

Lengthscales and Cooperativity in DNA Bubble Formation

Z. Rapti¹, A. Smerzi^{2,3}, K. Ø. Rasmussen² and A. R. Bishop²

¹ *Center for Nonlinear Studies, Los Alamos National Laboratory,
Los Alamos, New Mexico 87545, and School of Mathematics,
Institute for Advanced Study, Einstein Drive 1, Princeton, New Jersey 08540*

² *Theoretical Division and Center for Nonlinear Studies,*

Los Alamos National Laboratory, Los Alamos, New Mexico 87545 and

³ *Istituto Nazionale di Fisica per la Materia BEC-CRS, Università di Trento, I-38050 Povo, Italy*

C.H. Choi, and A. Usheva

*Beth Israel Deaconess Medical Center and Harvard Medical School,
99 Brookline Avenue, Boston, Massachusetts 02215*

(Dated: November 6, 2018)

It appears that thermally activated DNA bubbles of different sizes play central roles in important genetic processes. Here we show that the probability for the formation of such bubbles is regulated by the number of soft AT pairs in specific regions with lengths which at physiological temperatures are of the order of (but not equal to) the size of the bubble. The analysis is based on the Peyrard-Bishop-Dauxois model, whose equilibrium statistical properties have been accurately calculated here with a transfer integral approach.

The genetic code underlying all forms of life is encoded in the DNA molecule by the four bases guanine (G), thymine (T), adenine (A), and cytosine (C) strung along a sugar-phosphate backbone in a particular sequence. The four bases are, through hydrogen bonding, pairwise complementary (A-T and G-C) allowing the coding strand and its complement to form the characteristic double helical DNA macromolecule. Although this construct is extraordinarily stable, it is clearly necessary that the double strands be separated in biological processes, including gene transcription, where the code is read by the appropriate protein machinery in the cell. It has long been an experimental fact [1] that the DNA double-strand can be thermally destabilized locally to form temporary single stranded “bubbles” in the molecule. This local melting is made possible by the entropy gained by transitioning from the very rigid double-strand to the much more flexible single-strand, which already at biologically relevant temperatures can balance the energy cost of breaking a few base pairs. Considering this entropic effect together with the inherent energetic heterogeneity – GC base pairs are 25 % more strongly bound than the AT bases – of a DNA sequence, it is conceivable that certain regions (subsequences) are more prone to such thermal destabilization than others: This has been confirmed by model calculations as well as experiments. We have previously argued [2] that such regions may indeed experimentally coincide with transcription initiation and regulatory sites. In this way, the DNA molecule may help initiate its own transcription by containing bubble forming subsequences at the crucial positions in the sequence where the transcription machinery assembles and engages its operation. If a robust general link between the formation of large thermal bubbles and transcription initiation is sufficiently established, it becomes crucially important

to be able to accurately predict the subsequence of DNA with propensity for the formation of bubbles of appropriate sizes.

Here we show that the probabilities of finding bubbles extending over n sites do not depend on a specific DNA subsequences. Rather, such probabilities depend on the density of soft A/T base pairs within specific regions of length $L(\kappa)$. This characteristic length is of the order of the size n of the bubble at physiological temperatures, but it diverges as the DNA melting temperature is approached. Our results are based on a calculation of the thermal equilibrium statistical properties of the Peyrard-Bishop-Dauxois (PBD) model [3, 4] using a transfer integral operator (TIO) technique. This model constitutes a very powerful tool to not only predict bubble formation probability in a given sequence but also to understand the underlying physical mechanisms [5]. Our previous study of the PBD model has been performed using Langevin [2, 6] and Monte Carlo techniques [7]. However, since our interest is centered on a very small portion of the thermodynamical equilibrium state, namely on the formation of large bubbles, dynamical and iterative samplings as offered by these methods are not very efficient. Therefore, we have developed here a semi-analytic approach based on the TIO [8, 9] that allows us to efficiently calculate relevant thermodynamical probabilities.

