

# Hierarchical Assembly of Single-Stranded RNA

Lisa M. Pietrek, Lukas S. Stelzl, and Gerhard Hummer\*



Cite This: *J. Chem. Theory Comput.* 2024, 20, 2246–2260



Read Online

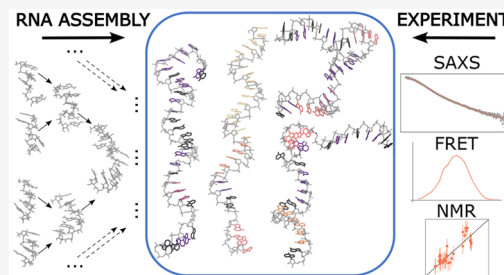
ACCESS |

Metrics & More

Article Recommendations

Supporting Information

**ABSTRACT:** Single-stranded RNA (ssRNA) plays a major role in the flow of genetic information—most notably, in the form of messenger RNA (mRNA)—and in the regulation of biological processes. The highly dynamic nature of chains of unpaired nucleobases challenges structural characterizations of ssRNA by experiments or molecular dynamics (MD) simulations alike. Here, we use hierarchical chain growth (HCG) to construct ensembles of ssRNA chains. HCG assembles the structures of protein and nucleic acid chains from fragment libraries created by MD simulations. Applied to homo- and heteropolymeric ssRNAs of different lengths, we find that HCG produces structural ensembles that overall are in good agreement with diverse experiments, including nuclear magnetic resonance (NMR), small-angle X-ray scattering (SAXS), and single-molecule Förster resonance energy transfer (FRET). The agreement can be further improved by ensemble refinement using Bayesian inference of ensembles (BioEn). HCG can also be used to assemble RNA structures that combine base-paired and base-unpaired regions, as illustrated for the 5′ untranslated region (UTR) of SARS-CoV-2 RNA.



## 1. INTRODUCTION

Single-stranded RNAs (ssRNAs) play important roles in many cellular processes, in particular, in the transmission of genetic information in the form of messenger RNA (mRNA). Noncoding stretches in mRNA or fully noncoding ssRNAs have key roles in the regulation of transcription and translation,<sup>1</sup> e.g., by acting as riboswitches<sup>2</sup> or by regulating the nuclear export of mRNA, and its stability and translation via polyadenylation.<sup>3</sup> In solution, ssRNAs can remain dynamically fully flexible and unstructured, transiently adopt secondary structures with paired bases, or form more stable secondary structures in complex with a binding partner.<sup>4,5</sup> mRNA in rapidly dividing cells was found to be substantially less structured than *in vitro*.<sup>6</sup> Therapeutics based on mRNA have long been explored,<sup>7</sup> which has recently led to the development of vaccines based on mRNA. In addition, nonbase-paired regions in RNA have emerged as promising drug targets.<sup>8</sup>

To improve our understanding of ssRNA and their functional mechanisms, we need to characterize their structural and dynamical features. However, experimentally investigating disordered ssRNA remains a challenging task. Nuclear magnetic resonance (NMR) techniques provide powerful tools to investigate local structure and dynamics with high-resolution in short disordered stretches of ssRNA shifts.<sup>9–15</sup> Small-angle X-ray scattering (SAXS) studies<sup>16–18</sup> or Förster resonance energy transfer (FRET) techniques<sup>16,19–21</sup> yield insight into the global structure of flexible biomolecules. The negatively charged ssRNA molecules have been shown to be strongly dependent on environmental buffer conditions,

including ion concentration and type,<sup>11,16–18</sup> an effect seen also in molecular dynamics (MD) simulations.<sup>9,21</sup>

Structural ensembles of ssRNA that capture the heterogeneity of these highly dynamic systems in atomic detail help the interpretation of data from experiments. Most experiments report ensemble averages. Such ensembles can, in principle, be obtained by performing MD simulations. However, MD simulations suffer from inaccuracies in the available force fields.<sup>22,23</sup> For RNA, special care has to be taken in setting the buffer conditions and choosing the ion force field parameters.<sup>24,25</sup> For biopolymers, small systematic errors in, say, backbone torsion potentials add up and result in major structural imbalances.<sup>26</sup> Inaccuracies in the energetics are amplified by the broad and shallow energy landscape of flexible biomolecules,<sup>13,21,27,28</sup> which requires extensive sampling. The sampling of ssRNA structural ensembles by MD simulations thus suffers both from systematic uncertainties due to inaccuracies in the force field and from statistical uncertainties due to the slow structural dynamics.

Fragment assembly is a promising approach to model RNA 3D structures. In early applications of RNA fragment assembly, Das et al. used FARFAR, a Rosetta-like fragment assembly approach to model noncanonical double-stranded (dsRNA) structure with atomistic detail.<sup>29</sup> Their fragment structures

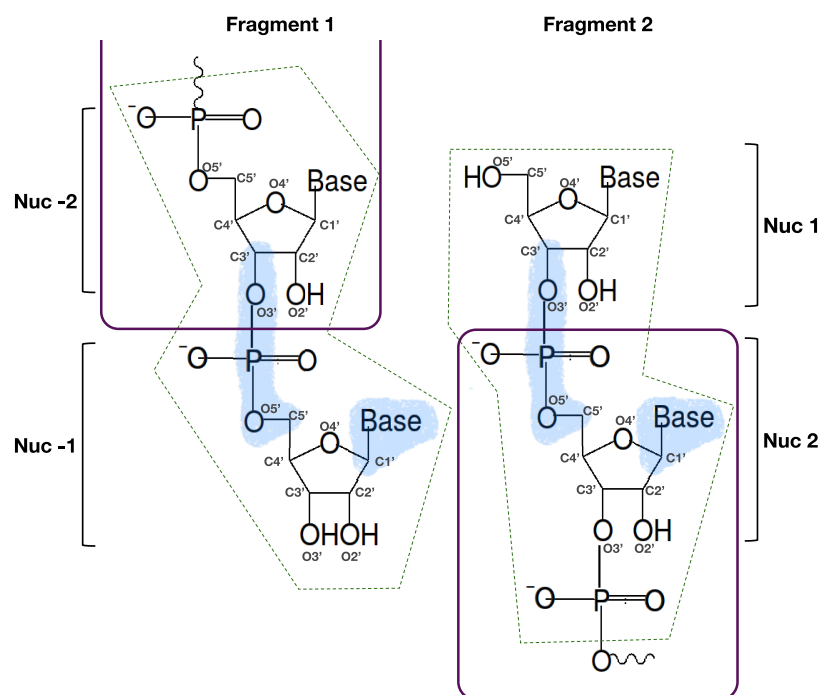
**Received:** September 22, 2023

**Revised:** December 9, 2023

**Accepted:** January 25, 2024

**Published:** February 16, 2024





**Figure 1.** Fragment alignment. Aligned heavy atoms are highlighted by blue shading. Heavy atoms within the dashed dark green boxes were excluded from the clash search. The purple lines indicate the regions of the two fragments that are included in the assembled chain.

were drawn from a library based on RNA structure in the large ribosomal subunit.<sup>29</sup> The more recent FARFAR2 approach<sup>30</sup> has been used to generate ensembles of short ssRNA polymers.<sup>14</sup> Chojnowski et al. developed a method to model 3D structures of short RNA polymers featuring base-paired strands as well as unpaired strands involved in loops by assembling RNA fragments from the PDB with the option to include experimental restraints.<sup>31</sup>

Ensembles of flexible biopolymers can be improved by integrating available experimental data. Approaches such as Bayesian/Maximum Entropy (BME)<sup>32–35</sup> and Bayesian inference<sup>36–40</sup> have been shown to work well in applications to ensembles of disordered biomolecules. For instance, Bottaro and co-workers refined tetrameric fragments according to NMR data using a BME approach, to improve their structural ensembles obtained via MD simulation, resulting in a more accurate description of the thermodynamic states.<sup>13</sup> In another example, Bergonzo et al. showed that a BME approach helped to improve conformational ensembles of a heteropolymeric oligonucleotide by integrating NMR and SAXS experimental data.<sup>14</sup> Alternatively, integration of experimental information can help to build models of observed molecules.<sup>18,41,42</sup>

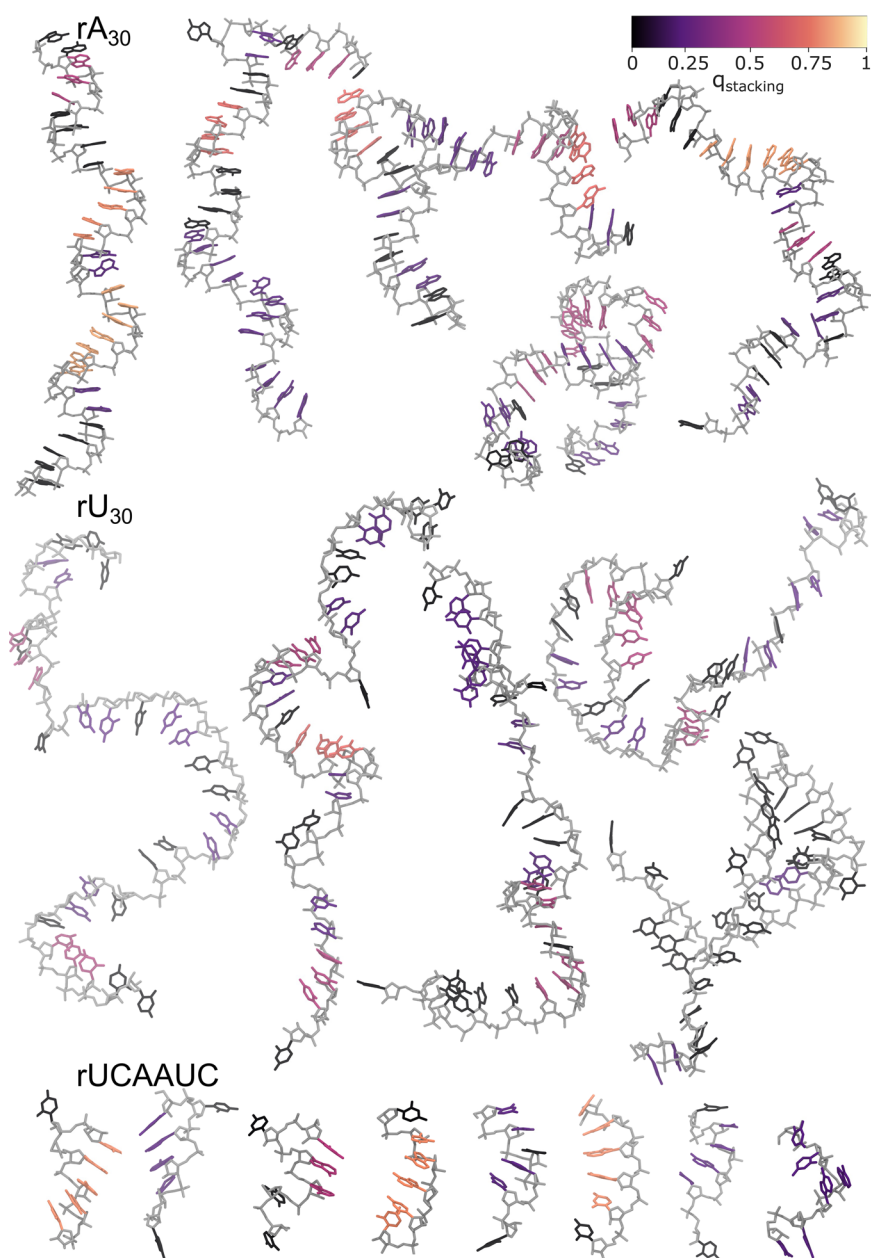
In previous work, we have introduced the hierarchical chain growth (HCG) method for disordered proteins.<sup>28,40,43</sup> We found that HCG is a robust approach well suited to efficiently growing broad structural ensembles of disordered proteins with atomic detail that are consistent with experimental findings. Here, we adapted HCG to model structural ensembles of disordered ssRNA. We focus on systems for which experimental data are available as reference:<sup>14,18,21</sup> homopolymeric adenosine monophosphate multimers ( $rA_n$  with  $n = 19, 30$ ), homopolymeric uridine monophosphate 30mer ( $rU_{30}$ ), and the short disordered heteropolymeric ssRNA rUCAAUC. We implemented ssRNA fragment assembly in the form of a Monte Carlo chain growth

algorithm, which we then used to assemble structural ensembles of conformations with atomic detail. We validated the modeled ensembles against diverse experimental data and could establish good agreement on average without refinement. We further improved the agreement with experimental observations by integrating experimental data using Bayesian inference of ensembles (BioEn) as a gentle ensemble refinement method.<sup>37</sup> As a proof of principle, we demonstrate that ssRNA chains grown with HCG can be combined with models of dsRNA, paving the way toward modeling short unstructured linkers, terminal untranslated regions (UTRs), or loops.

