTAAAGAACATT  
TATAGAACACTCCACCCAACAACAGCAGAATACACATTCTTCTCAAGTACACATGGAACA  
TACACCAAGATAGACCATATCCTAGGCCATAAAACAAACCTCAACAAATTTAAAAAAGG  
AAAAAATAAAAAAAGGATCTCCTACCACAACAAAAGAAAAAAGAAAAAACAACAAA  
AAAAAACAGGAAAATCTACAAACACGTGGAACTAAACAACACACTCCTAAAAAACCA  
GGGGCAAAAAAACCAAAAGAAAAATAAGAAAATACTTTGAGATGAATGAAAATGAA  
GACACAACATACCAAAATTTATGGGATGCAGCTAAAGCAGTGATTAGAGGAAAATTTATA  
GCTGTAAATGCCTATATTA AAAAAGAAGAAAGATCTCAAATCAATAACCTAACCTTCTAC  
CTTAAGACACTAAAAAAGAAGAGCAAACCTAAAGCAAGCAGAAGGAAGGAAATA  
ATAAGATTAGAGCAGAAATTAATGAAATAGAAGAAAAACAATAGAGAAAATCAATGAAA  
CCAAAAGCTGGTCTTTGAAAAGATCAACAAAATTGACAAACCTTTAGCTAGACTGACCA  
AGAAAAAGAGAAGACTCAAATTAATAAATCAGAAATGAAAGAGGGAACATTACTACTAA  
CCTTACAGAAATAAAAAGGATTATAAAGGAATACTATGAACAATTGTATGCCAATAAATT  
AAGATAACTTAGATGAAATGGACAAATTCCTAGAAAAAAGACACACAACTACAAAAAC  
TGACTCAAGAAGAAATAGAAAATCTGAATAGACCTATAAAAATAAAGAGATTGAATTAGT  
AATATAAAAACCTACCAACAAAAAAGCCAGACCCAGATGGCTTCACTGGTGAATTCTCC  
AAAAATTTAAAAAAGAATTAATACCAATTATTCACCTATTCCAAAAAATAGAAGAGGAGG  
AAAAACTACCAACTAATTCTATGAGGCCAGTATTATCCTGATACAAAACCAGACAAAG  
ACATAACAAAAGAAAAGAAAA

>3\_cons

CTTGGCACTTAATTCACCAACTCTGGGAGCAGAGGGGATTGGACATATGAAAAGAGAGTG  
ACGTCAGCAAAAAGGTAGTCTTTTAGCCAAGATGGCGGACTGGAGGAGCAGCCAGGGTCC  
GCCCCCCCAGGAGCAGCACAGAGAAAAACCAAAAAAATCGAGTTCCTGTTTCTCCCTG  
GAAGGGGTTTATGAAACATTCACACTGCGAACACTTGCCAACAAAATCATCTAAGTGAGA  
GCACTGGGGATTGAGCAAAGAGGGCAGCGCAACCAGGTGGAGCACCAGAGACTGGAAAAG  
ACCCAATGAAAGAGGAAGGAAGAAGACTGGATTACATTACCCACGGCACCCCTCCCCACCC  
GGGAACAGCACGAGAGCCAGGGGCACTCTCCCGGCCCAATGGCTCCTGCACGGGAAA  
GGGTGAGGGACTGTGCCATGAGGGAGCCCCAGCGCGGGGCCAGAGACCCACTTCCCCCA  
CGGACCTAGCTGCAATCCTGGCCACAGGAGACCCCCAGCCCCCGACCACCAGGGTA  
CCCACGCACAGAGCCGAAACAGACCCCAAGACAAAGCCTCTGGCCTGACACAGGGAGCA  
CAAGGTGTGCACTCCCCACCCCACTGACTCAAGCTGCTGCAGCACAGCCACCCGCTGCC  
CGGGAATGCTGCAGACTCCCTAGACCACAGGCAGAGCCAGAGCTCCAGCAGGCGGAGAGG  
AACATGCGCACCAGCCAAACCCACAGGCCACACGGCTACTGCGCTACGCCACCTCGGAA  
CTGCCCTTGACCCGGGGGAGCCTCAGCATCCCGGCAAAAAGACATGCAACCCCGCCAA  
AACC GAAGCTAAGACGCCCTGAACCACCCAGCCAGCGAGCCTAGCCCCGGAGCACGG  
CAAGCAACGCCATCCCCGCACTTCCCCCAAGGCGCGACAGCCGAGAACCCGCACCACC  
TACAACCTCCAGTCAAGCTACACAAACAGCGGCATACCCCCAGGACAGAGCTAAAGGG  
GGCTCTGGCCCCGCCAGCTGCCGCACCACGGGCACCCAGGCTGGCCCCCTGTGGAATCTT  
GCCTGCGGCACACGGAAGATCAAACCTGACCTTGTGGGGAAAGGATCCCCCAGCACAGCA  
CAGCTGCTTACAAAACGTGGCCAGACTGCTTCTTTAAGCGGGTCCCAGACTCAGCACC  
CCCCCCCCGGGCAAGACCACCTCCCCTGGGGCCAGAGCCTACCCCGCCAGCCCCACCG  
CAGCCTCGGAGTTGAAACCTCCCTGGGACGGAGCGCCCGGGGGGAGGCGCTGCACC  
TGCCCCTTTTCCCGGAAGTACACGAACCGCCATCCCCATGTCGGGGGCCTGAGCACAC  
CGGCCGCAAGCAACCCCCACCCACCCGGGCGGCAAAACTCTGGGCCAGTGGCCAGCCA  
ACCCCGTACCCGTATCACAGCCACAGCGCCACCCAAGGGGCCCGCCCTCAGGCCGAGG  
CCCAGACCCCTCCCCCAGAGTTCGACCACACACAAAGGGTAAACAGCCCAGAGCCG  
GGCCCCAGAAGAACC GGGAAGAGCCTCCCCAGTTCCCTGCACAGCCAGCACACCCC  
CCCCACCCCGGGCAGGCCGGTGCCAGTGGCCTGGGGGCGCAGCACAGCCCCACGC  
CCCCACAGCTGCGATCAAGCCAGGGCCCCAGACCCCTTTTCTGGGAAGGCCACAGGCC  
TCCCCACCCGGAACCTCCCTCCCCTGCCAGAGAGGCCCCACACGCGGGGACTCTCCCGC  
CGCCCCGTGACCCAGCACTGAGCCAGTTGCAGCTGCACCACCTCTGGGGCCGCACCCA

CAGAGAGGGGCCCCGAAACTCCAGCGGCCCTATCCCTGCCCCAGGCCAGCTGAAGGCTCCA  
GGTACTGGAAAATCCGAGGCGACTAGGGACTGGAGCGGGCCCCCAGCATACCGCAGCAGC  
CCTACGGAAAAGTGGCCAGACTGTTACGTGGGTGCCCGTTCCCATATCTCCTACCCGGGC  
AGGTCCTCCAGGCCTGGGCCTCCAGCCACCCCCGCCAGAGCTATCGAGCCAGTAGCAAC  
TCGGCAACTCCCTGGACAGAGCCTCCAGGGGCAACTGGAATGCCCTCTGCCCCACCCCA  
CAGCCGTCTGCCAGACAAAAAACCAGCCACCCAAGCAGCACACCAGCTATACACACCCA  
CGAACTACCCCCGACCAACGATCCAACGACAGTCACCCACCCCCCACCAAACCCCACT  
CACACCAAGGCCACCACGAGCCACCGTTCGCCCCGCCCGGACAACCTCCGTGCACCCCG  
GCCCCAGAGAGGCCACCCGCTGGCCCCGAGCCTCAGGCAGGAGCACCCTAAGGAAGGGGGG  
AGTGCAAACAGCCCTGCACCCCTGGCTGCCGGCCACAGGCGATAGCGGTCCCACCTATC  
GGAGGAGCACCAGATGCCACGCCAGGCGACCCCGCAGCACCCTACCTCTAACAACAAAA  
GAAAACACACCGGCCAGAGTAACAGGCCAGAGCGGGCCCCATTGCCCGGCCCGCTCCC  
CGGCTGCCACCAACAGGCCCAACAAAATAAACTGAGATCGCCCCAGAGCTGCAGTGG  
GCAGCCCAGGAGTGCCAAGCCACGATCTGCTGCCAGCAGACAGGCAACAAAGTAAACCCC  
CCCCAAAAAATAAAGAGCAAAGCTGACCCACATTAGACAAAAAATGGAACAAGCAT  
CTGACAGCCCTGACTCTTCCACAAGCGGCCGCTCAGGACAAAGACCCCAAACTTCAACT  
CTGCCACCCAAGCACACACCAGAGCACCAAGGCCAAGAACTGTGCTGACCCAGCCCC  
CCTGAAACCAAGGACAGGAAATTAGCCACAAATAAAGATCCTGCACAGAGCCTTGGCCCT  
CTGAAAGCACCAGAAAAGAAGCCAACAACTCAACCAACTTACACAGCCTTCCAAAAT  
CAAAGAAACCCCAAGGAGATCCAAGAATACAGCACAAACAAACAGGAAAGCGAGAAAAA  
AACTCCACCAAAAAGAACACACACTGCCTAAAAAAGACTGAAATGCTGAGCCCTGAGTAC  
CTGATGTTAAAACCCTCAAATGGCCCACCACTCCACTACTAGCTTCTAAGCAAACATCC  
ATAATGAGAAGAAACCAATCAGCAAAAGGGGCTCACACAAGGAATACAAAAAAGCAGGA  
GAGAGGAAACATGACACCTACCAAGGAACACAATAAATCTCCAGAAACGGCACAAAACC  
CTAAGGAAATTGAGATGTCTGGCTGAGTAAGAGCCTAATTGACAGACAAAGAATTCAAAA  
ATAACGATCATAAACAGAAGCTCAATGAGATACAGGACAAAAGCTATAGAAAACCAATTC  
AAGGAAACCAGGAAAACAATCCATGAAATAAATGAGAATTTCAACAAAACAGAAATAGAA  
ATTATAAAAAAGAACCACAAACAGAAAAATTCTAAGAGCTGAAAAATATACAATAACTGA  
ACCCAAAAAAAACCATTAATAAATTCAATAGAAAGTATCAACAGCGCGCGCAGAATAGATC  
AAGCAGAAGAAAGAATTCCCCTGAACCTGAAGACAGGTCATTTGAAATTATCCAGTCAG  
ACGAACAAATCCACAAGAAAAAAGAATGAAAAAAATTGAAGAAAGCCTCAGAGACATAT  
GGGACAACATCAAACAAGAGAACCAACATATGCATAATGGGAATCCAGAAAGGAGAAGAG  
AAAAAAGAGAAAAGGCTCAGAAAGCATATTTAAAGAAATAATGGCTGAAAACCTTCCAAA  
TCTGGCGAAAGAGATGGTATATAAATCAAAAGCCAATAACATCCAGATACAAGAAGCACA  
AAGAACCCCAATAGATTCAACACAAACAAAATAATCAACATATACACACAAAAAACC  
AATAAAAAAAGTCAACAAGAAAAAATAAATCTTAAAGGCAGCAAGAGAGAAAAGAAAGA  
TCACATACAAAGGAATCCCAATAAGACTAACAGCAGACTTCTCAGCAGAAACCTTACAAG  
CCAGAAGAGAGTGGGATGATATATTCAAAGTGCTGAAAGAAAAAATACTGTAGCCAAGA  
ATATTATACCCAGCAAAGCTATCCTTAGAAATGAAGAAATAAATCATCCACATACCAAC  
CAAGAATAATATATCCAGCAAAACTAACCTTCATAAATGAAGGAGAAATAAAAACCTTCC  
CAGACAAACAAAAACTAAAGGAATTCATATCAACTAGACCTGCATTACAAGAAAACCCAA  
AGGGAAAGTCTATAAACAGAAAAGAAAGAAAAAGAAATCAAACCACATAAAAAACAATAAAA  
CACAAAACCAAAAGATATAAAAAAATAAGACAAAAAATAACATAACAACCAGAAA  
ACAACAAAATGACAGAAATAAATCTCACATATCAATAATAACCTGAATGTAAATAGAA  
TAAACTCCCATAAAAAAATAAATGGATAGATAAAAAAAGAAACAACCTAAACA  
CATCCTAAAAAACAACCCCAACCCAAAAAGAAAAACATAGACGGAAACAAAAAATG  
AAAAAAGAAACACCATAAAAAAACAACCAAAAGGAGAAGAGTAGCAAGAAAGACACAA  
GACAAAAAAGCAAAAAACAAAACTAGAAAAAACAAGAAAAAAGACAAAAGAAAAA  
AAAAACGACCAATTGAGCAAGAGATATAATTGTAATATGTGTACAATTGTAAATATATAT  
GCACCCAAACACTAGAGCACCCAGATATATAAAGCAAATATTATTAGATCTAAAGGGAGA  
GATAGACCCCAATACAATAATAGTTGAGGACTTCAACCCCACTCTCAGCATTGGACAGAT  
CATCTAGACAGAAAATCAACAAAGAAACATGATTTAAACTGCACCATAGACCAATGGAC  
CTAAATAACAGACATTTACAGAACATTTACCCCAACAGCTGCAGAATACACATTCTTTTC  
ATCAGCACATGGAACATTTCTCCAGGATTGACCATATGTTAGGACACAAAAACAAGTCTCAA  
CAAATTTTATCAAGTATCTTATCTGACCACAATAGAATAAACTAGAAATCAATAACAAG  
AGGAACATTCAAACTATACAAATATATGGAAATTAACAACATGCTCCTGAATGACAAT  
GAGTGAAGAAGAAATTAAGAATGAAATTTAAAAATTCCTTGAAACAAATGAAAATAGAAA

CACAACATACCAAAACGGGAACAGCAAAAGCAGTTATTAAGAGGCAAGTTT

>4\_cons

GAAGGCGGAACAAGATGGCCGAATAGAAGACTCCACCGATCATCCTCCCTGCAGGAACAC  
CAAATTGAACAACACTATCCACACAAAAAATACCTTCATAAGAACCAAAAATCAGGTGAGC  
GATCACAGTACCTGGTTTTAACTTCATATCACTGAAAGAGGCACTGAAGAGGGTAGGAAA  
GACAGTCTTGAATTGCCGATGCCACCACTCCACCATACCCGGGCAGTGGCAGAGTGGTGT  
GGAGAGAGAATCTGTGCGCTTGGGGGAGGGAGAGTGCAGAGATTGTGAGACTTTGCATTG  
GAACTCAGTGCTGCCCTGTCACAGTAGAAAGCAAACCAGGCAGAACTCAGCTGGTGGCC  
ACGGAGGGAACATTTAGACCAGCCCTAGCCAGAGGGGAATCGCCTATCCCAGTGGTCCGA  
ACCTGAGTTCCGGCAAGCCTTGCCACCGCGGGCTAAAGTGCTCTGGGGATCTAAATAAAC  
TTGAAAGGCAGTCTAGGCCAAAAGGACTGCAATCCTAGGCAAGTCTAGTGCTGAACTG  
GGCTCAGAGACAGTGGACTTGGGGGACACATGACCTAAGGAGACACCAGCTGGGGCAGCA  
AAGGGAGTGCTTGACCACCCCTCACTAACTCCAGGCAGCACAGCTCACGGCTCCGAAA  
GAGACTCCTTCCTTCTGCTTGAGGAGAGGAGAGGGAAGAGTAAAGAGGACTTTGTCTTGC  
AACTTGATAACCAGCTCAGCCACAGTAGGATAGGGCACCAACAGAGTCATGAGGCCCCC  
ATTCAGGCCCTGGCTCCCGGACAACATTTCTAGACACACCCTGGGCCAGAAGAGAACC  
GCTGCCTTGAAGGGAAGGACCCAGTCTGGCAGGATACATCACCTGCTGACTAAAGAGCG  
CTTGGGCCCTGAATGATCAACAGCGATACCCAGGCAATACTCAATGTGGGCCTTGGGTGA  
GACTCAGAGACTTGCTGGCTTCAAGGTGTGACTCAGCACATTCCCAGCTGTGGTGGCTATG  
GGGAGAGACTCCTTATGCTTGAGAAAAGAAGAGGGAAAAGTAAAGGGGACTTTGTCTTGC  
ACCTTAGGTACCAGCTCGGCCACAGTGGGATAGAGCACCAAGTAGGCTCTTGGGGTCCCC  
GATTCCAGGACTTGGCTCTTGGATGGCATTCTGGACCTGCCCTGGGCCAGAGGAGAGCC  
CACTGTCTGAAGAGAGAGTCCCAGGCCTGGCAGCATTACCACAAGCTGACTGAAGAGC  
CCTTGGGCCTTGAGAGAACATTGGCGGTAGCCAGGCAGTACTCTCCATGGGCCTGGGATG  
GTGGTGGCCACAGGGAGCGACTCCTTTGCCTGTGGAAAGGGGAGGGAAGAGTGGGAAGGA  
CTTTGTCTCGTGGTTTTGGGTGCCAGCTCAGCCACGGTAGAATAGAGCACCAAGGTAGATTT  
CTAAGTTTTCTGACTCCAGGCCCTGGCTCCCGGATGACATCTCTGGACCTGCCTGAGGCC  
AGGGGGAACCTTACCACCCTGAAGGGAAGGACACAAGCCTGGCTAGCTTTTACAACCTGCTG  
ATTGTAGAGCCCTAGGGCCTTGAGCGAACATAGGCGGTAGCCAGGAAGTGGTTACAGCAG  
GCCTTGGGTGAGACCCAGTGCTATGCTGGCTTCAAGTCTGACCCAGCACAGTCCCAGTGG  
TGGTGGCCACAGGGGTGCTTGTGTACCACCTCCAGCTTCAAGCAGCTCAGAACAGAGAG  
AGAGACTCCATTTGTTTTGGGGGAAAGTAAGGGAAGAGAACAAGAGTCTCTGCCTGGTAAT  
CCAGAGAATTCTCCGGATCTTATCCAAGACCACCAAGGCAGTACCTCTATGAGTCTGCA  
AGAACCACAGTGTTAATGGGCTTGGGGTGGCCCTAAAGCAGATATGGCTACATGACCAA  
AACTTAGATCAAAACACCCAAGTCCATTCAAATACCTGGAAAGCCTTCCCAAGAAGAAT  
GGGTACAAACAAGCCAGACTGTGAAGACTACAATAAATACCTAACTCTTCAATGCCCAG  
ACACTGACGAACATCCACAAGCATCAAGACCTTCCAGGAAAACATGACCTCACCAAACGA  
ACTAAATAAGGCACCAGTGACCAATCCTGGAGAAACAGAGAGATATGTGAACTTTCAGAC  
AGAGAATTCAAATAGCTGTTTTGAGGAAACTCAAAGAAATTCAAGATAACACAGAGAAG  
GAATTCAGAATTCTATCAGATAAATCTAACAAGAGAATGAAATAATTAAGAAAGAATCAA  
GCAGAAATTCTGGAGCTGAAAAATGCAATTGGCATACTGAAGAATGCATCAGAGTCTATT  
AACAGCAGAATTGATCAAACAGAAAAAAGAATTAGTGAGCTTGAAGACAGGCTATTTGAA  
AATACACAGTCAGAGGAGACAAAAGAAAAAAGAATAAAAAACAATGAAGCATGCCTACAA  
GATCTAGAAAATAGCCTCAAAGGGCAAATCTAAGAGTTATTGGCCTTAAAGAGGAGGTA  
GAAAGAGAGATAGGGGTAGAAAGATTAATTCAAAGGATAATAACAGAGAACCTCCCAAAC  
CTAAAGAAAGATATCAATATTCAAGTACAAGAAGGTTATAGAACACCAAGCAGATTTAAC  
CCAAAGAAGACTACCTCAAGGCATTTAATAATCAAACCTCCCAAAGGTCAAGGATAAAGAA  
AGGATCCTAAAGCAGCAAGAGAAAAGAAATAACATGCAATAAAGCTCCAATACGTATGG  
CAGCAGACTTTTCAGTGGAACCTTACAGGCCAGGAGAGAGTGGCATGACATATTTAAAG  
TGCTGAAGGAAAAAACTTTTACCCTAGAATAGTATATCCAGTGAAAATATCCTTCAAACA  
TGAAAGAGAAATAAAGACTTTCCAGACAAAACAAAAGCTGAGGGATTTTCATCAACACCAG  
ACCTGTCTGCAAGAAATGCTAAAGGGAGTTCTTCAATCTGAAAGAAAAGGACGTTAATG  
AGCAATAAGAAATCATCTGAAGGTACAAAACCTCACTGGTAATAGTAAGTACACAGAAAAA  
CACAGAATATCGTAACACTGTAATTGTGGTATGTAAACTACTCATATCTTAAATAGAAAG  
ACTAAAAAATGAAACAATCAAAAATAATAACTACAACAATTTTCAAGACATAGACAGTAC  
AATAAGATATAAATAGAAACAACAAAAGTTAAAAAGAGAGGGGATGAAGTTAAAGTGTA



GAGTTTTTATTAGTTTTCGATTGTTTGCTTGTTTGTATGCAAACGGTGTTGTTATCAG  
CTTAAAATAATGGGTTATAAGATAATATTTGCAAGCCTCATGGTAACCTCAAATCAAAAA  
ACATACAACAGATACACAAAAAATAAAAAGCAAGAAATTAATCATACCACCAGAGAAAA  
TCACCTTCACTAAAAGAAAGACAGGAAGGAAGGAAAGAAGGAAGAGAAGACCACAAAAACA  
ACCAGAAAACAAATAACAAAATGGCAGGAGTAAATCCTTACTTATCAATAATAACATTGG  
AATGTAAATGGACTAACTCTAATCAAAAAGACATAGAGTGGCTGAATGGATAAAAAAAAC  
AAAACCCAATGATCTGATGCCTACAAGAAACACACTTCACCTATAAAGACACACATAGAC  
TGAAAATAAAGGGATGGAAAAAGATATTCATGCAAATAGAAACCAAAAAAGAGCAGGAG  
TAGCTATACTTATATCAGACAAAATAGAATTAAGACAAAAACTATAAGAAGAGACAAAG  
AATGTCATTTAATGATAAAGGGGTCAATTCAGCAAGAGGATATAACAATTTTATATATAT  
GCACCCAACACTGGAGCACCCAGATATATAAAGCAAATATTATTAGAGCTAAGAGAGAGA  
TAGACCCAATACAATAATAGCTGGAGACTTCAACACCTGTCTTTTAGCATTAGAAAAAT  
CATCCAGACAGAAAATCAACAAAGAAACATTGAACTTAATCTGCACTATAGACCAAATGA  
ACCTAATAGAAATTTACAGAACATTTAATCCAACAGCTGCAGAATACACATTCTTCTCCT  
CAGCACATGAATCATTCTCAAGGATAGACCATATATAAGGTCACAAAACAAATCTTAAAA  
CATTCAAAAAATTGAAATTATATCAAGAATCTTCTCTGACCACAATGGAATAAAACTAGA  
AATCAATAACAAGAGAAATTTTGGAAACTATACAAACACACGGAAATTAACAAAATGCT  
ACTGAATGACCAGTGAGTCAATGAAGAAATTAAGAAGGAAATTAAAAAATTTCTTGAAAC  
AAATGATAATGGAAACACAATATAACAAAACCTATGAGATACGGTAAAAGCAGTACTAAG  
AGGGAAAGTTTATAGCTGTAAGTGCCTACATCAAAAAAGAAGAAAAACTTCGAATAAACA  
ACCTAATGATGCATCTTAAAGAACTAGAAAAGCAAGAGCAAACCAACCCAAAATTAGTA  
GAAGAAAATAAATAAATAAAGATCAGAGCAGAAATAAATGAAATTGAAATAAAGAAAAACA  
TACAAAAGATAAATGAAACAAAAAGTTGGTTTTTTGAAAAGATAAACAAAATTGACAAAC  
CTTTAGCCAGAATAAGAAAAAAGAGAAGACCCAAATAAATAAATCAGAGATGAAAAAG  
GAGACATTACAACCTGATACCACAGAAATTCAAATGATCATTAGAGGCTACTATGAGCAAC  
TATATACCAATAAATCGGAAAACCTAGAAGAAATGGATAAATTCCTAGACACATACAACC  
TACAAAGATTGAACCATGAAGAAATCCAAAACCTGAACAGACCAATAACAAGTAACGAGA  
TCGAAGCCGTAATAAAAAGTCTCCCAGCAAAGAAAAGCCCGGGACCTGATGGCTTCACTG  
CTGAATTTTACCAAACATTTAAAGAATTAATACCAATCCTACTCAAATTTCTGAAAAA  
AGAGGAGGAAGGAATACTTACAACTCATTCTATGAGGCCAGTATAACCCTGATACCAAA  
ACCAGACAAAGACACATCAAAAAAATAAATAAAGGCAATATTCCTGATGAATATTGA  
TGCAAAAATCCTAAACAAAATACTAGCAAACCAATTCACAACACATTAAAAAAATCAT  
TCATCATGACCAAGTGGGATTTATCCCAGGGATGCAAGGATGGTTCAACATATGCAAATC  
AATCACAATCGATATGATACATCATATCAACAGAATGAAGGACAAAAACCATATGATAAT  
TTCAATTGATGCTGAAAAAGCATTGATAAAATTCACATCCCTTCATGATAAAAAACCT  
CAAAAACCTGGGTATAGAAGAACATAACTCACGACACAATAAAGCCATATACGACAGACA  
CACAGCTAGTATCAAATGAATGGGGAAAAACTGAAAGCCTTTCCATTAAGATCTGAAAC  
ATGACAAGGATGCCCACTTTCAACACTGTTATTCAACATAGTACTGGAAGTCTTAGTTAG  
AGCAATCAGACAAGAGAAAAGAAATAAAGGGCATCCAAATTGGAAAGGAAGAAGTCAAATT  
ATCCTTGTTTGCAGATGATATGATCTTATATTTGGAAAAACCTAAAGACTCCACCAAAAA  
ACTATTAGAACTGATAAACAAATTCAGTAAAGTTGCAGGATACAAAATTAACATACAAAA  
ATCAGTAGCATTCTATATGCCAACAGTGAACAATCTGAAAAAGAAATCAAGAAAGTAAT  
CCCATTTACAATAGCTACAAATAAAATTAATACCTAGGAATTAATATAAAAAAATAAAT  
GAAAGATATCTAAAATGAAAACATAAACAACACTGATGAAAGAAATGAAAGAGGACACAAA  
AAAATGGAAAGATATTCATGTTTATGGATTGGAAGAATCAATATTGTTAAAATGACCAT  
ACTACCCAAAGCAATCTACAGATTCAATGCAATCCCTATCAAAATACAAATGACATTTTT  
CACAGAAATAGAAAAACAATCCTAAAATTTATATGGAAACCACAAAAGACCCAGAATAGC  
CAAAGCTATCCTAAGCAAAAAGAACAACAACTGGAGGAATAACATTACCTGACTTCAAATT  
ATACTACAGAGATATAGTAACCAAAACGGCATGGTACTGGCATAAAAAACAGACACATAGA  
CCAATGGAACAGAATAGAGAACCAGAAACAAATCCATACATATACAGCGAACTCATTTT  
CGACAAAGGTGCCAAGAACATACTGAGGAAAAGACAATCTCTTCAATAAATGGTGCTG  
GGAAAACCTGGATATCCATATGCAGAAGAATGAAACTAGACCCCATCTCTCGCCATATACA  
AAAATCAAATCAAATAGATTAAGACTTAAATATAAGACCTCAAATATGAAACTACTA  
AAAGAAAACATTGGGGAAACTCTCCAGGACATTGGATTTGGGCAAAGATTTCTTGAGTAA  
TACCCCAACAAGCACAGGCAACCAAGCAAAAATGGACAAATGGGATCAAATCAAGTTAAA  
AAGCTTCTGCACAGCAAAGGAAACAATCAACAAAGTGAAGAGACAATACACAGAATGGGA  
GAAAATATTTGCAAACCTACCCATCTGACAAGGGACTAGTATCCAGAATATATAAAGAAT

