
Survival of the complex/Survival of the most flexible

Tom Röschinger¹, Simone Pompei², Michael Lässig^{3,*},

¹ Division of Chemistry and Chemical Engineering, California Institute of Technology, Pasadena, CA 91125

²

³ Institut für Biologische Physik, Universität zu Köln, Zùlpicherstr. 77, 50937, Köln, Germany

* mlaessig@uni-koeln.de

To Do List

- Supplementary Methods 3: Does the diversity within a population depend on the alphabet size, as seen from the analytic results? $\Delta E \sim \left(\frac{n}{n-1}\right)^2$. This is important for the results in this section.

-

Abstract

Following text is simply copied from thesis and not changed yet. Molecular recognition has a crucial role on multiple biological systems on various scales, from cell to cell recognition and the understanding of pathogen immunity to transcription. The stochastic nature of these processes are ideal circumstances for statistical models. The model I present consists of two components, the driver and the trailer, and the evolutionary dynamics depend on their mutual specificity. Quantitative evolutionary models for two component systems usually consider the evolution of one component only, the trailer, based on a fixed driver. The novelty of this approach is to explicitly include mutations of the driver that are caused from external events, and force the two component system out of the most beneficial interaction. In this system, complexity is introduced as a dynamic variable, and it is shown how neutral driver mutations and dynamic complexity can be included in models for quantitative trait evolution. As a result of coevolution, the complexity of the driver and trailer components increases.

Introduction

The complexity of a system is a crucial component in molecular evolution. Evolution of gene regulation is often assumed to be driven by the gain and loss of transcription factor binding sites (TFBS), rather than the evolution of a transcription factor (TF) itself [1]. This assumption is based on the pleiotropy of most TFs that have been observed, and that the pleiotropic effects of a mutation in such a transcription factor are deleterious for most interactions, such that no change is possible. However, these TF's are conserved across species to begin with [2]. Some classes of TFs display extensive levels of species diversity while maintaining structure and function [3] and it has been discussed that protein evolution is a crucial source of developmental variation [4], but the consequences of protein evolution remain unclear [5, 6]. The evolution of gene regulatory networks is assumed to be driven by modifications and the evolution of TFs [7, 8].

Methods

Fitness Model

We are considering systems where the fitness is proportional to the binding of a single molecule to its functional site, such as a transcription factor to its operator. In thermodynamic equilibrium, the binding probability for this system is [9]

$$p_+(E) = \frac{1}{1 + e^{\beta(E-F_0)}}, \quad (1)$$

where F_0 is the free energy of a random genome. This function has a sigmoid shape, which can be approximated as an exponential function close to one of the plateaus. In analytical computations we often approximate high binding probability plateau as $p_+(E) \sim (1 - e^{\beta(E-F_0)})$. In addition, we consider the fitness effect of the binding site length. Generally, we assume that longer binding sites come with an increased fitness cost, since genome size is under selection especially in prokaryotes (citation). Therefore, we include a linear fitness cost f_l per position of the binding site to the system,

$$F(E, l) = f_0 p_+(E, l) - f_l l, \quad (2)$$

where f_0 is the proportionality factor between binding probability to fitness. Note that the linear fitness cost for binding site length does not influence the dynamics of binding sites of fixed length, since the term maintains constant and selection coefficients are calculated as fitness differences, therefore canceling any constant additional term.

Binding Energy Model

We assume a minimal energy model, where each position contributes independently to the total binding energy of the sequence, which is called the independent nucleotide approximation and commonly used for Protein-DNA interactions [16, 17]. Therefore, we assume that minimal binding energy is achieved by a reference sequence, and each mismatch from that sequence brings a fixed cost to the binding energy of about $\epsilon\beta \approx 2 - 3$ [9]. In this work we fix the energy cost per mismatch to be $\epsilon\beta = 2$, independent of the actual nucleotide at the position. For specific examples, there are methods to obtain real energy matrices that make it possible to compute the actual binding energy of a transcription factor to its binding site (citation).

Due to the linearity of the model, the total binding energy will increase with the length of the binding site, for both unspecific and specific sequences. We assume that the total number of unspecific binding site maintains constant, hence, the free energy difference ΔE that is required to acquire specificity compared to off target binding sites maintains constant as well,

$$\Delta E = E_0(l) - E^*(l), \quad (3)$$

where we relabeled the free energy threshold in the sigmoid fitness landscape as $E^*(l)$ and $E_0(l) = 3/4l\epsilon$, since a random sequence has $1/4l$ just by chance. The free energy threshold is $\Delta E = 10k_B T$ (argue either from minimal length of binding site or other source).

Genetic Load and Length

One measure of adaption is the genetic load, which is computed as the difference of the mean fitness of a population to the maximal achievable fitness. Using the exponential approximation of the upper plateau of the binding term in the fitness landscape, the genetic load is simply given as the derivative of the fitness landscape evaluated at the mean energy $\mathcal{L} = \beta|f'(\Gamma)|$ [19]. The time evolution of the trait mean is given by the stochastic equation

$$\dot{\Gamma} = m^\gamma + \Delta_E f'(\Gamma) + \chi_\Gamma(t), \quad (4)$$

where

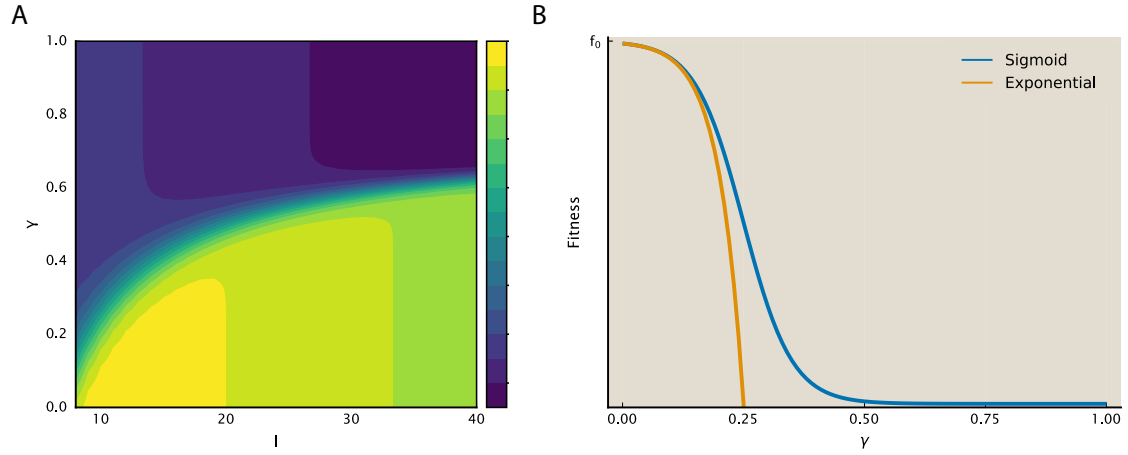


Fig. 1. (A) Two dimensional fitness landscape, which binding site length on the x-axis, and intensive binding energy $\gamma = E/l$ on the y-axis. (B) The fitness landscape for fixed binding site length (blue) and the exponential approximation (orange).

