

Excited state, non-adiabatic dynamics of large photoswitchable molecules using a chemically transferable machine learning potential

Simon Axelrod,^{1,2} Eugene Shakhnovich,¹ and Rafael Gómez-Bombarelli^{2,*}

¹*Department of Chemistry and Chemical Biology, Harvard University, Cambridge, MA, 02138*

²*Department of Materials Science and Engineering,
Massachusetts Institute of Technology, Cambridge, MA, 02139*

(Dated: March 18, 2022)

Light-induced chemical processes are ubiquitous in nature and have widespread technological applications. For example, photoisomerization can allow a drug with a photo-switchable scaffold such as azobenzene to be activated with light. In principle, photoswitches with desired photophysical properties like high isomerization quantum yields can be identified through virtual screening with reactive simulations. In practice, these simulations are rarely used for screening, since they require hundreds of trajectories and expensive quantum chemical methods to account for non-adiabatic excited state effects. Here we introduce a *diabatic artificial neural network* (DANN) based on diabatic states to accelerate such simulations for azobenzene derivatives. The network is six orders of magnitude faster than the quantum chemistry method used for training. DANN is transferable to azobenzene molecules outside the training set, predicting quantum yields for unseen species that are correlated with experiment. We use the model to virtually screen 3,100 hypothetical molecules, and identify novel species with extremely high predicted quantum yields. The model predictions are confirmed using high-accuracy non-adiabatic dynamics. Our results pave the way for fast and accurate virtual screening of photoactive compounds.

Light is a powerful tool for manipulating molecular systems. It can be controlled with high spatial, spectral and temporal precision to facilitate a variety of processes, including energy transfer, intermolecular reactions, and photoisomerization [1]. These processes are used in areas as diverse as synthesis, energy storage, display technology, biological imaging, diagnostics and medicine [1–3]. Photoactive drugs, for instance, are photoswitchable compounds whose bioactivity can be toggled through light-induced isomerization. Precise spatiotemporal control of bioactivity allows photoactive drugs to be delivered in high doses with minimal off-target activity and side effects. Such therapeutics are a promising path for the treatment of cancer, neurodegenerative diseases, bacterial infections, diabetes, and blindness [4, 5].

Theory plays a key role in explaining and predicting photochemistry because empirical heuristics learned from thermally activated ground state processes typically do not apply to excited states [3]. Computer simulations based on quantum mechanics can achieve impressive accuracy in the prediction of experimental observables. These include the isomerization efficiency and absorption spectrum of photoswitchable compounds [6, 7], which are key quantities in the design of photoactive drugs.

However, *ab initio* methods in photochemistry are severely limited by their computational cost [8]. In order to gather meaningful statistics for one molecule, hundreds of replicate simulations are needed, each of which involves thousands of electronic structure calculations performed in series with sub-femtosecond timesteps. The individual quantum chemical calculations are particularly demanding, requiring excited state gradients and some treatment of multireference effects. In some cases, both the ground- and excited-state gradients are required at each time step [9–11]. Using *ab initio* methods to compute photochemical properties of tens or hundreds molecules is impractical, and photodynamic simulations have not yet been used for large-scale virtual screening.

Among the most accurate and expensive electronic structure methods are multi-configuration perturbation techniques [12–16], but their cost and requirement for manual active space selection limit their use in virtual screening. The photochemistry community has made exciting developments over several years to overcome both of these hurdles. For example, reduced scaling techniques [17, 18] and graphics processing units [19] can significantly accelerate multi-reference calculations. The density matrix renormalization group (DMRG) [20, 21] and multi-reference density functional theory (DFT) methods [22–25] have expanded the size of systems that can be treated with high accuracy. DMRG has also been used to automate the

* Corresponding author: rafagb@mit.edu

selection of active spaces for multi-reference methods [26, 27]. Less accurate but more affordable black-box methods include spin-flip time-dependent DFT (SF-TDDFT) and hole-hole Tamm-Dancoff DFT [28], among others [29, 30]. Despite these developments, the cost of non-adiabatic simulations remains high. As discussed below, even SF-TDDFT is prohibitively expensive for virtual screening. Semi-empirical methods [31–33] are currently the only affordable approach for large-scale screening. They provide qualitatively correct results across many systems, but are ultimately bounded by their approximations, with average energy errors of 15 kcal/mol [32].

A different approach is to use data-driven models in place of quantum chemistry (QC) calculations. Machine learning (ML) models trained on quantum chemical data can now routinely predict ground state energies and forces with sub-chemical accuracy [34–36], and take only milliseconds to make predictions. These models have been successfully used in a variety of ground state simulations [35, 37, 38]. They have also been used to accelerate non-adiabatic simulations in a number of model systems [39–45]. However, excited state ML has not yet offered affordable photodynamics for hundreds of molecules of realistic size, which is the ultimate goal for predictive simulation in photopharmacology. Further, no excited-state interatomic potentials have been developed that are transferable to different compounds. They therefore require thousands of QC calculations for every new species to serve as training data.

Here we make significant progress toward affordable, large-scale photochemical simulations and virtual screening with ML. To develop a transferable potential we focus on molecules from the same chemical family, studying derivatives of azobenzene, a prototypical photoswitch. The derivatives studied here contain up to 100 atoms, making them the largest systems fit with excited-state ML potentials to date. Combining an equivariant neural network [35] and a physics-informed diabatic model, together with data generated by combinatorial exploration of chemical space, and configurational sampling through active learning, we produce a model that is transferable to large, unseen derivatives of azobenzene. This yields computational savings in excess of six orders of magnitude. Predicted isomerization quantum yields of unseen species are well-correlated with experimental values. The model is used to predict the quantum yield for over 3,100 hypothetical species, revealing rare molecules with extremely high *cis-to-trans* and *trans-to-cis* quantum yields.

Results

Azobenzene photoswitches. This work focuses on the photoswitching of azobenzene derivatives, but the methods are general and can be applied to other chemistries and other excited state processes. Azobenzene derivatives can exist as *cis* and *trans* conformers. The conformations are local minima in the ground state, but not in the excited state. Photoexcitation of either can therefore induce isomerization into the other (see the potential energy schematics in Figs. 1(a) and 2(b)). A key experimental observable is the quantum yield, defined as the probability that excitation leads to isomerization. The yield depends critically on the dynamics near conical intersections (CIs), configurations in which the excitation energy is zero. In these regions the electrons can return to the ground state with non-zero probability.

Many approaches have been developed over several decades to model such non-adiabatic transitions. These include *ab initio* multiple spawning [11] and cloning [46]; Ehrenfest dynamics [9, 10]; coherent switching with decay of mixing [47]; the variational multi-configurational Gaussian method [48]; exact factorization [49–53]; the multi-configuration time-dependent Hartree (MCTDH) method [54, 55]; Gaussian MCTDH [56]; and trajectory surface hopping [57]. A recent review of these methods can be found in Ref. [3]. Surface hopping is a popular approach because of its simplicity and efficiency. In this method, independent trajectories are simulated with stochastic hops between potential energy surfaces (PESs). Depending on the curvature of the PESs and the location of the hop, a trajectory can end in the original isomer or in a new isomer (Figs. 1(a) and 2(b)). The quantum yield is the proportion of trajectories that end in a new isomer. Our goal is to predict the quantum yield of azobenzene derivatives after excitation from the singlet ground state (S_0) to the first singlet excited state (S_1). This can be accomplished with the surface hopping approach described above, using a fast surrogate ML model to generate the PESs. The impact of considering only the first excited state is discussed in Supplementary Sec. IV.

ML architecture and training. Our model is based on the PaiNN neural network [35], which uses equivariant message-passing to predict molecular properties. In this approach, an initial feature vector is generated for each atom using its atomic number. The vector is then updated through a set of neural network

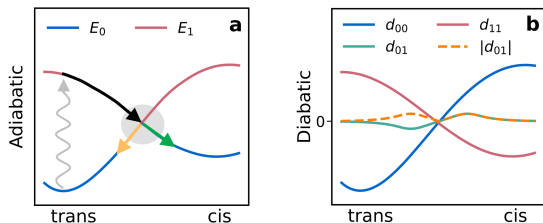


Figure 1. Depiction of the potential energy surfaces in azobenzene derivatives. (a) S_0 and S_1 adiabatic energies, with the CI region shaded in gray. Initial excitation is shown with a vertical zigzag line. Trajectories prior to hopping are shown in black. Reactive and unreactive trajectories after hopping are shown in green and yellow, respectively. (b) Diabatic energies $d_{nm} \equiv (\mathbf{H}_d)_{nm}$. The diagonal diabatic elements cross and become re-ordered along the isomerization coordinate. A CI occurs when the diagonal diabatic elements cross and the off-diagonal element becomes zero.

operations involving “messages”, which incorporate the distance, orientation, and features of atoms within a cutoff distance. A series of updates leads to information being aggregated from increasingly distant atoms. Once the updates are complete, the atomic features are mapped to molecular energies using a neural network.

This architecture can be used to predict energies and, through automatic differentiation, the forces for each state. However, models that predict adiabatic energies have a basic shortcoming for non-adiabatic molecular dynamics (NAMD). Since surface hopping is largely controlled by the energy gap when it is close to zero, small errors in the energies can lead to exponentially large errors in the hopping probability [58, 59]. This in turn can cause large errors in observable quantities like the quantum yield. This point is discussed in further detail in Supplementary Sec. II A. Further, since CIs are non-differentiable cusps in the energy gap, they are difficult to fit with neural networks. For N atoms in a molecule, the network must predict two different energies that are exactly equal in $3N - 8$ dimensions. We found this to be particularly challenging for *trans* species that are outside the training set. As shown in Supplementary Sec. VII, small errors in the gap lead to the incorrect prediction that many species never hop to the ground state.

To remedy this issue we introduce a model based on diabatic states, which we call DANN (*diabatic artificial neural network*; Fig. 2(a)). The approach builds on previous work using neural networks for diabatiza-

tion [60–62]. Much of the previous work could only be used for specific system types, such as semi-rigid molecules [61] and coupled monomers, and is thus not applicable to azobenzene. None of the methods have been used for large systems with significant conformational changes [60, 62], such as azobenzene derivatives. Further, our work uses diabatization to ease the fitting of adiabatic states across chemical space. In particular, it addresses the issue of gap overestimation near conical intersections of unseen species, as described in Supplementary Secs. II and VII. Our work uses diabatization to address this problem, whereas previous work developed diabatic states because of their favorable theoretical properties. We also note that gap overestimation in unseen species is both a newly-identified and newly-addressed problem, as previous work in ML-NAMD focused on single species only [39–45].

The diabatic energies form a non-diagonal Hamiltonian matrix, \mathbf{H}_d , which is diagonalized to yield adiabatic energies. When a 2×2 sub-block of \mathbf{H}_d has diagonal elements that cross, and off-diagonal elements that pass through zero, a CI cusp is generated (Fig. 1). The diabatic energies that generate the cusp are smooth, which makes them easier to fit with an interpolating function than the adiabatic energies. In the DANN architecture, smoothness is imposed through a loss function related to the non-adiabatic coupling vector (NACV). The loss minimizes the value that the NACV takes when it is rotated from the adiabatic basis (Eq. (3)) into the diabatic basis. The NACV measures the change in overlap between two wavefunctions after a small nuclear displacement. If the NACV between two states is zero, then their wavefunctions must change slowly in response to a nuclear perturbation. Therefore, their energies cannot form the cusp in Fig. 1(a), and must instead resemble the smooth energies in Fig. 1(b).

The DANN model was trained on SF-TDDFT [63] calculations for 567,037 geometries, using the 6-31G* basis [64] and BHLYP [65] exchange-correlation functional. Unlike traditional TDDFT [66], SF-TDDFT provides an accurate description of the CI region [67], and, unlike multi-reference methods, is fairly fast and requires no manual parameter selection. The configurations were sampled from 8,269 azobenzene derivatives, of which 164 were taken from the experimental literature. The remaining molecules were generated from combinatorial substitution using common literature patterns (Supplementary Tables S10 and S11).

The data generation process is shown in Fig. 2. Initial data was generated through *ab initio* NAMD with

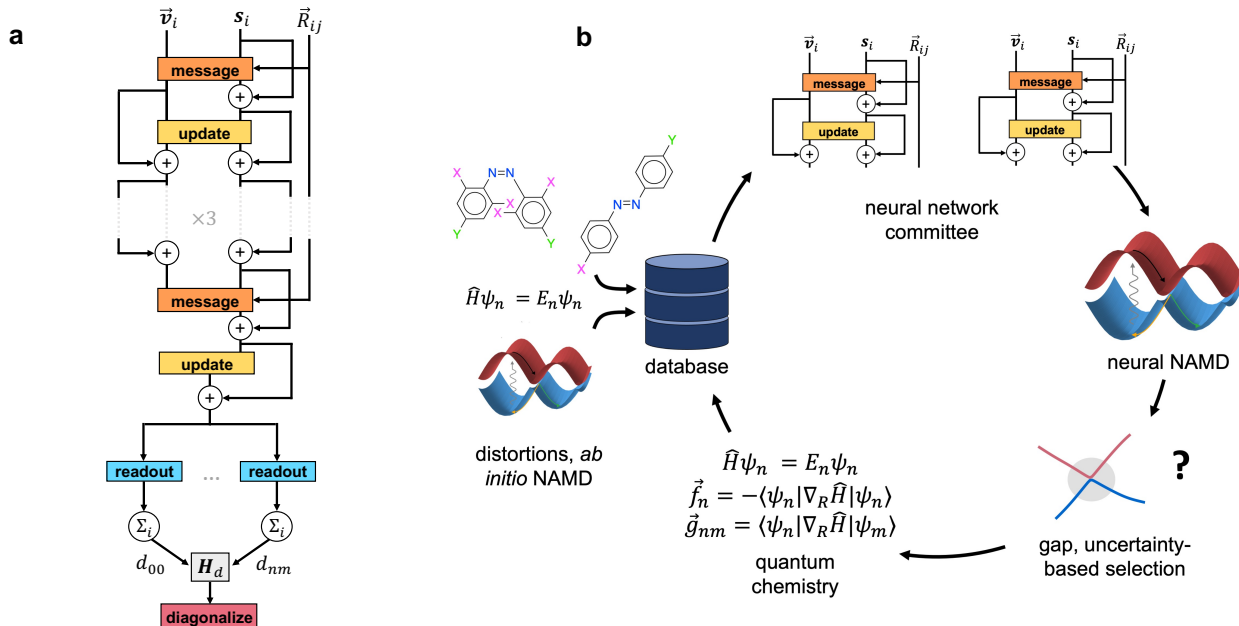


Figure 2. (a) Schematic of the DANN architecture, which is based on the PaiNN model. Scalar atomic features s_i and vectorial atomic features \vec{v}_i are updated through messages from neighboring atoms. The s_i are then mapped to atomic energies, which are summed to produce the diabatic Hamiltonian \mathbf{H}_d . The diabatic matrix is diagonalized to produce adiabatic quantities. (b) Schematic of the active learning loop. Geometries and QC data are first generated through *ab initio* NAMD, normal mode sampling, and inversion/rotation about the central N=N double bond. Two neural networks are then trained on the data and used to perform DANN-NAMD. Newly generated geometries with high committee variance and/or low predicted gaps receive QC calculations. The new calculations are added to the training data, the networks are retrained, and the cycle is repeated until convergence.