CCTACAACCTCAACAATAAAAAAAAAACAAACAACCCAATTAAAAAATGGGCAAAAAGACTTGA  
ATAGACATTTACAAAAGAAGATATACAAATGGCCAATAAGCATATGAAAAGATGCTCAA  
CATCATTAGTCATCAGGGAAATGCAAATTAACCACAATGAGATACCACTTCACATCGC  
CCATTAGAATGGCTAAAATCAAAAAGACAGACAATAACAAGTGTGGCGAGGATGTGGAG  
AAACGGGAACCTCATACTGCTGGTGGGAATGTAAAATGGTACAGCCACTTTGGAAAA  
CAGTTTGGCAGTTCTCAAAAAGCCTTAAACATACTTACCATATGACCCAGCAATTC  
CACTCCTAGGTATATACCCAAAAGAAATGAAAACAAGATGTTACAAAAGATACCTGTACA  
CGAATGATTCATAGCAGCATTAGATTCATAATAGCCAAATAATTGGAAAACAACCCAAATG  
TCCATCAACAGGTGAATGGATAAAACAAAATGTGGTATATACATACAATGGAATACTATTC  
AGCCATAAAAAGGAATGAACTACTGATACATGCAACAACATGGATGAATCTCGAAAACAT  
TATGTTTAAAGTGAATAAAGCCAGTCACAAGAAGGATACATACTGTATGATTCATTTATT  
AAATATGAAGTATCTAAATATAATCAAAAAGACAAGCAGGCAAAAATAAAAAATATATTG  
TTTAGGGATACATACATATGTGGTAAAATAAAGAAAAGCAAGGGACTCATAGAGACA  
GAAAGTAGCACAGGAAAGATAAAACAAGAACTAATAATAGTGGTTACCAGGGGGCTGGGAA  
GGGTAGTGGGGGGGAGGGGGAGAGGTGGGAGTGGATGGTAAATGGGAGTACTTTTCAAAA  
GGTTCATTTTGGGATGATGAAAATGTTCTATATAAACTATTCTGTATTTGATTGTGGTGG  
TGGTTGCACTGTAACAGTGTGAAATATACTTAATATAATATTGAAATTTGCCAAAACCCA  
CAGAACTTTACAGCATAAAGAGTGAACCTTAATGTATGAAAATTTTAAAAAAACAAACAA  
GAAAAAGGGAGATCACAAAATGAAATACAACTGAAACAAAAGAACCTAACTGTATTACA  
AATGAATAACATAACCTCACTGAAGGGAATAAGGAAAAAAAGTACTAACTAAATAACTT  
TAGAAATGAGTATTTTACTACACACTCTAAGACTAAAGACAAAAAGAACTGTTAAAAAT  
AACTGAACCCTAATGAATAAGCTTATTTTCCACAGGGGCACAGGTTAACAATTCTGACAC  
TACTATATATATATATAAATTAATAAAAAGGTAAATATATTGAAGATAATGGGAGCC  
AGGTTTCTCACTGTCGGAGTAGGGAGTTACAAATATGGAAAGGGAGAAGACTAGAATGAA  
CCCTGTGGTATTGGATTGGAATTGGAGGTATCAGTATGAACTCATGCCTTTTAAATATAAA  
TAGATATACAGACAGACAGATATAGAAATAGATATAGATATATATGTGTATGTGTATATG  
TGTATGTATATACGTACATATATTTTCTAGCTCTGTCCACTGAGAGGGCCTAGAAGCAAT  
GACACCCAGTAGCAATGAGCACACCTAGCACCATGATCTTGGTATGTAAATACCATTCT  
CCACTAAAAGGAACCAGGGCTCCTTGGAGAAATGGCTGATTCCAGGGCGCACAAAGAAAA  
GATACAAGATGAGCCTGGAACATCTTGCAGTGCCAGAAAATAAGGAAGTGCTCAAAAAAA  
CAAGGAGACAGGTATGCAAAAAGGAAAAAAGAACCAACTGAAAGAGCTCCCAATGGCCAA  
AGCTGGGACAATTTAAACAAAAAAGAAAATGCTTGAGGTGATGGATACCCATTACCC  
TGATAACACAGGTGATTATTATACATTCTATGCCTGTAACAAAACATCACTATGTAAATA  
AAGAAAAAGAGAAAACCTTCTTACAGTAGAATGCCAACTAATAAATGTAGAAGGAATG  
ATGGAATTAGAAAATCACCATTTGAAAAAATCACAGTAACAAAAGAATCACAAAAAAT  
CATCAAGAGATCCTAAAACCAAGGAAGGAAAGCGTAACGACCAACACCATAAGTACAGAA  
TCGAAAAAGAGACCCCAAAAAAAGCATAAAATACAAAGGAAAAAAGTAAATTTAC  
AGTAGAGAAACCTGACAAACACCACCTTAACCAAGTAATCAAAGTTAACATCAACAATAA  
TAAGACAAATCGACAGCATATACCCCTGAGATAAGACACGAAAAAACAACAACAAC  
TCCGGGGGCATCCCCCAAAAAACCAAAAACCAACCAACAATAAAAAAACATCAGAC  
AAACCCAACTGAGGGACATTCTACAAAATAACTGACCAGTACTCCTCAAAACTGTCAAG  
GTCCTCAAAAACAAAGAAAGACTGAGAACTGTCACAGACAAAAGGAGACTAAAGAGACA  
TGACAACATAAATGCAACGCGGGATCCTGGATGGGATCCTGGAACAGAATTTTTTTTGTCTA  
TAAAAGAAAATAAGGGAAAAACTAACGAAATCTGAATAAAGTATGGACATTAGATAATAA  
TAATGTATCAATATTAATTTACTAATTGTGACAAATGGACAGATGTTAAGAAAGAGAATG  
TCCTTGTTTTTAGGAAATACACACTGAAGTATTTAGGGGTAAAGGGGCATCATGTCTGCA  
ACTTACTCTCAAATGGTTTCAGAAAAAATAATGTATATGAAAACAGAGAATGATAAAGC  
AAATGTGGCAAAATGTAAACATTAATATATACAGATGAAGGGTATACGGGAACTCTCTG  
TACTATTTCTTCAAATTTTCTGTAAATATAAACTAAGCTTAAAAAAGTAAAT  
TTAAAAATAAAAA

LTR consensus sequences:

>1\_cons

GTTTTTGGCAACCAAGGAAGGGGGTTCGAGGTAGAGAAAATGCTAGGCATTCAAAAATCTC  
CTTTTCTTTTTGCTACAAACAGGAGATAACCTCACGCTCTAAGCTCAAAAATACTTTTCA

ACCTCGGACCAATGGGGCAAAGCGCCGAAACATGGAAGCAACTCCAGGTTTCTGGCCGTG  
GCCAGTGAAACTAAGGAGTTTCCATGTAGAGAAAACCAACCACCACCCCATTTACCC  
AGGGGTCAGGAGTCTTTTCTGTACTATTCCTCTCTTTTTTCGAGTTCAACCTGTTTCAA  
CACAAGAGAGGCGGCTTCTCTACCTCCATGAAAATGGAAGGCAAGTAGCTGGGGTCAC  
TCCCAGTACTGCCTGAAGGCCTAGGAATGAATGGGAATAATTGCCCTGCCCCGAAGGGG  
GAATGAAACTTTTATTTTTTTATCTTTTCCGAGTGTGGTCCCTGATCCCTACATGCGGCA  
CAGCTCAGAGCAAACCTCACACGTGTTTCAGGAGACTTAAACCTTCTTTTCTTATGCTAAA  
TTCTTCCCTTATCGTACTCAACTGGCTAAGGAACAAAAAGGCCACCCAGCATCCAGTTC  
CTATCATTACAGTTCATGGCTATAACTAATGGAATGGAAAGCACGGGAAAGCGTGGCGTT  
ATCAAATTATAAGAATGCTAAAAGATGAGGCCTCCATCCTGGAGCTAGAACACACCTCCA  
AAAGGGCACCAGAGGCATAGCACCCAAAGTGGCGATGGTAGCCACCTATCCTCAGAGTCA  
TCCCCAACCCCTAGGAAAGGAACCCAATTGGGGTTTGTCTTGGAGATAATCCCTATCGGAT  
AAAAATTAAGCTTAACTCTCTAGAAAAGAAACCATTAGCACAGAATAAAGATTAATAC  
CCAGCCACTCTTAAAATCTTCTTTTGAAGTTGCAATACTGTGTGGACCCCATATTGTTTG  
GAATCTGAAGTTTACTGTTGAATGAGAAAGCGAAACAACCATATAAGTTGCGAGACCATT  
GGGCTGCCTTAAACAACCAGGGGGCCTGGTGAACGTGTTAGATCTTCTTTTGTACACCTT  
GCCACCAGTTTACTTTTGAGTCTGACCAACTTTATCCTATAAAAAAAGTATTGCCAGGGA  
AGAAAAACGGCCAGCAAAAAGGAAAAAATTGAACTGTTTACTTTGACTACTATACCCGAA  
AGAGTGTACAAATCATATGTGAAATTTGTCTGTGGAGTCGCGAAACACCCACAAAATGTA  
ACAGTTTGGAGGGCTCAGAAAGAGAAAGAAAATGTTATGTTAGATTTTAGCTATAGAGAC  
AAAGCTGTAAGAGCTTTGGCTATGAATGTATACGGCTCGAAGCAAATTTATATGTAAG  
CACCCATATTGTTTTACGTTCTCACTGCCACAACCTGCATCCAAGTAAAAAAGCAAACAC  
AAATCAAGTAAATAAGTCAGAGCAATTTTCAAGTTCACATGACTATAAGTATAACTTTAC  
TAAACAAGCAGCTTTATAAATTATTGAGGAAATAAAAAATAGAAATGCCTTCAGAATTGCT  
AGCATAACGTTTTGATCTAAGTTTTAGTTTTAGTCTCTGCTACTTATGAAGAAGGATCATGG  
TTTGGCATAGAAAGTTATAAACTATAAACCCAGCCAAAACAAAATGATCTTGGTTTGCC  
TGCCCTTTTTTTTTTTTTGACAAATAGTAGCAATTTAACATGAACAATAAAGTCTGAGCTG  
TTGGCCAAAATATCTATGTACTTAACTTTGAGGCTCTTACTTAGGTTTGGGTGAGCACCT  
GATGTTCACTGGCTATTA AAAATGTGGTTAAAAAGGAAAAAAGA ACTTAAAAGAATAGT  
GTAAAAGAGCTCAATGGACAAAAGTAAGCTAAATAAAAAAGGCTAAAATCCAAAATGTAAG  
TACATATGAATGCGATAAATGTTTTAGGTAACTTTTTGTGTA AATTA AAAATCTTAAAAT  
TATATTTGATGCTCATTAAATATCTGGGTAATCTTCCATAAGGGAAGGGTTGCAATCTGG  
AGAAATACATGTTTGAATTTTGTGCAATGTTGATTCTCTTTAAGTCCCCGTAACGAATCG  
ACCAGAATTTCTCTCTAGACAGAGAAATTAATTCTAAGTTCTCAATAAAGTTTTAAGCCA  
CCAAGGACAAAATTCCAGTTAACACATAATGTACCAAACAAAACGAACCAGAAAAGGCTT  
CGATATTAATGAGAAAAAGAATAATTTTATCTAATTCAGAAGTTATATAAAAAGTTAGTTC  
AAATTACAGAATTA AAAAGGTTATTTATGAAAAAATG TACTAAGAAAAAATAAGTAGGG  
GAGAAAAATGTGGAAAAAGTTTAAATAATAAAAATATTCTTTAAAACCTGATAAAGAATTG  
GAGACATTTGACTAATTAACATTTTCATAGTTAAAGCTGTTAGTCTTGAATAATAAATA  
AGAAGTATTATATAAGAACCATAAACTCAGAAAATACTTCATACAAACACAGTTAA  
ACAATAAACTTAAATTTAGTCCTGAGCCAACTCAAAGAAAATAAATACACAGATTCCCCC  
CAGGTTTACTGTTTTGCGTGGATAGTGATAGAACTTGGCTGTTTGGAGAGAGATGATTTAG  
GAAAACACCTTAAATTGACTTTTTTGTGATTGAACAGGATGTATAATGATATTGGTAGACTTG  
AGGAAACAGAAAAGCGTACAAGGGGTGAAATGTGAGCCATGTGAGTTTTTTTTGGGGGACC  
CTAGGTAACACTATAGTCTCCAAGGTAAATTGAGTAGGAAAATTTAGGGTTGGTATCCTG  
TTTAAATTGATTTGCTTCAAATTTTCATTCGTTTGTGTTTATTTCTCCTCTGGCTTTAC  
TTGTGTGTGATCGTATATACATAAAACAATTGATTTTTTTGATAAGTTTTTTAGTTACTA  
GTGGAAGGCTTTTATTTGGTTCTGTGGATAGTTATTTTGTTCCTATGCTATATCATTAC  
TAGCAAGTCATCATTGTTCCATTTATCTGGAATGCAAAGAGACCCTGTTAGGCCCGC  
AAGAAATAATGTAGTACACAACTTTTCTATCCTATTATAAACTAACTTTTTGGAATTTA  
GGCTTCTGATACTTTAAGGAACTTAAAAATATTTTCGTAAAAAGAATTAAGCTCAATAT  
TTCTTTTCTCTGCCTAATTTCTCAAAAATGCATAAACTAGTTATGACTATTCTTAATTC  
ACCCTTATGTGATTCTTTGCATACACAGTCAAGCAGGGATGCTAGGGCAGCTCAGGGAGA  
GAGAACCAAGGGCATTAGGTTTGGAAAGAGAAAGAGTGAAAATTTCTTTCTAAATAGAC  
TCAGTCTCTTTCTCAATTTTATCATACCTCTCTAAATAAATTAATGAGAGAGTAGAAAT  
CACTGTTTAACTCATCTCTAAAACCACAAAACAATGAACCCCAAGAAATTTAAGCGGATA  
TTTAGCTATAAAAAATTTTCATCTGCTACATGTATCCTTCTATATTA AACTGTAGAAAAA

GGGGTACATTAGGATAAAATGCGTGCCTAGGACTCCATAGGCTTGCTGTTCAAGATGGCC  
CAGCAAACCTGGACAGTCATGTCCTTGGGAGCTTGACCTCGTAACCATGTGGCCATGCTTT  
CATACCCCTAGGTCATAAAACAGCTCAACAAAGAGTTCAATAACTAAATTAGGAAATTG  
GAAAACGAAATATCATAAAGCTACTGGGTCTTCTTCTATCTGTCTGTGTAATACGTCTAT  
GTATGTATGTGTTGTGTGTGTAATGTATCACTACTAAAAATACAAAAAAGAACCTAAATT  
GATCTGCATAATAAAAAAAAAAAAAAAAAATTTGAATCAAATATTTTATGAGGAGAGAAGAAAG  
AACGAGTCAAATGCTTTTTCAAGTTAATGAAGTACTTTATTAAGTTCATGTGACTTAAGT  
AATATTTAAGAGAAAGAGACAGCCTAAGGTTAATCCCCAAAGTATTAGAAAAAAGATATC  
AAAAATGTCTTAAAAAATGTAAAAATACATTTTGGTCTAAATTATACAGATCAAATACTT  
TCATACTTATCCCTGCCAAATACTATAAAGGTGTCAAAGGTTGGCTAAAATGTTTTAAGG  
TTATAAACCCAGCCCAAAACAGAATGATCTTTGCTTGTGTTGATTTTTAAAAATTATTA  
TTGATATTGGAATAATGAAAAAAGCTACATCTTGAATTTAGTAAGATTACCATAACTTCT  
AACCTTGTGGCTTTAGGCGATATTTAAATGATGACTATCGCAGTTTTTCATAAAGAATCTA  
GGTAAGCAATTAATAAAATAATTAGGTAAATGTAATGGGATAAATACCTGGAGACAACT  
TGTATAATTTAGAATATAAAGTTATATTAATTAATAATAGATAATTAATTATTTGAG  
TATTTTCCAATAAAAAATATATTGTAGGAAAACATTCTTACTTAAAAAAAAGTGTGTCCTT  
TTTTAAAAAATGGTGAATAAGTTTTGTCTAATTCAAAGCTTATTTAAAGGTTATATATA  
ACAAGGTAAAAGGAACCAGGAAATAAAAAAATGTAAAGAAAGTTATGAAAATAAACAC  
GTCACCTAGCTATGCAAAAAAGCTGAAAAAGAAAAAATCATATGAGAAAGAATCTTAT  
ATGGTAAATTCATAAGAAAGAATCTTATATGGTAAATCTTGTCTAAAATAAAAAAACA  
GGTTGTTAAAAAAGAGAGATGTTTAGGACAAATCAGAAAGTCCAAGCATATTATAAATGG  
TCTGTGTAATCATAAAAAAATTTACAAAAAAGAAATTTAAAAAATTTATATGATTAAG  
TTGGCTATAATTAAGAAATTTATTTATAATAGTCTTTCTAGAGATTGAAGTTTGATAT  
TAAAAATACACTAATACACTAAAAATATGCAAGAAAAACAAAACAGTCTAAAAGTAATGA  
GGCATCCAAAAAATACAAAAGATGTTAATGTGAATCCAAAAATTCAACATCTATTCAAC  
CTCACTGCTTTAAAGCTAGGTATACCCCGTAAGAACACATGAAATAACAACCTCCCCCCC  
AACTCTACTGTCAGCTCCTGTAATTTTTCTCAGGTTCTAACTGATGTTGTGTACTGA  
TGCTGGAAAGGGTCAAACCTAAAGGTCTAAAAGAAATGTTTTCTTCCAATATAACATTCT  
GTACTCAACTTTTCTTGATGTGTCTGAACTGCTCCATGAAACCAAAAAACCACACCTAGA  
ACACTGGAAACACTCTTCTTGTCTAATTAAAAAAACCACAACCTTACCAAGGTTTACATT  
AAAGTTAAAAGTCGACAGCAGTTCCATTATAACAGACAAAGGAACCTAGTGAAGAAAGA  
AAGCCTTGGAGAAAAGGCCAGGGTCCCTGTGACGGAAGTGGCCAAAGAAAAAGATTTTA  
TGTTTTATCAAAAAAAGTTTGGTGTTTTTAGGAACAGGTTATTAAGAAGCAAAGGAAAC  
TGAATTTTTAATTGTGCAAAAAAGGGTAAAAGCATCCATGTATCTTTCTGTATTGCTTTT  
AAAGTCCTTATTGTTTTAAGTTAGAGAACATAAAGCTAAAGGTTTAAACAGGTCGTGGAA  
GAATTGTAAACAATTAATCTTGAAAAAATTAAGCCACATCTTCAAGGCCGTAGAAGAT  
GCCAATCAAAATAAACTGCATTCCTGAGACACAGGAAATTAAGCTATTCAACTCCTCAA  
GGCCCAGGGACTATCCAGAAGAGGTGGGTATGTGTGAACATGATGTCTAATATCCAAAGA  
TAAAGTTATTTATGGTTTCTCTGTAAATTGAACATTGAAAGTTTCTCTATAAATTAATC  
ATTAATAATTAAGCACACTGATGCAAGACCAGCATATGGGCCCTGTGTCAGATTAACA  
AGGTTTTCTTAAAGCACTAATCTGCTCTTTAATAAAAAATTTATAAAGGGTTATAAAATGT  
TTACGAAAATCTCATCCTATGGTCAAACCTGATTAAGATCGGAAAGATTAATAATATAAGAG  
ATTATTTAAAAAATATTTCTGAGATTGACATTAATAGTACACTAATGCAAGGGTGAAATG  
TGGCTTTCTCTCCTGAACAAGATTTTCAAACAAAATTAAGAGACACCAAAAGATTTTTAT  
TAGCCTTTTGAATAAACTACCAACAAAAAAGAAGGGAAAGACAAGAGACAGATTGTTTG  
GAAAATAAGTCTTCCCTCTCTCAAAGAATGAAGGTTTTTGGCCTTTAAAAAAGTTTTCC  
TGGAGCAATCATTTTGGCTAAATGAATGACTTATTTAATGTAACCTGCAATTCTATTTT  
ATAATATCAAGTGTTTTAAACCTATAACATATCTCCTCAGTCTCCCCAAACCTTCAGTAC  
AGTCTATGTCTTTCTGACCAAAAATTGTCTTTTTCAGATATCAGGCTTCTTAGAAGCATCAG  
AAGGCCCCACGAAGAACCATCCAAAAGAGAGGTAAAAAGGATTATTTGACACATTTAGTT  
ACATTTCTTCCCTGCCAGAAAGCATTGACAAAAACGAAAAATGTTAATCTTCTTTAGGT  
TATATTTAATAAATAAGTTATTGATATATGTTCCAAAATTGTATGGGATTTCTAAAATT  
CTAAGATGTCTGAGTATATATTATCAATCATAATTAAGGTTATTATGTTAAATTATTGTA  
AACCACAGAAATAACAAAACCTTTTGATCTGTGTGAGTTGTGTTTTTAACTGTAACCTATT  
CTAAGAATTTTCCACAGTTATTCACAGACAATTGTTGTATTGTTTGAACCCGTTTCAAAG  
ATAGTTTATAATAAGCTATGGTGTCTTTTAGGAAGTTGATTAAAGGATGGAAAGAACTCA  
AAAAAGGGGGCTGAGATCCACACAAGGTCTCGGACAACGCGGTGGGAGATTGTAACATC

