# RNA Denaturation: Excluded Volume, Pseudoknots, and Transition Scenarios

M. Baiesi,[1] E. Orlandini,[1] and A. L. Stella[1,2]

[1]INFM–Dipartimento di Fisica, Università di Padova, I-35131 Padova, Italy
[2]Sezione INFN, Università di Padova, I-35131 Padova, Italy

A lattice model of RNA denaturation which fully accounts for the excluded volume effects among nucleotides is proposed. A numerical study shows that interactions forming pseudoknots must be included in order to get a sharp continuous transition. Otherwise a smooth crossover occurs from the swollen linear polymer behavior to highly ramified, almost compact conformations with secondary structures. In the latter scenario, which is appropriate when these structures are much more stable than pseudoknot links, probability distributions for the lengths of both loops and main branches obey scaling with nonclassical exponents.

    

In recent years, considerable efforts have been devoted to the description of secondary structure formation (base pairing map) in single molecular strands of RNA [1–8]. This is an important step within the general program of understanding how structure is encoded in the primary sequence of biopolymers. By disregarding excluded volume effects and pseudoknots, recent studies established the existence of a molten phase at relatively high temperatures $T$ for a RNA molecule in dilute solution [9–11]. In this phase the inhomogeneities associated with a specific primary sequence should be irrelevant for the large scale behavior and should allow the coexistence of many different secondary structures of comparable free energies. As $T$ increases, a long RNA molecule should pass from the molten phase to a regime in which secondary structures essentially disappear and the global behavior becomes that of a linear polymer in good solvent. Excluded volume should play a relevant role at such denaturation transition. Indeed, there the entropic free energy gain associated with the formation of hairpins or of more complicated branched structures with loops is comparable with the corresponding base pair (bp) binding and staking energies, and depends crucially on the excluded volume interactions. Recent studies have shown that the discontinuous nature of DNA denaturation is determined by these interactions [12–14].

In this Letter we propose a model of the large scale conformational behavior of RNA in the high-$T$ and molten phases. Our coarse-grained description applies to mesoscopic scales and treats RNA as a long homopolymer, disregarding many of its structural details. These simplifications are not expected to alter the large scale properties and allow one to fully take into account excluded volume effects and their interplay with pseudoknots in determining the denaturation transition.

Starting with the related pioneering work of de Gennes [15] on the statistics of branchings and hairpin helices in a periodic copolymer, excluded volume effects were never fully included in studies of RNA denaturation. This leaves open the problem of establishing the existence and of determining the character of this transition in the long chain limit. A realistic embedding of the system in space, taking into account excluded volume, is also necessary for discussing adequately the complex conformational constraints posed by pseudoknots. Pseudoknots occur, e.g., when two loops locally bind to each other forming helices in their interior [8]. Several recent papers have discussed secondary structure formation taking pseudoknots into account for RNA molecules with a specific base sequence [16–19]. The complexity of descriptions which are realistic up to the nucleotide scale did not allow these studies to properly account for excluded volume effects. On the other hand, they demonstrated the key importance of pseudoknots in determining secondary structures and folding pathways.

We model a conformation of the RNA strand as a two-tolerant trail of $N$ steps on the face centered cubic (fcc) lattice [20]. This is a random walk in which no more than two steps are allowed to overlap on a single lattice bond, forming what we call a contact. In addition, by giving an orientation to the trail, we impose that only pairs of antiparallel steps can form contacts, and whenever this happens a gain in energy $\epsilon < 0$ is counted. The orientation of the trail reflects the backbone directionality or RNA [8]. Because of its coarse-grained character, it is possible to establish only a rather rough correspondence between the geometry of our model and many microscopic structural features of RNA [2,16–19]. A single step on the lattice corresponds as order of magnitude to the Kuhn length of single stranded (ss) RNA ($\sim 2.5$ bases [18]). Helices correspond to sequences of double lattice steps, and a minimal helix of $\sim 3$ bp [18] can be assumed to be represented by a single double step. Consistent with the coarse-grained character and with the homopolymer approximation, $\epsilon$ is an effective parameter, averaging among several different energetic and entropic contributions at the smallest scales. To start with, a sequence of double steps has in our model the same flexibility as a

sequence of single steps. This is an approximation, because the persistence length of double stranded (ds) RNA should be much larger than that of ssRNA and comparable to that of dsDNA (60–150 bp [21]). On the other hand, the length of double step sequences should be limited by base matching requirements. The difference between ss and ds rigidity can be taken into account by introducing a suitable bending energy for double step sequences. In the case of DNA denaturation models such bending rigidity did not reveal the ability to alter the asymptotic transition properties [14,22].

