

Electronic Supplementary Information

For

**A first-principles exploration of the conformational space of neutral
and sodiated di-saccharides assisted by semi-empirical methods and
neural network potentials**

Huu Trong Phan,^{a,b,c} Pei-Kang Tsou,^a Po-Jen Hsu^a and Jer-Lai Kuo^{a,b,c,d}

- a. Institute of Atomic and Molecular Sciences, Academia Sinica, Taipei, 10617, Taiwan
- b. Molecular Science and Technology Program, Taiwan International Graduate Program, Academia Sinica, Taipei, 11529, Taiwan
- c. Department of Chemistry, National Tsing Hua University, Hsinchu 30013, Taiwan.
- d. International Graduate Program of Molecular Science and Technology (NTU-MST), National Taiwan University, Taipei 10617, Taiwan

*To whom correspondence should be addressed. E-mail:

jlkuo@pub.iams.sinica.edu.tw (Jer-Lai Kuo)

Contents

The sampling details

The analysis of vibrational IRPD spectra on 1-1, 1-2, 1-3, and 1-6 linkage

Table S1. The predictive performance of different NNP generations on the individual test set for of mono-saccharide (α/β -Glc) and 19 Glc-based di-saccharides. The unit for E-MAE and F-MAE is kJ/mol and kJ/mol/Å, respectively. The relevant information of sodiated α -Maltose (α 14 α) is highlighted in bold.

Figure S1. The evolution of the third quartile and the median of the energy from snapshots extracted from optimization trajectories carried out by the DFTB3 method with the initial geometries of sodiated α -Maltose are created from the attach-and-rotate sampling scheme. The energy is relative to the corresponding local minima of the respective trajectory.

Figure S2. The potential energy correlation between the NNP-0 and DFT M06-2X/6-311+G(d,p) on the test set of each individual sodiated di-saccharides. The zero of energy is set as the energy of the global minimum of sodiated α -Maltose evaluated at DFT level.

Figure S3. The correlation in the evaluation of atomic forces between the NNP-0 and DFT M06-2X/6-311+G(d,p) on the test set of each sodiated di-saccharide.

Figure S4. The histogram of RMSD values for pairs of minima in sodiated α -Maltose. Each pair consists of a NNP minima and a corresponding M06-2X minima, obtained via the reoptimization from NNP minima.

Figure S5. The temperature dependence of the relative population of low-energy conformers of sodiated di-saccharides 1-1, 1-2, 1-3, and 1-6 linkage. The number in the square bracket “[]” indicates the energy rank of the local minima sharing the same structural features. For example, “[0]” is used to denote the conformer with the lowest energy, while “[1]” is used to represent the conformer with the second lowest energy.

Figure S6. The accumulated HSA spectra of all 19 sodiated Glc-based disaccharides. The total accumulated spectra are shown in gray line, and the spectra from local minima sharing similar conformations are depicted in various colors and styles.

Figure S7. The vibrational spectra of sodiated α Glc-(1-1)- α Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S8. The vibrational spectra of sodiated α Glc-(1-1)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S9. The vibrational spectra of sodiated β Glc-(1-1)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S10. The vibrational spectra of sodiated α Glc-(1-2)- α Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S11. The vibrational spectra of sodiated α Glc-(1-2)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S12. The vibrational spectra of sodiated β Glc-(1-2)- α Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S13. The vibrational spectra of sodiated β Glc-(1-2)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S14. The vibrational spectra of sodiated α Glc-(1-3)- α Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S15. The vibrational spectra of sodiated α Glc-(1-3)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S16. The vibrational spectra of sodiated β Glc-(1-3)- α Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S17. The vibrational spectra of sodiated β Glc-(1-3)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S18. The vibrational spectra of sodiated α Glc-(1-4)- α Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S19. The vibrational spectra of sodiated α Glc-(1-4)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S20. The vibrational spectra of sodiated β Glc-(1-4)- α Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S21. The vibrational spectra of sodiated β Glc-(1-4)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S22. The vibrational spectra of sodiated α Glc-(1-6)- α Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S23. The vibrational spectra of sodiated α Glc-(1-6)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S24. The vibrational spectra of sodiated β Glc-(1-6)- α Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S25. The vibrational spectra of sodiated β Glc-(1-6)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

Figure S26. The probability distribution of CCS values (\AA^2) of the conformers of sodiated disaccharide with (a) 1-1 linkage, (b) 1-2 linkage, (c) 1-3 linkage, and (4) 1-6 linkage. The x-axis displays the CCS value (\AA^2), and the y-axis represents the HSA population probability. The width is set at 0.6\AA^2 .

The sampling details of the “attach-and-rotate” algorithm

The “attach-and-rotate” sampling algorithm has two primary operations: (1) the attachment of two mono-saccharide structures and followed by (2) the rotation of the two dihedral angles on the glycosidic bonds.

Two input structures are named as the base and seed molecules. In the attachment operation, for a given target di-saccharide molecule, the attachment of the two structures was performed by aligning the vector formed by the O-H bond of the base molecule with the O-C bond vector of the seed one. The H atoms on the O-H vector of the seed molecule and the O, H atoms attaching to the O-C bond of the seed molecule are discarded. Two remaining bodies are joined to form the initial guess of a di-saccharide molecule. The base and seed molecules are the non-reducing and reducing ends, respectively. In this work, the linkages sampled are 1-2, 1-3, 1-4, and 1-6, which corresponds to the O1 atom of the base molecule forming the linkages with C2, C3, C4, and C6 of the seed molecule. For the sampling of the sodiated di-saccharide molecule, one neutral conformer will combine with another sodiated one in two modes: the neutral conformer is used as the reducing end, the sodiated one is used as a non-reducing end, and another mode is vice versa.

The attachment of the two mono-saccharides and the rotation along the glycosidic bonds may result in unphysical geometries, with their atoms being placed too close to each other. In this work, the physical geometry is defined in a simple rule based on the distance between atoms. Specifically, the distance between heavy atoms (non-hydrogen atoms) should be at least 1.4 Å, and the distance between hydrogen and another (regardless of heavy and hydrogen atoms) should be at least 0.5 Å. With these criteria, the set of initial geometries was quickly refined by performing a fast screening to remove the unphysical geometries.

The analysis of IRPD spectra on 1-1, 1-2, 1-3, and 1-6 linkage

The accumulated spectra for the 1-1 linkage are detected in the blue shifts ($\sim 70\text{ cm}^{-1}$) of the lowest wavenumber, from $\alpha 11\alpha$ to $\alpha 11\beta$ and $\beta 11\beta$ (Figure S6). The 5-6-5'-6' is the favored coordination pattern with the dominant contribution to the accumulated vibrational spectra (Figure S7-S9). The lowest wavenumber peaks exhibit a hydrogen bond pattern akin to the interaction between OH2'...O2. The hydrogen bond length in $\beta 11\beta$ is longer than $\alpha 11\alpha$ due to structural restraints (the shortest distance is 2.12 Å vs. 1.89 Å in $\alpha 11\alpha$). The lowest wavenumber peak of OH2' stretching in $\alpha 11\beta$, is comparable to $\beta 11\beta$. Figure S6 shows a minor visible peak of $\alpha 11\beta$ from the conformer of pattern 4C1-4C1_5-6-5'-6' with relative energy +3.52 kJ/mol (Figure S8). In $\alpha 11\beta$, it is blue-shifted $\sim 30\text{ cm}^{-1}$ further to see a significant peak from the global minima of $\alpha 11\beta$ (Figure S6, S8).

In 1-2 linkage, OH stretching is lowest in $\alpha 12\alpha/\beta$, with peaks at 3297-3300 cm^{-1} , significantly lower than the lowest wavenumber peak in αGlc (3438 cm^{-1}). The weak population of these conformers, which are +4-5 kJ/mol above the most stable minimum (Figure S10, Figure S11), causes these peaks to appear as a small hump in the accumulated vibrational spectra, which may be imperceptible in experimental spectra (Figure S6). The visible intense peaks are from the global minimum of each kind of di-saccharide in 1-2 linkage. Since the lowest wavenumber peaks with

significant intensity are apparent at 3600 cm^{-1} , our simulation predicts $\beta 12\beta$ is easily identifiable with the absence of major features around $3450\text{-}3550\text{ cm}^{-1}$. The remaining types with the lowest wavenumber in that region are very close to each other. These types are challenging to assign, relying only on the left-most peak. Thus, the accurate assignment would also consider higher-wavenumber features.

For the 1-3 linkage, the low wavenumber peaks in the range $3400\text{-}3500\text{ cm}^{-1}$ are all observed in 4 types ($\alpha 13\alpha$, $\alpha 13\beta$, $\beta 13\alpha$, $\beta 13\beta$) (Figure S6). In $\alpha 13\alpha$, the expanded peak at 3447 cm^{-1} is mainly contributed by the second lowest conformers $4C1\text{-}4C1\text{-}5\text{-}6\text{-}2'$, with a relative energy of $+1.10\text{ kJ/mol}$ (Figure S6, S14). IR spectra in $\alpha 13\beta$ show four visible peaks from conformers with the same structural pattern as $1C4\text{-}OS2\text{-}g\text{-}2\text{-}1'\text{-}5'\text{-}6'$ (with relative energy $1.54, 2.25, 2.51\text{ kJ/mol}$). However, the experiment may not recognize these peaks well due to their small intensities. In $\beta 13\alpha$, the most substantial peak appears around 3443 cm^{-1} (Figure S6, S16), indicating a strong hydrogen bonding interaction comparable to sodiated $\alpha\text{-Glc}$. $OS2\text{-}4C1\text{-}g\text{-}3\text{-}5\text{-}6\text{-}4'$ contributes substantially to peak intensity with the energy $+1.22$ above the global minimum. Two apparent low wavenumber peaks in $\beta 13\beta$ (Figure S6, S17) can be observed around $3400\text{-}3500\text{ cm}^{-1}$. Both are associated with the same hydrogen bonding interaction ($OH4'\dots O6'$, $1.81\text{-}1.83\text{ \AA}$). The assignment of the exact type in this linkage, analogous to 1-2 linkage, is quite challenging as their lowest wave number peaks with major intensity are populating in $3400\text{-}3500\text{ cm}^{-1}$ for $\alpha 13\alpha$, $\beta 13\alpha$, $\beta 13\beta$. The $\alpha 13\beta$ type is easier to recognize as the major peaks are blue-shifted to $\sim 3580\text{ cm}^{-1}$.

Among four di-saccharides of 1-6 linkage (Figure S6), only $\alpha 16\beta$ has a small visible peak at around 3500 cm^{-1} , making it easiest to characterize. This peak is from the second most stable conformers with a relative energy of $+2.26\text{ kJ/mol}$ (Figure S23). In $\beta 16\alpha$, a peak at 3450 cm^{-1} is less intense due to the weak population of the $4C1\text{-}4C1\text{-}6\text{-}3'\text{-}4'$ conformer (relative energy $+5.17\text{ kJ/mol}$). Instead, the two most stable conformers with a similar conformational pattern of $OS2\text{-}4C1\text{-}3\text{-}5\text{-}6\text{-}1'\text{-}5'$ (relative energy $+0.0$ and $+0.03\text{ kJ/mol}$) contribute to the major peaks at 3550 cm^{-1} . In $\alpha 16\alpha$, the global minimum in $\alpha 16\alpha$ contributes to the intense peak at $\sim 3580\text{ cm}^{-1}$, along with three smaller visible peaks at $3500\text{-}3580\text{ cm}^{-1}$ from a $4C1\text{-}4C1\text{-}6\text{-}3'\text{-}4'$ conformer (relative energy $+2.07\text{ kJ/mol}$) and two $1C4\text{-}1C4\text{-}2\text{-}4\text{-}5\text{-}1'\text{-}5'$ conformers (relative energy $+2.7$ and $+3.3\text{ kJ/mol}$). On the other hand, the $\beta 16\beta$ shows the most of its features at 3600 cm^{-1} , with a tiny peak at 3580 cm^{-1} from a high-lying energy conformer ($OS2\text{-}4C1\text{-}3\text{-}5\text{-}6\text{-}1'\text{-}5'$) with a relative energy of $+5.32\text{ kJ/mol}$. Therefore, our analysis would expect peaks in the $3500\text{-}3550\text{ cm}^{-1}$ range for $\alpha 16\alpha/\beta$, whereas similar features are lacking in $\beta 16\alpha/\beta$.

Table S1. The predictive performance of different NNP generations on the individual test set for of mono-saccharide (α/β -Glc) and 19 Glc-based di-saccharides. The unit for E-MAE and F-MAE is kJ/mol and kJ/mol/Å, respectively. The relevant information of sodiated α -Maltose (α 14 α) is highlighted in bold.

		Test	NNP-0		NNP-1		NNP-2	
		set size	E MAE	F MAE	E MAE	F MAE	E MAE	F MAE
mono-		9233	1.35	1.55	1.43	1.56	1.45	1.60
saccharide								
neutral di-		25187	12.27	5.85	2.27	1.55	2.28	1.56
saccharides								
sodiated	α 11 α	3974	6.02	6.45	3.07	1.70	2.51	1.48
di-	α 11 β	3965	8.44	6.71	2.41	1.67	2.24	1.51
saccharides	β 11 β	4137	9.55	7.13	2.60	1.76	2.17	1.51
	α 12 α	3999	12.84	6.4	2.89	1.81	2.27	1.57
	α 12 β	3976	9.87	6.37	2.93	1.74	2.41	1.50
	β 12 α	4459	9.54	6.90	2.70	1.84	2.30	1.53
	β 12 β	4071	6.76	7.06	2.66	1.80	2.22	1.52
	α 13 α	4115	11.54	7.52	3.05	1.81	2.46	1.53
	α 13 β	3983	9.68	7.65	3.05	1.72	2.24	1.48
	β 13 α	4208	9.58	7.87	3.00	1.87	2.48	1.59
	β 13 β	4331	7.38	8.01	2.61	1.85	2.14	1.57
	α14α	3919	4.83	3.53	2.84	1.79	2.47	1.72
	α 14 β	3946	7.06	4.8	3.10	1.80	2.29	1.54
	β 14 α	3862	6.46	5.19	3.17	1.88	2.61	1.58
	β 14 β	4125	9.14	5.86	2.49	1.80	2.23	1.52
	α 16 α	3937	31.29	7.78	3.11	1.81	2.26	1.55
	α 16 β	4434	28.93	7.98	2.90	1.77	2.18	1.53
	β 16 α	3897	25.55	7.86	2.96	1.87	2.25	1.58
	β 16 β	4120	21.70	8.06	2.68	1.88	2.21	1.62

