Electronic Supplementary Information: Structure Prediction from Spectra amidst Dynamical Heterogeneity in Melanin

Arpan Choudhury,[†] Raghunathan Ramakrishnan,^{*,‡} and Debashree Ghosh^{*,†}

†School of Chemical Sciences, Indian Association for the Cultivation of Science, Kolkata 700032, India

‡Tata Institute of Fundamental Research Hyderabad, Hyderabad 500046, India

E-mail: ramakrishnan@tifrh.res.in; pcdg@iacs.res.in

Methods

In this Section, we outline the methodologies used for the ML-guided assignment of the valence excitation spectra of DHICA melanin. The steps involved in this process are depicted in the workflow diagram (Fig 2) in the main paper.

Data generation. Molecular dynamics simulations of DHICA, MKICA, and DKICA tetramers (see Supplementary Fig S1-S3) were carried out with a general AMBER force field (GAFF) as implemented in AMBER16.¹ The equilibrations in a canonical (NVT) ensemble were done with 0.5 fs of integration time step at 300K using the Langevin thermostat. The trajectories were run for a total of 1 ns, with the first 100 ps discarded, and configurations were sampled every 1 ps. Electronic excited state calculations for these configurations were carried out at the TDDFT level of theory using the CAM-B3LYP functional and 6-31G(d) basis set, utilizing the QChem software package.² The excitation energies and oscillator strengths for all singlet excitations were calculated within the wavelength range of 200-800 nm. Solvent effects have also been evaluated via continuum solvation models and we notice only small changes in the spectra upon solvation³ (see supplementary Fig S7).

Data clustering. The K-means algorithm was utilized to cluster the structures based on their three inter-monomer dihedral angles (θ), with a requirement of a standard deviation of ~ 5° within each cluster. Further sub-clustering was performed within each cluster based on the values of 8/7/6 ring-OH dihedral angles (ϕ) in DHICA/MKICA/DKICA melanin. These clusters are the atropisomers and their sub-clusters are the dynamic structures of statistically dominant conformational macrostates which are responsible for the population of each cluster. Finally, we take the mean of the structures and the spectra within each sub-cluster to generate unique ML training data for each cluster.

Machine learning architecture. In this study, we have developed a ML model using kernel-ridge regression $(KRR)^{4,5}$ to autonomously assign the DHICA melanin spectra. KRR-based ML approach has been shown to result in asymptotically vanishing errors for predicting molecular ground⁶ and excited state⁷ energies across the chemical space using structural

descriptors as the input. However, in this study, we have applied KRR for the inverse problem of structure assignment using the valence-excitation spectrum as the input, and the dihedral angles as the output. A novel aspect of our study is that our models are applicable for assigning the individual structure in a mixture of melanin.

For training the ML models, we obtained our data by randomly sampling a mixture of clusters (i.e., atropisomers) for each DHICA, MKICA and DKICA. In experiments, the objective is to identify the composition and structure of melanin that correspond to a particular spectral fingerprint within a limited wavelength range. To support these investigations, the full spectrum (200-800 nm) was divided into small 10 nm windows. Then the ML models were trained for each of these windows, with the intensities of all clusters in the mixture being combined into M bins. In other words, we train ML models for the entire UV-vis spectrum. Irrespective of the number of clusters in a mixture, the input to the ML models is a summed-intensity vector, \mathbf{p} , of length M.

The prediction of an angle for an individual configuration in a query mixture is computed by using the kernel-Ansatz

$$\theta^{pred.} = \sum_{i=1}^{N} \alpha_i k(\mathbf{p}_q, \mathbf{p}_i).$$
(1)

The equation states that the predicted angle, $\theta^{pred.}$, for a query mixture (q) is the sum of all α_i multiplied by the kernel function, $k(\mathbf{p}_q, \mathbf{p}_i)$, for each element (i) in the training set of size N. The kernel function $k(\cdot)$ measures the similarity between the query and training spectra, \mathbf{p}_q and \mathbf{p}_i . After evaluating the performance of different kernel functions, the Laplacian kernel was selected that is defined as $k(\mathbf{p}_q, \mathbf{p}_i) = \exp(-|\mathbf{p}_q - \mathbf{p}_i|_1/\sigma)$, where σ is a hyperparameter that determines the width of the kernel and $|\cdot|_1$ denotes the L_1 norm (see Supplementary Fig S8). We trained separate ML models for each dihedral angle by setting σ for all models constant according to the single-kernel formalism suitable for multi-property modeling.⁸ This value of σ obtained was found to be similar to the value obtained from a 5-fold cross-validated model (see Supplementary Fig S9). The regression coefficients in Eq.(1) are determined as the vector, α , for each angle by solving the linear equation

$$[\mathbf{K} + \lambda \mathbf{I}]\boldsymbol{\alpha} = \mathbf{x},\tag{2}$$

where \mathbf{x} is a vector containing the angle in question for all training instances, and \mathbf{K} is a square matrix that comprises all the pairwise kernel functions between the training instances. Due to the absence of outliers in our dataset, we have set the regularization strength (λ) to a constant value of 10⁻⁴ that is sufficient to make the kernel matrix positive definite.

DHICA



Fig S1. Initial structures of DHICA tetramers which are taken for MD simulations.

MKICA



Fig S2. Initial structures of MKICA tetramers which are taken for MD simulations.

