



## OPEN Telomeric G-quadruplex intermediates unveiled by complex Markov network analysis

A. Sáinz-Agost<sup>1,2</sup>, F. Falo<sup>1,2</sup> & A. Fiasconaro<sup>1,2</sup>✉

G-quadruplexes are secondary, non-canonical RNA/DNA structures formed by guanine-rich sequences assembled into four-stranded helical structures by the progressive stacking of G-Tetrads, planar arrangements of guanines stabilised by monovalent ions such as  $K^+$  or  $Na^+$ . Their stability plays a very important role in the prevention of DNA degradation, leading to the promotion or inhibition of specific biological pathways upon formation. In this work, we explore the occurrences of intermediates originating from the unfolding of these structures by using all-atom simulations, analyzing a small number of significant reaction coordinates to follow the evolution of the system by applying a mesoscopic simplification of the structures followed by two different dimensionality reduction techniques: Principal Component Analysis (PCA) and time-Independent Component Analysis (tICA). The data of the reduced trajectories are then encoded into a Complex Markov Network which, in conjunction with an Stochastic Steepest Descent, provides a hierarchical organization of the different nodes into basins of attraction. This procedure is able to reveal the main intermediates and the most relevant transitions the system undergoes in its denaturation path.

G-quadruplexes (G4s) are secondary, non-canonical structures arising from either one or multiple guanine rich DNA or RNA chains, and are very common in the human genome<sup>1,2</sup>. These conformations emerge from the gradual stacking of *G-tetrads*, planar arrangements of four guanines stabilized by Hoogsteen hydrogen bonds<sup>3–5</sup>. The stacking process is further aided by the presence of either monovalent or divalent cations along the central channel these tetrads define<sup>6</sup>, with typically  $K^+$  or  $Na^+$  being involved.

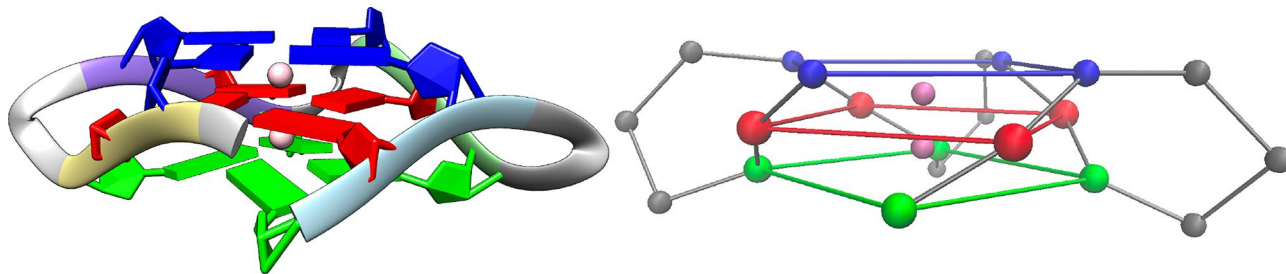
The relevance and recent interest in these structures comes from their biological function. G-quadruplexes have an important role in DNA stability<sup>7</sup> as well as in regulation processes<sup>8–10</sup>, the latter associated to either the over- or under-exposure of different binding sites of interest upon adoption of a G-quadruplex topology. The presence of these structures at the moment of DNA replication can lead to the interruption of this process, destabilizing the helicases<sup>11</sup> and, combining with other factors, can lead to epigenetic instability<sup>12</sup>. Similarly, overabundance of guanine in mRNA can lead to difficulties in protein translation in the ribosome. Several works have investigated their gene promotion and repression capabilities<sup>13,14</sup>, as well as its relevance as a potential therapeutic binding site for cancer treatment, given that the telomeric regions of the chromosomes have a tendency to form these structures<sup>15,16</sup>.

G-quadruplexes can adopt diverse topologies<sup>3</sup>, primarily influenced by two key factors: the number of single-stranded DNA/RNA chains participating in their formation (from either a single or several molecules, the latter called intermolecular G4s) and the particular twists and turns present in their backbones, dividing them into parallel (all guanine strands oriented in the same direction), antiparallel (neighboring strands oriented in opposite directions) or hybrid. Our research focuses specifically on unimolecular quadruplexes, formed by a single stranded DNA (ssDNA) chain. Within this context, we examined one of the multiple structures available: the *parallel* G-quadruplex, which has all guanine tracts are oriented in the same direction, in which the loops connect the top of a guanine track with the bottom of the next. This conformation is depicted on Fig. 1.

The mechanical stability of G-quadruplexes under force has been extensively documented<sup>6,17–19</sup>. Thus, the studies in this field have turned towards the study of their thermodynamical properties<sup>20,21</sup>. In particular, the observation and subsequent characterization of different conformations and folding intermediates of G-quadruplexes has become an object of large interest, due to the involvement of these structures in genetic regulation.<sup>22</sup>

The goal of this work is to focus into the thermal unfolding of G-quadruplexes and investigate the presence of intermediate structures in the unfolding pathway. The stability of the G4 is affected by several factors both

<sup>1</sup>Departamento de Física de la materia condensada, Universidad de Zaragoza, Zaragoza 50009, Spain. <sup>2</sup>Instituto de Biocomputación y Física de Sistemas Complejos, Universidad de Zaragoza, Zaragoza 50009, Spain. ✉email: afiascon@unizar.es



**Fig. 1.** On the left, parallel conformation sourced from the PDB (1KF1, ribbon representation in Chimera). On the right, the resulting structure after the coarse graining procedure. The different G-tetrads on the structures are marked in different colors, the loops in grey and the ions in pink. As concerns the PDB structure, the different guanine tracts are also represented in homogeneous colors: yellow (first tract, 5'-end), purple (second tract), lime (third tract) and ice blue (fourth tract, 3'-end), respectively.

external, such as ionic concentration<sup>23,24</sup>, presence of other elements in the solvent<sup>25,26</sup> and bath temperature, or internal, such as strand orientation (parallel, antiparallel or hybrid), length and structure of the loops<sup>24,27,28</sup>, and chain rigidity<sup>18</sup>.

Molecular Dynamics (MD) simulations have proven to be a valid and powerful tool in examining both the thermal and mechanical stability of biomolecular structures, also G-quadruplexes, revealing complex energy landscapes and characterizing their folding kinetics<sup>29,30</sup>. In the context of MD, we have used the *Replica Exchange Molecular Dynamics* (REMD)<sup>31,32</sup> method to study the evolution of our system through the parallel simulation of 8 copies of the human parallel G-quadruplex (PDB: 1KF1)<sup>33</sup> in a temperature range close to its reported melting point (65°), observing an unfolding event in only one of the eight replicas. To efficiently analyze the resulting high-dimensional data, a coarse-graining approach based on a consolidated mesoscopic model<sup>17,18</sup> was applied, followed by the dimensionality reduction techniques, specifically Principal Component Analysis (PCA)<sup>34,35</sup> and time-lagged Independent Component Analysis (tICA)<sup>36</sup>. These techniques allowed us to reduce the complexity of the data and retain only the most relevant information about the system. Using Complex Markov Networks (CMN)<sup>37–40</sup> combined with a stochastic steepest descent algorithm, we constructed a series of networks describing the unfolding process. These networks outlined the main conformations the system adopts as it unfolds, identifying both stable and transitory intermediates. The description provided by tICA proved to be clearer than that of PCA, sometimes identifying a temporal evolution of the denaturation pathway.

This work is structured as follows: in Section **Methods** we elaborate on the simulation setup, the different dimensionality reduction approaches, and the description and construction of the Complex Markov Networks. Section **Results** contains a discussion of the simulation findings, the effect of the presence of the cations inside the G4 structure on the unfolding dynamics, the results arising from the application of the two dimensionality reduction procedures considered, as well as the resulting CMNs describing the unfolding. The final conclusions are contained in Section **Discussion and conclusions**. Additional information regarding both methods and results can be found in the Supplementary Material associated to this publication.

## Methods

### All-atom simulations with Gromacs

The initial structure for the simulation, the intramolecular parallel DNA G-quadruplex (1KF1)<sup>33</sup>, formed by GGGTTA repeat units, as 5'-AGGGTTAGGGTTAGGGTTAGGG-3', was retrieved from the Protein Data Bank (PDB).

The simulations were carried out using the software Gromacs<sup>41,42</sup>, with the help of the force field amber modification Parmbsc0<sup>43</sup>, which has been used for both DNA<sup>44</sup> and G-quadruplex analysis<sup>6</sup> and has been proven to be the best option for our calculations<sup>45</sup>, although other options, such as OL15<sup>46</sup> or Parmbsc1<sup>47</sup>, could also be considered. Control measures to verify the stability of the complexes in normal conditions under the force field of interest were carried out, obtaining the expected results.

The structures were solvated using the Tip3P water model<sup>48</sup> in a periodic cubic box 1 nm larger than the diameter of the DNA body. Afterwards, more potassium K<sup>+</sup> ions were introduced in the simulation box to counteract the negative charges present in the DNA chains, resulting in an average potassium concentration of 0.14 M, with no negative ions added to the mix.

The energy of the solvent was minimized using a steepest descent algorithm for up to 50000 iterations, leaving the G4s frozen. This was followed by equilibration in the NVT and NPT ensembles, both lasting for 100 ps using a leap-frog integrator, with a Berendsen thermostat and a Parrinello-Rahman barostat. The electrostatic interactions were treated via the particle mesh Ewald method, establishing its cutoff, as well as the one for the Lennard-Jones interactions, at 1 nm. All simulations performed used a time step of 2 fs, and simulation data was recorded every 20 ps.

To observe the unfolding intermediates of the structures, our simulations have been performed at temperatures near the melting point of the parallel G4, being slightly higher than 65°C<sup>49</sup>. At these temperatures, the unfolding kinetics are notoriously slow.

To increase the unfolding probability of the G4s and avoid that the system may get stuck in a metastable state<sup>50</sup>, we applied in our simulations the *Replica Exchange Molecular Dynamics* (REMD) method, also called

Parallel Tempering<sup>31,32</sup>. This methodology consists in running concurrent simulations of multiple copies of the entire system (structure, solvent and ions), each of them at different temperatures. We refer to the simulation boxes containing a copy of the system at a given temperature as *replicas*. Once the simulation is started, after a sufficient time lapse (below described)<sup>51–53</sup>, attempts are made to exchange the content between replicas at different temperatures, re-scaling the velocities of the particles in the process. The probability of such exchanges is determined by a prescribed algorithm, in our case the Metropolis criterion<sup>54</sup>. This interchange of conformations at diverse temperatures aids the system in overcoming high energy barriers.

We constructed eight replicas, with assigned temperatures exceeding the nominal melting point of the structures (65°C). This adjustment was made in consideration of the ionic strength found in our simulations, which is higher than that of the experimental measurements, a fact that has been shown to increase the melting temperature  $T_m$ <sup>23</sup>. The concrete values were in the interval [343, 345, . . . , 357] K.

As refers to the exchange time between replicas, we followed a different criterion with respect to the standard procedure<sup>51–53</sup>. In fact, in order to guarantee the subsequent PCA and tICA analysis under a well determined temperature, we have used a thermalization condition after each replica exchange. This way, the minimum switching time has been set to a time interval of 200ns, time during which the energy potential can be considered thermalized, according to the relaxation of the potential energy visible in the Supplementary Material. The consequence of this choice lays in a reduced unfolding probability, but guarantees the correctness of the PCA and tICA procedure.

### Mesoscopic discretization of the model

The aim of this work is to extract the main conformations the human parallel G-quadruplex adopts along its unfolding trajectory. The method employed for the analysis relies partially in the diagonalization of different correlation matrices, whose dimension corresponds to the number of coordinates involved in the description of the all-atom trajectories of these structures, as explained in later paragraphs. In order to reduce the size of these matrices as well as to only focus on the conformations that lead to significant changes in the backbone of the G-quadruplexes, we decided to coarse-grain the structures before this analysis.

Specifically, the coarse-graining procedure, consistent with prior investigations in our research group<sup>17,18</sup>, involved the removal of water and ions from the simulation box, retaining only the DNA bases, which were subsequently reduced to unique identical beads. The positions of these beads were determined as the center of mass for the guanines and the position of the phosphates for the other nucleotides composing the loops. The transformation of gromacs coordinates into a coarse-grained system was executed using Python 3.9.18 in conjunction with the MDAnalysis package<sup>55,56</sup>. The resulting structure is depicted on the right side of Fig. 1.

### Dimensionality reduction techniques

Following the coarse-graining, our system now comprises 21 beads, each characterized by its position along the three Cartesian axes, resulting in a total of 63 degrees of freedom. Extensive research<sup>57,58</sup> has demonstrated that the majority of these degrees of freedom do not contribute significantly to our understanding of the system. Instead, the essential aspects can be effectively captured using a reduced set of reaction coordinates derived from combinations of the original ones.