Figure 1 sketches two possible configurations of our model. In both 1(a) and 1(b) a diagram on the right summarizes the corresponding contact map. A bridge in the diagrams connects each pair of steps forming a contact. Bridges are numbered in the order of appearance if one follows the trail orientation. A main bridge is a bridge which is not inscribed within other, larger bridges. Unlike 1(b), Fig. 1(a) shows a pseudoknot, indicated by the crossing of two bridges in the diagram. This crossing means that a step forming a loop overlaps with one outside the loop. We consider two variants of the model, which we refer to as (model) I and II. While in I configurations with pseudoknots are allowed, in addition to those without pseudoknots, in II the former are forbidden altogether.

Thermodynamic quantities and canonical averages are defined in terms of the partition function $Z = \sum_w \exp[-H(w)/T]$. The sum extends to all allowed configurations $w$ with $|w| = N$ steps, and $H = \epsilon N_c(w)$, $N_c(w)$ being the number of contacts in $w$. For both models, we sampled configurations by a multiple Markov chain Monte Carlo procedure [23] using several ($\approx 20$) temperatures satisfying $0 \leq \epsilon/T \leq 3.5$ [24]. We first

computed as a function of $T$ the specific heat of I and II for different $N$. At a continuous conformational transition with crossover exponent $\phi < 1$ one expects a singular behavior $C_{\max} \sim N^{2\phi-1}$ for the maximum of this quantity as $N \to \infty$ [26]. If $\phi < 1/2$ such singularity does not imply a divergence. For both models we find no evidence of a diverging $C_{\max}$. Hence, at this level we can only conclude that for both models the denaturation transition must be continuous and with $\phi < 1/2$, if it exists.

We also determined two geometrical radii of the configurations, namely, the end-to-end distance, $R_e$, and the radius of gyration with respect to the center of mass, $R_g$. The ratios of the averages of such radii in the $N \to \infty$ limit are expected to be universal numbers characteristic of the different regimes involved in the transition [26,27]. Figure 2 shows plots of $\langle R_e^2 \rangle / \langle R_g^2 \rangle$ for different $N$. For I the trend of the curves gives indication of a sharp transition at $\epsilon/T \approx 1.9$. Indeed, for high $T$ the ratio approaches from below the universal value 6.25(1) appropriate for the swollen, self-avoiding walk (SAW) regime [28]. On the other hand, at very low $T$'s the trail should fold in a double structure with a maximal number of contacts ($N_c \sim N/2$), in which $R_e$ necessarily approaches zero. This explains the trend towards zero (from above) of the curves at low $T$. Remarkable is the accumulation of intersection points for $\epsilon/T \approx 1.9$. These intersections mark a change of the trend of the curves for increasing $N$ and suggest the presence of a peculiar transition regime with universal ratio $\approx 4.8$. Hence, for I there is clear evidence of a second order transition at $\epsilon/T \approx 1.9$. For II there is no similar indication: the intersections are pushed towards lower and lower $T$'s as one compares curves corresponding to pairs of increasing $N$ values (Fig. 2, inset). This means that the larger $N$, the more the SAW regime extends in the low $T$ region. The whole pattern suggests for II a rather slow crossover, not a transition. Further insight is provided by the study of
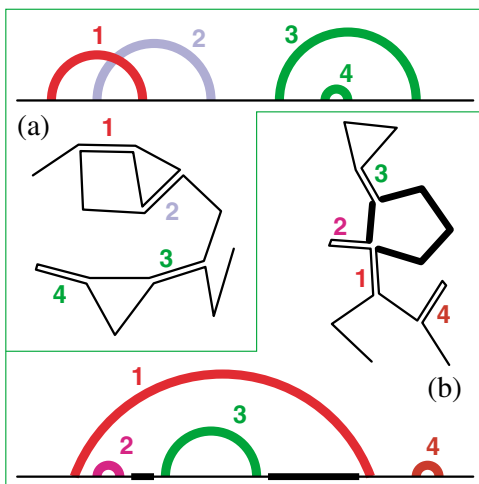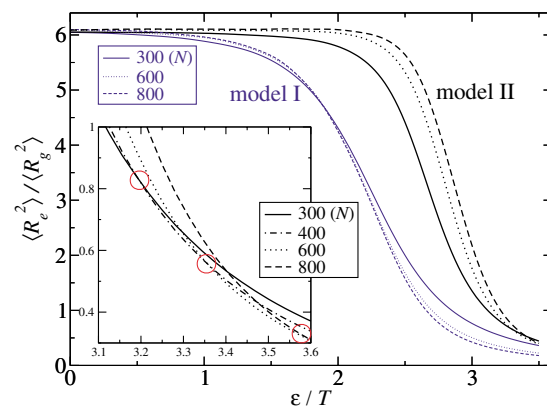