Mutation model

Results

Evolutionary Steady State in Non-equilibrium

We start by investigating how externally driven mutations that change the consensus sequence of the binding site force the system into a new state that is different from the steady state that can be observed when binding is only disrupted by mutations in the binding site. Therefore, we use the Fokker-Planck equation describing the evolution of the mean Γ of a quantitative trait E [11],

$$\frac{\partial}{\partial t} Q(\Gamma, t) = \left[\frac{\tilde{g}^{\Gamma\Gamma}}{2N} \frac{\partial^2}{\partial \Gamma^2} - \frac{\partial}{\partial \Gamma} (m^{\Gamma} + \tilde{g}^{\Gamma\Gamma} \tilde{s}_{\Gamma}) \right] Q(\Gamma, t). \quad (5)$$

The equation can be broken down into three evolutionary forces: genetic drift, selection and mutations. In the low mutation rate regime, the selection term can be written as $\tilde{s}_{\Gamma} = \partial_{\Gamma} f(\Gamma)$, i.e., the derivative of the fitness landscape at the value of the trait mean. The effective diffusion coefficient is given by $\tilde{g}^{\Gamma\Gamma} = \langle \Delta \rangle \sim \mu$ (include scaling with alphabet size here?), hence scaling the strength of selection with the variation of the trait within a population (mention Fisher's theorem here?). The mutational drift term is $m^{\Gamma} = -\frac{n}{n-1}(\Gamma - \Gamma_0)$, where n is the alphabet size of the binding site and Γ_0 is the mean trait in the absence of selection. The mutation term describes deterministic drift that sites undergo due to mutations. If this term dominates the deterministic behavior, then most sites will have binding energies similar to random sequences. Driver mutations have the same impact on the binding energy as trailer mutations. Hence, in the evolutionary description of traits, we can add the driver mutation rate ρ to the trailer mutation rate μ in the mutation term, $m^{\Gamma} \sim -(\rho + \mu)(\Gamma - \Gamma_0)$, for details see [methods oder supplementary](#). The selection term $\tilde{g}^{\Gamma\Gamma} \tilde{s}_{\Gamma}$ describes the selection coefficient of the trait mean value \tilde{s}_{Γ} and how it is amplified by the diversity within a population $\tilde{g}^{\Gamma\Gamma} = \langle \Delta_E \rangle$. For low mutation rates, the diversity within a populations is given by its value in the absence of selection, and very small. Most positions are monomorphic, thus a mutation in the driver sequence changes the binding energy for nearly all sites equally. Therefore the diversity within a population is not affected by driver mutations, and the selection term is invariant.

The non-deterministic term $\tilde{g}^{\Gamma\Gamma}/2N$ describes how mutations lead to a spread in mean binding energies across populations. In the original formulation this term also scales with the mean diversity within a population. However, driver mutations also lead to an increased diversity across populations, since they change the mean binding energy, while not effecting the mean diversity within a population. Hence, we have to distinguish between the two response terms $\tilde{g}^{\Gamma\Gamma}$. The response coefficient in the selection term, which only scales with the trailer mutation rate μ , is labeled $\tilde{g}_s^{\Gamma\Gamma} \sim \mu$. The response coefficient in the non deterministic term is scaling with both mutation rates and is labeled $\tilde{g}_d^{\Gamma\Gamma} \sim \mu + \rho$.

In the absence of selection, the steady state distribution $Q_0(\Gamma)$ is the same as in the absence of driving, since the scaling of both terms cancels out in the final solution (Methods). However, in a driven environment $\rho > 0$, the steady state solution becomes

$$Q(\Gamma) = \frac{1}{Z} Q_0(\Gamma) \exp \left[\frac{2N}{1+\kappa} F(\Gamma) \right], \quad (6)$$

where $\kappa = \rho/\mu$ is the ration of driver and trailer mutation rate, which is the key parameter describing the strength of non-equilibrium. A rescaling of the fitness landscape $\hat{F}(\Gamma) = F(\Gamma)/(1+\kappa)$ reproduces the steady state distribution in equilibrium with a weaker fitness landscape. This result is explained by the different scalings of the population genetic terms with non-equilibrium. The mutation drift towards smaller binding energies increases as well as diffusion across populations, while selection is not increased and therefore relatively weaker.

Note, that the rescaling of the fitness landscape effectively only applies to the term in the fitness landscape regarding the binding probability. The fitness length cost is absorbed in the normalization of the steady state distribution in 6. In the low mutation rate limit, an exact steady state distribution can be numerically computed from substitution rates (Methods or Supplementary). In the regime of small selection coefficients $Ns \sim 1$, an analytic steady state can be approximated (Methods or Supplementary), which is equal to the result in equation 6.

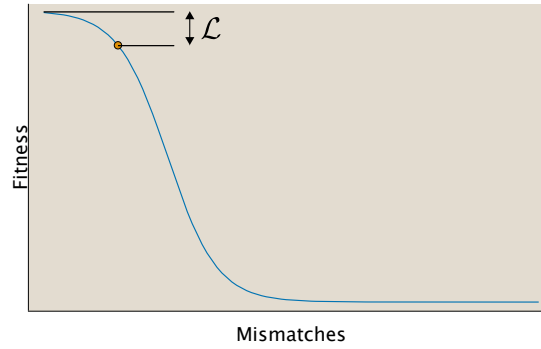


Fig. 2. This could be a figure to show what the genetic load is.