## 2. METHODS

**2.1. MD Fragment Library.** For poly adenine (A) RNA, an  $rA_4$  tetramer was modeled using the AMBER suite of programs.<sup>44</sup> The oxygen atoms of the terminal ribose groups at the 5' and 3' ends were protonated (Figure S1). For heteropolymeric ssRNA, we used heterotetrameric fragments rGXYZ. The nucleotide at the 5' position was fixed as guanine (G). We chose guanine as the headgroup, first, to mimic the interior of the ssRNA by providing a purine platform for stacking and, second, to facilitate the alignment with a relatively large base. For the following three nucleotides "XYZ", we used all  $4^3 = 64$  combinations of G, A, cytosine (C), and uracil (U). Each RNA fragment was placed in a dodecahedral box and solvated in TIP4P-D water<sup>45</sup> with 150 mM NaCl.<sup>46</sup> Charge neutrality was established with excess sodium ions. On average, the resulting systems contained about 6600 atoms in total. The RNA fragments were modeled with the DESRES<sup>23</sup> force field. We thus performed the fragment MD simulations using the same force field, water model, and ion parameters as described before.<sup>21</sup>

MD simulations were performed with GROMACS/2018.8.<sup>47</sup> Bonds including hydrogen atoms were constrained



**Figure 2.** Snapshots of ssRNA polymers grown with HCG. Representative renders of structures drawn at random from ensembles of (top) rA<sub>30</sub>, (center) rU<sub>30</sub>, and (bottom) rUCAAUC sampled by fragment assembly. The nucleic backbone is shown in light gray, and nucleobases are colored according to their base stacking factor<sup>55</sup>  $q_{\text{stacking}}$  (top: color code). Hydrogen atoms are omitted for clarity.

using the P-LINCS algorithm.<sup>48</sup> To maintain the pressure at a constant value of 1 bar, the Parrinello–Rahman barostat<sup>49</sup> was used. The cutoff distances for van der Waals and real-space electrostatic interactions were set to 1.2 nm. Electrostatic interactions were calculated using the particle mesh Ewald method<sup>50</sup> with the Fourier spacing set to 0.16 nm. The system was first energy minimized, followed by 400 ps of MD equilibration. The production REMD simulation was run in the NPT ensemble for 100 ns, with structures saved every 10 ps. For all tetramer fragments, we used 25 replicas that collectively spanned a temperature range of 300–431 K, as calculated using the algorithm by Patriksson and van der Spoel.<sup>51</sup> For each system, 10000 different structures collected at equally spaced time points from the replica simulated at 300 K were used for the respective fragment library.

**2.2. Hierarchical Chain Growth.** We adapted HCG<sup>40,43</sup> to grow full-length models of disordered homo- and heteropolymeric ssRNA chains from MD rA<sub>4</sub> and rGXYZ fragments. HCG was previously implemented and validated to model extensive ensembles of intrinsically disordered proteins (IDPs), displaying average properties that are in line with experimental observables.<sup>28,40,43</sup> HCG performs fragment assembly; i.e., a pool of fragment structures is combined at random into long polymers. The structural alignment of individual fragments and the rejection of poorly aligned or sterically clashing fragment pairs are critical for the quality of the resulting ensembles in terms of both local and global structural properties. We note that besides the root-mean-square distance (RMSD) alignment criterion and the steric exclusion we did not include any kind of attractive or repulsive interfragment interaction during the assembly. Thus, only

intrafragment electrostatic interactions are considered in HCG, as sampled in the fragment MD simulations. However, we integrate experimental data on a global level to account for possible discrepancies in assembled polymers.

We used a fragment alignment strategy that not only focuses on the nucleic acid backbone but also accounts for the position of the base. For two conformations of fragments adjacent in sequence and drawn at random from the respective pool, we performed a rigid-body superimposition of the O3' atom, phosphate atom, and OS' atom connecting nucleotides -2 and -1 in fragment 1 and nucleotides 1 and 2 in fragment 2, and all atoms of the nucleobase as well as the C1 atom of nucleotide -1 in fragment 1 and of nucleotide 2 in fragment 2 (Figure 1, light blue shaded area). For a successful alignment, we required the RMSD of the superimposed atoms to be below a given threshold,  $\text{RMSD} < 0.64 \text{ \AA}$ . In the superimposition, we doubled the weights of the aligned backbone atoms relative to the aligned atoms of the nucleobase to produce atom distances within the expected range.<sup>52</sup>

Alignment was followed by a search for steric clashes, defined as heavy atom distances below a cutoff of 2 Å. Note that we did not consider hydrogen atoms in clash detection. Atoms in the fragment-overlap region were excluded from the heavy atom clash search (Figure 1, dark green dashed boxes). Any steric clash resulted in the rejection of the fragment pair. Otherwise, the two fragments were merged. In merged fragments, nucleotide -1 from the first and nucleotide 1 from the second fragment were removed. In this way, the assembled chain featured only nucleotides sampled at the second and third positions (X and Y) of the rGXYZ fragments. The terminal nucleotides (G and Z) were treated as capping groups. We repeated this procedure in each hierarchical level of HCG until we reached the full-length sequence. For each polymer investigated in the present work, we grew ensembles with 10000 members. The RNA structure libraries (i.e., the MD fragment library as well as exemplary structures from the HCG ensembles discussed in this work) are available at <https://zenodo.org/record/8369324>. The HCG code to assemble ssRNA to the hierarchical chain growth is available at the GitHub repository <https://github.com/bio-phys/hierarchical-chain-growth/>.

We note that the ribose atoms are not included in the superimposition. We found that by not enforcing the sugar pucker configuration, we increased the diversity of grown structures and benefitted from diverse sugar pucker configurations sampled in fragment MD simulations. By including the nucleobase in the alignment, we improved the configuration of stacked bases, which is important to produce reasonable stacking also for longer sequences (Figure 2). The extent of base stacking in the assembled structures will to a significant degree be predetermined by the fragment library entering HCG and thus the MD simulation force field used to create the library.<sup>11</sup>

**2.3. Modeling the 5' UTR of SARS-CoV-2 RNA.** To build a structural ensemble of the 5' UTR of SARS-CoV-2 RNA by HCG, we combined fragments for the ssRNA segments with structural models for the stem loops using the secondary structure as input. The conformations of the structured stem-loop regions were randomly drawn from libraries filled with structures from MD trajectories published previously.<sup>53</sup> The connected disordered regions were grown using HCG as described above according to the sequence in ref 53. For the assembly, the same scheme as implemented in

HCG was used. In particular, the adjacent regions (fragments) were assembled in a hierarchical manner in subsequent levels. In the final level, the full-length model was assembled with a total of 233 nucleotides (sequence and secondary structure in Supplementary Figure S2). For the heavy atom superimposition, we set the RMSD cutoff to 1 Å, and the clash radius was kept at 2 Å. We grew only a small ensemble of 50 full-length chains. We note that the region spanned by nucleotides 162–200 was predicted to be structured and a part of stem-loop SL5.<sup>54</sup> However, to our knowledge for this region, there has been no structure solved so far. Therefore, we here modeled this region as a single-stranded region with HCG.

**2.4. Mapping of FRET Labels.** HCG is naturally suited to the inclusion of molecular labels, such as covalently attached fluorophores. To build a pool of dye-labeled rA<sub>4</sub> fragments, we used an MD library for the dyes Alexa Fluor 594 and Alexa Fluor 488 attached to dideoxyadenosinemonophosphate (dA<sub>2</sub>). The use of dA<sub>2</sub>-dye fragments to model fluorophores attached to both DNA and RNA chains has been validated by Grotz et al.<sup>21</sup> A random structure was drawn from the pool of rA<sub>4</sub> fragments and from the Alexa Fluor 594 or Alexa Fluor 488 MD library, to either label the 5' or 3' end, respectively, of the rA<sub>4</sub> fragments. We performed a rigid body alignment of heavy atoms of the sugar moiety and nucleobase from the terminal nucleotides. In particular, to attach Alexa Fluor 594 to rA<sub>4</sub>, we aligned the respective atoms from the terminal nucleotide at the 5' end of the rA<sub>4</sub> fragment with the respective atoms from the terminal nucleotide at the 3' end from the dA<sub>2</sub>-dye fragment. For Alexa Fluor 488, the same alignment was performed but at the 3' end of the rA<sub>4</sub> fragment and at the 5' end from the dA<sub>2</sub>-dye fragment. The RMSD cutoff for heavy atom distances was set to 0.8 Å. If the RMSD value was below the cutoff, we searched for clashing heavy atoms within a pair distance of 2.0 Å. If no clashing atoms were detected, then the dye molecules and the rA<sub>4</sub> fragment were assembled such that all atoms from the dA<sub>2</sub> fragment and terminal oxygens of rA<sub>4</sub> were excluded. In this way, we sampled a library of the FRET dyes mapped onto rA<sub>4</sub> fragments with 10000 conformations for each fluorophore, Alexa Fluor 594 and Alexa Fluor 488.

The fragment libraries used here for the fluorescent dyes contain only a short nucleic acid segment. Attractive interactions between dyes and nucleic acids are thus limited to the terminal bases. With more distant bases, only steric interactions are considered. However, MD simulations can result in excessive sticking of the dye to the nucleic acid with some force fields.<sup>21</sup> By contrast, experimentalists tend to exclude dyes that stick to the attached RNA or proteins, based on measurements for instance of the fluorescence anisotropy decay.<sup>21,40</sup> The fragment-based approach to modeling dyes taken here and by Grotz et al.<sup>21</sup> takes advantage of this strategy, being aimed at the modeling of experiments with nonsticky dyes.

**2.5. Ensemble Reweighting Using BioEn.** We refined the HCG ensembles of the ssRNA polymers investigated here against experimental SAXS or single-molecule FRET data by reweighting using BioEn.<sup>37,38</sup> We used uniform reference weights  $w_{0,i} = \text{const.}$  for the unbiased ensembles produced by HCG. The reference weights of the individual chains were then minimally adjusted such that the ensemble average better agrees with the experimental observable, while making sure that the refined ensemble was well-defined and converged. The confidence parameter  $\theta$  in BioEn<sup>37,38</sup> was chosen by L-curve

analysis. As a measure of the extent of reweighting, we used the Kullback–Leibler divergence  $S_{\text{KL}} = \sum_i w_i \ln(w_i/w_{0,i})$  between the reference weights  $w_{0,i}$  and the refined weights  $w_i$ , both normalized,  $\sum_i w_{0,i} = \sum_i w_i = 1$ . In addition, we inspected the cumulative distribution function (CDF) of rank-ordered weights  $w_i$ . A rapid initial rise indicates that few ensemble members carry a large fraction of the weight, which in turn indicates poor overlap between reference and refined ensembles.

**2.6. Analysis of ssRNA Conformations.** The python packages Barnaba,<sup>55</sup> MDTraj,<sup>56</sup> and MDAAnalysis<sup>57,58</sup> were used to perform analyses of the ssRNA conformations.

A cluster analysis was performed using the Barnaba software<sup>55</sup> exemplary of the rA<sub>4</sub> fragment conformations as sampled in 100 ns MD simulation trajectory. The sampled conformations within 10000 frames were assigned to 6 different clusters. In the cluster analysis, we first calculated the g-vectors describing the relative positions of the nucleotide pairs in each structure. We then performed a principal component analysis (PCA) of the g-vectors, projecting the data points onto the plane of the first and second principal component axes. The clustering was performed via a Barnaba wrapper of the DBSCAN function from the scikit-learn package. Here, we set the minimum distance for nearest neighbors to  $\text{eps} = 0.35$  and the minimum number of samples per cluster to 50.

**Quantification of Base Stacking.** We used the Barnaba software<sup>55</sup> to screen the assembled structures for stacked bases. For each base, we quantified the stacking by a factor  $q_{\text{stacking}}$ . For nucleobases not involved in any stack, we set  $q_{\text{stacking}} = 0$ . For stacks of  $n_{\text{stacked}} = 2$  and 3, we set  $q_{\text{stacking}} = 0.25c$  and  $q_{\text{stacking}} = 0.5c$ , respectively, where  $c = n/(n - 1)$  with  $n$  the total number of nucleobases in the ssRNA. For longer consecutive stacks with  $n_{\text{stacked}} \geq 4$  bases, we set  $q_{\text{stacking}} = c - c/(n_{\text{stacked}} - 1)$ . We then used  $q_{\text{stacking}}$  to color the bases in the structural visualizations (Figure 2).