		E_0	E_1	ΔE_{01}	$(\Delta E_{01})_{\text{small}}^a$	\vec{F}_0	\vec{F}_1	\vec{g}_{01}
Seen species	MAE (\downarrow)	0.86	1.01	0.75	0.47	1.00	1.17	0.87
	R^2 (\uparrow)	1.00	1.00	1.00	0.97	0.99	0.99	0.84
Unseen species	MAE (\downarrow)	3.06	3.77	1.89	0.97	1.72	2.31	1.36
	R^2 (\uparrow)	0.99	0.98	0.98	0.95	0.97	0.86	0.50

^a For these R^2 calculations, we computed the total sum of squares using $\text{mean}\{\Delta E_{01}\}$ instead of $\text{mean}\{(\Delta E_{01})_{\text{small}}\}$. The mean predictor should not know *a priori* which gaps are small, and hence should predict the mean of all gaps.

Table I. MAE and coefficient of determination (R^2) of the DANN model for various quantities. Units are kcal/mol for energies and kcal/mol/Å for forces and force couplings. E_i are energies, \vec{F}_i are forces, ΔE_{01} is the energy gap, and \vec{g}_{01} is the force NACV. $(\Delta E_{01})_{\text{small}}$ denotes the energy gap when it is under 4.6 kcal/mol (0.2 eV).

164 species from the literature, together with normal-mode sampling and distortions of the combinatorial species to near-CI regions. The remaining data was generated through active learning. In each cycle we trained a committee of models, used one model to perform NN-NAMD, and used the committee variance

and energy gap to choose NAMD geometries for new quantum chemistry calculations. The cycle was repeated five times in total; further details can be found in the Methods section.

Validation. To test whether the model could reproduce experimental results for unseen molecules, we evaluated it on species that were outside the training set. The test set contained 40 species (20 *cis/trans* pairs), including 33 with experimental S_1 quantum yields in non-polar solution. Non-polar solution was chosen because it is the closest to the gas-phase conditions simulated here. Solvent effects can be easily incorporated into the model through transfer learning to implicit solvent calculations. Previously this was shown to require new calculations for only a small proportion of the training set [37].

The performance of the model is summarized in Table I. Statistics are shown for both seen and unseen species. The former contains species that are in the training set, but geometries that are outside of it. The geometries were selected with the balanced sampling criteria described in Supplementary Sec. X. Geometries from unseen species were generated with DANN-NAMD using the final trained model. Half of the DANN-NAMD geometries were selected randomly from the full trajectory and half by proximity to a CI (Supplementary Eq. (S13)). 100 configurations were chosen for each molecule.

For species in the training set, all quantities are accurate to within approximately 1 kcal/mol(\AA). Apart from the NACV, all quantities have R^2 correlation coefficients close to 1. The R^2 of the NACV is 0.84. This may be somewhat low because diabatization cannot remove the curl component of the NACV in the diabatic basis [68, 69]. This would also explain the low R^2 value for the NACV in Ref. [42], which computed it as the gradient of a scalar. For molecules outside the training set, all quantities apart from the energies have an error below 3 kcal/mol(\AA). The energy gaps and ground state forces have R^2 correlation coefficients near 1. The gap error of 1.89 kcal/mol should be contrasted with the error of 15 kcal/mol in Ref. [32], which applied semi-empirical methods to azobenzene. The errors in the excited state forces are slightly larger, but still quite low. The correlation coefficient for the force NACV g_{01} is rather poor. As described in Supplementary Sec. VII, the yields of *trans* derivatives are better correlated with experiment when using Zhu-Nakamura surface hopping than Tully’s method. The latter uses the NACV and the former does not, so part of the difference may be explained by the high error in the force NACV. Nevertheless, there is still reasonable agreement between Tully’s method and experiment, suggesting that errors in the force NACV do not spoil the dynamics.

Figure 3(a) shows snapshots from an example

DANN-NAMD trajectory, and panel (b) shows random samples of the hopping geometries. Reactive hopping geometries are shown on top, and non-reactive ones are shown below. The molecule is the (aminomethyl)pyridine derivative **26**, with the species numbering given in Supplementary Tables S12 and S13. The overlays show *cis-trans* isomerization proceeding through inversion-assisted rotation, consistent with previous work [70]. The dominant motion is rotation, with the CNNC dihedral angle increasing in magnitude from -10° at equilibrium to -86° at the hopping points. Significant changes also occur in the CNN and NNC angles, with each transitioning from 123° to either 113° or 135° .

The predicted PES in the branching space (\vec{g}, \vec{h}) is shown beside the geometries. \vec{g} is the direction of the force coupling and $\vec{h} \propto \nabla_R(\Delta E_{01})$ is the direction of the gap gradient. Each vector was computed with automatic differentiation using Eq. (1). The diabatic energies, adiabatic energies, and gap are shown from top to bottom. We see that the model generates a true CI, in which the S_0 and S_1 energies are exactly equal. Further, the degeneracy is lifted in both the \vec{g} - and \vec{h} -directions, so that the S_1 energy and gap each form a characteristic cone. These hallmarks of CIs are built into the model because the adiabatic energies are eigenvalues of a diabatic matrix. For example, the cone emerges from the fact that $d_{11} - d_{00}$ and d_{01} each pass linearly through zero in different directions [71].

Figure 3(c) indicates that the predicted and experimental quantum yields of unseen species are correlated. The yields are for the 33 *cis* and *trans* species with experimental data in Supplementary Table S12. The R^2 value is 0.42, and the Spearman rank correlation coefficient ρ is 0.74. While the R^2 value is somewhat low, the Spearman rank correlation is high. The Spearman coefficient measures the accuracy with which the model ranks species by quantum yield. ρ only compares orderings, while R^2 compares the model error to the error of a mean predictor. This means that ρ is a more forgiving metric, and also a more relevant metric for virtual screening. Since *cis* isomers have yields 2 to 3 times higher than *trans* isomers, the high value of ρ means that the model properly separates the isomers into low- and high-yield groups.

Further, as shown in Supplementary Figs. S5 and S7, the model produces meaningful rankings among *trans* species. The correlation coefficients are $\rho = 0.32$ using Tully’s method [57] and $\rho = 0.57$ using the Zhu-Nakamura approach [72]. The model is largely able to differentiate between high- and low-yield *trans* deriva-

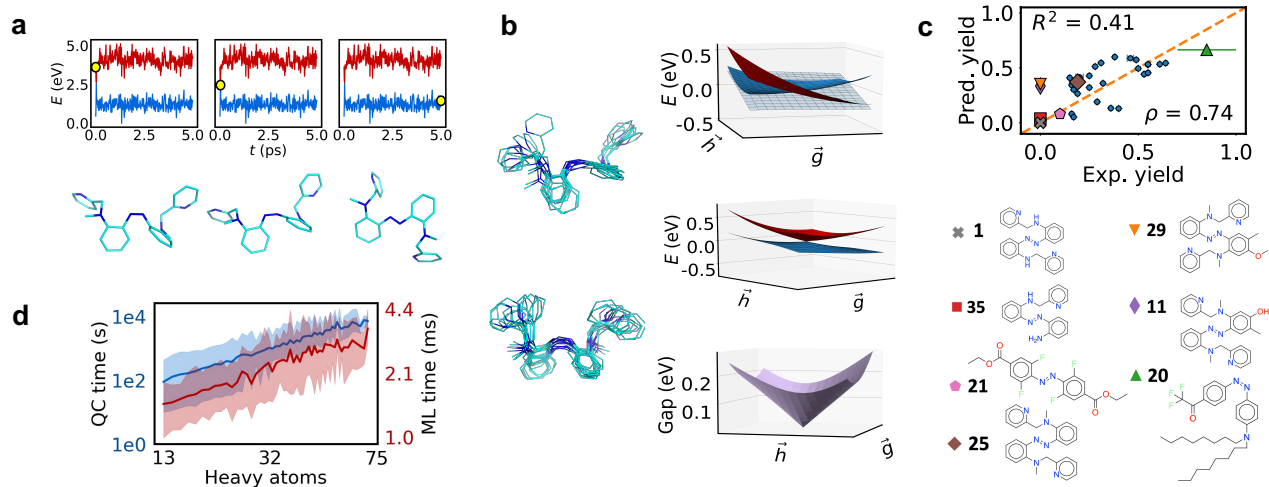


Figure 3. (a) Selected trajectory frames for a molecule outside the training set. The top panels show the S_0 and S_1 energy as a function of time. A yellow dot indicates the time at which the snapshot below was taken. (b) Left: Overlay of selected hopping geometries from reactive (top) and unreactive (bottom) trajectories. Right: PES as a function of branching plane coordinates at one of the reactive hopping geometries. Diabatic energies, adiabatic energies, and adiabatic gaps are shown from top to bottom. The diabatic coupling is shown in gray. (c) Predicted vs. experimental quantum yield for 33 species outside the training set. The R^2 value and Spearman rank correlation ρ are both shown. Color-coded data points are defined below. (d) Node time for QC and ML calculations.

tives. Several such molecules are of interest. They are color-coded in the plots, with the legend given below. A full list of predictions is given in Supplementary Table S12. We see, for example, that the (aminomethyl)pyridine derivatives **1** and **35** are both predicted to have near-zero yields. These species do not isomerize from *trans* to *cis*, because strong N-H hydrogen bonds lock the planar *trans* conformation in place [73]. Replacing the NH group in **1** with N – CH₃ gives species **25**. This molecule isomerizes because there is no hydrogen bonding. This, too, is predicted by the model. Further, the hepta-tert-butyl derivative **17** has an experimental and predicted yield of zero. This is likely because of steric interactions among the bulky tert-butyl groups. While able to account for these two different mechanisms, the model fails to predict the subtle electronic effects in species **11** and **29**. Resonance interactions between oxygen lone pairs and the azo group modify the PES, such that there is no rotational CI [74]. There is instead a concerted inversion CI, which occurs too early along the path between *trans* and *cis* to allow for isomerization. The changes in the PES may either be too small or too specific to the substituents for the model to predict without fine tuning. Finally, derivatives with high yields are

partly distinguished from those with low but non-zero yields. An example is **21**, whose experimental yield of 10% is half that of *trans*-azobenzene. The model properly identifies this molecule as having a low yield, but also mistakenly does the same for several high-yield species. The accuracy for unseen species could always be improved with transfer learning, in which the model is fine-tuned with a small number of calculations from a single molecule (discussed below). This would increase the computational cost, but would still be orders of magnitude less expensive than *ab initio* NAMD.

While meaningful correlations are produced for *trans* species, the same is not true of *cis* molecules ($\rho = 0.02$). This may be because there are no *cis* derivatives with zero yield. Nevertheless, the model properly identifies **20** as having the highest yield. Further, it does not mistakenly assign a zero yield to any derivative. This is noteworthy because, as shown in Fig. 4(a) and (b), several hypothetical *cis* species are predicted to have zero yield. Synthesis of non-switching *cis* derivatives and comparison to predictions could therefore be of interest in the future.

Overall, we observe moderate correlation between predicted and experimental yields. The Spearman cor-

relation is high when including both isomers, moderate for *trans* isomers, and low for *cis* isomers. The R^2 value, a measure of numerical error compared to that of a mean predictor, is moderate when including both isomers and near-zero when separating them. Indeed, the MAEs of the mean predictor are 9.5%, 10.3%, and 17.7% for *trans*, *cis*, and all species, respectively. The model MAEs before (after) subtracting the mean signed error are 14.4% (13.5%), 11.5% (11.2%) and 13.2% (13.0%). In addition to model error, sources of error include inaccuracies in SF-TDDFT, approximations in surface hopping, solvent effects, and experimental uncertainty. These are discussed in depth in Supplementary Sec. IV. Each source of error affects both R^2 and ρ , but is expected to have a larger effect on R^2 . The rank correlation with experiment is encouraging given the difficulty of the task, as captured by the sensitivity of the yield to model errors in the PES [72], and given the sources of error outside the model. Further, as discussed below, DANN provides an excellent starting point for fine-tuned, molecule-specific models that can be used for high-accuracy simulations of single species.

Figure 3(d) shows that DANN-NAMD is extremely fast. The plot shows the node time, defined as $t_{\text{calc}}/n_{\text{calc}}$, where t_{calc} is the calculation time per geometry, and n_{calc} is the number of parallel calculations that can be performed on a single node. We see that ML speeds up calculations by five to six orders of magnitude. The direct comparison of the pre-trained model node times and QC node times is appropriate because the model generalizes to unseen species. This means that it incurs no extra QC cost for any future simulations. The minimum speedup corresponds to the smallest molecules (14 heavy atoms or 24 total atoms), and the maximum to the largest molecules (70 heavy atoms or 99 total atoms). This reflects the different scaling of the QC and ML calculations. Empirically we see that DANN scales as $N^{0.49}$ for N heavy atoms, while SF-TDDFT scales as $N^{2.8}$. These values come from fitting the timings to $t = A \cdot N^x$, where t is the computational time, A and x are fitted constants, and N is the number of heavy atoms. DANN’s apparent sub-linear scaling is an artifact of diagonalizing \mathbf{H}_d ; when the diagonalization is removed, the scaling becomes linear. This is the expected scaling for a message-passing neural network with a fixed cut-off radius. Evidently diagonalizing \mathbf{H}_d introduces a large overhead with weak dependence on system size. Nevertheless, we see that DANN is still quite fast.