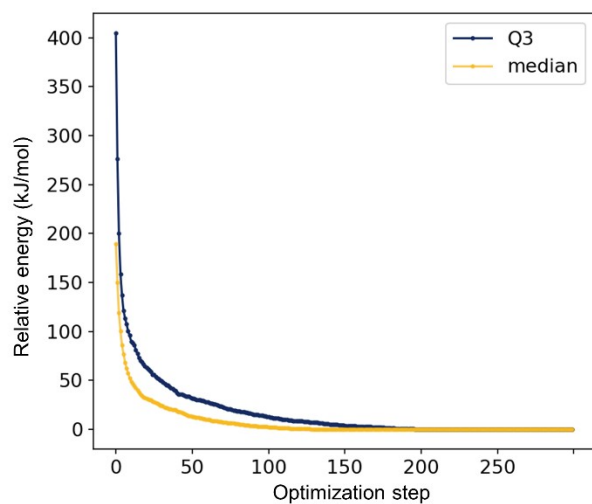
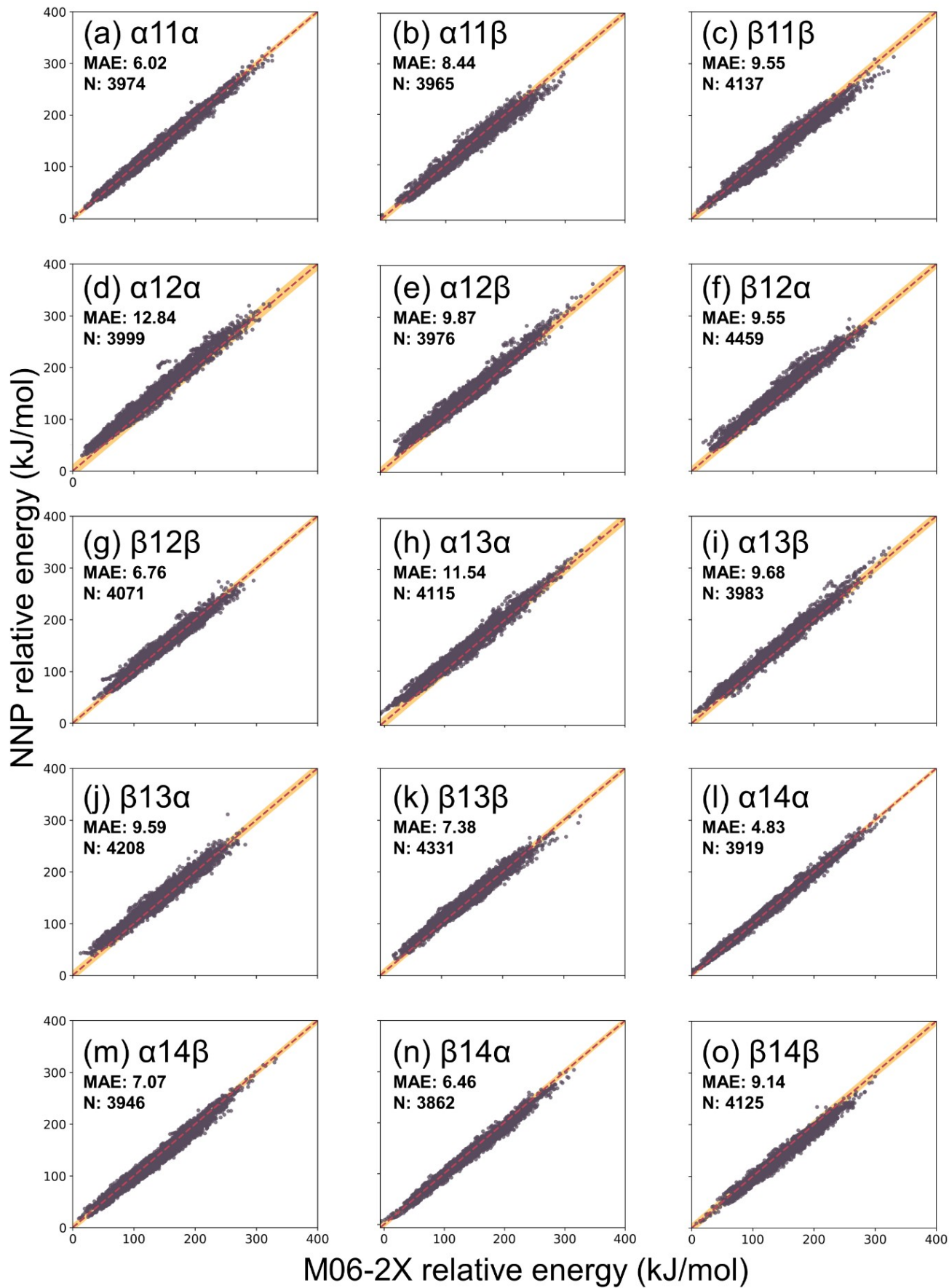


Figure S1. The evolution of the third quartile (Q3) and the median of the potential energy from snapshots extracted from optimization trajectories carried out by the DFTB3 method with the initial geometries of sodiated α -Maltose are created from the “attach-and-rotate” sampling scheme. The energy is relative to the corresponding local minima of the respective trajectory.



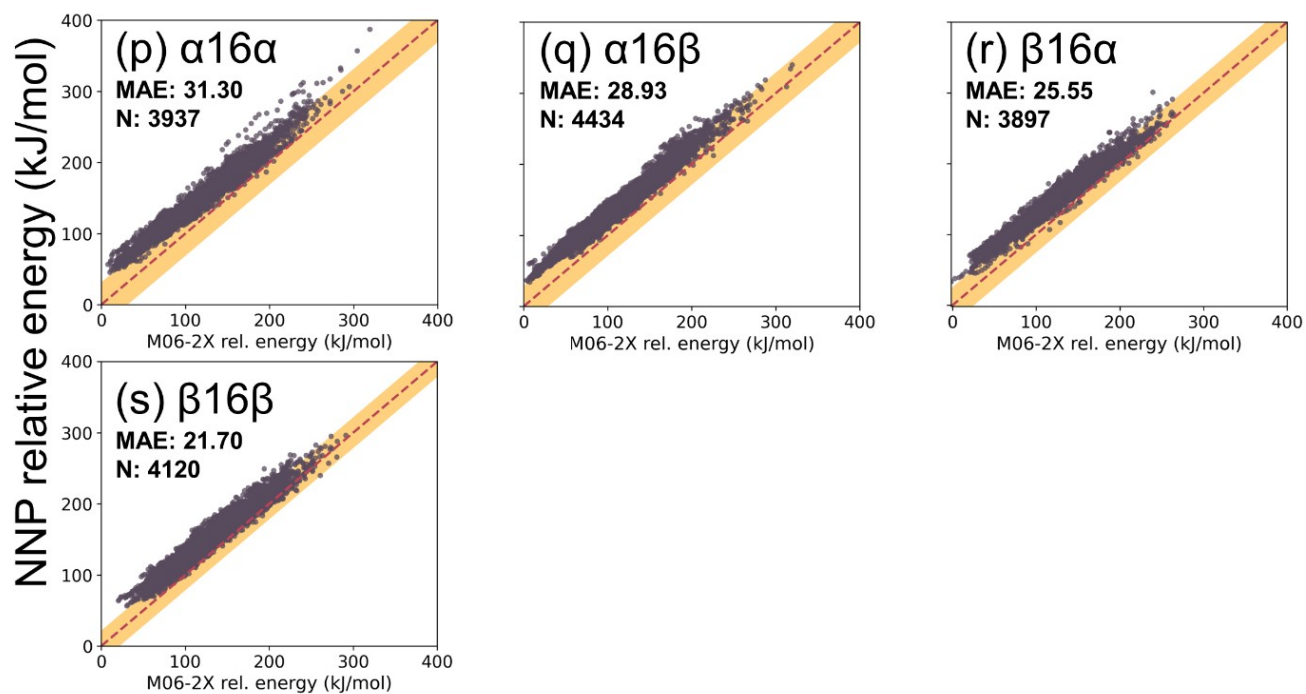
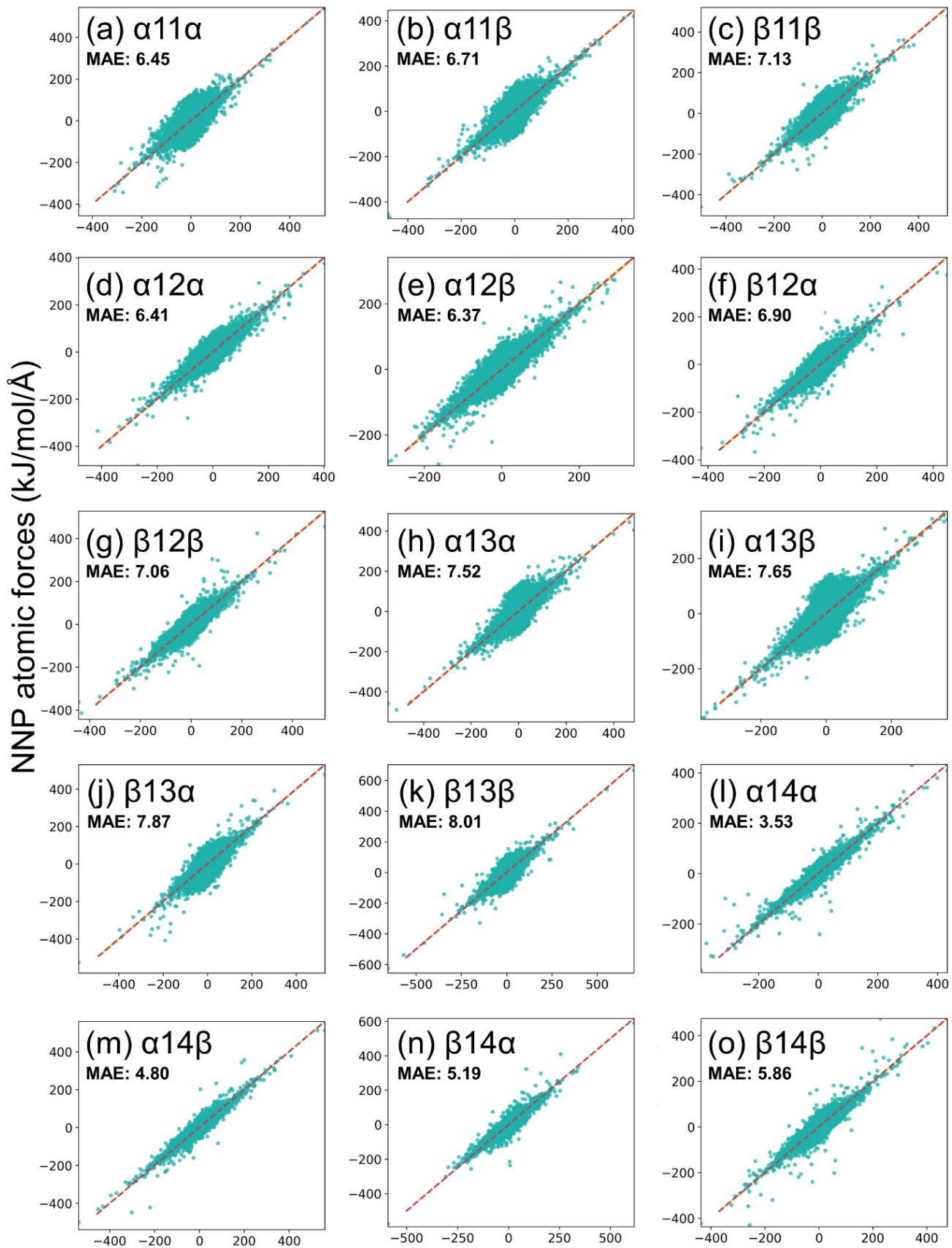


Figure S2. The potential energy correlation between the NNP-0 and DFT M06-2X/6-311+G(d,p) on the test set of each individual sodiated di-saccharides. The zero of energy is set as the energy of the global minimum of sodiated α -Maltose evaluated at DFT level.



M06-2X atomic forces (kJ/mol/Å)

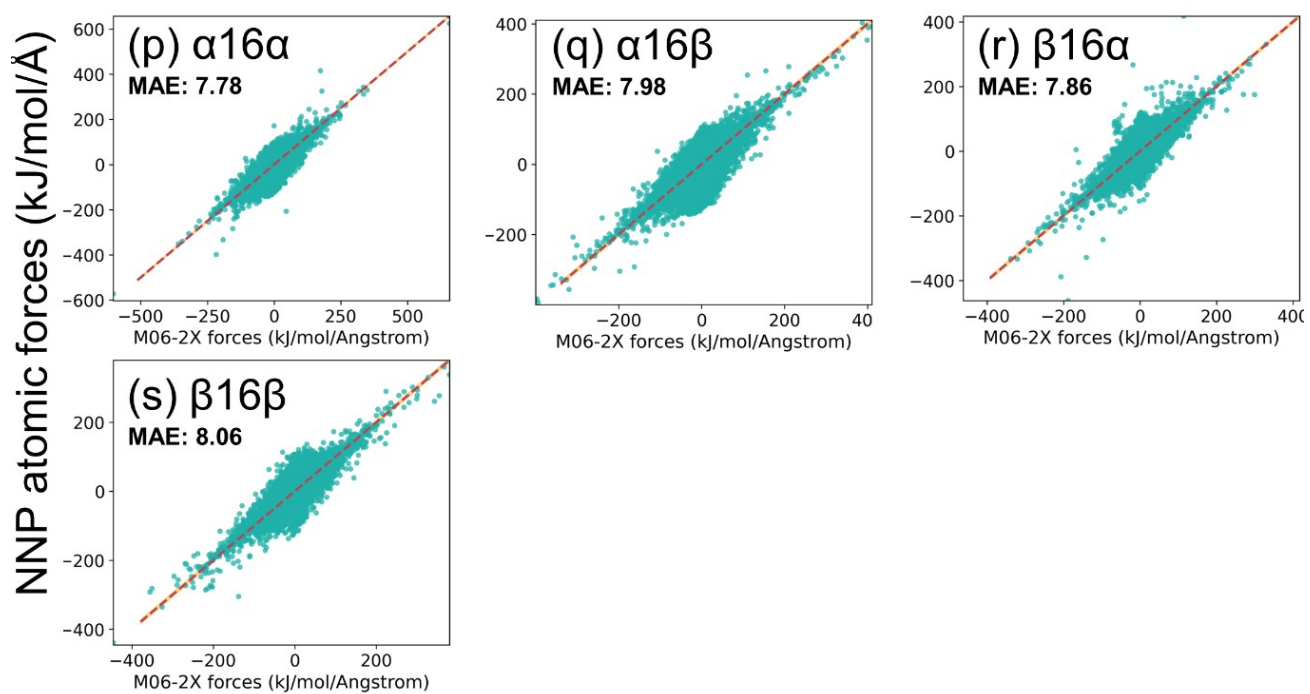


Figure S3. The correlation in the evaluation of atomic forces between the NNP-0 and DFT M06-2X/6-311+G(d,p) on the test set of each sodiated di-saccharide.