AGAAAAGAGAAAACACCTACAGGACCCCGAGAAGATCCAAGTGACTCAGAAAAATGCCTA  
AACCAAACCTCCAGCAAAAAGAAGCAGGAATTAAGAGCCAGCCCGTGAAGGTGACCAGGAGA  
GAACATGAAAAAATCTTTTTGACTTTTTGCTTAAAACATTGCTGATCCTTTGTTTTGTTT  
TTCAGAGTCAAGAAAACTTTTATTTTGAACATTTACAGCCTTTAACAATTGAGTAAAGT  
ATACTCCTATGAACAAAATTTGGAGCATGTTTGTCTCTCTGCCTGGTTCCCTGAGAAT  
TCGCCCTGATTAATACTTTGTTACGTCATTAAGAGAGCCAGTATTGGGGAGGGCACTGG  
GAACCCACTTCCTGCAATATAGTGACTTGCATAAAAAGACAATAAGAATCTAATTTTCATT  
TGCCACAGGACGTGATATGAGGGATTGAAGATTTGACTGGCCAATTTTATTTAGACCTTA  
AAGGGAAGGGTTTGTCTTCTGTAAGGAATCAATCTTGACATGTAGAGCCAATAAAAAGCC  
CTATGGAAAAACTGGCCTCATAACCTTATATACACAGTCCCTGTACAAGCTTTCTGACCA  
GTAATCAGCAAAGAATGTCACCTTTCTGACAGGCCAGGAGCCCCAAGTTTATCTTGGGAC  
CTCAAGAGGAGAGGGAATTCACCCAACCTCACAGGTATTTGAGGACACAAACCCATGGCTG  
GGCTCAGCTTTAAAAAAGTCTTATCTGAGATTCCTTCTATGGAAAAAAGTTCCATCAAAG  
CCAATCTAAAAAGACCATATATAAGAAATAATTATTTCTTGCTGCACCTTTATGCAAATAAT  
CAGGCCAAGTATAATAAGAATAAAAACCTATTTTACAAACAAATCAGTCCTAACATGATTT  
TTTTTTACAAAAATCAGGAAACTGGAGAGAGAAAAATTATGTTTCAAAAACTATAATACA  
ACTGTCATTAGATTCTAAACCCAGAAGTTGTTTTTAAGTTTTTGCCTACATGTTAGACTA  
ACCCTGCTTGTTCCTGTGAACCAACCAGCAATCTCCGGCTGCAACAAAAA  
GGGATCGGTACTGCGGACATTTGCGTTACAATTTTAATTCTCAACAATTTTCGCGAACAC  
CCTACAAAGGGATAATCTTATACTCTGATAGATAAAAGATGAAGACCCAGAATAAATAAG  
AAACCCAAAACCTACAGACGTCCCTCAGAAAAAGTAAGAAAAAGAAAACTAACCAAAGC  
CAAGTAACCTGAACACAAGTCTTAAAAAGAATAATTACAGCAACCAGTTATCTGGGTATG  
TCACAAGACATCCTCTTCACTCCCCTATAAGAGAAGGACATAATTTATCAGTTTTACCTT  
CCATTTTAGAGGATGAGAAAGAAAACATGCAACCACAAAAACAAACCCCTCACAACAAA  
CTCAATACCAAACCTCTACGGCAAAGCCCAAAGAAAACAAAACAGAATCTAAATTATCCTGA  
GAAAGATACTAAAAATAGAACTGGCCAAAGAGAAAAGCCACAATGCAGATGAGCACA  
GCCAAGGCCTACCGAGGAAAACAAGAGGCCAGTATCATAGAGAAACAGCAGGAGAGGATC  
CACAACAAAAATAAGGTCTAGGGAATTCAAGGCTACTGACAGCAGGGGAGATAGGGCATA  
AGTGAGTAGACCGCATAACTACCACACCCTAGGACCCACTGCATCATAGGTGCAAGCCGC  
TTTGACACCCATGGTGGCACCTGCCAAGGTCACGGGGACCCTGGAAAACAAAAACACAAG  
AAGAAAATAGTACCATCCTCCTTATATACATAACAATAGTCTACCCGGGGCATCGGATAC  
GAACAAAAGGGAAACAGGCATGCATGCAGCCATCTTTATAATGACAAATAGCCTCTCTTC  
TAATGTATGGACAACAGACGAATGAAACCTGTAACCACAAGAACACTGTAAAAAAGAAC  
GGCATGCTGACACCAGAAACCCAAAATCACGCTAAACGAAAAGGGAACCTAAGGCTCCAAA  
CTATGGGAAAACCCCCACAGAAGCAACAACCAAAACCAAAAAGAGGGAATTTTCCAAAACA  
AAAACGGGGTAGGCCATTGCTTAGGACTGAGCTCAGGCAAAAGACCTCATCAGATCAAAC  
CAAAACCAAAACAGACCGGCAAAAGCTAAAACCTTTAAGGCAATACATATGGATCTTATTA  
GTTTTATTTGTTTTTGTCTTCCAGAAAGCCCTATAACAAACATCCCCAGAGCAAAA  
GGATAAAACGCCTGAAGATCCTTGAGCTGTCTTACCCCTCCCATTTGTTTCATCATAAGA  
CATGAAATCTAATAACCTGATTAATATCATCTTTCCTAAAGACCATCAAACCTCAAATGA  
TCATGCAACAAAAGCCCCAGACAATGCCACCTGAAGACACCAACCCTGGTCATTAAGGAG  
CTACCCTGTCTCCATTAGAGAGAGCAGCTAAGAGAGATGTGACCCACAATTCCAAGAAA  
CAACACCCCCAGCCAGCAGGAAGCAGTTAAAAAAGGAGGACAGCCAGACCCTTTTCTTA  
TAAAAAGGTCAATAGAAATGACATCTGAAAAGGGGGGAA

>2\_cons

AAATTGCAGCGGCGAGCAGGGTCCGGCCGACAGAACCATGCCATAATGCCGAAGCCAGA  
AGGGCGATACCCCCAGTACCCGGCCACAGCTTGAGCCTAAGCTGCGCGCTCCATGGGA  
TTCAGACGTGTCTGGGCATGTCTGGAGAACCGGCCGTGAGAAGTGGGGCGGGCCCGCC  
CACCTCCCCGCAGATGTGAGGGGTCTATTTGCATGAGCGAATGAGGAGAGAAAGGAATG  
TGAAACAGGAGAGAGCAGCGCACGATCCCCTGGTTGTTGGCCTCCTAAGACAGAGCTGC  
AGACCCGGGCATCCTTGGGGAAGAGGGAGCCGGGAACAGCGAGCTCAGATAGACCCACAC  
GCATTGGTGCAGGAACAAGGGCCAGGCGTCGGAGGGTCCCTGAGGCATCTGAGAGGGAA  
CCAACACGTACAATGCGGGCCTTCACAGAAAGGGCGATTCACGAGCTAAAGGACAGATAT  
GGAGAGAGGACAGCACGGTGTCTGCTACTCCACAGGATTGGTCTCCCGTTTTGATAAAGG  
GACATTGCAGATCCAAATGTCTCAGGCTGAATGGCGCAGTTAGGAATGGGGCCAAGGTCT  
CCAGCCCCACTGAAGCCACCGCCCCTATACTCCTGTCAAGATAAAAGGAGAAAAGAA

AGGAATTCAGACTTTTCTGGGGTTCTTGGATACTGGAGCCCACATGACAACATGTCTGAG  
TCCCCTTAGGGGAAAAATTAACCTGATGACATCGGGAGGTTTGGGGACAAACATGGTGAC  
CCATGGTGCTTATTTGCTTATGGTGCTTATCTGCTTGTGGGTGGGGCCCTTTGGGCCATT  
TCGGGTGCCAGTGACCATGGTTCACCGCTGAGTGCATTATAGGCATTGACATTTTGGC  
TGCTTGTGGCACAGAACATCACCGCTGCCTGAGGGGGTATGCCCCCTCACAGCTAAGAAT  
TCGAGCCATAACAGCGCCGCAGACCCACAACCTGCCTGCCCCCTCAGCCTGCCAACTCCCT  
ATGGGTTATTCAACAAAAGCAGGACTGCAATCAAAGCACAGGCAAAAACCTAAAAGAGCTG  
ATTTCAGGACTTGCTACAGATAAAAATGTTACAAACCACCCTGTCACAATGTAACAGCCC  
AGCCCGGCCGGGCGAAAAACCCCGTGGGACATTGAGGACAACAATGGACTGTCGCAAGCA  
GCCTGCTGCTTTGGCCCCCTCCACACAAGCGGGCGCCGACATCACACAGTAATTGAACAC  
ATCATGGAGGCTTCCAACCAATAGTGTGACACAGTTATTGATCTGGCTAATGGATTCTTC  
TCAAGCCAGGGGGGGAGAGGGGCAGAGATCAATTTGTATTACATGGCAAAGTATACAA  
TATACATTTACAGTGCTGCCACAGGAGTATTTGAACTCACCTGCCATATGCCACCAGTGG  
GTAGGATGGGATTTCCGCACTGTGCTTTTGCCTAAAGTGGTCATGTGCATTACATTACATA  
GGTGACATCCTCAGCGCGGCCCTGCAGGAGCCCATCACAAAAACACCCTGGACGCCATG  
AACACAAGCACGAGACAAACAGACTGGGAAGTTAACCTAACAGTCTGGGATCAGCCAA  
ACTGGTGACCTTTTTCCCCCCTTTCGGCGCGGAGCCAAAGAGGCAGCGCAGCTTCGGT  
CAAGCAAAAATTGTTCCGCCCCGGGGCACCCACTAATAAAAAGGAGACCGGACAGCGGGT  
CGGCGCCCCGGGGGACTGCAGACAGCATATACCTCACCTGGGTGTTCTTTTGGCCCCCTT  
AGTCAAGGTGACCAACAAAGCCGCAACTTTGAATGGGGCCCCTGGGAGCAGCAGGCCCC  
GGAACCCAGCCAACAGACAGCGGCCAGCGACCAATCGGCACGGACTTCAGCGCCTTCGC  
GGACGGCAAAACCCATGGAATCACAGGTGTCCGCAACCTCCATGCATGCTGGCCGGCAGC  
AGTGGCAACGGAAAACCTGCCACTGGGGTGCACCAGCCTCTCAGATTTTGGACACATAAGT  
TGCTGAGGCAGCCACCAGATATACCTCTTTTGAATGGCAACTCCTTGCTTGCTATTGGG  
CACTGGTGGAGACTGAGCATCTTACGGCCGGAGCGCCACGTGTGACGCTGCAACCCGAAA  
CGCCAGTCTCACAGCGCGCTGCCAACCCACCAGCAAAACAGGACAGGCTCAACAGA  
GCTCAATTATCAAATGGAAATGGTACATTCAAGATCGGGCCCAGCCAGGACCCCAAGGGA  
CCAGCGGGCTCCATGAACAAATGGCTAGCTTACCAGAAGGGACCAAGCGACCCGTAGGGG  
ATGCTTTGGCTCCTCTGTGGCTACCTGGGGCCCAAGATTCAGAGACATGCCTACCGACG  
GTATGGCATGGGGTTTACTGACGGCTCTGCGAAACAACAAGCAAGAGGCTCCACCGGGCT  
GTGGACATCAACCAGCCAGTGGATGGCCATCTTTTACTGAGACTGGACATGGACGTTCT  
GCCAATGGGCCAAACTACATGCAGTGGTGTATGGCCATGCAGGCCGCCCTACCACCATA  
TCTTGCTACATTTTACTGACTCATGGGCCATTGCCAACAGCCTAGCCATCTGGTCAGGA  
GAATGGCAACTGAGTACTGGACTATTAAGGATCCCCTGTGTGGGGACAAGGACTATGG  
CAACAGCTTGCTGCCTGGAAGGGACAAATATATGTCATCATGTGGATGCTGGGACTACC  
ATGGCCACCCTTGAGAGGAATTTATGTCATGTTTTTGGATACCCCATGGGACTTCACTCT  
GACCAAGGAACATCCTTCACTGCCCAAGCAACATGACAATGGGCACACTCTCATGGAACA  
CGATGGACTTTCCATGCACCCTGTCATCCACAGGCCAATGGAGCTATTGAGCGCCAGAGC  
GGCCGAGGAAGAGCGCAACTGAAGAAAGGACATCAAGACGACCTGCTAGCGGAGAGGAAC  
CACCAACTCAAGAGGCCAATATCGACACAAAACACTGCACTCCAATGCAAGGGAAACACG  
GCACTGCAGCACACGTCGAGAAACACTGACCGCGGAGAGGAGCGAGACACACCAGGCAGC  
CGCCGAACTAGGCTACCCCTCAAAAAACCAATCTCAATCTTCCCAACCATCCACTTGCC  
TGTGTACCTCAACAGTCCACGACCCAGGACAAGAACGCGGGACAGGCCACCAAAGGACCC  
CACAAAAAGACCCCAACACAAAACAGGGAGGAAACACGCCCCCTGGGGACTCCCCTGCGA  
GCAGATCCCACAATAATTGGGGAAGAGACCAAGCCTTTGCAGGGGGGCTGGGGTCTGCGT  
GACGCAAGTGGGGGCAAGAACCCTGAAATTTGGGTTATTGCGGTTAATGTGACCCCAT  
TGTAAGTTTGTACAGGACGCCACCGCCCTCCCCGGACCCGACCCACAGGCTGAAAGCG  
TGAGGCCTGGGGAAACACCAAGGGCAATGGTGCCACAGAGGTAGTAGCCTCAGGACAGG  
GACAGACAGACTGGGTCGCTACACCAACTCAGCCCAACCCCTATCTGATAGGTAGGGAAC  
ACCTGACACCCTGGAAGAGAGCCGGCACCGGTGCCGGCCTGTCAGGCTGCTCCACCCGCA  
GCAACATGCAAGCAGCCGACTGGAAACAGAACCCCGACCCACGCTCGCCCAACACCC  
CCACCGCGCCGAACCTGACAAAACCTGGATCAGCCATCCAGACAGCAGCTAGCACAA  
ACCAAACGCACCCTCCCAGCCCGACAGGGCGAAACGACACCAGCAGGCCAAAGCAAAAAG  
CGGACAAAAACACAACCCGACACAGGGGTTACAGAGCAAAAAACCGACACCCCAACATC  
CCAGGAAACCGGAGGTGCCCTGTTTTAACTTAACTGACTTAAAGGTGGCAAAATGTCACGA  
CCACAACATAACAAAACCTTGGTGGGCTGGTACTTTGACGCACCACACTCCTTTGATTACA  
TGGACGAGAAGTGTCCAGTGGCGACGACGAAAACAAGGACCGGACTATTGCTAGCCCTC

TGTGTAGGGGCTTCATGAACAATATTGTATGGGGAAAACCTGAGCTCATGCAACTATGCCA  
TCAATGAGACTTGGCTGGTGAATGCCAATGCCTCCATACCCATGAATGGGTCCTGAACA  
ATAAAATGGGAAAGGGTGTGCTGTGTGCACCCGAGGGCTACATCTTTCTCTGTGGGCGGT  
CCGGGAGTGACCCAAATACGGGATGGGCAATGTCATGCCTGGAAAGCTGGCGGATGGTGG  
GATCCTGCACGTTGGGCGTGCTGGGGGTGCCCTGGATATCACCCCTGGGAATGAGATGC  
ACCATTGGGCCAGCAGCCTAAAGCTGTACACCAGGCTTACTAGGGACCTGCCAGGAGGTG  
TAACTGACTCTGGGTTTATGTCCTTTATGAGATCTTTGGTACCATACATAGGAGTCAGTG  
CTCATGAAAAAATGATAAGAAACCTGTCCCTGACCATGGCAGATATTGCTTCCCTCCACTG  
CCACTGCCTTGGCAGCCAGCAGACATCCCTCAACTCCCTTGGGAAGGTTGTTTTAGACA  
ACAGAATTGCTCTAGACTTTCTTTTAGCCCAACTGGGAGGAGTGTATGCAATTGCCAACA  
CCTCCTGCTGTACCTGGATAAACACCTCAGGTATCGTAGAAACACAAGTAGAGGAGATCC  
GGAAGCAGGTTCACTGGCTGCAGACAGTGGGGCCACCTGAAGGATCCTTCTTTGACCTCT  
TTAGCAACTTCTTACCTGGATCACTGGGATCCTGGGCTAGGTCACTGCTCCAGGCAGGCC  
TGATCATCCTGCTTGTGGTAGTAGTCCCTCCTGGGCCAGTGAAATGTATTCTGGCTATGG  
CTCAATGATGTTGCACTGAGATTGTGTCAAGGTGCTACATCAATCTGACAAGACAA  
ACCTCTGCCTCCAGATCCGGGGAGGTGCGTGGGCATATGAAATGGACTAGCTTTGCTAAG  
GGGGATATCTGGGTTGGGGG

>3\_cons

CAAAACCGCCCTGTTTGGAGAAAGAAGAAATGGACACCCTGAAAGCACCTCGCGACCCTAA  
CGCAAGACAAAAATACGTCAGGCAGACAGTGAAGGTAGGCCCCACAGAATGGGAAATGG  
CACCTTGGAGTCACTGAGGGTAGCCTCCCTTAGAAGGGTTGAGTTCTTTCAAACAGCCCA  
ACAGCCCACAGCCACAGACAGAAAATCTAAGTCACCCTGAGGTATGAAAATAAAAAAAAA  
AAACTCAGGCGGAACTTGACCCGGGGGGCTGGCCCAAGACAAGCAAACAAGGAAGGTCTC  
TCCAGATCAGCAAGGCTGCACAGCCGGGGCTAGCCAGAAGCCTTTTGTCTTTGTGTA  
ATTAACATGCCACAGGGGAAAATTCCCTCCCTTTTCAGACACATGCATAGTGGGCTCC  
AAAGGAACATAAACAAATATGGAGGAGCAATACCAACACCAAATAAGGGTCATACAAAC  
AAGAGAAGCAGCGCTTTGTGCCAACCGAGAGGCATCCTTACCAGAGCATAACAAAAATGG  
AGTGGAACAATCCTCCCTCAAAGAAACACTGCGCCTTATGCAATAAGAAGAGTCGCTCAG  
CAAGCGCCACACAGCACTAGAGGACCTTAAAATCTGCAGCACTAGTATTCAACAGCGACA  
CCATGTGCCAATGAAGAAAAATGGCAATAGTGGACTGATCCCACTTCCGGGCTCCCTCT  
CTGCTGCAGAGAGCTTTCCTCTCTCAAGAAGGAGACAAGACATGTTTATAAATAGCCAAA  
ATGCAGCTTAGTAACTTTCACTCCAACCTCACACTTCGAGGGAGAGAGAAATAAACGAA  
GCAAGAAGCCTCAAGAGGAAGAGATGTGTCTCCATGCTCATTAAATAACGCCGGCACTATA  
TTTCTGATTGTAAGAACAAGAACTCCGGATAATACCTCACAGAAGAAATTTAGAAAC  
AGCTAATTTGAGGAGCATTGGCGACACCTAACCTGTAACA

>4\_cons

TGTGGCAAGGATAATATTTTGGAGATATTAATTTATGTTTTGTTCTCCTCTTGGTAAACCT  
GTATTTTCCCCTTCCCACCTTCCCCATTATGGCCCAAGCAGGTAGCCAGGCCCTTGAT  
GACTCATTGCCTCGGGGAGGTATGTGCCGAGGGGAAGTCCATGCCAGAAAAGCAGAAA  
GAATGCTGCATGAAGTCAACAACCTTCTCTGGCTTTTGTTTTCCAAAAGCCTAAGCCCAT  
TGGAGGGAGTATGCTAGGAGACCTGGGGAGGAGGGGGAAGAGGTTAGAGGCTGAGAGGG  
AACCTGAGAGGAGGCCGGGCTCCCTCCCCCAGACAAAGAGGAAGAGATTCCCCTGGGC  
TGAAACCCAGGAAGGAAGGCGGGACAGGAGCTTTAAAAAAACCCAGATAGGTCTGGGGGA  
ATCCTGGGAGGAGAGCAGGGAAGGGGACCTGTGCCCTGCTTCCCGGCAGCGCAGCCCGG  
GAGGCGGCAAGACTCTCAGAGAGGCCCTGCATTTGGCTCGGCGCCACGTCCAACATGGCG  
CGAGAGCGGTAGAGCAGCAATGGCTGCGTGGTGTGTCTAGGCAGACGGGACCGAGGGCAG  
CCCACCGGGGCTCCAGCCTCCCATGGCCTGCGTGTGGCACGGAGGGGAGCCAAAGATT  
CCCGAGTGCCCAAGCATGGCACGGAAGGTGGGCTCAGAGAGAGCCGAGGCAGAAAGCAG  
CAGAGAGCCGCCGGCCCTGAAGGGGCCATCACAGAGGGGCAGAAGAGGGCCTACATGCCT  
GGGAGAAGGAGCCCGGCGACCCGAGGCCAAGAGGGAAATGGCCACCACCGCGGACACCAG  
CAGGGAGCGGAGGCGCATCACGCCAAGGACCAGACGGGACAAGGGACATCTCAGCGGTC  
ACCAGCACAGGGAGCAGACCAGACCAGCCACTCCGCAGCAGAGACCAGTGTGGATCCCGA  
TGACGCCACGAGGACCAGAGGCCCCCCCGATCCCCCAATGCCACGGCAACTGCGTAAG  
CCCACACTACCCCGGAACCTTGGCACAACCCTGGGGAGAAGGGGAGGAGAGGGGAGGGG  
AGAATCCTGAATTGACTGAGTATTTCCCCAAAAATGACTGAGTCATAAAAAAGAGACTAA

TTTACCTAAAAGAGACTGTTTAAATCACTGGAATTGACTGAGTTTACCTGGAAGTGACAA  
GATTAAGTCGTTCTCACGCTTCCC GCCAGCCCGGGCGGGATGGGGGCTCGGAATCAGAAAT  
TAAGTTGAGTTATAGAAAATAAAGAAATGTTACATTTTCCTTGCACACCTGAGTTTGTGG  
CGAGTAAGATTGCATACCCGCTACA

SVA consensus sequences:

>1\_cons

CTCTCCCTCTCCCTCTCCCTCTCCCTCTCCCATGGTCTCCCTCTCCCTGTCCCCTCTTT  
CCACGGTCTCCCTCTGATGCCGAGCCGAAGCTGGACTGTACTGCCGCCATCTCGGCTCAC  
TGCAACCTCCCTGCCTGATTCTCCTGCCTCAGCCTGCCGAGTGCCTGCGATTGCAGGCGC  
GCGCCGCCACGCCTGACTGGTTTTTCGTATTTTGTGGTGGAGACGGGGTTTTCGCTGTGTT  
GGCCGGGCTGGTCTCCAGCTCCTAACCGCGAGTGATCTGCCAGCCTCGGCCTCCCGAGGT  
GCCGGGATTGCAGACGGAGTCTCGTTCACTCAGTGCTCAATGTTGCCAGGCTGGAGTGC  
AGTGGCGTGATCTCGGCTCGCTACAACCTCCACCTCCAGCCGCCTGCCTTGGCCTCCCA  
AAGTGCCGAGATTGCAGCCTCTGCCCGGCCGCCACCCCGTCTGGGAAGTGAGGAGCGTCT  
CTGCCTGGCCGCCATCGTCTGGGATGTGAGGAGCCCCCTCTGCCCGGCCGCCATCGTCT  
GGGAAGTGAGGAGCGCCTCTGCCCGGCCGCCATCCCGTCTAGGAAGTGAGGAGCGTCTCT  
GCCCGGCCGCCATCGTCTGAGATGTGGGGAGCGCCTCTGCCTGGCAACCGCTCCATCTG  
AGAAGTGAGGAGCCCCCTCCGCCCGGCAGCCGCCCTGTCTGAGAAGTGAGGAGCCCCCTCCG  
CCCAGCAGCCACCTGGTCCGGGAGGGAGGTGGGGGGGTTCAGCCCCCGGCCGCCAGCCG  
CCCCGTCCGGGAGGGAGGTGGGGGGGTTCAGCCCCAGCCCGGCCGCCGCCCGTCTGGG  
ATGTGAGGAGCGCCTCTGCCCGGCCGCCCTACTGGGAAGTGAGGAGCCACTTTGCCCGG  
CCAGCCACTCTGTCCGGGAGGGAGGTGGGGGGGTTCAGCCCCCGGCCGCCAGCCGCCGCC  
GTCTGGGAGGGAGGTGGGGGGGTTCAGCCCCCGGCCGCCAGCCGCCCGTCCGGGAGGG  
AGGTGGGGGGGTTCAGCCCCCGGCCGCCAGCCGCCCGTCCGGGAGGTGAGGGGCGCCT  
CTGCCCGGCCGCCCTACTGGGAAGTGAGGAGCCCCCTCTGCCCGGCCACCACCCCGTCTG  
GGAGGTGTACCCAACAGCTCATTGAGAACGGGCCATGATGACGATGGCGGTTTTGTCTGAA  
TAGAAAAGGGGAAATGTGGGGAAAAGATAGAGAAATCAGATTGTTGCTGTGTCTGTGTA  
GAAAGAAGTAGACATAGGAGACTTTCCATTTTGTCTGTACTAAGAAAATTCTTCTGCC  
TTGGGATGCTGTTGATCTATGACCTTACCCCCAACCCCGTCTCTGAAACATGTGCTG  
TGTCCACTCAGGGTTAAATGGATTAAGGGCGGTGCAAGATGTGCTTTGTAAACAGATGC  
TTGAAGGCAGCATGCTCGTTAAGAGTCATCACCCTCCCTAATCTCAAGTACCCAGGGAC  
ACAAACACTGCGGAAGGCCGAGGGTCTCTGCCTAGGAAAACCAGAGACCTTTGTTTAC  
TTGTTTATCTGCTGACCTTCCCTCCACTATTGTCTATGACCCTGCCAAATCCCCCTCTG  
CGAGAAACACCCAAGAATGATCAATAAAAAAAAAAAAAAAAAAAAAA

>2\_cons

CTCTCCCTCTCCCTCTCCCTCTCCCTCTCCCTCTCCCTGTCCCCTCTTTCCACG  
GTCTCCCTCTGATGCCGAGCCGAAGCTGGACTGTACTGCTGCCATCTCGGCTCACTGCAA  
CCTCCCTGCCTGATTCTCCTGCCTCAGCCTGCCGAGTGCCTGCGATTGCAGGCGCGCGCC  
GCCACGCCTGACTGGTTTTTCGTATTTTGTGGTGGAGACGGGGTTTTCGCTGTGTTGGCCG  
GGCTGGTCTCCAGCTCCTAACCGCGAGTGATCCGCCAGCCTCGGCCTCCCGAGGTGCCGG  
GATTGCAGACGGAGTCTCGTTCACTCAGTGCTCAATGGTGCCAGGCTGGAGTGCAGTGG  
CGTGATCTCGGCTCGCTACAACCTACACCTCCAGCCGCCTGCCTTGGCCTCCCAAAGTG  
CCGAGATTGCAGCCTCTGCCCGGCCGCCACCCCGTCTGGGAAGTGAGGAGCGTCTCTGCC  
TGGCCGCCATCGTCTGGGATGTGAGGAGCCCCCTCTGCCCGGCCGCCAGTCTGGGAAGT  
GAGGAGCGCCTCTGCCCGGCCGCCATCCCGTCTAGGAAGTGAGGAGCGTCTCTGCCCGGC  
CGCCATCGTCTGAGATGTGGGGAGCGCCTCTGCCCGGCCGCCCGTCTGGGATGTGAGG  
AGCGCCTCTGCCCGGCCAGCCGCCCGTCTGGGAGGTGGGGGGGTTCAGCCCCCGGCCGG  
CCAGCCGCCCGTCCGGGAGGAGGTGGGGGGGTTCAGCCCCCGGCCGCCAGCCGCCCGG  
TCCGGGAGGTGAGGGGCGCCTCTGCCCGGCCGCCCTACTGGGAAGTGAGGAGCCCCCTCT  
GCCCGGCCACCACCCCGTCTGGGAGGTGTACCCAACAGCTCATTGAGAACGGGCCATGAT  
GACAATGGCGGTTTTGTGGAATAGAAAGGCGGGAAAGGTGGGGAAAAGATTGAGAAATCG  
GATGGTTGCCGTGTCTGTGTAGAAAGAAGTAGACATGGGAGACTTTTCATTTTGTCTGT  
ACTAAGAAAATTCTTCTGCCTTGGGATCCTGTTGATCTGTGACCTTACCCCCAACCCCTG