DKICA



Fig S3. Initial structures of DKICA tetramers which are taken for MD simulations.





Fig S4. Standard deviation of all the structural degrees of freedoms along the MD trajectory are shown for DHICA, MKICA and DKICA melanin tetramers. The standard deviation of bond distances (\mathbf{a}) and bond angles (\mathbf{b}) are very low. But for the dihedral angles (\mathbf{c}), only a few important parameters (i.e. inter-monomer and ring-OH dihedral angles) show significant standard deviation (marked by the dashed box). These are chosen as ML output vectors. All other degrees of freedom are considered as rigid.



Fig S5. Data distribution and performance of ML models for predicting the structures of different melanin configurations using an ensemble-averaged resultant spectrum as the input. (a) Distribution of cumulative spectral intensities for unique configurations of DHICA, MKICA, and DKICA melanin tetramers across different wavelength ranges. (b) Distribution of the dihedral angles $\{\theta, \phi\}$ for unique configurations of melanin. Learning curves for predicting structures with 5, 10, and 15 most important clusters for each DHICA, MKICA and DKICA using spectral intensities in the (c) 290-300 nm and (d) 560-570 nm wavelength range as the input. In (c) and (d), the mean absolute error (MAE) of the angles in the ML-reconstructed structures is plotted for various training set sizes. Predictions were made on a 10k holdout set. The inset displays error distribution for various ML models.



Fig S6. Error breakdown. Separate mean absolute error (MAE) for inter-monomer dihedral angle (θ) and ring-OH dihedral angle (ϕ) in comparison to the 'mean' of them are given for the intensities in the (a) 290-300 nm and (b) 560-570 nm wavelength range. It shows that the learning curves show much better convergence when predicting only the inter-monomer dihedral angles (i.e., the atropisomers).











Fig S7. Effect of solvent. Comparison of TDDFT stick spectra in solvent phase (PCM with dielectric constant of water) and gas phase for ten random molecules in the dataset.



Fig S8. Kernel function benchmarking. The learning curves for linear, gaussian and laplacian kernel functions are shown in the 290-300 nm wavelength region.



Fig S9. Hyperparameter benchmarking. Kernel width (σ) calculated with single kernel ansatz (vertical dashed lines) are compared with 5-fold cross-validation model (solid lines). Two different kernel functions for two training set sizes are used in the cross-validation. Single kernel σ are calculated with random 500 input representations from the training data.

Section S1. Validation of the ML models on an arbitrary artificial spectrum

To further validate our approach, we evaluated our ML models for an artificial model spectrum with a bimodal distribution (Fig S10). Our goal is to assign absorptions near the peak maxima in the wavelength ranges 290-300 nm and 690-700 nm. We found that for the 690-700 nm range, the only possible species was DKICA tetramers, while for the 290-300 nm range, all three forms were possible. The characteristic structures for the 5 most important clusters in these regions were reconstructed. They are shown in Fig S10 (blue denotes structures that absorb in the range 290-300 nm and green denotes those that absorb in the 690-700 nm range). The percentage contributions of each cluster to the resultant spectrum are also displayed next to the corresponding structure. The predicted inter-monomer dihedral angle values and standard deviation of the clusters are also shown in the figure. The full list of



predicted dihedral angles is shown in the supplementary Table S2 and S3.

Fig S10. Analysis of the predicted melanin structures for an artificial bimodal spectrum with peak intensities at 300 nm and 700 nm. Models were trained on mixtures of 5 clusters for each DHICA, MKICA and DKICA using inter-monomer $\{\theta\}$ and ring-OH $\{\phi\}$ dihedral angles as the targets. ML predictions were made for the spectral intensities in 290-300 nm and 690-700 nm. Reconstructed structures with their percentage contribution to the input spectrum are shown for each cluster. The individual ML-predicted angles $\{\theta\}$ and their uncertainty determined by the K-means clustering of the training set are also provided.

In the spectral range 690-700nm, the predicted structures display inter-monomer dihedral angles close to 0 degrees, resulting in unphysical planar structures. This is attributed to the low-intensity region of the tetramer spectra in the database (see supplementary Fig S5a). The low intensity in this range is indicative of the absence of low energy (stable or most probable) conformers that can absorb significantly in this wavelength region.

Overall, we have tested our model on a holdout validation set as well as on an artificial spectrum. The prediction on the artificial spectrum has shown expected outcomes of stable conformers in the high confidence region of the spectral window where the melanin absorbs significantly and null values for the structural parameters in the low confidence region where melanin does not absorb significantly. We have further confirmed the validity of our predictions in the 290-300 nm range by again calculating the TDDFT spectra of the ML predicted structures and their intensities are within the correct wavelength range (supplementary Table S4).

Table S1. Training set bias check. Inter-monomer dihedral angle (θ) predictions on a query spectrum for 5 different randomly chosen training sets of size 10k in the 290-300 nm wavelength range. All 5 training sets predict the same cluster for all the configurations in DHICA, MKICA and DKICA which indicates that there is no training set bias.