Virtual screening. Having shown that the model is fast and generalizes in the chemical and configurational space of azobenzenes, we next used it for virtual screening of hypothetical compounds. We first retrained the network on all available data, including species that were originally held out, for a total of 631,367 geometries in the training set. We then predicted the quantum yields of 3,100 combinatorial species generated through literature-informed substitution patterns, as in Ref. [75]. This screen served two purposes. The first was to gather statistics about the distribution of photophysical properties of azobenzenes at a scale not accessible to experiments or traditional simulations. The second was to identify molecules with rare desirable properties. In particular, we sought to find molecules with high $c \rightarrow t$ or $t \rightarrow c$ quantum yields and redshifted absorption spectra. The former is important because increasing the ratio $\text{QY}_{a \rightarrow b} / \text{QY}_{b \rightarrow a}$, where QY is the quantum yield, can lead to more complete $a \rightarrow b$ transformation under steady state illumination. This is critical for precise spatial control of drug activity when the two isomers have different biological effects [76]. Redshifting is a crucial requirement for photo-active drugs, since human tissue is transparent only in the near-IR [76].

The results are shown in Fig. 4. Panels (a) and (c) show the predicted yield vs. mean gap. For each species we averaged the gap over the configurations sampled during neural network ground state MD. The thermal averaging led to a typical blueshift of 0.2-0.3 eV relative to the gaps of single equilibrium geometries. The mean excitation energies are 2.95 eV for *cis* derivatives and 2.84 eV for *trans* species; the gaps are 2.98 eV and 2.97 eV for the respective unsubstituted compounds. The average gaps and their differences are similar to experimental measurements for azobenzene [77]. The average $c \rightarrow t$ and $t \rightarrow c$ yields are 54% and 24%, respectively, while those of the unsubstituted species are 59% and 37%. These are consistent with experimental results in non-polar solution, for which the base compound has yields of 44-55% and 23-28% [77]; the former is closer to 55% on average. However, the yield of the base *trans* compound is overestimated with respect to both theory and experiment [7, 72, 77]. The mean (median) proportion of trajectories ending in the ground state after 2 ps are 92% (100%) for *cis* species and 31% (17%) for *trans* species. The standard deviations are 25% and 30%, respectively.

Panels (b) and (d) show the yield plotted against the isomeric stability, defined as $E_{\text{trans}} - E_{\text{cis}}$ for *trans* isomers and $E_{\text{cis}} - E_{\text{trans}}$ for *cis* isomers. The energy E is the median value of the configurations sampled in

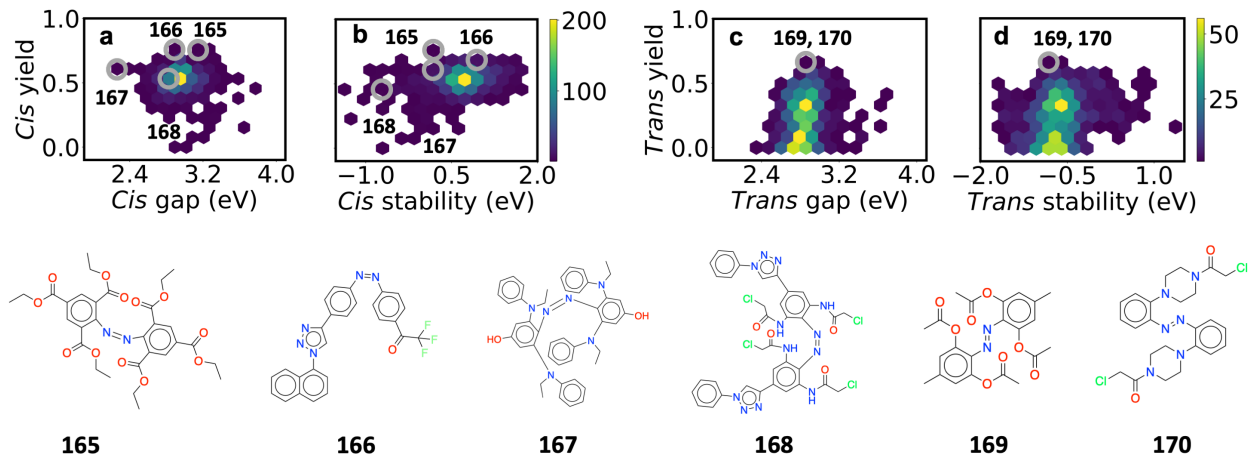


Figure 4. Results of virtual screening. Species of interest are circled in gray and shown below the plots. (a) Predicted yield vs. excitation energy for *cis* derivatives. (b) Predicted yield vs. stability for *cis* derivatives. (c)-(d) As in (a)-(b), but for *trans* derivatives.

the ground state; we used the median to reduce the effect of outlier geometries. On average the *trans* isomers are more stable than the *cis* isomers by 0.66 eV (15.3 kcal/mol), which is similar to experimental values over 10 kcal/mol for azobenzene [78]. The stability is of interest for three reasons. First, a large absolute value indicates that one isomer is dominant at room temperature. This is essential for photoactive drugs, and is the case for regular azobenzene. Second, an inverted stability, in which *cis* is more stable than *trans*, allows for stronger absorption at longer wavelengths. This is because the dipole-forbidden $n-\pi^*$ (S_1) transition is significantly stronger for *cis* than for *trans* [77]. Third, in photopharmacology, one often wants to deliver a drug in inactive form, and activate it with light in a localized region. If *trans* happens to be active and *cis* inactive, then localized activation is only possible if *cis* is more stable.

Several species of interest are shown in Fig. 4. The molecules **165** and **166** have predicted $c \rightarrow t$ yields of $75 \pm 6\%$ and $72 \pm 6\%$, respectively, which are well above the *cis* average of 55%. The species **169** and **170** have predicted $t \rightarrow c$ yields of $66 \pm 7\%$ and $63 \pm 10\%$, respectively, which are three times the average *trans* yield. Molecule **167** is highly redshifted, with a mean predicted gap of 2.26 eV (548 nm), and a standard deviation of 0.87 eV. QC calculations on the geometries sampled with DANN gave a gap of 2.26 ± 0.61 eV, in good agreement with predictions. The mean gap is lower than the median of 2.52 eV, which reflects the

presence of several ultra-low gap structures. The low gap and large variance mean that **167** may be able to absorb in the near IR. The redshifting is likely because of the six electron donating groups, which increase the HOMO energy, together with the crowding of the four *ortho* substituents. The latter distorts the molecule, leading to twisted configurations with smaller gaps. Finally, species **168** is more stable in *cis* form than *trans* form. The predicted *cis* stability is -0.79 eV (-18 kcal/mol), in good agreement with the QC prediction of -0.92 eV (-21 kcal/mol). As mentioned above, this inverted stability can be a desirable property for photopharmacology.

To validate the yield results, we performed DANN-NAMD using highly accurate species-specific models. As described in Supplementary Sec. XIV B, we generated a model for each species by refining the base network with data from that species alone. The data was generated through several active learning cycles, resulting in 1,200 to 2,500 training geometries for each compound. We used this approach in place of *ab initio* NAMD because of the latter’s prohibitive cost for large molecules. The QC computational cost for fine-tuning was at most 3% of that of an *ab initio* simulation, and hence far less demanding. The average gradient calculation for a molecule with 50 atoms took 58 minutes for two surfaces using 8 cores, and the average NACV calculation took 55 minutes. Fine-tuning with 2,000 geometries for a medium-sized molecule would thus take 30,000 core hours. For *ab initio* NAMD,

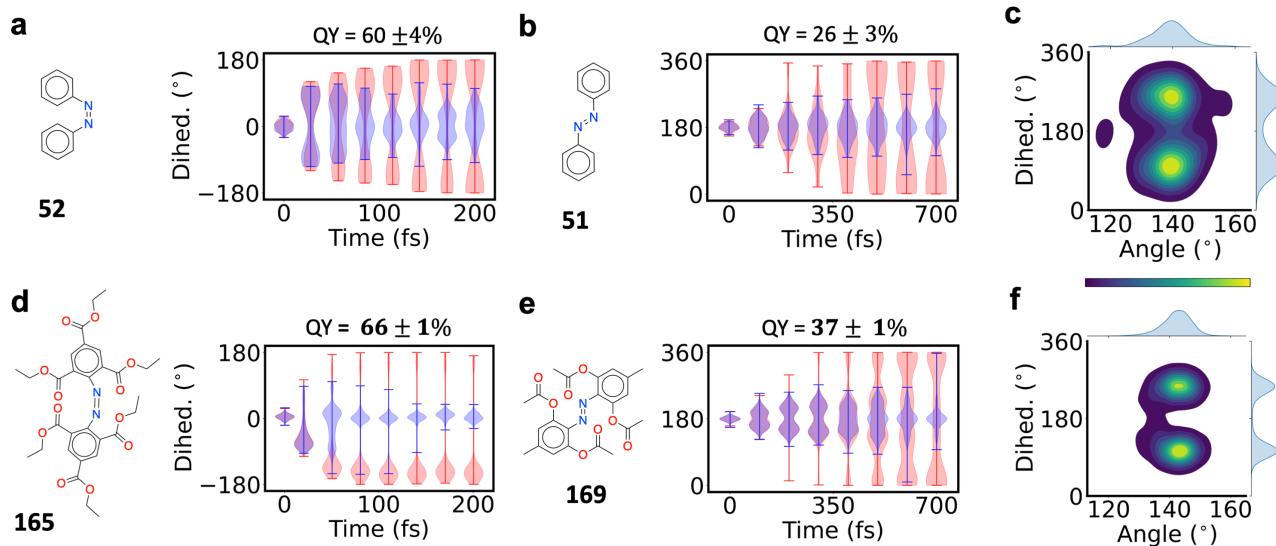


Figure 5. Visualization of high-accuracy non-adiabatic dynamics for several compounds of interest. (a)-(b), (d)-(e) Violin plots showing the CNNC dihedral angle vs. time. Reactive and non-reactive NAMD trajectories are shown in red and blue, respectively. The violin width at a given dihedral angle indicates the density of trajectories with that angle. The yield of each compound is shown above the plots. For ease of visualization we have used the range $[-180, 180]$ for *cis* dihedral angles and $[0, 360]$ for *trans* dihedral angles. (c) Distribution of hopping geometries for *trans* azobenzene. (f) As in (c), but for the derivative **169**. The density is visualized with kernel density estimation as a function of the CNNC dihedral and $\max(\alpha_{\text{CNN}}, \alpha_{\text{NNC}})$, where α is an angle. Yellow corresponds to the highest density and blue to the lowest. The marginal distributions over single coordinates are shown above and to the right of each plot.

a conservative estimate of 100 trajectories run for 1 ps each, with only one gradient computed per frame, would take 780,000 core hours.

We also computed the yields of *cis* and *trans* azobenzene for comparison. For these species we used full *ab initio* simulations, because of the central role of the unsubstituted compound as a reference point and because simulations were fairly affordable for such small molecules. Issues with spin contamination also hampered the fine-tuning process for these compounds (see Supplementary Sec. XIV B).

Initially we generated refined models for species **165**, **167**, **169** and **170**. It became clear early on that only **165** and **169** had high yields, and so we focused on those molecules thereafter. Using molecule-specific models, we computed the quantum yields of **165** and **169** to be $66 \pm 1\%$ and $37 \pm 1\%$, respectively. The computed yields for *cis* and *trans* azobenzene are $60 \pm 4\%$ and $26 \pm 3\%$, respectively, which are in excellent agreement with experiment [77]. Both of the new molecules have higher quantum yields than the associated base compounds. The improvement is particularly large for species **169**: its yield is 11 points higher than *trans*

azobenzene, which is a relative enhancement of 42 percent. We argue below that that this significant increase has an intuitive physical explanation.

The dynamics of the four molecules are shown in Fig. 5. Panels (a) and (b) show the CNNC dihedral angle vs. time for azobenzene, and panels (d) and (e) show the same for the derivatives. Both the substituted and unsubstituted *cis* isomers rapidly proceed through a rotational CI, but the derivative rotates much more quickly. Indeed, we see that the isomerization of the derivative is complete within 75 fs, while the base compound takes nearly 130 fs. The excited state lifetimes are 34.2 ± 0.3 fs and 63 ± 3 fs for the derivative and base compound, respectively, indicating that the former reaches the CI earlier than the latter. These observations may explain the enhanced yield, since a higher rotational velocity leads to more efficient isomerization [79]. We also note that the derivative rotates in only the counter-clockwise direction, while *cis* azobenzene rotates in both directions, but this is not expected to affect the yield.

The two *trans* molecules behave in qualitatively different ways. In *trans* azobenzene, the distribution of

dihedral angles slowly widens with time (Fig. 5(b)). This is consistent with a rotational barrier [7, 72], as different trajectories overcome the barrier at different times, and so the torsion angle becomes uniformly distributed. Additionally, as seen in the marginal dihedral distribution of Fig. 5(c), many of the geometries hop near 180° . This agrees with Ref. [7], which identified a non-reactive planar CI and a reactive twisted CI as the main hopping points for *trans* azobenzene. The non-reactive CI leads exclusively back to *trans*, while the reactive CI leads to *cis* and *trans* in different proportions. Using the method described in Supplementary Sec. XIV C, we found that 26% of the trajectories proceed through the planar CI and 74% through the rotational CI. This is the same distribution reported in Ref. [72]. Approximately 36% of the rotational trajectories generate *cis* azobenzene, giving an overall yield of 26%. This is in good agreement with previous computational and experimental values [7].

