

# Estimating Evolutionary Parameters for Protein Low Complexity Regions using an Approximate Bayesian Computation

Alexander Turco

December 5, 2022

# Overview

Background Information

Research Questions/Explorations

Experimental Approach

Future Work

## What are Low Complexity Regions?

## Saccharomyces cerevisiae SRP40 Protein LCRs

>CAA82171.1(25-125) complexity=0.92 (15/1.90/2.20)

```

sssssssssssssssssssssssgsssssssssssssdssdssdsessssssssss
sssssdsssssedssssgsssssssssdesssesede

```

>CAA82171.1(149-282) complexity=1.33 (15/1.90/2.20)

```

essssessssgsssssesgsgesdsdsssssssssdsestdesdsqsssssssdsss
dsdssssdsdssdsdssssssssssdsdsdsdsssdsssgssdsssssdsdssdestssds
dsdsdsdsgssse

```

>CAA82171.1(298-316) complexity=2.18 (15/1.90/2.20)

tpassnestpsassssan

# LCRs Present in Unique Ways

## Homorepeats

Consecutive iterations of a single residue



# LCRs Present in Unique Ways

## Homorepeats

Consecutive iterations of a single residue



## Direpeats

Consecutive iterations of two ordered, different residues



# LCRs Present in Unique Ways

## Homorepeats

Consecutive iterations of a single residue



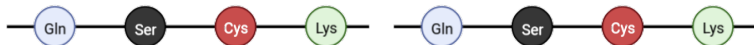
## Direpeats

Consecutive iterations of two ordered, different residues

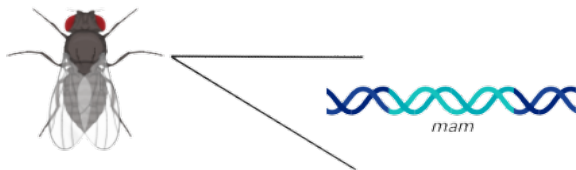


## Tandem Repeats

Sequence of residues which are repeated a number of times



# LCRs are Hypermutable

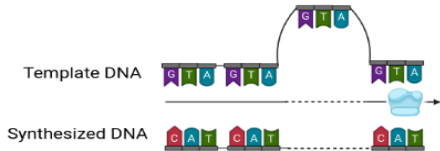


<b>mam domain</b>	<b>Size (bp)</b>	<b>Amino Acid Substitutions</b>	<b>Amino Acid/ Total Substitutions</b>
Unique	933	26	0.15
Repetitive	810	47	0.42

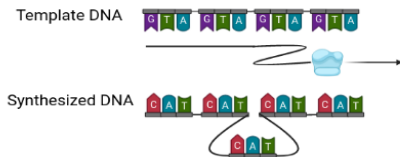
<sup>1</sup>Newfeld, Smoller, and Yedvobnick, 1991

# Proposed Mechanisms of LCR Evolution

## 1. *Polymerase Slippage/Slipped Strand Mismatching*



**Polymerase Slips Forward**



**Polymerase Slips Backwards**

<sup>2</sup>Levinson and Gutman, 1987; Sehn, 2015

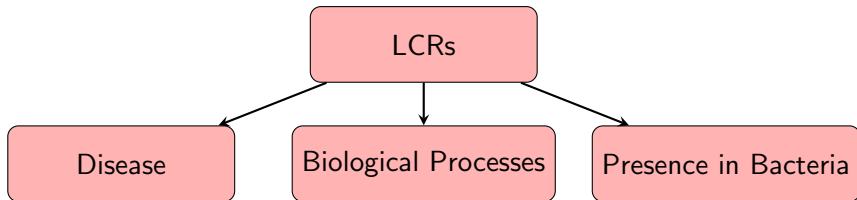


# Proposed Mechanisms of LCR Evolution

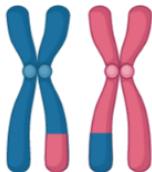
## 2. *Unequal Recombination*



# Why Care about LCRs and their Evolution?



*Huntington's Disease*



*Genetic Recombination*



*Neisseria meningitidis*

# What will this Study Explore?

- ▶ Estimation of evolutionary parameters (mutation rate, indel rates)
- ▶ Various models of insertions and deletions
- ▶ Summary statistics which best explain data