TGCTCTCTGAAACATGTGCTGTGTCCACTCAGGGTTAAATGGATTAAGGGCGGTGCAAGA  
TGTGCTTTGTAAACAGATGCTTGAAGGCAGCATGCTCGTTAAGAGTCATCACCCTCC  
TAATCTCAAGTACCCAGGGACACAAACACTGCGGAAGGCCGAGGGTCTCTGCCTAGGA  
AAACCAGAGACCTTTGTTCACTTGTATCTGCTGACCTTCCCTCCACTATTGTCCTATG  
ACCCTGCCAAATCCCCCTCTGCGAGAAACACCCAAGAATGATCAATAAAAAAAAAAAAAA  
AAAAAA

ERV consensus sequences:

>1\_cons

AATTTGGTGAGCCAGCCAGGAGCCGCTGGGACGGTGATGCATTCAGCGGCCAGCGGCTTG  
TGACGAGACAGTCTTCAGGAGACTCCCAGCAGCTGCTGGGTGAGATTATCCAGGGGACTC  
TCCTGAGGGCTGTCCCTTGGACAAAACCACACATCCCTCTCACTACTGGGGAAGAACGGA  
GGTCAGGAACGGACGTGCTCAAGGGGTGAGTAAACTGAACCTAATAAGGGACTTATTAT  
TTTCTATCTGGGCTTGTAAAGCCATTTGTCCGGTACCACCAGGGAGACAATAGGGCTCA  
TTTGTACACCCGCTTTGCGTTTGGTTAAAATCAGGTTTTGAGTTGGTTTTGAATCTGTTA  
TGTCAGCAAGACACCCAGGAGGTAGACCCATAGGTCCGGAATCATCCTACGTTTGTAAAT  
GGATCGCTCTAGGGATTGCGGTACCAGATAGGACCAACGGTGAGAAGACAAGATCCAGTG  
ACAGCAGGGACGACGGTGGACAGAAACCAGTCGAGAGGGGAAAAAGTAGCGATAGCAGCTG  
AGAGCGAAAAGGAAAGAAGAAGACCATTACACCGGCCGAGAGACAAGTAGGAATTAGTA  
ACTGCGCGAAGGGGTGTGAATGCACAAACCCTGAGGGATGCACCACTGTTTAAAGGTGGCC  
CCGGCCCAAGACAAGCCCATAGAAACCGACTCAAAAACCTACGCCACACACTGGCCTGA  
CGGCACCGCCCCGAAGGCAAGGAGCACAGACACGCCCGGCACGAGCCCCGGAACAACAGC  
GTAAACCCCACTTCAGCGGAGGGAAGGGAGGGACAACCACAACACACACCCCCAGTCCG  
TCCGAAGCCCCCTCCGGCTACTGCCCTGGCAACAGGCAGGGGGCCCCCCCCATGCCACGC  
CAGCATGAACGCCAACAGGATCCGGCAACGCAGGCGAGCCGACCCGCAGTCATAATGGCA  
CCCTGTGCTTACAGGGCCGACACTCAAACGCGCCAGCGCAGATCCGGACCCCTCGCCGGC  
CGACGTCAGCGGGAACATCCACGAGCCAAGGCGTTACCGGGGTGGGATGAGCCGGCCAGC  
AGGCGGGACTACTCCACCTTTCGACTTCCGAGTGAGGGCGCTGGCAGCCCTCCAGCCCC  
CGCACGCTGTCGGCGCCGACCCGCGTCCGGCGCACTCTGGCCGCGCTAGAGAAGCCCA  
CCCTATTCTGCAGGCCGGTGGATGCCTGCTCCCCCGCCTGGCCTGGCCCAGCACAAAG  
CCCGGGTCCGAGCTCGGAGGAAAGTTTGGGCCGAGCCGGGCGCAGTCACAACAAGAAC  
GGGTATGCACACACTCGTGGTAGCGCTGAAACCCAGACACCAGCCGCCTCGGCCACCTC  
CAGACATTTGGCGCCACCGAGCATGGGCGCCTGGCCAACGAGGGGCGGAGGGTAGCTAGG  
CGCAGGCCTCCAGGCGCACCTTGGCCTGAACAGCCAGGGTGCCATGAACAACACCACAAA  
ACTCACAATAGCGGCATTCGAGTCCAATGCACACATTAAGACTTCTCCGAATTA  
CCATGCTTACACCACCGCTTCTTGAAGGTTTAAAGGAGGCAGACAAATGCCTGGGCCCC  
GAAATCACTATAACGGGCAGCACGAACCTTCCACAGGGGGCCACCCCATGAGGCAGAAAGG  
AAAGTTAACATCAGGGCATCCCAAAGGAATGACACCCAAGCAAAGCTGTACTGCAGCAA  
CTTGAGACGGAATGAAAAAAAAAAGGCAGGGAAAAAACACTCATGGCAACAGCAGGAGGA  
AACCCAAGAGGGAAAGGGACCCGAAAAGCGGAAAGGGAAAAATAGAAAAGGATCAGTGC  
TACTGCTGGGAAACGGGGCATTGGAAAAGGATTGCCCAAAGTTAAGCCAGAAAGAACCA  
AGGCCAATAATGGCAGTTAAGCCCGGAATGAATCTGAGGAAGATTGAGGGTGCCCAAGA  
CTCCAGCAGCTCCAACCTATCTGACATCAAAATTTCCCACCGGGTCCCTTACTAGAC  
ATGGGGACAACCAGGAAAAGGCCACAGGACGCAAAAAAAAAAAGTGGACTTCAGACAAGACA  
AGCAAATAAAAATCAGCTAAGAAAAAAACAACCTAATGAACTATCGCAAAACCCCTGTCA  
GTATAATACGAGTAACCACAAAAGGGAATTCGGCGAAACCCACCTTCCATCCTGAGGAA  
AATATCGCAAAAATAGAATTACTCAAAAATACCCAAGAAAGATAGAAAGACCGCTTCTCT  
TTCCAGGCCACATAGGCCATCCAACCTAATAGTCTCATATCGGCTACCAGTCCGAGATAC  
GGCACCTTTTCCGCCAAAGACAGCCTGAAGCCTGAAAGCCAGGCTACCAGTTCCAGCTGA  
GACCCGCGGCCCAACAACTGATTGTGAGTTGATCAAACAGCAACACGGCGACAAGGTAA  
AACCAAGAGGTCTGGGCATGAGACCGACCCGGGAGGGCAATTAATGTGAGTCCGGTAAAAA  
TCAAACCTGAAGGAAGGGGCCAACCTATCCGGAAAAAAAATACCCCTTAAAGAGGGGAAG  
CCTTGGAAAGACATCCAGCCAGTCTTAGTCCAGTCTTGCAGTATGGCCTAGTGAGAACAG  
GCATTCCTTTTCTGCAATACTGTCTGGCCCCATTATCAGACAACAAGCCAGAATACCGG  
GCAAACAAGCCACCCTCAAATGAGAAAAGATCTACCAATTATCTAAAACATGATTTTGG

GAGTGCTTTCCTCTGAAACCCCCCTACTTCTCACCTATTATATGTACACCCCCCCCCGCA  
TGCTTCGTTGTTAATACTGTCTTTTTCCCCCCTGAATGGTGTCTAAGGCAGAACTAGAT  
ATTATAGATGACCCCTTTTACAAGGGCCACCTGTCTCTCAGGGTGAACAGCAACCGCCC  
CCATATAGCCCTTGCCAAGTGCTCCTGAGGCTAAAACCCAGGAGCAAACACCGGGGACC  
CTACTAAGTCCCCCTCACACTCGGAGGGGAACACCGTATTCAACCCTTTTTCTGGGCCTC  
AAGATAATAACCCCAATCCGAAAAAGAGAAATTGGAGGCAGCCACTCATGGCCGCCATG  
GAGACAAGACTAACCATCAGCCCGCACTCCCCCTCTAACCGATATACAACAATGCA  
AGGAAAAACTAGGAACTATTCGGAGAATCCAAGGAAGTTTAAAGACGAGTTAAATAAAT  
TGACCATGGCCTTTGATCTACCCTGAAGAGACCTACAATTAATTCTATAGGCCTCTGCA  
CACATGAGGAAAATCAACGGATTGTAGTTCCACAGCAACCGGACGCCAACCTGAGCCCA  
TAAAAACTGCCAGGCTCCACAATTCTCGTCATTTACTGTGCACCAGGAGCAAAAACATAC  
ACACCTAGACCAAACTAAGAATGCCAAAGAGACTACCAAGAGCTTGAAAAAAAAGGATA  
AATTACAACAAGATCAACATAAAGGAATGAGAAAAGTAAACAGGTAAGCAAATAAATTGGG  
ACAAAGGCAGAGAGCCAGACGGGGGAAGCCAGAGGAGCAGCCAGGCTGACAGGCAGAG  
AGACGACCGAACGTAGGAGAGATGCGAGGGAACGCCGGGGTAAGGCCAGAGGCGAGAGC  
TGCGGCGAGGGGGCCGACCCCCAAATCCACCAAGAGACAACCAGTAAGCTGCAAGAAA  
CAGAGAGGGGCACACAAAATACTAAAGACCAACATTAAAAAACCGCCTTTATCTTTTAAA  
AGAAACCGACCATGAAGAAATGAAGAGAAAAACAAAGAAAAGAAAAAAAACAGAACAAAT  
TTGTGGCAGCCTTAGCCCCCGCTCCCTTAAGGCCACATAGCTCCAGAGAGCCACGAAA  
AGTTTCCGTCCAAAAAGGACAGAAAAGAGTACCCCTTACCTAGCACAAAGGTGAGAATA  
AGTATAACAACCTAAGCATCATGGTCAAGGCACCAAGGATTTACCCGCCGGCACCAGTC  
AAAACGGAAAGCACCGAATCAATAACAGAGCTAACATCCTCCCAACAGCTCACTCAGACA  
AACTTTAGCAGCGAAAAAAACCAAATTACAAGAACATTGACCAAAAGGACAGCCGCCCA  
TAAACGAATACAAGGGCCAACGGCAACAATTGTCTGCAAAAATCCCCACGACCGAGCAGGT  
CCGCTAACAAAGAGGCCACACCAAAGGCCACCATTTAGAGAAAAGCCAGGCACACCCCCC  
CCACCAAGGTGAGGCAACGGCAGACTCTTCCCGCCTCTTCCGATATTCAGACCAACCG  
GGCCCAACAAAACCTCAAACAAGAAAACAGTCCAAAAGTTAACAAAAAAGGAGACGGCCT  
GGGAGTCGGACGAACAGCCCCAGACCACCAACATGGAAATAATAACTCATTATCTTCAGG  
GAAAGATCACCTTCATACCAGGCTGTCCCCTTCCGGGTCCCCTCTCCCTACCGCATGTTA  
AAACCGGCCCAAATGGTAAAGGGTGAAAATGAACGAACCCCAACCAAGGGCTTGAGTTAAA  
CGGAAGAAAAAGGGTACAAAAGATAAAAATATACCTCAGATATTTTATACATTGTA AAAACA  
ATCAGCTAAATACTGCATAAGATGGAAAAAATGCCAGCGGTTCGTTAAATTCGCATGGGA  
CCTTGTTCTTGAGAGGTGGTCTAAACATAGCTATTA AAAAGTAAAAATGTTAAACACATGT  
GAATGGAATAGATGCTTAAATAGTGAGCTTATTGTACGCCCTGAGAGCTGAACACTCGAT  
GAGACAACCCGCCTACAGATAAGATATACCTTCAAAGAAGCGGAACACCACCGGATCCCT  
ACCCCCAGGTAAGGGACGCCACAAACCGACACTTCCCTACCGGGGAGAAAGAATAACCAAT  
AAATAGGAACTTATAAAAAATAAGCTTCAAATATAAAGAAAACCAATACTGAAAAATCAG  
AAACATAACTAACACAACCGAGACAAACTGGATAAACGGCAAAAAGAAGCATAATAGAAAAG  
GAAGCGAAAAGCAGAAAATCTCTGGGAGACTAAATACTTGTA AACACACCAAGCATCTCT  
TCCTATCTGTCTCATCACTCAATAATCCGGGAAAAAAGAAAAACACCAACACCCACA ACT  
GACCCAAGACAAAACCATGGGCTCTTAAAACCTGCAACCCTCCTCCGCACAAACATTCC  
TCAAGCCTGTGATAAAGCTCTCCCGTTCAATAATTGTCTGCGGAAAACAAAGATAAGGAAA  
GGGAATCGGAATTTCTAAAGAGGGTCTTAGCAGAAGCAAACAACGCGGAGCATGTAAGTG  
ACCTCATGGGGACTTGAGAACTTTTTCAGCAAAAAGGTTTGAAAAGCTTTCCTCAACATC  
CTCCTAAACAACACTTGCATGACCTAAAACATGGAGATATTCTCTAAAGGACGCCGGACA  
AGAACAACCGAGCACAGACACAAAACGAAAAGATACCGGCTGCTCAAAAACCTAAAAACCC  
GAGATTACCTAAGGCATTTATTTATGAAATGGCAAAAAGGCTATGACTTCAAAGACAGTAC  
ACTACGACAGTTGAGCACTTCTCGAGGCTTTACGGAAAAGCCTTAATTCCTCTAATTGAT  
CAAAAGGGCTGGATGTTACCTACAGTTTTGGGAAAATGGCATGATGCTAGAGTATTTGG  
AATTTCGATCTTAAAAAGAACGGAACCTATAGAACATTATGAAGGCCATATATATAAACCTG  
TGCATAGCATTGAAGAAATGGAAAAACAAAGTGCATAAAAAGAAGTTAAAGGCGAAAGACC  
AATCGCTCACCTACATACAGAAATATTCAAGAGCGGTAGCCCGCAAATACACCTAGAAC  
TAGTACACCCGAGGTGTA AACTGAGGCCAAAGAACAATAAGAACTAGCTCAGGCCATGGC  
AGGCAAAAAGACGGTTGGAGAGGCAGCCAGATGCCACAATTATGGCTTATAGCATATCCAC  
AAAGGCCACCAACAACAGGGGATTA AAAAAAAGTACAGTAAAGACGAAGACCTTCTCAA  
GATGGAAAAAGGAATATCCAGACCGGCTCAAGGGAACCCACCAAAAAGGGACCACAAAGC  
AAGCTGTAAAACCTCCGCAACCTGCCCCCTGCCCTCAGGGAATAACTAGTTGGGTACAA

CTGTCCAGGATTAACCTGTTTCTTATGAGTCACAGGCACAAAAGGAGGACACCATGACC  
TACATCTGTGAACCTTTGGAAGACTTCCACTACCTATTTAAAAGAATCAACACTCAGCCA  
GAAGTGGTAACGTGATGCTGACGGAAGCAAAAGTAATAATAACCCTTCACTAACTACTGA  
ATTAACTAAATATGTAATAAAGAAGAAGGGAATAACATTGAAACTGCAAACACCGAAATT  
ACCATAATAACCCGGACAAACCGCCAAAATGACCAACAAGTTAAAAGAATCAACTCCAA  
ATCGTGTAAGCACACCCAGTTAAAACGAACCAATGTAATTGGAATCGAAAAACTCAAGAA  
AATAATACAAAACAAAAACAGAAAGAAATAACCCCAATGAAGTTGAATTAGTCATAAA  
CATAACCGGAAACAAGTGTACGATCACCGAAACGATCCCAAAAAACCAGCAAACAATTCA  
AAACAATCCCGATCAAACCAGAAAAACCCCTAACAATCTAAACACCTTTTCGTCAATGCG  
AACCGTCCCAAACCTCAGATGATTCACCCAGCCGTTCAAACTGATAGAGAAAGAACCTG  
GAAAGAACAGAAGCAACAGCTTCCCGAGCAAAAGGAGGAAAACAATTGGACCGACCGGAG  
TGCTGAAGTCTCGAAAGCCGAAGGGTCGACGCAATAGGAAGCCGCAACCCCGGAGCCCC  
GTGCCAGGGGGCGACACCCGCCCTGAAAAGGGGAAAGGGCCCAACAGAGCTGTAACCATT  
AAGCCGTCCGCGGACGGCAAAGCTAAAAGAGAACAACAACAGCCACACAGATCATAAATA  
AACATGCATTCTGGGTGATGCACATGCCAAAACTGCAAAGAAAGGGCACCACAAACTAC  
CAAACACTAACCTATAAAGCAAGTGAAGAAATTCCAAAGCCAGTGTAAAGCAATACAACA  
GCCACCCTCCGGACCTATAAAATATTAGAACAACACCCAATGTGTCTGGACCCGAATAGA  
GCCGGAAACAAAGAGGGCCACACGACCCAGCTGTGCCAACCCCTTGGGCCTCAGTTTGT  
GCCCAACTGGGCTCATTTTTAGGAGAACAATAAAAAGGCAAAAGCCCGCGGCCGCTGG  
CACCCCCCGTAGATGCCGCCTTTTGCCTCCCTGTGATCCTCGGCACGACCTCGCCGA  
CCATATGAAAAGAATGGAACAAACAGGAAACTAACCTGGAAACCCTAGACCCCGGGCTG  
GATAAGATAATAACCACGAAAAGAACTAATGAGTCAGAAAAACCCACAATCAATATTT  
CCAATTACGAAGAGCACCAAAACTCACCTCGGAAAAAAAAACAAAAGAAAACATAAGAA  
CCTGAGAACTCCCGTACCAGAGCCCAAGGCAGAAGCCAGACTGCCAGAACTCAAACC  
ACCAATATAAACCTGCTTCCGCCACGCTGGGCCAGCAGCAGCCTCCAGGGAGCCGGC  
ACCCGTGCCTGAACACACCTACCAACAGATGATACAAGAAAACAAAACAATAGATTGAGC  
CGGCGCTGGAGTCGCTCGCCCCACCACAGCAACCGACCCACACGGACAACGGCAGCAGC  
ACAACACCCAGACACAAAGACACCGAAGATACAAATATCAGTTCCAGAGACCTGACCATT  
TTTGGCGAGTGGTGGACACAAGCACTAACACAATAGTGCAAACCTTACGAGCTTAGCAAGG  
CGGAGGCCGGATAATTAATAGCACAATATGTCTCCATTGGAGACACGGAACATATATAA  
AAACGATACTAACTCAAAAAAGACAATACGGTGTCCCAAGAGTAGGCACATGAAAGTCT  
ATAGCAACTATGAACAAATGCCACCATGCTGCAGGTACACAGAAAACCTGGAACAACTA  
AAAAAACACACAAGACAAAAATTATCACGAAACAACAGAAACCACCAACACCGGACCCC  
CAAGGACCCTGCGAAAGAAACGAAAAACCCCGAAAGCAAAAACCTAACATAGACACAATTG  
GGACCTTATAACTAACGATATCCATAACAGACAGAAAAGTGACTTACAAAAACAAAACAG  
GCTTGGCATTATAACCAGCTGAGCATAGAACAATCTGCGCACAACCTAACTTCTGGCGCTCG  
CTTGATGCTCCCTCCCGCTTGGGGCCTAGCTCGGCCTCGCCAGTCGTGAAGCCAGCGAC  
GGTACGGGCACACAACTCCAAAAGAGATCGACTACTCAACCAGGGCATCACAAACCTCG  
GTCACCGAATTACATCTACTGGGAAGACATAACACAGGCAAAGAACTACAACCAATAG  
CCCGACACATCTTAAAAGAAGAAAGCCAGCAAAACAGGAAGCCCGGACCCAAAACCGGA  
TGGTCCCAGAAACAACAACACCTGCCAAACGGGGCAACACCGCAAGTACAGGAACACAAG  
ACAGCAAAAAGATAACTAAAAACCTACCAACTTTCTTGGTGTCTTTGTTGCATAGTTAC  
TGCGGGCTGAGCAAACCAGGCACCCCTGAGAATAACAATGACAGCCACAGGGGATCAGG  
GCCGGTCATATCCTAAGAGAGTTGGGCATACAAATATCCTGAACAATCAAAGGCTGTAGA  
CATCAGTGGGCCAAAAGAAACCAACGCAGCACAGCAAACTCAGAGAGATGTACCACTGA  
CCGCAAAGGTCTGCAGCTCGCTAATAGGACCAAAAAATCGGATAAAACAACTCAGTTCAAC  
TCGCCCACACACCATCAAGCGCCAAAAAATCAACCAAAAACGCTGACAACAACTCCTGC  
CAAGCACTTGCCACGTGACACACCACATGAAGACGATCCCCTGTGATAAGGGGGCGGGTGC  
ACAGCCACTGTCAGTCATCAACCAGGATCCAACAATGATCATCCAACAAAAGAACAAG  
GATCCTGCCCTTCCGGCGCGCGCCAAATCTGCCGAAAGTAAGTATAAAACAATGACACA  
AGAAAAACAACAGTAAACAGGACTCAAAAAAGAATTAACAATCACAAATCGACAGT  
GGGAAATCCTTCAGATCGGGATGAAGAATGACCAAAGGAGAAAT

>2\_cons

TACTTTGGTGCCGCTGACTCGGATACGTTCCCTAGTGGTAAGAAACCTCTATGCCTCGC  
CTTCTTTGGCTGGAGCGTTCAACCCCGTATGCAGTTTTCTTCTCCCCGGCACCCAC  
CGCGGACCAACCAACCCCAAAAAAAGCCGACCGCCACAGTCGCTCCGCTCCCCGAG

CCATAAACACACAACACCCAAAAGGAAAGAAAAATACTGGGAATGAGCATGGAAGCCA  
AACCGAAAAAGCCAGGGACACCAACAAAAGACACGGAGAGGAGGCTCACGGGAGAGGAAA  
GGGTAAAGCCAAAAGCATCCCAAAATCCGGGGTGCACCTCGTGGAGTCTAAAGAGAGTTGT  
CCCATGTAGCAAACCGGACACCACCTGGGAGCCAAGTCTTGACCAGCAGGGAACTCGAGG  
TCCCCTCGGGACCCAGCTCGCTGGCAGCAAAAAGAGCATCATCCTTTGTTCTCCCGCCG  
ATTCCTAACCTACGCGTATTGGTGATTGCCTTTATCCTTTTTTTTTTTAGCTCATCTCTC  
CTCCTAACTGGGATTAGGAGTCATGACAGTGCCTGGATTGGGAACAGCAACGGACCGGAT  
TTGGTAATTAACCGCCTTGGTATCAGTTAGAGGCCACTGAACCTCTGTGGGTCAGAAAA  
AAACTGGAAGTAAAAGAATTAGTTCAGGGGGTGTGACAGATTACAATAATTTACCAAAG  
ATTTTTCAAACCTTCCATTCCGACGAGAGAACTTCTCTTGTATCTTTGGGAGAAAACTTT  
CCCTCGGCAAATTCCCCACCTGCACCAAGTAAGGACCCACCTACCCGCCACCATGG  
AAGGAGCAAACATAGGCCAGGAAAGAAAAGGAGTCTTCGACATGAGCGATCAATTAAC  
CTAACAGCAAGCCCAATATTAGACCAAAAAATAAGTTCCAGGACTCACAAACACTGCCT  
TCTTTAAGGCAAACACGGCCACAAATGTTAGCACCTGAAGGCAGCGAAGCAATACAGAAA  
TTGCCACATCCCCTAAAGCTAAAATAAGGCGGAACGTACCGAATGCTTAAAAACTCGGT  
CTAGAGCATCTGGTAAAAAAGACTATCAACCAAGCCTCAGAGGGCAGCCCAAGAAACAAC  
TCCAAAACCTTAAGCCACACTTTTCGACATCAACATCAAAAAGGCTCATATAACAGCGTAGGG  
AAGACCAGCTCCTCAGTAGGAGTTAAATCCTAGATCTGGTTAGCTTTACATTACTCACA  
AGACTGTTCAAACCTAGGGAGAAAAGAACCAATATGGCTGACACACTGTCCCTGGTGGG  
AGTCACTTTTGTGCGAGGAAAACCAATGATTATGTGAGATCCCAGACCCCTCATTAAAAA  
CACAATATTCCAAATAGGCTGCACAATCACAGCTTTAAAGCTATCTTTTCTAAGACTCAG  
GCAGACATGAAGTGAAGTACTGATCATTATGACCAATTAATAATCAAACCTAAAAATGAACAA  
GGGGCCTAGGTAAGATCCTACCCCTCTCCAGCCCTAGAAAAAACATACTTAAAGTTTC  
CTCAAACCTAAGCGTCTCAAATAACTCAACTGAAGAAACCAATAAGGACCACTAAT  
GTAACCTCAGTCCATCAGACCCCTTTAAACAAGCCAACGCAACCCCTTAAATTCCAATG  
TATTCCTCGGGTCCCACACTCAAATCAGTAACCCTATCCATCAACCACCAGCACAAACC  
TAAAGAAAACCCATAAAAAACTCTGCACATTTTCACTCACAAACTGGAAAACATGAAAC  
AGCCATTCAAAGCCCTAACCCTCTAAACTACATTGGCCCCCGAAACCGACTCCTTCT  
GTCTCTAATAACAACCTAATTAATGCCCTTCTACACATAATCATTTACACTTAATTTT  
ACGCCAACACGAGAGAGACCCCGTATTGAAACACCGTCAGCCCATAGCATGGAATCAGT  
CATTTGGCAAACATTGAAAATCTATTAGACATGACTGGAGCCCCAGTCAGGAGGGCGCCG  
CAAGGTGGGTTCCCAAGGCAGCTAAGGAGAAAACCCAGAGAAACCTCCCGCGGGGGCAA  
AAGGCCAAGCATATCTTACCAATCCAGCTTCTCCACAAGATCCAACCATCCAGAAATC  
CTCATTTTAGGCAGGTCCCAGTCTCAATTCCCCTCCTATGCAATCTCGCCCATAGCAAAC  
CAAAACTCAAACCTTTAACCTTAAATAGCAATTTAGGCCCTAAAAAACTTGATTCTCTA  
TCCCATCTCATGAACTTAAGCAAATCAATAACGATCTTGTTATGTCTATGAATCACCCCG  
ATAAAGACATGTACCAGATCTGTAATGTAGCCATAGCATTTGAAGACACCTGGAAGGACA  
GATCAGTCATCTTGACTAAAACCTGACTAAAGAAGAAAATTAACCACCAGGAATGTGG  
CCCAAAATTTTGCAGATGAAATTCACACGACTAATCCTAATGATAATAGAGCTGGGGCAG  
AAGCTTTATCTAATTAAGAGTTCTACAATAAAAAATCCTTTATGAAAGCTAAAAAACAGC  
AGCTGCCTTAACCAAAAAGAGATTCTAGAGGACAAAAGAAGGGCCCTATATTGGATAAAG  
AAGGTGACTTAGGAAGCTGGAAGAGAAACCACATGATCGTGTGCTTCTCGAAGGAATAA  
AAATGGCTGGAATTAACCTTTATACTATGTAAAACTATCCAAAATAGATCAAGATCCTT  
ATGAAAACCTACTGCTTTCTCGCTCGTCTATGAGAGGCCATGAAAAATTACTCCACCT  
TTAACCTAACACGGAAAGAGAAAAAAACATCCTGCAATAACAGACAAGGCTCAACCCGG  
AACTAGGCAACGAGGACTCTCACAAGGAAATCTAAGAAACCCCTACCAATTGACTTCTAT  
TCCCTCAATCCTCCTCTTTGGCTCAGCCATCACCAGATGCACGACACAAACTACACTTCC  
TCCAACCCACCCCTTGAAGATTCTAGATTTTGTAGAACTACTTACCACCTCTTTGAAAAT  
ATCTCGTACACTCGTGGTTAAGTAATAACCTTAGTTGAGGCTTGTGGTTTCACCTGGGA  
GGTTACTTTTGGTAAAGTTCAAAGCCAGAAAATATTGGCCGTTTGGCCCGGCTAAAGTCG  
GGCAAAAAGCCAGCTGAACGGACTCAAACCCCACTACAGGACCTCCTAAAAGTAGCCTTC  
AAAGTCTTTAAAACCGGGACGAAAAAAATAAAAAAAGAAACCTAAAAGGAGAAGCAA  
AAATAAAAGCAGAAGCTCCAACAAAAGGCTGCCCTATGGGCACACAACCCAACTCCAGGT  
CGCCATCAGATCACCCCCCACTCAACAGCCACCCATGAGCGGGACCAGGCAAAAATGCA  
AAAGCCAACCTGCCCCCTTAAAGACCTCCAAAGAGCTATGGGCCGCGCAGTCGTCAAAG  
ACATACTAGCCTCCCGGATCCACCTGGTCCCTGTAGTAATAAAGGCCACCTGAAACACTA  
AGCTAAAGATGGCCCAAGTAATAGGAAGTCCCCAATTTTGGAGTTTGGGGGTATCAGAA