The potential energy of the PBD model, in its simplest form, reads

$$E = \sum_{\kappa=1}^N [V(y_n) + W(y_n, y_{n-1})] = \sum_{\kappa=1}^N \mathcal{E}(y_n, y_{n-1}), \quad (1)$$

where $V(y_n) = D_n(e^{-a_n y_n} - 1)^2$, represents the nonlinear hydrogen bonds between the bases. $W(y_n, y_{n-1}) = \frac{k}{2} (1 + \rho e^{-b(y_n + y_{n-1})}) (y_n - y_{n-1})^2$ is the nearest-neighbor coupling that represents the (nonlinear) stack-

ing interaction between adjacent base pairs: it is comprised of a harmonic coupling with a state depended coupling constant effectively modeling the change in stiffness as the double strand is opened (i.e. entropic effects). This nonlinear coupling results in long-range cooperative effects, leading to a sharp entropy-driven denaturation transition [3, 10]. The sum in Eq.(1) is over all base-pairs of the molecule and y_n denotes the relative displacement from equilibrium bases at the n^{th} base pair. The importance of the heterogeneity of the sequence is incorporated by assigning different values to the parameters of the Morse potential, depending on the the base-pair type. The parameter values we have used are those from Refs. [11, 12] chosen to reproduce a variety of thermodynamic properties.

Transfer Integral Method. All equilibrium, thermodynamic properties of the model (1) can be obtained through the partition function

$$\begin{aligned} \mathcal{Z} &= \int \prod_{\kappa=1}^N dy_n e^{-\beta\mathcal{E}(y_n, y_{n-1})} \\ &= \int \prod_{n=s}^{s+\kappa-1} dy_n Z_\kappa(s) e^{-\beta\mathcal{E}(y_n, y_{n-1})}, \end{aligned} \quad (2)$$

where the notation

$$Z_\kappa(s) = \int \prod_{n \neq s, \dots, s+\kappa-1}^N dy_n e^{-\beta\mathcal{E}(y_n, y_{n-1})}$$

has been introduced. $\beta = (k_B T)^{-1}$ is the Boltzmann factor. In order to evaluate the partition function (2) using the TIO method, we first symmetrize $e^{-\beta\mathcal{E}(x, y)}$ by introducing [10]

$$\begin{aligned} S(x, y) &= \exp\left(-\frac{\beta}{2}(V(x) + V(y) + 2W(x, y))\right) \\ &= S(y, x). \end{aligned}$$

Here the second equality holds only when x and y correspond to base-pairs of the same kind. Using Eq. (2) the expression for $Z_\kappa(s)$ is rewritten as

$$\begin{aligned} Z_\kappa(s) &= \int \left(\prod_{n \neq s, \dots, s+\kappa-1}^N dy_n S(y_n, y_{n-1}) \right) \\ &\quad \times dy_0 e^{-\frac{\beta}{2}V(y_0)} e^{-\frac{\beta}{2}V(y_N)}, \end{aligned} \quad (3)$$

where open boundary conditions at $n = 1$, and $n = N$ have been used. To proceed, a Fredholm integral equations with a real symmetric kernel

$$\int dy S(x, y) \phi(y) = \lambda \phi(x) \quad (4)$$

must be solved separately for the A/T and for the G/C base-pairs.

Since the eigenvalues are orthonormal and the eigenfunctions form a complete basis, Eq.(4) can be used sequentially to replace all integrals by matrix multiplications in Eq. (3). Whenever the sequence heterogeneity results in a non-symmetric $S(x, y)$, Eq.(4) cannot be used and we resort to a symmetrization technique, based on successive introduction of auxiliary integration variables, as explained in Ref. [13].

As noted, in order to quantify the sequence dependence on DNA's ability to form bubbles of different sizes, we have previously monitored the frequency of opening events using Langevin and Monte Carlo simulation techniques. Since the large openings constitute relatively rare events such techniques are not efficient (although essential for evaluating dynamical and non-equilibrium properties). It is much more effective to imply the probabilities of large bubbles at a given site in the sequence directly from the thermodynamic distributions using the TIO. Importantly, we have confirmed below that this equilibrium approach reproduces the bubble locations observed by Langevin simulations for the same sequences. This suggests that the bubbles – although large bubbles are rare events – are governed by equilibrium statistics.