Optimal Binding Site Length

First, we compute the binding site length which minimizes the genetic load, therefore being the fittest and is expected to be the steady state that binding site lengths are evolving towards. Non-equilibrium weakens the ability to find specific binding sites, while the length fitness cost is independent of the state of the binding. Hence, the fitness trade-off changes with the level of non-equilibrium. We can consider the fitness effects in the terms of genetic load. For weak non-equilibrium most sites are functional and are at the edge of the upper fitness plateau, which can be approximated by an exponential tail. Then the genetic load of a population is given by $\mathcal{L} \sim \left| \frac{\partial}{\partial l} f(\Gamma, l) \right| + f_l l$. The average load across populations is given by the load at the deterministic mutation-selection balance $\hat{\Gamma}$ (Methods or Supplementary) (Cite Torsten and Daniel Paper). We retrieve two terms for the genetic load, with different scalings of the binding length. The linear fitness cost for binding length, and the binding probability term with scales inversely with binding length,

$$\mathcal{L} = \xi \frac{l_0}{l} (1 + \kappa) + \lambda \frac{l}{l_0} = \mathcal{L}_\kappa + \mathcal{L}_\lambda, \quad (7)$$

Where ξ is a combination of various parameters, $\lambda = 2Nf_l/l_0 \sim 1$ is the scaled fitness cost for length, and l_0 is a length scaling parameter (supplementary methods). A derivation of this equation can be found in the supplementary methods. Due to the trade off between the cost of adding positions to the binding site and reducing load by adding positions, there is an optimal binding site length, which minimizes the total genetic load,

$$l_{\text{opt}} \sim l_0 \sqrt{\frac{1+\kappa}{\lambda}}. \quad (8)$$

In equilibrium, the optimal binding site length is close to the length scaling parameter, $l \sim l_0$. 145

Dynamical Evolution of Binding Site Length 146

In a real population, length adaptation is a dynamical process, driven by mutations which include a new position in the sequence in its recognition. We assume these mutations to be very rare compared to other type of mutations, such that the timescales separate. Therefore, we can assume that the system is in the respective steady state of binding energies each time a length mutation occurs. Assuming an alphabet size of 4 and using the match/mismatch energy model (see Methods), a newly added position will be a match with probability $1/4$, since it has not been selected for yet and hence is random compared to the driver sequence. Length increase mutations are on average neutral in respect to the adaption term in the fitness landscape in first order, and are only influenced by the fitness cost of length (supplementary methods), 155

$$\sigma^+ = -\frac{\lambda}{l_0} = -\frac{\partial \mathcal{L}_\lambda}{\partial l}. \quad (9)$$

Since $\lambda \sim 1$, length increase mutations are in the regime of weak selection (**Here either an explanation of what this means or a nice reference**). However, length decrease mutations are affected by the fact that the binding site is adapted, i.e., removing a position is likely to hit an interaction which is beneficial to the binding interaction. If the binding site has k matches, then the probability of removing a match is k/l , which is much larger than the probability of adding a match, which we showed to be $1/4$. The selection coefficient for removing a position has an additional term, 162

$$\sigma^- = \frac{\lambda}{l_0} - \xi \epsilon \beta \left(\frac{l_0}{l} \right)^2 (1 + \kappa) = \frac{\partial \mathcal{L}_\lambda}{\partial l} + l_0 \epsilon \beta \frac{\partial \mathcal{L}_\kappa}{\partial l}. \quad (10)$$

Unexpectedly, the adaptation term is scaled by the effective length scaling parameter, and the mutational effect. Therefore, the selection coefficient for length decrease mutations shows marginal selection, $\sigma^- \sim (l/l_0)^2 \sim 1$. Due to the separation of time scales, we can write down substitution rates for length mutations, which lead to a steady state distribution. In first order, we can compute the substitution rates as 167

$$u_{+/-} \sim 1 + \frac{\sigma^{+/-}}{2} \sim \exp \left[\frac{\sigma^{+/-}}{2} \right]. \quad (11)$$

Using detailed balance (**Maybe avoid this term here?**), we can compute the steady state distribution for binding site lengths. Using the rates above, we find that the resulting distribution is 169

$$P(l) \simeq \exp \left[-\frac{\epsilon \beta l_0}{2} \mathcal{L}_\kappa - \mathcal{L}_\lambda \right]. \quad (12)$$

When comparing the maximum of this distribution with the length which minimizes the load, we find that the length l_{opt}^* which optimizes the effective potential resulting from the dynamical calculation is larger by a factor $\sqrt{l_0}$, 172

$$l_{opt}^* \simeq \sqrt{l_0} l_{opt}. \quad (13)$$

This is an enormous difference, see Figure 3(A), and shows how important it is to consider the dynamics of the evolution of binding site lengths. Transcription factor binding sites have lengths around 10bp in prokaryotes (**cite "Why are TF binding sites 10bp long."**), so we can tune the fitness cost for length f_l such that the optimal length computed from the dynamical computation is around that value in equilibrium ($\kappa = 0$). Then, we observe an increase in binding site length with increasing non-equilibrium, due to the sites being under more selective pressure to maintain their binding specificity, which shifts the weights in the trade off with the cost of increased binding site length. 180

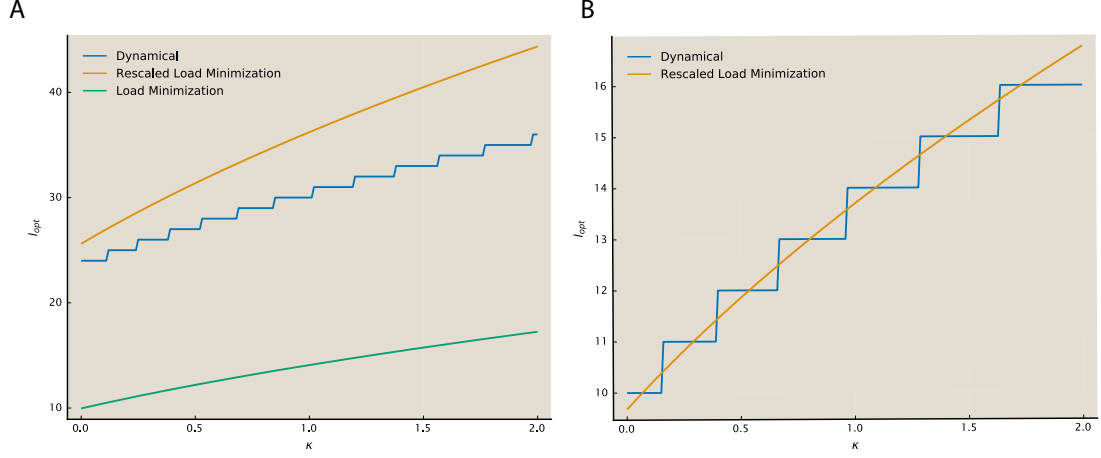


Fig. 3. (A) Optimal length computed from minimizing the genetic load (green), the minimum of the effective potential shown in equation (12) (blue), and the length minimizing the genetic load rescaled by $\sqrt{\ell_0}$ (orange) for parameters $f_0 = 50/2N$, $f_l = .2/2N$ and $N = 1000$. (B) Minimum of the effective potential shown in equation (12) (blue), and the length minimizing the genetic load rescaled by $\sqrt{\ell_0}$ (orange) for parameters $f_0 = 50/2N$, $f_l = 1.4/2N$ and $N = 1000$.