By contrast, nearly all trajectories of **169**, including non-reactive trajectories, rotate significantly. This can be seen in the marginal dihedral distribution in Fig. 5(f), in which the hops are tightly localized around $180 \pm 77^\circ$. Only 5% of the trajectories hop at the planar CI, which is five times lower than the base compound. Additionally, the yield of the rotational trajectories increases from 36% to 40%. The inhibition of the planar CI pathway, together with the enhancement of the rotational yield, leads to an overall yield increase from 26% to 37%. While the enhanced reactive yield does not have a simple explanation, the reason for the planar pathway inhibition can be clearly seen in Fig. 5(e). Whereas the rotation of **51** is stochastic, leading to a uniform distribution of angles, the rotation of **169** is initially concerted. Nearly all trajectories rotate in unison to a dihedral angle of $180 \pm 45^\circ$ at 300 fs. Past 300 fs, hopping begins and the trajectories separate from each other. Hence they proceed through the rotational reactive CI, and become distributed between 0° and 360° after hopping. The planar non-reactive CI is avoided because of the molecule’s initial rotation. This explanation is consistent with the presence of bulky *ortho* groups, which twist the equilibrium structure and hence weaken the N=N double bond. This lowers the excited state barrier to rotation, which leads to an initial torsion and hence increases the yield.

Discussion

The DANN model shows high accuracy and transferability among azobenzene derivatives. One limita-

tion is that the unseen species contained functional groups that were present to some degree in the training set. Model performance was generally higher for more highly represented functional groups, though some groups were highly represented yet difficult to fit, while others were weakly represented and well-fit (Supplementary Sec. V). Moreover, the model cannot be applied to other chemical families without additional training data. Further, as shown in Supplementary Sec. VII, it substantially overestimates the excited state lifetime for a number of *trans* derivatives. On the other hand, semi-empirical methods provide qualitatively correct predictions across a variety of chemistries, but cannot match DANN’s in-domain accuracy, and cannot be improved with more reference data. Adding features from semi-empirical calculations, as done in the OrbNet model [80], may therefore prove useful in the future. Recent developments accounting for non-local effects and spin states have improved neural network transferability [36], and could also be beneficial for excited states. The model could be further improved with high-accuracy multi-reference calculations, solvent effects, and the inclusion of the bright S_2 state. The use of spin-complete methods in particular is of crucial importance, since spin contamination prevented fine-tuning the model for the base compounds. It may also have affected the accuracy of the DANN model in general. Thus spin-complete, affordable alternatives are of particular interest [28]. Active learning could be accelerated through differentiable sampling with adversarial uncertainty attacks [81], which would improve the excited state lifetimes. Transfer learning could also be used to improve performance for specific molecules. Only a small number of *ab initio* calculations would be required to fine-tune the model for an individual species.

Diabatization may also prove to be useful for reactive ground states. Reaction barriers can often be understood as transitions from one diabatic state to another [82]. The diabatic basis may make reactive surfaces easier to fit with neural networks.

In conclusion, we have introduced a diabatic multi-state neural network potential trained on over 630,000 geometries at the SF-TDDFT BHHLYP/6-31G* level of theory, covering over 8,000 unique azobenzene molecules. We used DANN-NAMD to predict the isomerization quantum yields of derivatives outside the training set, and the results were well-correlated with experiment. We also identified several hypothetical compounds with high quantum yields, redshifted excitation energies, and inverted stabilities. The

network architecture, diabaticization approach, and chemical and configurational diversity of the training data allowed us to produce a robust and transferable potential. The model can be applied off-the-shelf to new molecules, producing results that replicate those of SF-TDDFT at orders of magnitude lower computational cost.

Methods

Network and training. As explained in Supplementary Sec. I, a unique challenge for non-adiabatic simulations is their sensitivity to the energy difference between states. Using a typical neural network to predict energies and forces for NAMD leads to some molecules becoming incorrectly trapped in the excited state. This is partly caused by overestimation of the gap and/or an incorrectly shaped PES in the vicinity of the CI. To address this issue we introduce an architecture based on *diabatic* states, whose smooth variation leads to more accurate neural network fitting (Fig. 1(b)).

In general diabatic states must satisfy [83]

$$(\mathbf{U}^\dagger [\nabla_R \mathbf{H}_d] \mathbf{U})_{nm} = \begin{cases} -\vec{f}_n, & \text{if } n = m \\ \vec{g}_{nm}, & \text{if } n \neq m. \end{cases} \quad (1)$$

where ∇_R is the gradient with respect to \vec{R} , \mathbf{U} diagonalizes the diabatic Hamiltonian through

$$(\mathbf{U}^\dagger \mathbf{H}_d \mathbf{U})_{nm} = E_n \delta_{nm}, \quad (2)$$

and $\vec{f}_n = -\nabla_R E_n$ is the adiabatic force for the n^{th} state. The dependence on \vec{R} has been suppressed for ease of notation. \vec{g}_{nm} is the force coupling,

$$\begin{aligned} \vec{g}_{nm}(\vec{R}) &= \left\langle \psi_n(\vec{r}; \vec{R}) \left| \nabla_R \hat{H}(\vec{r}, \vec{R}) \right| \psi_m(\vec{r}; \vec{R}) \right\rangle \\ &= (E_m(\vec{R}) - E_n(\vec{R})) \vec{k}_{nm}(\vec{R}), \end{aligned} \quad (3)$$

where $\hat{H}(\vec{r}, \vec{R})$ is the clamped nucleus Hamiltonian, $\psi_n(\vec{r}; \vec{R})$ is the n^{th} adiabatic wavefunction, and the matrix element is an integral over the electronic degrees of freedom \vec{r} . The vector $\vec{k}_{nm}(\vec{R})$ is the derivative coupling:

$$\vec{k}_{nm}(\vec{R}) = \left\langle \psi_n(\vec{r}; \vec{R}) \left| \nabla_R \psi_m(\vec{r}; \vec{R}) \right. \right\rangle \quad (4)$$

Combined with the following reference geometry conditions (Supplementary Sec. I),

$$(E_0, E_1) = \begin{cases} (d_{00}, d_{11}), & \text{if } \vec{R} \in \text{trans} \\ (d_{22}, d_{00}), & \text{if } \vec{R} \in \text{cis}, \end{cases} \quad (5)$$

we arrive at three sets of constraints, Eqs. (1), (2), and (5). In principle only Eqs. (1) and (2) are required for the states to be diabatic. However, we found the reference loss to provide a minor improvement in the gap near CIs (Supplementary Table S1).

We use a neural network to map the nuclear positions \vec{R}_i and charges Z_i to the diabatic matrix elements d_{nm} , and a loss function to impose Eqs. (1), (2) and (5). The adiabatic energies E_n are generated by diagonalizing \mathbf{H}_d , and the forces and couplings by applying Eq. (1) and using automatic differentiation. The design of the network is shown schematically in Fig. 2(a). The general form of the diabatic loss function is

$$\mathcal{L} = \mathcal{L}_{\text{core}} + \mathcal{L}_{\text{ref}} + \mathcal{L}_{\text{nacv}}. \quad (6)$$

Here $\mathcal{L}_{\text{core}}$ penalizes errors in the adiabatic energies, forces, and gaps, \mathcal{L}_{ref} imposes Eq. (5) and $\mathcal{L}_{\text{nacv}}$ imposes Eq. (1) for $n \neq m$. The terms are defined explicitly in Supplementary Eqs. (S1)-(S3).

For the network itself we adopt the PaiNN equivariant architecture [35]. In this approach a set of scalar and vector features

for each atom are iteratively updated through a series of convolutions (Fig. 2(a)). In the message block, the features of each atom gather information from atoms within a cutoff distance, using the interatomic displacements. The scalar and vector features for each atom are then mixed in the update phase. Hyperparameters can be found in Supplementary Table S4. Most were taken from Ref. [35], but some were modified based on experiments with azobenzene geometries. Further details of the PaiNN model can be found in Ref. [35]. Once the elements of \mathbf{H}_d are generated, the diabatic matrix is diagonalized to yield the transformation matrix \mathbf{U} and the adiabatic energies E_n . The vector quantities \vec{f}_n and \vec{g}_{nm} are given by Eq. (1). When non-adiabatic couplings are not required, the \vec{f}_n can be calculated by directly differentiating the E_n . This is more efficient than Eq. (1), since it requires only $M_{\text{ad}} = 2 < M_d(M_d + 1)/2 = 6$ gradient calculations. This approach was used for NAMD runs, which required only diabatic energies, adiabatic energies, and adiabatic forces, while Eq. (1) was used for training.

Data generation and active learning. Data was generated in two different ways. First, we searched the literature for azobenzene derivatives that had been synthesized and tested experimentally. This yielded 164 species (82 *cis* and 82 *trans*). For these species we performed *ab initio* NAMD, yielding geometries that densely sampled configurational space. Second, to enhance chemical diversity, we generated nearly 10,000 species through combinatorial azobenzene substitution. This was done using 48 common literature substituents and four common substitution patterns (Supplementary Tables S10 and S11). We then performed geometry optimizations, normal mode sampling, and inversion/rotation about the central N=N bond to generate configurations. QC calculations were performed on 25,212 combinatorial geometries. All calculations were performed with Q-Chem 5.3 [84], using SF-TDDFT [63] with the BHHLYP functional [65] and 6-31G* basis [64].

Two neural networks were trained on the initial data and used to perform DANN-NAMD. Initial positions and velocities for DANN-NAMD were generated from classical MD with the Nosé-Hoover thermostat [85, 86]. The initial trajectories were unstable because the networks had not been trained on high-energy configurations. To address this issue we used active learning [37, 38] to iteratively improve the network predictions (Fig. 2(b)). After each trajectory we performed new QC calculations on a sample of the generated geometries. For all but the last two rounds of active learning, geometries were selected according to the variance in predictions of two different networks, where the networks were initialized with different parameters and trained with different random batches. In the last two rounds, half the geometries were selected by network variance, and half by proximity to a CI. Further details are given in Supplementary Sec. XIII. The new data was then added to the training set and used to retrain the networks. The cycle was repeated three times with all species and another two times with azobenzene alone. In all, we computed ground-state gradients, excited-state gradients, and NACVs with SF-TDDFT for 641,367 geometries. 96% of the geometries were from the 164 literature species. 88% were generated through *ab initio* NAMD and 8% through active learning. The remaining 4% were from the combinatorial species. 1.5% were generated through geometry optimizations, 1.5% through inversion/rotation, and 1% through normal-mode sampling.

We initially set out to train a model using energies and forces alone. Since analytic NACVs are unavailable for many *ab initio* methods, an adiabatic architecture could have been used with a wider variety of methods. NACVs also add computational overhead, and so generating training data for an adiabatic model would have taken less time. To this end we initially used the Zhu-Nakamura (ZN) surface hopping method [79], which only requires adiabatic energies and forces. However, the issues with adiabatic models described in Supplementary Sec. VII led us to develop the diabatic approach. Since diabatic states can be used with any surface hopping method, we used the diabatic model to perform Tully’s fewest switches (FS) surface hopping [57] after the last round of active learning. All results in the main text use the FS method. A comparison of FS and ZN results is given in Supplementary Sec. VII.

Data availability

The quantum chemistry data is available at <https://doi.org/10.18126/unc8-336t>. A detailed description of how to load and interpret the data is given in the README file. Source data of experimental and predicted quantum yields are provided with this paper.

Code availability

Trained models and computer code are available in the Neural Force Field repository at <https://github.com/learningmatter-mit/NeuralForceField>. Notebook tutorials explain how to train the models and perform DANN-NAMD.

Acknowledgements

The authors thank Wujie Wang, Daniel Schwalbe-Koda, Shi Jun Ang (MIT), Kristof Schütt, and Oliver Unke (Technische Universität Berlin) for scientific discussions and access to computer code. Harvard Cannon cluster, MIT Engaging cluster, and MIT Lincoln Lab Supercloud cluster at MGHPC are gratefully acknowledged for computational resources and support. Financial support from DARPA (Award HR00111920025) and MIT-IBM Watson AI Lab is acknowledged.

Author Contributions

S.A. conceived the project, and developed the methodology with R.G.-B. and E.S. S.A. performed the calculations under the guidance of R.G.-B. S.A. wrote the first draft of the manuscript, and all authors contributed to the final version.

Competing interests

The authors declare no competing interests.

Supplementary Information

Contents

I. Extended methods	15
A. Relevance of adiabatic gap	15
B. Diabatization	15
C. Network loss	16
D. Data generation	17
II. Model accuracy and ablation studies	17
III. Experimental data	18
IV. Sources of error	19
A. Experimental error	19
B. Computational error	20
V. Influence of different functional groups	20
VI. Spin contamination	21
VII. Surface hopping results with different models	22
VIII. Architecture	23
IX. Training	25
X. Balanced sampling	26
XI. Dynamics	28
XII. Fewest switches implementation	29
XIII. Active learning	31
XIV. Validation	32
A. <i>Ab initio</i> NAMD	32
B. Transfer learning	32
C. Conical intersection pathways	35
XV. Figure details	35
A. Computational speed-up	35
B. Diabatic energies	35
XVI. Intensive and extensive quantities	36
XVII. Proof of diabaticity	37
XVIII. Training species	39

I. Extended methods

A. Relevance of adiabatic gap

Accurate prediction of the energy gap is crucial for generating reliable NAMD results. The gap controls the hopping rate, and in turn observables like the photoisomerization quantum yield, excited state lifetime and time-resolved photoelectron spectrum [87]. To understand the importance of the gap, consider that in most approaches, the hopping rate is determined by the derivative form of the NACV [57, 88]. The derivative coupling between states n and m can be written as $\vec{g}_{nm}/\Delta E_{nm}$, where ΔE_{nm} is the energy gap and \vec{g}_{nm} is the force coupling [83, 89]. The gap in the denominator leads to singular derivative coupling at CIs, and therefore to a guaranteed hop. The coupling can alternatively be obtained from the gap alone from through its first and second derivatives [44]. The energy difference also features prominently in the Zhu-Nakamura method, in which the hopping rate is approximately exponential in the square of the gap [79]. Therefore, it is crucial to accurately predict the energy gap in any NAMD simulation.