We evaluate the probabilities $P_\kappa(s)$, for a base-pair opening spanning κ base-pairs (our operational definition of a bubble of size κ), starting at base-pair s as

$$P_\kappa(s) = \mathcal{Z}^{-1} \int_t^\infty \prod_{n=s}^{s+\kappa-1} dy_n Z_\kappa(s) e^{-\beta\mathcal{E}(y_n, y_{n-1})}, \quad (5)$$

where t is the separation (which we have taken as 1.5 \AA) of the double strand above which we define the strand to be melted.

Numerical Results. Using this technique, we are able to systematically investigate the relation between a given sequence containing a (apparently) disordered mixture of A/T and G/C pairs and the probability of spontaneous, thermally activated, bubbles of various sizes. Our analysis begins with a thorough study of two viral promoter sequences, Adeno major late promoter (AMLP) and Adeno Associate viral promoter (AAV P5). We have previously investigated the dependence of the thermally induced large bubbles in these sequences [2, 6] and found that the opening profiles obtained through Langevin simulations of the PBD model agreed remarkably well with the local denaturation profiles indicated by S1 nuclease experiments (see Ref. [2] for details). Here we use the TIO to calculate the probabilities Eq. (5) for the thermal creation of bubbles of size 1,3,7, and 10 base-pairs for the AdMLP promoter at $T = 300 \text{ K}$ (Fig. 1). The significant feature of the sequence is the occurrence of a TATA-box at base-pair location -30 with 7 consecutive A/T base-pairs. Around $+1$ there is rich region containing ~ 12 A/T base pairs, which, however, are not located consecutively since a comparable amount of G/C pairs are alternately embedded among the A/T pairs. Since

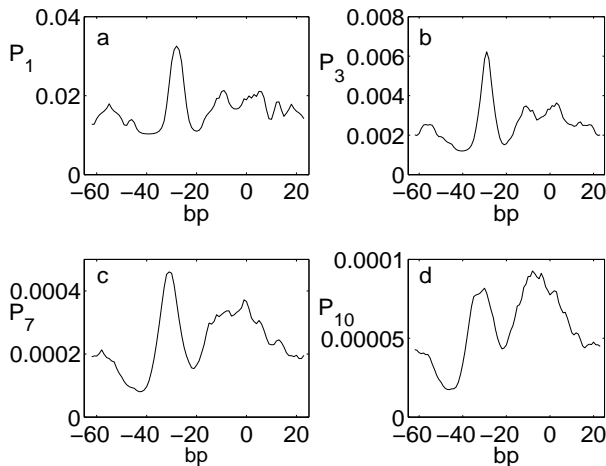


FIG. 1: Probabilities P_κ of creating bubbles spanning $\kappa = 1, 3, 7,$ and 10 -bp, respectively, for the Adeno major late promoter at $T=300\text{K}$.

A/T base-pairs are more weakly bound (softer) than G/C pairs, we could reasonably expect that bubbles have a predominant opening probability in the region -30 . This is indeed the case for small bubbles, Figs. 1a and 1b. However, surprisingly, this prediction breaks down when considering bubbles of larger sizes. The corresponding probability increases around $\text{bp} +1$ (Fig. 1c), up to the point that, for a bubble of size 10 bp, it becomes the highest (Fig. 1d). This finding illustrates the strong interplay between the sequence of base-pairs and the size of the bubble in the thermal activity of DNA (Indeed, it is likely that bubbles of different sizes may initiate different genetic processes).