Comment 29. 7., edited 8. 9. Does the slow dynamics of ℓ lead to a steady state that is localized at ℓ_{opt} ?

1. Our model has a conditional stationary distribution $Q(k|\ell)$ peaked at $k/\ell \equiv \gamma$ with

$$\gamma - \gamma_0 = \frac{\ell_0}{\ell} \quad (14)$$

and $\gamma_0 = 1/4$. This stationary state has a scaled load $\mathcal{L} \equiv (F - F_{\text{max}})/\tilde{\sigma}$ given by

$$\mathcal{L} = \mathcal{L}_\kappa + \mathcal{L}_\lambda = (\gamma - \gamma_0)(1 + \kappa) + \lambda \frac{\ell}{\ell_0}, \quad (15)$$

leading to an optimum length

$$\ell_{\text{opt}} = \ell_0 \sqrt{\frac{1 + \kappa}{\lambda}}. \quad (16)$$

2. In the exponential regime of the (scaled) fitness landscape $\mathcal{F}(k|\ell) \equiv F(k|\ell)/\tilde{\sigma}$, which is relevant for stable sites, the scaled load equals the scaled selection coefficient of site mutations,

$$s = \frac{\partial \mathcal{F}(k, \ell)}{\partial k} = \mathcal{L}_\kappa = (\gamma - \gamma_0)(1 + \kappa); \quad (17)$$

prefactors to be double-checked with Torsten-Daniel paper.

3. Perturbation theory in $\epsilon \equiv 1/\ell_0$. We can distinguish terms of order ϵ and of order $\ell/\ell_0 \sim 1$. We have $\gamma \sim \lambda \sim 1$ and $\mathcal{L} \sim 1$; the scaling of λ follows from the requirement $\ell_{\text{opt}} \approx \ell_0$ at $\kappa = 0$. In particular, we can distinguish selection coefficients $\sim 1/\ell_0$ (weak selection) and ~ 1 (marginal selection).
4. Adiabatic substitution dynamics of ℓ . Length changes are assumed to occur with much smaller rates than mutations; we assume they are emitted from the steady state of the k dynamics. Each length change is coupled to a stochastic change of k . Because the dynamics takes place in the regime of weak and marginal selection, we use a linear approximation for substitution probabilities. Hence, we can first compute average selection coefficients \bar{s}_+ and \bar{s}_- for changes of ℓ and then evaluate the corresponding substitution rates using these averages.

In this framework, we find the processes $\ell \rightarrow \ell + 1$ with average selection coefficient

$$\bar{s}_+ = \frac{3}{4} \frac{(-s)}{4} + \frac{1}{4} \frac{3s}{4} - \frac{\lambda}{\ell_0} \quad (18)$$

$$= -\frac{\lambda}{\ell_0} \quad (19)$$

(near-neutral evolution) and the processes $\ell \rightarrow \ell - 1$ with average scaled selection coefficient

$$\bar{s}_- = \gamma \frac{(-3s)}{4} + (1 - \gamma) \frac{s}{4} + \frac{\lambda}{\ell_0} \quad (20)$$

$$= -(\gamma - \gamma_0)s + \frac{\lambda}{\ell_0} \quad (21)$$

$$= -(\gamma - \gamma_0)^2(1 + \kappa) + \frac{\lambda}{\ell_0} \quad (22)$$

$$= \ell_0 \frac{\partial \mathcal{L}_\kappa}{\partial \ell} + \frac{\partial \mathcal{L}_\lambda}{\partial \ell} \quad (23)$$

(marginal selection). These selection coefficients are of different magnitude: length increases are constrained only by the weak linear term, length decreases are marginally constrained by conservation of site function.

5. These asymmetric dynamics define an adaptive ratchet, which acts to increase selection for complexity. A typical cycle $\ell \rightarrow \ell + 1 \rightarrow \ell$ has the following average scaled fitness balance:

- substitution $\ell \rightarrow \ell + 1$: weak fitness loss, $\bar{s}_+ = -\lambda/\ell_0 = -O(\epsilon)$;
- equilibration at $\ell + 1$: weak fitness gain, $\Delta \mathcal{L}_{\ell+1} = \mathcal{L}_\kappa(\ell) - \mathcal{L}_\kappa(\ell + 1) \sim \ell_0/\ell^2 = O(\epsilon)$;
- substitution $\ell + 1 \rightarrow \ell$: marginal fitness loss, $\bar{s}_- = -\ell_0^2/\ell^2 = -O(\epsilon^0)$;
- equilibration at ℓ : marginal fitness gain, $\Delta \mathcal{L}_\ell = -\mathcal{L}_\kappa(\ell) + \mathcal{L}_\kappa(\ell + 1) - \bar{s}_- = O(\epsilon^0)$.

This cycle could inform a Figure.

6. Effective fitness landscape for length changes. We compute the ratio of the forward and backward substitution rate in first-order approximation,

$$\frac{u_{\ell \rightarrow \ell+1}}{u_{\ell+1 \rightarrow \ell}} = \frac{\exp(\bar{s}_+/2)}{\exp(\bar{s}_-/2)} + O(s^2) \quad (24)$$

$$\simeq \exp \left[-\frac{\ell_0}{2} (\mathcal{L}_\kappa(\ell + 1) - \mathcal{L}_\kappa(\ell)) - (\mathcal{L}_\lambda(\ell + 1) - \mathcal{L}_\lambda(\ell)) \right] \quad (25)$$

$$\equiv \exp [\mathcal{F}_{\text{eff}}(\ell + 1) - \mathcal{F}_{\text{eff}}(\ell)] \quad (26)$$

with an effective scaled fitness potential

$$\mathcal{F}_{\text{eff}}(\ell) = -\frac{\ell_0}{2} \mathcal{L}_\kappa(\ell) - \mathcal{L}_\lambda(\ell). \quad (27)$$

These rates define an equilibrium distribution

$$Q(\ell) \sim \exp[\mathcal{F}_{\text{eff}}(\ell)], \quad (28)$$

which is peaked at

$$\ell^* = \ell_0 \sqrt{\frac{\ell_0(1 + \kappa)}{2\lambda}} = \ell_{\text{opt}} \sqrt{\frac{\ell_0}{2}}. \quad (29)$$

Compared to the naive distribution $Q(\ell) \sim \exp[-\mathcal{L}(\ell)]$, which would lead to a peak at ℓ_{opt} , the ratchet mechanism enhances the evolution of complexity.