**2.7. Calculation of Experimental Observables. FRET.** We calculated FRET efficiencies for the rA<sub>19</sub> ssRNA ensembles obtained by HCG with explicit dyes attached at the 5' and 3' ends. The interdye distance  $r$  was calculated as the geometric distance between the central oxygen atoms of the two FRET dye labels,<sup>21</sup> as determined using MDAAnalysis.<sup>57,58</sup> For the orientational factor  $\kappa^2$  in the Förster theory, we considered three models that differed in their assumptions on the dye dynamics. A similar approach has been employed before.<sup>59</sup>

In model 1, we set  $\kappa^2 = 2/3$ ,<sup>40,60</sup> assuming implicitly that dye rotation is isotropic and fast<sup>61,62</sup> compared to the fluorescence lifetime of the donor, which in the absence of the acceptor is  $\tau_D \approx 4$  ns.<sup>21</sup> The transfer efficiency  $E$  of an individual ssRNA conformation labeled with fluorophores at each end was then calculated as

$$E = \frac{1}{1 + (r/R_0)^6} \quad (1)$$

The Förster radius  $R_0$  was set to the experimentally determined value of  $R_0 = 5.4$  nm.<sup>21</sup> In model 2, we assumed also the dye linker dynamics to be fast and accordingly attached  $\approx 20$  conformers for the dye pairs to a given ssRNA conformation, averaged the interdye distance  $r$  over these conformers, and then calculated the FRET intensity according to eq 1 with the average  $r$ . By contrast, in model 3, we assumed the dye dynamics to be slow. Accordingly, we determined both  $r$  and  $\kappa^2$

explicitly for each dye-labeled ensemble member. We calculated  $\kappa^2$  as

$$\kappa^2 = (\hat{\mu}_D \cdot \hat{\mu}_A - 3(\hat{r} \cdot \hat{\mu}_A)(\hat{r} \cdot \hat{\mu}_D))^2 \quad (2)$$

where  $\hat{\mu}_D$  and  $\hat{\mu}_A$  are unit vectors in the direction of the transition dipole moments of donor and acceptor, respectively, and  $\hat{r}$  is a unit vector pointing in the direction between the central oxygen atoms of the two dyes. We then calculated the rate of energy transfer as  $k_T = (3/2)\kappa^2 k_D (R_0/r)^6$  and the FRET efficiency of each ssRNA conformation in the ensemble as<sup>63</sup>

$$E = \frac{1}{1 + 1/(k_T \tau_D)} \quad (3)$$

For all three models, the efficiency  $E$  was then averaged over the ssRNA conformations in the ensemble with their respective weights. The three models can be considered as extremes with respect to the assumed dye dynamics. Importantly, in all models, we assumed the ssRNA dynamics to be slow compared to the fluorescence lifetime  $\tau_D$ .

For model 2, we mapped FRET labels onto the full-length rA<sub>19</sub> grown for analysis with model 1, with labels integrated in the models at the fragment level. Particularly, we randomly picked an rA<sub>19</sub> conformation  $i$  and attempted to simultaneously replace both labels with randomly picked label conformations of Alexa Fluor 594  $j$  and Alexa Fluor 488  $k$ . Here, we followed the procedure for the heavy atom alignment and clash search as described above for the dye mapping on the fragments. In case of a steric clash or if the RMSD exceeded 0.8 Å in the alignment, both dye conformations were discarded, and a new pair of dyes was drawn. We attempted to replace the dye conformations 1000 times for each of the 10000 randomly drawn rA<sub>19</sub> conformations. The normalized acceptance rate for dye replacements determined for conformation  $i$  was then used as the weight for  $\langle r_i \rangle$  for each conformation in the ensemble.

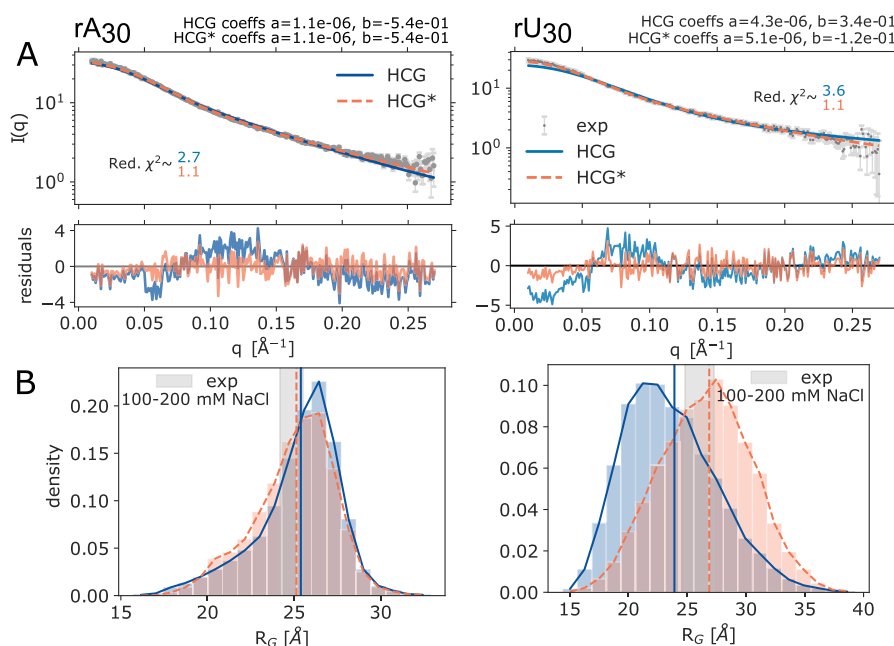
**SAXS.** For each ensemble member  $i$  of either the rA<sub>30</sub>, rU<sub>30</sub>, or rUCAAUC HCG ensembles, we calculated the SAXS scattering intensity  $I_i(q)$  at scattering vector  $q$  using Crystol<sup>64</sup> following ref 14. The calculated scattering intensities  $I_i(q)$  with normalized weights  $w_i$  (i.e.,  $\sum_i w_i = 1$ ) were averaged over the ensemble as  $I_{\text{sim}}(q) = \sum_i w_i I_i(q)$ . In the limit of  $q \rightarrow 0$ , the Guinier approximation becomes exact,  $I_i(q) \approx I_{0,i} \exp(-q^2 R_{G,i}^2/3)$ , where  $I_{0,i}$  is the intensity at  $q = 0$ , and  $R_{G,i}$  is the radius of gyration of ensemble member  $i$ . Accordingly, we calculated  $R_{G,i}$  from the slope of  $\ln I_i(q)$  with respect to  $q^2$  at  $q = 0$  as

$$R_{G,i}^2 = -3 \left. \frac{d \ln I_i(q)}{d(q^2)} \right|_{q=0} \quad (4)$$

We evaluated the slope as a numerical first difference. With  $I_i(q) \approx I_{0,i} \exp(-q^2 R_{G,i}^2/3)$  at small  $q$ ,  $I_{\text{sim}}(q) = \sum_i w_i I_i(q)$ , and  $I_{0,i}$  being nearly constant, we determined the root-mean-square (RMS)  $R_G$  by averaging over the ensemble

$$R_G^2 = \frac{\sum_i w_i I_{0,i} R_{G,i}^2}{\sum_i w_i I_{0,i}} \approx \sum_i w_i R_{G,i}^2 \quad (5)$$

We related the calculated intensity  $I_{\text{sim}}(q)$  to the measured intensity  $I(q)$  by performing a least-squares fit of  $I(q) = a I_{\text{sim}}(q) + b$  with an intensity scaling factor  $a$  and a constant background correction factor  $b$  as fit parameters. To assess the



**Figure 3.** Comparison of the rA<sub>30</sub> and rU<sub>30</sub> HCG structural ensemble to SAXS measurements<sup>18</sup> (left and right columns, respectively). (A) Top: Experimental SAXS profiles measured at 100 mM NaCl in gray and the average profile calculated using Crystol<sup>64</sup> for the unrefined HCG ensembles (blue), 10000 structures each, and the refined HCG\* ensembles (orange) with weights for  $\theta = 100$  and  $\theta = 46$  for rA<sub>30</sub> and rU<sub>30</sub>, respectively. Intensity scale factors  $a$  and background correction constants  $b$  determined by least-squares fitting are shown in the plots. Bottom: Residuals. (B) Distribution of  $R_G$  in the unrefined HCG and reweighted HCG\* ensembles (blue and orange, respectively). Vertical lines indicate the RMS  $R_G$  value of the HCG ensemble (solid blue) and the weighted RMS  $R_G$  value (HCG\*, dotted orange). The gray shaded area highlights the area spanned by the  $R_G$  value inferred from the SAXS profiles measured at 100 mM and 200 mM NaCl including the error range.

quality of the fit, we calculated the reduced chi-squared normalized by the number  $M$  of data points

$$\chi^2 = \frac{1}{M} \sum_{j=1}^M \frac{|aI_{\text{sim}}(q_j) + b - I(q_j)|^2}{\sigma_j^2} \quad (6)$$

with  $\sigma_j$  the reported experimental standard error of  $I(q_j)$ . Values of  $\chi^2 \lesssim 1$  indicate agreement within the experimental uncertainty. We used an implementation of BioEn specific for SAXS data that fits nuisance parameters globally for the ensemble average during the refinement, with the code available at [https://github.com/bio-phys/SAXS\\_BioEn/](https://github.com/bio-phys/SAXS_BioEn/). In brief, we performed an initial fit of the intensity scale factor  $a$  and background correction  $b$ , which were then updated using the refined weights until convergence was achieved.

We assessed the quality of the ensembles with the  $\chi^2$  statistic for the squared residuals and the hplusminus statistic for the sign-order of the residuals, calculating  $p$ -values for both tests individually and in combination.<sup>65</sup>

### 3. RESULTS AND DISCUSSION

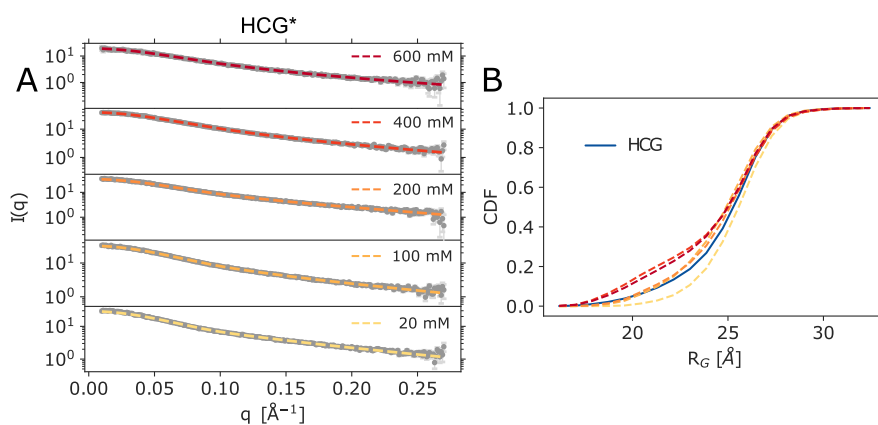
#### HCG Produces Broad Structural Ensembles of ssRNA.

We used HCG to sample structural ensembles of ssRNA polymers with four different sequences: rA<sub>30</sub>, rA<sub>19</sub>, rU<sub>30</sub>, and rUCAUC. For all four systems, we observed a combination of extended and compact conformations, as shown representatively in Figures 2 and S3. Compactness is associated with kinks in the ssRNA backbone (see, e.g., bottom center rA<sub>30</sub> structure in Figure 2). In the more extended structures, the chains retained features of the A-form helix (e.g., rA<sub>30</sub> top left in Figure 2). In particular, we observed stretches of continuously stacked nucleobases.