FIG. 1 (color online). RNA configurations and corresponding contact maps. Overlapped steps (contacts) are slightly split. In (a) a pseudoknot is present (crossing of bridge "1" with bridge "2"). In (b) a loop of length $\ell = 5$ is marked by a thicker line, both on the chain and in the contact map. Here, "1" and "4" are main bridges.



FIG. 2 (color online). $\langle R_e \rangle^2 / \langle R_g \rangle^2$ as a function of $\epsilon/T$ for three different values of $N$. Inset: detail of the crossings of four curves for model II. The circles enclose intersections between the curve pairs (300, 400), (400, 600), and (600, 800).
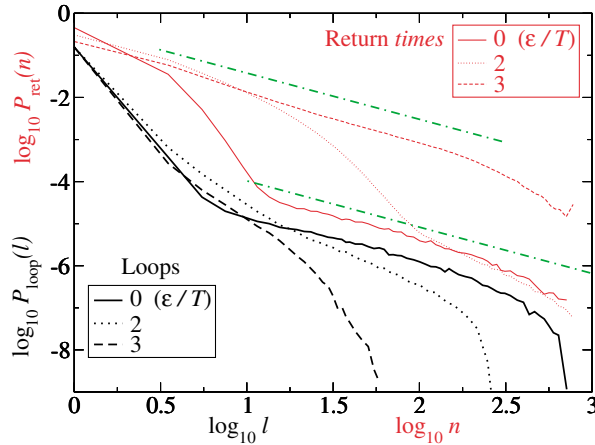
FIG. 3 (color online). Log-log plots for $P_{\text{loop}}(\ell)$ (thick lines, shifted down by 0.5 for clarity) and $P_{\text{ret}}(n)$ (thin lines), both for $N = 800$ and for different $T$ values. The dot-dashed lines have slope $-1.1$.

some scaling properties. The radius of gyration is expected to scale as $\langle R_g \rangle \sim N^\nu$ for large $N$ [26]. For both models we observe that at high $T$ the determinations of $\nu$ at finite $N$, for $N \rightarrow \infty$, approach a value $\approx 0.59$ appropriate for a SAW in $d = 3$ [26,28]. For I, at $T$'s sufficiently below the transition the $\nu$ estimates can be extrapolated to $\approx 0.35$ for large $N$. This indicates that the configurations are very close to compact in the low $T$, molten phase ($\nu \approx 1/d = 1/3$). For II at very low $T$'s we extrapolate $\nu \approx 0.4$ which is also not far from $\nu = 1/2$, as appropriate for branched polymers [26,29].

The different behaviors of the two models are due to the presence of pseudoknots in I, where they occur in almost all sampled configurations at sufficiently low $T$. Pseudoknots correspond to the formation of extra binding contacts and thus can lead to more compact configurations with respect to the case of II. These extra contacts trigger the sharp transition observed at $\epsilon/T \approx 1.9$. For II, if present, a transition should be located at much lower $T$'s, most likely below the range of applicability of the model. The possibility of forming pseudoknots is a driving factor in tertiary structure formation [4]. However, I somehow overamplifies this factor. Indeed, when pseudoknots form in real RNA, the high rigidity of the ds contact region between loops should stretch them, lowering their conformational entropy. This results in a destabilization of the links forming pseudoknots. One way to increase the binding free energies of these links is to introduce sufficiently high concentrations of divalent metal ions, like $Mg^{2+}$ in solution [4,31]. On the other hand, in II pseudoknot forming links would appear if one would look at the details of the configurations at a somewhat more coarse-grained level. This means that II gives essentially zero energy to such links. Thus, one could expect that the description of II is probably closer to real RNA in situations in which a low concentration of

$Mg^{2+}$ ions in solution enlarges the stability gap between secondary structure and pseudoknot forming links [4,31].