7. Summary.

- Non-equilibrium weakens the selection on point mutations at a given complexity, as expressed by an effective fitness landscape

$$\mathcal{F}_{\text{eff}}(k|\ell) \sim \frac{1}{1 + \kappa}. \quad (30)$$

- Non-equilibrium introduces ratchet selection for complexity, as expressed by an effective fitness landscape

$$\mathcal{F}_{\kappa, \text{eff}}(\ell) \sim \ell_0(1 + \kappa). \quad (31)$$

This landscape is proportional to ℓ_0 ; that is, complexity begets more complexity.

8. Questions:

- (1) Is this in qualitative accordance with the evaluation of the rates in (26) beyond first order?
- (2) How does the fitness balance compare with the numerics?
- (3) How does ℓ^* compare with the numerics?

Data

Discussion

Here we should discuss the meaning of our results. The two factors of increased binding site length. First, the dynamics of length evolution lead to a much higher binding site length in equilibrium, than observed from minimizing the genetic load. But. Also we should discuss which data could support our theory. And finally, the shortcoming of the theory. Like how the approximation of the distribution is not really good for higher selection coefficients. And what how it could be used to improve. And finally how we think experiments could be set up to test the hypothesis experimentally.

References

1. Tuğrul, M., Paixão, T., Barton, Nicholas H. & Tkačik, G. Dynamics of Transcription Factor Binding Site Evolution. *PLOS Genetics* **11** 1-28 (2015)
2. Schmidt, D. & et.al. Five-vertebrate ChIP-seq reveals the evolutionary dynamics of transcription factor binding. *Science* **328** 1036–1040 (2010)
3. Nowick, K. & Stubbs, L. Lineage-specific transcription factors and the evolution of gene regulatory networks. *Briefings in Functional Genomics* **9** 65–78 (2010)
4. Lynch, V. J. & Wagner, G. P. Resurrecting the role of transcription factor change in developmental evolution. *Evolution: International Journal of Organic Evolution* **62** 2131–2154 (2008)
5. Wagner, G. P. & Lynch, V. J. The gene regulatory logic of transcription factor evolution. *Trends in Ecology & Evolution* **23** 377 - 385 (2008)
6. Voordeckers, K., Pougach, K. & Verstrepen, K. J. How do regulatory networks evolve and expand throughout evolution?. *Current Opinion in Biotechnology* **34** 180–188 (2015)
7. Lozada-Chávez, I., Janga, S. C. & Collado-Vides, J. Bacterial regulatory networks are extremely flexible in evolution. *Nucleic Acids Research* **34** 3434–3445 (2006)
8. Perez, J. C. & Groisman, E. A. Evolution of Transcriptional Regulatory Circuits in Bacteria, *Cell* **138** 233–244 (2009)
9. Lässig, M. From biophysics to evolutionary genetics: statistical aspects of gene regulation. *BMC bioinformatics* **8**, 1 (2007).

10. Babu, M. M., Teichmann, S. A. & Aravind, L. Evolutionary Dynamics of Prokaryotic Transcriptional Regulatory Networks. *Journal of Molecular Biology* **358** 614 - 633 (2006).
11. Nourmohammad, A., Schiffels, S. & Lässig, M. Evolution of molecular phenotypes under stabilizing selection. *Journal of Statistical Mechanics: Theory and Experiment* **2013**, P01012 (2013).
12. Mustonen, V. & Lässig, M. Fitness flux and ubiquity of adaptive evolution. *Proceedings of the National Academy of Sciences* **107**, 4248–4253 (2010).
13. Berg, Otto G. & von Hippel, Peter H. Selection of DNA binding sites by regulatory proteins: Statistical-mechanical theory and application to operators and promoters. *Journal of molecular biology* **193**, 723–743 (1987).
14. Kimura, M. On the probability of fixation of mutant genes in a population. *Genetics* **47**, 713 (1962).
15. Berg, J. & Lässig, M. Stochastic evolution of transcription factor binding sites. *Biophysics* **48**, 36–44 (2003).
16. Stormo, Gary D. & Fields, Dana S. Specificity, free energy and information content in protein–DNA interactions. *Trends in biochemical sciences* **23**, 109–113 (1998).
17. Djordjevic, M. SELEX experiments: new prospects, applications and data analysis in inferring regulatory pathways. *Biomolecular engineering* **24** 179–189 (2007).
18. Stormo, Gary D. & Stekel, Dov J. Determining the specificity of protein–DNA interactions. *Nature Reviews Genetics* **11** 751 (2010).
19. Held, T., Klemmer, T., & Lässig, M. Survival of the Simplest in Microbial Evolution.. *Nature Communications* **10** (2019).

Acknowledgments

Author Contributions

Competing Interests

The authors declare no competing interests.

Supplementary Information

Supplementary Methods 1 Steady State Approximations.

In the regime of low mutation rates, $\mu N \ll 1$, a population is in a monomorphic state, i.e., dominated by a single genotype, for most of the time. A new mutation fixes in the population with probability $p(s)$ given by [14],

$$p(s) = \frac{1 - e^{-2s}}{1 - e^{-2Ns}}, \quad (32)$$

where s is the selection coefficient of the mutation and N the effective population size. Then, substitutions occur at a rate

$$u(s)_{E \rightarrow E'} = N\mu_{E \rightarrow E'}p(s), \quad (33)$$

where $\mu_{E \rightarrow E'}$ is the rate at which a mutation occurs which changing the binding energy from E to E' . Here, we assumed that the fixation of a mutation is under selective pressure, which is the case for mutations in the trailer. We investigate the effects of mutations in the driver, which occur with rate ρ . These mutations are results of external events, e.g., a mutation in the coding

sequence of a transcription factor leading to a missense mutation. In isolated system of driver and trailer, these mutations are not under selection pressure, but can change the binding energy nonetheless. Here we assume that these mutations happen for all driver sequences at the same time and therefore can be treated as substitutions as well. Thus, we add a term to the mutation rate in the form of

$$u(s)_{E \rightarrow E'} = N\mu_{E \rightarrow E'}p(s) + \rho_{E \rightarrow E'}, \quad (34)$$

where $\mu_{E \rightarrow E'}$ is the mutation rate with driver mutations change the binding energy from E to E' . Assume, both trailer and driver have the same alphabet size, then both mutation rates are given by a constant rate times a factor accounting for the entropic difference between both states, $\mu_{E \rightarrow E'} = \mu w_{E \rightarrow E'}$ and $\rho_{E \rightarrow E'} = \rho w_{E \rightarrow E'}$. Even if the assumption of equal alphabet sizes doesn't hold, e.g., in protein-DNA interaction, the differences in alphabet sizes can be accounted for by rescaling the mutation rates by a constant factor. The entropic term w is determined by the substitution rates in the absence of selection. We are looking for the distribution of binding energies at steady state, which is equivalent to