In this study, we employ two distinct methods to diminish the system's dimensionality while retaining a substantial amount of relevant information. These methodologies are recognized as “Principal Component Analysis” (PCA) and “time-Independent Component Analysis.” (tICA).

#### Principal component analysis

PCA is a dimensionality reduction procedure which produces a series of orthogonal reaction coordinates from linear combinations of the input data  $x_i$ ,  $i = 1, \dots, 63$ . Its aim is to retain and explain as much of the variance of the original data as possible. This procedure, first introduced in<sup>34,35</sup>, relies on the solution of an eigenvalue problem involving the covariance matrix of the input data

$$C(0)V = \Lambda V, \quad (1)$$

where  $C(0)$  is the covariance matrix  $C(0)_{ij} \propto \sum_t x_i(t)x_j(t)$  of the normalized data (mean 0, unit variance),  $V$  contains the eigenvectors  $v_i$  by columns, and  $\Lambda$  is the diagonal eigenvalue matrix. The *principal components* are therefore obtained by projecting our original data into the different eigenvectors. A more detailed look into the derivation of the method can be found in the Supplementary Material.

The relative magnitude of the resulting eigenvalues is a measure of the proportion in which the original variance is projected onto the corresponding PC. Thus, depending on the distribution of the eigenvalues  $\lambda_i$ , a subset of coordinates  $n \leq N$  can be selected for the ongoing system description. The appropriate value for  $n$  corresponds to a large reduction in the magnitude of the eigenvalue  $\lambda_{n+1}$  when compared to  $\lambda_n$ , meaning that most of the relevant information of the system (which corresponds to the majority of the variance) resides within the first  $n$  degrees of freedom.

In our specific context, PCA is employed to eliminate extraneous degrees of freedom while retaining those that contribute to insights into conformational changes within the G4 structure. It is crucial to note that this assumption holds true only when these transitions represent the most significant alterations in the variance of the system, a condition that may not always be met. To selectively filter changes based on their kinetics and retain only those with a timescale exceeding a predefined threshold, we employ tICA, here below presented.

### time-independent component analysis

tICA is a method, first introduced in<sup>36</sup>, which produces a series of reaction coordinates from linear combinations of the input data that maximise not the variance, but the autocorrelation of the projected data between times  $t$  and  $t + \tau$ , with  $\tau$  being a time window selected by the user, known as the *lag time*. The mathematical framework of tICA is similar in its nature to that one of PCA, and is presented in the Supplementary Material.

The method relies on the resolution of a generalized eigenvalue problem involving the covariance matrix of the data and the time-lagged correlation matrix  $C(\tau)_{ij} \propto \sum_t x_i(t)x_j(t + \tau)$ , similar to the former but correlating the data at times  $t$  and  $t + \tau$ :

$$C(\tau)W = C(0)\Gamma W, \quad (2)$$

where  $\Gamma$  is the diagonal eigenvalue matrix, and  $W$  the eigenvectors matrix, with each of the different  $w_i$  vectors as columns. This general equation (Eq. (2)) is typically not solvable due to the small value of the determinant of the matrices involved, leading to numerical errors in the calculations. The AMUSE algorithm<sup>59</sup> is typically used in its place, detailed in the Supplementary Material.

It has been shown<sup>60</sup> that the eigenvalues  $\gamma_i$  correspond to the value of autocorrelation of the  $i$ -th component and that the cross-correlation between the  $i$ -th and the  $j$ -th coordinates vanishes at the lag time  $t = \tau$ . Furthermore, if we assume the autocorrelation of the signal to have an exponential decay, the associated constant is given by the lag time and the eigenvalue as:

$$t_i = -\frac{\tau}{\log |\gamma_i|} \quad (3)$$

Therefore, the value of the lag-time acts as a threshold for the time-scales detected in our system: if a certain timescale is bigger than  $\tau$  it will be detected and included in our coordinates, otherwise it will be discarded, since its autocorrelation has had time to decay to 0: it is a kinetic filter. With this property in mind, tICA should be able to distinguish important conformational changes and intermediates that PCA could not consider if they are not affected by a remarkable variance value.

### Conformational Markov network and conformational basins

Conformational Markov Networks (CMNs) are a tool that has been previously introduced to study the Free Energy Landscape of a diverse array of physical systems<sup>37–40</sup>. It consists in building a series of nodes connected via directed links.

For the construction of a CMN, each of the  $n$  reduced coordinates, obtained from either PCA or tICA, are discretized into  $m$  intervals. Every possible combination of the intervals of the coordinates that has been occupied (i.e. the trajectory of our real system has been in that precise combination at an arbitrary time  $t$ ) will constitute a *node* in our system. Therefore, we have a number of nodes  $N_{\text{nodes}} \leq m^n$ . The weight of each node is given by the number of times its associated combination of intervals is visited by the trajectory, normalized with the total number of time frames the trajectory is divided in. The links between them are directed and described by  $P_{ij}$ , denoting the transition probability from node  $j$  to node  $i$ , normalized such that  $\sum_i P_{ij} = 1$ .

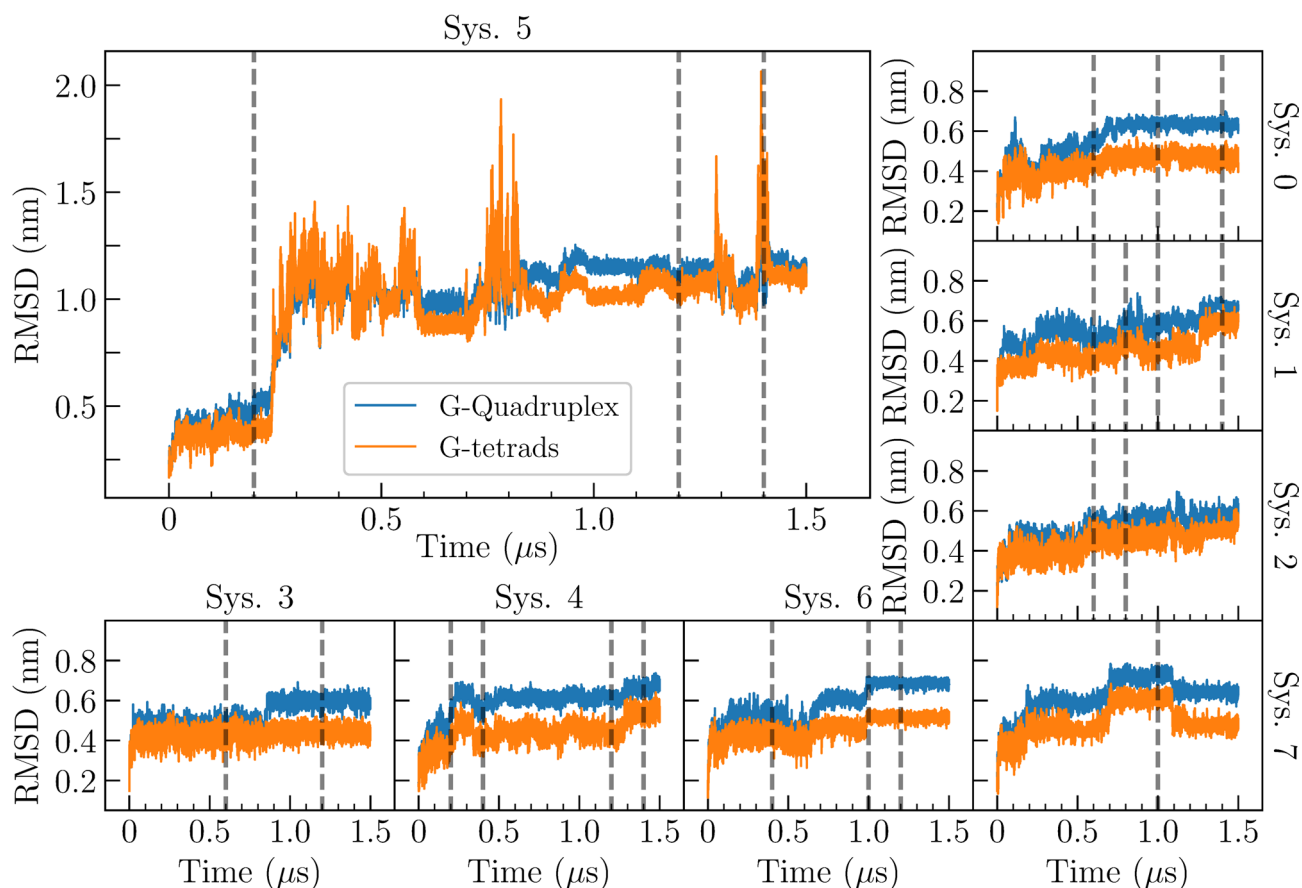
The resulting network size can become challenging to analyze depending on the chosen number of intervals  $m$ , which is directly correlated to the number of nodes. In our specific case,  $n = 3$  coordinates and  $m = 10$  intervals were chosen, resulting in up to 1000 distinct nodes. To reduce this number and eliminate potential redundancies where different nodes essentially represent the same conformation, a Stochastic Steepest Descent algorithm was employed. This algorithm, described in detail in Ref<sup>37</sup>, groups the nodes into basins, which represent different attractors of the trajectory. Additionally, a filtering based on a cutoff of the weight of the nodes has been introduced with the goal to retain only fundamental basins, and eliminating the ones less representative ( $P_i < 10^{-5}$ ).

## Results

### Direct GROMACS simulations

One of the magnitudes we investigated in our study is the *Root Mean Square Deviation* (RMSD) of the G-quadruplex structure: the difference at each time between the G4 positions and the equivalent native configuration, calculated along all replica's trajectories. This measure is the simplest way to indicate how much the structure changes during the time evolution, and it has been calculated both for the complete G4-system and the guanine tetrads only (i.e. the piled guanine structure without the external loops), the latter being the most important guideline to confirm that changes in the RMSD actually correspond to unfolding processes. In fact, a significant change in loop conformation could lead to an increase in the total RMSD, while not necessarily being associated to the denaturation of the G-quadruplex. Thus the confirmation of the unfolding relies on the G-tetrads structure only.

This information is contained in Fig. 2, that shows both the whole G-quadruplex and G-tetrads. To display the RMSD values properly, we have demultiplexed the trajectories: instead of tracking the behavior of a replica whose contents are continuously changing over time, we analyze the evolution of a specific G4 structure as it passes through between replicas. This approach prevents the misinterpretation of changes in RMSD that are actually due to exchanges between replicas as real unfolding events. We refer to these demuxed trajectories as *Systems* (Sys.), as they represent the same nucleic structure over all the replicas. The successful exchanges between replicas, when a Sys. changes temperature, are marked by gray vertical lines in Fig. 2. The RMSD curves focusing on the replicas can be found in the Supplementary Material, corresponding to Supplementary Figs. 1 & 2.



**Fig. 2.** RMSD calculated over the different starting configurations (Sys.) of the parallel G-quadruplexes. In blue, the RMSD of the whole structure. In orange, the RMSD taking into account only the guanines forming the planar arrangement. The vertical dashed lines correspond to successful exchanges between replicas. Sys. 5, since it is the one experiencing unfolding and used for the analysis, has been enhanced.

All but one of the Sys. present an RMSD that quickly stabilizes at either approximately 0.6 nm or 0.4 nm, depending either on the inclusion (complete G4, blue line in the panels) or absence (G-tetrads, orange line) of the loops in the calculation, respectively. Sys. 5 undergoes great deviations from the norm, reaching values up to 2 nm independently of the loops, thus confirming an unfolding event. These results further emphasize the correctness of the choice of the REMD method; in fact, the unfolding processes of the G-quadruplex have an average lifetime typically larger than the scope of the simulations, thus making them difficult to observe without the use of replicas at different temperatures.

The other G4 structures remain in a relatively stable configuration, though some of them present the loss of one ion from the central channel, occurrence described in the next Section [Effect of the monovalent ions in the stability of the system](#).