ATTACTTCGCATTATGAGAGAGCTTTGGTGTGTAATAACTAGGTAGGAAATACACATTTA  
GGGATGGCTAATGGCAGTTATGGGGGATACTCGGCTCTTTGCACATTTGGATAAGAGAAG  
CATGCTCTTGGCCACCTGGAAGGTATGAAAATGCCCTCCCCCTCTGTAGGATAGAAGAC  
CACTGGAAGGGCGACGGTGTCTACAGACAGGGTTAACGGGGGCCGGCTCCCGTTTCAATG  
AAAGCCCAGCAATAAGCTCTCCACCGTCTTCCCCAGAGGATTCACCACAACCTGGACAA  
CCACAAGCCAAAAGAAAAACACCACCAACAATATTCCGGGACTACGAGCAAGCCACCGAA  
ACCCACCCACTAGACCACAGAAGGGGCCCTTGAGGCCACAGACTGTCCCCCAACCCCTATT  
CCCTGTAAAACATAACCCATGAAGTGCGAAAGAGTTCTGAGCACAATCCCAATAGCAACTA  
AGAAACAAGTATAATTTCTTTAAATTTATACAAGAAACACCAAATCCGTTTTACACAATCA  
ATTGGGAAACAAGTTATTCTAATTAACCACAATAGTCGGAACCCATCCACAAGTCTCCAA  
ACCCCAAGCCAAGAGGCCAAAACCCCTCTACTCTTGGCTTCTCCAAACTAACCCCTTTTGA  
GCAAGCTCTTTAAATAATATAAAAAAATATTCTCTCAGCACTTCACTTAAAGCTGCTCA  
GCCATACTGCTCAATAATATATTCATAAATTTCCAGAAAACAAAAACGTATTGTCTGA  
AAAGGGAATGTTTCTTATATGGTTTAGCCAATACTTCAATCTCCACCCACAAACTTGAA  
TCTTAAATAAAAGTGAAATATACAAACAAATTAGCTGTTGGAAAAAATTCGAAACCTGCA  
GAAGCATCACCTTACACACCCGTCGACATAACCTTTAAAAATCCCACCAAGTTCCTAAAA  
AAAAACAATATCCCATGATCCCAGAAGCTAAAAAAGAGCTAAAACTATAATTTTTTGAT  
CTTTTAAAAAAGGATTACTCAGCCCAGGCAACATACCCTGTAATACACCTATTTTAACT  
GTAAAAAACCACAAAGAAATAGTAAAAAAGGAGATGTAATAATGGAAAGTGGATAG  
ATCATAACAGACTAGTTCAGGACCTTAGACTTATTAATAAAGCTGTAATACCAATTCATCC  
TGTAGTTCCTAACCCCTTATACATTACTTTCACATATTCCCTCAAGTACAAAATTATTTAC  
TGTACTTAAAGCTAAAGGATGCCTTTTTTTCTTAAACCCTTAAACACCATGACCCAAAATAT  
TTTTTTATTTAAAAATGGCAATAACAAAAACAAACACGTTTACACACAGATCACATGGTCT  
GTTCTACAAACAATTGAACAGGCACATAAAAGCTCTATAAGAAAGACATTTTAAAGGTCAA  
ACAAATAAAGATTGTATAATATTCTCTAAAAAATATTTAATAAGAATGCGAGCGAATATA  
ATAGTATTTGAACCAATAGATGAAGTTAAACAATCAATTTAAAAATGGAAAGCTCAGGGAT  
TCAGAGATAGCCCCACCTATTTGGACAGGCCCTTACAAAAGACCTCGCTGAAATACGCC  
TTACACAAGACACACTTCTTCAATAAGTCAATGACTTGCCTCTTTCTATGCCAAAAGGTT  
AAATAACCTTCCAAAACAGTAAGCATTCACTTAACTTAAATGGAATGAACAAAAAAAACCT  
AGAAAGCAACAAAAGAAAACCAGTCATCTGGCACCTCTAGCATCCTCACCAAGAGACACA  
AAACCCCTTGAACAAATCAGGTGTTGATATTTGGTCCACATCCAAGCTAGAAAGGACCAA  
AATACAGTATTAGTACTAAACAAACTTGCTGATTGTGGGTACAAAGTATCTCCTTCTAAG  
GTACAAACATCCACACAAAGAGTTCAATTTTGGGGCCCCACTAAAAAACACAAACAAAC  
ACCCCTCAAGCACCCATCTACAGAGGGCTCCTGAAGCATCCAAAAGAGAGAAAAGAAAG  
AATCTTATATAAGTCATGCCGCCACCCAAAACCTAAACAACAAAAAAGTGAATTTTTACCT  
ATTTCAGGATATTTTTGAATATGTATTCCAAACATTGGTTTAAATAGCTCAACCCTTATAT  
CCTGCATTACAGGTGTCTCAGTAGGTACAATCAATCACCCATAATATTAGAAAACAAGAAA  
AAAAAATGCATTACAAATAATAAAACATGCTTTTACATAACAGCCCCAGCACTGGCACTAC  
CAAACTCATTAACCCAATTAATTAGTTTGTACAAAAACAGTGTGGAAAACCTTTTCAAG  
CACTCAAAATCCTAGAATATGGTGTCTTCAAGGAGGTTTCATGAAAGGATGAAAAGGACCC  
CGAAAAGCACTCTTGAATACAGGTTTCTAATAACTTTAGAAATCACATCATTGGACTGGG  
TAAGAATTCCCGGAACCTAATGAAAAGACTGACTGGTTTATAAACTGCTAACCCAAGC  
AGAACAAAATTAATTGAATACCAAGAAAATACTTTGCCAGATTATCATGCTAAATCAGC  
CAATACTGAAATTGTTTAGATATACAATTTGAATGAACTCCATGGTCAAAGTCAAATGAC  
CTATGATAACCCATCAGTTATCAGGCCAAACAAACGGAGCCGCCAAAAAACTGTTGGCT  
AACTTAAACAAAAACTTAAATTTGTCTGTCGAAGGATGGCCAGCCTGTCTGAGAGCCTTGG  
CAGCCGTGGCTCTTTTAAATCCTAGAACTCAAAAACCTTTACATGCGAGAACACATAACTG  
TTGAAATAACACGTAACATAGAAGAAATGATAAATAAATCAGGTTACTACTTTAATTCT  
CTAAAAGACGTTTATACCAAACGGTAGTCTGCCTCTTCCAAATTCATGCCAAAACCATTG  
AAAGTTCCGCAATCTTTCCGGCCAACAAATAGGCCACTTAAACCTGCATACATGCACCACA  
ATCAAGTGTCTTATTCTGAGCATGTTTTCTCGTTCATTAATAATCTTACCCACCCAAGAGA  
ACCAATACCATAACACTACCTTCATTACAAAGAATGTAACCTTGTCTTGTAATGTAAGTT  
ACTCCTGAGATAAAGAAGGAAGACATACAGCTGATTATGCAATTGTAACCCAAAAACAAG  
ATATTGAAAGCACATATCTCCCATCAGACACCACGGCTCAAAAAGCTGAACTAATGGCAT  
TAACTAGGGCCCTTAAATTAAAAAAAGAAAAACAGTTAACATTTATACAGACTCTAAAT  
ATGCATTCCAAGTACTCCATTCATATGCCATTATCTGGAAAAGAAAGAGGTTTCTTAATA  
CAAAAGGAACTCCCATCAAAAATGGAAATCTCATATGCAAACCTATTTGCAGCGCTCTAAC

TACCACAGAAAGTGGCCATTATCAATTGTAAAGGACACCATAAACACGGCCAATCCAATAA  
CCGAAGGAAAACAGTTAACAAATTAAGCAACAAAAAAGCAGCACTAAATTCTGAGCAA  
GAGAGCCACTGAACTTCCCACCAACCCTAGAGACATGAGAATTAGAAGTAAAGCAACCAT  
TTTTAAACACACAAATTAATAACAACATGAGTAGTAAACAACGTGGGGTGGAACTTC  
CATTAAGAATAAAAAATCTGACGCATCCATTCACACTCTCGCCCCTGTACACCCCTTA  
ACAACCTTAAGTTAAACTTACCAGACAAATTCCTACCCAGAAAAAGATCGCATCCTTGA  
TGTGAACAGCTTATAACAATTAAGAAAAAAGAAAAAAGTGGCAGTAATAATAGAGACTTAT  
TCTTAAAAATTAAGACCCCAATGAACAACCCTGGGCATGATGTTTCTTTCTTTCCCCC  
ACCCCTGCCCTTTTCAAAGATATTAATTAGTCTTATTCTTTTTTCTACACAATAGTAACA  
AAGGAGGCCATTCGACTGCCTCCCTGGATCACACACTACGGCGCAGAAGAATCAACAGGG  
AACTTAGGCAAATAAACCTTACTTTTCGCCCTGGATGTAAGATGAAGGCCAATCTGGGT  
GAAAAACATTAAGGTATCTAGTATCTTCCGTAATTAGTAAAAACCATCACTGGGGGAT  
CCAAGCACTATCGAACCCCTCCCTCAAAGGAGACACTAAAACAATGCCTTGAATCTTCTA  
TATGTCTAATAATTCTTAGGCCTCCCTTCCATGGTGTCAAGCAAATGACCAAATGATTA  
GAAATAGATCCATCATAATAAGATCTATAGCAGATACAACCGTAAAGACTTTAGTTTTAC  
CAAATGCCCTAGATCTGCCCGTACTTTTAAACGGGAAGAAATACAAGTGCCACCGAATTCA  
AACTGACTCAACAGAGCCGTTTCTGTGCAGTTGAAGGAAATTTATCCCGTACGTGAACAA  
AAACATTTTCAGAACAAGGAACTGAGCTGGAAGAAATTAATAAAAAATTCCTTGAAGGGC  
ATTAAGCACAAGACCGTTTAAATAAACGAGCAGACTCTCCTTATAGCTCCTTCATTGACC  
CTTCTGACCCCAATCGTCTGGATAGGTTGGGAACCCGACCAAAGACCTCAATCCACACCT  
CAGAAAACAGTCTCAAGAACTCATAATCCCAAGCCTAAAAGTCCACTCTTACCCTCGA  
GAGACTCTACCAAGAACAATACGAGACGTCATAAAAAAACCTCAACCAAGCCTCCTGT  
CCAAGCGAATAGATGGGGGCCCCCCCTTTTTACACCAACTAATCCAAACGAGAAAAACAA  
ACCTTCTTTCCACACTGGATACACCATTCCAAACTAAAAGAGCACCAGATCCACATCC  
AGAAATTTCTCACCCCAAAATTATTCTCCCTCCCTCACAGGACCAACCTCACTGCCTT  
AACAAGAATTCCAGAAGTTGCCAATCCAGAACGCCCTGGTCCATAACACTCTCTGCCTCC  
AATTTCCAATCTTTTATCTCCTACTTTGTTTCAGATCTTTCCTGGTATCCCTTCCCCATG  
TCCCTGGATAGTCCACGCCAATTTCTCACACTAATCCAGGAGATATGGCTGCAGGGCACC  
TTCCAAAATTTCACTCCTACTCAAATCTCCTTTTTCTCCTAAGACCAAAACTGACTGTGC  
AAAAAAGTGGTCCACATAACCATTAATCGTGGAGATAAATGGCAGCCTTCAGACAAGAAA  
ATTACCTCTTAAATCAGTCACACTCCCCTCTTTCTTCCAAGTTAAATTTGTTTGTCCA  
CACAAATCCAGCAGTTCACAGCCCTTCCCTGTCAACCTCGCAACATAAACAGATACAAAA  
AAAAGCCGAAAAACAAAACAGACGATATACTCCTCCCTATGCCTACCTTAATATTTCTTC  
TTTTGAAATAGTTTTCTCTTAGTCACTGAGTTTTTGCATGCAGCCATCTCTAGAATGTC  
AAATGGTCTCTTCAACTGGAATGACAAGAGCTCAAAGAAATGTATGCCACAACGGCA  
CCCTAACCTATGACTAACACCCTAATACCAGAAACCCTAAACTATGATTCCTAAGAGTCG  
CACTCCAGCAATAACCTATGGACACCCTCTCAGCTATCAGAGAACCCTGACTAAAAATGGG  
GCAAACCTCAAAAACGCAACTTTGAGACAACAAAACCTTTGAACCAAGCTCGTGCACCACGC  
ATCCACAGACCTCAACAAAAAAGGACTAATACGTTCAAAAAACACAAATATAAAA  
AGATACCAACGAAAACAGATTCATACTCCCAAAGACCACCATTCCTTTACTTTCTGTCA  
GGGGGGCTAACCGGGCTAGCACACACCCAAATAGAAAAAATAAAAAGAAGAAAAAACA  
GGCGGTTCAAGAGACAACCACCAACCAACAGCCCCTGCCACCATAACCTACAACTTTTG  
TCTGTCCACCCCAAGTGTCTTCTCCTGTGTGGCACAACCTCCAGGTTTTGCCAGAAGGC  
GGGCCACCCTACACCAACAAGTACAAAAAACAATCACCAAAAGTCATCGTTCCACCCT  
GCAAAACCCAATGTTGACTGTTTTCCAGCAAATCTAATACCCCAATAACTATCGGAAA  
TTACACGATAACATAAAGCTCAAACAAAAATCGCTTACAGACATCCGAAACGGACCATAA  
CATCCCTTAAAAATCAGTAGGGATGGGTACCAAAAAAATTTCTCGTCTTGTACATCGGA  
TTTTCTCAATTAATTAACACACTGAAACTAAGAGACATTGTATTTCATCCATACTTCACTT  
AATGGAACCTCAGTTGCCTATACGCCCTTTTGAAGATAGGCACCAATAGACAGATTA  
GGTGAAAGCAGCCGACCCTCAGTATGGTCTGCTACAATGACAACCTGAATATCAGCCA  
CTCTCACAATAATCAGTGAAAGACCAAAAGGCGGAATTCAAAAACAATAACCACAACAC  
AACACCCAAACCAAAAATGAGCAGGCAAAATAACTTCCCTTGTCCCTGTTAAATTTACAA  
AAAAATAGAAGGGAAGAA

>3\_cons

TTTCTGGGGTCCCAGACCAGGAAGCGAGACAACAGATGGACTGTCTCCTTTGCCTGTGGG  
GAGGCAGCCCCTGAGCCGGCGGAGACCTGCCATCCGAAGCATCCCTGCGGGGAACCTCCAG

CCAGCAGGAGAGACCCGGACCCCCAGAACCCTCCCAGCCAGCCAATCACCCCAACGGAAC  
GCCTCAAGGGGCTACAGGACGATACCAGCAACAGCGTGCTAAAGAACC GCGGTAAGGAGC  
AAGGGCCCAAGGAAGGAGTGCCAGGCCATATGGACGAAAAAGGAACTGGATCAACACCC  
GGAAGGAACCGACAAATCCAACCCAGACATGCGAAACGTGGCAAGACTGGCCTGCGAAAC  
TGGACCAACCCCATACCCCTAACAAGTGAAAGTGGTTCACTGGTGGAGAAAATGGGCC  
GATAGAGCGGCAAGTCCAACAAAGAAGAGCTTGCTGGCAGGGTGGCAAGAGTGGCTTGCC  
ACCCCAACTGGGAGTGGGGAAGTGTGTGTGGACCTACCCAGGACATGAGAGAGGCTCGTT  
TCGTCCGATGAGGAGTCTGGGGTAGGAGTGGTGTGTGTATGTGTGTGAATGTGGGAGCC  
TAACTAGGCTCACCCGGGACACGGGAGAGGCCTATTACGTCCCAACAACAACCCGGCGAA  
AGGCAATGAGTGTGAACGTGGCCGTATCACACAAGGCCAACGTGGACAGGAACGTTGGGT  
GCTGTAGGTCACCTGGCTTGATTGTTTGCTCCGAGGGGAGTGGGTGGAACATCGGTTCT  
ACCAGCGAACATCCGCCAGGGCAATTACATATGGCTAACAGAAAAACAACA ACTTAGC  
GACCGAGGTGGCCTGCGTAAGGGGAAGGTACGTACCCAACAGCAAAATTGGGGGCTTTCC  
CAAGGATCACATTTTCGTCTGGGTCGGGTTTCTGGTAAAGTGTAAAAAGGGGAGTGCGC  
GCCCCACACCCCTTCTTTGGGTTGCGGTTTGCATCAGGAGGGAGCAAGCTGGATCCACCC  
AGAGGTTTCATGAAAAACACGAGAATAACACCAACTCAATAAGAACTGTGTTGAAAAATT  
GTTAAAAGGGTTCTGAAGGAGACTGGCGAATAAAAAATCACTGTACACAACTAAGGACGG  
AAAATCCAATAGAGTGGCCAACGTTTATCGTTGGTGGGCCGGCCGAAGGTAACATGGACG  
GGGCTGTAATTCGCCGTATGTTTAAGATATGGTATGAGGTATGTTGAAAGCCAGGGCACC  
CAGACCAGTTTCCGTACATAGACTCTTGCTACAGCTGGTCCTAAACCCCAACCCCGA  
AGTTGAAGATAATACAGTCTCCAAACGCCAAGCAGAACC AAAAAAAGACCTAGGAGAAA  
AGAGACAAAAGAAAAAGAAACAAAAGAAAAAAGCGTGGAGAAACACCAAAAGACAGACA  
GCAATGCCAGGCAAGAAATTGCAAGACGGAAGAAGCGAGCCAAACCGAGCCTCAGACAGA  
AGAAGAAGCAACTACGAGTAGGGGGAGAAAAGGAGACAGTAACCCTAGAAAAGACAACAA  
AAACCAGATTTGCAGAAACCAGCAAAGGAAAGACAAAGACTGGAACTTATAAGTCCAAA  
CGTAAAAGAAAAGAGAGAAAAGGGATGAAAAAAGGGATCCTATGACACACCTCCAGAAGCA  
AAAAAGATAGTCCCATAGGTTTAATGATTAACCTATGGAAGAATCATCAAAGCACTAATC  
ACAAGAATAAAAAACAGATAAAAACAATACTGTGCCCTTATCTGGCCACTAGGTCCCCTCC  
ACAGCCCCTCAATCTTCTAGCCTAAGTGGGGGTCAAATGAGGATGTAAGGGCACAAAATG  
TAACCCTATCCGATACTGAAAAAAGTAAAGGGGAAGAAAAAGAACAAGGTTACGCACTCT  
GTTGGATCTAGGAATTAACCGCCTTGAGCCTTCCACTTCAAGAGAAGAAGAAACCCAAA  
TACCATGAAAAAGTAGGAAGCCAAC TAGGAGGAATATCCCATTTAAAACCACCAGCGCTT  
GTGACCCCTGACCTTTTTTCTGGAAATACAGCCTCCATTCTACACCAAAAAGCCGCA  
GGCTCTCATCAAGCAGCTGCAATCGATTATGATGACACAACACCTCACCTGCCACACTG  
CCAACACCTTCTGCTTATAATCCTTAAAACAGAGGAGAGTAGCAGTGCGACCCAAGCAGT  
CCTTAACAGCCTATATACCATTACCCAAGAGGACACCCTCAACA ACTACA ACTATGTAAA  
AGCGATATTCACAAATACCGATCCCCATTTCGGACCAAAAAGAAATGAAGAAGAAGTGTATTC  
CCTATAAGGGGAGTGAAAATAAACTGGAAAAGGCTGAAATTAGTTTTATAAAAGGACTT  
TTAAATGAATCAGAAGTCAGAAATGCTTCAGAAGGAAGTAGAAAAGTCCCCGAAAAGTCCA  
TTAGAGAGGAGAACATAAAAATGAACAATTATAAGGGATCTGAGGCATACACTCGGTATG  
CCCCAATGGCCCTCGTGGACACTCAAAGTCAAGGTAAGGAAAGAAAGGAAGATGTGGAGG  
GAGGCCAGGAAGAGTCAGGAAAGTCAAAACCCCTCCCGCCAAAAAGTTCCAAAAGACATA  
CCCAACACAAATTCATTGCCAAATCGAGACATCAGGTGAGAGCTTCAAAAAGTCAGAGGGC  
TTTGCAAACATGAATGTCAGTAGGCAAATAATAATAGAAGGAAACGTAGAATCAAAACGC  
AAAAAACAAAATCTAATGAAAGCATTGGAAGACAACAAGAGAAAAGATGAAGAGCCAACA  
CGATGCCTACACAGACAAGAGGAGCAGATGAAGCAATAACTAGAGACCGATATGGAGGGC  
CCAGTCGAGGAAGGGATCAAAAAAGTCAAATTCATGAAGAAAAGTTGAACTGGCATAAAA  
ATAAAGGTGCAAAAAGTAGAGAAATAGGAAAATAGACCAATAAGTGAACAATTAAGGGAA  
AATAACAAGATATTCATAGGTAGAAAGAAAAGAAAGACAAAAACGAATGACAAAAGCAATG  
TTTCCCACTGCCTAACAGATAGCGACAAAATCGGCTTAACATAGACAAAGGATCAAGAGA  
ACCAGAAAAGATAAAATTGCCAGACCTTAACCACTCTAGGGACCCAGGACCCCTATAACA  
AAGGGAGGCACAAAGTCCATGGATATTATGCCCAATACAGTGCCTGAAAATTCAGTAGTA  
AATCTTCCAATAACCCAGCTATAAAAAAACCATAATACTTTTAAACAGGGCTACAGGGCAA  
TCAGAAAAAAGATTGCCCTTCAAGGAGAAATCTGATAGCAGGCCTAGA ACTCATG  
AATGTTCAAGTAGACCAGGAAAATCATCGGGTCCCTCTGTTTGAAAGAGAGTCCCACCAG  
AAACTGTTGACAAATATAGATTTTGCACCAGAACGTGATATAGTCACACTTTTGGTTAC  
ACAAGGGCTACTGTGTCCACTATTACTGTCCAGAATCTAATCCTGACTCCTCTGCAGAA