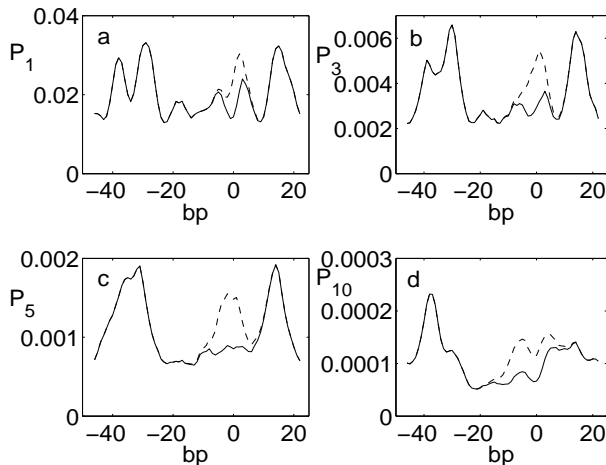


FIG. 2: Probabilities P_κ of creating bubbles of 1 bp ($\kappa = 1$) length (a), 3 bp (b), 5 bp (c) and 10 bp (d) length, for the wild (dashed line) and mutant (solid line) P5 promoter.

The replacement of 1 or 2 soft A/T with hard G/C base pairs in specific regions of the DNA can also hugely affect the probability for the formation of bubbles of given sizes.

This is illustrated with the AAV P5 promoter (Fig.2). This sequence regulates the AAV gene expression, and it has been shown [14] to bind the transcription initiator Yin Yang 1 (YY1) and to be active for TATA-Box protein (TBP)-independent transcription. The mutation of this promoter in which the two A/T bases at $+1$ and $+2$ are replaced by two G/C bases, is known to destroy the binding site for the YY1 initiator and thereby inhibit transcription. We have previously shown by Langevin simulations of the PBD model that this mutation also suppresses the formation of large bubbles around $\text{bp} +1$. Here we again calculate the probability to obtain bubbles of various sizes using the TIO. In Fig. 2 we show the probability of obtaining bubbles of sizes $n=1,3,5,$ and 10 for the wild (dashed line) and the mutated (solid line) AAV P5 promoter. The mutation causes a dramatic change in the double strand's ability to form large bubbles at and around the mutated region. However, the P_1 and the P_{10} probabilities are much less affected by the mutation. Notice for the wild type P5 AAV promoter, the region around the TATA-box has the largest probability for forming large bubbles (panel d). It is important to note that the AAV P5 promoter has four A/T rich regions: four consecutive A/T's around position -40 ; seven A/T's around position -30 ; five A/T's at the transcription start site $+1$; six A/T's around $+14$, and all these soft regions are clearly discernible in P_1 (panel a). From these results on the AAV P5 and AdMLP promot-

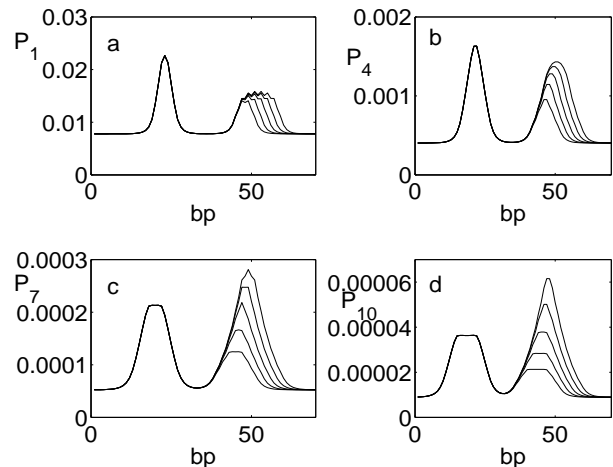


FIG. 3: Probability P_κ for the formation of bubbles of sizes $\kappa = 1, 4, 7, 10$ bps. The sequences are composed of 20 G/C, 5 A/T, 20 G/C followed by different sequences comprising $3,4,5,6,7$ A/T alternating with G/C bps. The last 20 bps are G/C.

ers we can speculate that the occurrence and intensity of a peak in the bubble probabilities does not depend on the specific composition of the DNA fragment. Rather, there is an essential interplay between the content of A/T

and G/C base pairs and the size of the bubble being examined. Understanding the mechanisms regulating the DNA openings are of great importance for predicting and engineering DNA processes, and we therefore now consider a series of simple (but experimentally realizable) DNA sequences where the effects discussed above are reproduced in detail. Our purpose is to isolate the underlying mechanisms. Our five sequences are all composed of 20 G/C, 5 A/T and 20 G/C base pairs. This is followed by a sequence that alternates A/T and G/C base pairs. We use 3, 4, 5, 6, and 7 A/T base-pairs in the five sequences. Finally, all five sequences are terminated with 20 G/C base-pairs.