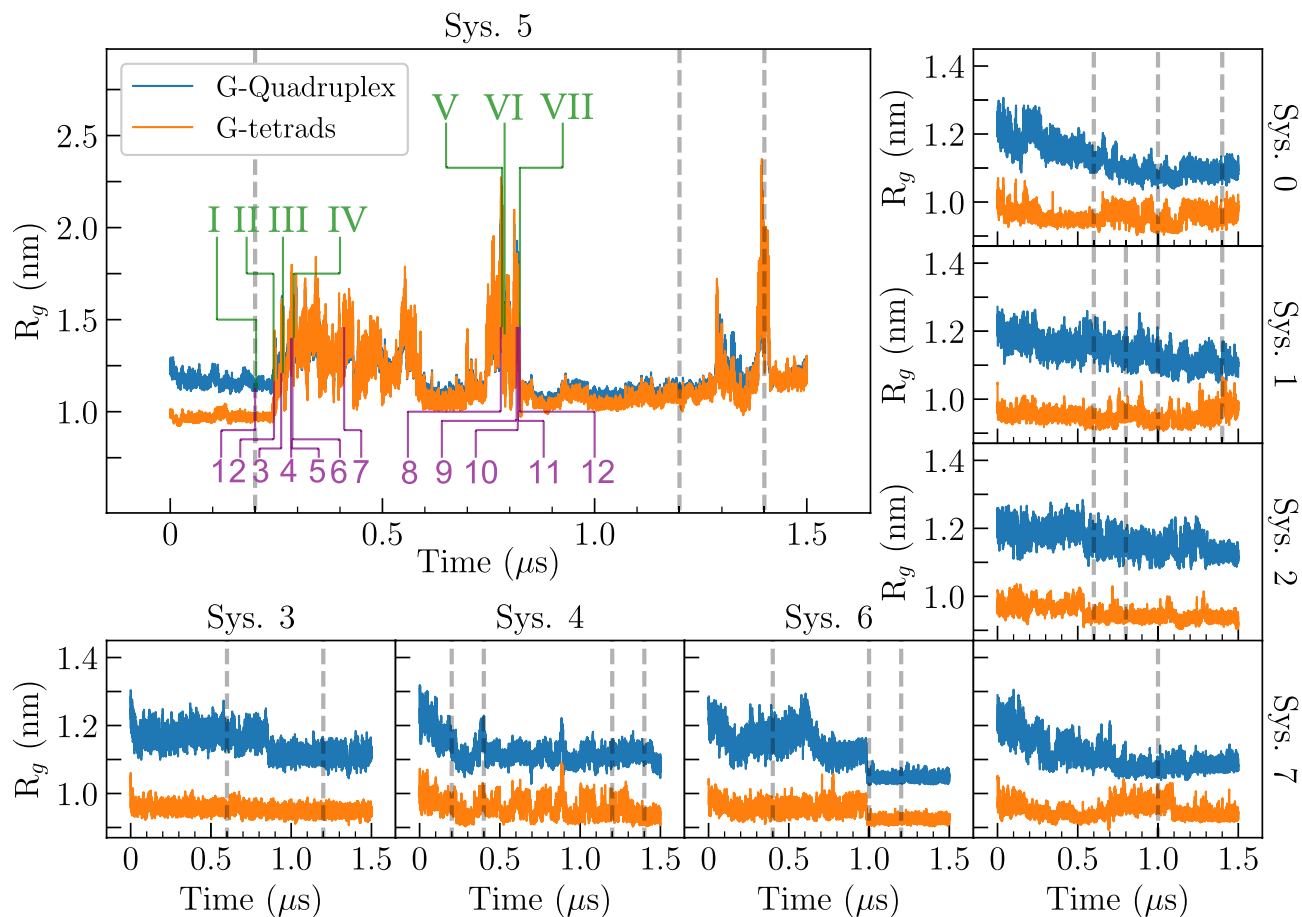
These same results can be understood through a second metric, the *radius of gyration* ( $R_g$ ), which measures the average square distance of the monomers respect to the center of mass of the structure. Fig. 3 shows the values of ( $R_g$ ), as a function of time, where we can observe its low variability in the majority of systems, while a rapid and sudden deviation with respect to the stable value occurs in Sys. 5, indicating the unfolding of the structure already seen in Fig. 2.

#### *Effect of the monovalent ions in the stability of the system*

The stability of G-quadruplex structures is highly influenced by the presence or absence of ions, typically monovalent, within the central channel<sup>6</sup>. These ions, which stably reside in the G4 because of their positive electrostatic interactions with the negatively charged G-tetrads and the spatial arrangement of guanines in the native conformation, may escape if a structural conformation with a sufficiently wide central channel occurs.

The replica analysis is once again inverted, focusing on a single structure as it traverses between replicas to prevent confusing exchanges between replicas with ion loss events.

To characterize the ions' position, two coordination numbers were employed: one for the site corresponding to the top and intermediate planes, denoted as  $Co_1$ , and the other for that of the intermediate and bottom planes, denoted as  $Co_2$ . These metrics quantify the number of bonds formed between each ion (filtered according to the closeness to the G4 structure) and the guanine bases constituting the tetrads, with values ranging from

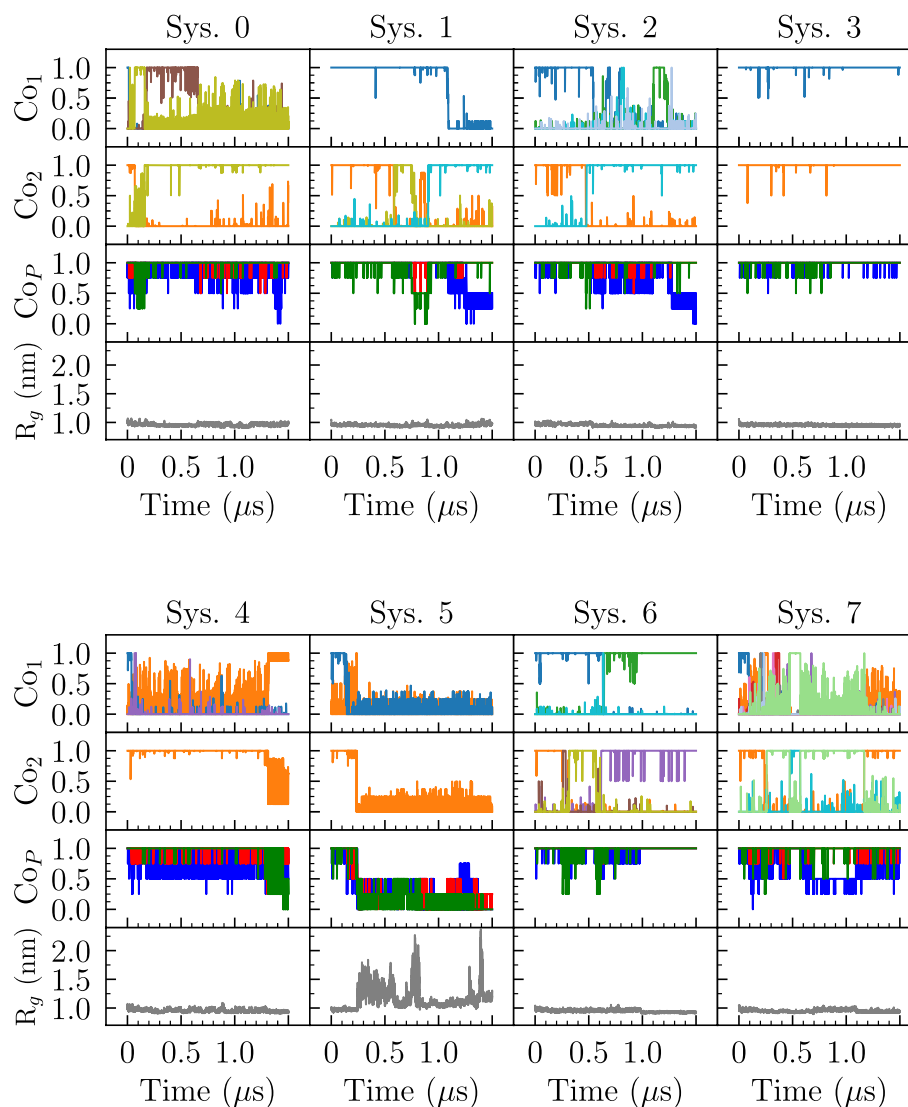


**Fig. 3.** Radius of gyration computed over the different starting configurations (Sys.) of the parallel G-quadruplex. In blue, the  $R_g$  of the whole structure. In orange, the  $R_g$  for the G-tetrads. Sys. 5 has been enhanced, since it is the one considered for posterior analysis due to the presence of unfolding. The vertical grey dashed lines correspond to successful exchanges between replicas. The purple and green colors of the numbers correspond to the states identified by PCA and tICA respectively, corresponding to the conformations reported in Figs. 10 & 11.

0.0 (indicating no bonds) to 1.0 (indicating all 8 possible bonds) in increments of 0.125. A bond is considered established if the distance between the  $O6'$  atom of guanine and the ion is less than 4.5 Å.

Since the middle tetrad is involved in the calculation of both coordination numbers, an ion exhibiting a coordination number of 1.0 in one of the two possible placements would also present a value of 0.5 in the other. A situation like this could be mistaken with the presence of two ions within the structure, one fully coordinated (hence  $Co_i = 1.0$ ) and another partially coordinated ( $Co_j \leq 1.0, j \neq i$ ). To prevent this specific confusion, once an ion achieves  $Co_i = 1.0$  in one position, its coordination number is set to 0 in the other position. Monitoring the width of the central channel itself involves defining an additional set of coordination numbers for the three guanine tetrads, referred to as  $Co_P$ . These numbers, similar to those used for ions, quantify the distances between the four guanines forming each tetrad, based on the Hoogsteen bonds present in the structure. The coordination number ranges from 0.0 to 1.0 in increments of 0.25, with a bond considered established if the distance between two neighboring guanines is less than 5.0 Å.

Fig. 4 illustrates the evolution of the two coordination numbers for the ions ( $Co_1$  and  $Co_2$  in the first two rows), along with the three coordination numbers of the G-tetrads (third row, each color represents a plane as depicted in Fig. 1: blue for the top plane, red for the middle plane and green for the bottom plane) and the radius of gyration of the G-tetrads (loops excluded). All systems, except Sys. 3, experience the loss of one or both ions at some point of their trajectories. Excluding Sys. 5, which unfolds and consequently lacks a definable central channel, the systems either recover both ions or maintain one ion while losing the other, as observed in Systems 0, 1, 2, 4 and 7. It should be noted that the recovered ions are not necessarily those initially lost; in fact all the ions in the simulation box are equivalent, thus color changes in Fig. 4 can eventually occur. The departure of an ion from the channel increases the system's susceptibility to destabilization. This phenomenon is reflected in Fig. 4, where the loss of an ion from either position in the channel leads to significant fluctuations in both the coordination number of the corresponding planes and the radius of gyration  $R_g$ . An example of this can be seen in Sys. 7 in the time window [0.75 - 1.0]  $\mu$ s, showing fluctuation in the plane coordination numbers and deviations in  $R_g$ . The same deviations can be seen in Rep. 7 in Figs. 2 & 3 during the same time interval.

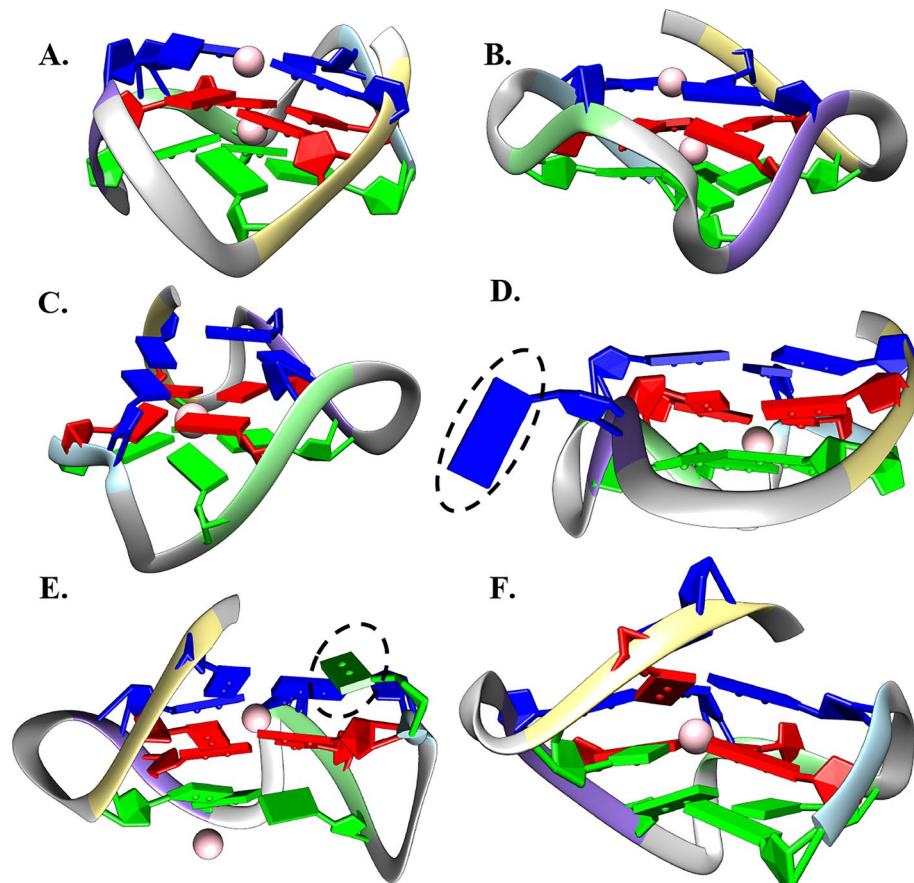


**Fig. 4.** Coordination numbers of the top ( $Co_1$ ) and bottom ( $Co_2$ ) positions available for the ions, coordination number of the three G-tetrads ( $Co_P$ , blue for the top plane, red for the middle plane and green for the bottom plane) and radius of gyration of the G-tetrads (loops excluded). Only ions which achieve a value of 1.0 in each system are represented.

