

# ESTIMATING EVOLUTIONARY PARAMETERS FOR LOW COMPLEXITY REGIONS USING AN APPROXIMATE BAYESIAN COMPUTATION

ALEXANDER TURCO

February 5, 2023

<sup>1</sup> Department of Biology, McMaster University, Hamilton, ON, Canada

# Contents

<b>Abstract</b>	<b>3</b>
<b>Literature Review/Proposal</b>	<b>4</b>
What are Low Complexity Regions? . . . . .	4
Characteristics and Types of LCRs . . . . .	4
Why care about LCRs? . . . . .	5
How do LCRs Evolve? . . . . .	6
What is an Approximate Bayesian Computation Markov chain Monte Carlo algorithm? . . . . .	7
How will we use an ABC-MCMC - Oct 28 First draft . . . . .	8
<b>Materials and Methods (mid-year stuff Jan 20)</b>	<b>9</b>
ABC-MCMC: The Algorithm . . . . .	9
Parameters and Summary Statistics . . . . .	10
Simulation Step: Creation and Mutation of Protein Sequences . . . . .	11
Normalization and Euclidean Distance Calculation . . . . .	12
<b>Results</b>	<b>12</b>
<b>Discussion</b>	<b>12</b>
<b>References</b>	<b>13</b>
<b>Appendix C++ simulation (so far)</b>	<b>16</b>

## Abstract

I want to finish the project in order to write the abstract, I want to be able to summarize everything concisely.

## Literature Review/Proposal

### What are Low Complexity Regions?

For decades, it was believed that peptide sequences which lack the ability to form stable three-dimensional structures also lack specific biological function (Haerty and Golding 2010b). Interestingly, among eukaryotic proteomes, the most commonly shared peptide sequences are found to be sequences with a low information content which lack a stable three-dimensional structure (Haerty and Golding 2010b; Marcotte et al. 1999; Bannen et al. 2007). These sequences have been termed ‘low-complexity regions’ (LCRs) due to their low information content and entropy, as well as their lack of diversity in amino acid composition (Wootton and Federhen 1993; Coletta et al. 2010). LCRs are found in DNA as well as protein sequences and can present in a variety of ways, all of which skew the composition of amino acids in a different manner (Wootton and Federhen 1993; Mier et al. 2020). Homorepeats, direpeats, tandem repeats, and imperfect repeats are common definitions of LCRs based on the periodicity of amino acids in a given sequence, but not every LCR is defined by a specific pattern (Mier et al. 2020). Most of the time, these patterns are found to occur in non-coding regions and evolve with minimal selective pressure (Kruglyak et al. 2000). Further research is being done in order to uncover the function of LCRs in protein coding regions, as well as the evolutionary background of these repetitive regions (Huntley and Golding 2006). To investigate the process of LCR evolution, this study proposes an Approximate Bayesian Computation (ABC) approach which will enable the prediction of evolutionary parameters such as mutation rates and insertion/deletion rates.

### Characteristics and Types of LCRs

Algorithms to detect the presence of low complexity regions in a sequence are available, and continue to be improved with further research into LCRs. Wootton and Federhen (1993) first developed an algorithm called SEG to find low complexity regions in protein sequences using information content (Huntley and Golding 2002). Information content is a common characteristic used to identify low complexity regions and in order to calculate the amount of information within a segment, the SEG algorithm implements Shannon’s entropy (Battistuzzi et al. 2016; Wootton and Federhen 1993). Shannon’s entropy (Shannon 1948) has been commonly used as a measure of complexity of a string of characters (Battistuzzi et al. 2016; Coletta et al. 2010; Wootton and Federhen 1993). This study will use the SEG algorithm and therefore Shannon’s entropy to assess the complexity of protein sequences. The less complex a sequence is (low variety of residues), the lower the entropy/information content of the sequence. Although LCRs are defined by their low information content, these regions have also been found to be hypermutable, and it is thought that throughout evolutionary history, they frequently gained and lost repeats (Marcotte et al. 1999; Kruglyak et al. 1998). In studying the *Drosophila melanogaster* gene mastermind, which encodes a highly repetitive nuclear protein, Newfeld et al. (1991) identified different patterns of evolutionary change between regions of high and low complexity. Repetitive regions were found to have a much higher rate of amino acid replacement, therefore the rate of evolution within these regions is higher than outside (Newfeld et al. 1991; Huntley and Golding 2000).

LCRs all share an overall low diversity of residues but present in unique ways as periodic or aperiodic repeats, which take on the form of homopolymers and heteropolymers (Wootton and Federhen 1993; Battistuzzi et al. 2016). Homopolymers/homorepeats are consecutive iterations of a single amino acid residue, and heteropolymers (direpeats, tandem repeats) are consecutive iterations of more than one residue that can be found in a variety of different patterns based on periodicity (Battistuzzi et al. 2016; Mier et al. 2020). Microsatellites, one of the best studied types of LCRs, commonly describe regions composed of tandem repeats that are typically made from anything between one to six nucleotides (Ellegren 2004). Although microsatellites normally refer to DNA sequences, it has been found that LCRs in proteins are comparable to microsatellites (DePristo et al. 2006). The molecular mechanisms involved in the process of evolution including slippage and unequal recombination are important for microsatellites and therefore protein LCRs as well (DePristo et al. 2006). A class of proteins which are related to, but slightly differ from LCRs are intrinsically disordered proteins (IDPs). IDPs are unable to form stable three dimensional structures and are characterized by low sequence complexity, biased amino acid composition, and high proportions of charged and hydrophilic residues (Wright and Dyson 2015). IDPs are composed of intrinsically disordered regions which are not necessarily defined by a low information content as LCRs are (Haerty and Golding 2010b; Dunker et al. 2002) (Haerty and Golding 2010b; Dunker et al. 2002). In this study, we propose a focus on low complexity regions.