The ssRNA structure in the HCG ensembles depends on the nucleotide sequence. Whereas the poly purine rA<sub>30</sub> tends to form relatively straight segments of stacked adenines, poly pyrimidine rU<sub>30</sub> is visually rather distorted with stretches of unstacked uridines (Figure 2). A stacking analysis using Barnaba showed about four fewer stacks in rU<sub>30</sub> than rA<sub>30</sub> on average (Figure S4A). Here, a single stack was defined as two nucleobases with particular distances and orientation to each other.<sup>55</sup> We further looked at consecutively stacked nucleobases, which we defined as four or more stacked nucleobases in a row and colored the nucleobases accordingly (Figures 2 and S4B). Key to retaining base stacking in the fragment assembly was the inclusion of the atoms of the nucleobase in the RMSD alignment, which ensured that the relative base–base orientation of the fragments was retained in the HCG assembly (Figure 1). The observed sequence dependence is in line with the behavior previously reported for the ssRNA polymers investigated here.<sup>14,16,18,21,66</sup>

The structure in long ssRNA chains is reflected in the fragment libraries used for HCG. We clustered the rA<sub>4</sub> fragment library according to their structure. In the largest clusters, stacked A-form like conformations dominate, either with perfectly stacked nucleobases (cluster 0 with 42%) or with nucleotide 4 inverted (cluster 1 with 39%; see Supplementary Figure S1). NMR studies support the presence of a substantial fraction of A-form like conformations for short single stranded RNA fragments.<sup>9,10,67</sup> The next largest clusters 2 and 3 are sparsely populated ( $\approx 1\%$ ), containing structures with A3 unstacked and A4 inverted (cluster 2), and all bases being unstacked (cluster 3).

HCG assembly also largely preserves the distribution of torsion angles in the fragment libraries, as shown for rA<sub>19</sub> in Figure S5. In the HCG assembly, the central two nucleotides at



**Figure 4.** SAXS measurements of  $rA_{30}$  at different salt concentrations. (A) SAXS profiles from the  $rA_{30}$  HCG ensemble refined against experimental profiles measured at 20, 100, 200, 400, and 600 mM NaCl. HCG\* is shown in yellow to dark red. Experimental profiles are shown in gray; errors are shown in light gray. (B) Cumulative distribution of  $R_G$  as predicted for HCG using Crysol<sup>64</sup> in blue and the weighted distributions using the refined weights.

positions 2 and 3 of the tetrameric fragments were retained. Figure S5 compares the distributions of the backbone torsions averaged over the 19 bases in each A2 and A3 as sampled in  $rA_{19}$  chains to the respective distributions in the fragments. First, we found that the distributions are essentially independent of the position in the HCG assembly, as would be expected for a long homopolymer. Second, we found that HCG largely retained the torsion angle distributions of the fragments. However, we observed that some populations were altered or completely vanished. The small differences reflect in part actual incompatibilities with longer chains yet also choices in HCG, in particular of the atoms to align and of the clash criteria. Differences as, e.g., in  $\alpha$  for A3 or  $\epsilon$  and  $\zeta$  for both A2 and A3 may be amplified by the different nature of the preceding nucleotide at the 5' position and the following nucleotide at the 3' position of A2 or A3, respectively. In particular, in the MD fragments these nucleotides were attached to terminal nucleotides, which may impact the distribution of torsional angles at the P - OS' bond ( $\alpha$ ), C3' - O3' bond ( $\epsilon$ ), and the O3' - P(+1) bond ( $\zeta$ ).

The recovery of the torsional distributions after assembly suggests that the local and global structural features observed in the assembled full-length chain arose from the local structure sampled in the fragments, as found before for tau K18.<sup>40</sup> In particular,  $rA_4$  fragments sampled in the DESRES force fields mostly exhibited an A-form helix-like conformation with a considerable population of conformations with one or two of the 4 nucleotides being unstacked (see Figure S1 clusters 1 and 3). In turn, this resulted in either pseudo A-form helix-like populations as well as populations of kinked or looplike structures for ssRNA polymers sampled with HCG, with some of the chains even featuring patterns with bulging nucleobases (Figure 2 and Figure S3).

**ssRNA from HCG Reproduces SAXS Data.** We compared the calculated SAXS intensity profiles obtained by averaging across the  $rA_{30}$  and  $rU_{30}$  HCG ensembles to experimental profiles measured at 100 mM NaCl<sup>18</sup> (Figure 3A). The  $rA_{30}$  and  $rU_{30}$  ensembles were assembled from  $rA_4$  and  $rGUUU$  fragment libraries, respectively (see Methods). Overall, the agreement was good, with reduced  $\chi^2$  errors (mean-squared residuals divided by the experimental error) of 2.7 and 3.6, respectively. However, the residuals revealed small but systematic deviations for both polymers. In the relatively

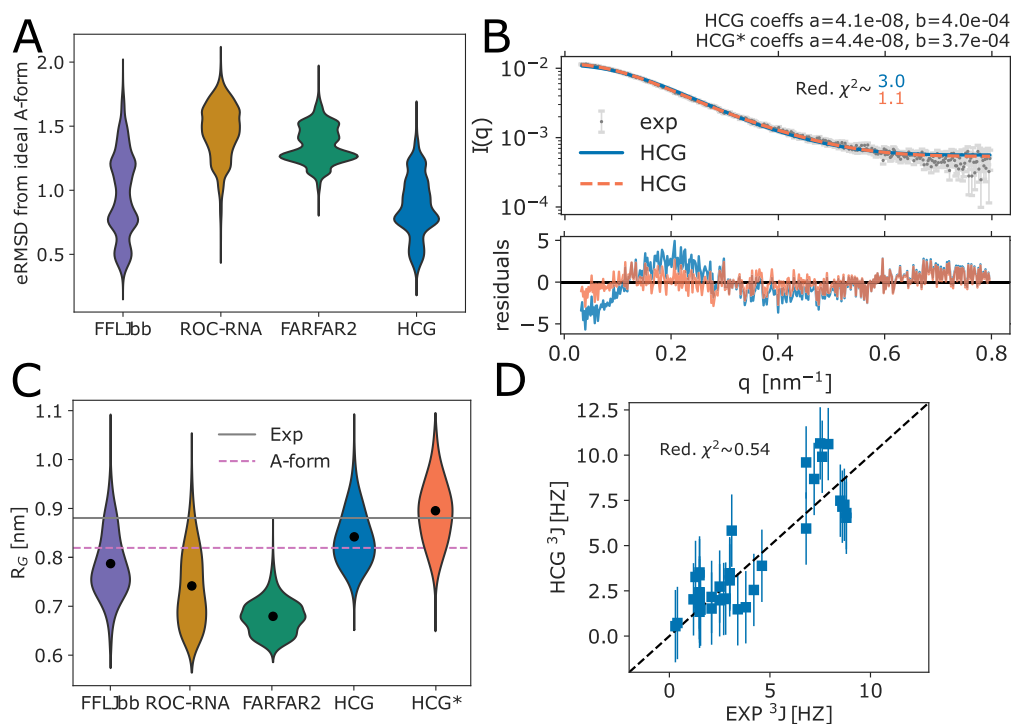
featureless intensity profiles, the residuals pointed to somewhat too extended structures for  $rA_{30}$  and too compact structures for  $rU_{30}$ .

Considering the fact that HCG does not account for long-range electrostatic interactions and salt screening effects beyond the scale of the fragments, the agreement between measured and calculated SAXS intensities at  $\sim 100$  mM NaCl is remarkably good. Solvent and, in particular, ions affect the global structure of the negatively charged nucleic acids polymers.<sup>9,11,16,17,21,24,68</sup> At the high concentrations of the SAXS experiments, interchain interactions may also be relevant.<sup>18,69</sup>

**Gentle Ensemble Reweighting Further Improves Agreement with SAXS Data.** We refined the HCG ensembles of  $rA_{30}$  and  $rU_{30}$  by performing BioEn reweighting<sup>37,38</sup> against the experimental scattering profile measured at 100 mM. Using a rather gentle bias, we adjusted the weights of the ensemble members to agree with the experimental profile with reduced  $\chi^2$  values of  $\approx 1.1$  for both polymers (HCG\* ensemble in Figures 3A, S6 light orange, and S7).

We also found good agreement between the measured and calculated values of the radius of gyration,  $R_G$ . We calculated the RMS  $R_G = \langle R_{G,i}^2 \rangle^{1/2}$  as an average over members  $i$  of the ensemble with their respective weights. For  $rA_{30}$ , the RMS average  $R_G$  over the HCG ensemble was within the uncertainty of the  $R_G$  measured by SAXS at 100 and 200 mM NaCl;<sup>18</sup> for  $rU_{30}$ , it was just below the expected range (Figure 3B). In this range, salt concentration was found to have only a small effect on  $R_G$ .<sup>18,69</sup> Reweighting by BioEn to match SAXS intensities  $I(q)$  also improved the agreement of the calculated and measured  $R_G$  values. Overall, HCG captured the global dimensions of  $rA_{30}$  and  $rU_{30}$ , and a gentle BioEn reweighting resulted in near-perfect agreement at higher NaCl concentrations for  $rA_{30}$  (see below).

**Dependence on Fragment Library.** For homopolymeric ssRNA, using a homopolymeric fragment (here:  $rA_4$ ) is a natural choice that facilitates the combination of fragments into longer chains, as the base overhangs at the 5' and 3' ends of the fragments to be combined are then identical. By contrast, to assemble generic ssRNA sequences by HCG, it is advantageous to fix the base at the 5' end to minimize the number of required fragment sequences. For the tetramer fragments used here, with sequence rGXYZ, we then have only



**Figure 5.** Characterization of structural ensembles of heteropolymeric rUCAUC sampled with MD and FARFAR2 from ref 14 and HCG. (A) Distribution of the eRMSD to ideal A-form. (B) Top: Average SAXS profile calculated for the HCG ensemble before (blue) and after ensemble refinement (HCG\* with weights for  $\theta = 100$  in orange) fitted to the experimental profile<sup>14</sup> (orange and gray, respectively). The intensity scale factors  $a$  and the background correction constants  $b$  as calculated by least-squares fitting are shown in the plot. Bottom: Residuals calculated for scattering profiles. (C) Distribution of  $R_G$  values in sampled ensembles, the refined HCG\* ensemble, experimental average as determined from SAXS, and for typical A-form. (D) Correlation plot of experimentally measured  $^3J$  couplings of the backbone and sugar moiety<sup>66</sup> and calculated values for HCG. Vertical bars indicate the estimated uncertainty of  $\pm 2$  Hz (one standard error) in calculating the  $^3J$ -couplings from RNA structures using approximate Karplus relations.<sup>13,66</sup>

$4^3 = 64$  sequences to consider. Here, having built both  $rA_4$  and  $rGA_3$  libraries, we can directly compare the two approaches to identify possible issues in HCG. We found that the distributions of  $R_G$  for  $rA_{30}$  chains grown by HCG with  $rA_4$  fragments and with  $rGA_3$  fragments show negligible differences (Kullback–Leibler divergence  $S_{KL}(rGA_3||rA_4) = 0.01$ ; see Figure S8), indicating that both libraries produce very similar results.

**$rA_{30}$  HCG Ensemble Reproduces SAXS Profiles Measured at Different Salt Concentrations after Gentle BioEn Refinement.** In their study on the salt dependence of ssRNA  $rA_{30}$  and  $rU_{30}$ , Plumridge et al.<sup>18</sup> have shown that the global structure of these highly charged polymers depends on (i) the concentration of ions in the solvent and (ii) the ion type. Here, we compared the  $rA_{30}$  HCG ensemble, grown from  $rA_4$  fragments simulated at 150 mM NaCl, to their SAXS profiles measured at 20, 100, 200, 400, and 600 mM NaCl concentrations. Despite the fact that we grew the polymer via HCG without taking into account long-range electrostatics beyond the fragment level, the HCG profile matched experimental SAXS profiles recorded at different salt concentrations reasonably well even without reweighting (Figures S9–S13A blue). For reference, we least-squares fitted the unrefined SAXS profile of the  $rU_{30}$  produced by HCG to the salt dependent experimental SAXS profiles of  $rA_{30}$ , adjusting only the intensity scale factors  $a$  and the background corrections  $b$ . For all salt concentrations, the HCG profiles for  $rA_{30}$  gave a better fit to the  $rA_{30}$  experiments than the HCG profiles for  $rU_{30}$ , albeit with only small differences in  $\chi^2$  (Figure

S14 and Supporting Table S1). Plumridge et al.<sup>18</sup> reported on the rather small differences in the global shape found for both polymers for salt concentrations  $> 20$  mM. For 100 mM NaCl, for instance, we found the experimental profiles of  $rU_{30}$  and  $rA_{30}$  to agree with reduced  $\chi^2 \approx 0.91$  after fitting  $a$  and  $b$  (Supporting Table S1).

BioEn reweighting of the HCG ensemble against the scattering profiles measured at the respective salt concentration established nearly perfect agreement with the experimental profiles (Figure 4A, Figures S9–S13A orange). For each salt concentration, we chose a set of weights for the regularization parameter  $\theta = 100$  resulting in reduced  $\chi^2 < 2$ . According to the L-curve analysis, with Kullback–Leibler divergences close to zero and the cumulative distribution functions (CDF) of rank-ordered weights staying close to uniform reference weights, all BioEn reweightings placed a rather gentle bias on the initial ensemble (Figure S6).

