



# Conformational Ensembles of Noncoding Elements in the SARS-CoV-2 Genome from Molecular Dynamics Simulations

Sandro Bottaro,\* Giovanni Bussi, and Kresten Lindorff-Larsen\*



Cite This: *J. Am. Chem. Soc.* 2021, 143, 8333–8343



Read Online

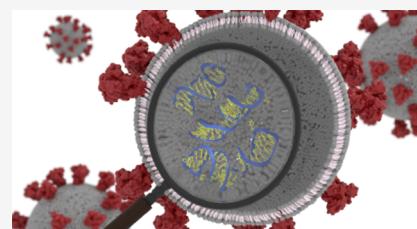
ACCESS |

Metrics & More

Article Recommendations

Supporting Information

**ABSTRACT:** The 5' untranslated region (UTR) of the severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) genome is a conserved, functional and structured genomic region consisting of several RNA stem-loop elements. While the secondary structure of such elements has been determined experimentally, their three-dimensional structures are not known yet. Here, we predict structure and dynamics of five RNA stem loops in the 5'-UTR of SARS-CoV-2 by extensive atomistic molecular dynamics simulations, more than 0.5 ms of aggregate simulation time, in combination with enhanced sampling techniques. We compare simulations with available experimental data, describe the resulting conformational ensembles, and identify the presence of specific structural rearrangements in apical and internal loops that may be functionally relevant. Our atomic-detailed structural predictions reveal a rich dynamics in these RNA molecules, could help the experimental characterization of these systems, and provide putative three-dimensional models for structure-based drug design studies.



## INTRODUCTION

The genome of the SARS-CoV-2 virus is a single, positive-stranded RNA consisting of approximately 30 thousand nucleotides. Coding regions of the genome are translated into several nonstructural proteins, as well as into the widely studied spike protein, envelope protein, membrane protein, and others.<sup>1</sup> The 5' and 3' untranslated regions (UTR) contain cis-acting sequences required for viral transcription and replication.<sup>2,3</sup> In particular, the 5'-UTR harbors conserved, functional elements that enhance viral transcription<sup>4–6</sup> and are involved in discontinuous transcription, leading to leader-body fusion.<sup>7</sup> Recent studies suggest that the 5'-UTR plays a key role in liquid–liquid phase separation phenomena with the SARS-CoV-2 nucleocapsid protein.<sup>8,9</sup> Proof-of-principle studies using locked nucleic acids targeting conserved structural motifs in the SARS-CoV-2 genome showed the potential for inhibiting growth via RNA-interacting molecules.<sup>10</sup>

Chemical probing experiments, NMR chemical shift measurements, and computational predictions have shown that the 5'-UTR is highly structured<sup>10–15</sup> and consists of several stem-loop elements (Figure 1). The secondary structures of the isolated elements are in close agreement with the full-length construct,<sup>11</sup> suggesting that they can fold independently. The three-dimensional structure of stem loop 2 (SL2) from SARS-CoV-1, which is completely conserved in SARS-CoV-2, has been determined by NMR spectroscopy,<sup>16</sup> while no experimental structures exist for the remaining elements in the 5'-UTR. In a recent computational study,<sup>17</sup> Rangan et al. predicted the three-dimensional fold of all stem-loop elements in the 5'-UTR and 3'-UTR and of the frameshifting stimulatory element using Rosetta's FARFAR2 algorithm.<sup>18,19</sup> The predicted model of SL2 in the 5'-UTR was

later refined using MD simulations,<sup>20</sup> leading to an improved agreement with available NMR data.<sup>16</sup>

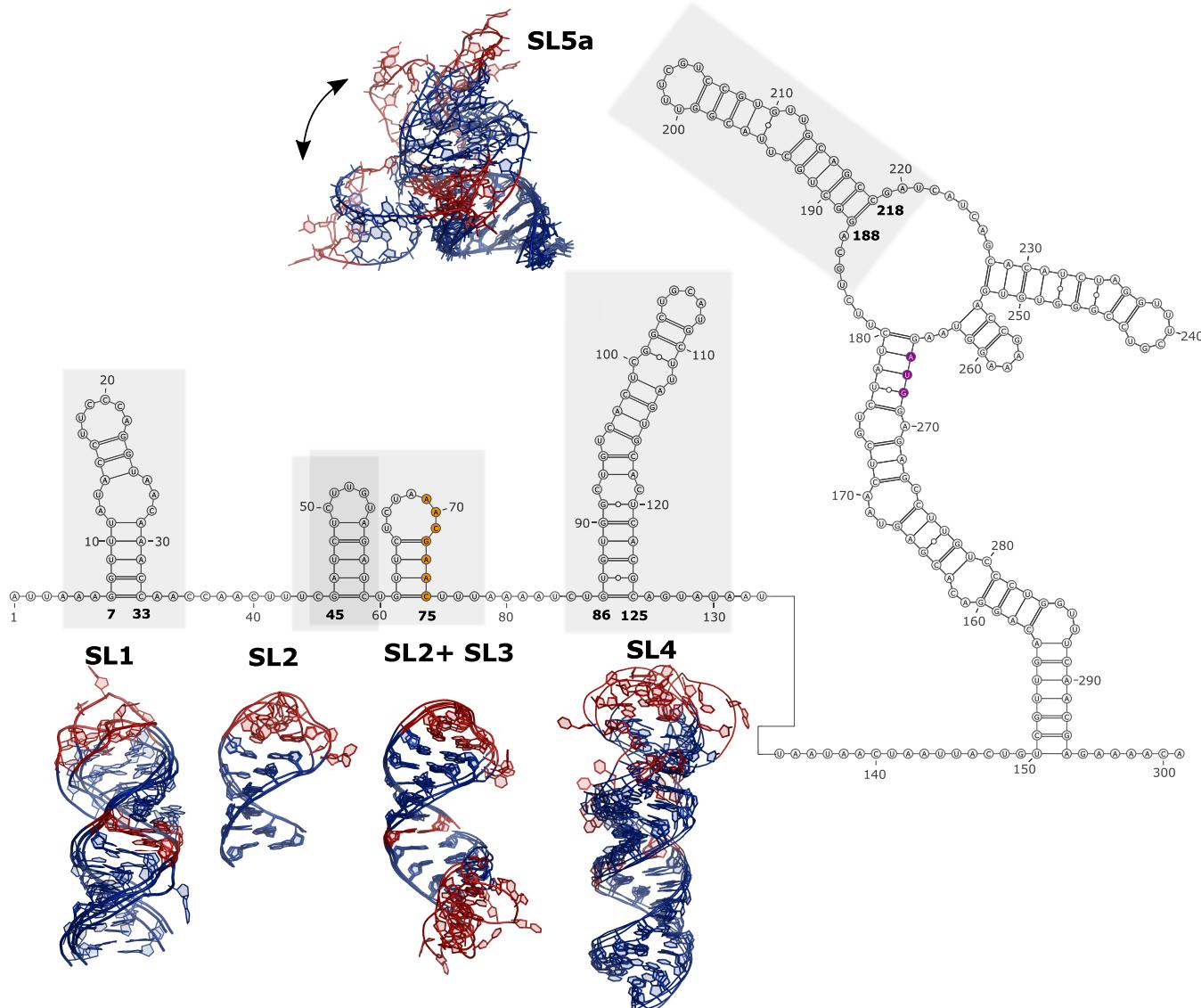
While there have been a plethora of structural studies of SARS-CoV-2 proteins, much less is known about the structure of the RNA genome at atomic resolution. RNA molecules often display a rich, complex dynamical behavior that is difficult to capture using experiments alone and that can play crucial roles in cellular functions.<sup>21,22</sup> In this work we have thus used atomistic molecular dynamics simulations to predict the structure and dynamics of five 5'-UTR stem loops. These elements were chosen because their high degree of conservation in betacoronaviruses, together with their vital role in viral replication, make them attractive, potential targets for small molecules. Experimental studies on the same constructs are currently being pursued within the COVID-19 NMR initiative (<https://covid19-nmr.de/>). The results presented in this work therefore represent a step in using computation and experiments in a synergistic manner. At the same time, the elements that we have chosen to model represent the largest structural elements for which it is currently feasible to obtain relatively converged simulations even with multi-microsecond enhanced sampling simulations.