AAACTGACAGTCTCAGAGGTTAAAAGGGGAAGGATTTACAGTGAAAATATTAATAAAAAACA  
GAAGTCAGATAAAGACCAGGATTACACAATCATATTCAGCGGTCATTCTTGTAAATCCC  
TGAAGCAGCAACTAATTTACTAGGAAGAGACTTAATGTTAAAGTTGGGATTAGGGCTACA  
AGTACGATAAACACAATTCCTAACTGTAATGGGCCTACACATAATGAGCGGTCTGCATAA  
CATGAATCCTGGTGTCTGGGCTGAAGAAGACAACCCTGGGATGGTGGAAAACCCCCCAC  
CCAGATAATAGAAATAAAACCACATGGAGAAGTAGTAAGAAGAAAACAATAACCCATTCC  
CAGAGAAGGCATTCAAGGGATAAAACCAATTATCGAGGGCCTCAATAAAGAAGGGATGAT  
AGTTCCTGCATGTCCCATATAAACTCCAATTCTGCCTGTGCAGAAGCCAGATGCCAA  
TGAAGTGTAGGTCATACAGGCTGGTACAGGACTTAAGGGCGATTAACCAAACGGTAAAGA  
CCATCCACCCTGTAGTACCTAACCAAGCCGCACACAAAAAAGCAAAGAATAAAATGC  
AGATGGTTTACACCCTTCTGCGAATGATTCCATCTAATGATAAATGGTTTACTGCAATAG  
ATCTGAAAGATGCTTTTTTTTCTTTCCCTCAGGATTCTACAGCCTGGAGTGCTAAAGA  
GAGCCAGGACCTATTTGCCTTTCACTGGGAAGATCCACACTCAGGGAGAAAACAACAATA  
TACCTGGACTGTCTGCCCAAGGGTTCATGAACTCCCCACCATCTTTACAGGGAGCAA  
ATTTTGGAGCAAGTACTTGAATAACTTTACCCTTCCAAAAGAAATCTGACTGGTACATTA  
CTTCGATGACATTCTTCTATTTGGAGACACTACGGAAGAGGTAGCAGACAAAGGAACACA  
AATTCTTAACCATCTCCTGGACAAAGGAGGGGATAAGAGTCTCAAGAAAAAGGCTCAGT  
GTTTAGATCCCGAAGTTAAATATCTAGGCTTCATGATAAGTGAAGGCAAGCGAAGAATAG  
GTCTGAAAGAAAGGAGGCAATTTTGTCTCTGCCACGCCTAAAACCAAGAAAGAAGTCA  
GAAAATTTTTAGGGGTAGTTGGATTCTGCAGGATATGGATTCCCAATTTGCACTAATAA  
CTAAGCCCTTATTTGAGAGGGTAAGAGATATCAAGAAAATGGAGAGGAAATCTCATGCAT  
GCAACCTATTATACGAGACTCAACATAAGAAAACCTTGACCAAATGAATAATGGAATTA  
AAACCTTAAAGCAACATCTCTATGATGCCCCAGGCCCTGGAGCTACCTTACTTAGA  
AAAGCCATTTGTACTGTTTGTGTCAGAAGGAGATAGGATGGCTGTTGGAGTCTTACTCA  
AGCCCTCAGAGGCTGGCGGCAGCCCGTGGCCTTCGTCTATCAAACTACTCGACCCGGTC  
ACCTGAGGATGGCCCAATGCCTTCAAGCCCTAGCAGCTACTGCCCTACTAGTACAAGAA  
AGAACTGAAAAACCCTCCTTCTTGGCCATCCTTGATAACCTTTGGCCAAAACCTAAA  
AATAAGCACACCCCATGCTCTCATGAAACTATTGCTTTCTGACCAATAAAGCAAAAAGGA  
CGGCGATGGCTAACAGACTCGAGAATTCTAAAGTATGAGGGGATTCTGTGATTGGGCCCG  
AGCAGGCCCTGAAGGCACAAGTGAATAATCATGATATAACAGTTGAACTTACAATACCCT  
GAACCCAGCCACCTTGCTATCTCTAGAACAGCATCCATTAGTGGCCTCAGAACAAAAATG  
TTTAGATTTAATTGAATACCAAACAAAAGTAAGACCAGATCTAAGAGAAACACCATTTCA  
TGACGGGTGATATCCATTTTTAGTTTTCTGCAGCACAGCCACAGGACCTCTTTATAGAT  
GGGTCTCCAGGTGATAGATGGAAAATGATACACTGATGGTTCTGCATTCAATGGAGAA  
AGACTCCCAAAACCAAAAAATAGAGTCATAAAAACCTTCTGGGTTGCCCCAGTATTAC  
AGGGGTGAAGGCAGGTCTGCTCAAACGTGAGAAGTGTGTTGCACTGAACCTAGACCTATAC  
CCCTTAAAAAATCAGGAAAGAACGCAGAACTAATTGCATTAACGAGAGCTCTGTTACTGA  
CTAAAGGCAAAAATGTGAACATATACTGATTCAAAGTATGCCTTTGCAGCGGTGCATA  
CATATGGAGAAATATGGAATGAAAGAGGCCTAATAAATACTAAGGGAAAAAGACATTATT  
TATAAAAAAGAAATCATGCAAGTATTAATGAACTACGGAAACCAGAAAAGATAGCTATT  
ATCCATGTAAGTGGACACCAAAAAAGGACCTCATTTGAAGGTGAAGGAAATAACAAGCA  
GATAAGGTGGCTAAACAAGCAGCTATGACTAAGCCTCAAACCTCCAGTCATACCCCTAACT  
CCCCATCTCCCTAAAGAACGACCTTAACCCAGTTTCTTAATCAATCATAAAGGAAAAAG  
AAGTATTTTCAAAAATAGGCGACAAAAAATAAAGAAGGATAATGGATGTTACCTGA  
CGGGAGAGAAACGATATCTAAACCCCTAACGACAAAATTAGTATAGCAACTACATATGTG  
CAAGGAACACACCCGGGGAAGCAAGGCTCTAGGTAATGCAGTAGTTCGGCATTATGGATTT  
ACACGGATTTATGCCCTCGCCAAAACGGTTACAGATAATTGCCTAATTTGTCAAAAAGAA  
AATCCAACACATGGTCCGAAATGTACCCCTGGAAGAAAAGATACTGGAACAAGACCATT  
TGAAAATCTGCAAGTTGACTTACAGAGATGCCTAAAATAGGAGGTCTACAATATTTACT  
AGTGCTCCTATTTACCTTCCCTGAAGAAGGGTAGAAGAAAAGATAGGAAATTTCTTTCCC  
CACCTGAAATGCAACTGCAAAAAAGATAATCAAAGCACTAATAAAAAACATTATACCCAG  
ATTTGGACTACCACAAAGCATTGATTACAGACAATGGAACCTCACTTTACTGCAAAAAGTAGT  
AAAGCAGTTGGCTCAAGTACTTAGAAAATAAATGGAAATGCACTAATCTGGTCTTACCA  
TACTCCCTACCATCCCCAGACATCAGGAAAGGTAGAAAGGATGAATCGGACTCTCAAAAA  
ACAATTAATAAATAAATTCAAGAACTTAATCATTGAAATGGGCCAAAAGTTCTCCAAT  
GGCCCTCCTTAGAATTAGATCTACTCCAATAAAGAAGAACTGGTTCTTCCACTTATGAGAT  
ATTATTCGAAAGGCCACCCCAATCATAAGTCAAATTAAGGTGATCTACAAGAGTTAGG



AGAAATAACATTA AATTATGTTTACCCTTAGAATAAGAGGCAAATGCAAAGCTTTAGGAA  
TAGCAATACAGAAAGTCCAAGGCTGGGTAAGAGAAAGAATACCTATAAGCCTAACAGACC  
CAGTACACCCATTCAAACCAGGGGACTCTGTCTGGGTCAAGAAATGGAATCCAACCACTT  
TGGGACCCATATGGGATGGGCCCATACTGTAATCATGTCCACTCCCCTGCTGTAAAG  
TTGCAGGAATCACACCTTGGATCCACCACAGCCGGCTGAAACCAGCAGCACCAGAAACCC  
CCGATGACGACAAGTGGACCAGCCAACAAGACCCAGACCACCCACCCGAATAATCCTAC  
GACGAAACCCAACCACCTAAGAGACGACAACAGCCCTGCTCCGACCACACCGGAAGCT  
GACCAGTCTACGCACGGCCGAAGCTTGAAGAGCCAACAAGCCCTGCTCTAGTCACACCCC  
GGAAGCTGACTAGTCTATGCATTGCCAAAGCTAAACTCCACTTCAAATGTCCCTAAAATG  
AAAATAAAAAATCTGATTCCCTATAACCTTCCCTGATAATACTAATTGTTCTACTATTACGC  
TACCACTGCAAATGCTGCAAACCTCCACCCCCAGAGAAAGACCTCCCATGACCCTGACTG  
GTATAAGCATGCTACTATTAACCAAAAGCTCACCACGGAAGAACTACACATAAACTGAT  
GATTTTACCATTACTGACTTCACGAAAAGAAGATGAAAATATACTTGGTTTATCACCTAT  
TATATCCTCCATCCAGCCTGGACAGCCATTAACAAAAATCGGCTGCTACCTAAGTTTTAC  
TATAAATGTACAGCAACCTTCCCATTAACCTGTATCTGCCTTAAGACCTTAGTTCCAGAA  
GCACCAACCTCCCATTCTGAGACCCAAAAAACACTAACCCAAATGTATGTTATGACCCT  
AATCTCTTACCATCAGACACCAGGTTTTGAAATAAAAAATAAAAACTAGCCAACCTCCGAGGG  
GATACAAATGGAATAAAAGAAGGAAAACCTTATAGATCAAACCAAAGAAGACCCTCCCTCC  
AATAAGAGGCCTATATCCTTGTACTTTGATGATTGCCATGCAACATATAAACATAAGCAT  
AACAAACCAAAGCAACCAGAAAATGTTTAAATCAAAAAACAACCTATCAATAGTGGCCCT  
AAAAATATAAGTTAAAAAAAACAGATTGCATCCCTGAACCGTTCCATTAACCTGGGAAAGG  
AGCTACAGAATAAAAAATAAATATGTTTGTAAAGAGTTCTGAGGTTCTGAGCTGGGCCCC  
CAGGTCTGAGCTCTGCTAAATCCGACCTGAAAATAACATTCAAAGGGCCAGTTTTTTCTA  
ACAAAAGGCAAAGACAAAAACTACTGTGCGAACAGCAGAAAGAGGCCGGAATAACCAAAC  
ACGAGTCAACCAACTCCTCCACAACGAAAAAAGGAAAAGAAGCAATGTAATAAAGCAAG  
TCCAGAGGCAAGCTGGAAGCATTTTTTTACTTTATGTTTGAACACTGCTGTTTTTTTT  
GTTTTATTTTCTTGAGTCAAGAAAACTTTTTCTTTTGACCAGACTCCCAGGTGAAACCAC  
TAATAAACGAAAACCTCAAAAAACCTAAAATGAACCAATATTCCAAATCTCTTATTATTA  
ACTAATTCTGCAAGCATCAGAAATTCTAAAAAATCAGAAAACCACATTTACAACAAACA  
TATTACTCGTTCTTCGTTTCTTCTCCAACCTAGCAGAAAATGGAGATTAATGTAACCAAG  
TCAATTCATGCTATGTATATGGAGTTTCAAGCATGGGAGATCAATCCCCGTGGGAAGTGA  
AGTTAATTTTGTAGAGCCAACAAAAAGTCCTTGGGAAAAACTGTCTGCACAATCCGGAGA  
GTCAATGACAGCTCTCTCAGTCAGAAAAGTACCTTTAGTCCTACTAACCCCTCCTCATGAA  
GGACAATAGAAAACCTACAACCATCTAGTCCTTAAAGAACTCCATAAAGGCAACATACTGC  
ATAGCTAGAGCCTTGAATGACCCCAACAAGGCCTATAGGAAGTATAACCTGAAGTCTCTGG  
CCTAGGCCAATAGTCTGTTGACGACGCAATAGAAAGAATGGACTGTCAAAAACAAGCTTTT  
ATATAATTACAAGCTACAGAACTAAAACAATAAATTCAATTAATCCAAGGGCCAACCCC  
CCCAAAAATTAATCCTACTAGAGGGACACCCAGACACAAACGGTGAATCTTTCTGAATAA  
AACCTATCCTTCTTTTTATCTTCTCTCAAACCCGCTTGGTACCACCCAGATTCTCAACCG  
GACTGGAAGGCTCCCTCTGGTATATACGGGATATGGTATAGGGGCAGAGGTTACAGGGAA  
AGACCCCATTTGGATTCTTTGAAAACCGTTTCTCAAATACACTAATCAAAAAATCAACCTC  
CAGCACCTGCCCCAAAAGATGTCTCCAACCACTCACAATAAAGAAAACCAAGCATTCC  
GCCACCATCTATCAAACGACCCGACCAAGATAACGATCGTAGAGGTTAAAGACTTAAAA  
CAAACCTTAGCAAGAAATTGAGACAGGATACCAAGATGTAATGCCTGGCTAAAATGGAT  
CAAATATTCCATCCGCACCTTGAACAAAAGGCATAGTTACGCTTGCGCACAAAGTAGACC  
GGAGGCTGTGTAATTGGCACAATTAATCCCTTTTCTTTATAATGCCCCCAAAAACAGGC  
GAACACATGCGGTACCTGTCTATGTTGAAAATAAAAAAAAAGAAAAACCAACCAGGAT  
AAGACAGCCTGGGGTGACAAGACATGCAAAGCTATAGCAATGCTGAATCACAAGAACGG  
CACCTGAAGGTAAGATCCCATGGGATAGCCACCTCCATCGCCACAAGAAAGATCAAGG  
GCACATCGTACCCCTATCTACAGGCACAAACCTATCATAACGGTTATAGGCGATCTGATTG  
AAATAATAACCAATGAATAAAGGAAAGCAGTGAGTTTACTAACTCGGAAAGACACACAAA  
TGAGTAAGCCCTTCTCATCCCCGAGCGGATGTATGGTGGTATTGTGGGAAATGCTACCT  
ATAAGAAAAAATTTGCCCTAGACCTCCTGCTGGACGTGTAACAGGTCCATTTGGAAATT  
AACTAACGTGTATTCACTTCTCACTGGTATTTGGTCACAAATAAAAAGAAAAAAGAAAGACA  
GTAAAAACAAGAGAAGCCACATATCAATCTTGGGACACGTACCTGTACATAGATGGCAT  
TGGAGTCCACGGGGACTACCTTTGGAATTTGGTGCCCAAAATAGATAGATGCAGAATT  
TTAATCAATGTTCTTCTGGTTGAAAAGGAACCTGCTTTATCAAACCTTTTTTTTTTTGTTT

TGATAAACTACATCTATTACAACCAACAGCGATTTATTA ACTACTAGAGATGCTCTTA  
AAGGAATAGCTGAACAATTAGGAGCTACTAGCCAGATGGCCTGGGAAAATAGGATAGCCT  
TAGACATGATATTAGCAGAAGAAGGAGGACTCTTTCTCATTATAAGAACCAAATGCTCTA  
CCTTCATAAACACACAACATCACGCCCAAACAATACCGCCCCTGATGGAAACATAACAAA  
AGCATTACAAGGTCTGACTGCTCTATCCAATGAGCTAGCCAAAAATTCTGGAATAAATGA  
CCCCTTTACAAAATGGCTAGAAAAGTGGTTCGGTAAATGGAAAGGAATAATAGCCTCAAT  
CCTTACATCTCTCGCAGTCGTAATAGGAGTACTTACTCTTGTAGGATGCTGTGTCATACC  
ATGTATCCGAGGATTAGTACAAAGGCTCATTGAAACAGCACTTACTAAAAAACCTCCCTT  
AACTATCCTCCACCTTATCCAATAAGCTGCTTACTTGCAACCTTTTAGCAGAATGAATC  
AGCATAATATTATCATTTGCATAAAAAATAAAAAATGTTTTAAAGAGATACGAAAAGAGAAA  
TAGTGAAAATGAGAATTACAACCAATGGTATAAATAGTAAAAAGGGGGGAAAAGGAAT

>4\_cons

GATGGTACCACAAGGGCCCCGAAAAAAAAGAAGGGCGCCAAGACGGGGATTAGGAGTCGAA  
ACAAACATTTTCAACCTCGCAATGGGGACTTCATCACACGTGGTAACTGCAGCATAGAAA  
GCATCCCACAAAACCTAGCACCCAGTCAACCTTGGGTACCGTGTGAGAGACACCGAAGG  
GACTCCCCCGCGAGGCAAGCGGCCCTCGCACCTCTCCCTACATGCAGTCTTTAATACC  
GCGGGACTCTCCCTCTAGAGAAGAACAGCTCACTATCCTGTCCGTGCCTGCCATAGAGAT  
GAAGGGCTCAAGATTGCTCAACTGCCCGCCAATTAGAATACCAGTGCCAAAATTAGGCTA  
GAAAAAGAAAAACTTCCAAGGACCTAAGCTATCGAACCTTCTTTTTTTCCTAAGCAACC  
ATGTGTTAGCGCCCTCCCTTAATCCAAGGTTAGGTGTCTTGATCAACCTTCGGTTGTGCT  
GTCAATTTAGGAGCTTTATAGTCGTTTCTATCCCTGGGGAAGGGCTCTTAACTATCCCC  
AACCTTTTCGCGTCTTAAGTAACGGTTTGCTACGAACAGGTTTTCTCTGCTCTTACTTTT  
CTGTGGAGCGATGTGCGCACCTGTGAGAAAAACATAACGTCTTTGAGTCACCCAGAAT  
TTCCAATTGAGCCAGAAGTTGAGCCCGATCTGTCCCTTGACATATCAGCTGACCTCTAG  
CGAGCTCGTACTATGGTTCCTCTGTCTTCTGCCTCCTGGGTCCAGATCATGACGAGCAC  
GAAAAGGTAGCCACAAAAGCACCAGGACCACCGAAGCGACAACCTCGCATCACAGACTGC  
CTGGAACCTAATGACTTCCCTCTCTCATCTCCTCTACAAGGTTATTCCTGCTAAGAAAAATC  
AAGAAGCCCTACGCAGAAGCCTTAAACACCTAGGCTAGGAACCAAAGATCCCTGTCCCT  
GGTGCCCTTCCAGATTTAGGCATAAGACTCAATTCAAGGGCAAATTTGAGGGACCAGTTC  
CCCACCATAGTGGACAGGCCCCACATCTGTAATGGCTAAGGGAAAAGAGAGACAGAGG  
AGAGAGAGAGAGACGGAGGAGAGAGAGAGAGACAGAGAGGAGAGAGAGACAGGAGAGA  
GACAAAACAAAGAGAGAGACATAGAGGAGAGAGAGAGAGTCAAAGAGAGAAAGAAAGAGA  
AAGAAATAGTAAAGAAAAACAGTGTGCCCTATTTTTTTAAAAACCACAGTAGCATTAGG  
GCCTATCATCAATTATTCCCAGAAAGACTTCCCCATAACACCAGGCCTCTCAAATACAA  
TCTTGTGTGTCAGTGTAACAAGGGCGTGGAGCAAGGGTACAGAGACCACAGACAATCAAT  
TGCTTTCCAATCAAAAATCCTTAACCCAGTAACCCGCGGATGGCCCAAATGCATTCAGTC  
AGTAGCGGCAACTGCTTTGCTAACAGAAGAAAGTACAAAATAACTATTAGAGGAAACCT  
CATGTGAGCACACCTCACCAGTTCAGAATTATTCTAAGTCAAAAAAGCAAAAAGGTAGC  
TACTAACTCAAAAATCTTAAAGTATGGGGCTATTCTGTTAGAAAAAGGTATAATAACTC  
CAACCACAGAAAACCTCCTTAACCCAGCAGATTTCTAACAGGGGATTTAAATCTTAATT  
ACCATACAAAGGTCCGACCAGACCTAGGAGGTACCCCTTCAGGACAGGGCGATAGATGG  
TACTCCTTCGTTATTGAGGAAAAACACCAAATAGGGAAACAAAAAGCAATAAAGACAC  
TCTCGCGAAAACAGAGGAAGAAAAATTGCCAAATAACAGCTATCAGCAAACACAGGAGAT  
GTTTCGCACACACCCTAACTTGAGAATACTACAACACCAATAAACTATAAAAAACATGAC  
GTAAATGATTCCAAATTCCTCTTTCAAGGTCAATTTCAATAAGAACAATAATTGATGGCC  
AGGCCTACAGAAAAGGCAGCTGAGCTGGCATGAATCATTTCGTACCCCCACCAGGACTA  
GTAACCACCTCGGAACGCAAAGCATGTCTCCATAAGGTGTAAGGACACCTAGAATCA  
AAAAATTCTAACGCCCTGTCTATGAGGCAAAATATGCGAGACAACAGGGACCAAACACTC  
TTGCAAATTCGTTTGCATGACTAAAGGCAGCTGATAAACAGATGATTGAAGGCAGTGAAG  
AAAACGTGTAAAAAATACACCTCAGGCCCATAGGGGACGCTCTAAGGGAACCTCAGACCC  
TAGACTAACTAAAAATGCGAGCATAACTATCTGTATACTTTTCAGATGAGAACTATACTC  
TCTGAGAACAATCTCCATTAATATGTATCCTTAACAATTGGGACAAATTCAAACCTGAAA  
TTTTAAAAAAAAGCGGCTGATATTCTTCTACAATACTGTCTGGCCCCAAAATTATCTTC  
TTTTAAGAATGAAAAACATGGCCACCTGACGGAAGTATTAATTATAACACTATTTTACAA  
ATAGACCTTTTTTTGTAAAAAAGAAGGCAAATGGAGTAACTCACATATGTCAATACAAAC  
TTTTCTTGCATTAATGACAATACTCAAATATGCAAAAAGGTTTTCCCTCTCGTCTTGT

ATGTCCTTGGGAGCTTGACTTTGTGACCATGTGGGGGTACTCTCTCTTGGTCTCCACCAT  
CCAGAAGAGGAAAAGAAGCAAAATAAGAAAAAGAACCAGACATAAAAAGAACAGCCAGTCT  
TGCAACTTATCCCCACAGACAGACCACCCAGCTCTCCCCCTTGCCCTAGTTCCTCCCA  
AACGCTCCTCCCCAACTAATGATGGTTCTTCTACAGCCCCCTTACCAAAGGGCCACAA  
GCAGATCGTCGATCGAATAAAATATAAGCCATACCCCTTCTAAAGACCCCCCTTTGGC  
CTAGGTGATCCTAAAGAAAACCTGTCAATCATAGAGAGTGGAATAAAAAAGCCTTAAAAA  
TTAAGGCAAAAGTAAAACATAAAACACCATAAAAAAGAAAAAGCACTTTAGAGAACTCATA  
GGCCAACTTTTACTACTTATGAAATATTATAGGTGTAATAATTGGCATAAGTTAAAAATAT  
CTAAAATTGTCGATTTTGTGTTGCGTGGAGAATTCAAGTGGTTTCATGCTTACACCAGAAG  
AAAAAATAAGACCTCAAAAAGAGCAACAAAGAACATAAAAACGACAAACACAGCCAAAA  
ACAATGGACAAAGAGCGTCTTATATGTGAAGGATAACAGGTACCTTAAAAACTGTTTAA  
AGCACAATAAGAAAAGGCAAGATAGAACACATATCCTAAAACCTGCCTCAACAAAAATGAT  
AGCCTCTCCCCTTGACATTGAAATTCAAAGAAAAAGCATAAAATACTAAGACAAAATCCT  
AAATAATAAAAAAAGCCTTCTAAGTGTCTTGTAACTTGTGTCAGTGCCAGAATAATAACTA  
AATTACCAGAAATTTATACCTCCCTTTTATCAAATTTGAAGAGATATAATTGTTTCAAGTTT  
TATCTTCAAACAAAAGACCAAAAAGAAACATGTATTCTATTTACAGAAACCTGAGAATTT  
AAAGATGCCTGGTATCTTAGTAAAGTCCATTATATACCTCAAGCCGAATAAAGAGTAACC  
TTAGAAACATTCAAAAAGTAAATCCCAAAAGTAGAACCCGACTGCCATATTGAATGAGAT  
AATGGTGATTGGAGCTGAAAAAATTAACCGAAAAAAGGTTAGTACTGTAAAAATAAGGT  
CACTCTCATTGTAGACGGGTTTATGAAAAATTGAAAAATCCTAAGAACTATTCAAAAAC  
AAACACTAATACTAAGGGAAAAGAAGAATATCCATTTGCTTTTTTTCGAAAAAACAAAGGA  
GCCGATGAGAAAAAATAGCAATCTGAAAACCTGAATCCATTGAATACTAAAAGATATAAAA  
AGATTAATCCAAAAGCCAATCAGCTGCATATATTATTAGTAACTTCAAAAAAATAATAT  
GAGCCCAAAAAAATTTTTTAAAAATAAAAGATTATTATTTCTACATTTAATGCCGTCTA  
GATTCTTGCCACCCTCAATGCCACAAGAGAGGACCTAAGGTAATTTCTGACAGCCTGGG  
ACTCCTTGGGAAAAACAGTGGAGGTGCCACAGACACTATTGAACCTGGCAACCTCGGTGT  
TCTATAATAGAGACCAAGAGGAACAGGCCAAAAGGAAAAAGCGAGACTAAGAAAAAGGCC  
GCAGCCTTAGTCATGGCCCTCAGACAAGCAGACTTTGGAGGCTCAGAGGGAACCAAAAGT  
GGAGCAGGACAATTGCATGGTAGGGCTTGTACCAGTGAGGTTTGCAAGGACACTTTAAA  
AAAGATTATCCAAGTAGAAACAACTGCCCCCTTACCCGTGTCCAATATGTAAAGGAAAT  
CACTGAAAGGAGCACTGCCCCAAGGAACAAAAGTCTCTAGGACAGAAGCAAACAACACTAG  
TTGTAATAACAGAAGAACTTTGGATATCTGGAGCAAAAAGATAAAAACAAGACATTACTATC  
ACAAAGACACTAGGATTCACTACCATTGAAGACCAGGAAAATGACTTCCCTCAGGGATACC  
AGTACGGCCTTATCAGTCTTACTCTAATATCATTGACCTCTGTCTTCTAAATCAATTACT  
ATCCTAGGGATATAAAGTAAACCCTTAAATGGTATTTCTCCAACCTCTTAAATTTGTAAT  
TAGGTAACCTTTGCTTTTATAAAAATCTTTTTCTTTTTATGTCTAAAAGTACCAACCTCTTA  
TTAGGGAAAGTTATAATAACATATCTTATAAAAAACAGTAGCTACTATCAAATGAACTTG  
TAGAAAAAATTAATAACTTAAAGTCTCTAACTTGAGAAAGGATATACTTTGTGATTTTA  
ATTTATATTTAATAAAAAGCTCTTATGGAAAACAAGGCTAAAAATATCAGTGAAGAAAAAT  
TCTAAAAGGCAAATAATGATCATCAAATTTCTATTTAAAATTAAGAATATATCACAATTC  
CTTAAAAAAAACAATAACTCCTTCTACCTAAAACCTTAAAAAAGATTTAAAAAGATAATAA  
AAAATATAATAACAAAGTATATTAAAAAAATTTAATTACATAAAATCTTACAAAATTAC  
AATCTTTTGAATACAAAATCTAACTGTAATCTATATGTCCTTAAAGTATAACATTATAA  
AAAACAGTATAAACTAAAAAAGTTACAAAAAAGCAGTACATTTATAACAAAAAACA  
AGATTTTAAAGATAAATATCTTAAAATAAAAAATAACAAAAATAAATTCTTGGACATAAAA  
TTCATATAATATTATAAAAAATGATTATATAAAAAATCTAAATTATGGAATACTTAAACCA  
TATACTCCGTTTTACATATAACAAAGAAAGTAGTATGGTTTATTATTAAAGATAATAAA  
GATTAATTTTTTATTATATACAAATATAATTACAATTAATAATTGTTTTATTTGAAACA  
TTGATTTGAAGATCATAAAAAACTGGCCTCTCATATAACTTGGAAAGTCTTACCACAAGA  
GTTGATAGATAGCCCAAACTAAATAGCCAACTTTAGTCCGGGATTTAAGCAAGTTCTC  
AAACATATTACCTAACCTCATATCATTCAATATGTGGATAACATACTTTTGGCTAACAG  
TTAAGAATAGACCTTGTACTTTCAACTACACATAGTAAACAAGTTACTCAAGAACCCTCT  
ATAAGAATAGAATCAACATTTGCTTAAAATTTCTAGATAATCAAGGGTACAAAATTTTAC  
TAAAAAGACTCAGATTTTCTTTTAAATAAGTAAATTATAAAGAATTAGGCTAATCGAAAG  
ACTTAAAGGGCCCTAAGCGAAGAAAGAATCAATCTTAAACTTGTTTATCCATGTTTCAAAA  
ACATGAAAAATTTTGAAGGGGATCTTTAAAAATAACTCAATTGTGTAACAATTATTCACA  
GGATACAGAGATATAATCCAGTCATTCTAGACCCTGAAAAAGGAAAAAATAAGCAAAT