## Why care about LCRs?

Proteins continue to be a large area of research due to their involvement in vital cellular processes and many human diseases. In protein sequence databases such as Swiss-Prot, the increase in the number of sequences and organisms represented has subsequently led to a decrease in the proportion of proteins containing LCRs (Coletta et al. 2010). On top of this, there is a lack of representation of LCRs in the protein data bank (Huntley and Golding 2002). Despite this underrepresentation of LCRs, they are known to be associated with several human neurodegenerative diseases and are thought to serve important biological functions (Coletta et al. 2010; Huntley and Golding 2006). It has also been found that the proteins of *Plasmodium falciparum* (the human malaria parasite) contain a high incidence of LCRs which has further highlighted the importance of both the evolution and function of LCRs (Gardner et al. 2002; DePristo et al. 2006). LCRs can appear as trinucleotide repeats which form repeated units of three nucleotides and are a well known form of deleterious mutation in humans (Ross et al. 1993). These are found to be the cause of diseases including fragile X syndrome, myotonic dystrophy, spinal atrophy, muscular atrophy, and Huntington's disease (Ross et al. 1993). These diseases can be broadly classified into two distinct groups, translated polyglutamine triplet repeat diseases and untranslated triplet repeat diseases (Everett and Wood 2004). Polyglutamine triplet repeat diseases, such as Huntington's disease, result in the formation of protein aggregates in the cell and occur due to expanded repeats being translated into expanded polyglutamine residues (Everett and Wood 2004). Untranslated triplet repeat diseases such as myotonic dystrophy and fragile X syndrome differ from polyglutamine repeat diseases as they contain trinucleotide repeats which are not translated into expansion within a mutant protein (Everett and Wood 2004).

The persistence of LCRs within genomes provides good evidence of their beneficial functions Verstrepen et al. 2005.

LCRs have been associated with important biological processes such as genetic recombination, antigen diversification, and protein-protein interactions (Karlin et al. 2002; Verstrepen et al. 2005; Kumari et al. 2015). The repetitive regions are thought to drive recombination events which alter genes and result in phenotypic variation (Verstrepen et al. 2005). In the genomes of *Haemophilus influenzae* and *Neisseria meningitidis*, LCRs are abundant and cause phase variation which gives the bacteria the ability to change their adherence patterns to host cells (Bayliss et al. 2001). This ultimately increases the fitness of the population and allows the bacteria to evade the host response (Bayliss et al. 2001). It was previously believed, based on structural evidence, that these hypermutable LCRs did not form stable structures but instead existed as solvent-exposed disordered coils (Wootton and Federhen 1993; Huntley and Golding 2002; DePristo et al. 2006). Using proteins from a non-redundant Protein Data Bank (PDB) dataset, Kumari et al. (2015) analyzed secondary structure content and surface accessibility and discovered that LCRs can form secondary structures within proteins. More specifically, in a large majority of identified LCRs, the analysis revealed the presence of more than one secondary structure, indicating that LCRs are found in regions where structure transition occurs (Kumari et al. 2015). Although more work is necessary to further understand the functions of LCRs, their role in genetic recombination, protein structure and function, and antigen diversity, highlight the importance of LCR research.

## How do LCRs Evolve?

Although research surrounding the evolution of LCRs is lacking, there are two proposed mechanisms of microsatellite evolution, which can be applied to many forms of LCRs. The first, polymerase slippage or slipped strand mispairing, involves loops being formed in either the coding or template strand, which causes a misalignment of strands and results in either the insertion or deletion of repetitive motifs (Levinson and Gutman 1987; Ellegren 2004). It is believed that slipped strand mispairing is the predominant mode of mutation of LCRs, specifically in homopolymer sequences (Levinson and Gutman 1987). The second mechanism, unequal recombination, occurs when repetitive regions in homologous chromosomes do not align properly during meiosis, which results in the repetitive region being expanded in one chromosome and contracted in the other (Warren et al. 1997; Mirkin 2007). In order to gain more insight into the evolutionary background of a variety of organisms, researchers have created models of events such as slippage in order to estimate mutation rate and other evolutionary parameters (Kruglyak et al. 2000). In a study of 10,844 parent/child allele transfers at nine short tandem repeat loci, Brinkmann et al. (1998) discovered 23 mutations, all of which were either gains or losses of repeats. Of the 23 mutations, 22 were due to single repeat mutations, which is why it has been common to use the stepwise mutation model of Ohta and Kimura (1973), that assumes repetitive regions expand or contract by 1 unit at a specific mutation rate (Kruglyak et al. 2000; Brinkmann et al. 1998). There are however, major drawbacks of the stepwise mutation model including that lengths can become negative, and the collection of repeat lengths in a sample will not have a stationary distribution (Kruglyak et al. 2000). It was thought that more complex models of LCR evolution were necessary to gain more accurate results, thus Kruglyak et al. (1998) proposed a model that incorporated length dependent slippage events. This differed from the stepwise mutation model in that the balance between slippage events and point mutations produced an equilibrium distribution of repeats (Kruglyak et al. 1998).