Ion escape and reabsorption are facilitated by slight changes in the central channel's width. The primary causes of these deformations includes the bending of G-tetrads, the motion of single guanines unbinding/drifting away from the structure. These deformations are not limited to the top and bottom tetrads. Disruptions in the intermediate plane can also result in the transfer of an ion from one position in the ionic channel to another, as observed in Sys. 0 at  $0.2 \mu s$ , Sys. 4 at  $1.3 \mu s$  and Sys. 7 at  $0.5 \mu s$ .

Fig. 5 shows various conformations associated with the dynamics of the ions in the different systems. Subfigure A illustrates the escape of an ion from the top position, corresponding in Fig. 4 to the fall of the blue line in Sys. 1 around  $1.1 \mu s$ , which triggers a decrease in  $Co_{P1}$ . Subfigure B shows the reverse process, where an ion is integrated into the structure, associated with the increase of the olive  $Co_1$  line at around  $0.2 \mu s$  in Sys. 0 and a corresponding rise in  $Co_{P2}$ . Subfigures C and D, taken from Sys. 2 and 7, respectively, depict the separation of a guanine from the top (blue) plane after ion escape around  $1.3 \mu s$  and  $0.7 \mu s$  respectively, leading to a significant drop in the coordination number of that plane,  $Co_{P1}$ . Subfigure E shows a similar phenomenon, but with the displaced guanine in the bottom plane. This separation occurs while an ion remains bound, eventually leading to its loss around  $0.76 \mu s$  in Sys. 2, accompanied by a decrease in  $Co_{P3}$ . Lastly, Subfigure F shows a slip-stranded conformation in Sys. 4 at  $1.3 \mu s$ . The slip-stranding effect reduces the maximum possible coordination number to 0.875 and widens the central channel, facilitating the transfer of the ion from the bottom to the top position, altering the trends of the coordination numbers. This slip-stranding also causes a decrease in the coordination number  $Co_{P3}$ , with a new maximum value of  $Co_{P3, \max} = 0.75$ .

From these results, it is evident that the presence of ions significantly influences the G-quadruplex stability, generally leading to partial guanine separation and other structural changes following their loss. In Sys. 5, which



**Fig. 5.** Relevant conformations associated to ion dynamics in the simulation. **A.** Escape of an ion from the top plane. **B.** Absorption of an ion through the top plane. **C.** Deformation of the top plane after ion loss. **D.** Separation of a single guanine from the tetrad of the top plane (indicated by black dashed outline). One ion has already left the structure. **E.** Separation of a single guanine from the bottom tetrad (indicated by black dashed outline), leading to ion loss. **F.** Slip-stranded conformation leading to the transference of an ion from the bottom to the top parts of the central channel. Slip-stranded contacts circled by dashed lines.

undergoes unfolding, both ions are lost before the process begins: one is lost early in the simulation, briefly replaced by another, which is lost again around  $0.1 \mu\text{s}$ , and the other ion is lost at  $0.24 \mu\text{s}$ , with unfolding commencing around  $0.26 \mu\text{s}$ .

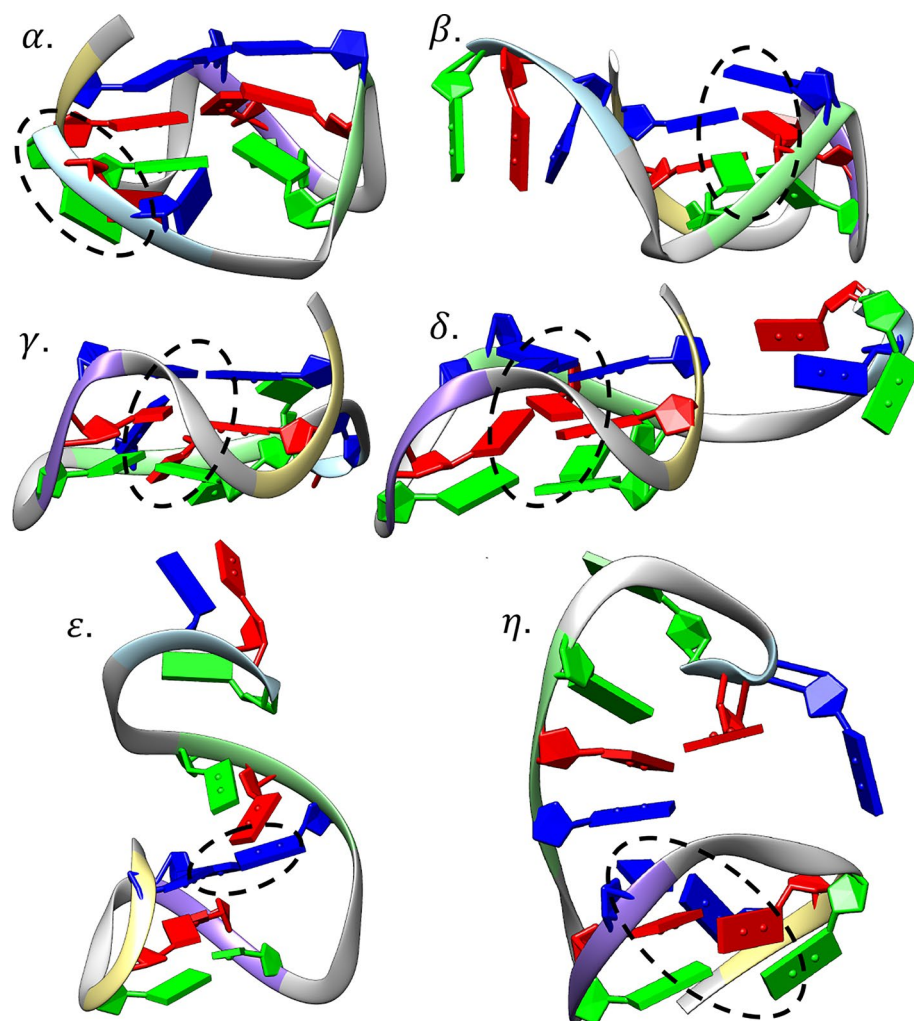
In the remaining systems, ion loss typically occurs in only one of the two available cavities. In this sense the loss of only one ion make the G4 structure less stable as it appears like a metastable state that maintains a certain stability. The G4 structure remains not completely compact, as visible in Fig. 5, where the structures are partially modified with the loss of one ion. We understand that these structures can maintain relative stability over time up to the loss of the second ion, a condition that definitely triggers the unfolding process. This observation may explain why only one system fully unfolded, as it required the loss of both ions before the unfolding process could initiate.

#### Unfolding process

The unfolding process occurs in Sys. 5, shown in Figs. 2 and 3. The identification of peculiar states in the trajectory will be largely developed with the analysis with tICA and PCA, later on presented.

In Sys. 5, the initial configuration after the replica exchange occurring at  $t = 0.2 \mu\text{s}$  corresponds to a folded G4 which has already lost one of its ions, as visible in Fig. 4. Afterwards, it quickly loses the other ion, at  $t = 0.24 \mu\text{s}$ , therefore triggering the unfolding itself, as reflected by the increase of RMSD and  $R_g$  in Figs. 2 & 3 at  $t = 0.25 \mu\text{s}$ .

Fig. 6 $\alpha$  presents a G4 conformation in which the fourth guanine tract (ice blue, see Fig. 1) separates from the structure with the remaining three segments, thus forming a G-triplex<sup>61,62</sup>, one common intermediates in G-quadruplex unfolding<sup>63,64</sup>. During the next 30 ns several attempts at refolding are made unsuccessfully. Moreover the dashed circle in Fig. 6 $\beta$  shows, in the remaining G-triplex, the so-called “slip-stranding”, *i.e.* a guanine tract that moves upwards or downwards respect to the others<sup>65,66</sup>. Other transient conformations can appear, such as G-hairpins<sup>63,67–69</sup>, *i.e.* the bond of two guanine tracts only (see Fig. 6 $\gamma$ ), or even temporary refolding events into a G-triplex (Fig. 6 $\delta$ ).



**Fig. 6.** Highlighted conformations achieved during the evolution of the system, with the dotted circles highlighting the features of interest.  $\alpha$ . Initial detachment of the fourth (icy blue) tract of guanines.  $\beta$ . Example of a G-triplex conformation with a slip-stranded effect between tracts 1 (yellow) and 3 (green).  $\gamma$ . Formation of a G-hairpin between tracts 1 (yellow) and 2 (magenta).  $\delta$ . Refolded triplex after re-attachment of tract 1.  $\epsilon$ . Stable intermediate state of G-hairpin (tracts 1 and 2) with stacking interaction of the third tract.  $\eta$ . Final state of the system, with tracts 1 and 2 in a cross hairpin conformation.

However, after the successful refolding into a triplex, the third G4 tract detaches again, rotates over itself and forms a single Hoogsteen bond with the top plane of the G-hairpin formed by the first and second tract, as shown in Fig. 6 $\epsilon$ . Remarkably, this latter hybrid between a G-triplex and a G-hairpin remains stable for approximately  $\Delta t \approx 0.6 \mu\text{s}$  of the simulation. The intense fluctuations observed in Fig. 3 in the interval  $[0.3, 0.9] \mu\text{s}$  correspond to changes in the distance between the third and fourth tracts. Eventually, the fourth strand rotates over itself and binds with the third strand through the bottom guanine (green), corresponding to the plateau between  $t = 0.6 \mu\text{s}$  and  $t = 0.75 \mu\text{s}$ . After instant  $t = 0.75 \mu\text{s}$ , the bond joining the G-hairpin with the third tract breaks. The chain elongates, pulling in the second tract and separating it from the first, breaking the G-hairpin. It is quickly reformed but with the two tracts rotated, forming a “cross-hairpin”<sup>69</sup>. During this time, the third and fourth tracts rearrange themselves into two main conformations: one characterized by the stacking interactions between both tracts leading to a high  $R_g$ , and another presenting the formation of Hoogsteen bonds between their guanines, with smaller  $R_g$ . The interchange between them forms the peaks observed around  $0.8 \mu\text{s}$ .

The final conformation of the system, stable and leading to the plateau in Fig. 3 starting at  $t = 0.8 \mu\text{s}$ , depicts the formation of a *cross hairpin* (an arrangement of guanine rich strands in a cruciform shape, as circled in Fig. 6 $\eta$ ) between the first and the second tract, with the third and fourth stacks remaining close, accompanied by the formation of transient Hoogsteen bonds between them.

### Dimensionality reduction

Both PCA and tICA are techniques quite sensible to the characteristics of the data of interest.

The application of the two methods in our trajectories has required to clean the data as follows: i. first of all, we corrected the effects of the periodic boundary conditions on the GROMACS trajectories by removing the

artificial discontinuities; ii. in order to prevent unrealistic distances of the trajectory coordinates emerging from global rotations and/or displacements of the structure as a whole, the coordinates of the structure have been rescaled by applying both a translation, superimposing the centers of mass with that of native structure, and rotations of the whole structure to recover the native orientation. A mean-square distance method, taking the initial state of the system as a reference, has been applied. iii. Finally, the mean was subtracted from each of the trajectory coordinates  $x_i$ , which have been used as input values for both PCA and tICA analysis.

#### Eigenvalues

The eigenvalues obtained from solving Eqs. (1) and (2) give the possibility to reduce the information provided by the complete degrees of freedom of the system into a smaller set of variables containing its most significant part.

In PCA, the relative magnitude of the eigenvalues (with respect to their total sum) is related to the percentage of the total variance projected onto the corresponding eigenvector which is combination of the initial coordinates. Thus, the larger the eigenvalues, the more relevant that particular direction is to represent the overall description of the dynamics. Ideally, one should find an eigenvalue or a small series of them with magnitudes clearly superior to the rest, indicating that the majority of information of the variance of the system is contained in one reaction coordinate. The outcomes of the diagonalization of the correlation matrix are contained in the left side of Fig. 7.

The first four eigenvalues are clearly larger than all the rest, specially the first. We choose to use the first three eigenvectors associated to them as the basis for the PCA procedure, accounting for 70.9% of the variance.