CGGTATTTAATAGAATGGAAACCTTCATAAACCACATTCAGATAAAAAACAGTCCCCTAA  
ACACAGAAACAGCCTATAAAAAAGCTTATATAAATCTTACCCTAAGATCCAACATAACAAA  
ATAATTAATAGTTTAAAGGTTTAAAGGTTTATTAACAAGCTTGAGCGATATAAGACTA  
CCCACAGGACAAAACCTTACCTTTAGAATACAGCCATCCACATGAGGTAAAACCAGCAAAG  
TACAATAATACACCTATTTTCAAAGAAATTACAGGTATAACCTTACATGAAGTTAAAATA  
AAAGTTATATAACGTCTTAATATCCAGTAAAATAAATTAGTAGAAAATTTAAAGTAATAG  
AAATGGGTTGGCTTCAATGTTTTAAAAATTGTTTCAACAGTTTCAATCTTAATAATTGATA  
CATTTAAAATAATTCATGAAAAATATCTTACTTTGAGGATATTTTACGATGTTAAAGGTA  
TATTAGATGCTAAAAATAATTTGTGGCACAAGATAACCAATTACTTAAATATCAGTCTC  
TACTCCTTGAAGGACCAATATTTCAAATATGAAATTGTAAAGCTCTTAACCATGCCACTT  
TTCTCCCAGAGGATAAGGAACCAATTAACCATGGCTGCCCAAAAATTTTTTCTCAGACAT  
AAACCATTTCGTGTGGATTTTTTGGACGTCCTTATAGATAATCTTAACATCAACTAGAAAA  
CTGTA AAAAAGTGCTTTTTGAATAAAAATTTAATGCAAAAATTAGAGTATGTTATAGTTAGTG  
ATACGAAAATAATTGAAAATTAGTAATGAACTGAAATTATCAAACATAACAAAAGGAACA  
GAAATATAGATAATAAAATGAACCAAAAAGAAAGCAGAATAATAATTAATTAATTTGTAC  
AAAAATAGTTTAAATCGGTGAATTTTAAATTATATTTTAAATAAAAAGTTTACATTTCTAT  
TTCTTTATGCTATAAGCATCAATATTTCTTATATGTTATAGAGTCAAAAATAACTATCTAA  
TAATTA AAAATCTGACAGAAATTAAGTATTATGCAAAGAAATTA AAAATAGAGCATTCCCTG  
ATCCCCTCTAGGGGAACACCCATTAACACCACAAAGAAATTATGAAATTATTGCAGGCA  
CTACAAAAACCTAAAAAGGTGGCAGTCTTACACTACCAGCGTCATCAAAAAGGTGAAGGA  
TAAAAAGTAGAAGGAAACCATCAAACAGAAAGCAAAAACCAAAATTGCTGAAAGGGAGAAA  
CTTCTTTTAAAAATAAATATCAAAGGACCCCTAGAATGGGACAAACCCCGCAAGGAAAAA  
AATCCACAATATTCACCAATAGAAATAAAATTAGGGAGACTTTTACAATAACATAATTAC  
AAAAATAAATAAACCTCAACAAGGTTAACAAAAGAAGAAGGAAAAAATATTTATTCCTGC  
AGCTTTCTCCAGAATATAGAACTATTCGTGAGCATTCTTATCTTATGGCAATGTAATT  
ATTTGCATAAATCAAATAAGAATATGTTTATTTTTGTAACAGGACATAATTGTA AAAACC  
GGTTAGAATGACCAAGGCTTTCCCTGCAAGACATTGGCAAGTGTTAGAGGAAAAGAGGAA  
AGTTTATCCGAAAGTGTTATTGAAATCCCCCGGAGAAGAGACAATGGAAAGGTCTGTTGG  
ATGTTCCAACCTTTAAATTA AAAATCAAATTATCATTAGAAGGCAATAAAAAAAAATAT  
GAAAGCCATAAAAATAAGGTATTACCCGGGAAAGGACAGGGAATCAGACCTTACGCAGATT  
ACCAAAGTCACGATAATACAGGACAGGCACATAAAAAAATTCAGTATTTATTAGTCTGTG  
TTGATACCTTTACTGGTTGGATAGAAGCCTTCCCCTGCAAGACAGAGAAGGCACAAGAAG  
TAATTAAGCACTAATTCATGAAATAATTCCTAGATTTGGGCTTCCCAAAAAGCTTACAGA  
GTGACAATGGTCCAGCTTTTAAAGCCACAATAACCCAAGGAATATCCACGGCGCTAGGAA  
TACAATATCACCTTCACTGCGCCTGGAGGCCACAATCCTCAAGGAAAGTTCGAAAAGGCAA  
ATGAAACACTCAAAGGCATTTAAGAAAATAACACAGGAAGCCCATCCCCCTTCTCTTT  
CTCTAAACCTCACGCAGAGAGCAATTTAACCAACTCATAAGTTTTTTTAGGGGAAAAAGAC  
ATGCCTTTCTTTAGCAGGCACAGGATCCAAAAGATGCCCATAGACTTAGTCAACTAAGAA  
TCCTACAAGTCTTTAAAAAGCTCTTATAAAAAATTCATTATAAGAAAAAAAATTCATTA  
AAAACAAATAAAAAGGCATATTGGAAAAATAAGTATTCTTGCTAAACTTTTTGCAAAAAA  
TAGGTTTAAATAGTAATAAATAAAGAATTGTATCTCAATCCAAAATTCTAAAGTTTATGT  
TCCAGCAAAGCAAACCTTAAAGAGCCTATGTGGTCAGTCACTATTCTTGCTGCATTTAT  
GTAAATAATCAGGCCAAGTCTAATGAGATCAGACTTATTTTGCAAGCAACCATTTTTAAA  
CTATTGACAACATGCTAAAACAGGAAAAAACAATGTGGACACAAAAAATATAGCACACC  
TGTTGTTAAATACTAGTATTGCCTAAAATTTTTCAAACCTCAGCACTGGCTTTGGGCAAA  
ATATCAAAAAAAAAGCAGGCAAAATCTTAAGACGAAAATAAGCCTCAAAGAAGACCAGA  
GTACAATCATTACTAATAAATGTTATTGACGTCTTAATATCAAACAACACTATAAAAAGC  
ATGTCACAGAGATAAAGGTATAATCCTTTTTCAAACCTTCAAATTAATCAATTC  
AATTCTCTAACCCTTCATCCCCTGGAATCCCCACAAAATTAGGAGTCCCCCTCCTCAGA  
TTTGATGCAGACAAGGAATAATCCAATCTACCCCAAAAACCTCTACGTAGAGATAAAGA  
CTCAAACCATGCAGTTAAAATCGAAGGAGTGGAATCATGAAAAAACACACCTGACTAAA  
ACATTAGACACTGCCTGAGATAACTAATAAAAACATCAACCCCGTAAAGACAAGGTCAGCG  
AAAGAAAAATCAATACACATGGAAACCAAGGAGGACCTGCAGCCTCCTTTGAGAAGGCA  
ACCTACATGAAGGAAAAAATATAGTCTGTAAAGGTTTTCTATTCTTTCTCGGTTAGCTGT  
CTTCTCCTAAAAGAAAAAGAAATAAACTTGCTTTTTCTTTTACCTAATTCTGAGAATAATT  
AACGGCAACCTTAAAAAAGAAAAAGGAAAAAACAAAAAAAACAAAAACCACATTAATC  
AGTTACAGCTCTCATTGCTCTCTTAAATGGAAGACTTCTACTATTCATCACATTATTA

GCAGCATACTAACCACACTCCTTATAATAGGACTATATACTGTAGCTCCTGCCAGGATGA  
AAATCCTAATCACATCAACCTTCTTTCTATCAGCCTTCCTTGTAACAGCAATTCCTCCT  
ACCTTTAACTCAGCCTGGAAAAAATGATGTCATCTTCCAGAGCACCCCCCTTATCTTTCT  
ATTTACTATTTCCCTAACAACACCTCATGATTCCTTTCACACTTCCTGCAATCACTCCTC  
CCCATACTCCAGCCCCTAACTCCCCTACAAAACACTCCACTTGACAATGTGCCTCCCCG  
GAAAAACAGAACAACCATTCTATAGGAAATTATTGAAGGCTGGGTAGCCCCCTACAAACC  
CCCACATATGTTGCCACTCACACTCACATGAAAAACAAATGCTATAAACTACAACCTC  
TGCACTCATAATAAAACAAAAACACCTTATAAGACAAAAAATAAAAACTAGCTGAC  
CTAAAAACGGGGAACCAACACTTGTGGACATACTATACCCATACAGGTATGTCTAACA  
AAGGAGAAGTCCAACCTTTACGTTAATAAGACTAAAAAATAACATATAAAACAAGTAATC  
AAAAACCTAAACCAAGTACACAAAACACCAAGACCATATAAAAAATTAGACCTCTCAAAA  
CTACAAAAACCCTCAGTTATCATACTCGCCCCTGGAGCCTATTTAACACCACCCTTACA  
ATTTAAACCCTACAACCTTCAAGCCCCAACTGATCATAGTAACCTTCTGAGTCACCCAAACAG  
CTCCATTCAAATGGTTTGTCTGCTCAGGGCCCCCAGAATACATGAGGTCTCCCCTAATAA  
ACCAACAAACTGTTGGATGTGACTCCCCATGCATTTCCAGCCATACATTTCAATCCCTGT  
CCCCAAACAGTGGAACAACCACTACTAGTGAATGCCCTCTAAATTTCTTTTCAGTCACT  
CTCTCAAAAAGTACAAAACACATCCAAATTAGTAGGTCCCATAGTTACCAATATAGAAAA  
CACAGAGACCTCAAACTCACATGCATAAACTTAAAGCAAGACTATATACAAAAACATCTC  
CCAATAAATTTTCATAGGAAACAACACAGTCACGAATCGCCTGGCTAACATAAAATAAAAA  
TAAAGAAAAATCAACCAAGGAATATTATTCTTCTGTTATAACACAACCTATCGATACCAG  
AATAGTATCTTAATAAAACAATGCACCTTCGAAGCTCAAAATGCTTACATTTAAATATAT  
AAAGATCACAAGCCTCACTTATATTCTTGTCTCTGACACCGAAAATTGATCTGAATAA  
CTGACTATGAGAAACACGAATTCTTTACAAAATCAGACTACCAAAATATCCACCTAGAAAC  
TTGCACAAGATCAAGGCCACCATTGAACCCTTGAAGAGCCCCCATTCTAAAATTTATTGT  
TTGAGCGCCAATGATAGATACAATCCAAGGGAATAGAAGTGCAGGAATAACAACCTCAAC  
ATAGACTAAACACAACTATCAAAAACAATTAAGGAAACAATAAAAAAATAACCAAACA  
ATAACAGACACTACAAACACACCTTAACTCTGTGGAACATGTACCCTTCCGAAACTGCAG  
AGCCATGAACCTTTCCTTGGTGGAACAAAGAAGAACCTTTTCCTTGTCCGATATAATATG  
GGGAGATTTGATTTATTAACACTACCAAAAACAATAAACATAACACATAAATTTAAAGAAAT  
TCAAGAACGCACACTAAGTAGACAAAAGGAGCCTCATAACACAGAATCCTTCAACCTCCT  
TAGCAGATGGATACCATCGCTTCTCCCCCTTTTATAGGTCCTGTAACAGCAATCATATTGCT  
ACTCGCCTTTGGGCCATGTATCTTTAACCTCCTTGTCAAATTTGTTTCCTCCAGAATCGA  
GGCCATCAAGCAACAAAGGTTCTCACAATGGAACCACAAATAAGCTCAACAACCAACTT  
CAGAACACCAGACCAAGAATACCAAGGACCCGTGGACCCACCCCTAGCCCCCTCACCG  
AACTGAATATCCAAAAACGATACCCTCTCGAGGACATAACAAAGGCAGGGCCACTTCTTC  
GCCCCTATTTCAGCAGGAAGTAGTTAGAAAAAATCCGCAACCTCCCCAACAGCATTTG  
GCATTTCTGTTTAGAAGGGGGACATGAGGAAGAATGGGGATA

>5\_cons

AGTGGCGTCCACGTGGGGGCTCTCTGGCGCCCAAAAGGGGTTTCTCACTCACTGAGGGAA  
AGTGTGGGAGTCCATTTAGCCTGTCTACCTCATCCTCTCACTGCACTCCGAGTGATGACT  
AAGGACGGCAGCTCTTAAGCATCGGGGAAGAGTGGCCACATAGGGTATTTCCGAAACTC  
CCCTCGTTAAACTTTTACACACCGGTTGTGTTGTAGAAGTATTGTATAAAAAACCTACTG  
AAAACACGAGAAATAACAGTTATTTATAACACTACTATATAAACTAAAGGAAAGACAGG  
GTCATATCCACGCGGCTTCCAACAATCTGAACCTTTTATTTAAAAAGATGGAAAAATACT  
GTCCTTAGTTTCCTGAAAAAGGAACCATGGAGCTAAAAGTATGGGACCGAGTTGGTGCAA  
CATTCCGGCAACGGGTCACAGCAGGTAATTATCTTCCATCACTATTTGGAGTGAATGGG  
CCCTAATACGTGTTGCCTTACTTCCATACCAGTCCAGTGAACCCCTACAACCTACCACAAC  
TTAACGCACATGGCGACCCGCAGCCTTTACCTCAGATATCCACCCCACTCCGACTTCAC  
TTTCTGATCACCAATACAATACAATCTACCTCTCCTACCTACCCAAAAGGAGGAATCTA  
TGAATAACTCCCGAAACATCCCCTTAACCTCACCACTGAATATCTTAAATCTTTTCAA  
CAGAGCTGCTACTCCGAAACCAGCGGAACAGACTCAGCCATCCTGTGAACATCTAAATC  
CTCATTCTTCTACCCCAATCATCAGCACCCCTACTCTAAGCCTACTCCTACTAGCAACG  
AGACCAAACAACATATTTATAAACTTATACTGCCCTCCCCAAAGACTACAGCCCCAC  
ACCCTCCTAACCTTTCGCTCATTACCCGGCCACTGTTCAACCCATCCAACTACTAATC  
AGCACGCAACTTAAACATGAAGACAATAATCACCAGGAAGTTTAGGCCCTCCAACAC  
CCACAACCCCAAGTCTCAAACTCCAATACCGGTCCGACCTCCTCAACCTCAGTTTCCT

TATCTACACATACTTTTCTGTCACCTTCTATGCCGACTCCGTCTCATGTGCCTGCTCTTG  
AACTTCCATGCAATGCTTATTACGCGTAGGTAGGAACTGTGCGTCAGCCTGGCTTGAA  
GTCTCTCAACTTCTATATTCAAACCCCAATTTCTTTTCTGCTTCTGGTCCAGTCACAAC  
TGCTGTTGCTACCCATAAGCAATAGGTTACATACATTCTGATAATGACACCCCTCTTAT  
GAGGGCCATTCTCAGGGCAAGGGAATACGGGGATCCCGAGGCATGGTGTCTGTTATTCT  
ACAATCTCCTATACCTGCTGCCCCATTCTAGCTGCCCTGCTCTGGTGTCAATGGATCA  
GCCACCACCTGCTGACCAAGTTCAGCAGGCAGCTGACGCCACTGCCTCTCCAGACCCGCA  
GCTCAGGGGATCAGGCTCCTCAGCCAGTGCAAGAAGGGCCTGATGTCCAGCAGAGCCAG  
TTCCTGAAATTTGATTTAGAGGCGTGGCAGTACCCCGTCACACTACACCCCCAGATAAA  
CAAGCAAAAGACATGCGACAATATGAACCTTTCCCTTTAAAATTCTAAAAGAATTTAAA  
GATGCTTGAAATCAGTATGGACCAAATTTCTCCTTATGTCAAAACAGTACTAAAACCTTT  
GCTACTGAAAAACGATTGGTTCTATTGACTGGGACATTCTAGCAAAAAGCTGTTCTAACT  
CCATCTCAATACTTACAATTTAAGACATGGTGGGCAGATGAGGCCAGATTCAAGCTCCG  
CTAAATCAGGAAAATGAAACTCAAATTAATGTGACTACTGACCAGCTTCTGGGAGGGGGC  
GATTGGGGCGGCTATAAGTAACCAACAAATAGCCTAGGATAAACCCACTTTAGATCAGGTT  
ACCAGAACAAGTTAGTTAGGAGGAAAGGAAAAAATCCCCTTTGAAGGTCTTGCCTTTTTG  
CAAATAACAGCTATTAACAGGGTCAAATGAACCATACCCCTGATTTTCATGGCTCAAAT  
ACAAGATGCTGTGAAAAATCTATTCTGATACGAATGCACAAGATATAGTCTGCAAAT  
GTTAGCTTTTGAATGCTAATCCAGAGTGTGAGGCTGCTATACAATCTGTCCAACGTAA  
AACCCAACCAGAAAATGATTTGACCACTACCTATATCAAAAATAGAGCAGGTGTTGGTAG  
AACATCACAACTCATAAAGTTAAGTAGAGTCTTCTAAAGCTATCTGTTTCTCTATTTCC  
TTTTCTGCCTGCTTTGAATCTGCTGTTATTAAGCTACCGGTGTTGAGATAAAACTCACTG  
TTTATGGTACCGCTAGCCTCAGCAAAGAATATAGCCCACTTGACAAAAGCATTCTTTGA  
GCACAGGCAATGAAAGAACCCAAACAAAATAAAGCAAATAATTCTTTTCCAGATCAACC  
TGTTACATCAGGAAACAAGGTCATACTCGACAAGATCAAAAAACTGTAGCCTAAAAGACC  
GAAAGAAACAAGTCTTTACTTAATGCTAATCCCCAACAACTGCTCCTCAGAGACGGA  
CAAAAACGAATACCTGTGTGGAATGTATGCCAAGATGGAAAAAAGGAAAAACATTGGACA  
AATCATTGCTACTCTAAATTCGATATAAATGGTAACCCGTTACCGCAAATTCAGATAAAC  
GGAAAGGGCAGGAAGCCCCAGCCACACAACCAACAGGAGGAGGCACCAGCCTCAAGCC  
CCAATCCAATCAGGGCTCCGGGGTTTTGCTCCACAACCTCCAGCACCTCCCACTAAAACC  
AAAACAGCATTCCCGCATCAACCAGTCCAAACAAAGCCACAGACACAACCTCAAATACTC  
AAACCACAACCATATGCGTCTCAGCCCCCTTCTCTTATCCCAGTACAATGCCCGTCCACCG  
CCACAACAGGAGGTGCCGCAGTAGATCTATGCAGTACTATACCTATGACCCTACTACCTG  
GGGAACCCCTAAAATTGTCCCCACAGGAGCCAATGGCCCTTTACCTGGAACTTTAACTA  
GATAAATTTTGGCCAACCCCTGCTTAGCAACAAAAGGTGTTAAAGTTCATACCGGACTCA  
TTGATTCTGATTACTCTGGGGAAATAAAAATTGTTATTTCTACTAAAGTTCCCTTTAAAA  
CTGAAGCAGGAGAATGAATTGCTCAACTTCTGCTTCTCCCGTAACTGAAAATCGGTACAA  
ATAAAGGTAAACAAACAAGAGGCCTTGGGAGTACCAATAAACAAGGAAAAGCCGCTTATT  
GGGTAAATAAAAATTTCTGATAAACGGTCCGTGACCTGAAGGTAGAGACACTATAAAGGGA  
AAGAACCTCCATGATTTTCTAGACAGAGGAACTGATTTTTTCTATAATTTCTCCTCAGCAA  
TGGCCTTCCACCTGGCCAAAACAACCCGCAAAAATCAAATTAGTGGGAGTTGGAAAAGCC  
CCGGAAGTTTATCAAAGCTCTTTTATTTTGCATTGTACAGGCCAGATGACCAAATGGGA  
ACAATTCAACCATATATAACTCTTTCCCATGTAATCCAGGGAACAGTGCTGCACTACAA  
CAATGGGGGGCGGACATGACCAAACACAAGAAAATCACAGGCTATGAAAATGGTTAACAG  
GAAAATAACTTATAATCCTGACCAGCTTTGTCTAAGGCCTTGGTTCTCAAAATCCACATA  
AATTAATACTAAGCAAAAACAAAATAAGTTAAAAAAAATGAGATATATGCCTGGAAGAG  
GACTAGGAGAAAATTGGCAAGGGATAAAAAGAACCCCTGCAACTCACCAAAAAACTTGACA  
ACTAAGGATTTGGATATCCTTTTTAGTGGCGGCCATTGTCAAGCCTCCAGACCCTATCCC  
TTTAAAATGGATATCTGATAAGCCAGTTTGGATAGAGCAGTGGCCGCTTCTTAAAAAAA  
ACTGGAGGCTTTAAATAAATTAGTTAATGAACAATTAGAAGATGGACACATTGAGCCATC  
TTTCTCTCCATGGAATTCACCTGTGTTTGTAAATACAAAAAAAATCAGCGGAAAATGGAGA  
ATGGTAACTGACTTAAGAGCCATTAATGCAGTAATTAACCTGGACGGTACCACCCAGG  
GGCACGTACAACCCGGCATGCCCTCCCCGCTATGATCCCTAAAATTTGGCCTCTAATAC  
TCATAGATCTTAAAGATTGCTTTTTTAATATTCCTTTAGACAAGCAAGACTGTGAAAAT  
TTGCTTTTACTGTACCTTCAATCAACAATCTGGAGCCTGCAACTCGTTATCAATGGAAAG  
TACTACCACAAGGAATGCTAAACAGTCTACAATTTGCCAGCCTTATGTTGGGCAAGTGC  
TTCAACCTGTCCGACATAAATTTCCACAGGGTTACATTCTTCATTATATGGATGATATAC