Models of LCR formation including replication slippage support the historical belief that LCRs evolve neutrally. More recently, there has been increasing evidence suggesting that LCRs are also acted upon by selective pressure (Haerty and Golding 2010b). Kimura (1983) proposed the neutral theory of molecular evolution which suggests that selection does not play a role in the genetic diversity within and between species, rather genetic diversity is neutral (Nevo 2001). Evidence for the neutral evolution of LCRs relies on a large number of factors including both their lack of stable structure and function (Dunker et al. 2002; Haerty and Golding 2010b), and ability to frequently gain or lose repeats through replication slippage (Marcotte et al. 1999; Kruglyak et al. 1998; Huntley and Golding 2000). Support for a selective model of LCR evolution comes from the non-random patterns of changes within LCRs, the deleterious effect of their expansion in humans (Karlin et al. 2002), and their enrichment in proteins involved in transcription, DNA, protein binding, reproduction, and development (Huntley and Clark 2007; Haerty and Golding 2010a; Battistuzzi et al. 2016). In a study of orthologous mouse and human genes, Mularoni et al. (2007) found a significant negative correlation between repeat number and gene nonsynonymous substitution rate, indicating that proteins acted upon by strong selective pressure contain a large number of repeats conserved between the two species (Mularoni et al. 2007). Interestingly, the study also reported a significant positive correlation between repeat size difference and protein nonsynonymous substitution rate, demonstrating that events such as slippage and substitutions occur in proteins which undergo neutral evolution (Mularoni et al. 2007). It was later revealed in a study by Battistuzzi et al. (2016) in which 11 representative Apicomplexa genomes were analyzed, that neutral mechanisms were found to act on highly repetitive LCRs (homopolymers) whereas selective pressures were influenced by the heterogeneity and length of the LCR (Battistuzzi et al. 2016). This work only begins to unravel the complexities of the evolutionary patterns associated with LCRs.

## What is an Approximate Bayesian Computation Markov chain Monte Carlo algorithm?

When studying molecular evolution, a common practice is to use model-based analyses of sets of DNA and amino acid sequences (Laurin-Lemay et al. 2022). This approach allows for the estimation of evolutionary genetic parameters such as mutation rates and insertion/deletion rates (Wu and Rodrigo 2015). Model-based statistical inference generally revolves around calculating the likelihood function, which represents the probability of the observed data under a chosen model (Sunnåker et al. 2013). The likelihood function therefore quantifies how well the data supports both the parameter values as well as the model (Sunnåker et al. 2013). However, due to an increase in the complexity and magnitude of available data, many current model-based analyses have become intractable by virtue of the likelihood function being difficult to calculate (Marjoram 2013). Approximate Bayesian computation (ABC) methods are rooted in Bayesian statistics and have been gaining popularity in areas such as genetics, as they bypass the calculation of the likelihood function (Sunnåker et al. 2013). The way in which they do this is by utilizing a simulation step in place of the calculation as a way to provide an estimate of the likelihood function (Marjoram 2013). Since there are many ways to approach a simulation, there are many different forms of ABCs. The more popular forms include ABC rejection methods, ABC Markov chain Monte Carlo methods (ABC-MCMC), and Sequential Monte Carlo ABC methods (ABC-SMC) (Marjoram 2013). This study proposes the use of an ABC-MCMC algorithm in order to estimate evolutionary parameters such as mutation and indel rates, and provide insight into the formation and evolution of protein LCRs.

The reason for proposing an ABC-MCMC in this study stems from the lack of a pre-existing model which explains how insertions and deletions work. Insertions and deletions alter the landscape of a sequence, making the likelihood calculation extremely challenging. Marjoram et al. (2003) originally proposed the algorithm for a MCMC method without the use of likelihoods. The algorithm first starts from a selected parameter value and proposes a move to a new parameter value based on a proposal distribution (Marjoram et al. 2003). Using this new parameter value, a dataset is then simulated and summary statistics are calculated, which makes it possible to quantitatively compare differences between the simulated dataset and the observed dataset (Marjoram et al. 2003). If the difference between summary statistics is small, the Hastings Ratio is calculated and the proposed parameter value can be accepted with a certain probability, then a new value is proposed and the process begins again (Marjoram et al. 2003). On the other hand if the difference in summary statistics between the observed and simulated data is very large, we propose a new parameter value and begin the algorithm again (Marjoram et al. 2003; Marjoram 2013). The use of this algorithm has enabled the analysis of complex problems which tend to arise in the areas of population genetics, ecology, epidemiology, and systems biology (Sunnåker et al. 2013). The group of Liepe et al. (2010) have been leaders in the use of ABCs for inference of genetic networks. This is evident through the creation of a software package they created called ABC SysBio which can implement ABC algorithms in a straightforward manner (Marjoram 2013; Liepe et al. 2010). Prior to the year 2000, there were essentially no papers published on ABCs (Marjoram 2013). As we enter into an era where larger and more complex data can be collected, the need for improved models is necessary, hence the large increase over the last decade in papers which mention ABC methods (Marjoram 2013).