In this spirit, we have generated three-dimensional configurations compatible with the secondary structure

Received: January 28, 2021

Published: May 27, 2021



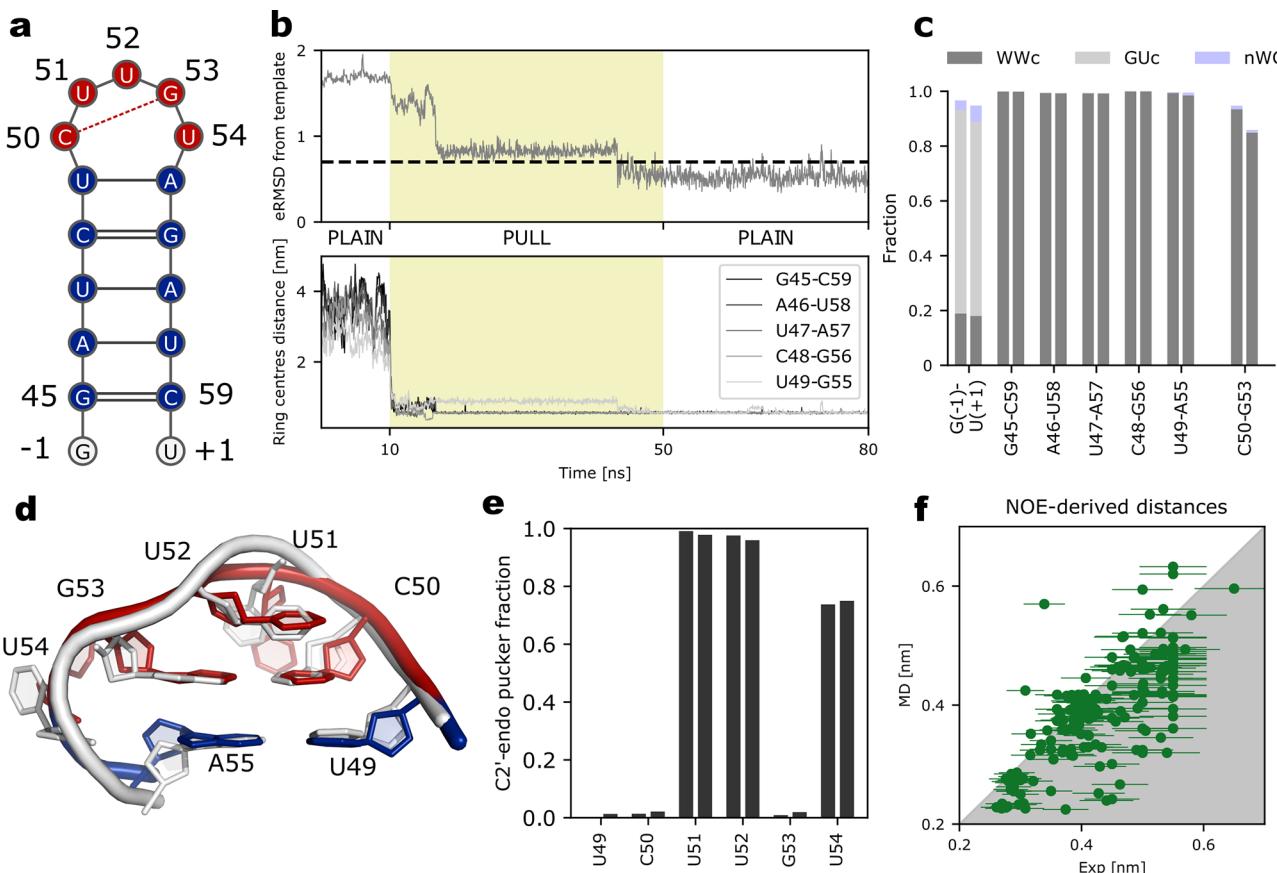


**Figure 1.** Reference sequence and secondary structure of the 5'-UTR in the SARS-CoV-2 genome. The numbering follows the convention in ref 11. The five structured elements considered in this work are labeled and highlighted in gray: SL1, SL2, SL2+SL3, SL4, and SL5a. The transcription regulatory sequence and the AUG start codon are highlighted in orange and purple, respectively. Selected centroids of the most populated clusters from MD simulations corresponding to each element are shown.

recently determined by chemical probing and NMR experiments.<sup>11</sup> While several fragment-based tools exist to perform such a task,<sup>23,24</sup> we here adapt an MD-based procedure.<sup>25</sup> Similar secondary structure restraining approaches have been employed in all-atom RNA folding simulations<sup>26,27</sup> and to study internal loop dynamics.<sup>28</sup> These initial structures are subject to extensive MD sampling, amounting to more than 0.5 ms of aggregate simulation time. More precisely, we enhance the sampling of apical and internal loop regions using a Hamiltonian replica-exchange scheme in which only specific portions of the RNA solute are affected.<sup>29–31</sup> For such regions we currently lack detailed structural information, and available experimental data suggest that conformational rearrangements are more likely to occur. Our simulations are thus a Boltzmann-weighted ensemble of three-dimensional conformations constituting our prediction of structure and conformational heterogeneity of these elements.

First, we apply the computational protocol to SL2 and validate the accuracy of our prediction against the three-dimensional structure of SL2 from SARS-CoV-1 and the NMR data used to generate this structure.<sup>16</sup> We then proceed by describing our predictions on four elements of increasing complexity (Figure 1) and discuss the agreement with available NMR measurements.<sup>32,33</sup> Stem loop 5a (SL5a) is part of the four-way junction SL5 containing the AUG starting codon, whose function is linked to RNA packaging.<sup>34</sup> Stem loop 1 (SL1) has been suggested to play a role in escaping the viral mechanism inhibiting mRNA translation.<sup>6,35</sup> The function of stem loop 4 (SL4) is still not completely understood and may act as a spacing element.<sup>36</sup> Finally, the SL2+SL3 construct consists of SL2 and stem loop 3 (SL3). The latter element contains the transcription regulatory sequence necessary for discontinuous RNA synthesis.<sup>37</sup>

Our predictions could help the atomic-detailed experimental characterization of such systems, as simulations can be easily



**Figure 2.** Description of the computational procedure and validation on SL2. (a) Sequence and secondary structure of SL2. The genomic numbering follows the convention in ref 11. The region used for constructing the template A-form three-dimensional structure is shown in blue. The sampling in the region shown in red is enhanced by partial tempering. The Watson–Crick base-pair between C50 and G53 predicted in partial tempering simulations is shown as a dashed red line. (b) Behavior of a successful pulling simulation. eRMSD (top panel) and distance between ring centers of the five WC base-pairs (bottom panel) during the three stages of the procedure: initial plain MD, pulling stage, and final plain MD stage, as labeled. The eRMSD threshold of 0.75 is shown as a dashed horizontal line. (c) Population of base-pairs from 20 replicas  $\times$  2.6  $\mu$ s partial tempering simulations. Following the Leontis–Westhof classification,<sup>46</sup> base-pairs are annotated as cis Watson–Crick/Watson–Crick (WWc), GU wobble cis (GUC), or non-Watson–Crick (nWC). For each interaction, the two bars show the statistics from two independent simulations. (d) Centroid of the most-populated cluster from MD simulations (blue/red) superposed onto the SL2 loop from the SARS-CoV-1 structure (white).<sup>16</sup> (e) Sugar pucker population in the loop region. (f) Calculated and experimental NOE-derived upper bound distances. The 10 of the 179 calculated NOEs that are significantly larger than the experimental values are located in the white upper-triangular region.

integrated *a posteriori* with NMR or SAXS data.<sup>38–40</sup> The predicted models provide a starting point for structure-based drug design studies<sup>41,42</sup> and could serve as additional benchmark systems to evaluate the accuracy of the atomistic RNA force field.<sup>43</sup> To maximize reproducibility and enable others to build on our work,<sup>44</sup> we have made trajectories, input files, and analysis scripts available via github at [https://github.com/KULL-Centre/papers/tree/master/2020/COVID\\_SUTR\\_MD](https://github.com/KULL-Centre/papers/tree/master/2020/COVID_SUTR_MD).