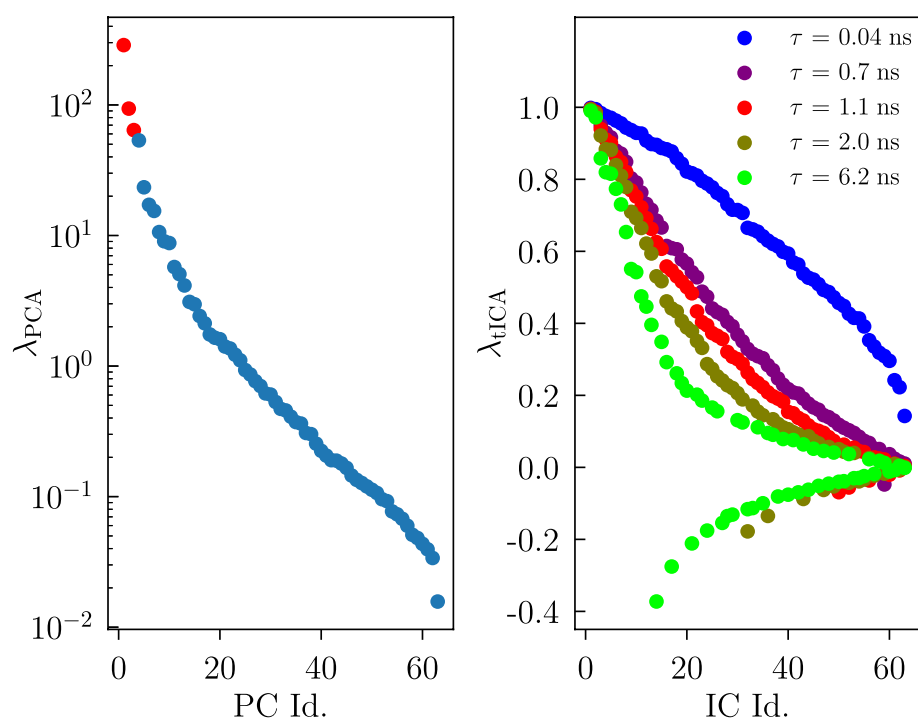
For tICA, the interpretation of the eigenvalues is not as straightforward; their magnitude indicates the minimum timescale its associated coordinate is able to discern, calculated as in Eq. (3).

The eigenvalues of tICA are constrained between -1 and 1, with the negative values only appearing for large values of the lag time  $\tau$  (see  $\tau = 2.04$  ns in Fig. 7). These eigenvalues indicate modes in the system (described by their corresponding eigenvectors) that decay over time, typically corresponding to fast transitions or back-and-forth fluctuations that do not contribute to the stable, meaningful changes tICA aims to capture. Thus, the appearance of negative eigenvalues at a certain lag-time  $\tau$  imposes an effective ceiling onto this magnitude,  $\tau = 1.1$  ns in our particular case. The right side of Fig. 7 contains the eigenvalues for tICA calculated at different lag-times, showing the appearance of negative eigenvalues from  $\tau = 1.1$  ns onwards.

The analysis of the trajectories was carried out with different values of  $\tau$  but, for the rest of this document, a lag-time of  $\tau = 0.7$  ns has been chosen. The reason for this particular value relies on the compromise we found which consists, on the one hand, in eliminating as many fast irrelevant motions of the system as possible that are translated into a few simulation instants and thus undetected by higher values of tICA, while, and on the other hand, in protecting the information contained in the longer-lived states.

#### Projected trajectories

Fig. 8 shows the trajectories of both methods, PCA and tICA, projected on the eigenvectors associated to the three chosen eigenvalues with, on the right hand side, the histogram of the corresponding coordinates occupation.



**Fig. 7.** On the left, PCA eigenvalues obtained from solving (1), with the vertical axis in log-scale. The three eigenvalues in red correspond to the combinations chosen for data projection. On the right, tICA eigenvalues obtained from solving (2) for different values of the time window  $\tau$ , seen in the legend.

The analysis of Fig. 8 reveals appreciable differences between the two methodologies. PCA yields trajectories characterized by histograms exhibiting broad Gaussian-like peaks, whereas tICA produces trajectories with narrower distributions, able to better distinguish between different states of the system with improved clarity. Commonalities are observed in both methods, visible, for example in the first coordinate PC1 and IC1 which show a similar trend: high value for the first interval between  $t = 0.25$  and  $t = 0.8 \mu\text{s}$ , and after that a narrow transition to another constant value. Both trajectories describe the general behavior observed in the root-mean-square deviation (RMSD) (see Fig. 2), with tICA exhibiting less fluctuations than PCA, that, instead, almost reproduces the same shape as the original trajectory.

The different capabilities of the methods in detecting independent and well-defined states throughout the system's trajectory become evident when plotting the coordinates against each other as in Fig. 9. The image depicts a 3D plot of the coordinates extracted from PCA (left) and tICA (right), with the color scale corresponding to the free energy differences in the trajectory determined by the relative occupation of specific coordinate combinations during evolution, and computed as  $-\log(P/P_0)$ , with  $P_0$  being the lowest occupation probability ( $\neq 0$ ) in the states ensemble. On the bottom of each 3D plot there are three heatmaps plotting the coordinates against each other by pairs.

In Fig. 9 is well visible that PCA exhibits greater variance in the evolution of its coordinates featuring bright states surrounded with clouds of relatively populated spots which, when 3D plotted, seem to collapse in a central globule. Conversely, tICA produces distinct bright spots that are well separated, with additional small areas of low population in between, leading to a 3D plot characterized by small cloud separated from one another.

### Complex Markov network results

Once a suitable value for  $\tau$  for tICA is chosen and the trajectories are projected onto the reduced coordinate space we can proceed with their encoding into a complex network. As previously explained, the trajectory undergoes segmentation into equidistant subsets, each constituting a node when populated by the system. Subsequently, a Stochastic Steepest Descent algorithm is applied to facilitate the classification of all nodes into distinct basins of attraction.

The resulting networks of basins depicting the unfolding process are illustrated in Figs. 10 and 11.

The results from tICA are presented first, since they provide a more straightforward depiction of the unfolding process. The PCA results are then analysed subsequently.

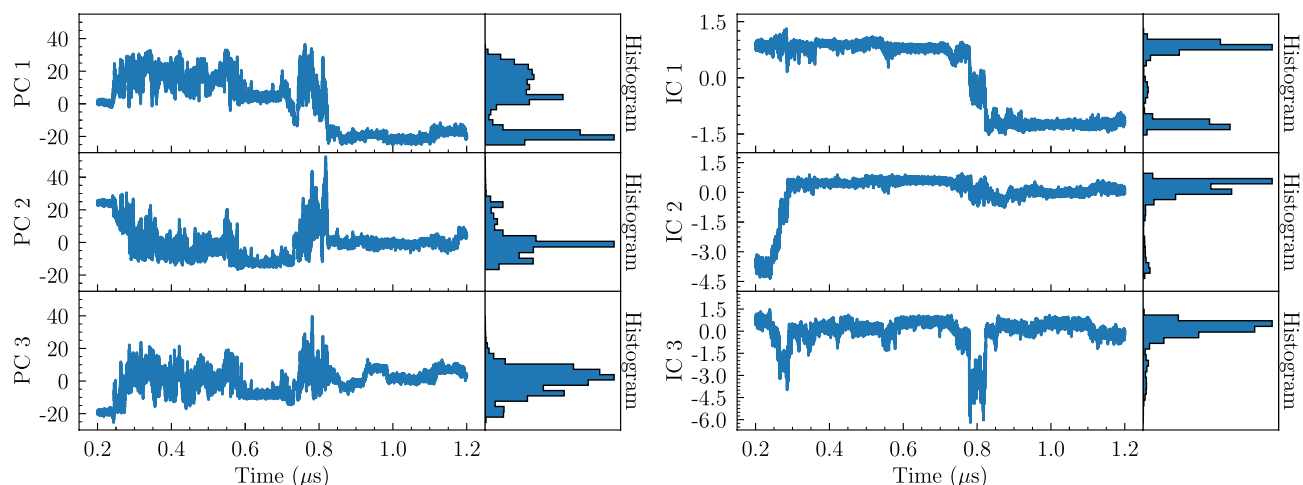
#### tICA analysis

The tICA analysis summarized in Fig. 10 shows a linear representation of the unfolding process, with 7 states identified as basins of attraction in the CMN analysis. The representative node of each basin is labeled from 1 to 7 in the figure, and its corresponding time of appearance in the system trajectory is marked in Fig. 3 with green lines.

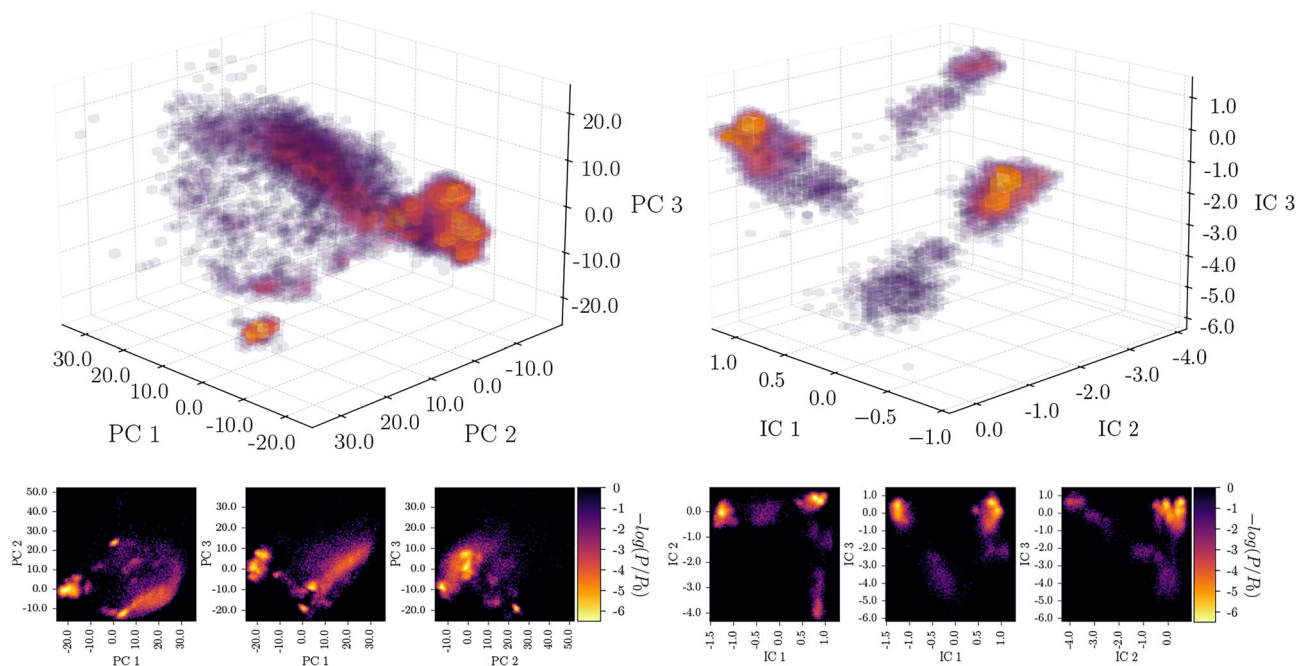
Starting from a fully folded conformation (basin #I), the next state (basin #II) portrays the system as already having severed its Hoogsteen bonds with the fourth guanine tract at the 3'-end of the chain. In this basin tICA is able to recognize the "slip-stranded" conformations between the remaining tracts in the G-triplex, reminiscing of Fig. 6 $\beta$ .

Basin #III depicts the moment after which the fourth guanine tract, previously disconnected from the main triplex but still proximal to it due to crowding effects, has separated. Moreover, another detail can be observed, consisting in the slip-stranding of tract 2 from the remaining triplex (slip-stranded contacts circled in dashed line in Fig. 10, basin #3).

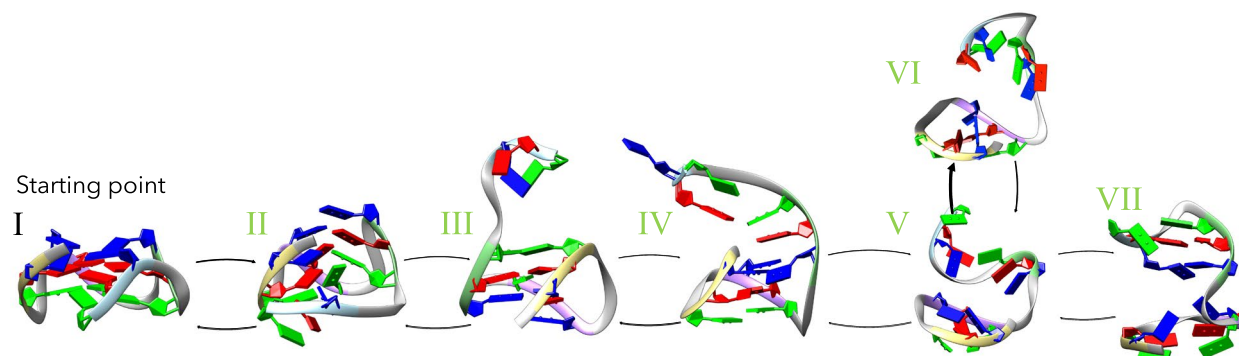
In Basin #IV the third tract has separated from the triplex, rotating over itself and forming a single Hoogsteen bond with the top guanines (blue) of the first and second tracts, as described in Fig. 6 $\epsilon$ . This basin represents



**Fig. 8.** Projected trajectory from Sys. 5 between 0.2 and 1.2  $\mu\text{s}$  onto the first three eigenvectors (left), along with the histogram of the coordinates (right). PCA on the left, tICA with  $\tau = 0.7 \text{ ns}$  on the right.