We tested the possible effect of ds vs ss rigidity by introducing in our Hamiltonian a term $\epsilon_b N_b(w)$, where $\epsilon_b > 0$ is a bending energy and $N_b(w)$ is the total number of pairs of successive double steps which are not parallel. For $\epsilon_b = -4\epsilon$, we verified that the average number of junctions with bending double steps remains always almost equal to zero, realizing a limit situation of infinitely rigid double helices. This limit is a reasonable approximation in view of the relatively large ds persistence length. With such rigidity included, we can reach less asymptotic values of $N$ ($N \lesssim 500$) in our simulation. The transition pattern of I and the differences between I and II remain essentially unaltered. The only appreciable change is a slight shift of the transition in I toward lower $T$'s, which is revealed by the positions of the specific heat maxima ($\epsilon/T$ roughly increases by 10%). This shift is consistent with the expected entropic destabilization of pseudoknot links due to ds rigidity. As far as the universal features and the trends observed in Fig. 2 are concerned, we cannot detect significant changes. Also the low $T$ scalings remain essentially unaltered.

RNA denaturation corresponds to a substantial suppression with increasing $T$ of the highly ramified structure of loops and branches characterizing the molten phase. The analysis of the loops is feasible and particularly instructive in II. For DNA, the distribution of the lengths of denaturated loops, corresponding to openings of the double helix, follows a power law whose exponent $c$ determines the character of the transition [13,14]. A simple example of a loop in RNA is given by the closure of an isolated hairpin. In this case the loop is connected to the rest of the structure by a single branch of double steps. Of course, more complicated situations may occur [thick loop in Fig. 1(b)]. Even at high $T$ an extensive number of minor spikelike branches is present along the RNA backbone. Thus, loops with a fixed number of branches naturally have length distributions with rather sharp cutoff. Therefore, we decided to sample the length of all loops identified in the various configurations, irrespective of the number of outgoing double step branches. The various lengths can be obtained from the contact map of each configuration, by using a recursive algorithm that identifies all the loops inside each main bridge in the diagram. Another interesting quantity is the return *time*, i.e., the total number of steps comprised within a main bridge. In the assumed planar topology of model II this time is the arc length corresponding to each departure of the configuration from a contact-free, linear polymer behavior. Probability distributions of the loop lengths, $\ell$, and of the return times, $n$, are plotted in Fig. 3 for different $T$'s and for $N = 800$. At high $T$, after transients both distributions behave as power laws with approximately identical exponents: $P_{\text{loop}}(\ell) \sim \ell^{-c_\ell}$, $P_{\text{ret}}(n) \sim n^{-c_r}$, with $c_\ell \simeq c_r = 1.1(1)$. The peaks at small $\ell$ and $n$ in the distributions

indicate that loops at high $T$ mostly occur within isolated small hairpins in II. The identity of $c_\ell$ and $c_r$ means that almost all large bridges are also main bridges. Thus, the return time essentially coincides with the loop length at high $T$. At lower $T$, while $P_{loop}$ becomes shorter and shorter ranged for decreasing $T$, the behavior of $P_{ret}$ remains of power law type at large arguments. The value of $c_r$ remains stable and close to that estimated for $\epsilon/T = 0$. This means that as the RNA molecule enters deeper and deeper into the molten phase with developed secondary structures, the loops become shorter and shorter. On the other hand, the main branches departing from the contact-free backbone encompass all accessible length scales, as appropriate for a branched polymer. This could also explain why in this range of $T$ the exponent $\nu$ discussed above is not far from $1/2$, as for branched polymers [26]. The exponent $c_r$ obtained here definitely deviates from the mean field value $3/2$ [9]. The introduction of ds rigidity in II does not alter to a significant extent the scenario discussed above and, in particular, the exponents $c_r$ and $c_\ell$.