Training set 1	Training set 2	Training set 3	Training set 4	Training set 5
DHICA	DHICA	DHICA	DHICA	DHICA
$\theta_1 = 54.82$	$\theta_1 = 54.60$	$\theta_1 = 54.41$	$\theta_1 = 54.57$	$\theta_1 = 54.49$
$\theta_2 = 78.13$	$\theta_2 = 77.92$	$\theta_2 = 77.80$	$\theta_2 = 77.91$	$\theta_2 = 77.79$
$\theta_3 = 98.51$	$\theta_3 = 98.29$	$\theta_3 = 98.16$	$\theta_3 = 98.26$	$\theta_3 = 98.15$
(Cluster $id = 59$)				
$\theta_1 = 84.42$	$\theta_1 = 84.24$	$\theta_1 = 83.86$	$\theta_1 = 84.34$	$\theta_1 = 83.90$
$\theta_2 = 88.97$	$\theta_2 = 88.79$	$\theta_2 = 88.56$	$\theta_2 = 88.76$	$\theta_2 = 88.47$
$\theta_3 = 47.78$	$\theta_3 = 47.66$	$\theta_3 = 47.64$	$\theta_3 = 47.51$	$\theta_3 = 47.65$
(Cluster $id = 50$)				
$\theta_1 = 103.93$	$\theta_1 = 103.66$	$\theta_1 = 103.55$	$\theta_1 = 103.60$	$\theta_1 = 103.54$
$\theta_2 = 91.81$	$\theta_2 = 91.62$	$\theta_2 = 91.49$	$\theta_2 = 91.62$	$\theta_2 = 91.46$
$\theta_3 = 93.98$	$\theta_3 = 93.80$	$\theta_3 = 93.65$	$\theta_3 = 93.78$	$\theta_3 = 93.67$
(Cluster id $= 39$)				
$\theta_1 = 99.29$	$\theta_1 = 99.80$	$\theta_1 = 99.01$	$\theta_1 = 99.93$	$\theta_1 = 98.55$
$\theta_2 = 78.82$	$\theta_2 = 80.13$	$\theta_2 = 78.87$	$\theta_2 = 80.32$	$\theta_2 = 77.96$
$\theta_3 = 68.07$	$\theta_3 = 68.47$	$\theta_3 = 67.87$	$\theta_3 = 68.43$	$\theta_3 = 67.72$
(Cluster id = 44)	(Cluster id $= 44$)			
$\theta_1 = 58.34$	$\theta_1 = 58.29$	$\theta_1 = 58.20$	$\theta_1 = 58.22$	$\theta_1 = 58.19$
$\theta_2 = 87.13$	$\theta_2 = 86.88$	$\theta_2 = 86.72$	$\theta_2 = 86.87$	$\theta_2 = 86.66$
$\theta_3 = 50.24$	$\theta_3 = 50.00$	$\theta_3 = 50.00$	$\theta_3 = 50.05$	$\theta_3 = 49.88$
(Cluster id $= 25$)				
MKICA	MKICA	MKICA	MKICA	MKICA
$\theta_1 = 88.41$	$\theta_1 = 88.20$	$\theta_1 = 88.08$	$\theta_1 = 88.10$	$\theta_1 = 87.98$
$\theta_2 = 64.73$	$\theta_2 = 64.61$	$\theta_2 = 64.47$	$\theta_2 = 64.53$	$\theta_2 = 64.47$
$\theta_3 = 86.86$	$\theta_3 = 86.74$	$\theta_3 = 86.51$	$\theta_3 = 86.61$	$\theta_3 = 86.51$
(Cluster id $= 53$)				
$\theta_1 = 81.84$	$\theta_1 = 81.60$	$\theta_1 = 81.53$	$\theta_1 = 81.54$	$\theta_1 = 81.46$
$\theta_2 = 95.95$	$\theta_2 = 95.68$	$\theta_2 = 95.58$	$\theta_2 = 95.56$	$\theta_2 = 95.49$
$\theta_3 = 90.57$	$\theta_3 = 90.35$	$\theta_3 = 90.24$	$\theta_3 = 90.40$	$\theta_3 = 90.22$
(Cluster id = 54)	(Cluster id = 54)	(Cluster id $= 54$)	(Cluster id = 54)	(Cluster id $= 54$)