## RESULTS AND DISCUSSION

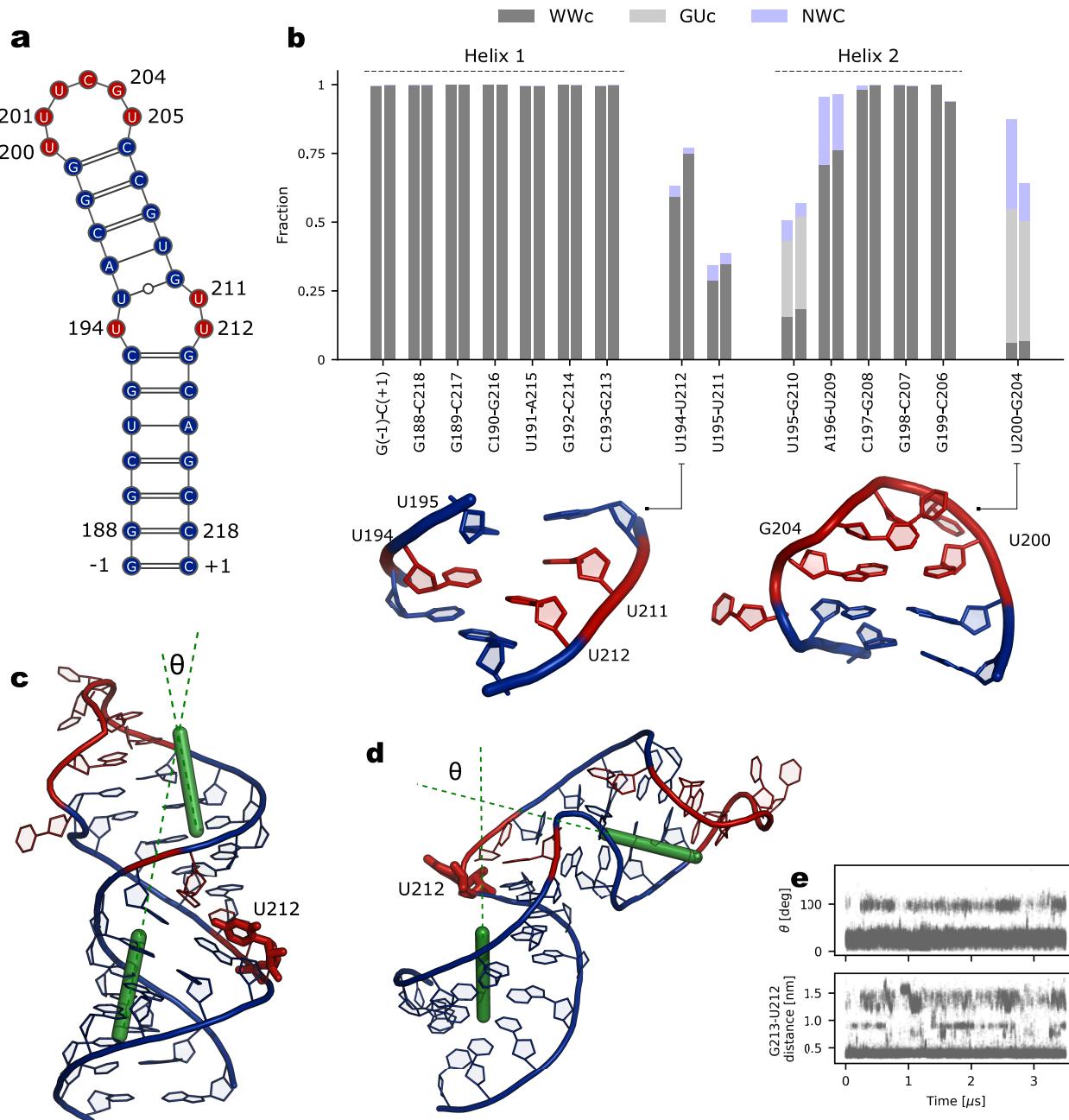
In this section we describe the conformational ensembles obtained from MD simulations. We first show the results obtained on SL2, the simplest system considered in this study. Crucially, in this case it is possible to validate the procedure by comparing with available structural data and NMR measurements.

**Stem Loop 2.** We run unbiased MD simulations of a single-stranded, extended RNA with sequence GGAUCUCUUGUAGAUCU. After 10 ns, we apply a biasing force (in a process we term “pulling”) to promote the

formation of an ideal A-form RNA structure in the region where the secondary structure is known (blue nucleotides in Figure 2a). By construction the structural dissimilarity from the template, here evaluated using a nucleic-acids-specific distance called ellipsoidal root-mean-square distance (eRMSD),<sup>45</sup> decreases during pulling (Figure 2b, top panel). After 40 ns the external bias is removed, and the system fluctuates around a stable structure in which all five Watson–Crick (WC) base-pairs are formed (Figure 2b, bottom panel).

While pulling simulations make it possible to quickly generate conformations with the desired stem structure, these runs are not sufficiently long to allow rearrangements in the loop region. Previous simulation studies have shown these motions to occur on long time scales ( $\mu$ s to ms) that require extensive simulations.<sup>47</sup> We thus enhance the sampling in the loop region (red in Figure 2a) by running extensive partial tempering simulations (20 replicas, each 2.6  $\mu$ s long) initialized from a successful pulling run (see Methods).

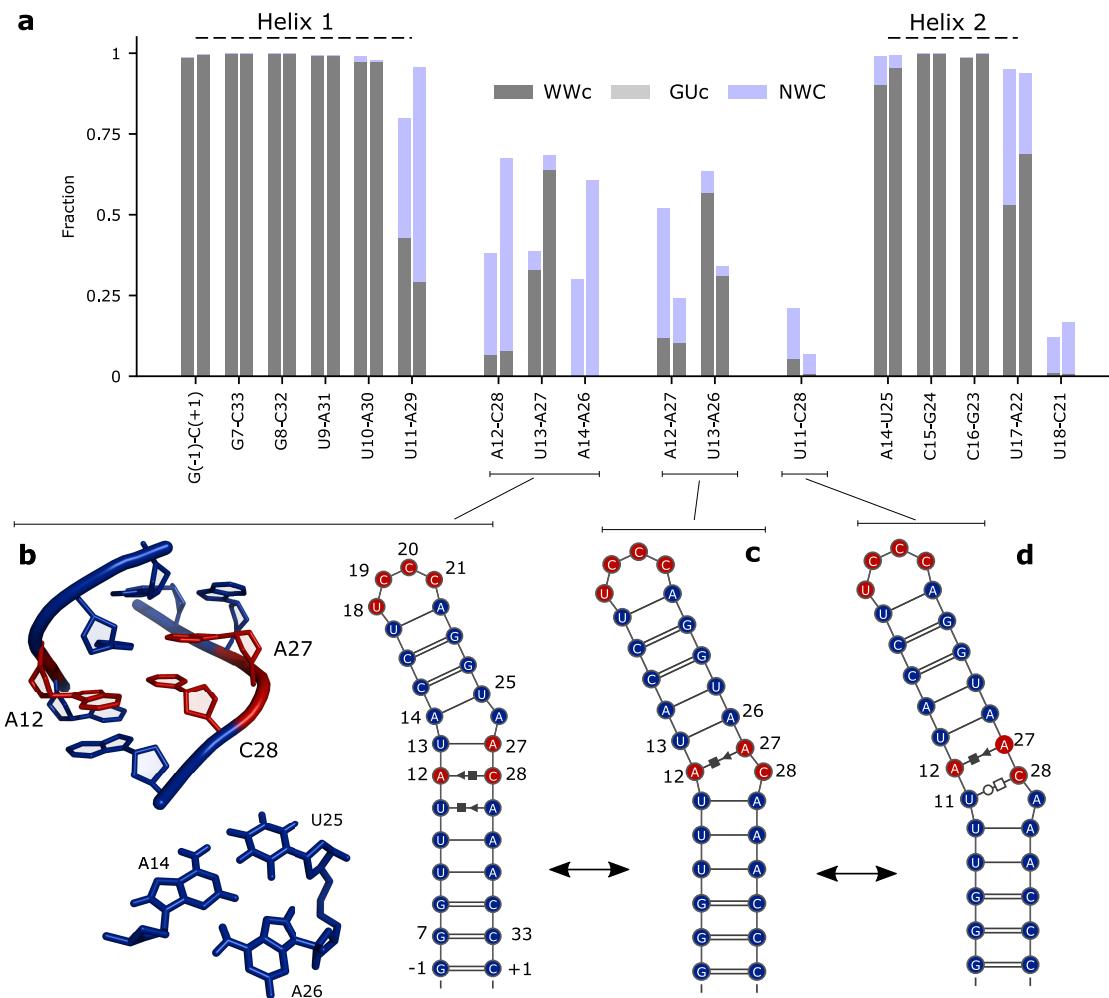
Analyses of the enhanced simulations show that all five base-pairs in the stem are stable, and a comparison of two independent simulations shows that our simulations are well



**Figure 3.** SL5a results. (a) Sequence and secondary structure of SL5a. The genomic numbering follows the convention in ref 11. The region used for constructing the A-form templates is shown in blue. The sampling in the region shown in red is enhanced by partial tempering. (b) Population of base-pairs from 20 replicas  $\times$  3.6  $\mu$ s partial tempering simulations. Centroids of the most populated internal and apical loop structures discussed in the text are shown. (c) Representative straight conformation. The axes defined by helix 1 and 2, together with the angle  $\theta$  between them, are shown in green. The helical axes are calculated as the vector normal to the base-pair plane, averaged over all base-pairs within the helix. (d) Representative kinked conformation. (e) Time series of the angle between helices,  $\theta$  (top panel), and of the distance between U212 and G213 (bottom panel) in the neutral replica. We note that due to the replica-exchange setup used in our simulations, these simulations do not reflect the time scales needed to switch between the two conformations.

converged (Figure 2c). In addition, we observe the formation of a stable terminal GU wobble base-pair, as well as of a C50–G53 Watson–Crick interaction in the loop, resulting in a structure resembling a CUUG tetraloop<sup>48,49</sup> with a bulged U at position 54. The interaction between C50 and G53 is also compatible with the decreased sensitivity to chemical probing of these two bases relative to the remaining three in the loop.<sup>10</sup> A three-dimensional cluster analysis<sup>50</sup> reveals that the largest cluster ( $\sim 90\%$ ) is remarkably similar to the available NMR

loop structure from SARS-CoV-1,<sup>16</sup> whose sequence differs from our construct in the two terminal base-pairs. For example, the centroid of the most populated cluster from MD simulations superimposes extremely well to the first NMR model in PDB 2L6I (Figure 2d). In the loop, C50 forms a canonical base-pair with G53, U52 stacks on top of C50, and both U51/U54 are solvent exposed. The median distance between the experimental structure and the MD simulation is



**Figure 4.** SL1 results. (a) Population of base-pairs from partial tempering simulations. (b) Three-dimensional structure of the most populated internal loop configuration showing the U11–A29 Hogsteen–sugar and A12–C28 sugar–Hoogsteen base-pairs, as well as the A14–U25–U26 base triple. The corresponding secondary structure is shown on the right. The genomic numbering follows the convention in ref 11. The base-paired regions selected for creating initial configurations are shown in blue, while partial tempering acts in regions shown in red. (c, d) Secondary structures of two alternative internal loop arrangements.