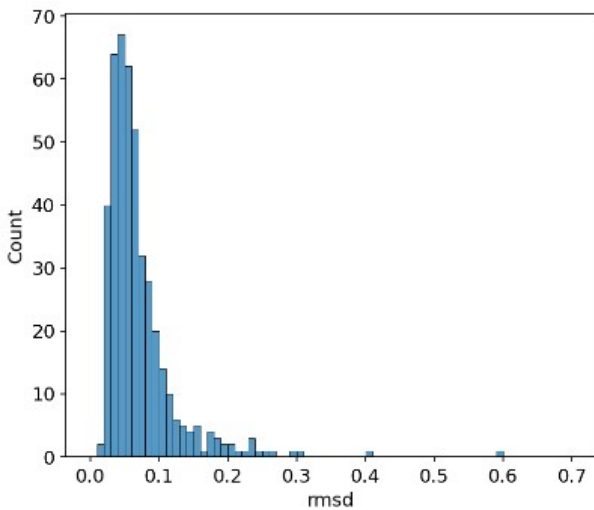
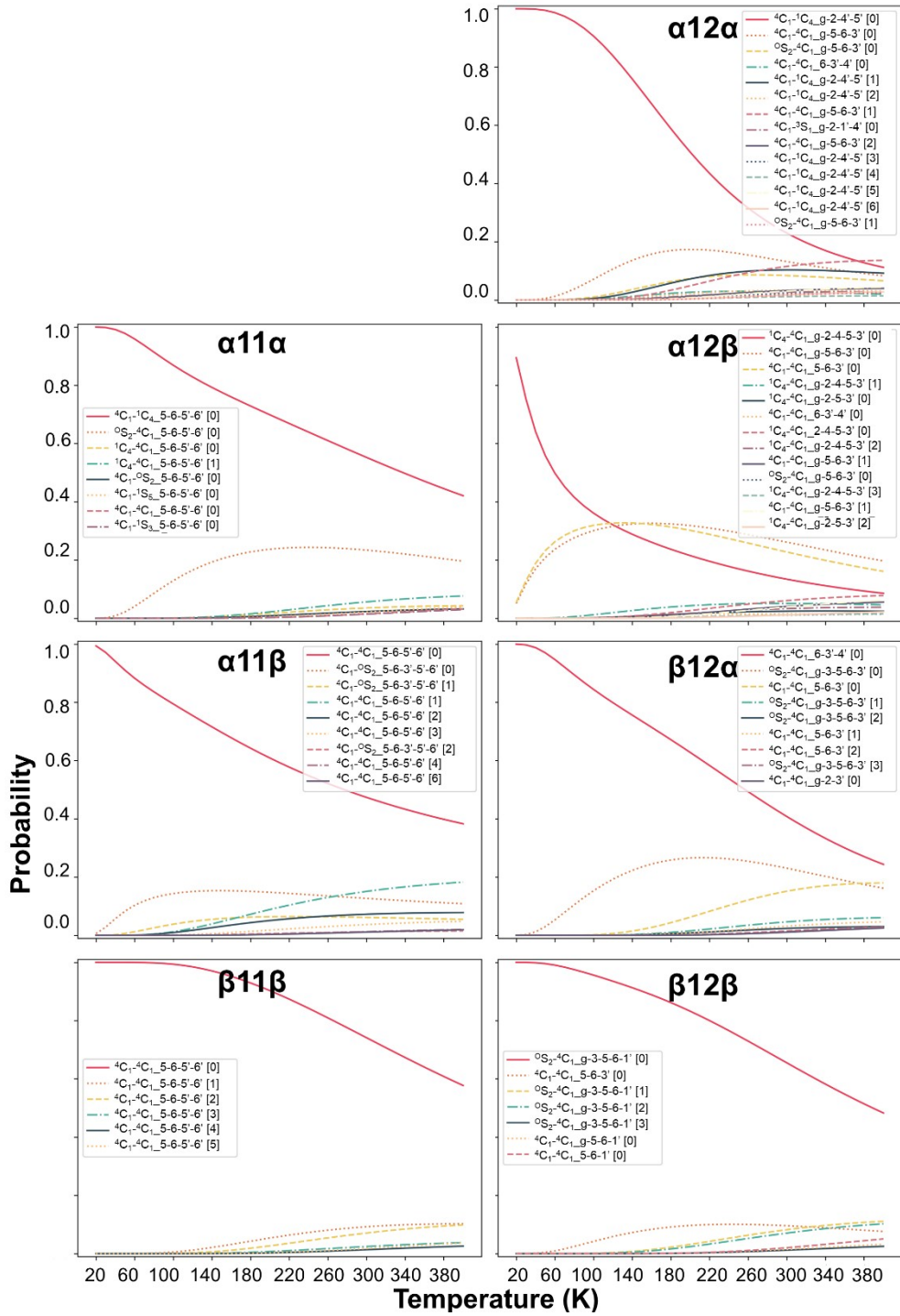


Figure S4. The histogram of RMSD values for pairs of minima in sodiated α -Maltose. Each pair consists of an NNP minima and a corresponding M06-2X minima which obtained via the re-optimization from NNP minima.



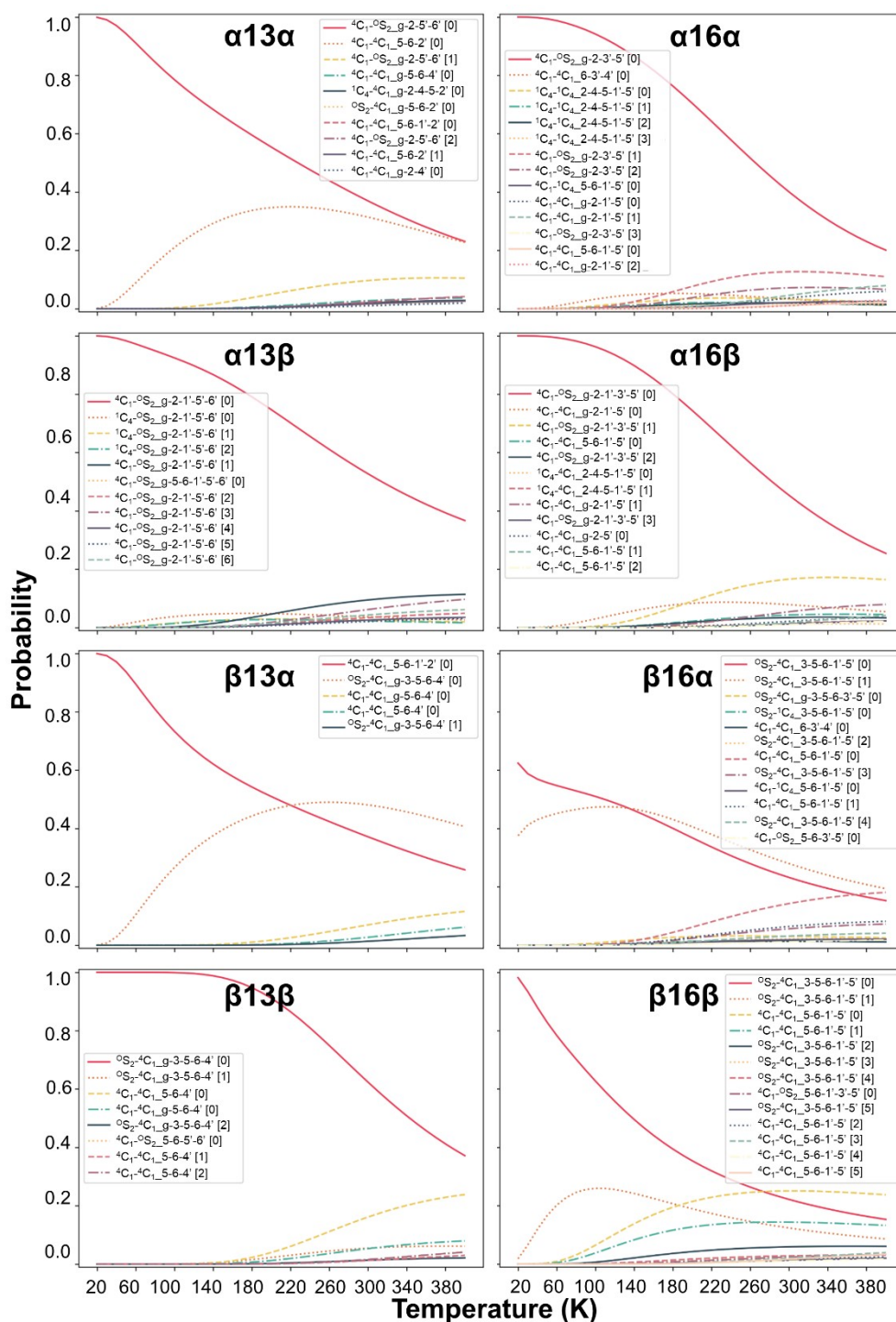
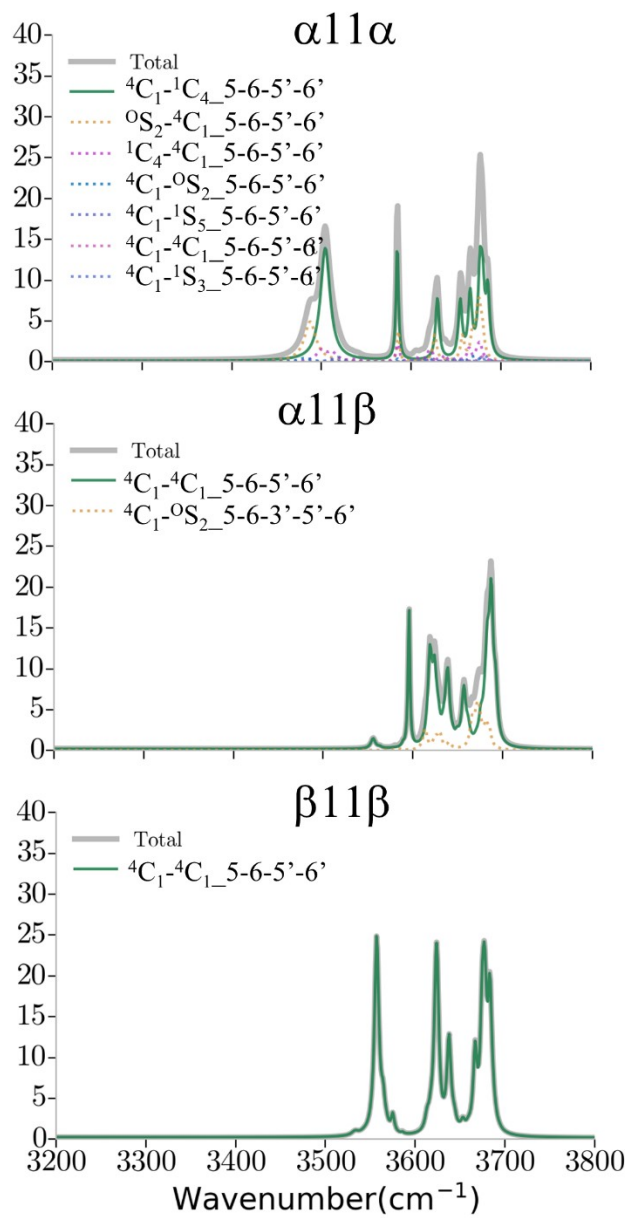
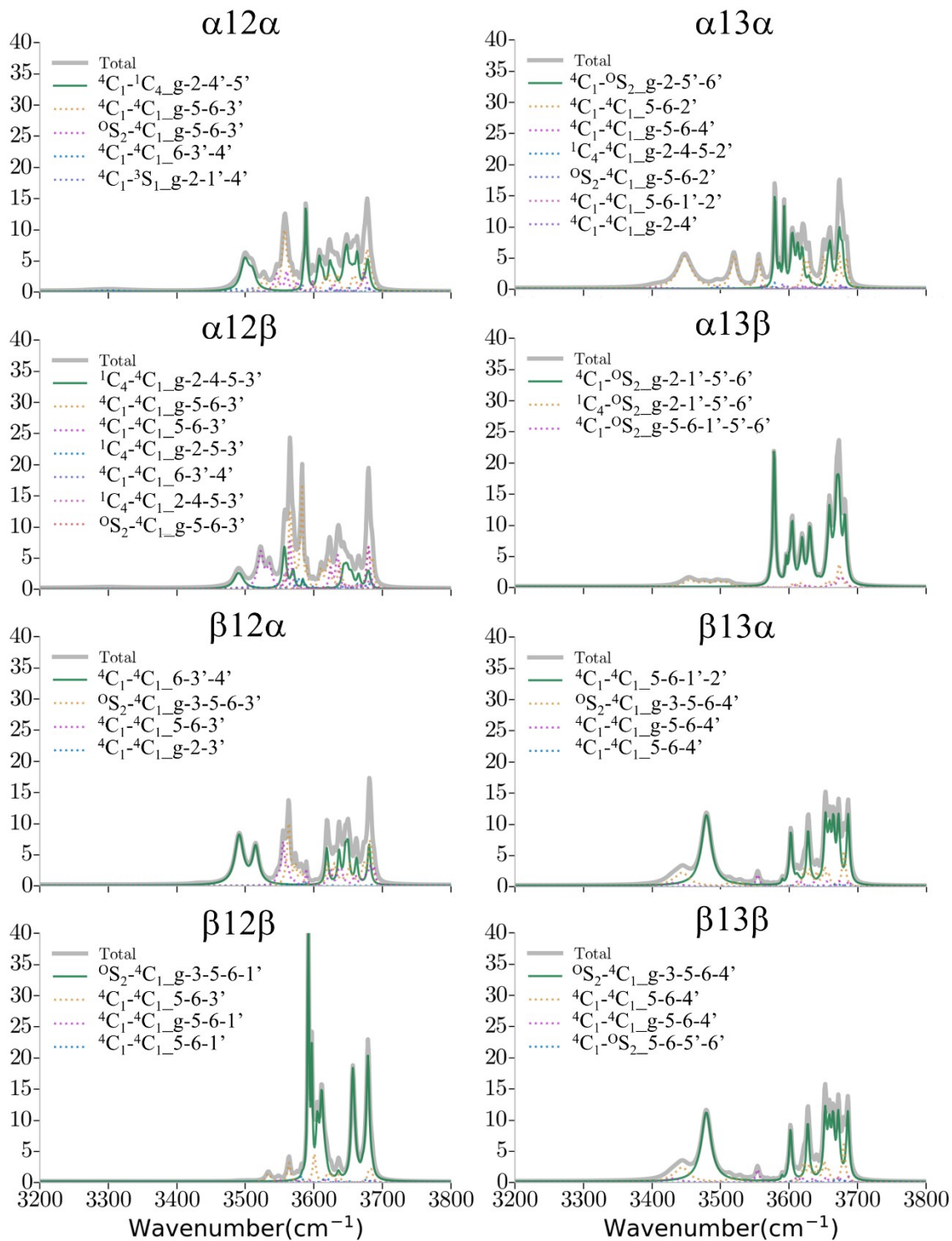


Figure S5. The temperature dependence of the relative population of low-energy conformers of sodiated di-saccharides 1-1, 1-2, 1-3, and 1-6 linkage. The number in the square bracket “[]” indicates the energy rank of the local minima sharing the same structural features. For example, “[0]” is used

denote the conformer with the lowest energy, while “[1]” is used to represent the conformer with the second lowest energy.





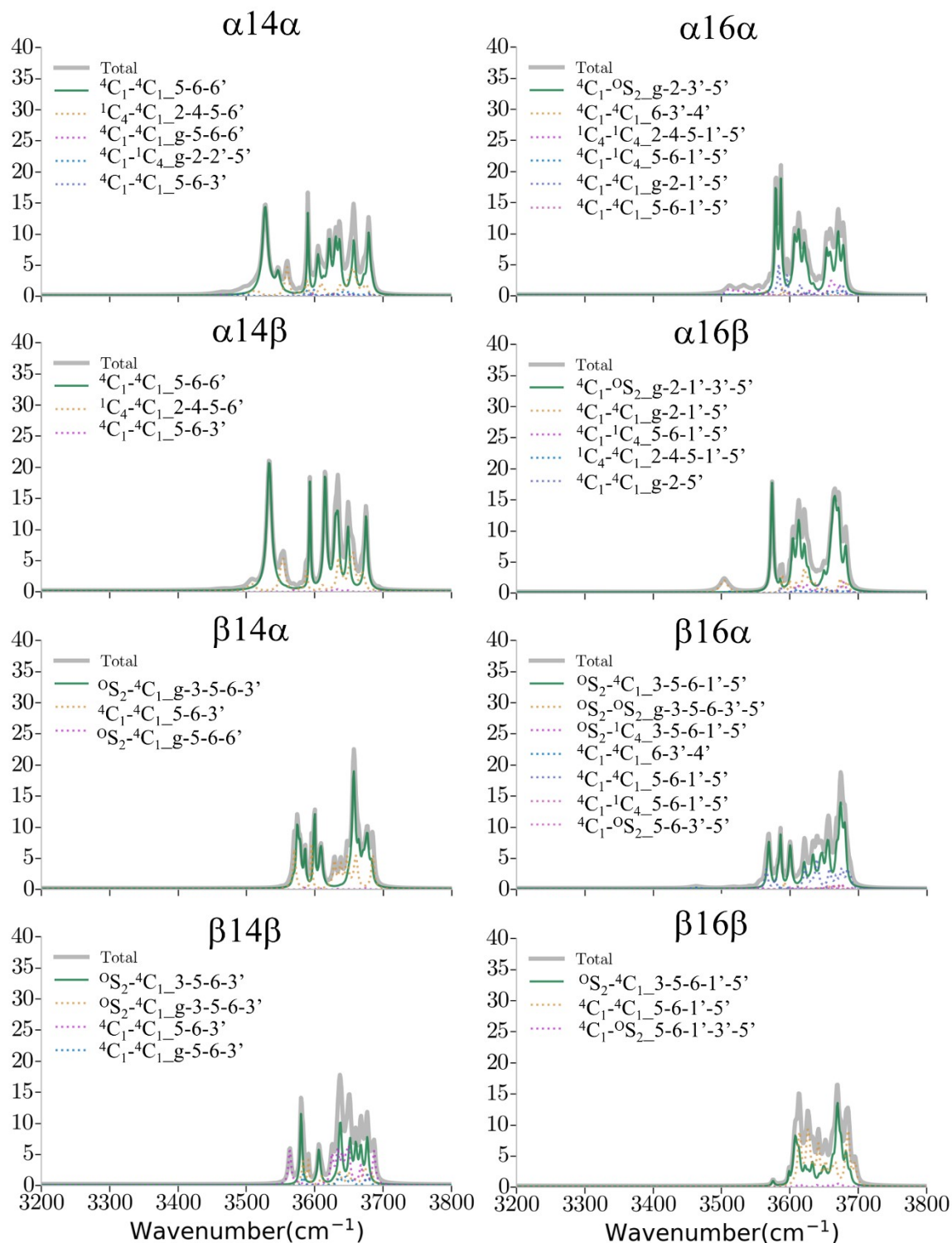


Figure S6. The accumulated Q-HSA spectra at 300K of all 19 sodiated Glc-based disaccharides. The total accumulated spectra are shown in gray line, and the spectra from local minima sharing similar conformations are depicted in various colors and styles.