**Fig. 9.** Colormaps depicting the correlation of the different coordinates emerging from PCA (left) and tICA (right). Within each subfigure, the bottom three panels depict heatmaps illustrating the pairwise relationships between the projected coordinates, with the color representing the relative occupancy of the corresponding coordinate pairs throughout the trajectory. The three-dimensional plots illustrate the combined coordination of the coordinates. Each data point in this plot corresponds to a specific conformation sampled during the trajectory, with the color representing the relative occupancy of a particular combination during the simulation.

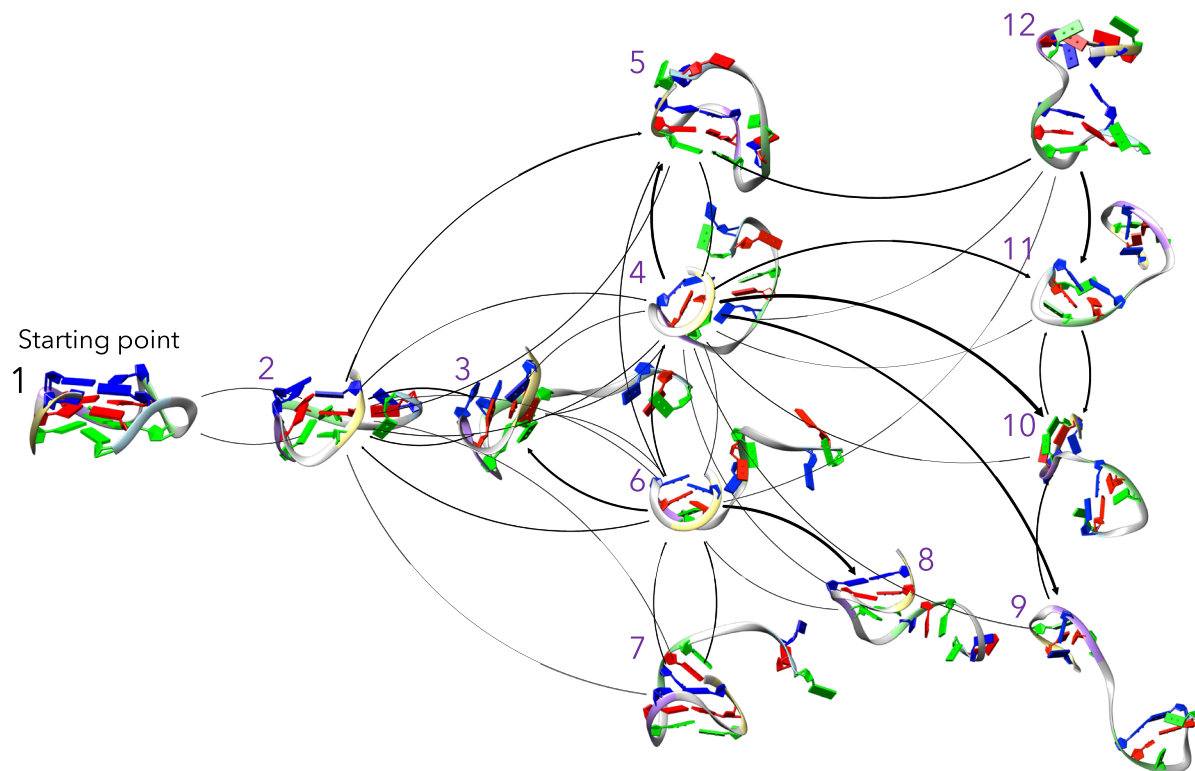


**Fig. 10.** Network of basins found by tICA,  $\tau = 0.7$  ns. The green numbers refer to the equally colored lines in Fig. 3, showing the position of the basins along the trajectory. The files corresponding to the 7 structures, in .pdb format, as well as the videos of the unfolding, can be downloaded from the Supplementary Material.

the first relatively stable state of the system, corresponding to the time interval  $[0.3-0.7] \mu\text{s}$  in Figs. 2 & 3. The fluctuations of these magnitudes in that time interval correspond to waving motions between the third and fourth guanine tracts.

Basin #V and #VI describe very similar situations. The last bond between the third tract and the remaining G-hairpin breaks. The third tract drifts away and pulls the second tract, breaking the G-hairpin and making it into the so called “cross-hairpin” (Fig. 6 $\eta$ ). There is a slight difference between the two states, consisting in the relative position of tracts 3 and 4, which in basin #VI are a little closer due to a stacking interaction created by the top guanines (green) of both tracts (dashed circle in Fig. 10).

Basin #VII corresponds to the final state of the system: tracts 1 and 2 remain in a cross-hairpin, while the third and fourth remain close to each other, with two Hoogsteen bonds formed between them, almost constituting a complete hairpin. This configuration is very stable, with the presence of fluctuations due to the separation of the two blocks (see Figs. 2 and 3 from  $t = 0.8 \mu\text{s}$  onwards).



**Fig. 11.** Network of basins found by PCA. The purple numbers refer to the equally colored lines of Fig. 3, indicating the position of the basins along the trajectory.

With these results in mind, it becomes evident that tICA is able to provide us, upon selection of a proper value for the time window  $\tau$ , a clear denaturation path. The states identified by means of the basins exhibit a high degree of dissimilarity, capturing the significant structural transformations and avoiding the oversampling of short-lived conformations.

In the figure, the links between states are bidirectional (two arrows) and weighted. They represent multiple transitions between the different basins with no information about a particular time sequence. This is due to the fact that, in principle, the system investigated lies in a dynamical equilibrium condition. This means that if the trajectories were long enough, they could be analyzed independently of time, with the occurrence of the same intermediates.

To study the stability of the basins found by tICA, the average escape time of each state can be calculated. These times reflect the average time the G-quadruplex remains in a given conformation before undergoing a transition. By analyzing escape times, insights into the stability and transition kinetics of different conformational states can be gained. The limitation of our analysis relies on the fact that we have a single unfolding trajectory, thus preventing us from performing broad statistics on the occupation of the different states. Nevertheless, we calculate the escape times by using two approaches: one based on Complex Markov Network self-loops and another analyzing the time spent in a given conformation, providing an estimation of the free energy difference between the different basins. The detailed methodology, as well as the results from these calculations, can be found in the Supplementary Material; they reveal that the structure depicted in basin #IV has a larger lifetime than the rest, thus classifying it as a stable intermediate of the unfolding.

Additionally, the above analysis has been performed under a mesoscopic reduction of data, as previously explained. Nevertheless, the same tICA procedure can be applied without that intermediate step. In that case an even broader spectrum of states is resolved, with the resulting network still describing a clear denaturation path. This analysis is also included in the Supplementary Material.

#### PCA results

Fig. 11 shows the network of states generated by the PCA procedure. It is evident that the interpretation of the basins becomes severely more challenging here when compared to the one produced by tICA. PCA aims to build coordinates containing the maximum possible variance which, when applied to the identification of unique basins of attraction, leads to the labeling of some equivalent states as different, overestimating the differences between equivalent configurations.

The initial state corresponds to the same starting conformation as in tICA, designated as basin #1 in Fig. 11. Subsequently, basin #2 marks the point at which the 3'-end disengages from the primary structure, forming the G-triplex. Following this event, a number of tightly connected states is revealed (from basin #3 to #8). Basin #3

shows a slip-stranded G-triplex, with the fourth guanine tract away from the structure. Up to this point, the states analyzed by PCA are equivalent to those of tICA, namely basins #I, #II and #III.

Basin #4 shows a G-hairpin formed by the first and second tracts, with one of the guanines from the third tract in contact with the hairpin, while in basin #5 that contact has broken, leaving a G-hairpin and two drifting guanine tracts (third and fourth).

The remaining basins (numbers #6, #7 and #8) in this heavily interconnected section of the network are classified as erroneously different by PCA respect to tICA. In particular, basins #6 and #7 depict a G-hairpin accompanied by an additional third guanine tract forming a single Hoogsteen bond. The difference between these conformations fundamentally lies in the relative orientation of bonds on the guanine columns forming the hairpin. Basin #8, on the other hand, analogous to basin #5, represents a G-hairpin no longer in contact with the third tract. All these states (basins #5, #6, #7 and #8) identified by PCA are thus equivalent to basin #IV in tICA. Remarkably, the state identified by PCA only (basin #7) is then completely superfluous in describing the unfolding process because it is contained in basin #IV in tICA.

States #9, #10, and #11 and #12 show the final conformation of the system. They correspond to basins #V, #VI and #VII in tICA. They reveal the formation of transient and relatively short-lived bonds reminiscent of G-hairpins, resulting each of them in the division of the G-quadruplex into two well-separated guanine tracts: a cross-hairpin (tracts 1 and 2) and a disordered clump (tracts 3 and 4). However, while the three states found by tICA were distinguishable by the presence or absence of certain interactions (stacking between guanines in basin #VI and guanine-guanine bonds in basin #VII), in the case of PCA they are functionally equivalent: basins #10, #11 and #12 show the cross-hairpin in the first two tracts, along with an additional Hoogsteen bond between the third and fourth, while basin #9 presents an additional bond between them, with no stacking interaction being captured.

To conclude, our simulations reveal that PCA falls short in providing a minimal unfolding path akin to that offered by tICA, instead providing a multitude of interconnected states, some of them functionally equivalent.

Analogously as done with tICA analysis, the relative stability of the basins identified can be studied through the escape times, whose results can be found in the Supplementary Material. Furthermore, the PCA analysis without the mesoscopic reduction is also found in the Supplementary Material, leading to a higher number of basins with increased degeneration between them. They correspond to Supplementary Figures 6 & 7.

## Discussion and conclusions

The study of unfolding pathways in biological systems is essential to understand their stability and functional roles within their environments. In this work, we combined molecular dynamics simulations and multiple analytical techniques to explore the unfolding process of the human G-quadruplex in its parallel configuration (PDB: 1KF1). By characterizing key structural changes during unfolding, we have gained insight into the underlying factors that contribute to the stability and flexibility of this biologically relevant structure.

Replica Exchange Molecular Dynamics (REMD) simulations were employed to enhance conformational sampling. A distinct unfolding transition was observed in one of the eight replicas at temperatures near the denaturation point. Ion dynamics were found to play a fundamental role in this process: total ion loss appears to be a necessary step for unfolding, consistent with previous mechanical denaturation studies<sup>17</sup>. Systems that experienced only partial ion loss largely maintained their structural integrity, even reabsorbing ions from the environment and returning to their initial configurations, as evidenced in Fig. 4.

A mesoscopic reduction of the system coordinates<sup>18</sup> allowed us to simplify the analysis while preserving key dynamics. In fact, comparison with the full-atom representation showed no significant differences in the observed trends (Supplementary Figs. 6-7). This opens the door for the creation and development of realistic coarse-grained models capable of reproducing the unfolding of these structures, similar to those applied to DNA chains<sup>70</sup>.

Dimensionality reduction techniques, specifically PCA and tICA, were then applied to the coarse-grained trajectories to identify the most relevant collective motions. The resulting output trajectories were used to construct Complex Markov Networks, revealing the main states and transitions that define the G-quadruplex unfolding pathway.

Our analysis of the simulations reveals a clear unfolding sequence. Firstly, partial strand separation is accompanied by ion loss from the central channel, allowing the unfolding to commence. Afterwards, one strand detaches, producing a G-triplex<sup>22,61,62,71,72</sup>, then rapidly reorganizing into a partially folded G-hairpin<sup>22,67-69,71-74</sup>, which from then onward remains as a fundamental structural motif, while exhibiting sometimes slip-stranding<sup>65,66</sup> or accompanied by a cross-hairpin<sup>69</sup>. The transition sequence *parallel* G4 → triplex-like → hairpin, represents the dominant unfolding route captured in our simulations. Interestingly, the hairpin intermediate exhibits significantly longer persistence ( $\sim 0.4 \mu\text{s}$ ) when compared to the triplex intermediate ( $\sim 40 \text{ ns}$ ), indicating its relatively higher thermodynamic stability.