1.6 Å heavy-atom root-mean-square deviation and 0.5 eRMSD, as shown in Supporting Information 1 (SI1).

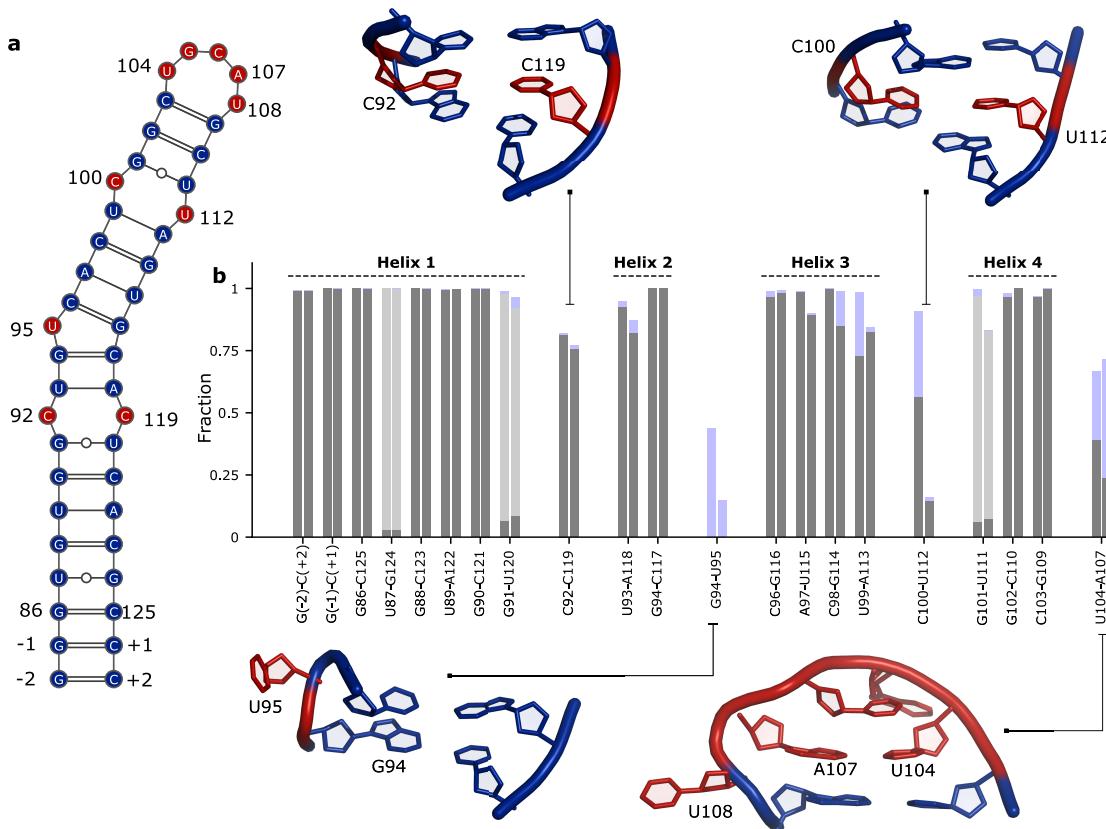
The sugar pucksers of US1, US2, and US4 are preferably in C2'-endo conformation, in full agreement with NMR measurements<sup>16</sup> (Figure 2e). The great majority (169 out of 179) of upper-bound NOE-derived distances are satisfied in MD simulations (Figure 2f). The most significant discrepancies consist of short proton–proton distances between G53–US4 that are not sampled in our MD runs. The full list of violated NOEs is shown in Figure SI2.

**Stem Loop 5a.** SL5a consists of two helices (blue nucleotides in Figure 3a) connected by an internal loop and capped by a hexaloop with sequence UUUCGU. We use the same procedure as described above for SL2 to generate initial configurations and to sample the ensemble of SL5a by enhancing regions where secondary structure is not known (red regions in Figure 3a). Helix 1 is highly stable throughout the 20 replicas × 3.6 μs partial tempering simulations (Figure 3b).

In the internal loop in SL5a we observe the formation of a U194–U212 base-pair with a population of ~60% (Figure 3b). The wobble GU base-pair between U195 and G210 in helix 2

is in equilibrium with an alternative WC/WC base-pair between U195 and U211. U212 is highly dynamic (see below), while other nucleotides in the internal loop, including G210, remain in the stack throughout MD simulations. Our predictions are overall compatible with NMR data, suggesting U194 to be base-paired and U195–G210 at least partially stable.<sup>32</sup>

The most frequent conformation of the apical loop is characterized by a wobble interaction between U200 and G204, with U205 completely solvent exposed. Sugar pucker conformations for nucleotides 201 to 204 are predominantly in C2'-endo conformation, at variance with NMR data, which indicate C2'-endo conformations only at positions U202 and C203.<sup>32</sup> In one of the two simulations we observe the transient formation of a U200–U205 base-pair capped by a canonical UUCG tetraloop fold.<sup>31</sup> The low population of this state is not surprising, as the canonical UUCG tetraloop fold has been shown not to be stabilized to a sufficient degree within this force field.<sup>52</sup> While the detailed architecture of the hexaloop has not yet been solved experimentally,<sup>11,32</sup> our computational prediction and available NMR data cannot be easily reconciled in this case. We suggest that the simulations might be



**Figure 5.** SL4 results. (a) Sequence and secondary structure of SL4. The genomic numbering follows the convention in ref 11. The region used for constructing the A-form templates is shown in blue. The sampling in the region shown in red is enhanced by partial tempering. (b) Population of base-pairs from 20 replicas  $\times$  1.6  $\mu$ s partial tempering simulations. Centroid of the most populated internal and apical loops discussed in the text is shown.

integrated with structural data collected in the future through a reweighting procedure.<sup>40</sup>

In addition to the analysis of individual base-pairs, we identify in the simulations frequent transitions from straight to kinked conformations (Figure 3c,d). The latter structure is reminiscent of a kink-turn motif, but lacks its signature interactions and sequence propensities.<sup>53</sup> This large-scale motion does not affect the structures of the individual stems, but only their relative orientation. We also note that kinked conformations are observed only when U212 is solvent exposed and therefore not engaged in base-pairing interactions with U194 or U195 (Figure 3e) and suggest that the relative populations could be modulated by binding of proteins or small molecules to this motif.

**Stem Loop 1.** Similar to SL5a, SL1 consists of two helices connected by an asymmetrical internal loop (Figure 4a). We use the pulling procedure to initialize our runs from the previously determined secondary structure.<sup>11</sup> While U13–A26/U17–A22 base pairs could not be confirmed by NMR, we have included these two interactions for generating the initial structures, as they were present in previous consensus secondary structure assignments.<sup>14</sup> Note that secondary structure restraints are absent during partial tempering simulations, and we indeed observe a rich dynamics involving the loop regions as well as the neighboring residues.