## How will we use an ABC-MCMC - Oct 28 First draft

Using the algorithm mentioned above for an ABC-MCMC, this study aims to better understand the evolutionary background/formation of protein LCRs. An ABC-MCMC will enable the prediction of two important evolutionary parameters, mutation rate and indel rate. There is possibility for the estimation of other parameters which will be explored upon investigating the first two. We will utilize amino acid sequences in this study, one being the SRP40 protein found in *Saccharomyces cerevisiae*, which is extremely biased in composition. This protein sequence will act as our observed data and we will compare this observed data to our simulated data.

In terms of a simulation, we will use C++ to first generate a random amino acid sequence of a certain length. This randomly generated sequence will then be mutated over a number of generations in a two-step process. The first process is to choose a random poisson deviate with a mean that is equal to the mutation rate multiplied by the total number of sites. A poisson distribution is used here because mutation is a rare event and rare events can be modelled using this distribution. The value of the poisson deviate yields the total number of sites in the simulated sequence which should be mutated at random. The second mutation process deals with amino acid expansion and in this case we iterate through each residue in the simulated protein sequence and scan for repeats. If a residue is part of a repeat, we take the total length of the repeat, multiply it by the mutation rate and use this value as the mean of a random exponential deviate. We use the exponential distribution as it models



waiting times between events. Based on the random exponential deviates assigned, we select the lowest value which represents the residue that will change fastest, and we alter that residue to either delete or insert a repeat at that position.

Once we simulate a protein sequence for a number of generations under certain parameter values, we need to obtain a set of summary statistics and compare the summary statistics of the observed and simulated data. We have proposed summary statistics based off notable characteristics such as protein length, number of LCRs, and the average entropy of the LCRs. There is the possibility for additional summary statistic characteristics upon exploration of the initially proposed characteristics. To quantitatively compare the differences between the observed and simulated data, we propose using a distance measure between the two vectors of summary statistics. This distance is just the norm of the vector observed-simulated. Along with this, we also propose the use of a threshold as a way to assess how close the two datasets are. If the distance between the two vectors of summary statistics is larger than this threshold, we can not accept the proposed parameter value and the algorithm begins again. On the other hand, if the distance is very small, we can move forward in the algorithm and potentially accept the new parameter value.

We intend to run the simulation under the same parameters many times and take the average of the produced summary statistic vectors before calculating the distance between observed and simulated data. It is also worth noting that each time we begin the algorithm again, new parameter proposals will be selected using random normal deviates. We hope to see the distance between summary statistics being minimized upon every iteration of the algorithm as this means the simulated protein closely resembles the observed protein under specific parameters.

## Materials and Methods (mid-year stuff Jan 20)

Custom scripts and commands utilized in this analysis can be found on GitHub at <https://github.com/opticrom/abcmcmc-thesis4c12>.

### ABC-MCMC: The Algorithm

This study utilized the ABC Markov chain Monte Carlo algorithm, originally proposed by Marjoram et al. (2003). The algorithm begins from a randomly selected parameter value and follows the steps below.

1. If now at  $\theta$ , propose a move to  $\theta'$  according to a proposal distribution  $q(\theta, \theta')$ .
2. Simulate a dataset,  $D'$  using  $\theta'$ .
3. If  $D' \sim D$  proceed to step 4; else, output  $\theta$  and return to 1.
4. Calculate the Hastings Ratio.

5. Accept, and output, the new  $\theta'$  with probability  $h$ . Else return to, and output,  $\theta$ , Go back to 1.

A custom C++ script was written to iterate through the algorithm for a desired number of simulations. The Normal distribution was used to control how new parameter values were proposed **SHOULD I MENTION THE MEAN AND STDEV OF THE NORMAL HERE? WHAT IF I WANT TO TEST MANY DIFFERENT KINDS**. We simulated a dataset under the newly proposed parameter value, which consisted of a randomly generated protein sequence. The simulated protein was compared to a protein of known low complexity called SRP40, which is found in the model organism *Saccharomyces cerevisiae*. The protein sequence for SRP40 was downloaded from the NCBI database. To quantitatively compare similarities between simulated and observed protein sequences, vectors of summary statistics (characteristics that describe the sequences) were created for each sequence. The Euclidean distance between the two vectors was calculated and compared to a threshold value to determine if the newly proposed parameter value could be accepted.

**NEED HELP:** In order to create a posterior distribution, parameter values need to be accepted with a certain probability. In the algorithm proposed by Marjoram et al. (2003), the Hastings Ratio is calculated to determine this acceptance probability. Currently having issues implementing the Hastings algorithm and struggling to understand why I can not use it or figure it out. What are other alternatives in this case for calculating acceptance probabilities? Also having issues with the threshold value, how do I know what an appropriate value to set it at is? why are my distances not decreasing with every iteration of the simulation? why am I either accepting too many or not enough, this is something we NEED to talk about.