To summarize, for I we could establish the existence of a sharp denaturation transition which is due to the presence of pseudoknots. Unlike DNA denaturation, this is a second order transition. The ds rigidity does not alter its universal features. II is a more adequate description of RNA in situations with a very wide stability gap between secondary and tertiary folding levels [4,31]. In this model there is no sharp transition and denaturation occurs as a crossover from linear to branched-compact polymer behavior. The geometry of this crossover is well described by the distributions $P_{loop}$ and $P_{ret}$ and by their exponents, whose nonclassical values are a further consequence of excluded volume.

[1] M. Zuker, Science **244**, 48 (1989).
[2] I. L. Hofacker *et al.*, Monatsh. Chem. **125**, 167 (1994), URL: http://www.tbi.univie.ac.at/~ivo/RNA/.
[3] P. G. Higgs, Phys. Rev. Lett. **76**, 704 (1996).
[4] I. Tinoco, Jr. and C. Bustamante, J. Mol. Biol. **293**, 271 (1999).
[5] R. Bundschuh and T. Hwa, Phys. Rev. Lett. **83**, 1479 (1999).
[6] A. Pagnani, G. Parisi, and F. Ricci-Tersenghi, Phys. Rev. Lett. **84**, 2026 (2000).
[7] S.-J. Chen and K. A. Dill, Proc. Natl. Acad. Sci. U.S.A. **97**, 646 (2000).
[8] P. G. Higgs, Q. Rev. Biophys. **33**, 199 (2000).
[9] R. Bundschuh and T. Hwa, Phys. Rev. E **65**, 031903 (2002).
[10] R. Bundschuh and T. Hwa, Europhys. Lett. **59**, 903 (2002).
[11] M. Müller, Phys. Rev. E **67**, 021914 (2003).
[12] M. S. Causo, B. Coluzzi, and P. Grassberger, Phys. Rev. E **62**, 3958 (2000).
[13] Y. Kafri, D. Mukamel, and L. Peliti, Phys. Rev. Lett. **85**, 4988 (2000).
[14] E. Carlon, E. Orlandini, and A. L. Stella, Phys. Rev. Lett. **88**, 198101 (2002).
[15] P.-G. de Gennes, Biopolymers **6**, 715 (1968).
[16] E. Rivas and S. R. Eddy, J. Mol. Biol. **285**, 2053 (1999).
[17] E. Rivas and S. R. Eddy, Bioinformatics **16**, 334 (2000).
[18] H. Isambert and E. D. Siggia, Proc. Natl. Acad. Sci. U.S.A. **97**, 6515 (2000).
[19] R. B. Lyngso and C. N. S. Pedersen, J. Comput. Biol. **7**, 409 (2000).
[20] The fcc lattice allows for closed loops also of odd length and thus enriches their sampling.
[21] S. B. Smith, Y. J. Cui, and C. Bustamante, Science **271**, 795 (1996).
[22] Y. Kafri, D. Mukamel, and L. Peliti, Phys. Rev. Lett. **90**, 159802 (2003).
[23] M. C. Tesi, E. J. Janse van Rensburg, E. Orlandini, and S. G. Whittington, J. Stat. Phys. **82**, 155 (1996).
[24] To increase the mobility of the Markov chain at low $T$, where branched structures are expected, in addition to the usual set of local and pivot moves (see [23] and references therein), we implemented a variant of the pivot move, in which only one arm of the branched structure is rotated. In order to obtain smooth plots, data were processed by means of the multiple histogram method [25].
[25] A. M. Ferrenberg and R. H. Swendsen, Phys. Rev. Lett. **61**, 2635 (1988).
[26] C. Vanderzande, *Lattice Models of Polymers* (Cambridge University Press, Cambridge, U.K., 1998).
[27] V. Privman, P. C. Hohenberg, and A. Aharony, in *Phase Transitions and Critical Phenomena*, edited by C. Domb and J. L. Lebowitz (Academic Press, New York, 1991), Vol. 14, p. 1.
[28] B. Li, N. Madras, and A. D. Sokal, J. Stat. Phys. **80**, 661 (1995).
[29] Two tolerant trails have been used to model collapse from SAW to branched polymer behavior [30]. See also P. Leoni and C. Vanderzande, cond-mat/0303421.
[30] E. Orlandini, F. Seno, A. L. Stella, and C. Tesi, Phys. Rev. Lett. **68**, 488 (1992).
[31] V. K. Misra and D. E. Draper, Biopolymers **48**, 113 (1998).