The reweighted RMS  $R_G$  values for the  $rA_{30}$  ensemble were shifted toward the experimental values for each concentration (Figures S9–S13C). The shape of the  $R_G$  distributions was minimally modified when we applied the refined weights for 20, 100, and 200 mM (Figures 4B left column and S9, S10, and S11C). In fact, for 100 and 200 mM NaCl, the RNA conformations of  $rA_{30}$  with  $R_G < 20$  Å lost weight against conformations with  $20$  Å  $< R_G < 25$  Å (Figures 3B left column and S11C). For the highest salt concentrations of 400 and 600 mM NaCl, a distinct shoulder developed in the reweighted  $R_G$  distribution at  $R_G \approx 20$  Å (Figures S12 and S13C). The diminished role of electrostatic repulsion between ssRNA



phosphate groups at high salt concentration may explain this trend to compaction.

We further assessed the quality of the  $I(q)$  fits using the hplusminus test statistic for ordered data.<sup>65</sup> Applied to the scattering intensities  $I(q)$ , it tends to pick up indications for systematic errors, e.g., as a result of deviations in the global size and shape of the ensemble members. Here, going in the HCG\* ensembles from low to high salt concentration, we found that systematic deviations at small  $q$  values decreased, as judged by the residuals and a screening of their signs. In return, we found improved  $p$ -values for the reweighted ensembles at higher salt concentrations (Figures S9–S13A and B).

Despite the overall efficient reweighting of the HCG ensemble resulting in almost perfect agreement with experiment, we emphasize that the refined ensemble lacks information about electrostatics and other interaction. For a more detailed assessment, one could perform additional MD simulations at the respective salt condition using a small subset of the models sampled in the HCG ensemble as start structure and choosing a reasonable force-field.<sup>28,43</sup> Such simulations would provide information about electrostatic interactions and the solvent layer.

**HCG Ensembles of rUCAAUC Capture SAXS and NMR Experiments.** Using the heterotetramer fragment library for HCG, we are able to grow heteropolymeric ssRNAs of an arbitrary sequence. In the following, we show results for the rUCAAUC hexamer, which has been investigated previously by experiments and MD simulations.<sup>14,66</sup> Visually, the structures appeared rather extended, albeit with populations of structures in which one or two nucleotides were unstacked, similar to what was observed by Bergonzo et al.<sup>14</sup> (Figures 2 and 5A).

Judging from the comparison to the published SAXS data,<sup>14</sup> the HCG ensemble of rUCAAUC chains captured the global dimensions without any refinement. The average SAXS profile calculated for the HCG ensemble was in good agreement with the experimental profile, with a reduced  $\chi^2$  of about 3.0 and small deviations at small  $q$  (Figure 5B). For reference, Bergonzo et al.<sup>14</sup> found profiles of similar quality in their MD simulations of full-length rUCAAUC. For the LJbb force field, their agreement with the SAXS experiments was slightly better, with  $\chi^2 \approx 2.4$  and the deviations at small  $q$  being less pronounced.

The distribution of  $R_G$  in the HCG ensemble was in line with the distribution in the MD ensemble sampled with the LJbb force field (FFLJbb) by Bergonzo et al.<sup>14</sup> (Figure 5C). The RMS  $R_G$  as calculated for the HCG ensemble was slightly closer to the experimentally determined value than that of the ensemble from full MD simulations. Interestingly, the RMS  $R_G$  is close to that of an rUCAAUC polymer in an ideal A-form helix conformation. Overall, the conformations sampled with HCG seemed to resemble a typical A-form to a larger extent than conformations sampled with the other approaches shown here, with the eRMSD from a typical A-form being smaller on average (Figure 5A).

The analysis of NMR  $^3J$  couplings of the backbone and sugar moiety revealed that overall HCG sampled local properties, as reflected in the torsion angles of the sugar moiety as well as the nucleic acid backbone, in excellent agreement with the experiments (Figure 5D). Small deviations in the calculated  $^3J$  couplings were within their predicted uncertainty ( $\approx 2$  Hz for the Karplus relation used to calculate the scalar coupling<sup>55</sup>), resulting in a reduced  $\chi^2$  value of  $\approx 0.54$ . Thus, we found our

ensemble to agree with experiment as good as the MD ensemble (LJbb force field) from Bergonzo et al.<sup>14</sup> We note that for the experimental values, Zhao et al. suggested an error of 2 Hz as well due to the deviations of measured values in multiple independent measurements (see ref 66 and Table S4), which we did not consider here.

We reweighted the SAXS profiles calculated for HCG against the experimentally measured scattering profile. Using a small bias with weights for  $\theta = 100$ , we found almost perfect agreement with the experimental profile with reduced  $\chi^2 \approx 1.1$  and  $S_{\text{KL}} \ll 1$  (Figure S15A, B). Deviations we observed for the refined profile were within the experimental error range, and only very small deviations for  $q < 0.1$  nm. The refined weights were used to calculate weighted distributions and averages for properties we analyzed here. The weighted distribution of  $R_G$  values was shifted toward larger values with the weighted RMS  $R_G$  value being in perfect agreement with the experimental value. Interestingly, we observed only small deviations from ideal A-form and scalar couplings close to experiment (Figure S15C, D).

**HCG Produces Ensembles with a Large Conformational Variability.** We observed a high diversity of global dimensions (e.g., for rU<sub>30</sub> in Figure 3, right panel). Using larger fragment sizes to prepare a fragment library, e.g., pentamers with the central trimers and the 3' terminal capping nucleotide being flexible and the 5' terminal cap fixed, would still be computationally feasible, with  $4^4 = 256$  fragments. It is interesting to speculate if we may be able to sample a higher population of structures that feature important local motifs, by sampling more local interactions within the input MD fragments.

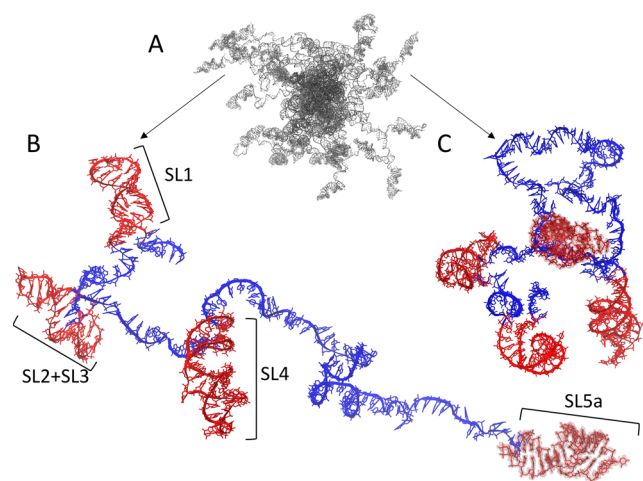
We compared the sampled structural diversity in HCG and MD ensembles in terms of pairwise RMSDs, calculated by using all heavy atoms within a polymer. For short chains (tetramers and hexamers), the distribution of the pairwise RMSD within MD ensembles was slightly larger than within the HCG ensemble. This suggests that a slightly larger conformational variability was sampled with MD (Figure S16), at least in terms of the pairwise RMSD. The pairwise RMSDs between MD and HCG were distributed around 4 nm (Figure S16, rose), similar to what we observed for two independent HCG ensembles of the same polymer (Figure S16C, dashed gray). For rA<sub>4</sub>, all distributions were shifted toward smaller pairwise RMSDs, probably due to the large population of A-form like helix conformations (Figures S1 and S16A).

Importantly, we do not know the actual extent of structural variability or the expected distribution of pairwise RMSDs for a native ssRNA ensemble. For ensemble refinement, however, it is advantageous to have a broad sampling that covers the relevant conformation space. By integrating experimental information, ensemble refinement methods such as BioEn<sup>37,38</sup> then down-weight conformations with low statistical relevance. By contrast, if the starting ensemble does not cover the relevant conformation space, conformations in this region would have to be added by biased sampling for a proper ensemble refinement.

In general, efficient comparisons of structurally heterogeneous ensembles are difficult.<sup>27</sup> Several algorithms exist to cluster ensemble members according to different properties, often accompanied by machine learning techniques. However, finding appropriate collective variables that really capture the important properties needed to display the differences between ensembles is not straightforward. Recently, a tool to compare

structural ensembles of IDPs by determining differences in distributions of local and global properties of the conformations based on a Wasserstein metric was introduced.<sup>70</sup>

**ssRNA Polymers Grown with HCG Can Be Combined with dsRNA.** Conformations sampled with HCG can be easily combined with structured dsRNA or other ordered structures. Here, we exemplarily modeled a region of the 5' UTR of the SARS-CoV-2 genomic RNA (sequence shown in Figure S2). Structured stem-loops were taken from earlier studies using FARFAR2<sup>53,71</sup> and from NMR studies.<sup>54</sup> To model the 5' UTR, we used MD trajectories of five stem-loops provided by Bottaro et al.<sup>53</sup> as input ensembles for the structured parts (see Methods). The stem-loop, highlighted in red, was then connected by single-stranded regions grown with HCG, highlighted in blue, resulting in a model containing 233 nucleotides in total. The final ensemble with about 44 different conformations features models with more extended and more compact single-stranded regions, dictating the overall global dimensions of the modeled 5' UTR (see Figure 6A). For illustration, we randomly chose two structures (Figure 6B and C, respectively).



**Figure 6.** Model of the 5' UTR region of SARS-CoV-2 genomic RNA built by HCG. (A) Ensemble overview. (B) and (C) show two representatives of extended and compact conformations, respectively, with more detail. Structures of the five stem-loop regions drawn at random from MD trajectories<sup>53</sup> are shown in red. The connecting single-stranded RNA is shown in blue. The structures are shown with atomic detail. Hydrogen atoms are omitted for clarity. The backbone atoms are shown in a cartoon representation except for the last stem-loop at the 3' end, which is highlighted as a surface. The full-length structure shown here covers 233 nucleotides.

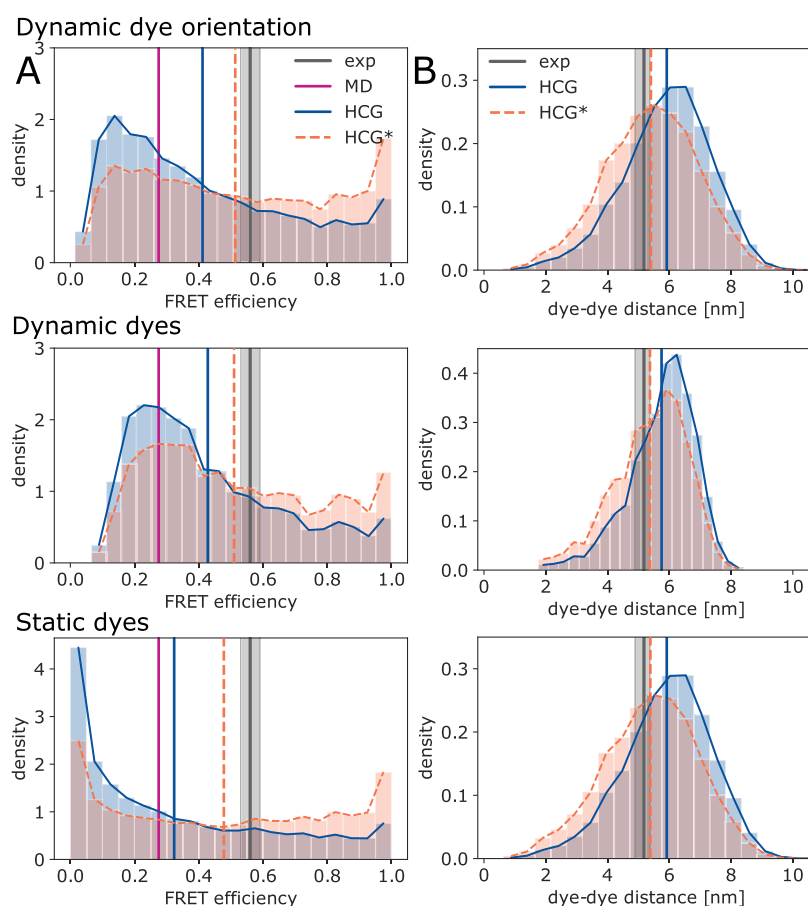
We demonstrated here a possible application of HCG to model structures of RNA molecules that combined structured and unstructured regions, such as mRNA molecules. More generally, a similar scheme may be applied to model any kind of biomolecule featuring unstructured parts, e.g., by adding a polyA tail to mRNA. Importantly, since HCG is modular, we can either add additional assembly steps and assemble the different regions after the initial growth or grow the flexible chain with HCG directly at the structured biomolecule. In particular, such models can be used for further analysis, e.g., as an initial structure for MD simulations.