$$\frac{u(0)_{E \rightarrow E'}}{u(0)_{E' \rightarrow E}} = \frac{\mu_{E \rightarrow E'}}{\mu_{E' \rightarrow E}} = \frac{w_{E \rightarrow E'}}{w_{E' \rightarrow E}} = \frac{Q_0(E')}{Q_0(E)}, \quad (35)$$

where $Q_0(E)$ is the distribution of binding energies in the absence of fitness. Even though there is non-equilibrium on the level of genotypes, we can look for a steady state distribution of the binding energies. Since the dynamics are one dimensional and fully described by the substitution rates, we can find the steady state distribution $Q(k)$ by imposing detailed balance and solving

$$\frac{Q(E')}{Q(E)} = \frac{u(s)_{E \rightarrow E'}}{u(-s)_{E' \rightarrow E}} = \frac{N\mu_{E \rightarrow E'}p(s) + \rho_{E \rightarrow E'}}{N\mu_{E' \rightarrow E}p(s) + \rho_{E' \rightarrow E}}. \quad (36)$$

In equilibrium, $\rho = 0$, the fraction has an exact solution [15],

$$Q(k) = Q_0(k) \exp[2NF(k)]. \quad (37)$$

In non-equilibrium $\rho > 0$, the fraction does not simplify to an exact term as equation 37, but we can compute the exact steady state distribution numerically. However, in the limit of small selection coefficients $Ns \sim 1$, the can expand the fraction in orders of the selection coefficient and retrieve,

$$\begin{aligned} \frac{Q_0(E')}{Q_0(E)} \frac{N\mu p(s) + \rho}{N\mu p(s) + \rho} &\approx \frac{Q_0(E')}{Q_0(E)} \left[1 + \frac{2s(N-1)}{1 + \frac{\rho}{\mu}} + \frac{2s^2(N-1)^2}{(1 + \frac{\rho}{\mu})^2} + \mathcal{O}(s^3) \right], \\ &\approx \frac{Q_0(E')}{Q_0(E)} \exp \left[\frac{2N}{1 + \kappa} F(E) \right]. \end{aligned} \quad (38)$$

Populations are usually large, so $N-1 \approx N$, and up to second order the series is equal to an exponential function. The resulting distribution is therefore given by

$$Q(E) = Q_0(E) \exp \left[\frac{2N}{1 + \kappa} F(E) \right]. \quad (39)$$

Supplementary Methods 2 Generic Mutation Drift and Driver mutations.

Here we use the population genetic notation of quantitative traits [11], and give a general form of the mutation drift. In addition we compute the contribution of driver mutations to the mutation drift.

The frequency of nucleotide i at locus k is given by y_k^i . The mutation drift of trailer mutations is therefore given by

$$m^\Gamma|_a = \sum_{k=1}^l \sum_{i>j}^n \mu_{j \rightarrow i} \Delta_{j \rightarrow i} (y_k^i - y_k^j), \quad (40)$$

where $\mu_{j \rightarrow i}$ is the mutation rate from nucleotide j to nucleotide i and $\Delta_{j \rightarrow i}$ is the energy difference of nucleotide j to i , $\Delta_{j \rightarrow i} = \epsilon_j - \epsilon_i$. In the absence of any mutational bias, the mutation rates are equal, $\mu_{j \rightarrow i} = \mu/n - 1$, where μ is the total mutation rate per position. Note that self mutations are excluded by definition. Inserting into the expression,

$$m^\Gamma|_a = \frac{\mu}{n-1} \sum_{k=1}^l \sum_{i>j}^n (\epsilon_k^j - \epsilon_k^i) (y_k^i - y_k^j), \quad (41)$$

we can multiply out the brackets to

$$m^\Gamma|_a = \frac{\mu}{n-1} \sum_{k=1}^l \left((1-n) \sum_{j=1}^n \epsilon_k^j y_k^j + \sum_{i=1}^n \epsilon_k^i \sum_{j \neq i}^n y_k^j \right). \quad (42)$$

Now we use $\sum_i y_k^i = 1$, such that we can write $\sum_{j \neq i} y_k^j = 1 - y_k^i$,

$$\begin{aligned} m^\Gamma|_a &= \frac{\mu}{n-1} \sum_{k=1}^l \left((1-n) \sum_{i=1}^n \epsilon_k^i y_k^i + \sum_{i=1}^n \epsilon_k^i (1 - y_k^i) \right), \\ &= \frac{\mu}{n-1} \sum_{k=1}^l \left(-n \sum_{i=1}^n \epsilon_k^i y_k^i + \sum_{i=1}^n \epsilon_k^i \right). \end{aligned} \quad (43)$$

The mean trait value is given by $\Gamma = \sum_{k=1}^l \sum_{i=1}^n E_i y_k^i$, and the neutral mean value by $\Gamma_0 = \frac{1}{n} \sum_{k=1}^l \sum_{i=1}^n E_i$, therefore the mutation drift from trailer mutations is given by

$$m^\Gamma|_a = -\mu \frac{n}{n-1} (\Gamma - \Gamma_0) \quad (44)$$

In our model driver mutations occur as substitutions, therefore the nucleotide frequencies in the trailer do not change. In general, with the new driver element, each possible nucleotide at that position can have a different contribution to the binding energy. This difference is weighted with the allele frequency of each nucleotide. The energy of nucleotide i with driver element α is denoted by E_i^α , the energy at position before the mutation by E_i^* . The number of possible driver elements is given by β , and the mutation rate of driver element $*$ to α by $\rho_{* \rightarrow \alpha}$. Then the contribution to the mutation drift from driver mutations is given by

$$m^\Gamma|_b = \sum_{k=1}^l \sum_{\alpha=1}^\beta \sum_{i=1}^n \rho_{* \rightarrow \alpha} (E_i^\alpha - E_i^*) y_k^i. \quad (45)$$

Again, we assume that there is no mutational bias, $\rho_{* \rightarrow \alpha} = \rho/(\beta - 1)$. We can transform the sum to

$$\begin{aligned} &= -\frac{\rho}{\beta-1} (\beta \Gamma - \sum_{k=1}^l \sum_{\alpha=1}^\beta \sum_{j=1}^n \epsilon_\alpha^j y_k^j), \\ &= -\rho \frac{\beta}{\beta-1} (\Gamma - \Gamma'_0), \end{aligned} \quad (46)$$

where Γ'_0 is the average over all possible binding energies with given trailer sequence. In general, 339
this average is close to the neutral mean binding energy Γ_0 . Then, the total mutation drift sums 340
to 341

$$m^\Gamma = - \left(\mu \frac{n}{n-1} + \rho \frac{\beta}{\beta-1} \right) (\Gamma - \Gamma_0). \quad (47)$$

In the case if equal alphabet sizes of trailer and driver, the mutation drift takes the simple form 342

$$m^\Gamma = - \frac{n}{n-1} \mu (1 + \kappa) (\Gamma - \Gamma_0), \quad (48)$$

with the non-equilibrium parameter $\kappa = \rho/\mu$. If there are different alphabet sizes, then rescaling 343
one of the mutation rates will give the same result. 344

Supplementary Methods 3 Mutation-Selection-Balance and Non-Equilibrium.