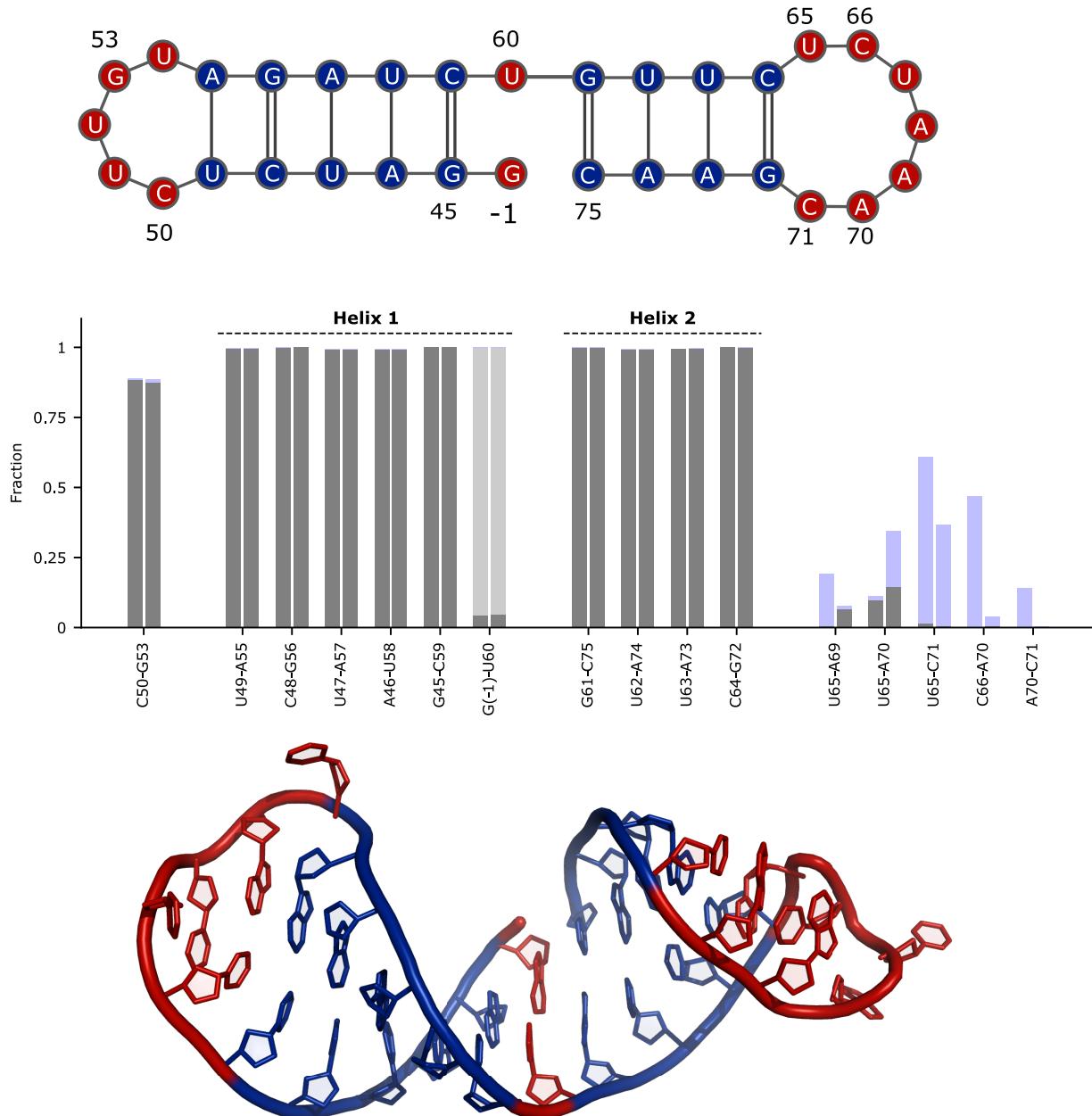
While helix 1 is stable, the terminal base-pair U11–A29 interacts through either the Watson–Crick/Watson–Crick or Hoogsteen–sugar edges (Figure 4b). The presence of this

noncanonical interaction allows the formation of a A12–C28 sugar–Hoogsteen interaction as well as of a U13–A27 base-pair. A26 is not completely solvent-exposed, but forms a base triple together with A14 and U25. In our simulations we also observe two alternative, lowly populated secondary structures (Figure 4c,d) where A12–A27/U13–A26 are base-paired, and either C28 or A29 is partially excluded from the stack. Qualitatively, the nucleobase arrangements shown in Figure 4b–d all preserve the coaxial stacking of the two helices.

The apical loop is highly dynamic and lacks a major, stable conformation with specific base–base interactions. The U17–A22 base-pair interconverts from Watson–Crick/Watson–Crick to sugar/Hoogsteen, and U18 can engage a sugar–Hogsteen interaction with C21, albeit with a small population (Figure S13).

**Stem Loop 4.** The SL4 element has a more complex organization, and consists of four A-form helices of different lengths linked by short pyrimidine internal loops and capped by a UGCAU pentaloop (Figure 5a). The experimentally derived secondary structure is generally preserved in the partial tempering simulations (Figure 5b). The A-form geometry between helix 1 and helix 2 is stabilized by a WC/WC C92–C119 base-pair (Figure 5b). U95 is solvent-exposed and transiently engages in noncanonical interactions with G94 (Figure 5b).

While the interactions observed in the first and second loops are relatively consistent between the two independent simulations, this is not the case for the third internal loop.



**Figure 6.** SL2+SL3 results. Top: Sequence and secondary structure of SL2+SL3. The genomic numbering follows the convention in ref 11. The region used for constructing the A-form templates is shown in blue. The sampling in the region shown in red is enhanced by partial tempering. Middle panel: Population of base-pairs from partial tempering simulations. Bottom panel: Centroid of the most populated cluster.

Here, the two simulations behave differently, suggesting slow motions that are not fully sampled in our simulations. In the first run C100 preferentially base-pairs with U112, while in the second run additional configurations are sampled, where either C100 or U112 is not in the stack. The partial overlap between the two runs, that we ascribe to insufficient sampling, is further confirmed by a principal component analysis of this internal loop region (SI4), as well as from the low number of round trips in replica space in run 1 (SIS). Experimental NMR measurements on SL4 agree with A-form helix compatible arrangement of the opposing residues C100 and U112,<sup>11</sup> thus

supporting the presence of pairing between the two bases. The apical loop structure resembles a tetraloop with U108 solvent-exposed and U104:A107 forming either a Watson–Crick/Watson–Crick or a sugar/Hogsteen noncanonical interaction, as shown in Figure 5b.

**SL2+SL3.** We here report the results on the SL2+SL3 construct, consisting of two consecutive stem-loop structures: SL2 (described above) and SL3, which contains the important transcription-regulating sequence required for subgenomic viral RNA synthesis.<sup>37</sup> Although we did not enforce this geometry, we find structures where the two helices are stacked

coaxially tail-to-tail, thereby forming a unique A-form-like stem (Figure 6).

The SL2 loop conformations are similar to those observed for SL2 alone, with a stable C50–G53 base-pair (see above). The UCUAAAC heptaloop in SL3 is highly flexible and adopts several conformations that are partially structured in terms of base-pairing (Figure 6). A cluster analysis of SL3 conformations reveals the presence of a loop structure where U65 forms a stacking interaction with C66, followed by the strand inversion and by two consecutive stackings between A68–A69–A70. This specific loop architecture can be considered the most RNA-like sampled in simulations, as we found it to be remarkably similar to the helix 58 loop in the archeal large ribosomal subunit<sup>54</sup> as well as to the anticodon loop in a nonproteinogenic tRNA<sup>55</sup> (Figure SI6).

Helix 2, containing the transcription regulatory sequence ACGAAC, is fully stable in our simulations. NMR measurements show that helix 2 is less stable than helix 1 in the absence of magnesium.<sup>11</sup> Such decreased stability of helix 2 is in line with a model for the body-to-leader fusion mechanism, in which the leader transcription regulatory sequence, “hidden” inside SL3, is required to unfold to be active.<sup>3</sup> The tail-to-tail structure observed in simulations, which may not form in the context of the full 5'-UTR construct, may overstabilize helix 2 in this construct. Note that quantifying duplex stability is a challenging problem that is not considered here, as it usually requires dedicated simulation protocols.<sup>47</sup>

## CONCLUSIONS

We report the prediction of the structure and conformational heterogeneity of five stem-loop elements, ranging in size between 17 and 44 nucleotides, located at the 5'-UTR in the SARS-CoV-2 genome. As a starting point for our study we use the secondary structure determined by NMR experiments.<sup>11</sup>

Structures of several of these elements have also been predicted using FARFAR2.<sup>17</sup> A principal component analysis of the structures generated for SL1 and SL2 shows small overlap between FARFAR2 and MD ensembles, suggesting the two methods to be complementary (SI7). In the case of SL2, the MD predictions mostly agree with the available NMR data, although small discrepancies remain. For this structure FARFAR2 does not predict the experimental loop architecture with the C50–G53 base-pair (SI1), and it thus appears that the ensemble generated by MD provides a more accurate description in this case. Given the scarcity of available experimental data on the 5'-UTR elements, it is not possible at this stage to compare the accuracy of different computational methods. Indeed, FARFAR2 has been shown to better predict experimental data compared to MD on an RNA internal loop.<sup>28</sup>

As in most MD studies, the results presented in this work depend on the amount of sampling. We took care in identifying and discussing states that are visited multiple times during one simulation and in both replicates, indicating these conformations to correspond to *bona fide*, highly populated metastable states. Moreover, the overall high agreement between the two independent simulations suggests that the conformational ensembles are relatively well converged. Our simulations of SL4, consisting of 44 nucleotides, represent an important exception with greater deviations between the two simulations, in particular of the third internal loop (SI4, SI5), indicating that our approach

would require larger computational resources for molecules of this size.

The modeling protocol used in the present study implicitly considers helix formation and internal loop dynamics as partially decoupled events. While this is a relatively common assumption in computational studies of RNA,<sup>18,27,31</sup> we can envisage situations in which the rearrangement of a loop requires one or more base pairs to break transiently. In our partial tempering MD simulations the secondary structure is not fixed: these rare events can thus be observed, albeit with a low frequency due to the absence of auxiliary replicas enhancing the breakage and formation of stem regions.

Despite recent improvements of RNA force fields,<sup>43,47,56</sup> the accuracy of the simulation results still depends on the accuracy of the physical models used, and we have validated our simulations against available experimental data. As more data become available, we note that it should be straightforward to integrate our extensive conformational sampling with data from future solution measurements (e.g., NMR or SAXS)<sup>38–40</sup> to refine the atomic-level description of SARS-CoV-2 RNA structure and dynamics.