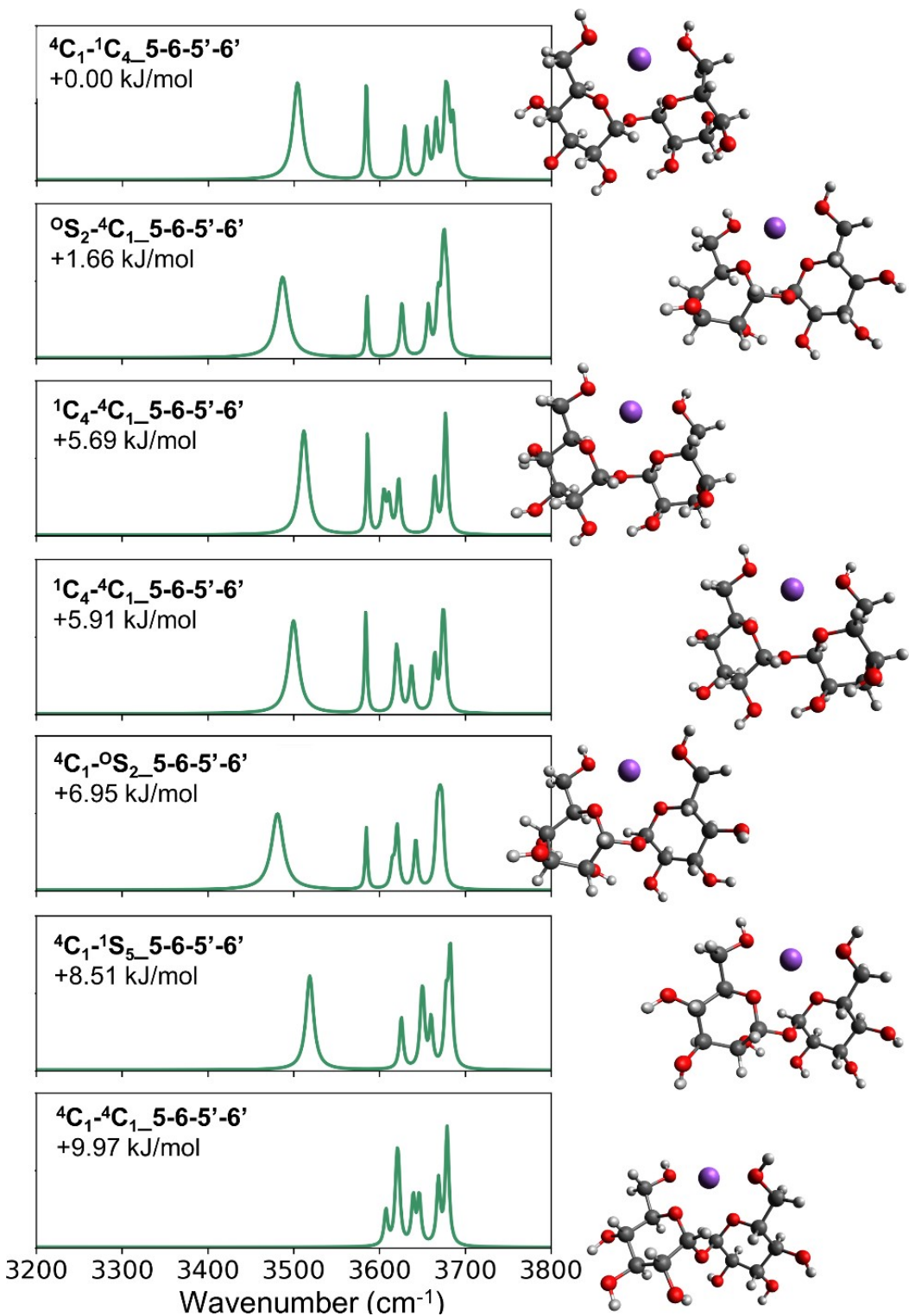


Figure S7. The vibrational spectra of sodiated α Glc-(1-1)- α Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

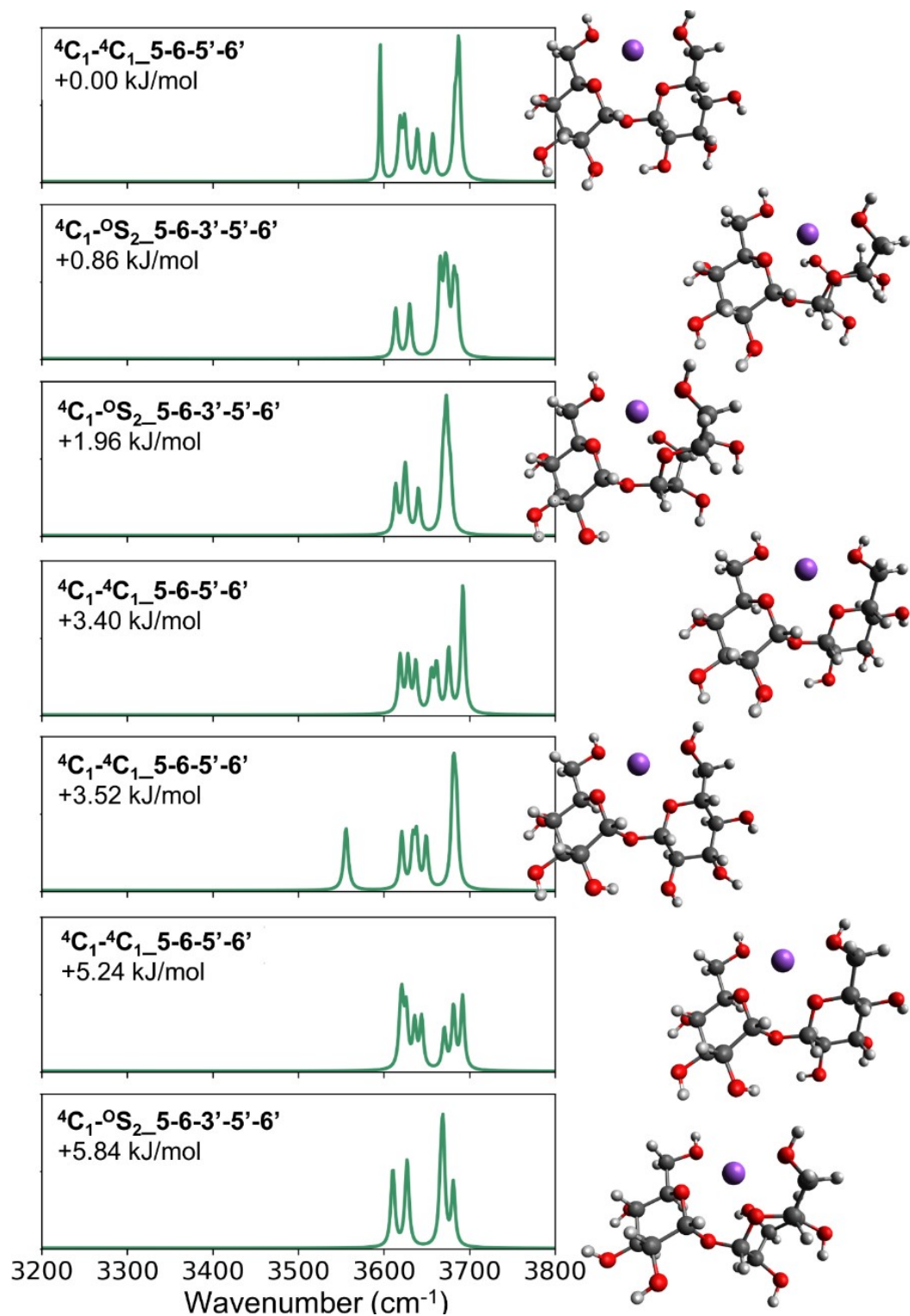


Figure S8. The vibrational spectra of sodiated $\alpha\text{Glc-(1-1)-}\beta\text{Glc}$ conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

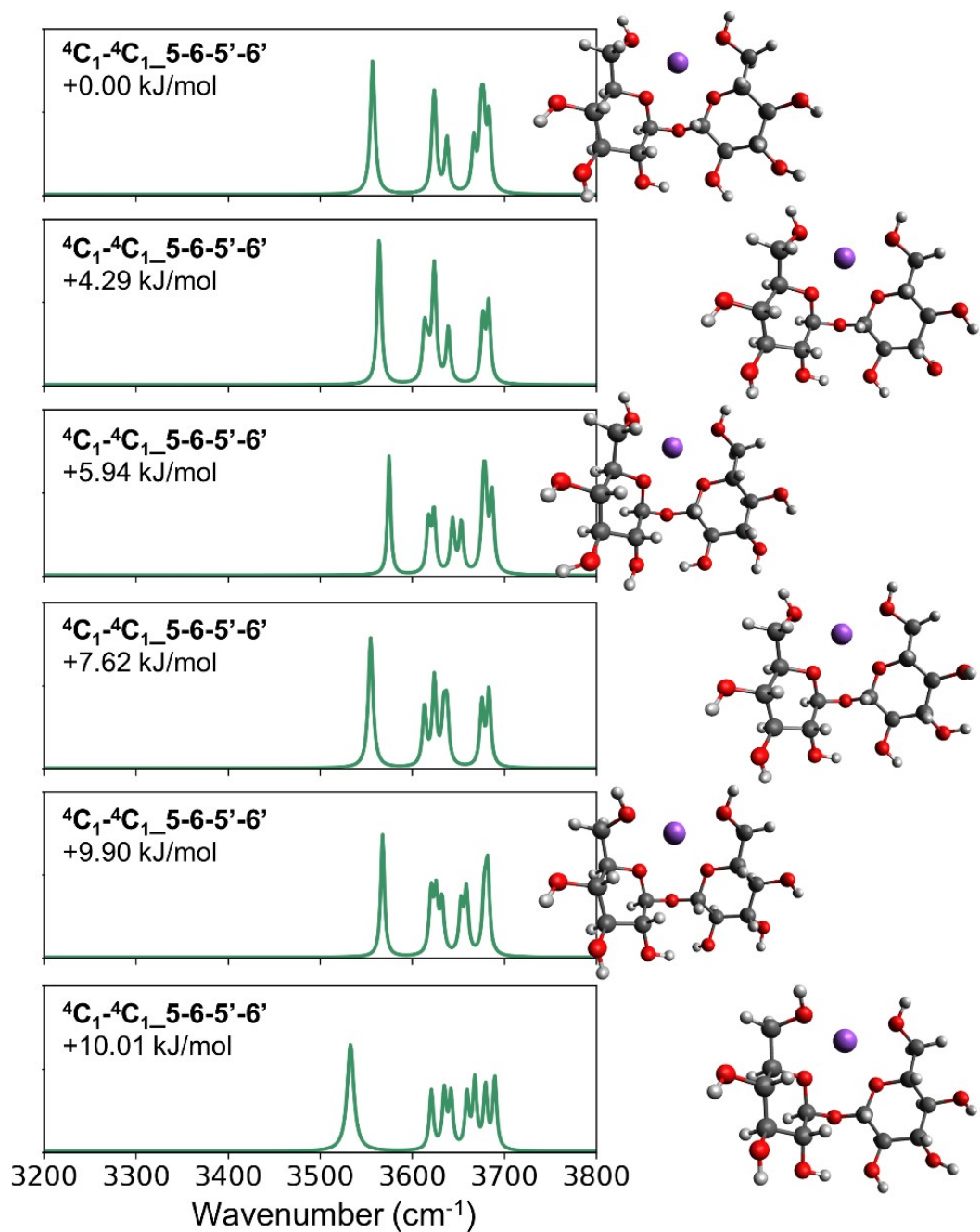


Figure S9. The vibrational spectra of sodiated $\beta\text{Glc-(1-1)-}\beta\text{Glc}$ conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

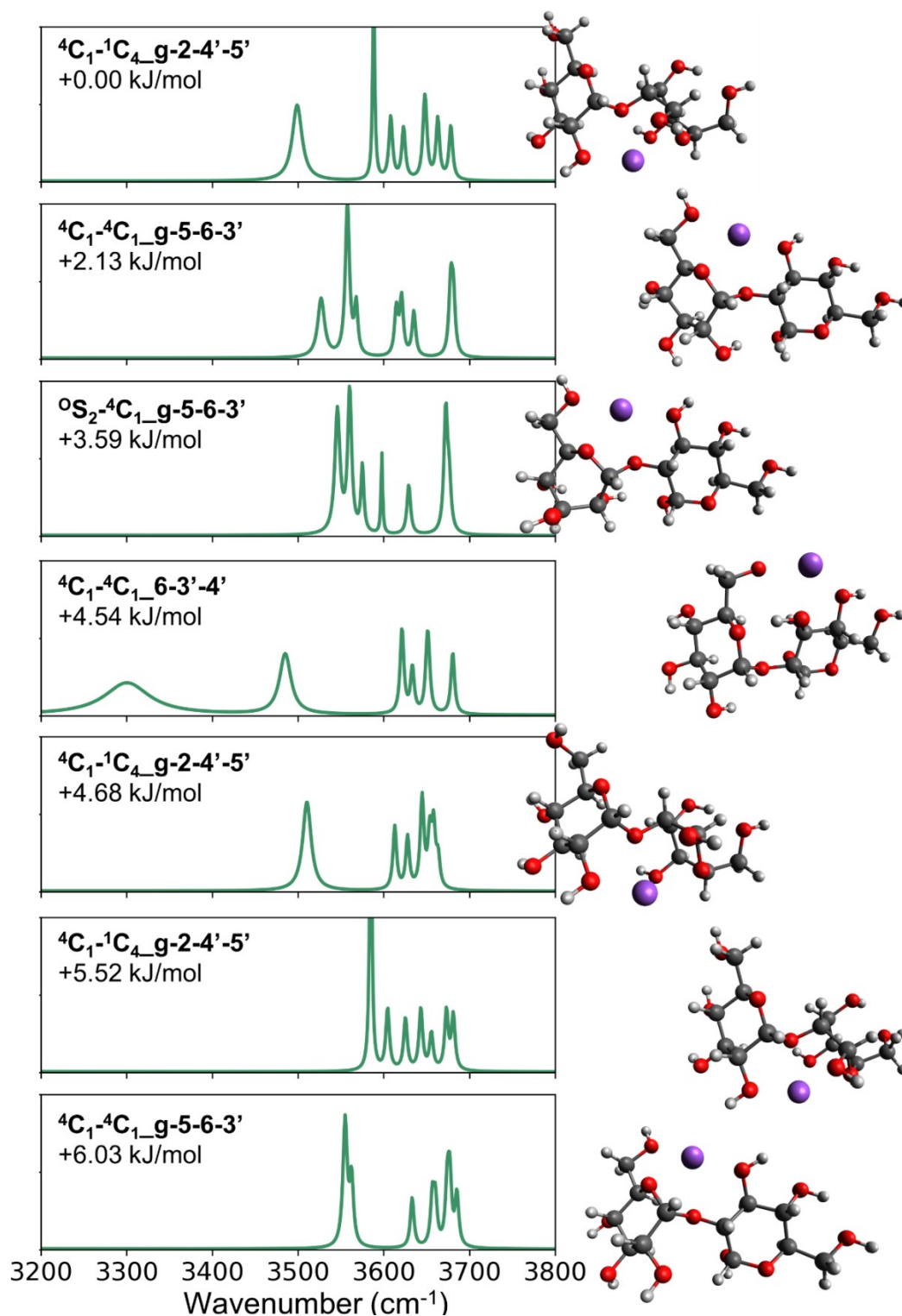


Figure S10. The vibrational spectra of sodiated α Glc-(1-2)- α Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