The time evolution of the mean binding energy is given by [19]

$$\dot{\Gamma} = m^{\Gamma} + \Delta_E f'(\Gamma) + \chi_{\Gamma}(t), \quad (49)$$

with $m^{\Gamma} = -\frac{n}{n-1}\mu(1+\kappa)(\Gamma - \Gamma_0)$, where the mutation rate μ is the absolute mutation per position. The mutation selection balance (MSB) is the deterministic fix point of this equation $\dot{\Gamma}_{\text{MSB}} = 0$. In the low mutation rate regime, the diversity within a population Δ_E is given by the neutral expectation, $\Delta_E = \mu Nl$ (This is from Armita's paper. Is this true for all Alphabet sizes, or do we need to rescale this somehow? For this calculation, I will assume that this term scales with the alphabet size the same way the mutational drift does. However, the scalings are perfect, if $\Delta_E \sim \left(\frac{n}{n-1}\right)^2$. So that is what I am going to do for now in these notes, but please double check this). With the mutation drift from equation 48, we get (Added β^{-2} to have matching units, maybe refine something here such that this is not an issue.)

$$0 = -\frac{n}{n-1}\mu(1+\kappa)(\Gamma_{\text{MSB}} - \Gamma_0) + \left(\frac{n}{n-1}\right)^2 \mu Nl f'(\Gamma_{\text{MSB}}) \beta^{-2}. \quad (50)$$

Note that the MSB is invariant under the limit of a vanishing mutation rate $\mu \rightarrow 0$. At close to the upper plateau of the sigmoid fitness landscape the landscape can be approximated by an exponential landscape with derivative

$$\partial_{\Gamma} f(\Gamma, l) = -f_0 \beta \exp[\beta(\Gamma - (\Gamma_0 - \Delta_E))], \quad (51)$$

where Γ_0 is the mean energy in a population in the absence of selection. In the model of matches and mismatches, this is simply $\Gamma_0 = \epsilon 3l/4$ in a four letter alphabet, where ϵ is the energy contribution of a mismatch.

Then, we can rewrite equation 50 to retrieve the transcendental equation

$$\Gamma_{\text{MSB}} = \Gamma_0 + \frac{1}{\beta} \frac{n}{n-1} \frac{f_0 N}{1+\kappa} \exp[-\beta(\Gamma_{\text{MSB}} - (\Gamma_0 - \Delta_E))], \quad (52)$$

which is solved by the ProductLog,

$$\Gamma_{\text{MSB}}(l) = \Gamma_0 - \frac{1}{\beta} \text{ProductLog}\left(\frac{n}{n-1} \frac{N f_0 l e^{\beta \Delta_E}}{1+\kappa}\right). \quad (53)$$

We can rewrite this result in terms of relative number of mismatches $\gamma/\epsilon l$,

$$\gamma_{\text{MSB}}(l) - \gamma_0 = -\frac{1}{l \epsilon \beta} \text{ProductLog}\left(\frac{n}{n-1} \frac{N f_0 l e^{\beta \Delta_E}}{1+\kappa}\right). \quad (54)$$

Although the ProductLog scales with the binding site length l and the non-equilibrium ratio κ , we can treat it as a constant (refer to figure.).

Defining the constant

$$l_0 = \frac{1}{\beta \epsilon} \text{ProductLog}\left(\frac{n}{n-1} \frac{N f_0 l e^{\beta \Delta_E}}{1+\kappa}\right), \quad (55)$$

we can write

$$\gamma_{\text{MSB}}(l) - \gamma_0 \approx -\frac{l_0}{l}. \quad (56)$$

Having equation (53) in hand, we can solve equation (50) for the fitness derivative at the Mutation-Selection-Balance, $f'(\Gamma_{\text{MSB}}(l)) = \hat{f}(l)$,

$$\hat{f}(l) = \beta \frac{n-1}{n} \frac{1+\kappa}{Nl} \text{ProductLog}\left(\frac{n}{n-1} \frac{N f_0 l \beta e^{\beta \Delta_E}}{1+\kappa}\right), \quad (57)$$

which we can also write as

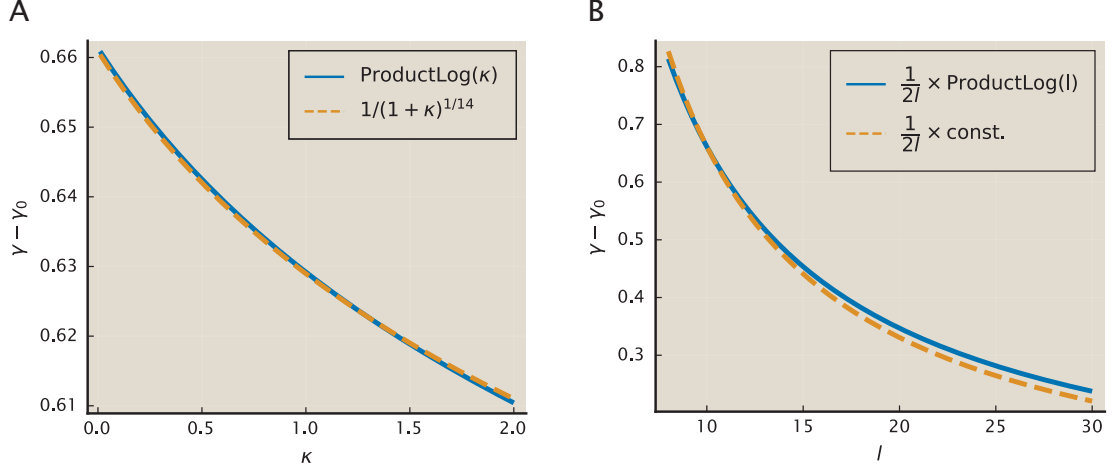


Fig. 4. Comparing the scaling of equation (54). The Productlog does not influence the scaling by a lot.