B. Diabatization

The adiabatic energies of a system, $\{E_n(\vec{R})\}$, are the energies produced by a QC calculation. They depend on the nuclear coordinates \vec{R} , and form the usual Born-Oppenheimer PESs. In the adiabatic basis the nuclear kinetic energy operator, related to $\nabla_{\vec{R}}^2$, is non-diagonal. Its off-diagonal elements are related to the NACVs, which generate transitions between adiabatic PESs [90]. The derivative form of the NACV also becomes singular at CIs, which is an undesirable numerical property. The diabatic basis is designed to remove this singularity [91]. The diabatic Hamiltonian $\mathbf{H}_d(\vec{R})$ is a rotation of the adiabatic energies into a new basis, given by $\mathbf{H}_d(\vec{R}) = \mathbf{U}(\vec{R}) \text{diag}(\{E_n(\vec{R})\}) \mathbf{U}^\dagger(\vec{R})$. Here diag denotes a diagonal matrix, and $\mathbf{U}(\vec{R})$ is a rotation matrix that depends on the nuclear coordinates (see Eq. (2)). \mathbf{H}_d can also be viewed as the clamped nucleus Hamiltonian expressed in the basis of diabatic wave functions. These wave functions are given by $\psi_{d,n}(\vec{r}; \vec{R}) = \sum_m U_{nm}(\vec{R}) \psi_{\text{ad},m}(\vec{r}; \vec{R})$, where $\psi_{\text{ad},m}$ is the m^{th} adiabatic wave function and \vec{r} denotes the electronic coordinates.

The diabatic states are defined such that the nuclear kinetic energy is approximately diagonal [91]. The states are instead coupled through off-diagonal elements in the potential energy matrix, known as diabatic couplings. The diabatic couplings are smooth functions of the nuclear coordinates. The diagonal elements are also smooth, maintaining their orbital character and switching energy ordering through a CI (Fig. 1(b)). In many applications, diabatic states are preferred over adiabatic ones because the singular coupling is removed. Here we prefer them because diabatic energies are easier to fit than adiabatic energies. That is, even if one is only interested in adiabatic energies and not NACVs, it is easier to learn the diabatic energies and then diagonalize \mathbf{H}_d than to learn the adiabatic energies directly. This is because the diabatic states possess no CI cusps or avoided crossings (Fig. 1), and are thus more easily fit by interpolating functions such as neural networks. As discussed below, diabatic states improve the model accuracy for species outside the training set.

While many diabatization methods exist, the most common ones cannot be straightforwardly applied to the current problem. Property-based methods were developed for charge-transfer type problems [92], while orbital-based approaches [93] are not implemented for TDDFT. Diabatic models that are parameterized by electronic structure data [61, 83, 94] are not designed for systems undergoing large distortions. Approaches that solve for the adiabatic-to-diabatic transformation matrix [95] are difficult to implement in practice, because the matrix varies rapidly near a CI.

Recent work introduced neural network diabatization based on reference geometries [60, 96]. In this procedure one assumes that the diabatic Hamiltonian is diagonal at a set of known reference geometries. At such geometries the elements of \mathbf{H}_d are equal to the adiabatic energies, but possibly reordered. For example, for two states and two reference geometries, one would have $\mathbf{H}_d = \text{diag}(E_0, E_1)$ at the first geometry and $\mathbf{H}_d = \text{diag}(E_1, E_0)$ at

the second (Fig. 1(b)). A neural network is then trained to produce \mathbf{H}_d , such that its eigenvalues are always equal to the adiabatic energies, and its elements are as above at the reference geometries. These two constraints generate \mathbf{H}_d at all intermediate geometries.

This method was successfully applied to thiophenol dissociation, yielding results consistent with the fourfold way [60]. However, as shown by the NACV error in Table S1, this method alone does not generate true diabatic states for azobenzene. To understand why, consider that near a CI, the true diabatic coupling d_{01} must closely resemble the coupling shown in Fig. 1(b) (d_{nm} is shorthand for $(\mathbf{H}_d)_{nm}$). This is because d_{01} must be linear in displacements about a CI, and quadratic only for Renner–Teller type intersections [71]. However, for two diabatic states, only the square of the diabatic coupling $|d_{01}|^2$ enters into the expression for the adiabatic energies. The model might then generate an off-diagonal element similar to $|d_{01}|$. This error would incur no penalty on the network, because the diagonal elements would properly switch ordering and the adiabatic energies would be correct. Hence the model would not generate smooth diabatic states.

To remedy this issue we used the combined loss described in the Methods section. The NACV component of this loss was also used in Ref. [62]. We note that the number of diabatic states and choice of orderings in Eq. (5) may depend on the system. For azobenzene we were interested in fitting $M_{\text{ad}} = 2$ adiabatic states, and in general one should use $M_d \geq M_{\text{ad}}$ diabatic states. In this work chose $M_d = 3$ because it significantly improved on the results of $M_d = 2$. In fact, with $M_d = 2$, the R^2 correlation of the force NACV was negative. The orderings in Eq. (5) were chosen by taking a small sample of azobenzene configurations, training several small models with different diabatic orderings, and picking the one with the best results. Thus no system knowledge is required to choose the orderings. The only system knowledge required is the set of reference geometries¹, but the reference loss can be omitted with negligible impact on model performance.

The diabatic model leads to true CIs, where the ground and excited state energies are exactly equal. To see why, consider Fig. 3(b), which shows the two diagonal diabatic elements (red and blue) and the off-diagonal coupling (light gray). The off-diagonal element passes linearly through zero in the \vec{h} direction. The diagonal elements cross in the \vec{g} direction, and so $\Delta \equiv d_{11} - d_{00}$ passes linearly through zero. Therefore, one can begin at a geometry for which $\Delta = 0$, and move in the \vec{h} direction until $d_{01} = 0$ without changing Δ . The final geometry will therefore be a CI. The fact that both d_{01} and Δ are locally linear, and hence must pass through zero, is known theoretically [71] and properly predicted by the model.

C. Network loss

The neural network loss terms are defined as:

$$\begin{aligned} \mathcal{L}_{\text{core}} = & \sum_n \rho_{E_n} \cdot \text{mse}(E_n) + \rho_{f_n} \cdot \text{mse}(\vec{f}_n) \\ & + \sum_{n>m} \rho_{\Delta E_{nm}} \cdot \text{mse}(\Delta E_{nm}) \end{aligned} \quad (\text{S1})$$

$$\mathcal{L}_{\text{ref}} = \sum_n \rho_{\text{ref}} \cdot \text{mse}(d_{nn}(\vec{R})) \Big|_{\vec{R} \in \{\text{cis}, \text{trans}\}} \quad (\text{S2})$$

$$\mathcal{L}_{\text{nacv}} = \sum_{n>m} \rho_{\text{nacv}} \cdot \text{mse}(A_n(\vec{R})A_m(\vec{R})\vec{g}_{nm}(\vec{R})), \quad (\text{S3})$$

¹ Typically the reference geometries are reactants and products. For reactions in which the product is not known *a priori*, one might use the following approach. First, train an adiabatic network on a single molecule. We have found that adiabatic models match or outperform diabatic ones for single species. Then use the model to simulate dynamics. By using active learning to improve the model’s coverage of configurational

space, the simulations can eventually discover the product. If derivatives of the molecule are expected to have similar reactions, then their products are also now known. With knowledge of the reactants and products, and hence the reference geometries, one can now build a diabatic model. If this is not possible, then one can train the model without \mathcal{L}_{ref} . Table S1 shows that this will likely decrease the model performance near CIs by a small amount.

where $\Delta E_{nm} = E_n - E_m$ is the energy gap. Each loss function is a sum over mean squared error (MSE) terms scaled by different weights ρ . For scalar molecular quantities X the loss is given by $\text{MSE}(X) = \sum_{j=1}^M (X_j - \hat{X}_j)^2 / M$, where \hat{X} is the predicted quantity and the sum is over M geometries in a batch. For atomic vectorial quantities the mean is additionally over the $3L_j$ vector components, where L_j is the number of atoms in the j^{th} geometry.

The loss term $\mathcal{L}_{\text{core}}$ contains the usual energy and force losses, plus an additional term for gap errors. While the gap term was not used in previous work, we found it to be crucial for the systems studied here. Indeed, without this loss term, one would expect the gap MAE in adiabatic models to be the geometric sum of the energy MAEs. This is what we found when using $\mathcal{L}_{\text{core}}$ without the gap loss. Table S1 shows that, when using the full core loss, the gap MAE is actually lower than each individual energy MAE.

The reference loss, \mathcal{L}_{ref} , is a sum over geometries \vec{R} which are considered to be *cis* or *trans*. At these geometries the target d_{nn} are given by Eq. (5). A geometry is considered to be *cis* or *trans* if its central CNNC atoms deviate from those of the equilibrium structure by $< 0.15 \text{ \AA}$. The distance is computed as the root-mean-square deviation (RMSD) among the atoms after alignment.

The NACV loss imposes diabaticity. It involves the force NACV \vec{g}_{nm} , and a phase correction $A_n(\vec{R}) = \pm 1$. The phase correction is chosen separately for each geometry to minimize the error between the predicted and target force coupling. This factor is necessary because each adiabatic wavefunction can have an arbitrary sign [97]. The signs cancel for diagonal terms like E_n and \vec{f}_n , but not for off-diagonal terms like \vec{g}_{nm} . The A_n account for these arbitrary sign changes.

D. Data generation

To train a useful model one must generate reliable QC data. TDDFT [98] typically offers a good compromise between speed and accuracy for modeling excited states. However, it has known instabilities near CIs [66], and, as a result of the Brillouin theorem, produces the wrong branching space dimensionality for S_0/S_1 intersections [99]. These issues, which can be traced back to the single-reference description of the excitation, can be partially alleviated with SF-TDDFT [63]. In this approach the excitation is performed with respect to a high-spin reference state. This yields some transitions that have double-excitation character with respect to the singlet ground state. Here we used SF-TDDFT with the BHHLYP functional [65] and the 6-31G* basis [64]. SF-TDDFT/BHHLYP is well-benchmarked for CIs in a number of molecules [67, 100]. Because SF is not spin complete, the singlet states must be identified based on their square spins $\langle S^2 \rangle$. We used the approach of Ref. [72], which identifies singlets as the two states with the lowest $\langle S^2 \rangle$ from the three excitations of lowest energy. Calculations were performed with the Q-Chem package [84].

Near-CI regions of the combinatorial species were sampled by setting the central CNNC dihedral angle of relaxed geometries to ± 90 degrees and/or the CNN/NNC angles to 180 degrees. The other internal coordinates were not changed. The former corresponds to the rotation pathway and the latter to inversion. This led to 11 possible combinations of rotation and inversion, including pure rotation, pure inversion, concerted inversion, and inversion-assisted rotation [77].

II. Model accuracy and ablation studies

Here we compare the DANN model to a model trained without a NACV loss (“ $-\mathcal{L}_{\text{nacv}}$ ”), a model trained without a reference geometry loss (“ $-\mathcal{L}_{\text{ref}}$ ”), an adiabatic model, and a median predictor. The latter predicts the median ground- and excited-state energy for each species, and the median value among all species for other quantities. The MAEs for unseen species are compared in Table S1. Half the geometries were sampled randomly and half by proximity to a CI, as described in Supplementary Sec. IX. Of particular interest are $(\Delta E)_{\text{small}}$, the MAE of the gap when it is under 0.2 eV, and $\text{sgn}(\Delta E)_{\text{small}}$, the average overestimation of small gaps by the model. DANN actually *underestimates* the small gaps, indicating that it does not “miss” conical intersections. Using an adiabatic model or removing the NACV loss increases both the error and overestimation of $(\Delta E)_{\text{small}}$.

	E_0	E_1	ΔE_{01}	$(\Delta E)_{\text{small}}$	$\text{sgn}(\Delta E)_{\text{small}}$	\vec{F}_0	\vec{F}_1	\vec{g}_{01}
DANN	3.06	3.77	1.89	0.97	-0.29	1.72	2.31	1.36
$-\mathcal{L}_{\text{nacv}}$	2.31	2.49	1.65	1.22	0.52	1.72	2.39	2.21
$-\mathcal{L}_{\text{ref}}$	2.21	2.58	1.68	1.06	0.16	1.72	2.32	1.37
adiabat	2.11	3.24	2.43	1.24	0.60	1.69	2.31	—
adiabat (st. 1)	3.28	2.99	1.88	1.49	0.95	1.67	2.30	—
median	16.86	19.04	22.80	36.04	-36.04	17.27	17.18	2.05

Table S1. Performance of different ablated models. Each column, apart from $\text{sgn}(\Delta E)_{\text{small}}$, shows the MAE of a different quantity. $\text{sgn}(\Delta E)_{\text{small}}$ is the mean signed error of $(\Delta E)_{\text{small}}$, given by $\text{mean}\{(E)_{\text{small}}^{\text{pred}} - (\Delta E)_{\text{small}}^{\text{target}}\}$. “st. 1” indicates the first stage of training of the adiabatic model, as described in Supplementary Sec. IX. The best score in each category is shown in bold.

	E_0	E_1	ΔE_{01}	$(\Delta E_{01})_{\text{small}}$	\vec{F}_0	\vec{F}_1	\vec{g}_{01}
MAE (\downarrow)	0.75	0.96	0.72	0.40	0.94	1.13	0.87
R^2 (\uparrow)	1.00	0.99	0.99	0.97	1.00	0.99	0.88

Table S2. Test set MAE for the DANN model trained on all species.

The same is true of removing \mathcal{L}_{ref} , but the effect is smaller. While the changes in $(\Delta E)_{\text{small}}$ appear minor, we show in Supplementary Sec. VII that they lead to noticeable quantum yield differences for a number of species.

Also of note is the difference in NACV error among the different models. Removing $\mathcal{L}_{\text{nacv}}$ substantially increases the NACV error, even when \mathcal{L}_{ref} is used. Since $\mathcal{L}_{\text{nacv}}$ imposes diabaticity, this shows that the reference loss alone does not give accurate diabatic states. We also see that the adiabatic model is much worse at predicting ΔE away from the CI than the diabatic model. This is because of the large loss weight used for $(\Delta E)_{\text{small}}$ in the second stage of adiabatic training (see Supplementary Sec. IX). Indeed, the results after the first stage are much better for ΔE , but far worse for $(\Delta E)_{\text{small}}$. Finally, the models without reference or NACV losses have far lower energy errors than DANN. It may be possible to decrease DANN’s energy errors by increasing the weight of the energy loss.

A key benefit of the diabatic model is that it accurately predicts the gap near CIs. It also produces the NACVs, adiabatic energies, and forces all with one model. Without the diabatic model, one might learn the force coupling as a separate property. This has been done in other work by taking the gradient of a learnable scalar [42, 44]. However, one would still have to divide by the adiabatic gap to obtain the derivative coupling. Thus gap prediction errors from an adiabatic model would still be problematic. A separate option would be to compute the NACV through the Hessian of the gap [44]. This is quite slow, and would have similar problems with the gap.