**HCG Ensembles of  $rA_{19}$  with Mapped Dyes Are Somewhat too Extended on Average as Judged by**

**Experimental FRET Efficiencies.** We calculated FRET efficiencies for the HCG ensemble of  $rA_{19}$  and compared them to the measured mean FRET efficiency  $\langle E \rangle = 0.56 \pm 0.03$  obtained in single-molecule FRET experiments at 150 mM NaCl concentration.<sup>21</sup> For model 1 with fast and isotropic averaging for the dye orientations about fixed dye positions (eq 1), we obtained a mean efficiency of  $\langle E \rangle \approx 0.41$ . In model 2, we start from the ensemble in model 1 but with multiple dye pair conformations placed onto every conformation  $i$  in the  $rA_{19}$  ensemble. By averaging the FRET efficiency over these dye pairs and their orientations, we effectively assumed dynamic dyes in model 2, which we consider to be more realistic than model 1. We observed a mean FRET efficiency of  $\langle E \rangle \approx 0.42$ . In the less realistic model 3 with fully static dyes, we fixed interdye distances and determined  $\kappa^2$  explicitly from the dye conformations (eq 3). For model 3, we obtained  $\langle E \rangle \approx 0.32$  with a high population of conformations having  $E < 0.1$ . By comparison, MD simulations of full-length  $rA_{19}$  using the same force field as in our fragment MD simulation gave FRET efficiencies of  $\langle E \rangle \approx 0.3$ ,<sup>21</sup> calculated with explicit dyes mapped onto the sampled conformers and  $\kappa^2 = 2/3$  fixed, as in our model 1. Differences we observed for FRET efficiencies calculated from the MD and the HCG ensemble using model 1 may indicate that the MD simulation was too short, with 7  $\mu$ s of sampling in aggregate. Alternatively, we may have a favorable compensation of errors in chain growth by accounting primarily for the local structure.

Using BioEn,<sup>37,38</sup> we then gently reweighted the ensembles of models 1, 2, and 3 to match the experimental mean FRET efficiency. For model 1 and model 2, we obtained reduced  $\chi^2 \approx 1.4$  for a BioEn confidence parameter  $\theta = 40$ , and for model 3, we obtained  $\chi^2 \approx 2.3$  for  $\theta = 60$  (see Figure 7A orange top row, middle row, and bottom row, respectively). The ensemble refinement assigned higher weights to the tail of the ensemble, i.e., to more compact chains with  $E > 0.6$ . In turn, the weighted distributions and the average interdye distance were shifted to shorter distances within the range inferred from the single-molecule FRET experiment using a worm-like chain polymer model (see Figure 7B). Here, either the structures were more compact, the mapped dyes featured less extended linkers, and/or the mapped dyes pointed toward each other (Figure S3). This shift in population toward more compact structures is qualitatively consistent with what we found in the BioEn reweighting for  $rA_{30}$  according to the SAXS data (Figure 3B). However, the shift there was considerably smaller, as the  $R_G$  value had already agreed with the measurements within the uncertainty. We note that the  $r$  distributions in the HCG\* ensembles for models 1 and 3 are nearly identical (Figure 7B). For model 2, the distance distribution in HCG was narrower, and for HCG\*, the peak of the distribution was slightly shifted toward larger interdye distances. A small shoulder at around 4 nm dye–dye distance, present in HCG\* for all models, was more pronounced.

The shape of the reweighted distributions of the FRET efficiency may indicate slight overfitting. However, judging from the L-curve analysis and the CDF of rank-ordered weights (Figure S17, orange, dark green, and dark red), the set of weights we chose seemed to impose a rather gentle bias with  $S_{KL} < 0.2$  for all three models. An important point to consider is how to properly perform the ensemble reweighting for polymers with attached labels.<sup>39,72</sup> In approaches such as FRETpredict,<sup>73</sup> dyes are placed onto proteins using a rotamer library approach (RLA) to predict FRET efficiencies with



**Figure 7.** HCG ensembles of  $rA_{19}$ , compared to single-molecule FRET experiments. (A) Distribution of FRET efficiencies in the HCG ensemble (blue) and in the reweighted HCG\* ensemble (orange) calculated with model 1 (top) with dynamic dye orientations and  $\kappa^2 = 2/3$ , model 2 (middle) with dynamic dyes and  $\kappa^2 = 2/3$ , and model 3 (bottom) with static dyes. For HCG\*, we chose refined weights for  $\theta = 40$  (models 1 and 2) and  $\theta = 60$  (model 3). Vertical lines indicate the mean FRET efficiency measured in experiment (black, with gray shading indicating  $\pm$  SEM), sampled in MD simulations of full-length  $rA_{19}$  with the same force field as used here to build fragment libraries<sup>21</sup> (magenta), and calculated for the HCG (blue) and HCG\* ensembles (orange). (B) Distributions of the interdy distance as determined for the HCG ensemble (blue) and the reweighted HCG\* ensemble (orange) with models 1 (top), 2 (middle), and 3 (bottom). The mean distances for experiment (black), MD simulation (magenta), HCG (blue), and HCG\* (orange) are shown as vertical solid and dotted lines. The experimental mean interdy distance was inferred from experimental single-molecule FRET efficiencies using a worm-like-chain model for the distance probability density function.<sup>21</sup>

individual statistical weights. In the present study, we reweighted the whole molecule, i.e., polymer chain plus the attached fluorophore molecules. As an alternative, one could reweight the chain and dye separately.

**HCG Compared to MD Simulations of  $rU_{40}$  Using the Anton Supercomputer.** We compared the distributions of the radius of gyration  $R_G$  and the  $OS'-O3'$  distance in an HCG ensemble of  $rU_{40}$  to those sampled in  $\sim 100 \mu s$  MD simulation runs of full-length  $rU_{40}$  with the DESRES (DES-Amber0.9) force field<sup>23</sup> at 0.05, 0.1, 0.5, and 1.0 M NaCl and a newer, modified version DES-Amber3.2<sup>74</sup> at 0.05, 0.1, 0.2, 0.4, and 0.5 M NaCl. For both force fields,  $rU_{40}$  transitioned between an extended state and a more compact state with folded-back conformations.<sup>23,74</sup> The population of the compact state increased with increasing salt concentration. For the  $rU_{40}$  polymers assembled by HCG from MD fragments sampled at 0.15 M NaCl with the DESRES force field, the distributions of  $R_G$  and the  $OS'-O3'$  distance are intermediate between the distributions in full MD simulations with the DES-Amber0.9<sup>23</sup> and DES-Amber3.2<sup>74</sup> at the closest NaCl concentrations below (0.1 M) and above (1.0 M, 0.2, and 0.4 M, respectively; see Figure S18). Overall, the HCG ensemble covers a range of  $R_G$

values similar to that of the MD simulations. Together with our results for  $rU_{30}$  compared to the experimental SAXS data<sup>18</sup> (see Figure 3, right column), we conclude that HCG performs well for polymeric  $rU_n$  chains.

The HCG ensemble for  $rU_{40}$  also compares well to the experiments without adjustments. Chen et al.<sup>16</sup> have determined mean end-to-end distances of  $\approx 66 \text{ \AA}$  and  $\approx 64 \text{ \AA}$  for 100 and 200 mM NaCl, respectively, from FRET measurements. The mean end-to-end distance in the HCG ensemble of  $rU_{40}$  at 150 mM NaCl is  $\approx 67 \text{ \AA}$ , very close to the experimental values.

The consistency of HCG and full MD simulations for  $rU_{40}$  with state-of-the-art force fields is reassuring because there are important differences between the two, even if they are conducted with the same force field. Whereas HCG uses RNA fragment libraries built by MD simulations, only steric interactions are considered between distant fragments so that ssRNA structures folding back onto themselves have no stabilizing interactions. Structures such as the stem-loops in the SARS-CoV-2 5' UTR (Figure 6) can be included in HCG as fragments, as shown here. However, folded-back structures of ssRNA can also be an artifact of the MD simulation force field,

which was a major driver for the development of HCG fragment-based methods, e.g., for the modeling of fluorophore labeled ssRNA.<sup>21</sup> In terms of sampling efficiency, HCG makes it possible to build models of arbitrary sequence with prebuilt fragment libraries, given the limited four-letter alphabet of RNA sequences. For proteins, where sampling of full-length chains is feasible in principle, HCG was found to provide substantially larger coverage of configuration space than MD simulations on an  $\approx 2 \mu\text{s}$  time scale.<sup>43</sup>

#### 4. CONCLUSIONS

Single-stranded RNA appears prominently in many cellular regulation processes, e.g., in mRNA and its poly-A tail but also in loops and linkers. The structural modeling of a flexible nucleic acid with unpaired nucleobases poses formidable challenges. Here, we showed that hierarchical chain growth, previously introduced for disordered proteins,<sup>43</sup> can be used to produce structural ensembles of ssRNA with atomic detail, starting from fragments sampled in MD simulations. The resulting structural ensembles feature highly diverse conformations (Figures 2 and S3), in good agreement with NMR experiments probing the local structure (Figure 5D). Also SAXS and FRET experiments probing the global structure are reproduced well. Overall, we found the HCG ensembles to agree with experiments about as well or better than MD simulations of full-length ssRNA (Figures 3–5, 7).

HCG relies on a number of simplifying assumptions. Most importantly, it assumes that the relevant local structure of disordered biopolymers is sampled properly in short fragments and that these fragments can be assembled into full length chains subject to only steric interactions. In particular, in its simplest form, HCG does not account for long-range electrostatic interactions. It is therefore remarkable that we obtained excellent agreement for  $r_{A_{30}}$  SAXS data over a wide range of salt concentrations, in particular for 100–200 mM NaCl (Figure 4).

The computational efficiency of HCG makes it possible to construct large ensembles with diverse conformations, sampling also significant populations of rare but relevant conformations. This broad coverage of conformation space enables ensemble reweighting schemes to match a wide range of experiments within expected uncertainties (Figures 3–5, 7, S6, S7, S9–S13, S15A and B, and S17).

HCG is implemented as a Monte Carlo chain growth algorithm with a well-defined ensemble and partition function.<sup>43</sup> Therefore, HCG can easily be combined with other Monte Carlo sampling techniques, e.g., to perform importance sampling as in the reweighted hierarchical chain growth (RHCG).<sup>40</sup> In RHCG, one uses a fragment library that is refined against experimental data prior to fragment assembly to improve the sampling of local properties in the grown full-length ensemble. An exciting perspective is to adopt an RHCG-like sampling scheme to include information on interfragment interactions during chain growth or to grow loop structures. This task may be turned into a machine learning problem. Methods based on artificial intelligence (AI) have been proven to reliably predict tertiary structure of folded double-stranded and also single-stranded RNA.<sup>30,71,75,76</sup> Query sequences that require modeling of both structured and disordered regions may be excellent targets for AI-guided applications of Monte Carlo techniques. We have shown that HCG is suited to model segments of mRNA that feature structured and unstructured regions (see Figure 6). HCG in

combination with machine learning approaches could prove useful for modeling more complicated mRNA or long-noncoding RNA (lncRNA) with internal short disordered loops. To improve the grown structures, one can include experimentally derived information<sup>40</sup> and information from secondary structure prediction tools. Fragment libraries for 3D structures of RNA secondary structure motifs<sup>77</sup> can be used as input for HCG of more complex RNA folds. One could also use coarse-grained RNA simulation models<sup>78</sup> to build fragment libraries for the assembly of large RNA structures.

HCG can also be combined with MD simulations to gain insight on inter- and intramolecular interactions and the dynamics of ssRNA. In previous work, we have shown that the conformations of IDPs sampled with HCG are well-suited as starting structures for parallel but independent MD simulations with atomic detail.<sup>43</sup> Similarly, this could be done with the ssRNA conformations as modeled here, either the fully flexible single chains or molecules with structured and flexible regions. Starting from a multitude of reasonable initial structures will facilitate exploration of the relevant conformational space and dynamics.