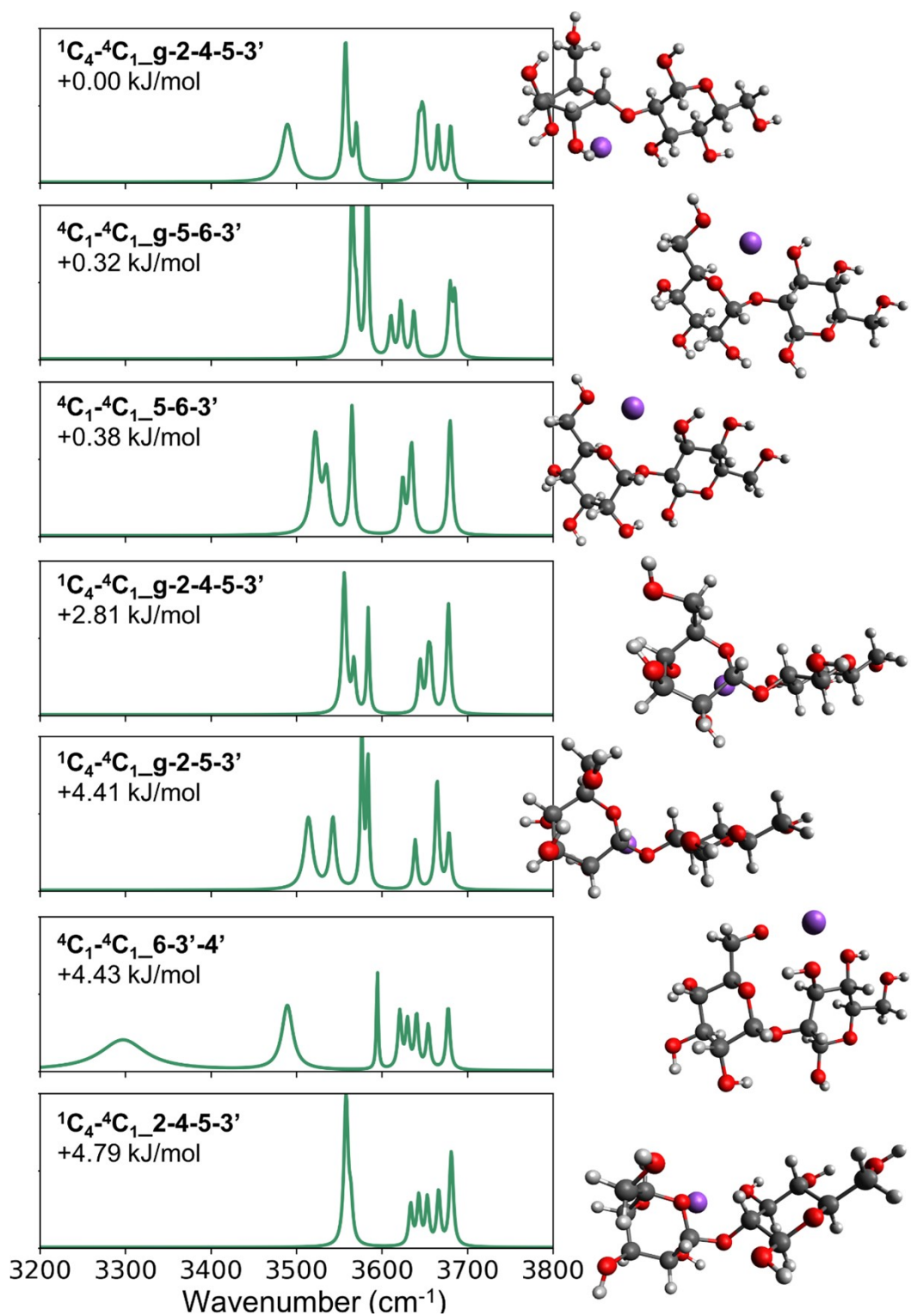


Figure S11. The vibrational spectra of sodiated α Glc-(1-2)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

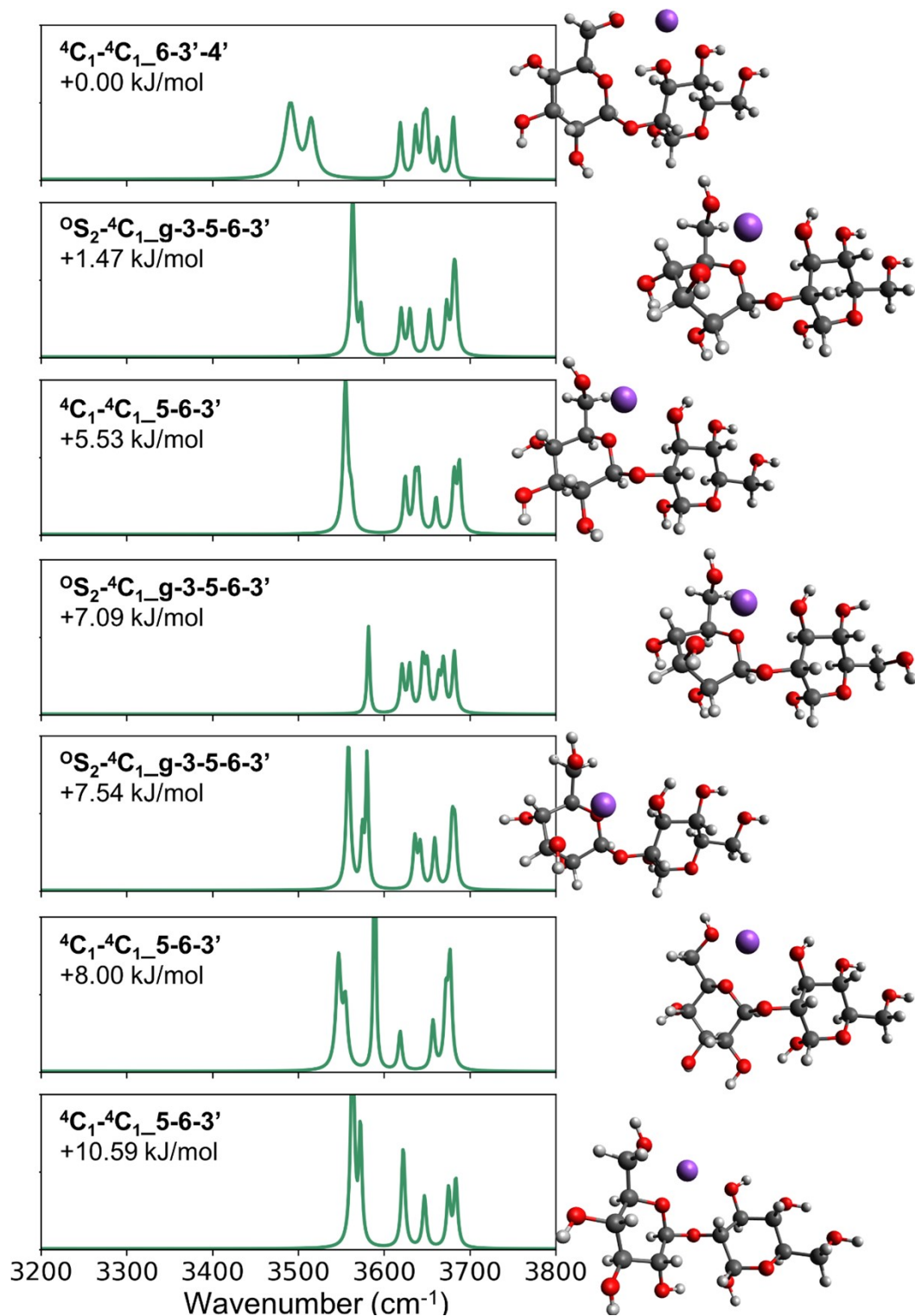


Figure S12. The vibrational spectra of sodiated β Glc-(1-2)- α Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

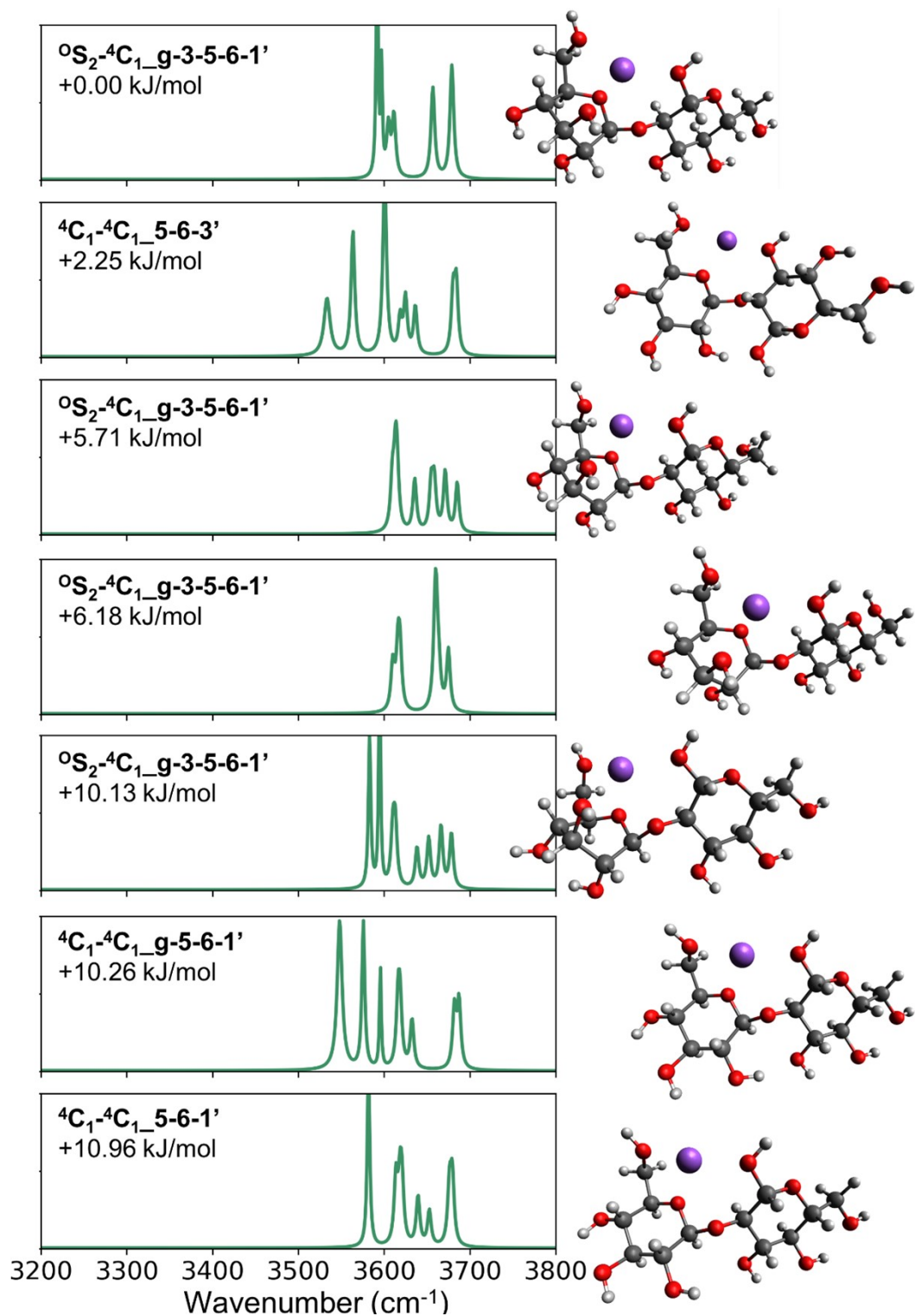


Figure S13. The vibrational spectra of sodiated β Glc-(1-2)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

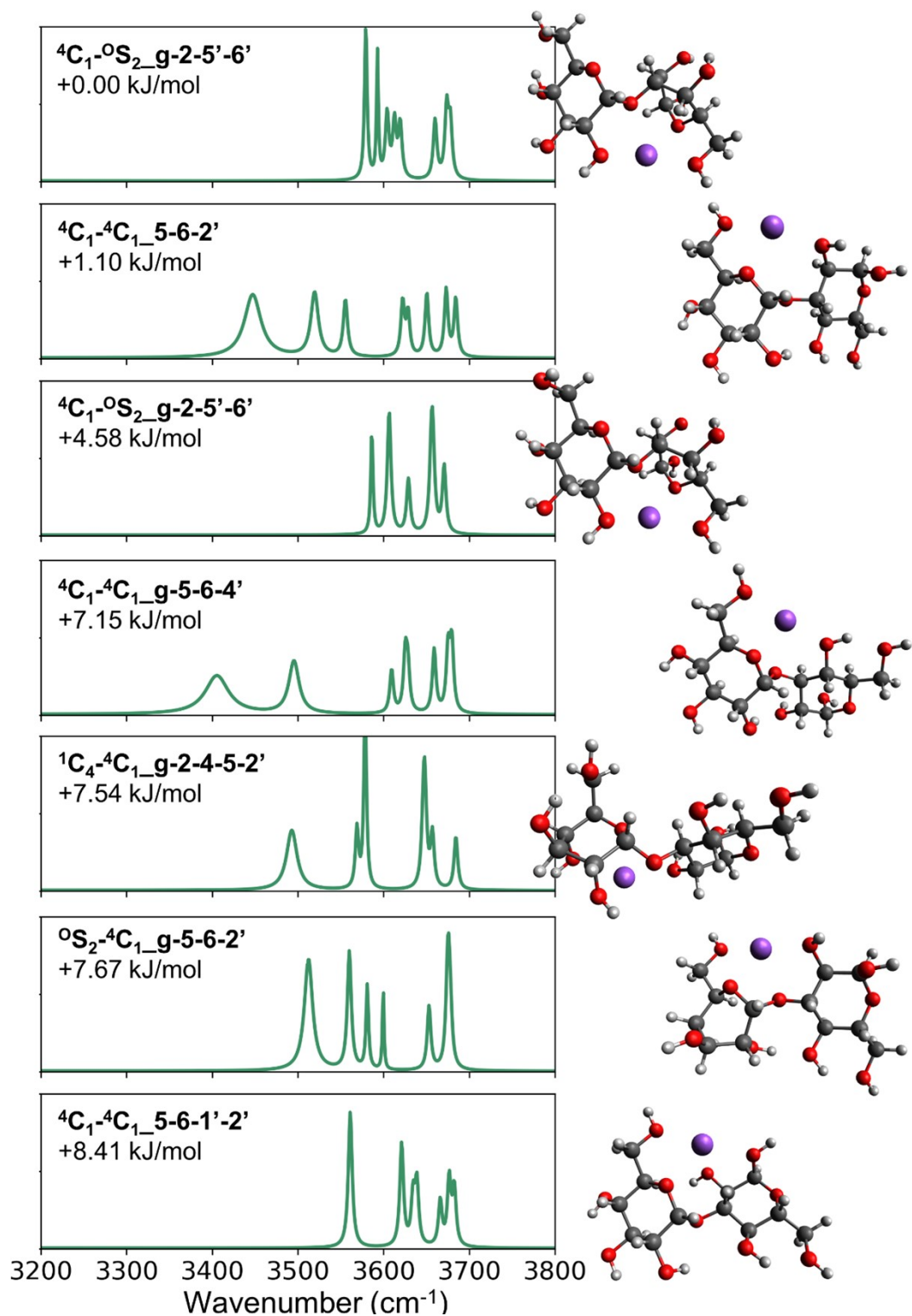


Figure S14. The vibrational spectra of sodiated α Glc-(1-3)- α Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