TTTGTGCTGCCCCACTGAAGAAGAATTAATTCAGTGTTCCTTCTTGAAACAAGCCA  
TTTCAGAGGCTGGATTAACATAGCTCCAGATAAAATTCAAAATACCACTCCTTTTCAAT  
ATTTGGGAATGCAGGTAGAAGACAAACCCATTAAGCCACAAAAAGTCCAACCTAGTAGAG  
ATAATTTAAAAACCTTAAATGACTTTCAAAAATTACTAGGTGACATTAATTAGATAAGAC  
CTACTTTAGGCATCCCTACATATGCGATGTCTAACCTGTTTGCACACTATGTGGAGATC  
CAAATCTAAACAGTCAAAGGCCTCTAACAGAAACCGCAGACTAAAAGAGGCTAAACCAGA  
GTTGCAATTGATGGAAAAAAGAGTCCAAAAGGCTCAAGTAACTAGAATAGATCCAAATTA  
GCCTTTACATTTTCTAATTTTTTCCAACCTCAGCACTCTCCTACGGGACTAATAGTTCAACA  
GCATGATCTAGTTGAATGGGGTTTTTCTTCTCATTCCACTTCAAAAACCTAACTATTTA  
TCTGGACCAAATCGCCACCATAATTGGGCAAGCAAGATCTCATATTATTAATAATTTACGG  
ATATGATCCTAAAAAATTATAGTCCCTTTAAAAACAACAACAATAACAACAAGCCTTTAC  
AAATTCTCTTACTTGGCAAATAAATTTGGCTGACTTTATTGGCATTATTGATAATCATT  
GCCTAAAAAAAATTGTTTCAATTTCTAAAAATAACTTCTTGGATTCTACCTAAAAATAAC  
CAAAGATAAACCAATTACAGGAGCCGTTACAATGTTCACTGATGGGTCCAGTAATGGAAA  
AGCGGTCTACGTCCACCAAAAACCAAGCAATCCACACAACATCTGCCTCCTTTGAAAT  
ATCATATATAATAAACATAGAAGAAAGAAGGGGGGGCTTACTGCCGTTTCTGAGCCATTC  
AAGGAGATTAATATACCCCTAAAAATTGTCTCTGATTCTGCATATGTAGTACATGCCACT  
AAGAAAATAGAAACAGCTACCATCAAATATATTGCTGATGAAAAACTGATTTCTTTATTT  
CCAAGGTTACAAACGGGACCTAGGAACCTTAGTCACCACCCCTTAAGCCGCCACTAAA  
AACCTGCCCCATACCCATCCGCCCGAAACCGGTCTGCTGGCAATCATAAAGCTGATGCT  
CTAGTCTCTTCCGCAATTAAGAAGCAGCAGACTTTCATAATCTCACTCATGTCAATGCC  
GCAGGACTCAAACACAAATACCCTCTCACATGGAAAGAAGCTAAACATATTGTACAGCGC  
TGTTACATTGCAAAGAGAGAATGGGAAAAACAACAAGCAGGCAATGAAGCATAACAACC  
ACCACCTTCCAGACTCAACACCAAAACAACCCCGCCACGAGTCAATCCCACAAACAATTC  
CGTCAAGTGCTAATCCTACCAACTCTGGCTCCAGGAGTTAATCCCAGAGGCTTGGCACCT  
AACGCTCTTTGGCAAATGGATGTCACCCATGTTCCATCTTTTGGAAAGACTAGCTTATGTA  
CATGTATCAGTAGACACCTTTTTACATTTTATCTGGGCTACATGCCAAACAGGAGAAGGC  
ACTGCCCGTGTTAAAAGACATATGTCTTCTGTTTTGCGGTTATGGGCATTCCACCTCAG  
ATTA AACAGACAACGCCCCAGGCTATAACCAGCAAAGCTTTTAAAAAATTTATTCAACAA  
TGGAATATTAACCGCACTACTGGAATCCCTTATAAGCCCCAAGGACAGGCTCCAGTAGGA  
GTGAGCAAATAACTTCCAAAAACAACAGTTACAAAAACAGAAAGAAAGAAAAAAGGA  
ATTAAGTACCCCCACAAGCAATTAATCTGGCACTTCTGACTCTGAATCCTTCCATTTT  
GTCAAAGCTCCGTCCTCTAATGGCAGCCGAACAACACTATAACAGGCAATAAATTTTTTCG  
GCACCTAAACAAGAAATTAAGAAACAAGCACAAAAAAGAAAGAAAAACAATGACGCAC  
TGGAAGAACAAGAAACA AAAAGTTGGCAAATAGCTAAAAAAGAAAGTTCCGGGCCACTG  
CCGGTGTTCTCTCCTTGATGTCTGGGAGCAGACCAGATTGGTAGACAAACACGAAAACCA  
AAAACCTAGGCAAGAGGAACAATCAGAACACGAGAAACAGGGTATGCTACAGGTCCACCAG  
AAGAAAATCAATCCCCTGTTTGGGTCCCTACTAGAAATCCCTGAGTCCGTCTGAAGAATG  
ACAATGAAAAACAACAAGAAAAAACAATCAGGGCCACAAACCACCCGCAAACATAGCCAAA  
ACTGGGCAAAAAGAACAAAAACAGACACGACAAACCAAAAACCATATAACCCAACAAGAC  
CAGAGCACAAAGAAAACATAGAAATTTAATCCCAAAACCCCAAATTCGTCACGATCAAC  
ACACGCAAACATAAAAAACTAACCTAGGAAAAGACCAAGAAACATACACAGCTCAAAGCA  
ATCCACCAAACCAAAGCAAAGAATCTCGCCCCAACGATGATTGTCATGATCACTCTGATA  
CGCATAATCAACACTGCAGTAACTCTCCCTTACACCAAGCTGCATACAACAATAAATCTG  
TCTCAGTGGACTTCTTTGCCTTTTCTCCTCACTTATTTCGACCCATCACATGGATGGATGCT  
CCTGTAGAAGTCTATACTAACGATAGTGCTTGCATGCCTGGATCTATAGATGACCGTTGT  
CCTGCTCAACCAGGAGAAGAAGGAACGCCTTTTAATGTTACCATTGGATATAAATATCCA  
CCTTTGTGCCTGGGACATGCACCTGGTTGTATCCCATAGATAATCAAAATTGGCTGGCG  
ACACTACCAGCCGGCAACACTGATACGAAATAGGGACATATGGTCTCAGATCTCACAATT  
AAACCTTTAAGATATACTATTACGGGTGTGGCAGACTCACTCAAAAATCTCAATATAAG  
CCAATAGGAACCACGCCAGAGCAGACGAACCTTGGCTCGCAGTGCCAGACCCCTAAAAAGAC  
CAAAAAAAGGGAATAAAAAACTAAAAATTTAATATGGAAAGATTGCATTAACGCACAAGC  
AGAAGTGCTAAAAAATGATTCCCACAGAATCATTATTGACTGGGCCCAAGGGGCATTT  
TAGGAATAATTGCTCTGCTCAGCAAACACAATGTCAGGAGGCTACCTATTTTATTGCTTA  
TTAAGAGAATAGCGACCACCCTCACATATTAAGGAAAGGTTGACCACATTCTGTCCCTC  
TAATTGGAAAAATAAAGGCATTGCCTGCATGAGACCAGGAGCCAGGGTCCGGCTCCGAGG  
GACAGAAAAAATAGAACCAAAAAAGAACCCTCAACAGTTAGAAATATGGAAATTGGCTAT

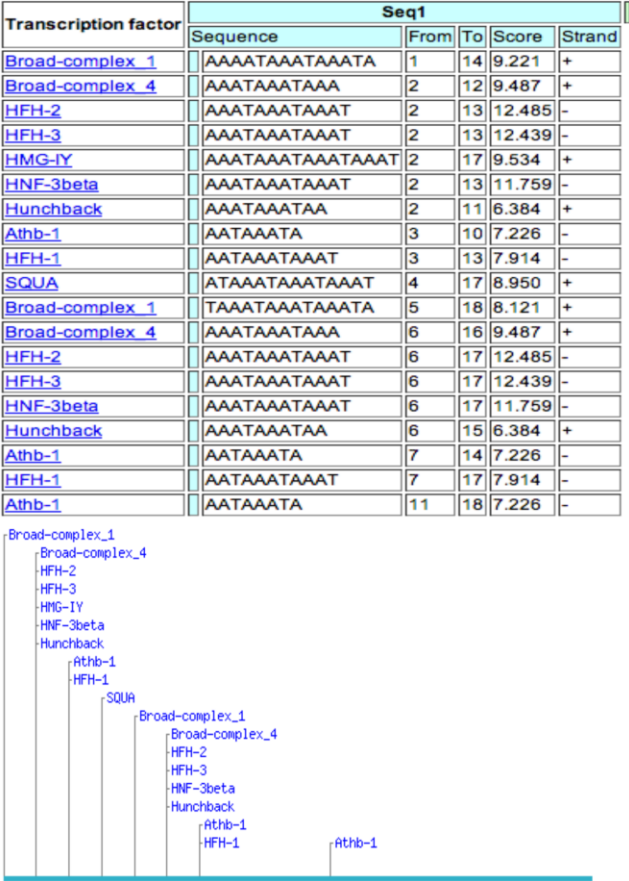
AGCCATATCCGGAATCAGAGTATGGGAAGGTGATAATAATAAATCTATTATAACAACAACTAA  
AACACAAAAAAAACAGATTTCTTTACACTATGATAAACACACACCCAAAAAATATCACTAT  
GGCAAACACCAAAAACGCCAATACAAAAATCCGACAGGAAAGACGATAAGGACTAACAACC  
CACCGATGCCAATCTACCCCCACAACATCCCAAACAACCTAGGACCCCCACCCACAGTC  
AAAAAAAAGAGGAATTGCCACCGCCGCTCCTCTCCCTCAGTATCAACCCAAAAACAGACA  
TACTACCTAGATGAACTCCAACAAAACCTATACCAATAAAAAAGTTGTGTTAAACCACCATA  
TATGTTATTAGTAGGAAACATAAATATTAGCACAAAAAATCAAACCTATTAATGCATTAA  
TTGTAAATTGTATACTTGTATTGACTCAACATTTGATCCAAAAAAAAGTGTATAATGGT  
CAGAGCCAGAGAAGGAATATGGATACCAGTAACTTTACACAGACCTTGGGAATCCTCCCC  
TTCAATCCATTTTATTAAAAAAATTCTACAGAGAATTCTTAAAAGAACTAAGAGATTTAT  
TTTTACTTTAATTGCAGTGATAATGGGCTTAATTGCTGTTACTGCAACAGCCGCTACTGC  
TGGAGTAGCATTACATCAATCTATTCAAACCGCTCATTTTGTGGATAAATGGCAAAAAAA  
TTCCACCCGAATGTGGAATTCTCAGCCAGGCATTGATCAAAAATTGGCCAATCAAATTA  
TGATCTAAGACAGACTGTTATATGGATGGGAGATAGGATAATGAGTTTAGAACATCGAAT  
GCAAATGCAATGTGATTGGAATACTTCTGATTATTGTATAACACCATATAGATATAATGA  
GAATCACCACAGTTGGGAAACAGTAAAAAGCCATCTACAAGGAAGCGATGATAATTTATC  
CTTAGACATAACAAAACATAAAGAACAATTTTTGAAGCCTCCCAAGCTCACTTAACTAC  
TATACCCGGAGCTGAAGTGTTTGAAGGAATCGCAGAAAGATTATCTGATCTAAACCCCAT  
TAAATGGATAAAATCTCTTGGAGGCTCCATTATTGTAAATATTGTAAGTATTTAATCTG  
TTTTATTTGTTTGTGTTTTAGTCTGCAGAACTCGACACAGAATCCTACGATAAAATCGTGA  
CCAGGACCAAGCCATCATCGCAATTGTTGACTTAGCAAAAAAGAAACGGCGACAGAATAT  
AAGAGCTGGGGAACGGGGTGGAT



Supplementary Material S3:

Figure 18: The regulatory motifs of 17bp subsequence of LTR in lincRNA.

LTR



A

Figure 19: The regulatory motifs from two subsequence of ERV in lincRNA.

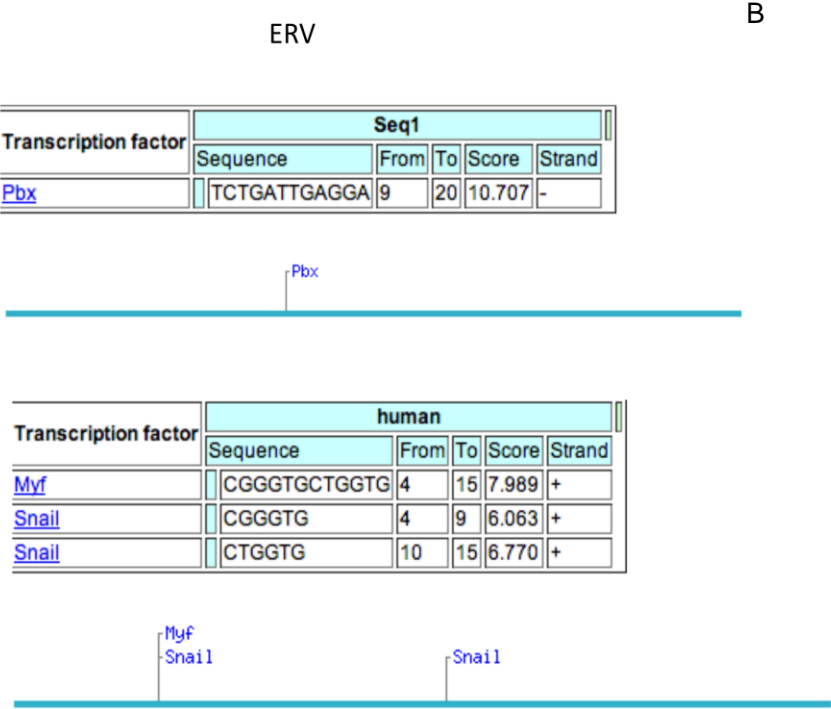


Figure 20: The regulatory motifs of 150bp subsequence of MIR in lincRNA.

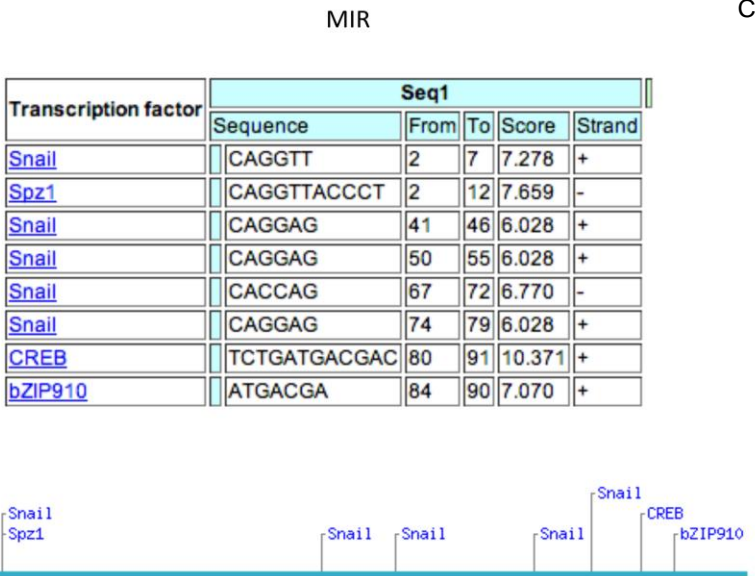
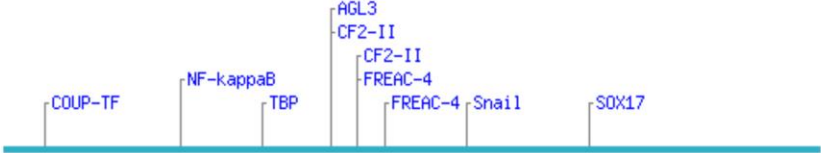


Figure 21: The regulatory motifs of 100bp subsequence of SVA in lincRNA.

SVA

D

Transcription factor	Seq1				
	Sequence	From	To	Score	Strand
<a href="#">COUP-TF</a>	CTGTTCAAGGGGCC	4	17	8.853	-
<a href="#">NF-kappaB</a>	GGCCAGTCCC	14	23	7.528	-
<a href="#">TBP</a>	TCCCTCCATGTATAC	20	34	7.528	-
<a href="#">AGL3</a>	CCATGTATAC	25	34	5.939	+
<a href="#">CF2-II</a>	CCATGTATAC	25	34	8.425	-
<a href="#">CF2-II</a>	ATGTATACAA	27	36	6.185	+
<a href="#">FREAC-4</a>	ATGTATAC	27	34	8.269	-
<a href="#">FREAC-4</a>	GTATACAA	29	36	8.102	+
<a href="#">Snail</a>	AACCTG	35	40	7.278	-
<a href="#">SOX17</a>	ACCATGGTC	44	52	6.411	+



# Reference

1. Belancio, V.P., D.J. Hedges, and P. Deininger, *Mammalian non-LTR retrotransposons: for better or worse, in sickness and in health*. *Genome Res*, 2008. **18**(3): p. 343-58.
2. Mirkin, S.M., *Expandable DNA repeats and human disease*. *Nature*, 2007. **447**(7147): p. 932-40.
3. Kunarso, G., et al., *Transposable elements have rewired the core regulatory network of human embryonic stem cells*. *Nat Genet*, 2010. **42**(7): p. 631-4.
4. Lynch, V.J., et al., *Transposon-mediated rewiring of gene regulatory networks contributed to the evolution of pregnancy in mammals*. *Nat Genet*, 2011. **43**(11): p. 1154-9.
5. Cowley, M. and R.J. Oakey, *Transposable elements re-wire and fine-tune the transcriptome*. *PLoS Genet*, 2013. **9**(1): p. e1003234.
6. Pereira, V., D. Enard, and A. Eyre-Walker, *The effect of transposable element insertions on gene expression evolution in rodents*. *PLoS One*, 2009. **4**(2): p. e4321.
7. van de Lagemaat, L.N., et al., *Transposable elements in mammals promote regulatory variation and diversification of genes with specialized functions*. *Trends Genet*, 2003. **19**(10): p. 530-6.
8. Rebollo, R., M.T. Romanish, and D.L. Mager, *Transposable elements: an abundant and natural source of regulatory sequences for host genes*. *Annu Rev Genet*, 2012. **46**: p. 21-42.
9. Jacques, P.E., J. Jeyakani, and G. Bourque, *The majority of primate-specific regulatory sequences are derived from transposable elements*. *PLoS Genet*, 2013. **9**(5): p. e1003504.
10. Conley, A.B., J. Piriyaopongsa, and I.K. Jordan, *Retroviral promoters in the human genome*. *Bioinformatics*, 2008. **24**(14): p. 1563-7.
11. Medstrand, P., J.R. Landry, and D.L. Mager, *Long terminal repeats are used as alternative promoters for the endothelin B receptor and apolipoprotein C-I genes in humans*. *J Biol Chem*, 2001. **276**(3): p. 1896-903.
12. Franchini, L.F., et al., *Convergent evolution of two mammalian neuronal enhancers by sequential exaptation of unrelated retroposons*. *Proc Natl Acad Sci U S A*, 2011. **108**(37): p. 15270-5.
13. Sela, N., et al., *Comparative analysis of transposed element insertion within human and mouse genomes reveals Alu's unique role in shaping the human transcriptome*. *Genome Biol*, 2007. **8**(6): p. R127.
14. Schoeftner, S. and M.A. Blasco, *Developmentally regulated transcription of mammalian telomeres by DNA-dependent RNA polymerase II*. *Nat Cell Biol*, 2008. **10**(2): p. 228-36.
15. Azzalin, C.M., et al., *Telomeric repeat containing RNA and RNA surveillance factors at mammalian chromosome ends*. *Science*, 2007. **318**(5851): p. 798-801.
16. Nilsson, M.A., et al., *Expansion of CORE-SINEs in the genome of the Tasmanian devil*. *BMC Genomics*, 2012. **13**: p. 172.
17. Kelley, D. and J. Rinn, *Transposable elements reveal a stem cell-specific class of long noncoding RNAs*. *Genome Biol*, 2012. **13**(11): p. R107.
18. Mc, C.B., *The origin and behavior of mutable loci in maize*. *Proc Natl Acad Sci U S A*, 1950. **36**(6): p. 344-55.

19. Kramerov, D.A. and N.S. Vassetzky, *Short retroposons in eukaryotic genomes*. Int Rev Cytol, 2005. **247**: p. 165-221.
20. McCarthy, E.M. and J.F. McDonald, *Long terminal repeat retrotransposons of Mus musculus*. Genome Biol, 2004. **5**(3): p. R14.
21. Singer, M.F., *SINEs and LINEs: highly repeated short and long interspersed sequences in mammalian genomes*. Cell, 1982. **28**(3): p. 433-4.
22. King, R.C. and W.D. Stansfield, *A dictionary of genetics*. 5th ed. 1997, New York: Oxford University press. vii, 439 p.
23. Kajikawa, M. and N. Okada, *LINEs mobilize SINEs in the eel through a shared 3' sequence*. Cell, 2002. **111**(3): p. 433-44.
24. Jurka, J., *Sequence patterns indicate an enzymatic involvement in integration of mammalian retroposons*. Proc Natl Acad Sci U S A, 1997. **94**(5): p. 1872-7.
25. Cordaux, R. and M.A. Batzer, *The impact of retrotransposons on human genome evolution*. Nat Rev Genet, 2009. **10**(10): p. 691-703.
26. Chen, C., T. Ara, and D. Gautheret, *Using Alu elements as polyadenylation sites: A case of retroposon exaptation*. Mol Biol Evol, 2009. **26**(2): p. 327-34.
27. Roy-Engel, A.M., et al., *Human retroelements may introduce intragenic polyadenylation signals*. Cytogenet Genome Res, 2005. **110**(1-4): p. 365-71.
28. Shen, S., et al., *Widespread establishment and regulatory impact of Alu exons in human genes*. Proc Natl Acad Sci U S A, 2011. **108**(7): p. 2837-42.
29. Sela, N., et al., *Characteristics of transposable element exonization within human and mouse*. PLoS One, 2010. **5**(6): p. e10907.
30. Vorechovsky, I., *Transposable elements in disease-associated cryptic exons*. Hum Genet, 2010. **127**(2): p. 135-54.
31. Chen, L.L., J.N. DeCervo, and G.G. Carmichael, *Alu element-mediated gene silencing*. EMBO J, 2008. **27**(12): p. 1694-705.
32. Levanon, E.Y., et al., *Systematic identification of abundant A-to-I editing sites in the human transcriptome*. Nat Biotechnol, 2004. **22**(8): p. 1001-5.
33. Hohjoh, H. and M.F. Singer, *Sequence-specific single-strand RNA binding protein encoded by the human LINE-1 retrotransposon*. EMBO J, 1997. **16**(19): p. 6034-43.
34. Kolosha, V.O. and S.L. Martin, *In vitro properties of the first ORF protein from mouse LINE-1 support its role in ribonucleoprotein particle formation during retrotransposition*. Proc Natl Acad Sci U S A, 1997. **94**(19): p. 10155-60.
35. Feng, Q., et al., *Human LI retrotransposon encodes a conserved endonuclease required for retrotransposition*. Cell, 1996. **87**(5): p. 905-16.
36. Mathias, S.L., et al., *Reverse transcriptase encoded by a human transposable element*. Science, 1991. **254**(5039): p. 1808-10.
37. Fanning, T. and M. Singer, *The LINE-1 DNA sequences in four mammalian orders predict proteins that conserve homologies to retrovirus proteins*. Nucleic Acids Res, 1987. **15**(5): p. 2251-60.
38. Weiner, A.M., *SINEs and LINEs: the art of biting the hand that feeds you*. Curr Opin Cell Biol, 2002. **14**(3): p. 343-50.
39. Deininger, P., *Alu elements: know the SINEs*. Genome Biol, 2011. **12**(12): p. 236.
40. Team, R.D.C., *R : A language and environment for statistical computing*. 2010, R Foundation for Statistical Computing.
41. Vance , A., *Data Analysts Captivated by R's Power*. New York Times, 2009.
42. Dreszer, T.R., et al., *The UCSC Genome Browser database: extensions and updates 2011*. Nucleic Acids Res, 2012. **40**(Database issue): p. D918-23.
43. Kent, W.J., et al., *The human genome browser at UCSC*. Genome Res, 2002. **12**(6): p. 996-1006.

44. Pruitt, K.D., T. Tatusova, and D.R. Maglott, *NCBI reference sequences (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins*. Nucleic Acids Res, 2007. **35**(Database issue): p. D61-5.
45. Ernst, J., et al., *Mapping and analysis of chromatin state dynamics in nine human cell types*. Nature, 2011. **473**(7345): p. 43-9.
46. Meyer, L.R., et al., *The UCSC Genome Browser database: extensions and updates 2013*. Nucleic Acids Res, 2013. **41**(Database issue): p. D64-9.
47. Edgar, R.C., *MUSCLE: multiple sequence alignment with high accuracy and high throughput*. Nucleic Acids Res, 2004. **32**(5): p. 1792-7.
48. Price, M.N., P.S. Dehal, and A.P. Arkin, *FastTree: computing large minimum evolution trees with profiles instead of a distance matrix*. Mol Biol Evol, 2009. **26**(7): p. 1641-50.
49. Han, M.V. and C.M. Zmasek, *phyloXML: XML for evolutionary biology and comparative genomics*. BMC Bioinformatics, 2009. **10**: p. 356.
50. Talavera, G. and J. Castresana, *Improvement of phylogenies after removing divergent and ambiguously aligned blocks from protein sequence alignments*. Syst Biol, 2007. **56**(4): p. 564-77.
51. Huang da, W., et al., *The DAVID Gene Functional Classification Tool: a novel biological module-centric algorithm to functionally analyze large gene lists*. Genome Biol, 2007. **8**(9): p. R183.
52. McCue, A.D. and R.K. Slotkin, *Transposable element small RNAs as regulators of gene expression*. Trends Genet, 2012. **28**(12): p. 616-23.
53. Cabili, M.N., et al., *Integrative annotation of human large intergenic noncoding RNAs reveals global properties and specific subclasses*. Genes Dev, 2011. **25**(18): p. 1915-27.
54. Altschul, S.F., et al., *Gapped BLAST and PSI-BLAST: a new generation of protein database search programs*. Nucleic Acids Res, 1997. **25**(17): p. 3389-402.
55. Belancio, V.P., D.J. Hedges, and P. Deininger, *LINE-1 RNA splicing and influences on mammalian gene expression*. Nucleic Acids Res, 2006. **34**(5): p. 1512-21.
56. Belancio, V.P., A.M. Roy-Engel, and P. Deininger, *The impact of multiple splice sites in human L1 elements*. Gene, 2008. **411**(1-2): p. 38-45.
57. Cenik, C., et al., *Genome analysis reveals interplay between 5'UTR introns and nuclear mRNA export for secretory and mitochondrial genes*. PLoS Genet, 2011. **7**(4): p. e1001366.
58. Cenik, C., et al., *Genome-wide functional analysis of human 5' untranslated region introns*. Genome Biol, 2010. **11**(3): p. R29.
59. Barrett, L.W., S. Fletcher, and S.D. Wilton, *Regulation of eukaryotic gene expression by the untranslated gene regions and other non-coding elements*. Cell Mol Life Sci, 2012. **69**(21): p. 3613-34.
60. Faulkner, G.J., et al., *The regulated retrotransposon transcriptome of mammalian cells*. Nat Genet, 2009. **41**(5): p. 563-71.
61. Kim, D.D., et al., *Widespread RNA editing of embedded alu elements in the human transcriptome*. Genome Res, 2004. **14**(9): p. 1719-25.
62. Cohen, C.J., W.M. Lock, and D.L. Mager, *Endogenous retroviral LTRs as promoters for human genes: a critical assessment*. Gene, 2009. **448**(2): p. 105-14.
63. Moran, J.V., R.J. DeBerardinis, and H.H. Kazazian, Jr., *Exon shuffling by L1 retrotransposition*. Science, 1999. **283**(5407): p. 1530-4.
64. Costa, F.F., *Non-coding RNAs, epigenetics and complexity*. Gene, 2008. **410**(1): p. 9-17.

65. Martianov, I., et al., *Repression of the human dihydrofolate reductase gene by a non-coding interfering transcript*. Nature, 2007. **445**(7128): p. 666-70.
66. Caretti, G., et al., *The RNA helicases p68/p72 and the noncoding RNA SRA are coregulators of MyoD and skeletal muscle differentiation*. Dev Cell, 2006. **11**(4): p. 547-60.
67. Siomi, M.C., et al., *PIWI-interacting small RNAs: the vanguard of genome defence*. Nat Rev Mol Cell Biol, 2011. **12**(4): p. 246-58.