#### ■ ASSOCIATED CONTENT

##### Data Availability Statement

An implementation of BioEn specific for SAXS data that fits nuisance parameters globally for the ensemble average during the refinement was used, with the code available at [https://github.com/bio-phys/SAXS\\_BioEn/](https://github.com/bio-phys/SAXS_BioEn/). The RNA structure libraries (i.e., the MD fragment library as well as exemplary structures from the HCG ensembles discussed in this work) are available at <https://zenodo.org/record/8369324>. The HCG code to assemble ssRNA to the hierarchical chain growth is available at the GitHub repository <https://github.com/bio-phys/hierarchical-chain-growth/>.

##### SI Supporting Information

The Supporting Information is available free of charge at <https://pubs.acs.org/doi/10.1021/acs.jctc.3c01049>.

Supporting table of SAXS fit parameters. Supporting figures of Cluster analysis of the MD trajectory of a  $r_{A_4}$  fragment; sequence of the SARS-CoV-2 5' UTR as modeled here; snapshots of  $r_{A_{19}}$  models from HCG with attached FRET labels; stacking analysis of poly A and poly U 30mers from HCG; backbone dihedral angle distributions sampled in  $r_{A_4}$  MD trajectories and  $r_{A_{19}}$  HCG models; BioEn reweighting of  $r_{A_{30}}$  HCG ensemble against SAXS data at different salt concentrations; BioEn reweighting of  $r_{U_{30}}$  HCG ensemble against experimental SAXS data; distribution of the radius of gyration for HCG  $r_{A_{30}}$  ensembles from different fragment sequences; refinement of  $r_{A_{30}}$  ensemble against SAXS data at 0.1, 0.02, 0.2, 0.4, and 0.6 M NaCl; comparison of SAXS intensities calculated for  $r_{A_{30}}$  and  $r_{U_{30}}$  against experiments for  $r_{A_{30}}$  at different salt conditions; BioEn reweighting of  $r_{U\text{-CAAUC}}$  HCG ensemble against experimental SAXS data; CDF of pairwise RMSDs in (RE)MD ensembles and HCG ensembles of ssRNA polymers with different lengths; BioEn reweighting of  $r_{A_{19}}$  HCG ensemble against experimental FRET data; comparison of HCG to earlier MD simulations of full-length  $r_{U_{40}}$  on Anton (PDF)

## AUTHOR INFORMATION

### Corresponding Author

Gerhard Hummer – Department of Theoretical Biophysics, Max Planck Institute of Biophysics, 60438 Frankfurt am Main, Germany; Institute for Biophysics, Goethe University, 60438 Frankfurt am Main, Germany; [orcid.org/0000-0001-7768-746X](https://orcid.org/0000-0001-7768-746X); Phone: +49 69 6303-2501; Email: [gerhard.hummer@biophys.mpg.de](mailto:gerhard.hummer@biophys.mpg.de)

### Authors

Lisa M. Pietrek – Department of Theoretical Biophysics, Max Planck Institute of Biophysics, 60438 Frankfurt am Main, Germany; [orcid.org/0000-0002-4494-4041](https://orcid.org/0000-0002-4494-4041)

Lukas S. Stelzl – Faculty of Biology, Johannes Gutenberg University Mainz, 55128 Mainz, Germany; KOMET 1, Institute of Physics, Johannes Gutenberg University Mainz, 55099 Mainz, Germany; Institute of Molecular Biology (IMB), 55128 Mainz, Germany; [orcid.org/0000-0002-5348-0277](https://orcid.org/0000-0002-5348-0277)

Complete contact information is available at: <https://pubs.acs.org/10.1021/acs.jctc.3c01049>

### Funding

This project was funded by the Deutsche Forschungsgemeinschaft (DFG) project number 161793742 (CRC 902: Molecular Principles of RNA Based Regulation) and by the Max Planck Society. Open access funded by Max Planck Society.

### Notes

The authors declare no competing financial interest.

## ACKNOWLEDGMENTS

We thank Kara Grotz for sharing her data and for insightful discussions. We thank Mark Nüesch and Benjamin Schuler for comments and detailed exchange on the analysis of the single-molecule FRET data. We thank Jürgen Köfinger for fruitful discussions and for sharing his experience with the analysis of scattering data, ensemble refinement via Bayesian Inference of Ensembles, and the hplusminus analysis.

## REFERENCES

- (1) Cech, T. R.; Steitz, J. A. The noncoding RNA revolution—trashing old rules to forge new ones. *Cell* **2014**, *157*, 77–94.
- (2) Micura, R.; Höbartner, C. Fundamental studies of functional nucleic acids aptamers riboswitches ribozymes and DNazymes. *Chem. Soc. Rev.* **2020**, *49*, 7331–7353.
- (3) Passmore, L. A.; Collier, J. Roles of mRNA poly(A) tails in regulation of eukaryotic gene expression. *Nat. Rev. Mol. Cell Biol.* **2022**, *23*, 93–106.
- (4) Leulliot, N.; Varani, G. Current topics in RNA-protein recognition: control of specificity and biological function through induced fit and conformational capture. *Biochemistry* **2001**, *40*, 7947–7956.
- (5) 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.
- (6) Rouskin, S.; Zubradt, M.; Washietl, S.; Kellis, M.; Weissman, J. S. Genome-wide probing of RNA structure reveals active unfolding of mRNA structures in vivo. *Nature* **2014**, *505*, 701–705.
- (7) Sahin, U.; Karikó, K.; Türeci, Ö. mRNA-based therapeutics—developing a new class of drugs. *Nat. Rev. Drug Discovery* **2014**, *13*, 759–780.
- (8) Sreeramulu, S.; Richter, C.; Berg, H.; Wirtz Martin, M. A.; Ceylan, B.; Matzel, T.; Adam, J.; Altincekic, N.; Azzaoui, K.; Bains, J.

K.; Blommers, M. J. J.; Ferner, J.; Fürtig, B.; Göbel, M.; Grün, J. T.; Hengesbach, M.; Hohmann, K. F.; Hymon, D.; Knezic, B.; Martins, J. N.; Mertinkus, K. R.; Niesteruk, A.; Peter, S. A.; Pyper, D. J.; Qureshi, N. S.; Scheffer, U.; Schlundt, A.; Schnieders, R.; Stiral, E.; Sudakov, A.; Tröster, A.; Vögele, J.; Wacker, A.; Weigand, J. E.; Wimmer-Bartoschek, J.; Wöhnert, J.; Schwalbe, H. Exploring the druggability of conserved RNA regulatory elements in the SARS-CoV-2 genome. *Angew. Chem. Int. Ed.* **2021**, *60*, 19191–19200.

(9) Eichhorn, C. D.; Feng, J.; Suddala, K. C.; Walter, N. G.; Brooks, C. L.; Al-Hashimi, H. M. Unraveling the structural complexity in a single-stranded RNA tail: Implications for efficient ligand binding in the prequeuosine riboswitch. *Nucleic Acids Res.* **2012**, *40*, 1345–1355.

(10) Tubbs, J. D.; Condon, D. E.; Kennedy, S. D.; Hauser, M.; Bevilacqua, P. C.; Turner, D. H. The nuclear magnetic resonance of CCCC RNA reveals a right-handed helix, and revised parameters for AMBER force field torsions improve structural predictions from molecular dynamics. *Biochemistry* **2013**, *52*, 996–1010.

(11) Condon, D. E.; Kennedy, S. D.; Mort, B. C.; Kierzek, R.; Yildirim, I.; Turner, D. H. Stacking in RNA: NMR of four tetramers benchmark molecular dynamics. *J. Chem. Theory Comput.* **2015**, *11*, 2729–2742.

(12) Sponer, J.; Bussi, G.; Krepl, M.; Banás, P.; Bottaro, S.; Cunha, R. A.; Gil-Ley, A.; Pinamonti, G.; Poblete, S.; Jureacka, P.; Walter, N. G.; Otyepka, M. RNA structural dynamics as captured by molecular simulations: A comprehensive overview. *Chem. Rev.* **2018**, *118*, 4177–4338.

(13) Bottaro, S.; Bussi, G.; Kennedy, S. D.; Turner, D. H.; Lindorff-Larsen, K. Conformational ensembles of RNA oligonucleotides from integrating NMR and molecular simulations. *Sci. Adv.* **2018**, *4*, earr8521.

(14) Bergonzo, C.; Grishaev, A.; Bottaro, S. Conformational heterogeneity of UCAAUC RNA oligonucleotide from molecular dynamics simulations, SAXS, and NMR experiments. *RNA* **2022**, *28*, 937–946.

(15) Liu, B.; Shi, H.; Al-Hashimi, H. M. Developments in solution-state NMR yield broader and deeper views of the dynamic ensembles of nucleic acids. *Curr. Opin. Struct. Biol.* **2021**, *70*, 16–25.

(16) Chen, H.; Meisburger, S. P.; Pabit, S. A.; Sutton, J. L.; Webb, W. W.; Pollack, L. Ionic strength-dependent persistence lengths of single-stranded RNA and DNA. *Proc. Natl. Acad. Sci. U.S.A.* **2012**, *109*, 799–804.

(17) Plumridge, A.; Meisburger, S. P.; Andresen, K.; Pollack, L. The impact of base stacking on the conformations and electrostatics of single-stranded DNA. *Nucleic Acids Res.* **2017**, *45*, 3932–3943.

(18) Plumridge, A.; Andresen, K.; Pollack, L. Visualizing disordered single-stranded RNA: connecting sequence, structure, and electrostatics. *J. Am. Chem. Soc.* **2020**, *142*, 109–119.

(19) Schuler, B.; Eaton, W. A. Protein folding studied by single-molecule FRET. *Curr. Opin. Struct. Biol.* **2008**, *18*, 16–26.

(20) Zheng, W.; Borgia, A.; Buholzer, K.; Grishaev, A.; Schuler, B.; Best, R. B. Probing the action of chemical denaturant on an intrinsically disordered protein by simulation and experiment. *J. Am. Chem. Soc.* **2016**, *138*, 11702–11713.

(21) Grotz, K. K.; Nuuesch, M. F.; Holmstrom, E. D.; Heinz, M.; Stelzl, L. S.; Schuler, B.; Hummer, G. Dispersion correction alleviates dye stacking of single-stranded DNA and RNA in simulations of single-molecule fluorescence experiments. *J. Phys. Chem. B* **2018**, *122*, 11626–11639.

(22) Kührová, P.; Mlýnský, V.; Zgarbová, M.; Krepl, M.; Bussi, G.; Best, R. B.; Otyepka, M.; Šponer, J.; Banás, P. Improving the performance of the amber RNA force field by tuning the hydrogen-bonding interactions. *J. Chem. Theory Comput.* **2019**, *15*, 3288–3305.

(23) Tan, D.; Piana, S.; Dirks, R. M.; Shaw, D. E. RNA force field with accuracy comparable to state-of-the-art protein force fields. *Proc. Natl. Acad. Sci. U.S.A.* **2018**, *115*, E1346–E1355.

(24) Grotz, K. K.; Cruz-León, S.; Schwierz, N. Optimized magnesium force field parameters for biomolecular simulations with accurate solvation, ion-binding, and water-exchange properties. *J. Chem. Theory Comput.* **2021**, *17*, 2530–2540.