## METHODS

MD simulations were conducted using the Amber ff99SB force field with parmbsc0<sup>57</sup> and OL3<sup>58</sup> corrections in conjunction with the OPC water model.<sup>59</sup> The initial conformations were minimized in a vacuum and then solvated in a dodecahedral box with a minimum distance from the box edge of 1.2 nm. Potassium ions were added to neutralize the box.<sup>60</sup> The systems were equilibrated in the NVT and NPT ensembles at 1 bar for 1 ns. Production runs were performed in the canonical ensemble using a stochastic velocity rescaling thermostat.<sup>61</sup> All bonds were constrained using the LINCS<sup>62</sup> algorithm, and equations of motion were integrated with a time step of 2 fs. Simulations were performed using GROMACS 2019.4<sup>63</sup> patched with PLUMED 2.5.3.<sup>64</sup> Force-field parameters are available at [http://github.com/srnas/ff/tree/opc-water-model/amber\\_na.ff](http://github.com/srnas/ff/tree/opc-water-model/amber_na.ff).

**Pulling Simulations.** The aim of the pulling simulations is quickly to generate three-dimensional RNA conformations with a specific input secondary structure. We adapt a procedure originally introduced for reconstructing atomistic structures from coarse-grained models.<sup>25</sup> Our protocol consists of three stages, as illustrated in Figure 2a,b for SL2:

1. Plain MD simulation initialized from extended, single-stranded, A-form RNA structures.
2. Pulling. We construct a reference structure with an ideal A-form double-stranded RNA template for regions with known secondary structure (shown in blue in Figure 2a) and perform pulling simulations to minimize the eRMSD<sup>45</sup> distance between the RNA molecule and the reference. We here use adiabatic bias MD simulations,<sup>65</sup> which introduce a biasing potential damping the fluctuations only when the system moves further away from the template. In this way the pulling procedure follows the thermal fluctuations of the system and does not require a predefined scheduling for the moving restraint.<sup>25</sup> When more than one stem is present, the bias acts on multiple templates simultaneously. The constant of the harmonic damping bias is set to 2 kJ/mol.
3. Plain MD. After 40 ns of pulling simulation, the external bias is removed and the system freely fluctuates without external forces.

At the end of the procedure we extract samples from the final plain MD stage satisfying the secondary structure constraints, i.e., with an eRMSD from the template of <0.75. This criterion guarantees that all base-pairs corresponding to those in the target structure are correctly formed (Figure 2b, bottom panel). Note that individual pulling simulations may fail: stems are sometimes not formed correctly and

unfold during the final plain MD stage (Figure S18). The success rate of this procedure ranges from ~80% for SL2 to ~10% for SL4, which consists of four-stem structures, consistent with previously reported benchmarks.<sup>25</sup> During preliminary tests we have observed that the formation of the correct secondary structure is hampered by early forming non-native contacts. For this reason we conduct pulling simulations at 340 K.

**Partial Tempering.** We resolute and minimize the RNA structures obtained from pulling as described above. We use partial tempering to enhance sampling in regions with unknown secondary structure.<sup>36</sup> One reference replica exchanges via the Metropolis–Hastings algorithm with parallel simulations that are conducted with a modified Hamiltonian. The modification is such that interactions in the cold region (blue and gray in Figure 2a, together with ions and water molecules) are kept at the reference temperature  $T = 300$  K. Interactions in the hot region (shown in red in Figure 2a) are rescaled by a factor of  $1/\lambda$ , while the interactions between cold and hot regions are at an intermediate effective temperature of  $T/(\sqrt{\lambda})$ , with  $0 < \lambda \leq 1$ . The system was then neutralized by adding a uniform compensating background. In this study we use 20 geometrically distributed replicas in the effective temperature range of 300 to 1000 K and attempting exchanges every 240 steps. Hot regions are shown in red in Figures 2–6 and include the phosphate group at the 3'-end of the selected region. We have found this small addition to be beneficial in preliminary test runs. System length, simulation time per replica, average acceptance rate, and number of round trips in replica space are reported in Table 1, and additional statistics are available in Figure

**Table 1. Simulated Systems and Related Statistics**

system	length (nt) <sup>a</sup>	simulation time [μs] <sup>b</sup>	acceptance rate [%] <sup>c</sup>	round trips <sup>d</sup>
SL1	29	3.6	15–41	44–191
SL2	17	2.6	26–51	220–435
SL2+SL3	32	1.6	4–23	0–10
SL4	44	1.9	7–30	0–20
SL5a	33	3.6	12–34	6–33

<sup>a</sup>Number of nucleotides (nt). <sup>b</sup>Per-replica simulation time. Each simulation consists of 20 replicas, and two independent runs for each system are carried out. The first 100 ns of each run is discarded during analysis. <sup>c</sup>Minimum and maximum acceptance rate across all replicas and runs. <sup>d</sup>Minimum and maximum number of round trips in replica space across all replicas and runs.

**S15.** To help the identification of slow degrees of freedom and assess convergence, we conduct two independent runs with different initial conformations obtained from two separate pulling simulations (Figure S19).

**Analysis and Visualization.** Ideal A-form structures were constructed using the nucleic acids builder (NAB) software in Ambertools.<sup>67</sup> Trajectories were analyzed using PLUMED<sup>64</sup> and MDtraj.<sup>68</sup> Base-pair annotation and puckering angles were calculated with Barnaba.<sup>50</sup>

## ASSOCIATED CONTENT

### Supporting Information

The Supporting Information is available free of charge at <https://pubs.acs.org/doi/10.1021/jacs.1c01094>.

RMSD between PDB: 2L6I, MD, and FARFAR2 simulations; NOE measurements violated in partial tempering simulations; 3D centroids of the first five clusters in the SL1 apical loop; partial tempering simulations of SL4 projected onto the first two principal components; average acceptance rate and number of round trips in replica space in partial tempering simulations; SL3 structure similar to experimental loop structures in the PDB database; principal component

analysis performed on FARFAR2 and MD ensembles for SL1 and SL2; example of unsuccessful pulling simulation; initial configurations in partial tempering simulations (PDF)

## AUTHOR INFORMATION

### Corresponding Authors

Sandro Bottaro – *Structural Biology and NMR Laboratory & Linderstrøm-Lang Centre for Protein Science, Department of Biology, University of Copenhagen, DK-2200 Copenhagen N, Denmark;*  [orcid.org/0000-0003-1606-890X](https://orcid.org/0000-0003-1606-890X); Email: [sandro.bottaro@bio.ku.dk](mailto:sandro.bottaro@bio.ku.dk)

Kresten Lindorff-Larsen – *Structural Biology and NMR Laboratory & Linderstrøm-Lang Centre for Protein Science, Department of Biology, University of Copenhagen, DK-2200 Copenhagen N, Denmark;*  [orcid.org/0000-0002-4750-6039](https://orcid.org/0000-0002-4750-6039); Email: [lindorff@bio.ku.dk](mailto:lindorff@bio.ku.dk)

### Author

Giovanni Bussi – *Scuola Internazionale Superiore di Studi Avanzati (SISSA), 34136 Trieste, Italy;*  [orcid.org/0000-0001-9216-5782](https://orcid.org/0000-0001-9216-5782)

Complete contact information is available at:

<https://pubs.acs.org/10.1021/jacs.1c01094>

### Notes

The authors declare no competing financial interest.

Analysis scripts, MD trajectories, and input files are available via [https://github.com/KULL-Centre/papers/tree/master/2020/COVID\\_SUTR\\_MD](https://github.com/KULL-Centre/papers/tree/master/2020/COVID_SUTR_MD) PLUMED input files are available in the PLUMED NEST<sup>69</sup> under the accession code 20.034, <https://www.plumed-nest.org/eggs/20/034/>.

## ACKNOWLEDGMENTS

We acknowledge Drs. Harald Schwalbe, Anna Wacker, Julia E. Weigand, Andreas Schlundt, and other members of the Covid19-NMR consortium for discussions and insights. S.B. and K.L.L. acknowledge funding from the Lundbeck Foundation BRAINSTRUC structural biology initiative (R155-2015-2666). We acknowledge access to computational resources from PRACE for the COVID-RNA project (COVID19-72). We thank Matteo Cagiada for preparing the TOC graphic.