## Parameters and Summary Statistics

Two parameters were estimated in this study, mutation rate and insertion/deletion (indel) rate. Mutation rate referred to the rate at which a single amino acid in a protein sequence changed into a different amino acid. Indel rate referred to the rate at which an amino acid was deleted or inserted from a protein sequence. For the purpose of this study, to determine if the length of a repetitive region played a role in insertions, any amino acid that was inserted into the simulated protein sequence was the same unit as the previous amino acid in the sequence **(Should I explain why here, not too sure, going to talk more about it in the next part).**

Summary statistics were utilized to capture important information about simulated and observed protein sequences in order to assess how similar the sequences were. Three summary statistics were used which included, the length of the protein sequence, the number of LCRs in the sequence, and the average entropy of the LCRs. To identify LCRs and their corresponding entropies, the *Seg* algorithm was implemented (Wootton and Federhen 1993). The following *Seg* parameters were utilized to search for LCRs in proteins; a window length of 15, a trigger segment complexity of 1.9, and an extension segment complexity of 2.2. We selected these due to previous research which demonstrated that these parameter values would better detect regions of low complexity in eukaryotes which are typically longer and contain more repetitive repeats (Huntley and Golding 2000;

Huntley and Golding 2002). In the 2002 paper, the 2000 paper is cited for this fact so I was not sure which to cite, put both for now.

Summary statistics were stored in vectors and normalized in order to prevent large values (like the length of the sequence) from dominating when calculating the Euclidean distance.

**Need Help:** We currently do not know how well these summary statistics explain the data. Weighting summary statistics could be beneficial, but determining how we weight each one appropriately is challenging at the moment. There are potential weighting techniques, such as one method used by Hamilton et al. (2005) in which summary statistics are regressed on retained parameters, allowing for a weight to be calculated for each summary statistic and subsequently used when calculating the Euclidean Distance. This method is very difficult to implement and I am struggling. Maybe other alternatives for this weighting?

### Simulation Step: Creation and Mutation of Protein Sequences

To bypass calculation of the likelihood function, a custom C++ script was written to simulate the generation and mutation of protein sequences over numerous generations. A random protein of desired length was generated first. We generated a protein sequence 400 amino acids in length, similar to the length of the SRP40 protein in *S. cerevisiae*. This simulated protein sequence was then mutated in the following ways.

We iterated over the simulated protein sequence and assigned exponential deviates to each amino acid. We utilized the exponential distribution because it is commonly used to model waiting times between events. Mutation and indel rates served to act as the scale parameter ( $\beta$ ), or mean of the distribution, indicative of the mean time until mutation occurs. In the case of mutation rate, the same rate was utilized across all sites, with the assumption made that repeats do not play a role in point mutation. In the case of the indel rate, we scanned for repeats, and if found, we multiplied the length of the repetitive segment by the indel rate, and utilized this new value as the scale parameter ( $\beta$ ) for the exponential distribution. **THINK ABOUT WHAT THIS MEANS IN TERMS OF THE EXPONENTIAL DISTRIBUTION, THE LONGER THE LENGTH OF REPEAT, HOW DOES THIS AFFECT THE EXPONENTIAL DISTRIBUTION, DOES IT MAKE SENSE?**

Exponential deviates were stored in two vectors, one for deviates generated using mutation rate, and the other for deviates generated using the indel rate. We then identified a single deviate with the lowest value (based on both vectors), which represented the residue that mutated quickest. Depending on which vector the lowest deviate came from, we either altered the corresponding amino acid or inserted/deleted a repetitive amino acid.

We ran the simulation ten times for each newly proposed parameter value and took the average of all ten vectors of summary statistics prior to calculating the distance. On each simulation, we mutated the protein sequence for 2000 generations before obtaining summary statistics **If I want to test with a different number of iterations should I highlight that in here to? Or should I just mention alternate simulation parameters were used.**

**Need Help:** After running the simulations, I want to visualize the data. I want to see the posterior distribution but this is hard without acceptance probabilities. Currently, I utilized gnuplot to try and visualize the change in parameter values over the number of simulations and I am not sure how accurate it is. This goes back to the issue of the threshold value and acceptance probability.

## Normalization and Euclidean Distance Calculation

To determine the similarity between simulated and observed protein sequences, the distance between the vectors of summary statistics was calculated by employing a Euclidean distance measure (2). Before calculating this distance, vectors of summary statistics were normalized (1) to prevent large values (such as length) from dominating the distance measure.

$$length = \sqrt{(x * x) + (y * y) + (z * z)} \quad (1)$$

$$normalized = [x/length, y/length, z/length]$$

$$d(P1, P2) = \sqrt{(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2} \quad (2)$$