- (25) Cruz-León, S.; Vanderlinden, W.; Müller, P.; Forster, T.; Staudt, G.; Lin, Y.-Y.; Lipfert, J.; Schwierz, N. Twisting DNA by salt. *Nucleic Acids Res.* **2022**, *50*, 5726–5738.
- (26) Best, R. B.; Hummer, G. Optimized molecular dynamics force fields applied to the helix-coil transition of polypeptides. *J. Phys. Chem. B* **2009**, *113*, 9004–9015.
- (27) Ahmed, M. C.; Skaanning, L. K.; Jussupow, A.; Newcombe, E. A.; Kragelund, B. B.; Camilloni, C.; Langkilde, A. E.; Lindorff-Larsen, K. Refinement of  $\alpha$ -synuclein ensembles against SAXS data: comparison of force fields and methods. *Front. Mol. Biosci.* **2021**, *8*, 654333.
- (28) Pietrek, L. M.; Stelzl, L. S.; Hummer, G. Structural ensembles of disordered proteins from hierarchical chain growth and simulation. *Curr. Opin. Struct. Biol.* **2023**, *78*, 102501.
- (29) Das, R.; Karanicolas, J.; Baker, D. Atomic accuracy in predicting and designing noncanonical RNA structure. *Nat. Methods* **2010**, *7*, 291–294.
- (30) Watkins, A. M.; Rangan, R.; Das, R. FARFAR2: Improved de novo rosetta prediction of complex global RNA folds. *Structure* **2020**, *28*, 963–976.e6.
- (31) Chojnowski, G.; Zaborowski, R.; Magnus, M.; Mukherjee, S.; Bujnicki, J. M. RNA 3D structure modeling by fragment assembly with small-angle X-ray scattering restraints. *Bioinformatics* **2023**, *39*, btad527.
- (32) Różycki, B.; Kim, Y. C.; Hummer, G. SAXS ensemble refinement of ESCRT-III CHMP3 conformational transitions. *Structure* **2011**, *19*, 109–116.
- (33) Boomsma, W.; Tian, P. F.; Frellsen, J.; Ferkinghoff-Borg, J.; Hamelryck, T.; Lindorff-Larsen, K.; Vendruscolo, M. Equilibrium simulations of proteins using molecular fragment replacement and NMR chemical shifts. *Proc. Natl. Acad. Sci. U.S.A.* **2014**, *111*, 13852–13857.
- (34) Larsen, A. H.; Wang, Y.; Bottaro, S.; Grudin, S.; Arleth, L.; Lindorff-Larsen, K. Combining molecular dynamics simulations with small-angle X-ray and neutron scattering data to study multi-domain proteins in solution. *PLoS Comput. Biol.* **2020**, *16*, No. e1007870.
- (35) Bottaro, S.; Bengtson, T.; Lindorff-Larsen, K. In *Struct. Bioinforma. Methods protoc.*; Gáspári, Z., Ed.; Springer US: New York, NY, 2020; pp 219–240.
- (36) Borgia, A.; Zheng, W.; Buholzer, K.; Borgia, M. B.; Schüler, A.; Hofmann, H.; Soranno, A.; Nettels, D.; Gast, K.; Grishaev, A.; et al. Consistent view of polypeptide chain expansion in chemical denaturants from multiple experimental methods. *J. Am. Chem. Soc.* **2016**, *138*, 11714–11726.
- (37) Hummer, G.; Köfinger, J. Bayesian ensemble refinement by replica simulations and reweighting. *J. Chem. Phys.* **2015**, *143*, 243150.
- (38) Köfinger, J.; Stelzl, L. S.; Reuter, K.; Allande, C.; Reichel, K.; Hummer, G. Efficient ensemble refinement by reweighting. *J. Chem. Theory Comput.* **2019**, *15*, 3390–3401.
- (39) Reichel, K.; Stelzl, L. S.; Koefinger, J.; Hummer, G. Precision DEER distances from spin-label ensemble refinement. *J. Phys. Chem. Lett.* **2018**, *9*, 5748.
- (40) Stelzl, L. S.; Pietrek, L. M.; Holla, A.; Oroz, J.; Sikora, M.; Köfinger, J.; Schuler, B.; Zweckstetter, M.; Hummer, G. Global structure of the intrinsically disordered protein tau emerges from its local structure. *JACS Au* **2022**, *2*, 673–686.
- (41) Plumridge, A.; Meisburger, S. P.; Pollack, L. Visualizing single-stranded nucleic acids in solution. *Nucleic Acids Res.* **2017**, *45*, No. e66.
- (42) Shi, H.; Rangadurai, A.; Assi, H. A.; 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*, 5531.
- (43) Pietrek, L. M.; Stelzl, L. S.; Hummer, G. Hierarchical ensembles of intrinsically disordered proteins at atomic resolution in molecular dynamics simulations. *J. Chem. Theory Comput.* **2020**, *16*, 725–737.
- (44) Salomon-Ferrer, R.; Case, D. A.; Walker, R. C. An overview of the Amber biomolecular simulation package. *Wiley Interdiscip. Rev. Comput. Mol. Sci.* **2013**, *3*, 198–210.
- (45) Piana, S.; Donchev, A. G.; Robustelli, P.; Shaw, D. E. Water dispersion interactions strongly influence simulated structural properties of disordered protein states. *J. Phys. Chem. B* **2015**, *119*, 5113–5123.
- (46) Joung, I. S.; Cheatham, T. E. Determination of alkali and halide monovalent ion parameters for use in explicitly solvated biomolecular simulations. *J. Phys. Chem. B* **2008**, *112*, 9020–9041.
- (47) Abraham, M. J.; Murtola, T.; Schulz, R.; Páll, S.; Smith, J. C.; Hess, B.; Lindahl, E. GROMACS: High performance molecular simulations through multi-level parallelism from laptops to supercomputers. *SoftwareX* **2015**, *1–2*, 19–25.
- (48) Hess, P.; LINCS, B. A parallel linear constraint solver for molecular simulation. *J. Chem. Theory Comput.* **2008**, *4*, 116–122.
- (49) Parrinello, M.; Rahman, A. Polymorphic transitions in single crystals: A new molecular dynamics method. *J. Appl. Phys.* **1981**, *52*, 7182–7190.
- (50) Darden, T.; York, D.; Pedersen, L. Particle mesh Ewald: An  $N \log(N)$  method for Ewald sums in large systems. *J. Chem. Phys.* **1993**, *98*, 10089–10092.
- (51) Patriksson, A.; van der Spoel, D. A temperature predictor for parallel tempering simulations. *Phys. Chem. Chem. Phys.* **2008**, *10*, 2073–2077.
- (52) Cordero, B.; Gómez, V.; Platero-Prats, A. E.; Revés, M.; Echeverría, J.; Cremades, E.; Barragán, F.; Alvarez, S. Covalent radii revisited. *Dalton Trans* **2008**, 2832.
- (53) Bottaro, S.; Bussi, G.; Lindorff-Larsen, K. Conformational ensembles of noncoding elements in the SARS-CoV-2 genome from molecular dynamics simulations. *J. Am. Chem. Soc.* **2021**, *143*, 8333–8343.
- (54) Wacker, A.; Weigand, J. E.; Akabayov, S. R.; Altincekic, N.; Bains, J. K.; Banijamali, E.; Binas, 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*, 12415–12435.
- (55) Bottaro, S.; Bussi, G.; Pinamonti, G.; Reiber, S.; Boomsma, W.; Lindorff-Larsen, K. Barnaba: software for analysis of nucleic acid structures and trajectories. *RNA* **2019**, *25*, 219–231.
- (56) 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.
- (57) Michaud-Agrawal, N.; Denning, E. J.; Woolf, T. B.; Beckstein, O. MDAAnalysis: A toolkit for the analysis of molecular dynamics simulations. *J. Comput. Chem.* **2011**, *32*, 2319–2327.
- (58) Gowers, R.; Linke, M.; Barnoud, J.; Reddy, T.; Melo, M.; Seyler, S.; Domański, J.; Dotson, D.; Buchoux, S.; Kenney, I. MDAAnalysis: A python package for the rapid analysis of molecular dynamics simulations. *Proc. 15th Python Sci. Conf.*; 2016; pp 98–105.
- (59) Schuler, B.; Lipman, E. A.; Steinbach, P. J.; Kumke, M.; Eaton, W. A. Polyproline and the “spectroscopic ruler” revisited with single-molecule fluorescence. *Proc. Natl. Acad. Sci. U.S.A.* **2005**, *102*, 2754–2759.
- (60) Hellenkamp, B.; Schmid, S.; Doroshenko, O.; Opanasyuk, O.; Kühnemuth, R.; Rezaei Adariani, S.; Ambrose, B.; Aznauryan, M.; Barth, A.; Birkedal, V.; et al. Precision and accuracy of single-molecule FRET measurements—a multi-laboratory benchmark study. *Nat. Methods* **2018**, *15*, 669–676.
- (61) Best, R. B.; Merchant, K. A.; Gopich, I. V.; Schuler, B.; Bax, A.; Eaton, W. A. Effect of flexibility and cis residues in single-molecule FRET studies of polyproline. *Proc. Natl. Acad. Sci. U.S.A.* **2007**, *104*, 18964–18969.
- (62) Holmstrom, E. D.; Holla, A.; Zheng, W.; Nettels, D.; Best, R. B.; Schuler, B. Accurate transfer efficiencies, distance distributions, and ensembles of unfolded and intrinsically disordered proteins from single-molecule FRET. *Methods Enzymol* **2018**, *611*, 287–325.
- (63) Hummer, G.; Szabo, A. Dynamics of the orientational factor in fluorescence resonance energy transfer. *J. Phys. Chem. B* **2017**, *121*, 3331–3339.

(64) Svergun, D.; Barberato, C.; Koch, M. H. J. CRYSOLE – a program to evaluate X-ray solution scattering of biological macromolecules from atomic coordinates. *J. Appl. Crystallogr.* **1995**, *28*, 768–773.

(65) Köfinger, J.; Hummer, G.; Köfinger, J. Powerful statistical tests for ordered data. *ChemRxiv* **2021**, DOI: [10.26434/chemrxiv-2021-mdt29-v3](https://doi.org/10.26434/chemrxiv-2021-mdt29-v3).

(66) Zhao, J.; Kennedy, S. D.; Berger, K. D.; Turner, D. H. Nuclear magnetic resonance of single-stranded RNAs and DNAs of CAAU and UCAAUC as benchmarks for molecular dynamics simulations. *J. Chem. Theory Comput.* **2020**, *16*, 1968–1984.

(67) Sripakdeevong, P.; Kladwang, W.; Das, R. An enumerative stepwise ansatz enables atomic-accuracy RNA loop modeling. *Proc. Natl. Acad. Sci. U.S.A.* **2011**, *108*, 20573–20578.

(68) Cruz-León, S.; Grotz, K. K.; Schwierz, N. Extended magnesium and calcium force field parameters for accurate ion–nucleic acid interactions in biomolecular simulations. *J. Chem. Phys.* **2021**, *154*, 171102.

(69) Chen, Y.; Pollack, L. SAXS studies of RNA: structures, dynamics, and interactions with partners. *Wiley Interdiscip. Rev. RNA* **2016**, *7*, 512–526.

(70) González-Delgado, J.; Sagar, A.; Zanon, C.; Lindorff-Larsen, K.; Bernadó, P.; Neuvial, P.; Cortés, J. WASCO: A Wasserstein-based statistical tool to compare conformational ensembles of intrinsically disordered proteins. *J. Mol. Biol.* **2023**, *435*, 168053.

(71) Zhang, K.; Zheludev, I. N.; Hagey, R. J.; Haslecker, R.; Hou, Y. J.; Kretschnig, R.; Pintilie, G. D.; Rangan, R.; Kladwang, W.; Li, S.; et al. Cryo-EM and antisense targeting of the 28-kDa frameshift stimulation element from the SARS-CoV-2 RNA genome. *Nat. Struct. Mol. Biol.* **2021**, *28*, 747–754.

(72) Tessei, G.; Martins, J. M.; Kunze, M. B. A.; Wang, Y.; Crehuet, R.; Lindorff-Larsen, K. DEER-PREdict: Software for efficient calculation of spin-labeling EPR and NMR data from conformational ensembles. *PLOS Comput. Biol.* **2021**, *17*, No. e1008551.

(73) Montepietra, D.; Tessei, G.; Martins, J. M.; Kunze, M. B. A.; Best, R. B.; Lindorff-Larsen, K. FRETpredict: A Python package for FRET efficiency predictions using rotamer libraries. *bioRxiv* **2023**, DOI: [10.1101/2023.01.27.525885](https://doi.org/10.1101/2023.01.27.525885).

(74) Tucker, M. R.; Piana, S.; Tan, D.; LeVine, M. V.; Shaw, D. E. Development of force field parameters for the simulation of single- and double-stranded DNA molecules and DNA–protein complexes. *J. Phys. Chem. B* **2022**, *126*, 4442–4457.

(75) Townshend, R. J. L.; Eismann, S.; Watkins, A. M.; Rangan, R.; Karelina, M.; Das, R.; Dror, R. O. Geometric deep learning of RNA structure. *Science* **2021**, *373*, 1047–1051.

(76) Baek, M.; McHugh, R.; Anishchenko, I.; Jiang, H.; Baker, D.; DiMaio, F. Accurate prediction of protein–nucleic acid complexes using RoseTTAFoldNA. *Nat. Methods* **2024**, *21*, 117.

(77) Gan, H. H.; Pasquali, S.; Schlick, T. Exploring the repertoire of RNA secondary motifs using graph theory; implications for RNA design. *Nucleic Acids Res.* **2003**, *31*, 2926–2943.

(78) Cragolini, T.; Laurin, Y.; Derreumaux, P.; Pasquali, S. Coarse-grained HiRE-RNA model for ab initio RNA folding beyond simple molecules, including noncanonical and multiple base pairings. *J. Chem. Theory Comput.* **2015**, *11*, 3510–3522.