## REFERENCES

- Kim, D.; Lee, J.-Y.; Yang, J.-S.; Kim, J. W.; Kim, V. N.; Chang, H. The architecture of SARS-CoV-2 transcriptome. *Cell* **2020**, *181*, 914–921.
- Lal, S. K. *Molecular Biology of the SARS-CoV-2 Virus*; Springer Science & Business Media, 2010.
- Yang, D.; Leibowitz, J. L. The structure and functions of coronavirus genomic 3' and 5' ends. *Virus Res.* **2015**, *206*, 120–133.
- Chen, S.-C.; Olsztyn, R. C. Group-specific structural features of the 5'-proximal sequences of coronavirus genomic RNAs. *Virology* **2010**, *401*, 29–41.
- Madhugiri, R.; Karl, N.; Petersen, D.; Lamkiewicz, K.; Fricke, M.; Wend, U.; Scheuer, R.; Marz, M.; Ziebuhr, J. Structural and functional conservation of cis-acting RNA elements in coronavirus 5'-terminal genome regions. *Virology* **2018**, *517*, 44–55.
- Schubert, K.; Karousis, E. D.; Jomaa, A.; Scaiola, A.; Echeverria, B.; Gurzeler, L.A.; Leibundgut, M.; Thiel, V.; Mühlmann, O.; Ban, N. SARS-CoV-2 Nsp1 binds the ribosomal mRNA channel to inhibit translation. *Nat. Struct. Mol. Biol.* **2020**, *27*, 959–966.

- (7) Sola, I.; Almazan, F.; Zuniga, S.; Enjuanes, L. Continuous and discontinuous RNA synthesis in coronaviruses. *Annu. Rev. Virol.* **2015**, *2*, 265–288.
- (8) Isberman, C.; Roden, C. A.; Boerneke, M. A.; Sealfon, R. S.; McLaughlin, G. A.; Jungreis, I.; Fritch, E. J.; Hou, Y. J.; Ekena, J.; Weidmann, C. A.; et al. Genomic RNA elements drive phase separation of the SARS-CoV-2 nucleocapsid. *Mol. Cell* **2020**, *80*, 1078–1091.
- (9) Carlson, C. R.; Afshari, J. B.; Ghent, C. M.; Howard, C. J.; Hartooni, N.; Safari, M.; Frankel, A. D.; Morgan, D. O. Phosphoregulation of phase separation by the SARS-CoV-2 N protein suggests a biophysical basis for its dual functions. *Mol. Cell* **2020**, *80*, 1092–110.
- (10) Huston, N. C.; Wan, H.; Strine, M. S.; Tavares, R. d. C. A.; Wilen, C. B.; Pyle, A. M. Comprehensive in vivo secondary structure of the SARS-CoV-2 genome reveals novel regulatory motifs and mechanisms. *Mol. Cell* **2021**, *81*, 584–598.
- (11) Wacker, A.; Weigand, J. E.; Akabayov, S. R.; Altinçekic, N.; Bains, J. K.; Banijamali, E.; Binns, O.; Castillo-Martinez, J.; Cetiner, E.; Ceylan, B.; et al. Secondary structure determination of conserved SARS-CoV-2 RNA elements by NMR spectroscopy. *Nucleic Acids Res.* **2020**, *48*, 12415.
- (12) Andrews, R. J.; Peterson, J. M.; Haniff, H. S.; Chen, J.; Williams, C.; Grefe, M.; Disney, M. D.; Moss, W. N. An in silico map of the SARS-CoV-2 RNA Structurome. *BioRxiv*. 2020, DOI: 10.1101/2020.04.17.045161 (accessed on 17/05/2021).
- (13) Manfredonia, I.; Nithin, C.; Ponce-Salvatierra, A.; Ghosh, P.; Wirecki, T. K.; Marinus, T.; Ogando, N. S.; Snijder, E. J.; van Hemert, M. J.; Bujnicki, J. M.; et al. Genome-wide mapping of SARS-CoV-2 RNA structures identifies therapeutically-relevant elements. *Nucleic Acids Res.* **2020**, *48*, 12436.
- (14) Rangan, R.; Zheludev, I. N.; Hagey, R. J.; Pham, E. A.; Wayment-Steele, H. K.; Glenn, J. S.; Das, R. RNA genome conservation and secondary structure in SARS-CoV-2 and SARS-related viruses: a first look. *RNA* **2020**, *26*, 937–959.
- (15) Miao, Z.; Tidu, A.; Eriani, G.; Martin, F. Secondary structure of the SARS-CoV-2 5'-UTR. *RNA Biol.* **2021**, *18*, 1–10.
- (16) Lee, C. W.; Li, L.; Giedroc, D. P. The solution structure of coronaviral stem-loop 2 (SL2) reveals a canonical CUYG tetraloop fold. *FEBS Lett.* **2011**, *585*, 1049–1053.
- (17) Rangan, R.; Watkins, A. M.; Chacon, J.; Kretsch, R.; Kladwang, W.; Zheludev, I. N.; Townley, J.; Rynge, M.; Thain, G.; Das, R. De novo 3D models of SARS-CoV-2 RNA elements from consensus experimental secondary structures. *Nucleic Acids Res.* **2021**, *49*, 3092.
- (18) Watkins, A. M.; Rangan, R.; Das, R. FARFAR2: Improved de novo Rosetta prediction of complex global RNA folds. *Structure* **2020**, *28*, 963–976.
- (19) Zhang, K.; Zheludev, I. N.; Hagey, R. J.; Wu, M. T.-P.; Haslecker, R.; Hou, Y. J.; Kretsch, R.; Pintilie, G. D.; Rangan, R.; Kladwang, W.; et al. Cryo-electron microscopy and exploratory antisense targeting of the 28-kDa frameshift stimulation element from the SARS-CoV-2 RNA genome. *Biorxiv*. 2020, DOI: 10.1101/2020.07.18.209270 (accessed on 17/05/2021).
- (20) Bergonzo, C.; Szakal, A. L. Using All-Atom Potentials to Refine RNA Structure Predictions of SARS-CoV-2 Stem Loops. *Int. J. Mol. Sci.* **2020**, *21*, 6188.
- (21) Ganser, L. R.; Kelly, M. L.; Herschlag, D.; Al-Hashimi, H. M. The roles of structural dynamics in the cellular functions of RNAs. *Nat. Rev. Mol. Cell Biol.* **2019**, *20*, 474–489.
- (22) Alderson, T. R.; Kay, L. E. NMR spectroscopy captures the essential role of dynamics in regulating biomolecular function. *Cell* **2021**, *184*, 577–595.
- (23) Parisien, M.; Major, F. The MC-Fold and MC-Sym pipeline infers RNA structure from sequence data. *Nature* **2008**, *452*, 51–55.
- (24) Cheng, C. Y.; Chou, F.-C.; Das, R. *Methods in Enzymology*; Elsevier, 2015; Vol. 553; pp 35–64.
- (25) Poblete, S.; Bottaro, S.; Bussi, G. Effects and limitations of a nucleobase-driven backmapping procedure for nucleic acids using steered molecular dynamics. *Biochem. Biophys. Res. Commun.* **2018**, *498*, 352–358.
- (26) Ebrahimi, P.; Kaur, S.; Baronti, L.; Petzold, K.; Chen, A. A. A two-dimensional replica-exchange molecular dynamics method for simulating RNA folding using sparse experimental restraints. *Methods* **2019**, *162*, 96–107.
- (27) Angela, M. Y.; Gasper, P. M.; Cheng, L.; Lai, L. B.; Kaur, S.; Gopalan, V.; Chen, A. A.; Lucks, J. B. Computationally reconstructing cotranscriptional RNA folding from experimental data reveals rearrangement of non-native folding intermediates. *Mol. Cell* **2021**, *81*, 870–883.
- (28) Shi, H.; Rangadurai, A.; Abou Assi, H.; Roy, R.; Case, D. A.; Herschlag, D.; Yesselman, J. D.; Al-Hashimi, H. M. Rapid and accurate determination of atomistic RNA dynamic ensemble models using NMR and structure prediction. *Nat. Commun.* **2020**, *11*, 1–14.
- (29) Wang, L.; Friesner, R. A.; Berne, B. Replica exchange with solute scaling: a more efficient version of replica exchange with solute tempering (REST2). *J. Phys. Chem. B* **2011**, *115*, 9431–9438.
- (30) Kamiya, M.; Sugita, Y. Flexible selection of the solute region in replica exchange with solute tempering: Application to protein-folding simulations. *J. Chem. Phys.* **2018**, *149*, 072304.
- (31) Chyzy, P.; Kulik, M.; Re, S.; Sugita, Y.; Trylska, J. Mutations of N1 riboswitch affect its dynamics and recognition by neomycin through conformational selection. *Frontiers in molecular biosciences* **2021**, *8*, 12.
- (32) Schnieders, R.; Peter, S. A.; Banijamali, E.; Riad, M.; Altinçekic, N.; Bains, J. K.; Ceylan, B.; Fürtg, B.; Grün, J. T.; Hengesbach, M.; et al. 1 H, 13 C and 15 N chemical shift assignment of the stem-loop Sa from the 5'-UTR of SARS-CoV-2. *Biomol. NMR Assignments* **2021**, *15*, 1–9.
- (33) Novakovic, M.; Kupce, E.; Scherf, T.; Oxenfarth, A.; Schnieders, R.; Grün, T.; Wirmer-Bartoschek, J.; Richter, C.; Schwalbe, H.; Frydman, L. Magnetization transfer to enhance NOE cross-peaks among labile protons: Applications to imino-imino sequential walks in SARS-CoV-2-derived RNAs. *Angew. Chem., Int. Ed.* **2021**, *60*, 11884.
- (34) Escors, D.; Izeta, A.; Capiscol, C.; Enjuanes, L. Transmissible gastroenteritis coronavirus packaging signal is located at the 5' end of the virus genome. *J. Virol.* **2003**, *77*, 7890–7902.
- (35) Tidu, A.; Janvier, A.; Schaeffer, L.; Sosnowski, P.; Kuhn, L.; Hammann, P.; Westhof, E.; Eriani, G.; Martin, F. The viral protein NSP1 acts as a ribosome gatekeeper for shutting down host translation and fostering SARS-CoV-2 translation. *RNA* **2021**, *27*, 253–264.
- (36) Yang, D.; Liu, P.; Giedroc, D. P.; Leibowitz, J. Mouse hepatitis virus stem-loop 4 functions as a spacer element required to drive subgenomic RNA synthesis. *J. Virol.* **2011**, *85*, 9199–9209.
- (37) Dufour, D.; Mateos-Gomez, P. A.; Enjuanes, L.; Gallego, J.; Sola, I. Structure and functional relevance of a transcription-regulating sequence involved in coronavirus discontinuous RNA synthesis. *J. Virol.* **2011**, *85*, 4963–4973.
- (38) Cesari, A.; Reißer, S.; Bussi, G. Using the maximum entropy principle to combine simulations and solution experiments. *Computation* **2018**, *6*, 15.
- (39) Bottaro, S.; Lindorff-Larsen, K. Biophysical experiments and biomolecular simulations: A perfect match? *Science* **2018**, *361*, 355–360.
- (40) Bottaro, S.; Bengtsen, T.; Lindorff-Larsen, K. *Structural Bioinformatics*; Springer, 2020; pp 219–240.
- (41) Stelzer, A. C.; Frank, A. T.; Kratz, J. D.; Swanson, M. D.; Gonzalez-Hernandez, M. J.; Lee, J.; Andricioaei, I.; Markovitz, D. M.; Al-Hashimi, H. M. Discovery of selective bioactive small molecules by targeting an RNA dynamic ensemble. *Nat. Chem. Biol.* **2011**, *7*, 553.
- (42) Costales, M. G.; Childs-Disney, J. L.; Haniff, H. S.; Disney, M. D. How we think about targeting RNA with small molecules. *J. Med. Chem.* **2020**, *63*, 8880–8900.
- (43) Dans, P. D.; Gallego, D.; Balaceanu, A.; Darré, L.; Gómez, H.; Orozco, M. Modeling, simulations, and bioinformatics at the service of RNA structure. *Chem.* **2019**, *5*, 51–73.