# What Approach will be Taken?

## Bayesian Statistics: Model-based statistical inference

$$p(D|\theta) \tag{1}$$

*Likelihood*

# What Approach will be Taken?

## Bayesian Statistics: Model-based statistical inference

$$p(D|\theta) \quad (1)$$

*Likelihood*

$$p(\theta|D) \quad (2)$$

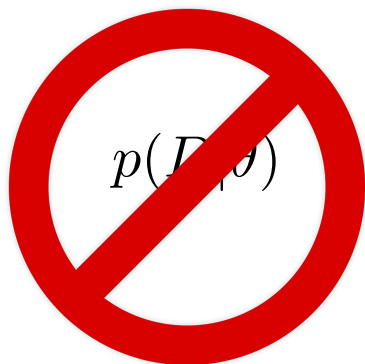
*Posterior*

# Why use an ABC-MCMC

- ▶ The increasing complexity and magnitude of available data can make the likelihood difficult to calculate

$$p(D|\theta)$$

# Why use an ABC-MCMC



- ▶ Calculation of the likelihood is replaced with a simulation step

# MCMC for ABC

- 1 Propose a move from  $\theta$  to  $\theta'$  according to a transition kernel  $q(\theta, \theta')$ .
- 2 Generate simulated dataset  $D'$  using  $\theta'$  and calculate  $S'$ .
- 3 If  $\rho(S', S) \leq \epsilon$  continue to 4, otherwise remain at  $\theta$  and go to 1.
- 4 Calculate

$$\alpha(\theta, \theta') = \min(1, \frac{\pi(\theta')q(\theta', \theta)}{\pi(\theta)q(\theta, \theta')})$$

- 5 Accept  $\theta'$  with probability  $\alpha$ , otherwise stay at  $\theta$ .
- 6 Return to 1.

---

<sup>4</sup> Marjoram et al., 2003



# MCMC for ABC: Modified Algorithm

- 1 Propose a move from  $\theta$  to  $\theta'$  according to the normal distribution

# MCMC for ABC: Modified Algorithm

- 1 Propose a move from  $\theta$  to  $\theta'$  according to the normal distribution
- 2 Create a random protein sequence and mutate over many generations using  $\theta'$  to generate simulated Dataset  $D'$

# MCMC for ABC: Modified Algorithm

- 1 Propose a move from  $\theta$  to  $\theta'$  according to the normal distribution
- 2 Create a random protein sequence and mutate over many generations using  $\theta'$  to generate simulated Dataset  $D'$
- 3 Calculate summary statistics for simulated dataset  $D'$

# MCMC for ABC: Modified Algorithm

- 1 Propose a move from  $\theta$  to  $\theta'$  according to the normal distribution
- 2 Create a random protein sequence and mutate over many generations using  $\theta'$  to generate simulated Dataset  $D'$
- 3 Calculate summary statistics for simulated dataset  $D'$
- 4 If  $d(S', S) \leq \epsilon$ , go to next step, otherwise stay at  $\theta$  and return to 1

# MCMC for ABC: Modified Algorithm

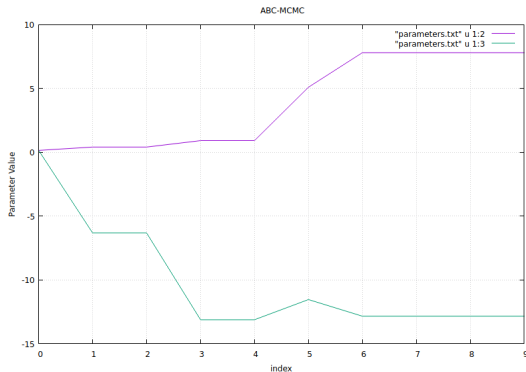
- 1 Propose a move from  $\theta$  to  $\theta'$  according to the normal distribution
- 2 Create a random protein sequence and mutate over many generations using  $\theta'$  to generate simulated Dataset  $D'$
- 3 Calculate summary statistics for simulated dataset  $D'$
- 4 If  $d(S', S) \leq \epsilon$ , go to next step, otherwise stay at  $\theta$  and return to 1
- 5 Accept  $\theta'$

# MCMC for ABC: Modified Algorithm

- 1 Propose a move from  $\theta$  to  $\theta'$  according to the normal distribution
- 2 Create a random protein sequence and mutate over many generations using  $\theta'$  to generate simulated Dataset  $D'$
- 3 Calculate summary statistics for simulated dataset  $D'$
- 4 If  $d(S', S) \leq \epsilon$ , go to next step, otherwise stay at  $\theta$  and return to 1
- 5 Accept  $\theta'$
- 6 Return to step 1

# Future Work

- ▶ Graphical representations of simulation iteration versus parameter values
- ▶ Implementation of weighted summary statistics in distance calculation
- ▶ Adjustment of values such as mean and standard deviation of the proposal distribution



# Acknowledgements

- ▶ Dr. Brian Golding
- ▶ Sam Long
- ▶ Zachery Dickson
- ▶ Johanna Enright