As shown in Fig. 3a, the largest 1 bp opening probability is localized at +20, a region that contains five consecutive A/T bases, and is therefore expected to be more susceptible to open than the region localized around +50, containing A/T's alternating with G/C's. However, this simple picture changes dramatically as we move to larger bubbles, Figs. 3b, 3c, and 3d. In all these cases, the height of the second peak increases as compared to the peak at +20. With 3 and 4 A/T's, the peaks saturate at a value lower than the first peak. However, the height of the two peaks for the sequence with 5 A/T's, becomes equal in P_7 , and remains so for larger bubbles. The sequence with 6 A/T's shows an inversion of the opening probability, similar to that observed in the AdMLP sequence: at P_{10} the most probable 10 base-pair opening occurs around the base-pair location 50. These data indicate that the opening probability of a bubble of a given length does not trivially depend on the number of consecutive A/T's in the DNA sub-sequence. Instead, bubbles of sizes n form with higher probabilities in regions where the number of A/Ts over some characteristic length $L(\kappa)$ is higher, even if the A/Ts are mixed with G/C pairs.

To confirm this hypothesis we have extracted the characteristic lengths $L(\kappa)$ as a function of the bubble size n from the probability distributions of the simple sequences considered in Fig. 3. For instance, for $\kappa = 1$, Fig. 3a, we have obtained $L(1) = 4$ sites, while for $\kappa = 5$, Fig. 3b, we have $L(5) = 10$ sites. We have therefore considered the AAV P5 promoter DNA sequence of Fig. 2. Starting from each site s of the sequence, we count the number $N_\kappa(s)$ of A/T pairs contained over the corresponding next $L(\kappa)$ sites. In Fig. 4a we show the results for bubbles of size $\kappa = 1$, which can be compared with Fig. 2a. The small difference between the mutant and the wild opening probability for the sites around located at 0 is well reproduced. The difference between the wild and the mutant sequence is most pronounced for $\kappa = 5$, Fig. 4b, which is also in agreement with the TIO calculation shown in Fig. 2c.

We have also demonstrated that the opening probability does not depend on the specific distribution of the AT pairs contained in the characteristic regions of length $L(\kappa)$. This is shown in Fig. 5, where we have calculated

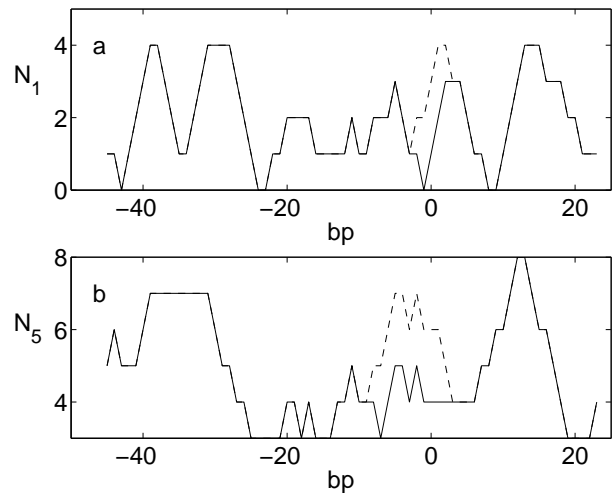


FIG. 4: Number of A/T bps contained in the characteristic length $L(\kappa)$, where $\kappa = 1$ (top panel), and 5 (bottom panel) for the wild (dashed line) and mutant (solid line) P5 promoter.