$\theta_1 = 67.16$	$\theta_1 = 67.08$	$\theta_1 = 66.96$	$\theta_1 = 67.05$	$\theta_1 = 66.91$
$\theta_2 = 75.89$	$\theta_2 = 75.81$	$\theta_2 = 75.69$	$\theta_2 = 75.78$	$\theta_2 = 75.59$
$\theta_3 = 60.08$	$\theta_3 =$	$59.96\theta_3 = 59.90$	$\theta_3 = 59.88$	$\theta_{3} = 59.85$
(Cluster $id = 24$)	(Cluster $id = 24$)	(Cluster $id = 24$)	(Cluster $id = 24$)	(Cluster $id = 24$)
$\theta_1 = 42.89$	$\theta_1 = 42.51$	$\theta_1 = 42.79$	$\theta_1 = 42.53$	$\theta_1 = 42.24$
$\theta_2 = 68.29$	$\theta_2 = 68.11$	$\theta_2 = 68.01$	$\theta_2 = 68.11$	$\theta_2 = 67.97$
$\theta_3 = 102.84$	$\theta_3 = 102.53$	$\theta_3 = 102.47$	$\theta_3 = 102.47$	$\theta_3 = 102.30$
(Cluster $id = 21$)	(Cluster $id = 21$)	(Cluster $id = 21$)	(Cluster $id = 21$)	(Cluster $id = 21$)
$\theta_1 = 60.29$	$\theta_1 = 60.21$	$\theta_1 = 60.08$	$\theta_1 = 60.13$	$\theta_1 = 60.10$
$\theta_2 = 113.46$	$\theta_2 = 113.16$	$\theta_2 = 113.02$	$\theta_2 = 113.07$	$\theta_2 = 112.87$
$\theta_3 = 105.70$	$\theta_3 = 105.43$	$\theta_3 = 105.21$	$\theta_3 = 105.49$	$\theta_3 = 105.26$
(Cluster id $= 33$)	(Cluster id $= 33$)	(Cluster id $= 33$)	(Cluster id $= 33$)	(Cluster id $= 33$)
DKICA	DKICA	DKICA	DKICA	DKICA
$\theta_1 = 70.96$	$\theta_1 = 70.79$	$\theta_1 = 70.70$	$\theta_1 = 70.77$	$\theta_1 = 70.66$
$\theta_2 = 68.69$	$\theta_2 = 68.53$	$\theta_2 = 68.45$	$\theta_2 = 68.54$	$\theta_2 = 68.42$
$\theta_2 = 83.30$	$\theta_2 = 83.09$	$\theta_2 = 83.00$	$\theta_2 = 83.06$	$\theta_2 = 82.95$
(Cluster $id = 20$)	(Cluster $id = 20$)	(Cluster $id = 20$)	(Cluster $id = 20$)	(Cluster $id = 20$)
	(010000110 20)	(010000110 20)	(010000110 20)	(010000110 20)
$\theta_1 = 63.52$	$\theta_1 = 63.37$	$\theta_1 = 63.27$	$\theta_1 = 63.35$	$\theta_1 = 63.26$
$\theta_1 = 61.00$	$\theta_1 = 60.85$	$\theta_1 = 60.27$	$\theta_1 = 60.84$	$\theta_1 = 60.20$
$\theta_2 = 71.00$	$\theta_2 = 00.00$ $\theta_3 = 71.08$	$\theta_2 = 30.09$	$\theta_2 = 00.01$ $\theta_3 = 71.05$	$\theta_2 = 30.19$ $\theta_3 = 70.94$
(Cluster id - 46)	(Cluster id - 46)	(Cluster id - 46)	(Cluster id - 46)	(Cluster id - 46)
$\left(\text{Cluster id} = 40 \right)$	(Cluster Id = 40)	(Cluster Id = 40)	$\left(\text{Cluster id} = 40 \right)$	(Cluster Id = 40)
$\theta_{1} = 96.72$	$\theta_{1} = 96.46$	$\theta_{1} = 96.36$	$\theta_{1} = 96.40$	$\theta_{1} = 96.34$
$\theta_1 = 50.12$ $\theta_2 = 75.72$	$\theta_1 = 50.40$ $\theta_2 = 75.53$	$\theta_1 = 50.50$ $\theta_2 = 75.43$	$\theta_1 = 50.40$ $\theta_2 = 75.50$	$\theta_1 = 50.34$ $\theta_2 = 75.38$
$\theta_2 = 10.12$ $\theta_2 = 102.84$	$\theta_2 = 10.00$ $\theta_2 = 102.62$	$\theta_2 = 10.45$ $\theta_2 = 102.45$	$\theta_2 = 102.63$	$\theta_2 = 102.45$
$0_3 = 102.04$ (Cluster id = 27)	$0_3 = 102.02$	(Cluster id - 27)	$0_3 = 102.05$	(Cluster id - 27)
(Cluster Id = 27)	(Cluster Id - 21)	(Cluster Id - 21)	(Cluster Id = 27)	(Cluster Id - 21)
$\theta_{1} = 65.02$	A 65.66	$\theta_{1} = 65.80$	$A_{1} = 65.64$	$\theta_{1} = 65.87$
$0_1 = 05.92$ $\theta_1 = 07.51$	$0_1 = 05.00$ $\theta_1 = 07.31$	$v_1 = 05.80$ $\theta_1 = 07.16$	$0_1 = 05.04$ $\theta_1 = 07.22$	$b_1 = 05.87$ $\theta_1 = 07.10$
$b_2 = 97.51$ $a_1 = 00.18$	$b_2 = 97.31$ $a_1 = 80.00$	$b_2 = 97.10$ $a_1 = 80.78$	$b_2 = 91.22$ $A_1 = 80.02$	$b_2 = 97.10$ $\theta_1 = 80.60$
$0_3 = 90.10$	$0_3 = 09.99$	$0_3 = 09.10$	$0_3 = 09.93$	$0_3 = 09.09$
(Oruster Id = 23)	(Oruster Id = 23)	(Oruster Id = 23)	(Oruster Id = 23)	(Oruster Id = 23)
A = 111.07	A = 111.91	A = 110.97	A = 111.17	A = 110.96
$\sigma_1 = 111.07$	$\sigma_1 = 111.31$	$\sigma_1 = 110.87$	$\sigma_1 = 111.17$	$\sigma_1 = 110.80$
$\sigma_2 = 81.12$	$\sigma_2 = 78.19$	$\sigma_2 = 80.05$	$\sigma_2 = 79.01$	$\sigma_2 = 79.81$
$\sigma_3 = (2.59)$	$\sigma_3 = (4.80)$	$\sigma_3 = (3.05)$	$\sigma_3 = (4.08)$	$\sigma_3 = (3.18)$
(Cluster 1d = 2)	(Cluster $id = 2$)	(Cluster $1d = 2$)	(Cluster 1d = 2)	(Cluster $1d = 2$)