Experimental studies have also provided evidence for the existence of these kind of intermediates. Circular dichroism (CD) and single-molecule FRET (smFRET) measurements<sup>63,64,74</sup> demonstrated that *antiparallel* and *hybrid* G4 structures, in both  $\text{K}^+$  and  $\text{Na}^+$  environments, can unfold through transient G-triplex and G-hairpin-like conformations. Although triplex-like intermediates are frequently observed, several single-molecule and spectroscopic studies have also reported alternative pathways that may bypass a well-defined triplex, involving partially folded or hairpin-like conformations instead<sup>64,73-75</sup>. Similar findings were obtained using UV-resonant Raman spectroscopy<sup>76</sup>, while time-resolved optical studies<sup>77</sup> identified G-hairpin formation as a key early event in folding, consistent with the final conformations observed in our unfolding trajectories. Other techniques, including DNA origami<sup>22,72</sup> and microfluidic mixing<sup>71</sup>, have revealed similar transient structures.

The main distinction between these experimental observations and our present results lies in the G4 topology. We focus here on the unfolding of a *parallel* G-quadruplex, in which all detected intermediates maintain parallel

orientation. In contrast, the experimentally identified G-triplexes correspond to either *antiparallel* or *hybrid* structures<sup>78</sup>, so direct evidence of parallel G-triplex intermediates remains elusive. Aznauryan *et al.*<sup>75</sup> noted that fully parallel G4s are rarely stabilized under typical experimental conditions, complicating the observation of their unfolding pathways. Photochemical trapping studies<sup>64,79</sup> also failed to identify long-lived parallel G-triplex species. The most plausible explanation is that parallel G-triplex intermediates possess extremely short lifetimes, making their detection experimentally challenging.

Our simulations are consistent with this interpretation. The distinct lifetimes registered for the intermediates can be coded into occupation probabilities and relative free energy differences. The calculations, though affected by the lack of a broad statistics, are detailed in the “Basin escape times” section of the Supplementary Material. The final state of the system (basin #VII in tICA, basins #9–#12 in PCA) includes a combination of G-hairpin and cross-hairpin features, confirming these as the most stable intermediates emerging from the parallel G4 structure.

Recent computational studies have further clarified the atomistic mechanisms of G-quadruplex folding and unfolding. Enhanced-sampling simulations reveal a highly rugged free-energy landscape populated by metastable triplex-like, hairpin, and slipped intermediates that interconvert through multiple pathways<sup>50,80,81</sup>. Kim *et al.*<sup>81</sup> identified coexisting triplex-like and hairpin states during human telomeric G4 folding, while Pokorná *et al.*<sup>50</sup> showed that parallel G4s can also fold via alternative routes that bypass a well-defined triplex. Janeček *et al.*<sup>80</sup> likewise demonstrated multiple competing pathways in the folding of a parallel G4 from a single strand. Taken together, these results depict G4 unfolding/folding transitions as complex multi-pathway processes where specific linear series of events can probabilistically occur. Within this computational context, we believe our work contributes by uncovering one of the plausible unfolding routes of the parallel G-quadruplex, characterized by a sequential transition from the folded structure through triplex-like and hairpin intermediates.

Regarding the scope of our studies, the simulations were conducted at biologically relevant temperatures slightly above the melting point, thus ensuring that the observed conformations reflect thermally accessible states separated by realistic energy barriers. However, under these thermal conditions, only one replica showed complete unfolding. The pathway observed reproduces intermediate steps that have been reported in both computational and experimental studies, the latter with different loops’ topologies, lending confidence to its biological relevance. Thus, more unfolding events are needed to uncover additional alternative unfolding pathways<sup>50,80,81</sup>, which would require extremely long trajectories in the conditions used in this work.

The sequence here investigated (5 -AGGGTTAGGGTTAGGGTTAGGG-3') represents the minimal telomeric repeat capable of forming a parallel G-quadruplex. Extensions to longer sequences, as well as variations in loop size and rigidity<sup>18,83,84</sup>, cooperative interactions between adjacent G4 units<sup>85,86</sup>, molecular crowding<sup>6,87,88</sup> and sequence constraints<sup>89</sup>, could influence both stability and kinetics. Besides all the many parameters that can be analyzed in the different physiological and biological contexts, our model focuses into the essential unfolding mechanisms followed by the parallel G-quadruplex. Similarly, the incorporation of other G4 topologies, such as antiparallel or hybrid, could lead to a more complete view of the G-quadruplex landscape. We have already made several attempts for both of these conformations, under the same conditions as the parallel simulations, without finding any significant results up to now.

With respect to the analysis methods, both tICA and PCA yielded Markov networks that effectively captured the relevant unfolding transitions. tICA proved particularly efficient in distinguishing kinetically distinct conformations, whereas PCA identified a larger number of states, several of which were structurally equivalent. Future studies may benefit from exploring non-linear dimensionality reduction and machine-learning-based techniques<sup>90,91</sup>, which could refine the description of conformational space, albeit at increased computational cost due to the need for larger ensembles of unfolding trajectories.

Finally, the last configurations obtained (basin #VII in tICA) suggest that refolding may not necessarily restore the original parallel topology. Instead, the unfolded state could evolve into alternative conformations such as antiparallel or hybrid G4s, as reported in other studies, showing a conformational switch between different G4 topologies, that warrants further exploration.

## Data availability

All data generated or analyzed during this study are included in the published article. The structure used in the simulation corresponds to the human telomeric parallel G-quadruplex, found in the Protein Data Bank under the code 1KF1. The files of the structures identified in Figure 10 extracted from the simulations, in .pdb format, as well as all the gromacs parameter files .mdp, are available at the GitHub repository at the address: <https://github.com/AlejandroSainzAgost/GromacsGQuadruplex>.

Received: 11 August 2025; Accepted: 20 November 2025

Published online: 11 February 2026

## References

1. Hans J. Lipps & D. Rhodes. G-quadruplex Structures: In Vivo Evidence and Function. *Trends Cell Biol.* **19**, 414 (2009)
2. Lam, E. Y., Beraldi, D., Tannahill, D. & Balasubramanian, S. G-quadruplex structures are stable and detectable in human genomic DNA. *Nat. Commun.* **4**, 1796 (2013).
3. Burge, S., Parkinson, G. N., Hazel, P., Todd, A. K. & Neidle, S. Quadruplex DNA: sequence, topology and structure. *Nucleic Acids Res.* **34**, 5402 (2006).
4. Louit, G., Hocquet, A., Ghomi, M., Meyer, M. & Sühnel, J. Are guanine tetrads stabilised by bifurcated hydrogen bonds? An AIM topological analysis of the electronic density. *PhysChemComm* **5**, 94 (2002).
5. Li, J., Correia, J. J., Wang, L., Trent, J. O. & Chaires, J. B. Not so crystal clear: the structure of the human telomere G-quadruplex in solution differs from that present in a crystal. *Nucleic Acids Res.* **33**, 4649 (2005).

6. Bergues-Pupo, A. E., Arias-Gonzalez, J. R., Morón, M. C., Fiasconaro, A. & Falo, F. Role of the central cations in the mechanical unfolding of DNA and RNA G-quadruplexes. *Nucleic Acids Res.* **43**, 7638 (2015).
7. Murat, P. & Balasubramanian, S. Existence and Consequences of G-quadruplex Structures in DNA. *Curr. Opin. Genet. Dev.* **25**, 22 (2014).
8. Hans, J. L. & Rhodes, D. G-quadruplexes and Their Regulatory Roles in Biology. *Nucleic Acids Res.* **43**, 8627 (2015).
9. Hänsel-Hertsch, R., Di Antonio, M. & Balasubramanian, S. DNA G-quadruplexes in the human genome: detection, functions and therapeutic potential. *Nat. Rev. Mol. Cell Biol.* **18**, 279 (2017).
10. Endoh, T. & Sugimoto, N. Mechanical insights into ribosomal progression overcoming RNA G-quadruplex from periodical translation suppression in cells. *Sci. Rep.* **6**, 22719 (2016).
11. Lopes, J. et al. G-quadruplex-induced instability during leading-strand replication. *EMBO J.* **30**, 4033 (2011).
12. Papadopoulou, C., Guilbaud, G., Schiavone, D. & Sale, J. E. Nucleotide pool depletion induces G-quadruplex-dependent perturbation of gene expression. *Cell Rep.* **13**, 2491 (2015).
13. Siddiqui-Jain, A., Grand, C. L., Bearss, D. J. & Hurley, L. H. Direct evidence for a G-quadruplex in a promoter region and its targeting with a small molecule to repress c-MYC transcription. *Proc. Natl. Acad. Sci. USA* **99**, 11593 (2002).
14. Tian, T., Chen, Y., Wang, S. & Zhou, X. G-quadruplex: a regulator of gene expression and its chemical targeting. *Chem* **4**, 1314 (2018).
15. Kosiol, N., Juranek, S., Brossart, P., Heine, A. & Paeschke, K. G-quadruplexes: a promising target for cancer therapy. *Mol. Cancer* **20**, 40 (2021).
16. Sen, D. & Gilbert, W. Formation of parallel four-stranded complexes by guanine-rich motifs in DNA and its implications for meiosis. *Nature* **334**, 364 (1988).
17. Bergues-Pupo, A. E., Gutiérrez, I., Arias-Gonzalez, J. R., Falo, F. & Fiasconaro, A. Mesoscopic model for DNA G-quadruplex unfolding. *Sci. Rep.* **7**, 11756 (2017).
18. Bergues-Pupo, A. E., Falo, F. & Fiasconaro, A. Modelling the DNA Topology: The Effect of the Loop Bending on G-quadruplex Stability. *J. Stat. Mech.: Theory Exp.* **9**, 094004 (2019).
19. Cheng, Y., Zhang, Y. & Huijuan, Y. Characterization of G-quadruplexes Folding/Unfolding Dynamics and Interactions with Proteins from Single-Molecule Force Spectroscopy. *Biomolecules* **11**, 1579 (2021).
20. Jana, J. & Weisz, K. Thermodynamic stability of G-quadruplexes: impact of sequence and environment. *ChemBioChem* **22**, 2848 (2021).
21. Lane, A. N., Chaires, J. B., Gray, R. D. & Trent, J. O. Stability and kinetics of G-quadruplex structures. *Nucleic Acids Res.* **36**, 5482 (2008).
22. Rajendran, A. et al. Small molecule binding to a G-hairpin and a G-Triplex: A new insight into anticancer drug design targeting G-rich regions. *Chem. Commun.* **51**, 9181 (2015).
23. Schildkraut, C. & Lifson, S. Dependence of the melting temperature of DNA on salt concentration. *Biopolymers* **3**, 195 (1965).
24. Nguyen, T. Q. N., Lim, K. W. & Phan, A. T. Folding kinetics of G-quadruplexes: Duplex stem loops drive and accelerate G-quadruplex folding. *J. Phys. Chem. B* **124**, 5122 (2020).
25. Han, H., Langle, D. R., Rangan, A. & Hurley, L. H. Selective interactions of cationic porphyrins with G-quadruplex structures. *J. Am. Chem. Soc.* **123**, 8902-13 (2001).
26. Miyoshi, D., Karimata, H. & Sugimoto, N. Hydration regulates thermodynamics of G-quadruplex formation under molecular crowding conditions. *J. Am. Chem. Soc.* **128**, 7957 (2006).
27. Phan, A. T., Luu, K. N. & Patel, D. J. Different loop arrangements of intramolecular human telomeric (3+1) G-quadruplexes in K<sup>+</sup> Solution. *Nucleic Acids Res.* **34**, 5715 (2006).
28. Nakata, M., Kosaka, N., Kawachi, K. & Miyoshi, D. Quantitative effects of the loop region on Topology, Thermodynamics, and cation binding of DNA G-quadruplexes. *ACS Omega* **9**, 35028 (2024).
29. Haider, S., Neidle, S. (2010). Molecular Modeling and Simulation of G-Quadruplexes and Quadruplex-Ligand Complexes. In: Baumann, P. (eds) G-Quadruplex DNA. Methods in Molecular Biology, vol 608. Humana Press.
30. Zhu, H., Xiao, S. & Liang, H. Structural dynamics of human telomeric G-quadruplex loops studied by molecular dynamics simulations. *PLoS one* **8**, e71380 (2013).
31. Hukushima, K. & Nemoto, K. Exchange Monte Carlo method and application to spin glass simulations. *J. Phys. Soc. Jpn.* **65**, 1604 (1996).
32. Sugita, Y. & Okamoto, Y. Replica-exchange molecular dynamics method for protein folding. *Chem. Phys. Lett.* **314**, 141 (1999).
33. Parkinson, G. N., Lee, M. P. H. & Neidle, S. Crystal Structure of Parallel Quadruplexes from Human Telomeric DNA. *Nature* **417**, 876 (2002).
34. Pearson, K. LIII. On lines and planes of closest fit to systems of points in space. *The London, Edinburgh, and Dublin Philosophical Magazine and Journal of Science* **11**, 559 (1901).
35. Hotelling, H. Analysis of a complex of statistical variables into principal components. *J. Educ. Psychol.* **24**, 417 (1933).
36. Molgedey, L. & Schuster, H. G. Separation of a mixture of independent signals using time delayed correlations. *Phys. Rev. Lett.* **72**, 3634 (1994).
37. Prada-Gracia, D., Gómez-Gardeñes, J., Echenique, P. & Falo, F. Exploring the free energy landscape: from dynamics to networks and back. *PLoS Comput. Biol.* **5**, e1000415 (2009).
38. Wales, D. J. Energy Landscapes: Some New Horizons. *Curr. Opin. Struct. Biol.* **20**, 3 (2010).
39. Wales, D. J. Exploring Energy Landscapes. *Annu. Rev. Phys. Chem.* **69**, 401 (2018).
40. Zeng, X. et al. Unfolding mechanism of thrombin-binding aptamer revealed by molecular dynamics simulation and Markov state model. *Sci. Rep.* **6**, 24065 (2016).
41. Bekker, H. et al. Gromacs: A parallel computer for molecular dynamics simulations. *Physics computing* **92**, 252 (1993).
42. Berendsen, H. J. C., van der Spoel, D. & van Drunen, R. GROMACS: A message-passing parallel molecular dynamics implementation. *Comp. Phys. Comm.* **91**, 43 (1995).
43. Pérez, A. et al. Refinement of the AMBER Force Field for Nucleic Acids: Improving the Description of  $\alpha\gamma$  Conformers. *Biophys. J.* **92**, 3817 (2007).
44. Šponer, J., Cang, X. & Cheatham, T. E. Molecular dynamics simulations of G-DNA and perspectives on the simulation of nucleic acid structures. *Methods* **57**, 25 (2012).
45. Li, N., Gao, Y., Qiu, F. & Zhu, T. Benchmark force fields for the molecular dynamic simulation of G-quadruplexes. *Molecules* **26**, 5379 (2021).
46. Zgarbová, M. et al. Refinement of the sugar-phosphate backbone Torsion Beta for AMBER force fields improves the description of Z- and B-DNA. *J. Chem. Theory Comput.* **11**, 5723 (2015).
47. Ivani, I. et al. Parmbsc1: A refined force field for DNA simulations. *Nat. Methods* **13**, 1355 (2016).
48. Jorgensen, W. L., Chandrasekhar, J., Madura, J. D., Impey, R. W. & Klein, M. L. Comparison of simple potential functions for simulating liquid water. *J. Chem. Phys.* **79**, 926 (1983).
49. Mergny, J. L., Phan, A. T. & Lacroix, L. Following G-quartet formation by UV-spectroscopy. *FEBS letters* **435**, 74 (1998).
50. Pokorná, P., Mlýnský, V., Bussi, G., Šponer, J. & Stadlbauer, P. Parallel G-quadruplex folds via multiple paths involving G-tract stacking and structuring from coil ensemble. *Preprint, bioRxiv* <https://doi.org/10.1101/2023.09.09.556957> (2023)
51. Qi, R., Wei, G., Ma, B. & Nussinov, R. *Replica Exchange Molecular Dynamics: A Practical Application Protocol with Solutions to Common Problems and a Peptide Aggregation and Self-Assembly Example* (In Peptide Self-Assembly. Springer, New York, 2018).