$$\hat{f}(l) = (1 + \kappa) \frac{n-1}{n} \frac{l_0}{Nl} \beta^2 \epsilon, \quad (58)$$

In the exponential fitness landscape, the scaled genetic load is given by the derivative of the fitness landscape in respect of the energy and the length fitness cost, 373
374

$$\mathcal{L} = 2N \left(-\frac{f'(\Gamma)}{\beta} + f_l l. \right). \quad (59)$$

Also we know, that the deterministic load is equal to the average load, 375

$$\mathcal{L} = 2N \left(-\frac{\hat{f}(l)}{\beta} + f_l l \right) \quad (60)$$

Finally, we can write the genetic load as a function of the driving parameter κ and the binding length, 376
377

$$\mathcal{L}(l, \kappa) = 2(1 + \kappa) \epsilon \beta \frac{n-1}{n} \frac{l_0}{l} + \lambda \frac{l}{l_0}, \quad (61)$$

with the scaled fitness cost $\lambda = 2N f_l / l_0$. 378

This load is minimized for given parameters at length 379

$$l_{opt} = l_0 \sqrt{\frac{2\epsilon\beta(n-1)}{n\lambda}} (1 + \kappa). \quad (62)$$

Supplementary Methods 4 Dynamical Length Evolution.

The slow dynamics of binding site length evolution leads to an asymmetry. As shown in equation (56), the conditional stationary distribution $Q(k|l)$ is peaked at $\gamma - \gamma_0 = l_0/l$. We compute the selection coefficients of length mutations in first order, using the expansion of the exponential fitness landscape

$$F(\Gamma, l) = F(\Gamma', l') + (\Gamma - \Gamma')f'(\Gamma', l') + (l - l') \left(\frac{3}{4}\epsilon f'(\Gamma', l') - f_l \right), \quad (63)$$

where we used $f'(\Gamma', l') = \partial_\Gamma F(\Gamma', l') = 4/3\epsilon \partial_l F(\Gamma', l')$. The selection coefficients for length mutations are the given by,

$$s_+^M = F(\Gamma, l+1) - F(\Gamma, l) = -f_l - \frac{3}{4}\epsilon f'(\Gamma, l), \quad (64)$$

$$s_+^{MM} = F(\Gamma + \epsilon, l+1) - F(\Gamma, l) = -f_l + \frac{1}{4}\epsilon f'(\Gamma, l), \quad (65)$$

$$s_-^M = F(\Gamma, l) - F(\Gamma, l+1) = f_l + \frac{1}{4}\epsilon f'(\Gamma, l), \quad (66)$$

$$s_-^{MM} = F(\Gamma - \epsilon, l) - F(\Gamma, l+1) = f_l - \frac{3}{4}\epsilon f'(\Gamma, l), \quad (67)$$

where the subscript is indicating increase (+) or decrease mutation (-) and the superscript is indicating whether the mutation is regarding a match(M) or mismatch(MM) When adding a position, a match is added with probability 1/4, and mismatch with probability 3/4. The average selection coefficient of a length increase mutation is then given by

$$\bar{s}_+ = \frac{1}{4}s_+^M + \frac{3}{4}s_+^{MM} = -f_l. \quad (68)$$

When removing a position, the probability of removing a match is equal to the ratio of matches in the binding site $1 - \gamma$, hence the average selection coefficient for a decrease mutation is given by

$$\bar{s}_- = \gamma s_-^{MM} + (1 - \gamma)s_-^M = f_l - (\gamma - \gamma_0)f'(\Gamma, l)\epsilon. \quad (69)$$

Now we can insert equations (56) and (58). Then, the scaled selection coefficients $\sigma = 2Ns$ are given by

$$\sigma_+ = -\frac{\lambda}{l_0}, \quad (70)$$

$$\sigma_- = \frac{\lambda}{l_0} - 2 \left(\frac{l_0}{l} \right)^2 \frac{n-1}{n} (1 + \kappa) \beta^2 \epsilon^2. \quad (71)$$

Here we can see the asymmetry of length mutations. Length increase mutations are only constrained by the small fitness cost of each position, with $\sigma_- \sim 1/l_0$ (near-neutral evolution), while length decrease mutations are marginally constrained by conservation of site function with $\sigma_+ \sim l_0^2/l^2 \sim 1$ (marginal selection). We can write the selection coefficients as derivatives of the load,

$$\sigma_+ = -\frac{\partial \mathcal{L}_l}{\partial l}, \quad (72)$$

$$\sigma_- = l_0 \frac{\partial \mathcal{L}_\kappa}{\partial l} \epsilon \beta + \frac{\partial \mathcal{L}_l}{\partial l}. \quad (73)$$

We have this extra factor of $\epsilon\beta$, which is a factor of 2. Let's see if we can get rid of it somehow? Equation 56 might be the right equation to do so. Due to the separation of timescales, we can compute a steady state distribution of binding lengths $P(l)$. This distribution is obeying detailed balance, therefore we can compute

$$\frac{P(l)}{P(l+1)} = \frac{u_-(l+1)}{u_+(l)}, \quad (74)$$

where u_+ is the length increase substitution rate and u_- is the length decrease mutation rate. In first order we can rate the ratio of the to rates as

$$\frac{u_-(l+1)}{u_+(l)} = \frac{\exp(\sigma_+/2)}{\exp(\sigma_-/2)} + \mathcal{O}(\sigma^2), \quad (75)$$

$$\simeq \exp \left[-\frac{l_0 \beta \epsilon}{2} (\mathcal{L}_\kappa(l+1) - \mathcal{L}_\kappa(l)) - (\mathcal{L}_\lambda(l+1) - \mathcal{L}_\lambda(l)) \right], \quad (76)$$

$$= \exp [\mathcal{F}_{\text{eff}}(l+1) - \mathcal{F}_{\text{eff}}(l)], \quad (77)$$

which the scaled effective potential

$$\mathcal{F}_{\text{eff}}(l) = -\epsilon \beta \frac{l_0}{2} \mathcal{L}_\kappa(l) - \mathcal{L}_\lambda(l). \quad (78)$$

These rates define an equilibrium distribution

$$Q(l) \sim \exp [\mathcal{F}_{\text{eff}}(l)], \quad (79)$$

which is peaked at

$$l_{\text{opt}} = l_0 \epsilon \beta \sqrt{\frac{l_0}{2} \frac{2(n-1)}{n\lambda} (1 + \kappa)}. \quad (80)$$