Predicted dihedral angles (in deg.) of DHICA Configuration 1 Configuration 2 Configuration 3 Configuration 4 Configuration 5 $\theta_1 = 54.82$ $\theta_1 = 99.29$ $\theta_1 = 84.42$ $\theta_1 = 103.93$ $\theta_1 = 58.34$ $\theta_2 = 78.13$ $\theta_2 = 88.97$ $\theta_2 = 91.81$ $\theta_2 = 78.82$ $\theta_2 = 87.13$ $\theta_3 = 98.51$ $\theta_3 = 47.78$ $\theta_3 = 93.98$ $\theta_3 = 68.07$ $\theta_3 = 50.24$ $\phi_1 = -4.82$ $\phi_1 = 25.38$ $\phi_1 = -7.72$ $\phi_1 = -313.63$ $\phi_1 = 8.00$ $\phi_2 = 160.28$ $\phi_2 = 176.95$ $\phi_2 = 164.21$ $\phi_2 = 495.45$ $\phi_2 = 168.38$ $\phi_3 = 71.89$ $\phi_3 = 183.97$ $\phi_3 = 93.17$ $\phi_3 = 72.95$ $\phi_3 = 28.16$ $\phi_4 = 148.06$ $\phi_4 = 13.70$ $\phi_4 = 86.94$ $\phi_4 = 87.61$ $\phi_4 = 139.05$ $\phi_5 = 72.42$ $\phi_5 = 70.45$ $\phi_5 = -85.11$ $\phi_5 = 121.93$ $\phi_5 = -173.72$ $\phi_6 = 152.26$ $\phi_6 = 216.35$ $\phi_6 = 57.90$ $\phi_6 = 278.49$ $\phi_6 = 106.17$ $\phi_7 = 159.84$ $\phi_7 = 102.77$ $\phi_7 = 139.41$ $\phi_7 = 76.45$ $\phi_7 = 150.34$ $\phi_8 = 6.65$ $\phi_8 = 60.23$ $\phi_8 = 20.92$ $\phi_8 = -12.53$ $\phi_8 = 77.98$ Predicted dihedral angles (in deg.) of MKICA Configuration 1 Configuration 3 Configuration 5 Configuration 2 Configuration 4 $\theta_1 = 88.41$ $\theta_1 = 81.84$ $\theta_1 = 67.16$ $\theta_1 = 42.89$ $\theta_1 = 60.29$ $\theta_2 = 64.73$ $\theta_2 = 95.95$ $\theta_2 = 75.89$ $\theta_2 = 68.29$ $\theta_2 = 113.46$ $\theta_3 = 86.86$ $\theta_3 = 90.57$ $\theta_3 = 60.08$ $\theta_3 = 102.84$ $\theta_3 = 105.70$ $\phi_1 = 155.50$ $\phi_1 = 169.33$ $\phi_1 = 134.17$ $\phi_1 = 201.59$ $\phi_1 = 160.55$ $\phi_2 = 121.67$ $\phi_2 = 65.23$ $\phi_2 = -166.22$ $\phi_2 = 500.60$ $\phi_2 = 6.09$ $\phi_3 = 84.68$ $\phi_3 = 97.48$ $\phi_3 = 177.83$ $\phi_3 = -321.00$ $\phi_3 = 164.56$ $\phi_4 = 170.79$ $\phi_4 = 99.12$ $\phi_4 = 81.89$ $\phi_4 = 170.64$ $\phi_4 = 113.69$ $\phi_5 = 90.49$ $\phi_5 = 75.48$ $\phi_5 = 265.65$ $\phi_5 = 521.45$ $\phi_5 = 153.18$ $\phi_6 = 160.81$ $\phi_6 = 162.98$ $\phi_6 = -151.34$ $\phi_6 = 179.04$ $\phi_6 = 143.46$ $\phi_7 = -4.80$ $\phi_7 = -42.59$ $\phi_7 = 3.02$ $\phi_7 = 321.42$ $\phi_7 = 45.24$ Predicted dihedral angles (in deg.) of DKICA Configuration 3 Configuration 1 Configuration 2 Configuration 4 Configuration 5 $\theta_1 = 70.96$ $\theta_1 = 63.52$ $\theta_1 = 96.72$ $\theta_1 = 65.92$ $\theta_1 = 111.07$ $\theta_2 = 68.69$ $\theta_2 = 61.00$ $\theta_2 = 75.72$ $\theta_2 = 97.51$ $\theta_2 = 81.12$ $\theta_3 = 83.30$ $\theta_3 = 71.25$ $\theta_3 = 102.84$ $\theta_3 = 90.18$ $\theta_3 = 72.59$ $\phi_1 = 164.09$ $\phi_1 = 165.58$ $\phi_1 = 127.68$ $\phi_1 = 9.51$ $\phi_1 = -5.71$ $\phi_2 = 2.86$ $\phi_2 = 153.71$ $\phi_2 = 155.68$ $\phi_2 = 198.27$ $\phi_2 = 132.62$ $\phi_3 = -26.20$ $\phi_3 = -4.59$ $\phi_3 = 33.28$ $\phi_3 = 259.94$ $\phi_3 = -22.82$ $\phi_4 = 165.45$ $\phi_4 = 165.62$ $\phi_4 = 160.24$ $\phi_4 = 197.30$ $\phi_4 = 142.56$ $\phi_5 = 222.96$ $\phi_5 = 147.02$ $\phi_5 = 149.21$ $\phi_5 = 225.04$ $\phi_5 = -12.37$ $\phi_6 = -29.54$ $\phi_6 = 9.75$ $\phi_6 = 4.15$ $\phi_6 = -40.28$ $\phi_6 = 194.45$