Table S2 shows the test set statistics for the DANN model trained on all species. This model was used for virtual screening, as described in the main text. Here the test set simply consists of 5,000 held-out geometries. Unlike in Table S1, the test set contains mostly seen species, and the geometries were not generated from DANN-NAMD with the trained model. The results are quite similar to the “seen species” rows in Table I.

III. Experimental data

We performed an extensive literature search to find quantum yields of azobenzene derivatives. To best compare to the vacuum conditions simulated here, we chose quantum yields measured in non-polar solvents. When multiple results were available in different non-polar solvents, we computed the prediction error relative to the mean experimental value. To compare to the simulated S_1 dynamics, we selected yields measured at the

peak of the highest-wavelength absorption band, which is a dipole-forbidden $n - \pi^*$ transition in unsubstituted azobenzene. All sources used here specifically reported yields at wavelengths close to the $n - \pi^*$ and $\pi - \pi^*$ absorption peaks, so it was straightforward to choose the appropriate wavelength. Lastly, we only selected yields at or around room temperature, since yields can have a strong temperature dependence [101].

IV. Sources of error

A. Experimental error

There are two main sources of uncertainty in the experimental quantum yield results. The first is the error in the absorption coefficients ε of the two isomers. The absorption strengths are used in standard rate equations to compute the yield; see, for example, Refs. [102, 103] and the supplementary of Refs. [104, 105]. The error can be traced back to the overlapping absorption spectra of the two isomers, which must be disentangled. The magnitude of the error depends on the method used to compute ε , which itself is related to the year of publication. For example, Ref. [102], published in 1988, indicated that the absorption coefficients were the primary source of error, but did not give an estimate for their uncertainty. They isolated the different isomers using chromatography and computed their absorption coefficients separately. When the isomers could not be isolated in sufficient amounts, the authors used the time-dependent absorption approach of Ref. [106]. This method had relative errors of $\pm 25\%$ for $n - \pi^*$ quantum yields [106] (i.e., $\text{yield} \rightarrow \text{yield} \times (1 \pm 0.25)$). Ref. [103], published in 2004, reported a yield of 0.7 to 1.0 for compound **20**, a large range stemming from the overlap of the isomers' spectra and the method used to separate them [107]. More modern methods have lower errors; for example, Ref. [105], published in 2015, reported a relative error of only $\pm 10\%$ in the supplementary.

The second important source of error is the overlap of the $n - \pi^*$ and $\pi - \pi^*$ absorption bands. The more strongly the two bands overlap, the higher the error of using only the S_1 state in the dynamics. Moreover, many experiments do not irradiate at the precise maximum of the $n - \pi^*$ band. Often a representative $n - \pi^*$ or $\pi - \pi^*$ wavelength is chosen, and then used for each derivative [101, 104], even though the derivatives have different absorption peaks.

To get an idea of the range of errors this could introduce, consider the results of Ref. [108] from 1958. This work computed the quantum yield of unsubstituted azobenzene over several wavelengths. The *trans* to *cis* yield was measured as 21% at 405 nm and 27% at 436 nm (10^{-3} M solution). The difference is because 405 nm light excites more of the $\pi - \pi^*$ band (313 nm) than 436 nm light. 405 nm excitation leads to a lower yield because the $\pi - \pi^*$ yield is only 11%. The difference is even more pronounced for the *cis* quantum yield, which drops from 55% to 40% moving from 436 nm to 546 nm. This result does not have a simple explanation, as there is no lower energy band beyond 436 nm. In principle, the error due to overlapping bands could be mitigated with an explicit dipole-electric field coupling term [109], which would excite multiple states in different proportions. For derivatives with completely overlapping $n - \pi^*$ and $\pi - \pi^*$ bands, the S_1 approximation would incur a maximum error of the S_2 yield minus the S_1 yield. This is because the S_2 state is much brighter than S_1 , and would therefore dominate the experimental yield. The yield difference is about 10 percentage points for both *cis* and *trans* unsubstituted azobenzene [110, 111]. Therefore, in the worst case scenario, we would expect an error of about 10 percentage points from the S_1 approximation. The error would likely be closer to the 6% reported in Ref. [108] at wavelengths of moderate overlap.

We can further quantify these errors by computing the range of results from different studies. Measurements of the $t \rightarrow c$ azobenzene yield range from 20% to 28% between 1979 and 1987 [106, 110, 112], though this could also be related to solvent effects (Sec. IV B). The yields of compounds **9** and **10** range from 16% to 24% and 44% to 50%, respectively, with references from 1962 and 1988 [101, 102]. Aggregating the results of all three compounds gives an average yield range of 7.3 percentage points, and an average relative error of $\pm 14\%$.

A final source of uncertainty is specific to Refs. [73, 74]. Several species were reported to have zero yield over a large range of irradiation wavelengths. However, as noted in the footnote to Table S12, the dipole-allowed S_2 transition was highly redshifted and thus overlapped strongly with the S_1 transition. The absorbance changes after irradiation were quite small, indicating a near-zero yield, but could have been due to the S_1 transition. Hence it is possible that the S_1 yield is not precisely zero.

B. Computational error

There are several sources of computational error. The first is error in the PES, which can be decomposed into error from SF-TDDFT and error from the model. While it is intractable to perform SF-TDDFT for all species with experimental data, our results for unsubstituted azobenzene suggest that it is rather accurate. As reported in the main text, we computed yields of $60 \pm 4\%$ and $26 \pm 3\%$ for $c \rightarrow t$ and $t \rightarrow c$, respectively. Experimental measurements between 1979 to 1987 gave $t \rightarrow c$ yields between 20% and 28% [106, 110–112]. Experimental measurements in 1974 and 1979 gave $c \rightarrow t$ yields of 55% and 56% [110, 111], while a measurement in 1958 gave $48\% \pm 5\%$ [108]. All measurements were performed in non-polar solvents. The yields computed with the original model were 59% and 37% for $c \rightarrow t$ and $t \rightarrow c$, respectively. Hence the main source of error seems to be the model, rather than SF-TDDFT.

Solvent effects may also affect the yield. These effects can be decomposed into a systematic term and a non-systematic term. It has been argued that non-polar solvents systematically reduce the quantum yield relative to vacuum [113]. However, careful analysis of the experimental data reported in [114] and cited in [113] does not support this conclusion. Indeed, given the good agreement between non-polar experimental results and both SF-TDDFT and hh-TDA DFT [7], systematic effects of non-polar solvents are likely to be small. Solvent-specific effects are part of the range of reported experimental values, since different works often use different non-polar solvents. Typical ranges were discussed in Sec. IV A.

A final source of error is the approximations in surface hopping. A large body of literature has examined surface hopping’s accuracy; see, for example, Refs. [115–118]. In this work, we can roughly evaluate its performance by comparing its results to those of *ab initio* multiple spawning (AIMS) [7]. AIMS is a fully quantum mechanical method and can thus be treated as a benchmark. We computed the $t \rightarrow c$ yield as $26 \pm 3\%$ using surface hopping with SF-TDDFT, and Ref. [7] computed the yield as $24\% \pm 6\%$ using AIMS with hh-TDA DFT. While there may be some error cancellation from the different quantum chemistry methods, the good agreement is still encouraging.

V. Influence of different functional groups

Here we discuss how the model transferability depends on the functional groups in a compound. To answer this question, we analyzed each of the 40 unseen species in the test set, and computed the model error for each geometry. As described in the main text, the geometries were taken from DANN-NAMD with the trained model; half were selected by proximity to a CI and half selected randomly. For each functional group, we aggregated the errors of all geometries that contained the group, and then computed the mean.

The functional groups with the lowest gap errors are shown in Fig. S1(a), and those with the highest are shown in panel (b). Table S3 shows the number of times that each functional group appears in the training set, together with the errors for all properties. Of the groups with the lowest errors, we see that both **B** and **C** are well-represented in the training set, which may explain their accurate results. The other groups each have about 1,000 samples in the training set. This is smaller than that of **B** and **C**, but not negligible. The compounds are also rather simple, which may partly explain their low error.

Of the groups with the highest gap error, only the nitro group **G** is well-represented in the training set. There are almost 10,000 training geometries with nitro groups, yet the error is still quite high. This may be because NO_2 is a strong electron-withdrawing group, and can thus have a significant effect on the electronic structure of azobenzene. Of the remaining substituents, both **H** and **J** are rather complicated, and each has only about 200 samples in the training set. Groups **F** and **I** both have about 1,000 samples in the training set, and are thus moderately represented. **I** may have large errors because of the electron-withdrawing effects of the three fluorines, or because it is found together with **J** in compound **20**. **F**, the *tert*-butyl group, likely has large errors simply because it is bulky and thus leads to distorted geometries.

We see that the transferability of the model depends on how well-represented a functional group is in the training set, and how complicated its electronic effects are. The former is supported by Fig. S2, which shows that the functional group error is anti-correlated with its prevalence in the training set. However, the effect is

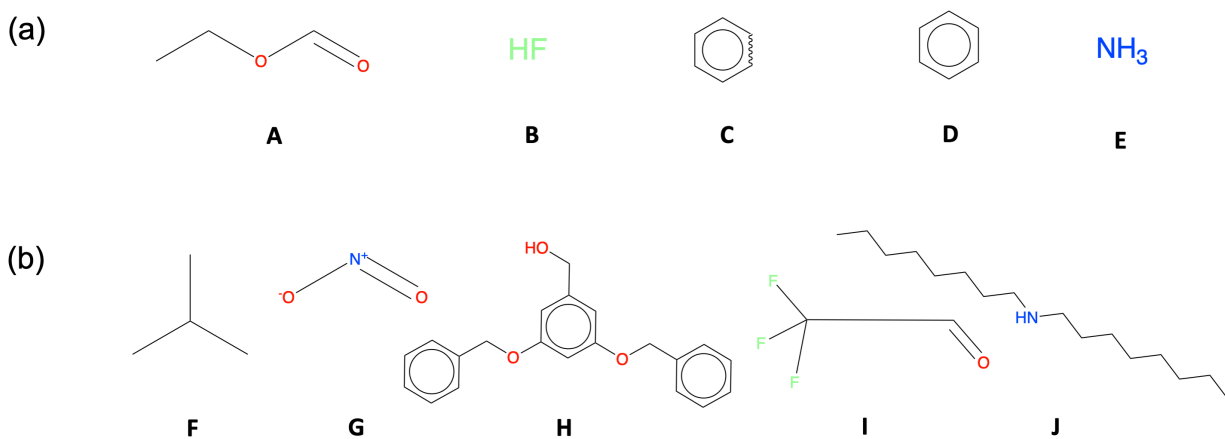


Figure S1. Functional groups in the test set of unseen species. (a) The five functional groups with the lowest gap error. Error increases from left to right. **C** is fused to a benzene ring in azobenzene. (b) The five groups with the highest gap error. Error decreases from left to right.

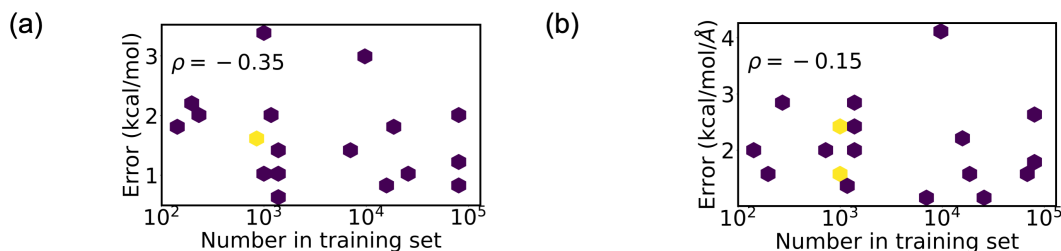


Figure S2. Test set error by functional group, plotted against the number of times the group appears in the training set. The Spearman rank correlation ρ is also shown. (a) Gap error. (b) Excited state force error.

stronger for the gap than the excited state forces ($\rho = -0.35$ and $\rho = -0.15$, respectively), and neither fully explains the error. The remaining error can best be explained by analyzing the groups themselves.

VI. Spin contamination

Here we discuss the amount of spin contamination in the training set. Figure S3 shows the square spin, $\langle S^2 \rangle$, for both the ground and excited states. The spin contamination is quite low for the ground state, with $\langle S^2 \rangle$ under 1.0 for 96% of geometries. It is much higher for the excited state, with 76%, 82%, and 93% of geometries having $\langle S^2 \rangle$ under 1.0, 1.2, and 1.4, respectively. Most of the geometries with excited-state $\langle S^2 \rangle \gtrsim 1.2$ were near the non-reactive S_1/S_2 CI, characterized by near-planarity and CNN angles near 108° [7]. The S_0/S_1 CIs, on the other hand, did not have severe spin contamination. The spin-contamination near the S_1/S_2 CI led to difficulty in fine-tuning the model for unsubstituted azobenzene. We also visually inspected a sample of geometries with extreme spin contamination, $\langle S^2 \rangle > 1.8$, and found that nearly all had broken apart. We kept

	Num. train	E_0	E_1	ΔE_{01}	$(\Delta E)_{\text{small}}$	$\text{sgn}(\Delta E)_{\text{small}}$	\vec{F}_0	\vec{F}_1	\vec{g}_{01}
A	1,219	1.05	0.94	0.63	0.61	-0.04	1.24	1.32	1.19
B	71,940	0.95	1.02	0.81	0.52	-0.01	1.33	1.66	1.62
C	17,624	1.08	1.14	0.85	0.37	0.12	1.24	1.52	1.62
D	994	0.73	1.00	0.91	0.44	-0.07	1.18	1.70	1.54
E	1,154	1.09	1.21	0.95	0.72	-0.04	1.61	2.42	1.41
F	979	2.28	3.36	3.39	1.14	-0.31	1.53	1.53	1.96
G	9,578	0.94	3.49	2.90	0.68	-0.16	1.59	4.11	1.83
H	200	1.62	3.17	2.18	0.44	-0.13	0.99	1.67	0.68
I	1,225	1.33	2.90	2.10	0.78	-0.17	1.63	2.90	0.91
J	236	1.33	2.90	2.10	0.78	-0.17	1.63	2.90	0.91

Table S3. Test set error by functional group. The groups are divided into those with the lowest gap errors (top) and those with the highest (bottom). “Num. train” refers to the number of geometries in the training set that contain the functional group. While each of the functional groups is contained in the training set, the precise combinations and arrangements of groups (and hence the molecular graphs) are unique to the test set.