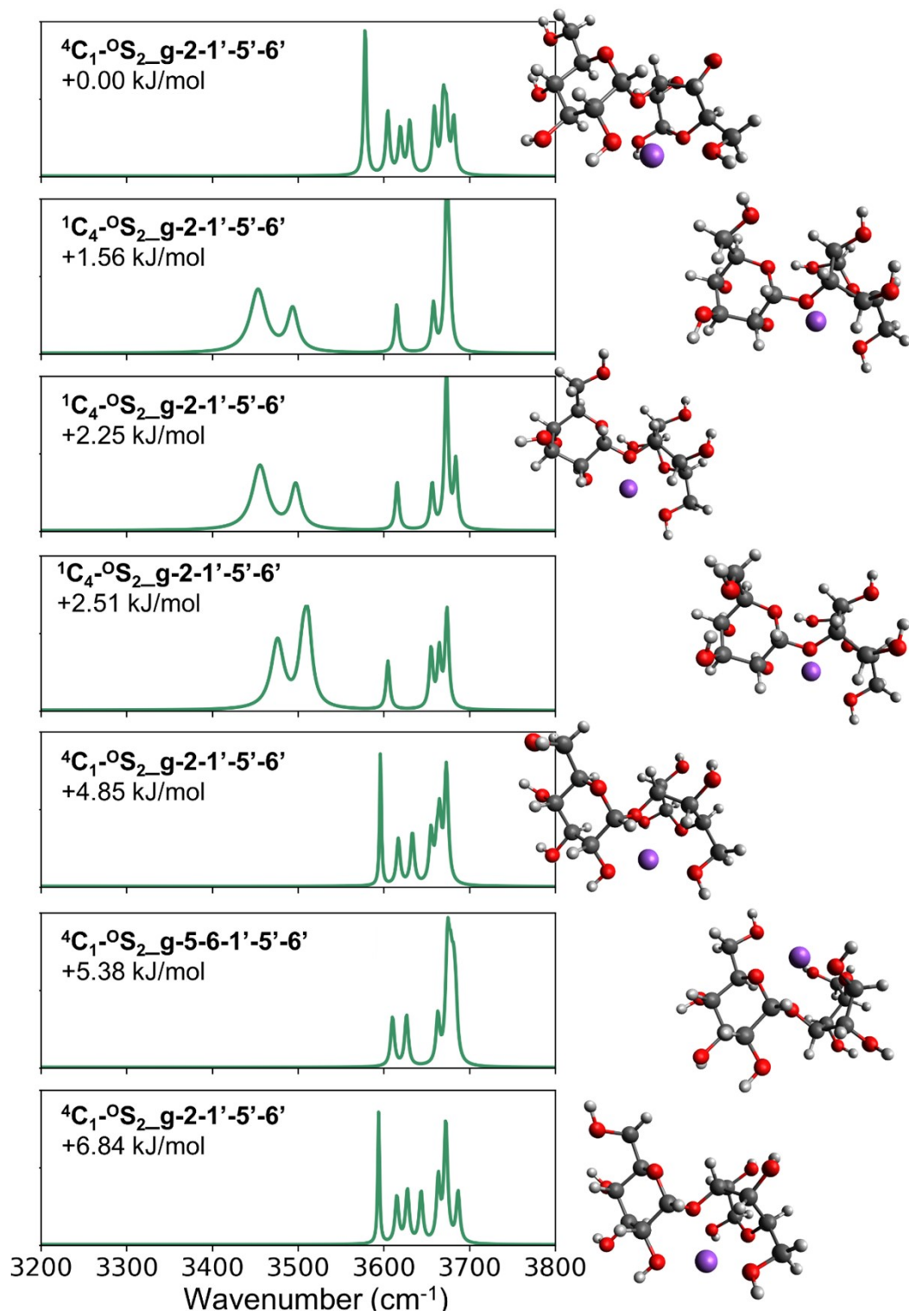


Figure S15. The vibrational spectra of sodiated α Glc-(1-3)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

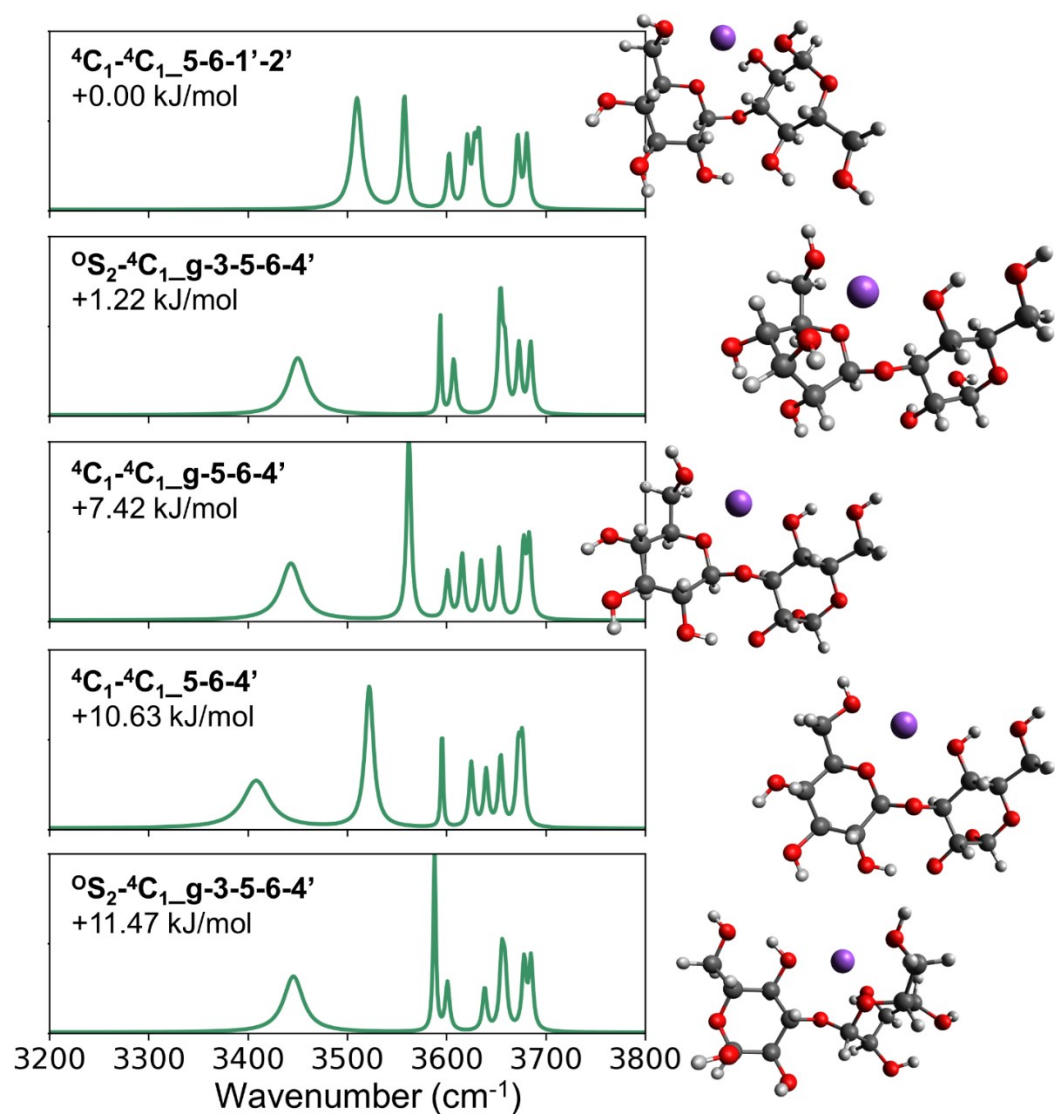


Figure S16. The vibrational spectra of sodiated β Glc-(1-3)- α Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

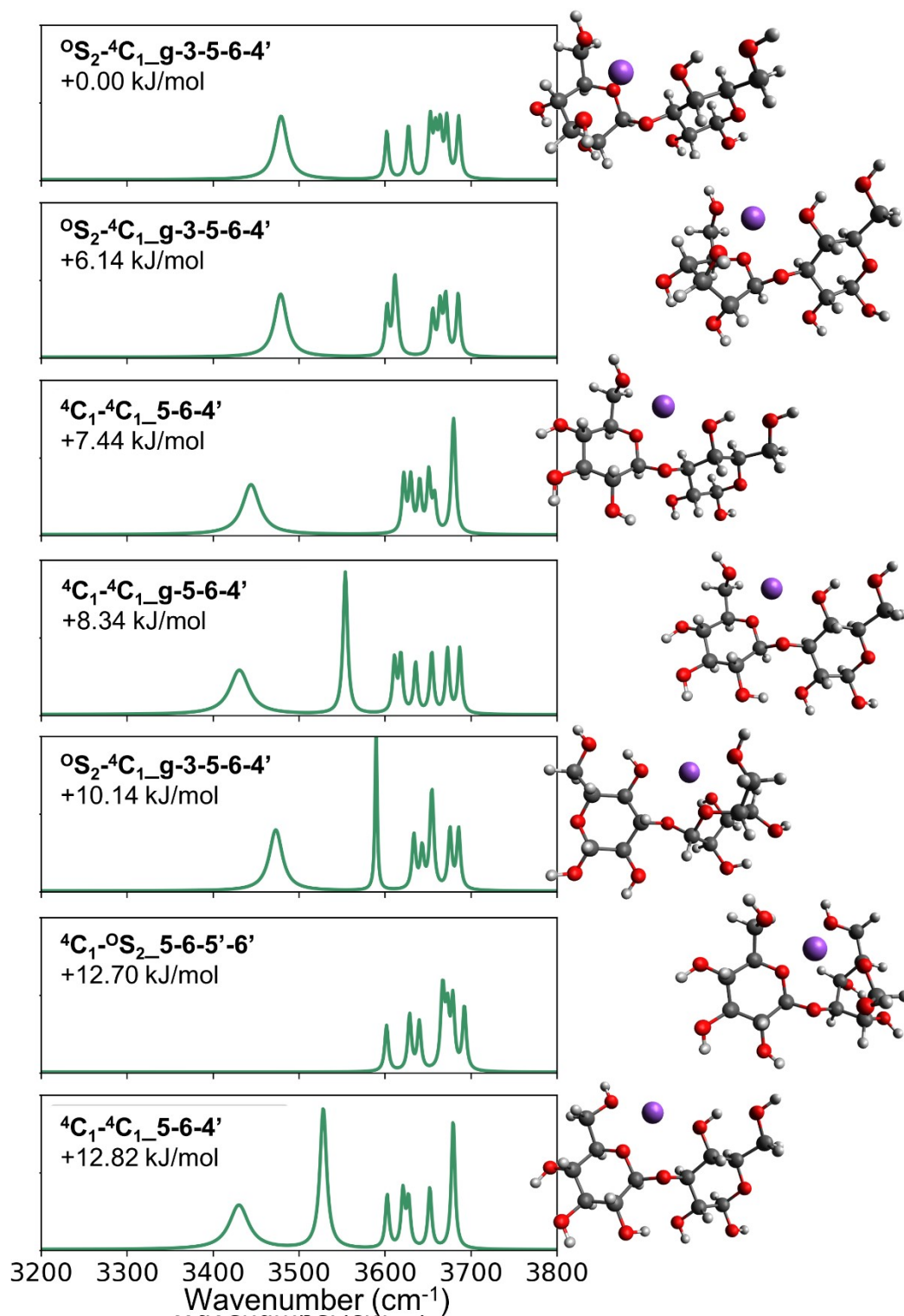


Figure S17. The vibrational spectra of sodiated β Glc-(1-3)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

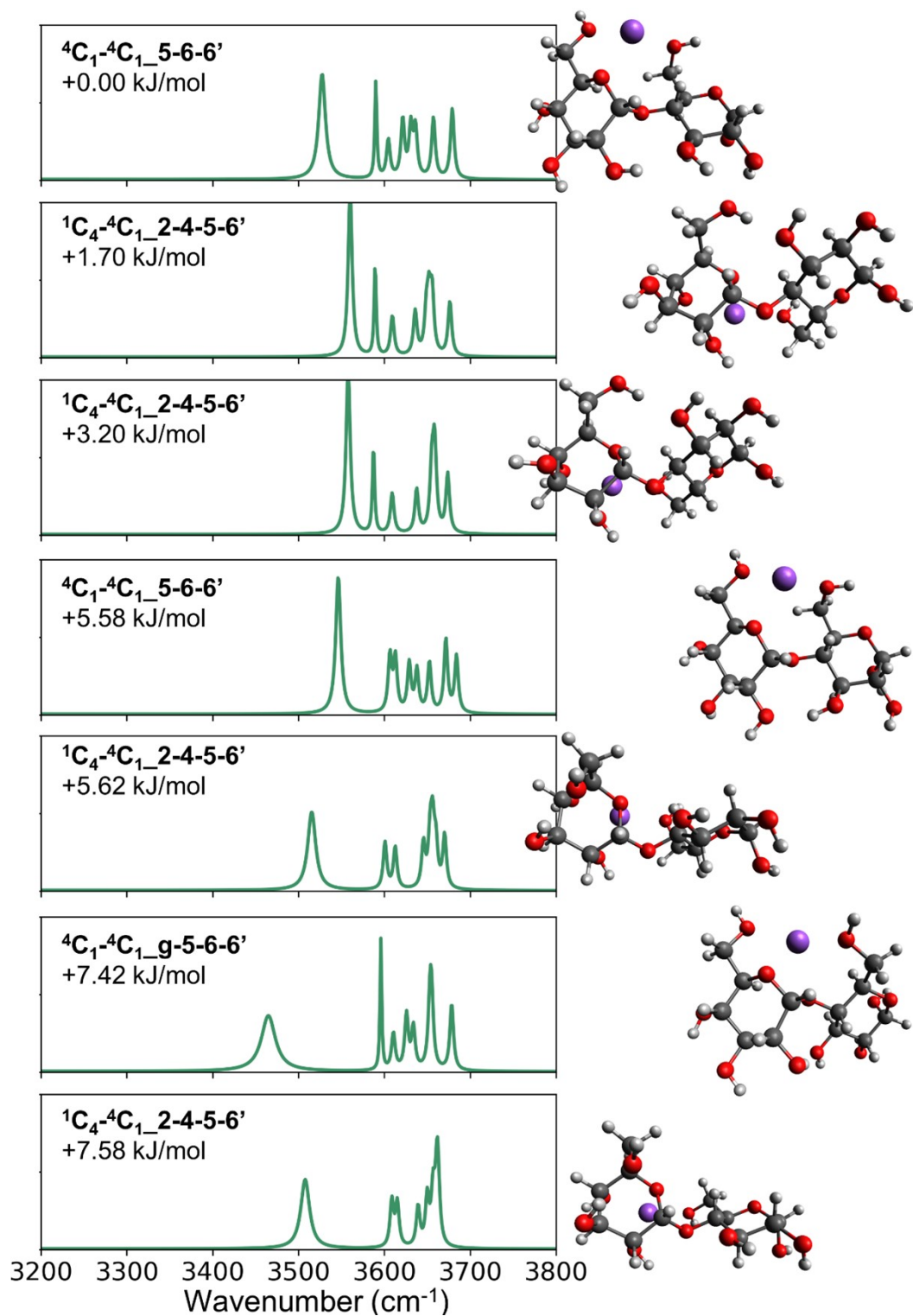


Figure S18. The vibrational spectra of sodiated α Glc-(1-4)- α Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