Table S2. All predicted inter-monomer $\{\theta\}$ and ring-OH $\{\phi\}$ dihedral angles for the artificial spectrum in the 290-300 nm wavelength range.

Table S3. All predicted inter-monomer $\{\theta\}$ and ring-OH $\{\phi\}$ dihedral angles for the artificial spectrum in the 690-700 nm wavelength range.

Predicted dihedral angles (in deg.) of DKICA				
Configuration 1	Configuration 2	Configuration 3	Configuration 4	Configuration 5
$\theta_1 = 0.027$	$\theta_1 = 0.033$	$\theta_1 = 0.041$	$\theta_1 = 0.039$	$\theta_1 = 0.039$
$\theta_2 = 0.028$	$\theta_2 = 0.029$	$\theta_2 = 0.042$	$\theta_2 = 0.041$	$\theta_2 = 0.026$
$\theta_3 = 0.023$	$\theta_3 = 0.031$	$\theta_3 = 0.040$	$\theta_3 = 0.029$	$\theta_3 = 0.040$
$\phi_1 = 0.022$	$\phi_1 = 0.002$	$\phi_1 = 0.028$	$\phi_1 = -0.002$	$\phi_1 = 0.002$
$\phi_2 = 0.060$	$\phi_2 = 0.063$	$\phi_2 = 0.066$	$\phi_2 = 0.064$	$\phi_2 = 0.063$
$\phi_3 = 0.033$	$\phi_3 = 0.028$	$\phi_3 = 0.039$	$\phi_3 = 0.002$	$\phi_3 = 0.018$
$\phi_4 = 0.065$	$\phi_4 = 0.065$	$\phi_4 = 0.064$	$\phi_4 = 0.055$	$\phi_4 = 0.068$
$\phi_5 = 0.044$	$\phi_5 = 0.031$	$\phi_5 = 0.024$	$\phi_5 = 0.039$	$\phi_5 = 0.036$
$\phi_6 = 0.031$	$\phi_6 = 0.027$	$\phi_6 = 0.034$	$\phi_6 = 0.062$	$\phi_6 = 0.033$

Table S4. TDDFT calculations of the ML predicted structures in the 290-300 nm wavelength range of the artificial spectrum.

DHICA		MKICA		DKICA	
Wavelength	Intensity	Wavelength	Intensity	Wavelength	Intensity
(nm)	(a.u.)	(nm)	(a.u.)	(nm)	(a.u.)
296	0.100	290	0.116	291	0.674
285	0.081	285	0.599		
295	0.832	290	0.081	301.27	0.591
				301.16	0.088
				295	0.163
295	0.236	302	0.489	288	0.567
288	0.057	291	0.070		
		288	0.170		
296	0.620	304	0.487	286.61	0.388
288	0.027	293	0.137	286.28	0.053
		287	0.188		
303	0.649	304	0.553	291	0.664
293	0.111	292	0.190		
		286	0.119		