## Results

## Discussion

## References

- Bannen R M, Bingman C A, and Phillips G N (2007). Effect of low-complexity regions on protein structure determination. *Journal of Structural and Functional Genomics* 8(4), 217–226.
- Battistuzzi F U, Schneider K A, Spencer M K, Fisher D, Chaudhry S, and Escalante A A (2016). Profiles of low complexity regions in Apicomplexa. *BMC evolutionary biology* 16(1), 1–12.
- Bayliss C D, Field D, Moxon E R, et al. (2001). The simple sequence contingency loci of *Haemophilus influenzae* and *Neisseria meningitidis*. *The Journal of clinical investigation* 107(6), 657–666.
- Brinkmann B, Klintschar M, Neuhuber F, Hühne J, and Rolf B (1998). Mutation rate in human microsatellites: influence of the structure and length of the tandem repeat. *The American Journal of Human Genetics* 62(6), 1408–1415.
- Coletta A, Pinney J W, Solís D Y W, Marsh J, Pettifer S R, and Attwood T K (2010). Low-complexity regions within protein sequences have position-dependent roles. *BMC systems biology* 4(1), 1–13.
- DePristo M A, Zilversmit M M, and Hartl D L (2006). On the abundance, amino acid composition, and evolutionary dynamics of low-complexity regions in proteins. *Gene* 378, 19–30.
- Dunker A K, Brown C J, Lawson J D, Iakoucheva L M, and Obradović Z (2002). Intrinsic disorder and protein function. *Biochemistry* 41(21), 6573–6582.
- Ellegren H (2004). Microsatellites: simple sequences with complex evolution. *Nature reviews genetics* 5(6), 435–445.
- Everett C and Wood N (2004). Trinucleotide repeats and neurodegenerative disease. *Brain* 127(11), 2385–2405.
- Gardner M J, Hall N, Fung E, White O, Berriman M, Hyman R W, Carlton J M, Pain A, Nelson K E, Bowman S, et al. (2002). Genome sequence of the human malaria parasite *Plasmodium falciparum*. *Nature* 419(6906), 498–511.
- Haerty W and Golding G B (2010a). Genome-wide evidence for selection acting on single amino acid repeats. *Genome research* 20(6), 755–760.
- Haerty W and Golding G B (2010b). Low-complexity sequences and single amino acid repeats: not just “junk” peptide sequences. *Genome* 53(10), 753–762.
- Hamilton G, Currat M, Ray N, Heckel G, Beaumont M, and Excoffier L (2005). Bayesian estimation of recent migration rates after a spatial expansion. *Genetics* 170(1), 409–417.
- Huntley M and Golding G B (2000). Evolution of simple sequence in proteins. *Journal of molecular evolution* 51(2), 131–140.
- Huntley M A and Clark A G (2007). Evolutionary analysis of amino acid repeats across the genomes of 12 *Drosophila* species. *Molecular biology and evolution* 24(12), 2598–2609.
- Huntley M A and Golding G B (2002). Simple sequences are rare in the Protein Data Bank. *Proteins: Structure, Function, and Bioinformatics* 48(1), 134–140.
- Huntley M A and Golding G B (2006). Selection and slippage creating serine homopolymers. *Molecular biology and evolution* 23(11), 2017–2025.

- Karlin S, Brocchieri L, Bergman A, Mrázek J, and Gentles A J (2002). Amino acid runs in eukaryotic proteomes and disease associations. *Proceedings of the National Academy of Sciences* 99(1), 333–338.
- Kimura M (1983). *The neutral theory of molecular evolution*. Cambridge University Press.
- Kruglyak S, Durrett R, Schug M D, and Aquadro C F (2000). Distribution and abundance of microsatellites in the yeast genome can be explained by a balance between slippage events and point mutations. *Molecular Biology and Evolution* 17(8), 1210–1219.
- Kruglyak S, Durrett R T, Schug M D, and Aquadro C F (1998). Equilibrium distributions of microsatellite repeat length resulting from a balance between slippage events and point mutations. *Proceedings of the National Academy of Sciences* 95(18), 10774–10778.
- Kumari B, Kumar R, and Kumar M (2015). Low complexity and disordered regions of proteins have different structural and amino acid preferences. *Molecular BioSystems* 11(2), 585–594.
- Laurin-Lemay S, Dickson K, and Rodrigue N (2022). Jump-Chain Simulation of Markov Substitution Processes Over Phylogenies. *Journal of Molecular Evolution*, 1–5.
- Levinson G and Gutman G A (1987). Slipped-strand mispairing: a major mechanism for DNA sequence evolution. *Molecular biology and evolution* 4(3), 203–221.
- Liepe J, Barnes C, Cule E, Erguler K, Kirk P, Toni T, and Stumpf M P (2010). ABC-SysBio—approximate Bayesian computation in Python with GPU support. *Bioinformatics* 26(14), 1797–1799.
- Marcotte E M, Pellegrini M, Yeates T O, and Eisenberg D (1999). A census of protein repeats. *Journal of molecular biology* 293(1), 151–160.
- Marjoram P (2013). Approximation bayesian computation. *OA genetics* 1(3), 853.
- Marjoram P, Molitor J, Plagnol V, and Tavaré S (2003). Markov chain Monte Carlo without likelihoods. *Proceedings of the National Academy of Sciences* 100(26), 15324–15328.
- Mier P, Paladin L, Tamana S, Petrosian S, Hajdu-Soltész B, Urbanek A, Gruca A, Plewczynski D, Grynberg M, Bernadó P, et al. (2020). Disentangling the complexity of low complexity proteins. *Briefings in Bioinformatics* 21(2), 458–472.
- Mirkin S M (2007). Expandable DNA repeats and human disease. *Nature* 447(7147), 932–940.
- Mularoni L, Veitia R A, and Albà M M (2007). Highly constrained proteins contain an unexpectedly large number of amino acid tandem repeats. *Genomics* 89(3), 316–325.
- Nevo E (2001). Genetic diversity.
- Newfeld S J, Smoller D A, and Yedvobnick B (1991). Interspecific comparison of the unusually repetitive *Drosophila* locus-mastermind. *Journal of molecular evolution* 32(5), 415–420.
- Newfeld S J, Tachida H, and Yedvobnick B (1994). Drive-selection equilibrium: homopolymer evolution in the *Drosophila* gene mastermind. *Journal of molecular evolution* 38(6), 637–641.
- Ohta T and Kimura M (1973). A model of mutation appropriate to estimate the number of electrophoretically detectable alleles in a finite population. *Genetics Research* 22(2), 201–204.