- (44) Amaro, R. E.; Mulholland, A. J. A community letter regarding sharing biomolecular simulation data for COVID-19. *J. Chem. Inf. Model.* **2020**, *60*, 2653–2656.
- (45) Bottaro, S.; Di Palma, F.; Bussi, G. The role of nucleobase interactions in RNA structure and dynamics. *Nucleic Acids Res.* **2014**, *42*, 13306–13314.
- (46) Leontis, N. B.; Westhof, E. Geometric nomenclature and classification of RNA base pairs. *RNA* **2001**, *7*, 499–512.
- (47) Šponer, J.; Bussi, G.; Krepl, M.; Banáš, P.; Bottaro, S.; Cunha, R. A.; Gil-Ley, A.; Pinamonti, G.; Poblete, S.; Jurečka, P.; Walter, N. G.; Otyepka, M. RNA structural dynamics as captured by molecular simulations: a comprehensive overview. *Chem. Rev.* **2018**, *118*, 4177–4338.
- (48) Jucker, F. M.; Pardi, A. Solution structure of the CUUG hairpin loop: a novel RNA tetraloop motif. *Biochemistry* **1995**, *34*, 14416–14427.
- (49) Bottaro, S.; Lindorff-Larsen, K. Mapping the universe of RNA tetraloop folds. *Biophys. J.* **2017**, *113*, 257–267.
- (50) Bottaro, S.; Bussi, G.; Pinamonti, G.; Reißer, S.; Boomsma, W.; Lindorff-Larsen, K. Barnaba: software for analysis of nucleic acid structures and trajectories. *RNA* **2019**, *25*, 219–231.
- (51) Cheong, C.; Varani, G.; Tinoco, I. Solution structure of an unusually stable RNA hairpin, 5GGAC (UUCG) GUCC. *Nature* **1990**, *346*, 680–682.
- (52) Mráziková, K.; Mlynšký, V.; Kuřrova, P.; Pokorná, P.; Kruse, H.; Krepl, M.; Otyepka, M.; Banas, P.; Šponer, J. UNCG RNA tetraloop as a formidable force-field challenge for MD simulations. *J. Chem. Theory Comput.* **2020**, *16*, 7601–7617.
- (53) Klein, D.; Schmeing, T.; Moore, P.; Steitz, T. The kink-turn: a new RNA secondary structure motif. *EMBO journal* **2001**, *20*, 4214–4221.
- (54) Ban, N.; Nissen, P.; Hansen, J.; Moore, P. B.; Steitz, T. A. The complete atomic structure of the large ribosomal subunit at 2.4 Å resolution. *Science* **2000**, *289*, 905–920.
- (55) Chang, A. T.; Nikonorowicz, E. P. Solution nuclear magnetic resonance analyses of the anticodon arms of proteinogenic and nonproteinogenic tRNAGly. *Biochemistry* **2012**, *51*, 3662–3674.
- (56) Vangaveti, S.; Ranganathan, S. V.; Chen, A. A. Advances in RNA molecular dynamics: a simulator's guide to RNA force fields. *Wiley Interdisciplinary Reviews: RNA* **2017**, *8*, e1396.
- (57) Pérez, A.; Marchán, I.; Svozil, D.; Šponer, J.; Cheatham, T. E., III; Laughton, C. A.; Orozco, M. Refinement of the AMBER Force Field for Nucleic Acids: Improving the Description of alpha gamma Conformers. *Biophys. J.* **2007**, *92*, 3817–3829.
- (58) Zgarbová, M.; Otyepka, M.; Šponer, J.; Mlaadek, A.; Banas, P.; Cheatham, T. E., III; Jurecka, P. Refinement of the Cornell et al. nucleic acids force field based on reference quantum chemical calculations of glycosidic torsion profiles. *J. Chem. Theory Comput.* **2011**, *7*, 2886–2902.
- (59) Izadi, S.; Anandakrishnan, R.; Onufriev, A. V. Building water models: a different approach. *J. Phys. Chem. Lett.* **2014**, *5*, 3863–3871.
- (60) Joung, I. S.; Cheatham, T. E., III Determination of alkali and halide monovalent ion parameters for use in explicitly solvated biomolecular simulations. *J. Phys. Chem. B* **2008**, *112*, 9020–9041.
- (61) Bussi, G.; Donadio, D.; Parrinello, M. Canonical sampling through velocity rescaling. *J. Chem. Phys.* **2007**, *126*, 014101.
- (62) Hess, B.; Bekker, H.; Berendsen, H. J.; Fraaije, J. G. LINCS: a linear constraint solver for molecular simulations. *J. Comput. Chem.* **1997**, *18*, 1463–1472.
- (63) Pronk, S.; Páll, S.; Schulz, R.; Larsson, P.; Bjelkmar, P.; Apostolov, R.; Shirts, M. R.; Smith, J. C.; Kasson, P. M.; van der Spoel, D.; Hess, B.; Lindahl, E. GROMACS 4.5: a high-throughput and highly parallel open source molecular simulation toolkit. *Bioinformatics* **2013**, *29*, 845–854.
- (64) Tribello, G. A.; Bonomi, M.; Branduardi, D.; Camilloni, C.; Bussi, G. PLUMED 2: New feathers for an old bird. *Comput. Phys. Commun.* **2014**, *185*, 604–613.
- (65) Marchi, M.; Ballone, P. Adiabatic bias molecular dynamics: a method to navigate the conformational space of complex molecular systems. *J. Chem. Phys.* **1999**, *110*, 3697–3702.
- (66) Bussi, G. Hamiltonian replica exchange in GROMACS: a flexible implementation. *Mol. Phys.* **2014**, *112*, 379–384.
- (67) Macke, T. J.; Case, D. A. *Molecular Modeling of Nucleic Acids*; ACS Publications, 1998; pp 379–393.
- (68) McGibbon, R. T.; Beauchamp, K. A.; Harrigan, M. P.; Klein, C.; Swails, J. M.; Hernández, C. X.; Schwantes, C. R.; Wang, L.-P.; Lane, T. J.; Pande, V. S. MDTraj: a modern open library for the analysis of molecular dynamics trajectories. *Biophys. J.* **2015**, *109*, 1528–1532.
- (69) Bonomi, M.; Bussi, G.; Camilloni, C.; Tribello, G. A.; Banáš, P.; Barducci, A.; Bernetti, M.; Bolhuis, P. G.; Bottaro, S.; Branduardi, D.; et al. Promoting transparency and reproducibility in enhanced molecular simulations. *Nat. Methods* **2019**, *16*, 670–673.