Predicted dihedral angles (in deg.) of DHICA				
Configuration 1	Configuration 2	Configuration 3	Configuration 4	Configuration 5
$\theta_1 = 84.49$	$\theta_1 = 56.67$	$\theta_1 = 102.37$	$\theta_1 = 110.14$	$\theta_1 = 84.17$
$\theta_2 = 91.44$	$\theta_2 = 69.75$	$\theta_2 = 91.00$	$\theta_2 = 89.41$	$\theta_2 = 60.40$
$\theta_3 = 109.04$	$\theta_3 = 64.55$	$\theta_3 = 74.40$	$\theta_3 = 111.56$	$\theta_3 = 69.77$
$\phi_1 = -2.40$	$\phi_1 = -6.47$	$\phi_1 = 88.55$	$\phi_1 = 6.72$	$\phi_1 = -5.35$
$\phi_2 = 184.38$	$\phi_2 = 170.31$	$\phi_2 = 91.04$	$\phi_2 = 167.35$	$\phi_2 = 180.38$
$\phi_3 = 123.60$	$\phi_3 = 125.28$	$\phi_3 = 26.81$	$\phi_3 = 158.86$	$\phi_3 = 26.74$
$\phi_4 = 50.46$	$\phi_4 = 86.77$	$\phi_4 = 192.58$	$\phi_4 = 35.76$	$\phi_4 = 45.07$
$\phi_5 = 7.46$	$\phi_5 = 83.32$	$\phi_5 = 87.25$	$\phi_5 = 144.33$	$\phi_5 = -0.37$
$\phi_6 = 153.55$	$\phi_6 = 91.66$	$\phi_6 = -39.57$	$\phi_6 = 33.92$	$\phi_6 = 219.89$
$\phi_7 = 159.93$	$\phi_7 = 167.42$	$\phi_7 = 265.89$	$\phi_7 = 157.19$	$\phi_7 = -33.53$
$\phi_8 = -8.41$	$\phi_8 = 6.46$	$\phi_8 = -70.13$	$\phi_8 = 4.25$	$\phi_8 = 141.50$
	Predicted dihe	dral angles (in de	eg.) of MKICA	
Configuration 1	Configuration 2	Configuration 3	Configuration 4	Configuration 5
$\theta_1 = 73.71$	$\theta_1 = 100.15$	$\theta_1 = 85.95$	$\theta_1 = 95.32$	$\theta_1 = 77.62$
$\theta_2 = 87.03$	$\theta_2 = 85.36$	$\theta_2 = 113.50$	$\theta_2 = 62.56$	$\theta_2 = 99.24$
$\theta_3 = 89.01$	$\theta_3 = 92.38$	$\theta_3 = 77.84$	$\theta_3 = 107.96$	$\theta_3 = 87.00$
$\phi_1 = 135.83$	$\phi_1 = 172.82$	$\phi_1 = 100.76$	$\phi_1 = -45.46$	$\phi_1 = 62.24$
$\phi_2 = 111.78$	$\phi_2 = 182.22$	$\phi_2 = 2.41$	$\phi_2 = 90.21$	$\phi_2 = 35.96$
$\phi_3 = 203.20$	$\phi_3 = 91.24$	$\phi_3 = 173.97$	$\phi_3 = 192.35$	$\phi_3 = 137.99$
$\phi_4 = 108.26$	$\phi_4 = 30.34$	$\phi_4 = -82.39$	$\phi_4 = 323.72$	$\phi_4 = -2.98$
$\phi_5 = 112.53$	$\phi_5 = 182.55$	$\phi_5 = 100.54$	$\phi_5 = -183.77$	$\phi_5 = 94.46$
$\phi_6 = 155.67$	$\phi_6 = 169.66$	$\phi_6 = 128.65$	$\phi_6 = 63.95$	$\phi_6 = 175.25$
$\phi_7 = 24.04$	$\phi_7 = 2.88$	$\phi_7 = 131.21$	$\phi_7 = 89.74$	$\phi_7 = 7.69$
Predicted dihedral angles (in deg.) of DKICA				
Configuration 1	Configuration 2	Configuration 3	Configuration 4	Configuration 5
$\theta_1 = 65.95$	$\theta_1 = 69.50$	$\theta_1 = 84.13$	$\theta_1 = 82.41$	$\theta_1 = 105.88$
$\theta_2 = 80.19$	$\theta_2 = 90.62$	$\theta_2 = 94.37$	$\theta_2 = 69.18$	$\theta_2 = 83.83$
$\theta_3 = 100.09$	$\theta_3 = 70.66$	$\theta_3 = 100.92$	$\theta_3 = 96.90$	$\theta_3 = 66.08$
$\phi_1 = -4.83$	$\phi_1 = 0.47$	$\phi_1 = 5.69$	$\phi_1 = -5.05$	$\phi_1 = -41.32$
$\phi_2 = 201.58$	$\phi_2 = 176.57$	$\phi_2 = 177.28$	$\phi_2 = 157.23$	$\phi_2 = 186.02$
$\phi_3 = 255.72$	$\phi_3 = 65.91$	$\phi_3 = 269.98$	$\phi_3 = -3.51$	$\phi_3 = 193.82$
$\phi_4 = 153.85$	$\phi_4 = 166.14$	$\phi_4 = 173.79$	$\phi_4 = 171.53$	$\phi_4 = 163.81$
$\phi_5 = 289.77$	$\phi_5 = 123.24$	$\phi_5 = 145.94$	$\phi_5 = 165.04$	$\phi_5 = -83.45$
$\phi_6 = -96.50$	$\phi_6 = 100.83$	$\phi_6 = 0.98$	$\phi_6 = 3.01$	$\phi_6 = 259.15$

Table S5. All predicted inter-monomer $\{\theta\}$ and ring-OH $\{\phi\}$ dihedral angles for the experimental spectrum in the 280-290 nm wavelength range.