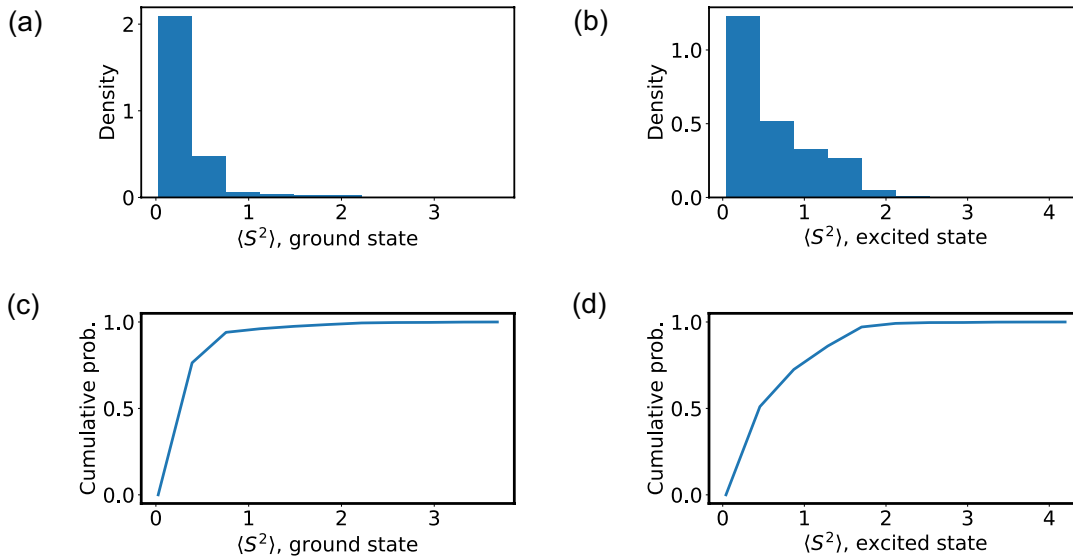


Figure S3. Analysis of spin contamination in the training set. (a) Distribution of $\langle S^2 \rangle$ in the ground state. (b) As in (a), but for the excited state. (c) Cumulative probability as a function of $\langle S^2 \rangle$ in the ground state. (d) As in (c), but for the excited state.

these geometries in the training set so that the model could learn the high energy of bond breaking.

VII. Surface hopping results with different models

Figure S4 compares ZN hopping statistics of adiabatic and diabatic models for unseen species. Panel (a) shows the distribution of hopping percentages among species. The hopping percentage is defined as the percent

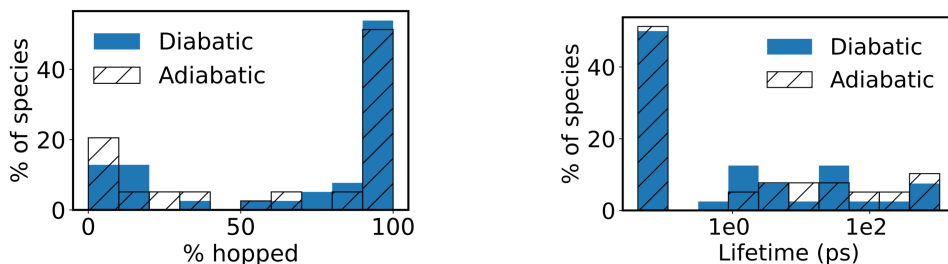


Figure S4. Comparison of ZN hopping statistics in diabatic and adiabatic models. (a) Proportion of trajectories in each species that hopped to the ground state. (b) Excited state lifetime for each species.

of trajectories for a given species that end in the ground state. The y -axis is the percent of all species that correspond to each bin. Panel (b) has an analogous plot, but with the hopping percentage replaced by the lifetime. The lifetime was estimated from an exponential fit of S_1 population, $p = \exp(-(t - t_{\text{on}})/\tau) \Theta(t - t_{\text{on}}) + \Theta(t_{\text{on}} - t)$. Here $p \in [0, 1]$ is the S_1 population, t is the time, t_{on} is the fitted turn-on time, τ is the fitted lifetime, and Θ is the Heaviside step function. Trajectories that did not contain any hops were assigned the maximum lifetime of all the other trajectories.

The *cis* derivatives have lifetimes around 50 fs, consistent with computational results for *cis* azobenzene [72]. Nearly 100% of all *cis* trajectories ended in the ground state. The *trans* derivatives have a wide distribution of lifetimes. Some are between 1 and 2 ps, which is similar to *trans* azobenzene [72]. Others are in the range of tens to hundreds of ps, which reflects the high proportion of trajectories that never returned to the ground state. These are very likely incorrect. They may be because of barriers between S_1 minima and S_0/S_1 CIs, which are known to exist for *trans* azobenzene [7]. Relatively small errors in the barrier may lead to large over-estimations of the lifetime.

Excited state barriers make it even more important to have accurate PESs near CIs. Since trajectories spend little time near crossing regions, trapping becomes even more severe when the gap near CIs is overestimated. The diabatic model helps to address this problem: in each plot we can see that the diabatic model leads to more hopping. For example, using the adiabatic model, 21% of species have hopping percentages under 10%. This number is reduced to 13% for the diabatic model.

The diabatic model also improves the quantum yield. Figure S5 shows the ZN quantum yields with the diabatic model, and Fig. S7 shows the diabatic FS yields. The Spearman rank correlation for the *trans* species is fairly high with both the ZN and FS methods, with the model accurately predicting low yields for a number of species. Figure S6 shows the ZN yield with the adiabatic model. The correlation among all species is reasonable, but for *trans* species is rather low. This is because of three molecules with zero predicted yield, all of which became trapped in the excited state. The diabatic model only predicts zero yield for one of these. Further, the diabatic model properly predicts zero yield for species **35**, while the adiabatic model does not. For these reasons the diabatic model has a fairly high correlation with experiment. Still, it is clear that excited state trapping of *trans* species is not a fully solved problem. Preferential sampling of excited state barriers may help to further address this issue in the future.

VIII. Architecture

All models were implemented in PyTorch [119]. The model hyperparameters are given in Table S4, and an in-depth explanation of the PaiNN parameters can be found in Ref. [35]. Note that we used five convolutions instead of the three used originally, as this substantially improved model performance. Following DimeNet

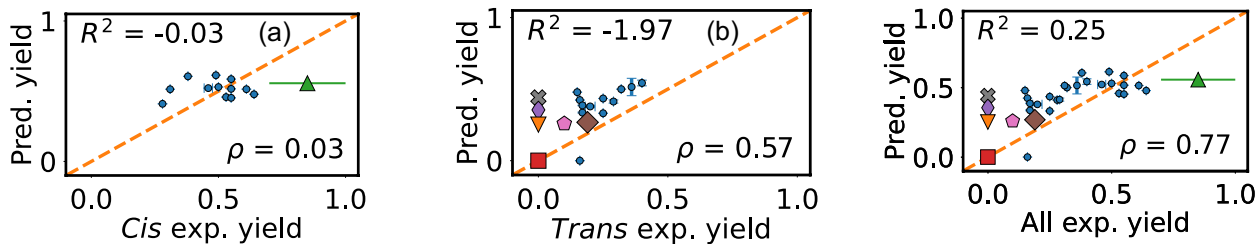


Figure S5. Experimental vs. predicted ZN yields using the diabatic model. (a) *Cis* isomers. (b) *Trans* isomers. (c) All species.

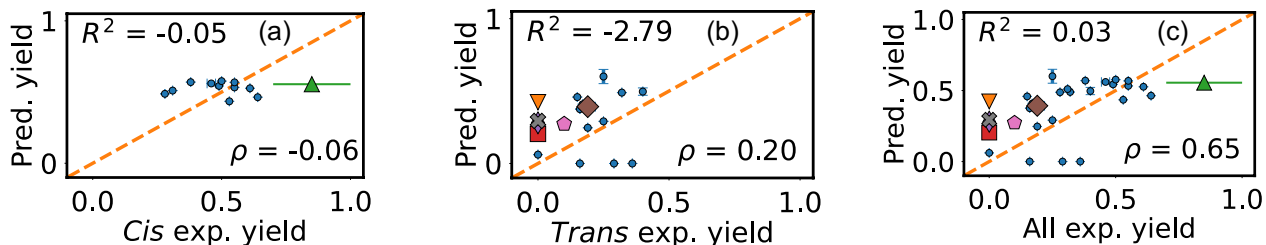


Figure S6. Experimental vs. predicted ZN yields using the adiabatic model. (a) *Cis* isomers. (b) *Trans* isomers. (c) All species.

[120], we allowed the k values in the radial basis functions to be updated during training. For the adiabatic model, we predicted each property as a sum over per-convolution properties, which was also used in DimeNet. In particular, each convolution had a readout network to convert the atomic features to an output. The final property was obtained by summing each of these outputs.

Several variations on the architecture were tested. For example, we trained both two- and three-state diabatic models on 5,000 azobenzene configurations. We found that adding a third diabatic state significantly decreased the error for all properties. We then trained models using all possible three-state reference orderings and chose the best one. We also experimented with intensive pooling for off-diagonal energies (see Supplementary Sec. XVI). In particular, we generated a molecular fingerprint through an attention-weighted average of atomic fingerprints. We then mapped this fingerprint to the d_{nm} for $n \neq m$. Even though the d_{nm} are intensive for $n \neq m$, this approach did not improve model predictions.

We also experimented with several adiabatic models. The model in the main text predicted E_0 and E_1 directly. This approach was used in previous work for single-molecule non-adiabatic dynamics [39–45] and for the prediction of absorption spectra across chemical space [121]. We also examined three models that predicted E_0 directly and E_1 as the sum of E_0 and a learned gap. We tried three different pooling methods for the learned gap ΔE : taking the mean over atomwise gaps, taking an attention-weighted average over atomwise gaps, and applying a dense readout network to an attention-weighted sum of atomic fingerprints. All approaches performed quite poorly compared to learning E_0 and E_1 separately as summed atomwise energies. This is problematic, because only adiabatic models that predict E_1 as $E_0 + \Delta E$ can guarantee the positivity of the gap (e.g. by squaring ΔE or applying a softplus function). Hence the adiabatic model in the main text sometimes generated predictions with $E_1 < E_0$. This is impossible in the diabatic model by construction.

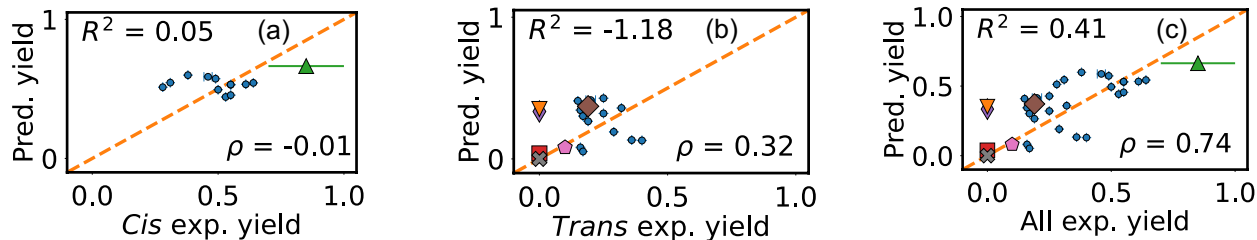


Figure S7. Experimental vs. predicted FS yields using the diabatic model. (a) *Cis* isomers. (b) *Trans* isomers. (c) All species.

Hyperparameter	Meaning	Value or name
F	dimension of hidden atomic features	128
n_{conv}	number of convolutions	5
n_{RBF}	number of radial basis functions (RBF)	20
R_{cut}	cutoff distance for convolutions	5.0 Å
activation function	activation used in message-passing and readout	Swish
learnable k	whether k parameters in RBF are learnable	true
skip	output is sum of per-convolution outputs	true only for adiabatic
M_d	number of diabatic states	3
M_{ad}	number of ground truth adiabatic states	2

Table S4. Model hyperparameters

IX. Training

The training set contained 562,037 geometries from 8,197 species. 5,000 geometries from 308 species were used for validation. The remaining 74,322 geometries from 40 species were held out for testing, so that the predicted yields of unseen species could be compared with experiment. A different random seed was used to determine the training and validation splits for each committee model, and also to initialize the different models. After training, we ran FS DANN-NAMD on 40 holdout species using the trained diabatic model. For each species we selected 50 geometries randomly and 50 by CI proximity (Eq. (S13)), for a total of 4,000 geometries. These geometries were used as the test set, giving the “unseen” statistics in Table I. The DANN-NAMD geometries from the diabatic network were used for *all* models, including the adiabatic and ablated ones in Table S1. The “seen” statistics were generated using the validation set. For all statistics, a phase correction was applied to \bar{g}_{01} to minimize the prediction error, as in Eq. (S3).

The model for screening new azobenzene derivatives was trained on all species. Inclusion of the holdout species provided an additional 74,322 geometries. We used 631,367 geometries for training, 5,000 for validation, and 5,000 for testing. This corresponded to 8,215 training species, 332 validation species, and 303 test species.

Training was performed over energies and forces/force couplings in units of kcal/mol and kcal/mol/Å, respectively. Per-species reference energies were subtracted from each energy. These were obtained by summing atomic reference energies, computed using multi-variable linear regression from (atom type, count) to relaxed geometry energy. Configurations with 10- σ energy and force outliers were removed prior to training. Those with forces or energies ≥ 450 kcal/mol/(Å) from the mean were also removed. We found that more stringent removal of outliers led to less stable trajectories. For example, removing 3- σ outliers and maximal energies/forces of

ρ_{E_0}	ρ_{E_1}	$\rho_{\Delta E_{01}}$	ρ_{f_0}	ρ_{f_1}	ρ_{ref}	ρ_{nacv}
0.2	0.1	0.5	1	1	0.01	1

Table S5. Loss parameters used to train all diabatic models.

	lr	lr _{min}	$\rho_{\Delta E_{01}}$	ρ_{small}
Stage 1	10^{-4}	10^{-5}	0.5	0
Stage 2	10^{-5}	10^{-6}	1.0	100

Table S6. Variable loss parameters used for training the adiabatic model. lr is the starting learning rate and lr_{min} is the minimum learning rate, at which point training is stopped.