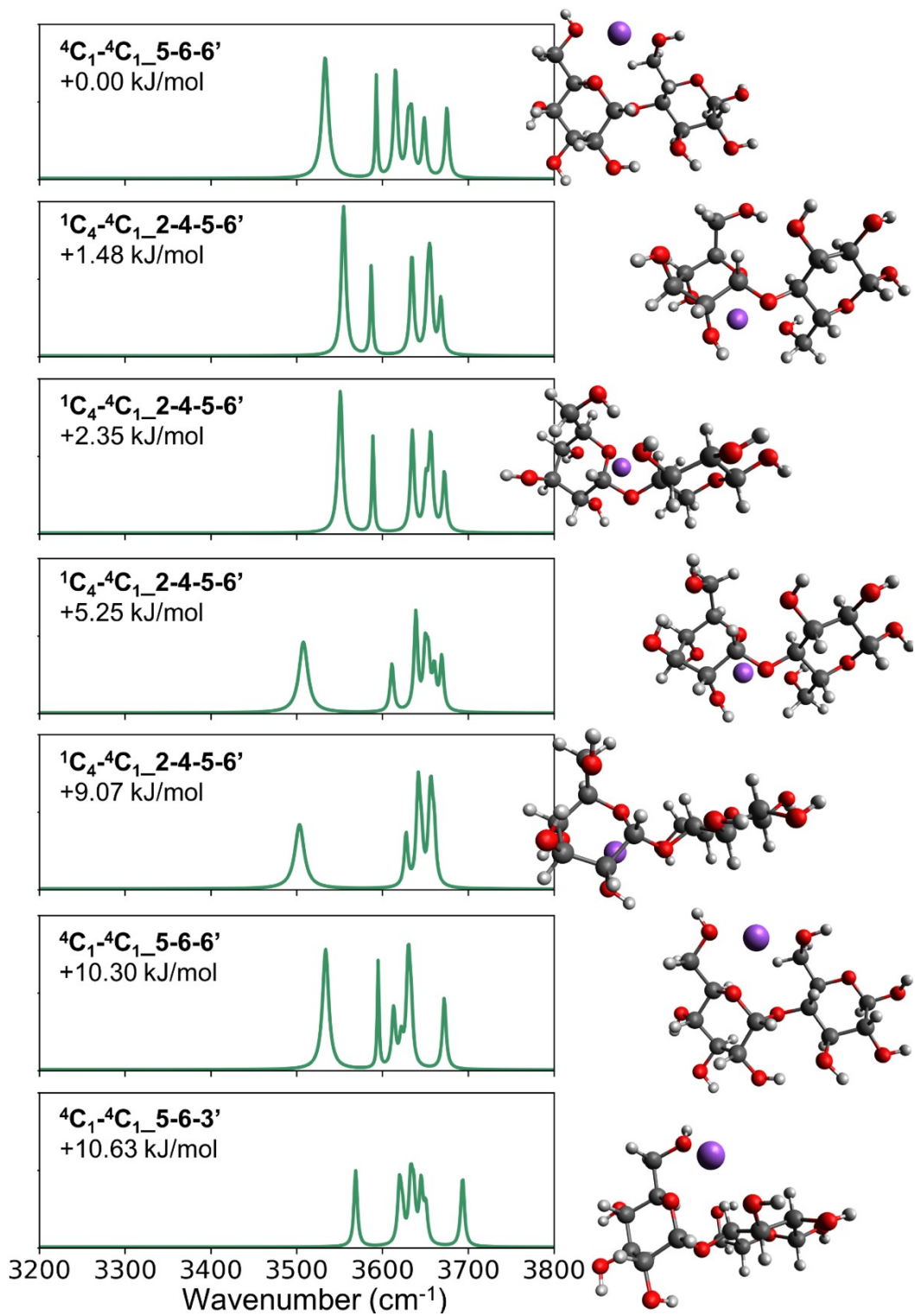


Figure S19. The vibrational spectra of sodiated α Glc-(1-4)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

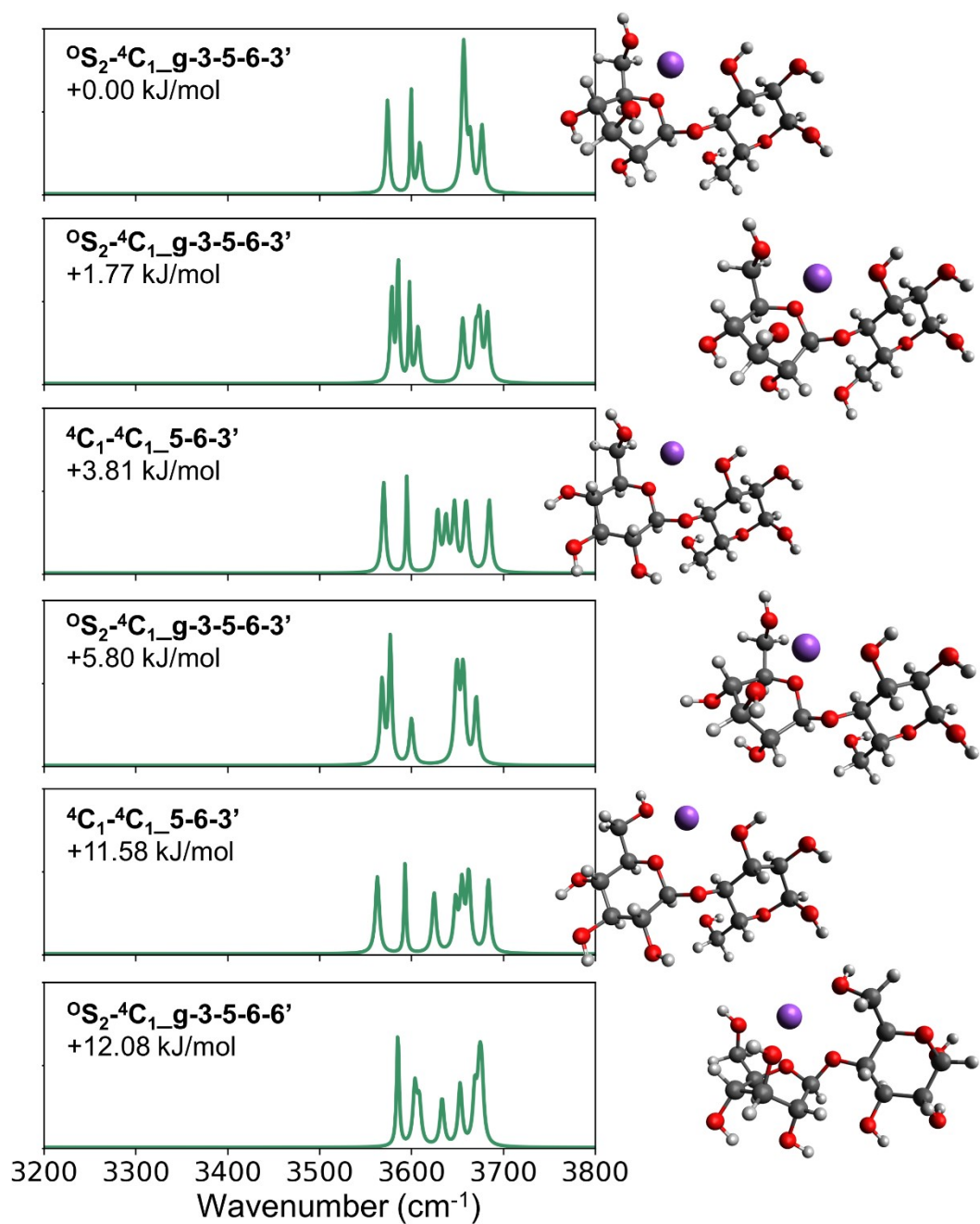


Figure S20. The vibrational spectra of sodiated $\beta\text{Glc-(1-4)-}\alpha\text{Glc}$ conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

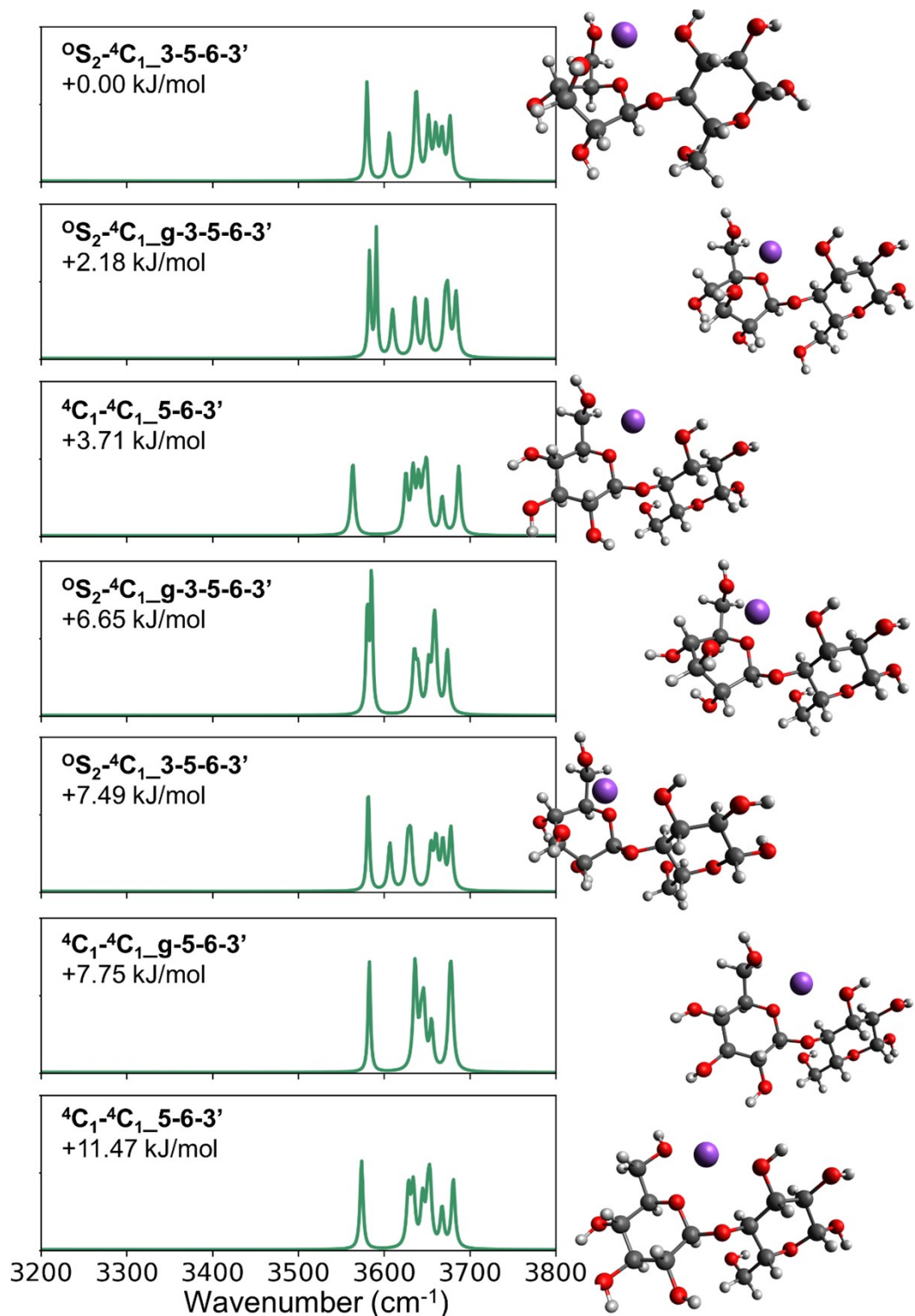


Figure S21. The vibrational spectra of sodiated β Glc-(1-4)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

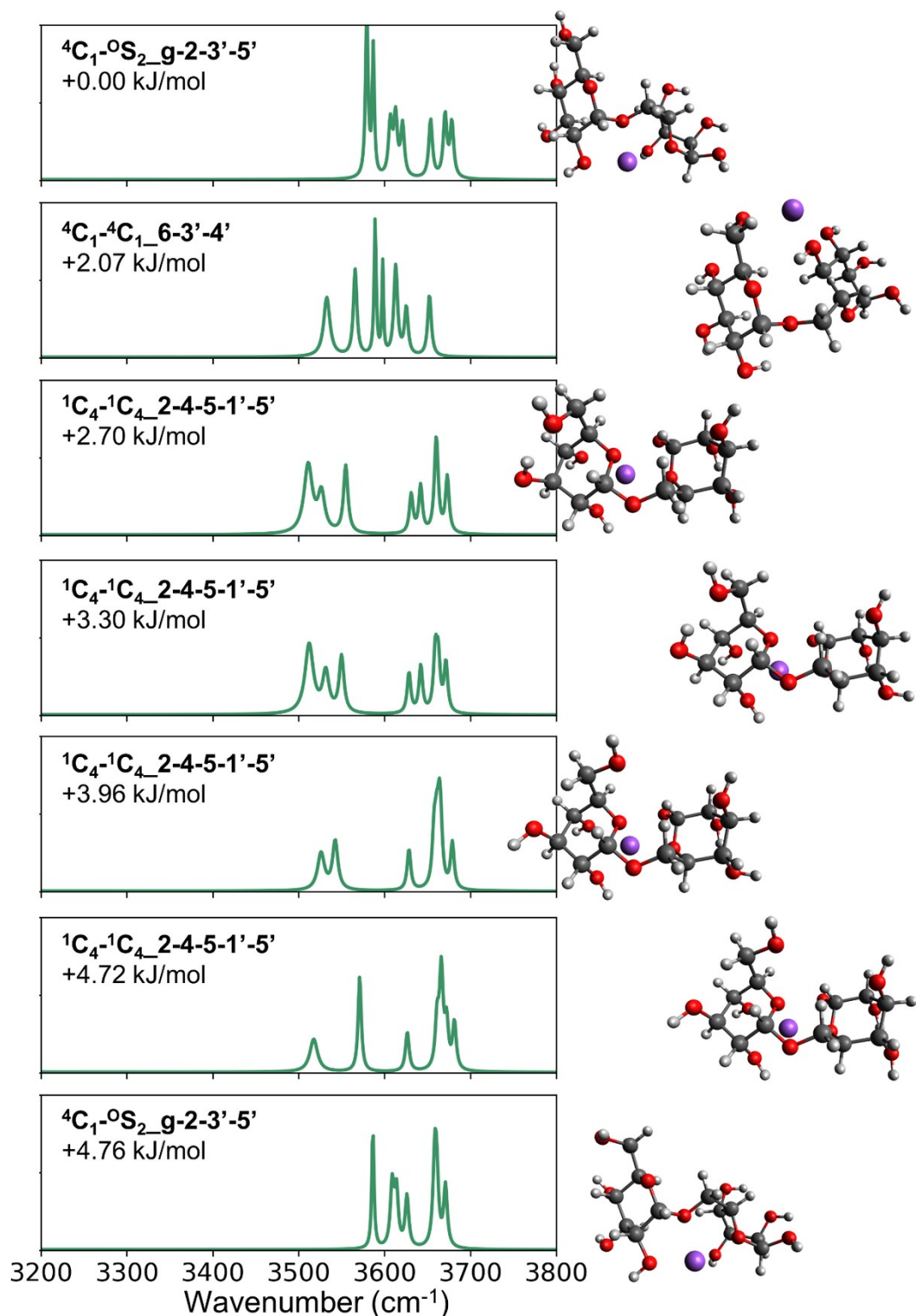


Figure S22. The vibrational spectra of sodiated $\alpha\text{Glc-(1-6)-}\alpha\text{Glc}$ conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