- Ross C A, McInnis M G, Margolis R L, and Li S.-H (1993). Genes with triplet repeats: candidate mediators of neuropsychiatric disorders. *Trends in neurosciences* 16(7), 254–260.
- Shannon C E (1948). A mathematical theory of communication. *The Bell system technical journal* 27(3), 379–423.
- Sunnåker M, Busetto A G, Numminen E, Corander J, Foll M, and Dessimoz C (2013). Approximate bayesian computation. *PLoS computational biology* 9(1), e1002803.
- Verstrepen K J, Jansen A, Lewitter F, and Fink G R (2005). Intragenic tandem repeats generate functional variability. *Nature genetics* 37(9), 986–990.
- Warren S T, Muragaki Y, Mundlos S, Upton J, and Olsen B R (1997). Polyalanine expansion in synpolydactyly might result from unequal crossing-over of HOXD13. *Science* 275(5298), 408–409.
- Wootton J C and Federhen S (1993). Statistics of local complexity in amino acid sequences and sequence databases. *Computers & chemistry* 17(2), 149–163.
- Wright P E and Dyson H J (2015). Intrinsically disordered proteins in cellular signalling and regulation. *Nature reviews Molecular cell biology* 16(1), 18–29.
- Wu S H and Rodrigo A G (2015). Estimation of evolutionary parameters using short, random and partial sequences from mixed samples of anonymous individuals. *BMC bioinformatics* 16(1), 1–12.

## Appendix C++ simulation (so far)

```

1 #include "functions.cpp"
2 #include "getindex.cpp"
3 #include <bits/stdc++.h>
4 #include <iostream>
5 #include <vector>
6 #include <string>
7 #include <ctime>
8 #include <algorithm> //this is to get min element stuff, cool library
9 #define numAA 20
10 using namespace std;
11
12
13 Ran myran(time(NULL)); //We will use 21 as the random seed right now, used in Poissondev too
14
15 ///////////////////////////////////////////////////
16 // FIRST FUNCTION TAKES AN INTEGER VALUE AND GENERATES A RANDOM/
17 // AMINO ACID SEQUENCE OF THAT LENGTH /
18 ///////////////////////////////////////////////////
19
20 std::string createSeq(int n){
21
22     char aminoAcids[numAA] = { 'G', 'A', 'L', 'M', 'F', 'W', 'K', 'Q', 'E', 'S', 'P', 'V',
23                                'I', 'C', 'Y', 'H', 'R', 'N', 'D', 'T' };
24
25     std::string protein = "";
26     for (int i = 0; i < n; i++){
27         protein += aminoAcids[myran.int64() % numAA];} //this rand() % 20 means in the range 0-19
28
29     //std::cout << protein << "\n" << "\n" ;
30     return protein;
31 }
32
33
34 /*THIS FUNCTION WILL MUTATE THE SIMULATED PROTEIN SEQUENCE
35 BY CHOOSING A RANDOM EXPONENTIAL DEVIATE (WITH MEAN = MUTATION RATE)
36 FOR EACH AMINO ACID IN THE SEQUENCE AND SUBSEQUENTLY SELECTING
37 THE AMINO ACID WITH THE LOWEST NUMBER (QUICKEST TO MUTATE) AND
38 MUTATING IT RANDOMLY, THIS IS DONE SUCCESSIVELY TO PRODUCE A
39 PROTEIN AND CREATE A VECTOR OF VALUES SIMILAR TO ABOVE (SEP 21)
40 SEP 26 - PROCESS CHANGE, THIS FUNCTION WILL BE USED FOR AMINO ACID EXPANSION
41 UPDATE NOV 15 - WE ARE MAKING THIS NEW FUNCTION THAT ASSIGNS DEVIATES