52. Abraham, M. J. & Gready, J. E. Ensuring mixing efficiency of replica-exchange molecular dynamics simulations. *J. Chem. Theory Comput.* **4**, 1119 (2008).
53. Periolo, X. & Mark, A. E. Convergence and sampling efficiency in replica exchange simulations of peptide folding in explicit solvent. *J. Chem. Phys.* **126**, 014903 (2007).
54. Robert, C.P. & Casella, G. Monte Carlo Statistical Methods. *Springer Texts in Statistics*. New York, NY: Springer New York (2004)
55. Gowers, R. et al. MDAnalysis: A Python package for the rapid analysis of molecular dynamics simulations. Proc. of the 15th PYTHON in SCIENCE CONF 98 (2016)
56. Michaud-Agrawal, N., Denning, E. J., Woolf, T. B. & Beckstein, O. MDAnalysis: A toolkit for the analysis of molecular dynamics simulations. *J. Comput. Chem.* **32**, 2319 (2011).
57. Amadei, A., Linssen, A. B. & Berendsen, H. J. Essential dynamics of proteins. *Proteins: Structure, Function, and Bioinformatics* **17**, 412 (1993).
58. Ali, M. U., Ahmed, S., Ferzund, J., Mehmood, A. & Rehman, A. Using PCA and factor analysis for dimensionality reduction of bio-informatics data. *arXiv preprint arXiv:1707.07189* (2017)
59. Hyvarinen, A., Karhunen, J. & Oja, E. Independent component analysis. *John Wiley & Sons* (2001).
60. Pérez-Hernández, G., Paul, F., Giorgino, T., De Fabritiis, G. & Noé, F. Identification of slow molecular order parameters for Markov model construction. *J. Chem. Phys.* **139**, 015102 (2013).
61. Stadlbauer, P., Trantírek, L., Cheatham, T. E., Koča, J. & Šponer, J. Triplex intermediates in folding of human telomeric quadruplexes probed by microsecond-scale molecular dynamics simulations. *Biochimie* **105**, 22 (2014).
62. Bončina, M., Lah, J., Prislán, I. & Vesnaver, G. Energetic basis of human telomeric DNA folding into G-quadruplex structures. *J. Am. Chem. Soc.* **134**, 9657 (2012).
63. Hou, X. et al. Involvement of G-Triplex and G-Hairpin in the multi-pathway folding of human telomeric G-quadruplex. *Nucleic Acids Res.* **45**, 11401 (2017).
64. Tassilo Grün, J. & Schwalbe, H. Folding dynamics of polymorphic G-quadruplex structures. *Biopolymers* **113**, e23477 (2022).
65. Stadlbauer, P., Krepl, M., Cheatham, T. E., Koca, J. & Šponer, J. Structural dynamics of possible late-stage intermediates in folding of quadruplex DNA studied by molecular simulations. *Nucleic Acids Res.* **41**, 7128 (2013).
66. Kejnovská, I. et al. G-Quadruplex formation by DNA sequences deficient in Guanines: Two tetrad parallel quadruplexes do not fold intramolecularly. *Chem. Eur. J.* **27**, 12115 (2021).
67. Gajarský, M. et al. Structure of a stable G-Hairpin. *J. Am. Chem. Soc.* **139**, 3591 (2017).
68. Stadlbauer, P. et al. Hairpins participating in folding of human telomeric sequence quadruplexes studied by standard and T-REMD simulations. *Nucleic Acids Res.* **43**, 9626 (2015)
69. Stadlbauer, P. et al. Parallel G-Triplexes and G-Hairpins as potential transitory ensembles in the folding of parallel-stranded DNA G-quadruplexes. *Nucleic Acids Res.* **47**, 7276 (2019).
70. Ouldrige, T. E., Ard, L. A. & Doye, J. P. K. Structural, mechanical, and thermodynamic properties of a coarse-grained DNA model. *J. Chem. Phys.* **134**, 085101 (2011).
71. Li, Y., Liu, C., Feng, X., Xu, Y. & Liu, B.-F. Ultrafast microfluidic mixer for tracking the early folding kinetics of human telomere G-quadruplex. *Anal. Chem.* **86**, 4333 (2014).
72. Rajendran, A., Endo, M., Hidaka, K. & Sugiyama, H. Direct and single-molecule visualization of the solution-state structures of G-Hairpin and G-Triplex intermediates. *Angewandte Chemie International Edition* **53**, 4107 (2014).
73. Mashimo, T., Yagi, H., Sannohe, Y., Rajendran, A. & Sugiyama, H. Folding pathways of human telomeric Type-1 and Type-2 G-quadruplex structures. *J. Am. Chem. Soc.* **132**, 14910 (2010).
74. Gray, R. D., Trent, J. O. & Chaires, J. B. Folding and unfolding pathways of the human telomeric G-quadruplex. *J. Mol. Biol.* **426**, 1629 (2014).
75. Aznauryan, M., Sondergaard, S., Noer, S. L., Schiott, B. & Birkedal, V. A direct view of the complex multi-pathway folding of telomeric G-quadruplexes. *Nucleic Acids Res.* **44**, 11024 (2016).
76. Di Fonzo, S., Amato, J. & Marzano, S. et al. Investigating temperature-induced unfolding pathways of DNA G-quadruplexes via 2D UV resonant raman spectroscopy. *Int. J. Biol. Macromol.* **323**, 147141 (2025)
77. Monsen, R. C., Sabo, T. M., Gray, R., Hopkins, J. B. & Chaires, J. B. Early events in G-quadruplex folding captured by time-resolved small-angle X-ray scattering. *Nucleic Acids Res.* **53**, gkaf043 (2025).
78. Cerofolini, L. et al. G-Triplex structure and formation propensity. *Nucleic Acids Res.* **42**, 13393 (2014).
79. Grün, J. T., Blümler, A., Burkhart, I., Wirmer-Bartoschek, J., Heckel, A. & Schwalbe, H. Unraveling the kinetics of spare-tire DNA G-quadruplex folding. *J. Am. Chem. Soc.* **143**, 6185 (2021)
80. Janeček, M. et al. Computer folding of parallel DNA G-Quadruplex: Hitchhiker's Guide to the conformational space. *J. Comput. Chem.* **46**, e27535 (2025)
81. Kim, H., Eunae, K. & Youngshang, P. Computational probing of the folding mechanism of human telomeric G-quadruplex DNA. *J. Chem. Inf. Model.* **63**, 6366 (2023).
82. Lim, K. W., Khong, Z. J. & Phan, A. T. Thermal stability of DNA quadruplex-duplex hybrids. *Biochemistry*. **53**, 247 (2014).
83. Tucker, B. A. et al. Stability of the Na<sup>+</sup> form of the human telomeric G-quadruplex: Role of adenines in stabilizing G-quadruplex structure. *ACS Omega* **3**, 844 (2018)
84. Jurkowski, M., Kogut, M., Olewniczak, M., Glinko, J. & Czub, J. Large-scale conformational analysis explains G-quadruplex topological landscape. *J. Phys. Chem. B* **129**, 9622 (2025).
85. Kristoffersen, E. L., Coletta, A., Lund, L. M., Schiott, B. & Birkedal, V. Inhibited complete folding of consecutive human telomeric G-quadruplexes. *Nucleic Acids Res.* **51**, 1571 (2023)
86. Saintomé, C., Amrane, S., Mergny, J. L. & Alberti, P. The exception that confirms the rule: a higher-order telomeric G-quadruplex structure more stable in sodium than in potassium. *Nucleic Acids Res.* **44**, 2926 (2016).
87. Gao, C. et al. Effects of molecular crowding on the structure, stability, and interaction with ligands of G-quadruplexes. *ACS Omega* **8**, 14342 (2023)
88. Lacer, A. N., Symasek, A., Gunter, A. & Lee, H. T. Slow G-quadruplex conformation rearrangement and accessibility change induced by potassium in human telomeric single-stranded DNA. *J. Phys. Chem. B* **128**, 5950 (2024).
89. Bugaut, A. & Alberti, P. Understanding the stability of DNA G-quadruplex units in long human telomeric strands. *Biochimie* **113**, 125 (2015).
90. Moses, A. Statistical modeling and machine learning for molecular biology. *Chapman and Hall/CRC* (2017)
91. Mehta, P., Bukov, M., Wang, C., Day, A.G.R., Richardson, C., Fisher, C.K. & Schwab, D. J. A high-bias, low-variance introduction to machine learning for physicists. *Phys. Rep.* **810**, 1–124 (2019)

## Author contributions

AF & FF conceived the project. AS-A conducted the simulations and the data analysis. AF, FF & AS-A wrote and reviewed the manuscript.

## Funding

The authors acknowledge the Grant No. PID2020-113582GB-I00 and the Grant No. PID2023-147734NB-I00

funded by MCIN/AEI/10.13039/501100011033, the support of the Aragon Government to the Recognized group 'E36\_23R Física Estadística y no-lineal (FENOL)'. AS-A also acknowledges the support of the predoctoral FPI fellowship PRE2021-100456 funded by MCIN/AEI/10.13039/501100011033.

## Declarations

### Competing interests

The authors declare no competing interests.

### Additional information

**Supplementary Information** The online version contains supplementary material available at <https://doi.org/10.1038/s41598-025-29993-1>.

**Correspondence** and requests for materials should be addressed to A.F.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

**Open Access** This article is licensed under a Creative Commons Attribution-NonCommercial-NoDerivatives 4.0 International License, which permits any non-commercial use, sharing, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if you modified the licensed material. You do not have permission under this licence to share adapted material derived from this article or parts of it. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by-nc-nd/4.0/>.

© The Author(s) 2026