

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
sssssdssssesdssssgsssssssssdesssesede

```

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

```

essssessssgsssssesgsgesdsdsssssssssdsestdesdsqsssssssdsss
dsdssssdsdssdsdssssssssssdsdsdsdsssdssdgssdsssssdsdssdestssds
dsdssdsdsgssse

```

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

tpassnestpsasssssan

LCRs Present in Unique Ways

Homorepeats

Consecutive iterations of a single residue



Direpeats

Consecutive iterations of two ordered, different residues

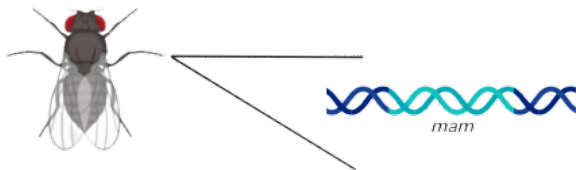


Imperfect Repeats

Regions in which the repeat units are not the same



LCRs are Hypermutable

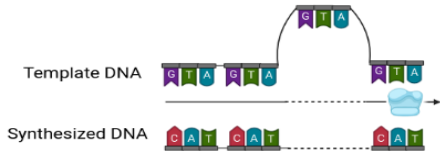


mam domain	Size (bp)	Amino Acid Substitutions	Amino Acid/ Total Substitutions
Unique	933	26	0.15
Repetitive	810	47	0.42

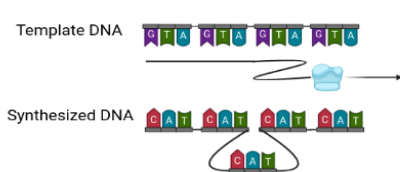
¹Newfeld, Smoller, and Yedvobnick, 1991

Proposed Mechanisms of LCR Evolution

1. *Polymerase Slippage/Slipped Strand Mismatching*



Polymerase Slips Forward



Polymerase Slips Backwards

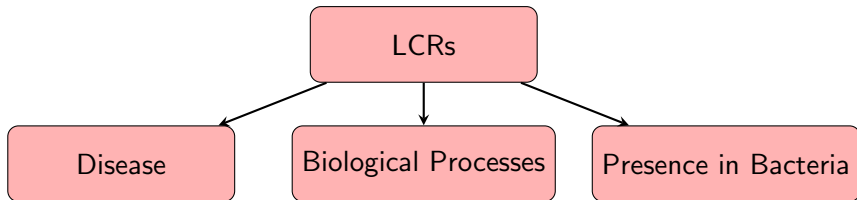
² Levinson and Gutman, 1987; Sehn, 2015

Proposed Mechanisms of LCR Evolution

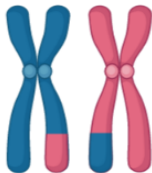
2. Unequal Recombination



Why Care about LCRs and their Evolution?



Huntington's Disease



Genetic Recombination



Neisseria meningitidis

What Did we Do in this Study?

- ▶ Programmed an ABC-MCMC using C++ which consisted of a simulation step where amino acid sequences were altered by point mutations and insertions/deletions
- ▶ Utilized the exponential distribution with ($\beta = \text{length} * \text{indel_rate}$) to see if the length of a repeat plays a role in its mutation
- ▶ Attempted to estimate evolutionary parameters such as mutation rates and insertion/deletion rates
- ▶ Explored summary statistics that could quantitatively explain characteristics of LCRs

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

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

Likelihood

$$p(\theta|D) \tag{2}$$

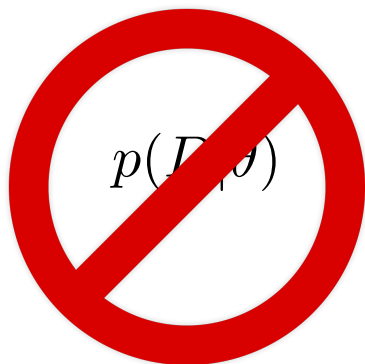
Posterior

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

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

- 5 Accept θ' with probability α , otherwise stay at θ .
- 6 Return to 1.

⁴ Marjoram et al., 2003

MCMC for ABC: Modified Algorithm

- 1 Propose a move from θ to θ' according to the normal distribution

MCMC for ABC: Modified Algorithm

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

MCMC for ABC: Modified Algorithm

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

MCMC for ABC: Modified Algorithm

- 1 Propose a move from θ to θ' according to the normal distribution
- 2 Create a random protein sequence and mutate over many generations using θ' 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 θ and return to 1

MCMC for ABC: Modified Algorithm

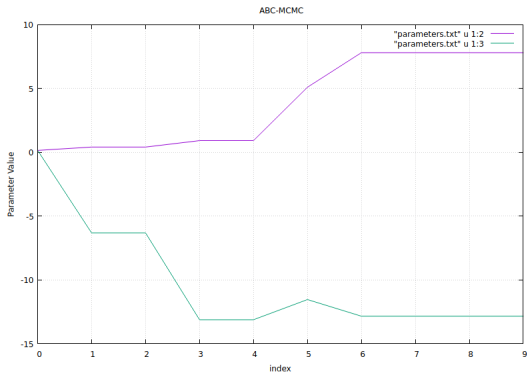
- 1 Propose a move from θ to θ' according to the normal distribution
- 2 Create a random protein sequence and mutate over many generations using θ' to generate simulated Dataset D'
- 3 Calculate summary statistics for simulated dataset D'
- 4 If $d(S', S) \leq \epsilon$, go to next step, otherwise stay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

MCMC for ABC: Modified Algorithm

- 1 Propose a move from θ to θ' according to the normal distribution
- 2 Create a random protein sequence and mutate over many generations using θ' 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 θ and return to 1
- 5 Accept θ'
- 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