```



```

42 BASED ON BOTH MUTATION RATES AND INDEL RATES AND THEN OUT OF BOTH VECTORS
43 WE WILL FIND THE LOWEST AND EITHER MUTATE IT OR DO AN INS/DEL DEPENDING
44 ON WHICH VECTOR THE LOWEST DEVIATE CAME FROM*/
45
46 std::string mutateSeqExp(std::string simulated_protein){
47
48     // Setting up the vectors
49     std::vector<double> exp_deviates_vtr_ind; // Creating a vector to hold the values of the deviates
        for indel rate
50     std::vector<double> exp_deviates_vtr_mut; // Creating a vector to hold the values of the deviates for
        mutation rate
51     std::vector<double> smallest_vtr; // Creating a vector to store the smallest element of each of the 2
        vectors
52
53     //std::cout << "before mutateseqEXP:\t" << simulated_protein << "\n"; // Initially printing the non-
        mutated strin.
54
55     // Mutation and indel rate set here now
56     float mutation_rate = 0.14;
57     float indel_rate = 0.14;
58
59     // First loop will assign deviates based on mutation rates
60     for (int i = 0; i < simulated_protein.length(); i++) {
61
62         float betal = mutation_rate ; // 1 will always be used here because the length if no repeats is 1
63         Expondev myexp(betal,myran.int64());
64         double deviate = myexp.dev();//here we choose exp_deviates(mean of beta)
65         exp_deviates_vtr_mut.push_back(deviate) ; //Here we are storing the exponential deviates
66     }
67     // Traversing the string
68     for (int i = 0; i < simulated_protein.length(); i++) {
69
70         int counter = 1 ;
71
72         //Code to scan back and forth to find repeats
73         if (simulated_protein[i] != simulated_protein[i+1] && simulated_protein[i] != simulated_protein[i
        -1]) {
74             float beta2 = indel_rate ; // 1 will always be used here because the length if no repeats is
        1
75             Expondev myexp(beta2,myran.int64());
76             double deviate = myexp.dev();//here we choose exp_deviates(mean of beta)
77             exp_deviates_vtr_ind.push_back(deviate) ; //Here we are storing the exponential deviates
78         } else {

```

```

79     int x = 1 ;
80     int y = 1 ;
81
82     //Be careful in these while loops, for i-y, when i is 0
83     //and y is 1, how does it not throw error
84     //Looking forward for repeats
85     while (simulated_protein[i] == simulated_protein[i + x]) {
86         counter += 1 ;
87         x++;
88     }
89     //Looking backwards for repeats
90     while (simulated_protein[i] == simulated_protein[i - y]) {
91         counter += 1 ;
92         y++;
93     }
94
95     float beta3 = indel_rate * counter ;
96     Expondev myexp(beta3,myran.int64());
97     double deviate = myexp.dev();
98     exp_deviates_vtr_ind.push_back(deviate);
99 }
100 }
101
102 //selecting the lowest deviate from both vectors
103 double min_mut = *min_element(exp_deviates_vtr_mut.begin(), exp_deviates_vtr_mut.end());
104 double min_ind = *min_element(exp_deviates_vtr_ind.begin(), exp_deviates_vtr_ind.end());
105
106 // Append the 2 minimums to the new vector
107 smallest_vtr.push_back(min_mut);
108 smallest_vtr.push_back(min_ind);
109
110 // Get the smallest of the small numbers
111 double smallest_num = *min_element(smallest_vtr.begin(), smallest_vtr.end());
112
113 // Getting the index of the smallest of the small numbers
114 int position = getIndex(smallest_vtr, smallest_num);
115
116 // If index = 0, mutation If index = 1, indel
117 if (position == 0){
118
119     for (int i = 0; i < 1; i++) {
120
121         char aminoAcids[20] = { 'G', 'A', 'L', 'M', 'F', 'W', 'K', 'Q', 'E', 'S', 'P', 'V',

```

```

122         'I', 'C', 'Y', 'H', 'R', 'N', 'D', 'T' };
123
124         char random_AA = aminoAcids[myran.int64() % numAA]; // sets up the random amino acid, same
125         used in first function to createSeq
126         int position2 = getIndex(exp_deviates_vtr_mut, min_mut);
127         simulated_protein[position2] = random_AA; // indexes the simulated protein at a random spot
128         and replaces the existing AA with a new random one
129         //std::cout << position2 << "\n";
130     }
131
132     } else {
133
134         int position3 = getIndex(exp_deviates_vtr_ind, min_ind);
135         char aa_index = simulated_protein[position3];
136         float random_number = myran.doub();
137
138         //Inserting or deleting a repeat
139         if (random_number < 0.5){
140             simulated_protein.erase(position3, 1);
141         } else {
142             simulated_protein.insert(position3+1,1, aa_index);
143         }
144     }
145
146     //printing out this stuff to check its working
147     //std::cout << min_mut << "\n" << min_ind << "\n" << smallest_num << "\n" << position << "\n";
148
149     // THIS IS JUST TO PRINT THE VECTOR
150     //for (int x = 0; x < exp_deviates_vtr_ind.size(); x++) {
151     //    std::cout << exp_deviates_vtr_ind[x] << ' ';
152     //}
153
154     //std::cout << "\n" << "after mutateSeqEXP:\t" << simulated_protein << "\n" << "\n" ;
155     return simulated_protein;
156 }
157
158
159 double getNormalDev(double mu, double stdev) {
160     Normaldev mynorm(mu, stdev, myran.int64());
161     double dev = mynorm.dev();
162     //std::cout << dev << "\n";

```

```
163     return dev;
164 }
165
166 //int main() {
167 //    std::string x = "MKNHCHKISAKHHHHHAM";
168 //    mutateSeqExp(x);
169 //}
```