Predicted dihedral angles (in deg.) of MKICA				
Configuration 1	Configuration 2	Configuration 3	Configuration 4	Configuration 5
$\theta_1 = 103.45$	$\theta_1 = 49.17$	$\theta_1 = 106.07$	$\theta_1 = 106.00$	$\theta_1 = 53.96$
$\theta_2 = 92.05$	$\theta_2 = 107.63$	$\theta_2 = 48.50$	$\theta_2 = 73.02$	$\theta_2 = 106.34$
$\theta_3 = 56.97$	$\theta_3 = 50.99$	$\theta_3 = 90.01$	$\theta_3 = 93.77$	$\theta_3 = 100.31$
$\phi_1 = 142.78$	$\phi_1 = 133.99$	$\phi_1 = 158.81$	$\phi_1 = 159.10$	$\phi_1 = 106.39$
$\phi_2 = 89.25$	$\phi_2 = 30.90$	$\phi_2 = 75.95$	$\phi_2 = 133.91$	$\phi_2 = -3.09$
$\phi_3 = 102.63$	$\phi_3 = 116.26$	$\phi_3 = 145.01$	$\phi_3 = 94.75$	$\phi_3 = 152.62$
$\phi_4 = 116.35$	$\phi_4 = 81.97$	$\phi_4 = 42.07$	$\phi_4 = 131.68$	$\phi_4 = -4.22$
$\phi_5 = 98.99$	$\phi_5 = 115.89$	$\phi_5 = 130.41$	$\phi_5 = 136.12$	$\phi_5 = 150.84$
$\phi_6 = 161.02$	$\phi_6 = 119.63$	$\phi_6 = 102.10$	$\phi_6 = 153.61$	$\phi_6 = 118.58$
$\phi_7 = -2.30$	$\phi_7 = 29.87$	$\phi_7 = 59.28$	$\phi_7 = 9.64$	$\phi_7 = 92.12$
	Predicted dihe	dral angles (in d	eg.) of DKICA	
Configuration 1	Configuration 2	Configuration 3	Configuration 4	Configuration 5
$\theta_1 = 65.45$	$\theta_1 = 61.11$	$\theta_1 = 81.83$	$\theta_1 = 82.27$	$\theta_1 = 65.55$
$\theta_2 = 79.03$	$\theta_2 = 58.86$	$\theta_2 = 90.90$	$\theta_2 = 109.16$	$\theta_2 = 123.50$
$\theta_3 = 86.23$	$\theta_3 = 69.04$	$\theta_3 = 69.34$	$\theta_3 = 96.98$	$\theta_3 = 73.23$
$\phi_1 = 15.49$	$\phi_1 = 133.13$	$\phi_1 = 2.12$	$\phi_1 = 3.08$	$\phi_1 = 8.18$
$\phi_2 = 158.89$	$\phi_2 = 148.44$	$\phi_2 = 161.31$	$\phi_2 = 159.12$	$\phi_2 = 131.77$
$\phi_3 = 70.57$	$\phi_3 = -0.71$	$\phi_3 = -0.55$	$\phi_3 = 14.26$	$\phi_3 = 2.70$
$\phi_4 = 147.13$	$\phi_4 = 161.21$	$\phi_4 = 166.67$	$\phi_4 = 157.94$	$\phi_4 = 179.86$
$\phi_5 = -3.63$	$\phi_5 = 145.89$	$\phi_5 = 92.01$	$\phi_5 = 148.88$	$\phi_5 = -3.43$
$\phi_6 = 186.07$	$\phi_6 = 8.51$	$\phi_6 = 161.73$	$\phi_6 = 5.46$	$\phi_6 = 168.04$

Table S6. All predicted inter-monomer $\{\theta\}$ and ring-OH $\{\phi\}$ dihedral angles for the experimental spectrum in the 560-570 nm wavelength range.

Table S7. MAE (in degree) with variable spectral wavelength range for 10k training set. Below results show that the 10 nm wavelength range offers the best error. MAEs are calculated on a 10k hold-out set.

Wavelength range (nm)	MAE in inter-monomer	MAE in ring-OH dihedral	
	dihedral angle (θ)	angle (ϕ)	
299-301	1.57	21.23	
298-302	0.30	5.61	
295-305	0.09	1.85	
290-310	0.11	1.95	
275-325	0.23	3.76	

References

- Case, D.; Berryman, J.; Betz, R.; Cerutti, D.; Cheatham III, T.; Darden, T.; Duke, R.; Giese, T.; Gohlke, H.; Goetz, A.; others AMBER 2016; University of California: San Francisco; 2016.
- (2) Shao, Y.; Gan, Z.; Epifanovsky, E.; Gilbert, A. T.; Wormit, M.; Kussmann, J.; Lange, A. W.; Behn, A.; Deng, J.; Feng, X.; others Advances in molecular quantum chemistry contained in the Q-Chem 4 program package. *Molecular Physics* 2015, 113, 184–215.
- (3) Gauden, M.; Pezzella, A.; Panzella, L.; Neves-Petersen, M.; Skovsen, E.; Petersen, S. B.; Mullen, K.; Napolitano, A.; d'Ischia, M.; Sundstrom, V. Role of solvent, pH, and molecular size in excited-state deactivation of key eumelanin building blocks: implications for melanin pigment photostability. *Journal of the American Chemical Society* 2008, 130, 17038–17043.
- (4) Murphy, K. P. Machine learning: a probabilistic perspective; MIT press, 2012.
- (5) Schölkopf, B.; Tsuda, K.; Vert, J.-P.; others Kernel methods in computational biology; MIT press, 2004.
- (6) Rupp, M.; Tkatchenko, A.; Müller, K.-R.; von Lilienfeld, O. A. Fast and accurate modeling of molecular atomization energies with machine learning. *Physical review letters* 2012, 108, 058301.
- (7) Ramakrishnan, R.; Hartmann, M.; Tapavicza, E.; von Lilienfeld, O. A. Electronic spectra from TDDFT and machine learning in chemical space. *The Journal of chemical physics* 2015, 143, 084111.
- (8) Ramakrishnan, R.; von Lilienfeld, O. A. Many molecular properties from one kernel in chemical space. CHIMIA International Journal for Chemistry 2015, 69, 182–186.