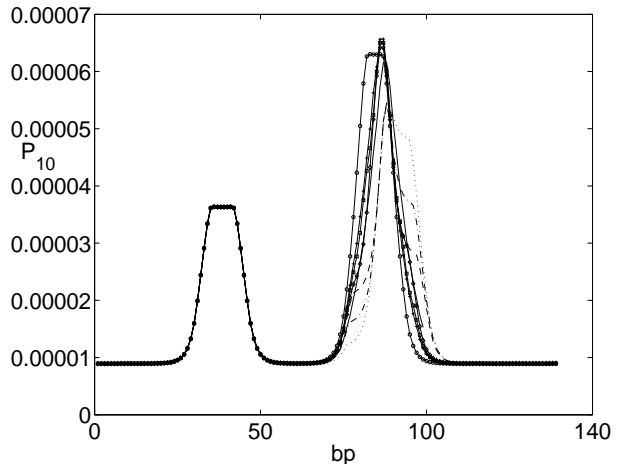


FIG. 5: Probability P_{10} for the formation of a bubble of 10 consecutive bps. The sequences consist of 40 G/C, 5 A/T and 20 G/C followed by 14 bps containing different random combinations of 7 A/T and 7 G/C. All sequences end with 40 G/C bps.

with the TIO the probability for formation of bubbles with size P_{10} for a sequence where the second A/T rich regions always contains 7 A/T distributed in several different combinations, but always over a 20-base region. We see that independently of the distribution of A/T base-pairs, the probability of 10 base-pair bubbles is always largest in the right-most region. Physically, we interpret this as the nonlinear coherence dominating (smoothing out the effect of) the base-pair disorder.

In summary, we have developed a semi-analytical techniques (TIO) that allows the efficient prediction of a given sequence for thermally induced bubbles of given sizes. We have found that large thermally induced bubbles arise

through a subtle interplay between length scales inherent in the nonlinear dynamics, and the sequence disorder. Our results provide new understandings that can help to not only identify new protein coding genes, but also enable reverse-engineering for use in future gene therapeutic applications.

This work at Los Alamos National Laboratory is supported by the US Department of Energy (contract No. W-7405-ENG-36) and by a NIH grant for A.U. (Grant Number: R01 GM073911).

-
- [1] M. Gueron, M. Kochoyan, and J.L. Leroy *Nature* **328**, 89 (1987); M. Frank-Kamenetskii *Nature* **328** 17 (1987).
 [2] C.H. Choi, G. Kalosakas, K.Ø. Rasmussen, M. Hiromura, A. R. Bishop and A. Usheva, *Nucleic Acids Res.* **32**, 1584, (2004).
 [3] T. Dauxois, M. Peyrard and A. R. Bishop, *Phys. Rev. E* **47** R44 (1993).
 [4] M. Peyrard and A. R. Bishop, *Phys. Rev. Lett.*, **62**, 2755, (1989).
 [5] M. Peyrard, *Nonlinearity* **17**, R1 (2004).
 [6] G. Kalosakas, K.Ø. Rasmussen, A. R. Bishop, C.H. Choi, and A. Usheva, *Europhys. Lett.* **68**, 127, (2004).
 [7] S. Ares, N.K. Voulgarakis, K.Ø. Rasmussen and A.R. Bishop, *Phys. Rev. Lett.*, **94**, 035504, (2005).
 [8] D. J. Scalapino, M. Sears and R.A. Ferrell, *Phys. Rev. B*, **6**, 3409, (1972).
 [9] Y. Zhang, W.-M. Zheng, J.-X. Liu, and Y. Z. Chen, *Phys. Rev. E*, **56**, 7100, (1997).
 [10] T. Dauxois and M. Peyrard, *Phys. Rev. E* **51** 4027 (1995).
 [11] A. Campa and A. Giansanti, *Phys. Rev. E*, **58**, 3585, (1998).
 [12] The parameters were chosen in Ref. [11] to fit thermodynamic properties of DNA: $k = 0.025eV/A^2$, $\rho = 2$, $\beta = 0.35A^{-1}$ for the inter-site coupling; for the Morse potential $D_{GC} = 0.075eV$, $a_{GC} = 6.9A^{-1}$ for a G-C bp, $D_{AT} = 0.05eV$, $a_{AT} = 4.2A^{-1}$ for the A-T bp.
 [13] M. B. Fogel, *Nonlinear Order Parameter Fields: I. Soliton Dynamics, II. Thermodynamics of a Model Impure System*, Ph.D. Thesis, Cornell University, (1977).
 [14] A. Usheva and Shenk, *Cell* **76**, 1115 (1994)