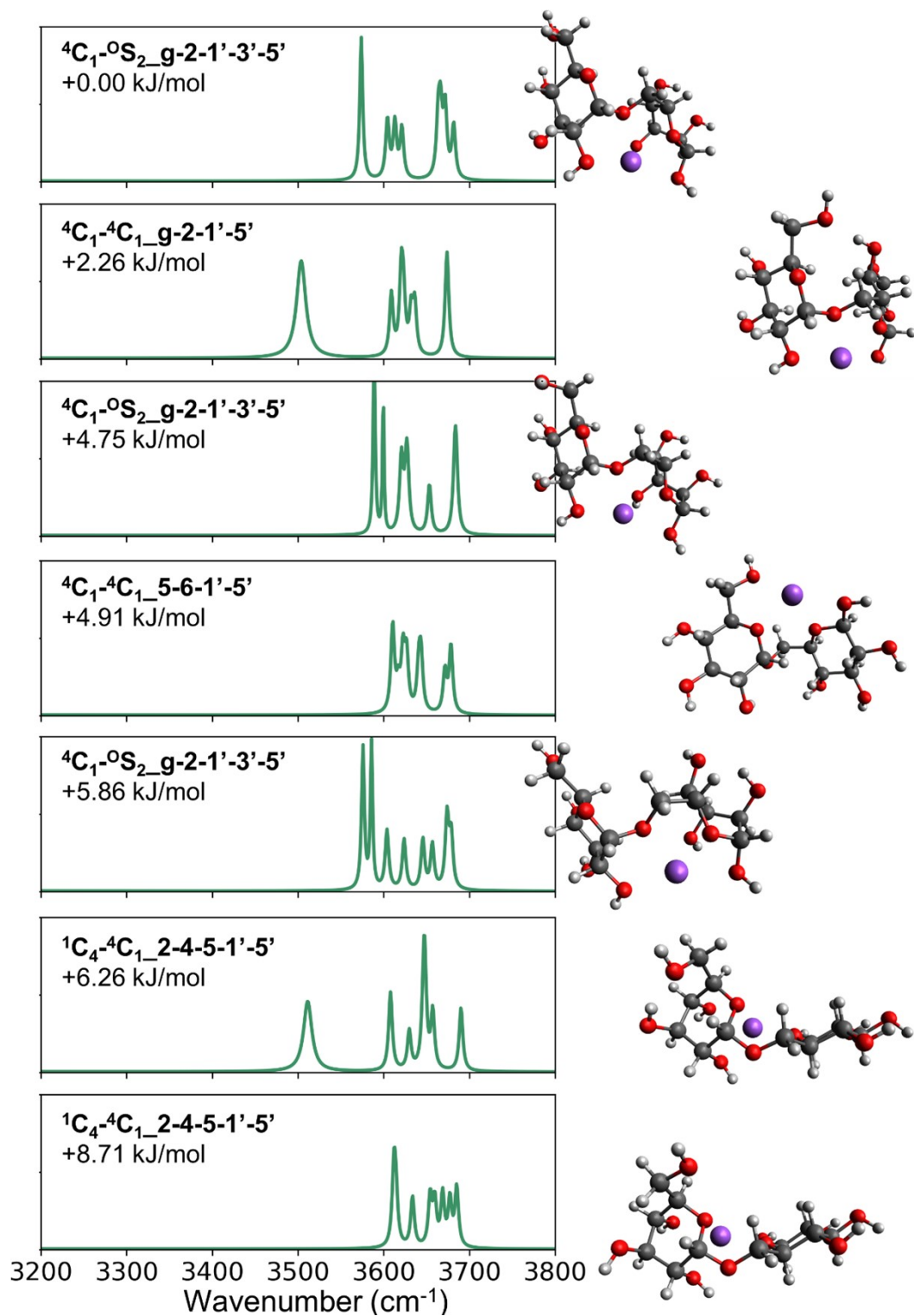


Figure S23. The vibrational spectra of sodiated α Glc-(1-6)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

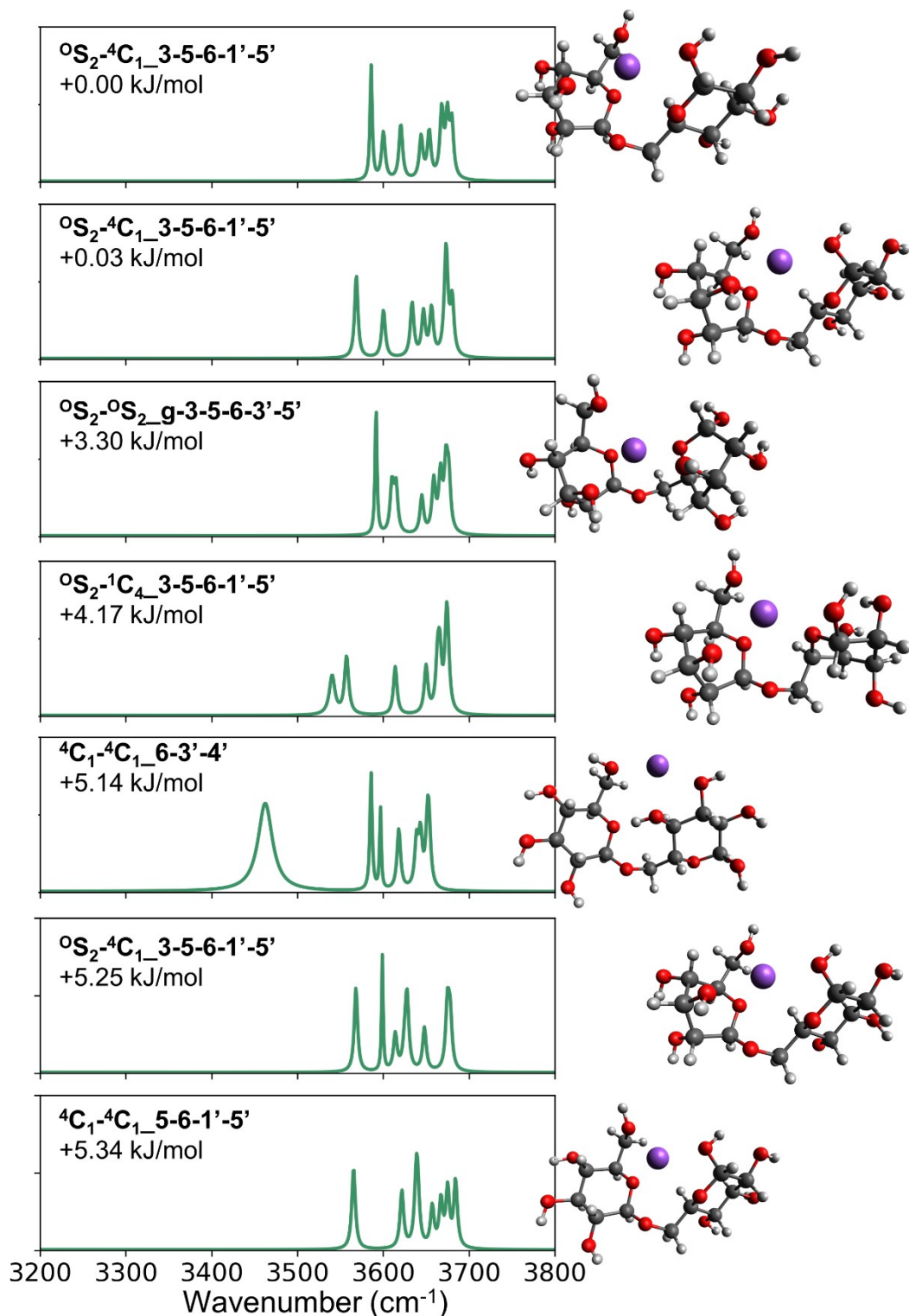


Figure S24. The vibrational spectra of sodiated β Glc-(1-6)- α Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

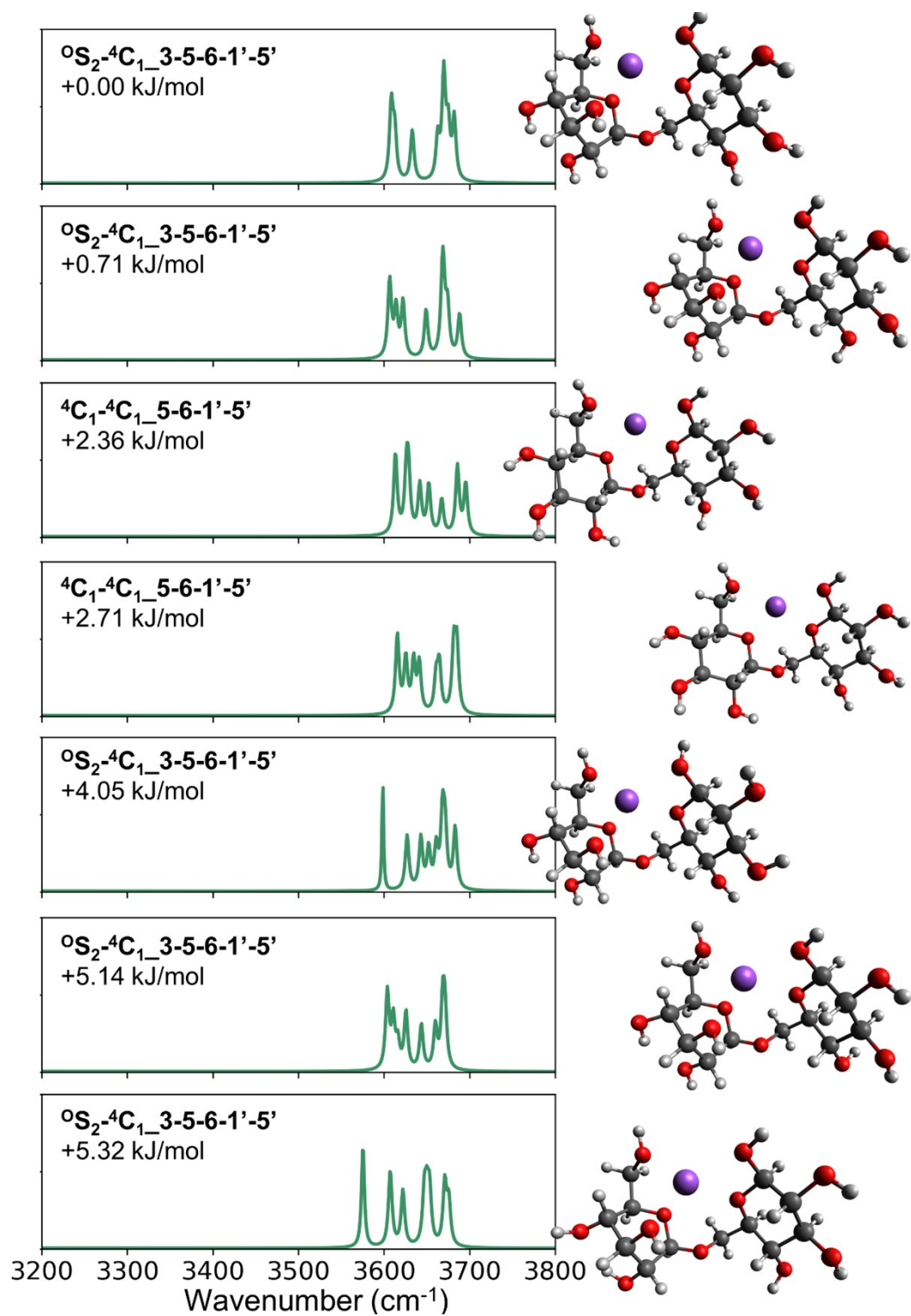


Figure S25. The vibrational spectra of sodiated β Glc-(1-6)- β Glc conformers of significant population derived by Q-HSA analysis at 300K, with conformation name and the relative energy values indicated at the top left corner of each spectrum.

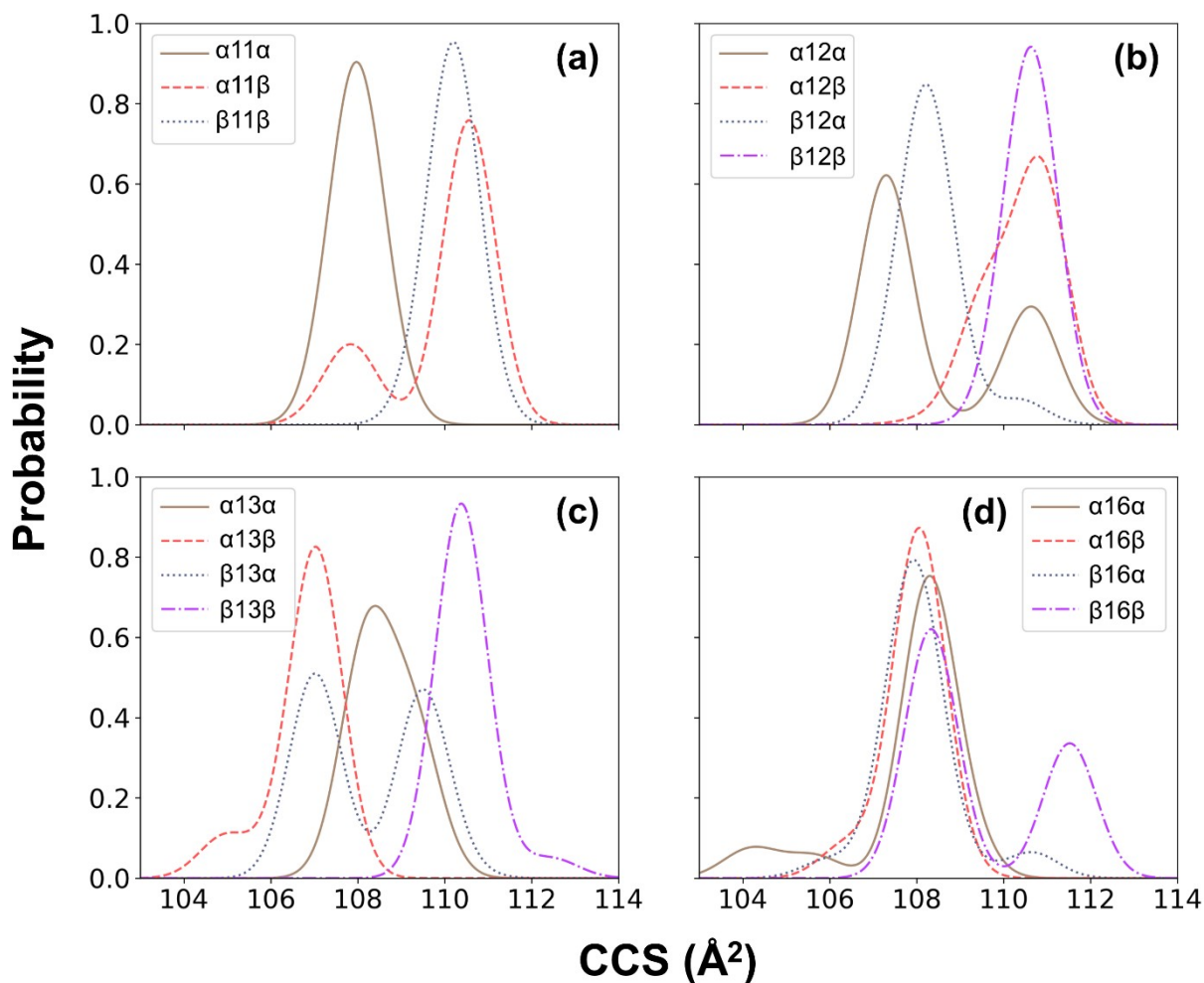


Figure S26. The probability distribution of CCS values (\AA^2) of the conformers of sodiated disaccharide with (a) 1-1 linkage, (b) 1-2 linkage, (c) 1-3 linkage, and (d) 1-6 linkage. The x-axis displays the CCS value (\AA^2), and the y-axis represents the HSA population probability. The width is set at 0.6\AA^2 .