≥ 300 kcal/mol(\AA) led to unstable ground state trajectories for six of the 40 unseen species. The 10- σ and 450 kcal/mol(\AA) criteria led to no diverging ground state trajectories. A small proportion of excited state trajectories were still unstable. We discarded all excited state trajectories that produced NaN geometries or energies, which were at most a few percent of the trajectories for a given species.

Models were trained with the Adam algorithm using a batch size of 60. Geometries were sampled for each batch using Eq. (S12), as described below. The loss was given by Eq. (6), using parameters in Table S5. Note that the range of NACVs is approximately 10 times smaller than the range of forces. This means that $\rho_{\text{nacv}} = 1$, $\rho_{f_i} = 1$ gives much higher weight to the forces. We experimented with a range of NACV coefficients between 0.1 and 10, and found that $\rho_{\text{nacv}} = 1$ gave the best performance.

The learning rate was initialized to 10^{-4} and reduced by a factor of two if the validation loss had not improved in 10 epochs. The final model was selected as the one with the lowest validation loss. Training was performed on a single 32 GB Nvidia Volta V100 GPU, and took 13 days to complete.

For diabatic models, the training was stopped once the learning rate reached 10^{-5} . Adiabatic models were trained in a two-step process, using different loss functions in each stage. Each step ended once the learning rate fell below a certain value. The following loss function was used with different loss trade-offs at different stages:

$$\begin{aligned} \mathcal{L} &= \mathcal{L}_{\text{core}} + \mathcal{L}_{\text{small}}, \\ \mathcal{L}_{\text{small}} &= \sum_{n>m} \rho_{\text{small}} \cdot \text{mse}(\Delta E_{nm}^{\text{small}}), \end{aligned} \tag{S4}$$

where $\mathcal{L}_{\text{small}}$ penalizes errors in gaps under 0.2 eV. The parameters for each stage are given in Table S6, and the ρ_E and ρ_f coefficients are the same as in Table S5. The first stage emphasized energy gaps and gradients, while the second stage emphasized small gaps to fine-tune the model near conical intersections. We also experimented with scheduled training and using $\mathcal{L}_{\text{small}}$ for diabatic models, but did not find any improvements.

X. Balanced sampling

A custom data sampler was used during training because the dataset was imbalanced in the following ways. First, the combinatorial species only had a few geometries each, and so would be very rarely sampled during training. Second, there were more equilibrium *trans* geometries than *cis* geometries. Third, there were more equilibrium geometries in general than near-CI configurations. Our sampler addressed these imbalances by giving higher sampling probability to underrepresented species and/or configurations.

To explain the sampling procedure, let us define two types of sampling weights. Sampling weights that are balanced by species are denoted by w , and those that are not are denoted by v . For example, the weights $w_{\text{cluster}}(g_{i,A})$ and $v_{\text{cluster}}(g_{i,A})$ are both assigned to the i^{th} geometry in species A , denoted $g_{i,A}$, based on the

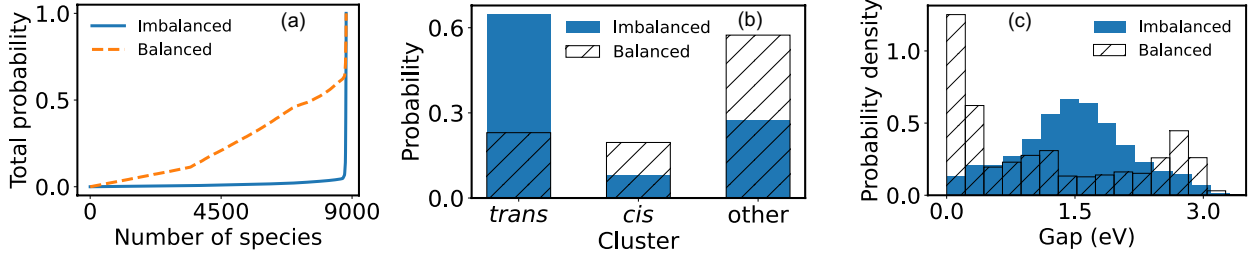


Figure S8. Sampling probabilities for different geometries in the training set, with and without the custom sampler. (a) Cumulative probability $p(i)$ of sampling any species $\leq i$. The species are ordered from fewest to most geometries. (b) Probability of sampling a geometry in each of the three clusters. (c) Probability density of sampling a geometry with a given energy gap.

cluster of configurations that it belongs to. This cluster is denoted k , where $k \in (\textit{cis}, \textit{trans}, \textit{other})$. For a geometry $g_{i,A}$ in cluster k , the weights are given as follows:

$$\begin{aligned} w_{\text{cluster}}(g_{i,A}) &= \left(\frac{1}{n_{\text{spec}} \cdot n_{\text{clusters}}} \right) \left(\frac{1}{n_{k,A}} \right) \\ v_{\text{cluster}}(g_{i,A}) &= \left(\frac{1}{n_{\text{geoms}} \cdot n_{\text{clusters}}} \right) \left(\frac{n_A}{n_{k,A}} \right) \end{aligned} \quad (\text{S5})$$

Here $n_{k,A}$ is the number of geometries of species A in cluster k , and n_A is the number of geometries in species A . n_{spec} is the total number of species, n_{clusters} is the total number of clusters, and n_{geoms} is the total number of geometries. The above definitions ensure that there is an equal probability P_k of sampling any cluster k :

$$\text{Using } w : P_k = \sum_A \sum_{g_{i,A} \in k} w_{\text{cluster}}(g_{i,A}) = \sum_A \frac{1}{n_{\text{spec}} \cdot n_{\text{clusters}}} = \frac{1}{n_{\text{clusters}}} \quad (\text{S6})$$

$$\text{Using } v : P_k = \sum_A \sum_{g_{i,A} \in k} v_{\text{cluster}}(g_{i,A}) = \sum_A \frac{n_A}{n_{\text{geoms}} \cdot n_{\text{clusters}}} = \frac{1}{n_{\text{clusters}}} \quad (\text{S7})$$

Here we have used the fact that $\sum_{g_i \in k,A} 1 = n_{k,A}$, $\sum_A 1 = n_{\text{spec}}$, and $\sum_A n_A = n_{\text{geoms}}$. The difference between the two weights is that w leads to balanced sampling among species, while v does not:

$$\text{Using } w : P_A = \sum_k \sum_{g_{i,A} \in k} w_{\text{cluster}}(g_{i,A}) = \sum_k \frac{1}{n_{\text{spec}} \cdot n_{\text{clusters}}} = \frac{1}{n_{\text{spec}}} \quad (\text{S8})$$

$$\text{Using } v : P_A = \sum_k \sum_{g_{i,A} \in k} v_{\text{cluster}}(g_{i,A}) = \sum_k \frac{n_A}{n_{\text{geoms}} \cdot n_{\text{clusters}}} = \frac{n_A}{n_{\text{geoms}}} \quad (\text{S9})$$

Here we have used the fact that $\sum_k 1 = n_{\text{cluster}}$. We see that w gives equal weights for each species, whereas v gives weights proportional to the number of geometries in that species.

Similarly to the cluster weights, we define Zhu-Nakamura weights w_{ZN} and v_{ZN} , such that geometries with a

higher hopping probability p_{ZN} are sampled more often. The corresponding expressions are:

$$w_{\text{ZN}}(g_{i,A}) = \left(\frac{1}{n_{\text{spec}} \cdot \sum_j p_{\text{ZN}}(g_{j,A})} \right) p_{\text{ZN}}(g_{i,A}) \quad (\text{S10})$$

$$v_{\text{ZN}}(g_{i,A}) = \left(\frac{1}{\sum_{j,B} p_{\text{ZN}}(g_{j,B})} \right) p_{\text{ZN}}(g_{i,A}). \quad (\text{S11})$$

The overall sampling weight for a given geometry is determined by each w and v , together with two user-defined weights: $P_{\text{spec}} \in [0, 1]$, the importance of species balance, and $P_{\text{ZN}} \in [0, 1]$, the importance of sampling geometries with high hopping rates. $P_{\text{cluster}} = 1 - P_{\text{ZN}}$ is the importance of sampling different clusters in a balanced way. The final sampling weight is then given by

$$p(g_{i,A}) = P_{\text{ZN}} [P_{\text{spec}} \cdot w_{\text{ZN}}(g_{i,A}) + (1 - P_{\text{spec}}) \cdot v_{\text{ZN}}(g_{i,A})] + P_{\text{cluster}} [P_{\text{spec}} \cdot w_{\text{cluster}}(g_{i,A}) + (1 - P_{\text{spec}}) \cdot v_{\text{cluster}}(g_{i,A})]. \quad (\text{S12})$$

In applying Eq. (S12) we used $P_{\text{spec}} = 0.6$ and $P_{\text{cluster}} = 0.5$. We defined a geometry as *cis* or *trans* if its RMSD with respect to the corresponding reference structure was ≤ 0.25 Å, and “other” otherwise. In contrast to diabatic reference geometries, the RMSD was computed using all atoms, not just the central CNNC atoms. Note that while our clustering approach was specific to *cis-trans* isomerization, many black-box clustering methods exist, such as hierarchical clustering [81]. Any of these methods could have been used in place of our user-defined clusters.

In principle p_{ZN} should depend on trajectory-specific factors such as velocity, and is therefore not a function of geometry alone. However, during simulations, hops almost always occurred below 0.5 eV. Further, since p_{ZN} is approximately exponential in the square of the gap, we approximated it with

$$p_{\text{ZN}}(\Delta E) \approx \exp(-\Delta E^2 / (2\Delta E_0^2)), \quad (\text{S13})$$

where ΔE is the energy gap and $\Delta E_0 = 0.15$ eV. This choice of ΔE_0 gave $p(\Delta E = 0.5 \text{ eV}) \approx 3 \times 10^{-3}$. This meant that geometries were usually selected only if their gap was under 0.5 eV.

The sampling probabilities for the training set are shown in Fig. S8. Panel (a) shows the cumulative probability of sampling different species, with the compounds ordered from fewest to most geometries. Without the custom sampler the probability was quite small for most molecules. This was because the majority were combinatorially generated and had few geometries. 90% of the probability was contained in the last 41 species. This probability was reduced to 35% when using the balanced sampler. Panel (b) shows that most geometries were *trans* isomers and few were *cis*. The custom sampler gave approximately equal probability to the two isomers, and the highest probability to the *other* group. Most geometries in this group had small gaps, and hence were highly weighted by w_{ZN} and v_{ZN} . This is demonstrated in panel (c), which shows that near-CI geometries had the highest probabilities when using the balanced sampler.

XI. Dynamics

Simulations were performed in two stages. First the starting geometries and velocities were generated for each NAMD trajectory. For *ab initio* simulations we optimized the ground state geometry of each species and computed its normal modes. These modes were used to sample geometries and velocities from the harmonic oscillator Wigner distribution at temperature $T = 300$ K [122]. For neural trajectories we ran classical MD for 15 ps with the Nosé-Hoover thermostat [85, 86], and selected random samples to start the NAMD trajectories. Classical MD was implemented with ASE [123]. The effective mass Q was set to $(3N - 6) \cdot \tau^2 k_{\text{B}} T$, where N is the number of atoms, k_{B} is Boltzmann’s constant, $T = 300$ K, and $\tau = 25$ fs is the relaxation time. Initial atomic velocities were sampled from a Maxwell-Boltzmann distribution at 300 K. Optimized geometries were used for the initial atomic positions. The total linear and angular momenta of the system were set to zero, and

the target kinetic energy was set to $(3N - 6) \cdot k_B T / 2$. The time step of the simulation was set to 0.5 fs, and the neighbor list of the system was updated every 10 time steps (5 fs). All pairs of atoms within with 7 Å were considered neighbors. In each step the distance between all pairs of neighbors was computed, and only pairs within 5 Å of each other were used in the model. This procedure meant that atoms entering the 5 Å cutoff between neighbor list updates would not be missed (i.e., an extra 2 Å “skin” was added). All ground state and excited state trajectories in this work were propagated with the velocity Verlet algorithm [124]. Frames for NAMD were taken only after the first 1 ps of ground state MD to allow for equilibration.

Next we ran ZN dynamics [72, 79, 125] using N_{trj} trajectories, where N_{trj} was 10 during the active learning cycle and 500 during final inference. Hops were restricted to gaps ≤ 0.5 eV to avoid unphysical transitions, the time step was set to 0.5 fs [72, 79, 125], and the neighbor list was again updated every 10 steps. Excited state trajectories were run for 1.5 ps during active learning and 5 ps for final inference. All other details of the implementation can be found in Refs. [72, 79, 125]. For each species we performed inference in parallel over 100-250 geometries at a time, depending on the number of atoms in the molecule. This was done by batching together geometries from 100-250 trajectories and evaluating the model on the batch. The batch size depended on the size of the molecule, which determined the GPU memory consumption. FS simulations were also performed for final inference. The ZN parameters above were used for FS DANN-NAMD of molecules in the holdout set. For virtual screening we used 100 trajectories and 2 ps of excited state dynamics. Species with a hopping rate under 10%, or with molecular graphs that changed during ground state dynamics, were deemed unreliable. These molecules were excluded from Fig. 4, but not from the average hopping percentage reported in the main text.

The quantum yield was calculated as $Y = n_R / n_T$, where n_R is the number of reactive trajectories and n_T is the total number of trajectories. The uncertainty was computed as the standard deviation of 1,000 bootstrapped samples. To identify reactivity, we computed the RMSD between the CNNC atoms in the final frame and the corresponding atoms in the optimized *cis/trans* geometries. A trajectory was considered reactive if it started as *cis* and ended closer to *trans*, or vice-versa. Trajectories that ended in the excited state were excluded from n_R and n_T . The yield was reported as 0 ± 0 if all trajectories ended in the excited state.

XII. Fewest switches implementation

Tully’s surface hopping propagates the nuclei on one electronic surface at a time, called the active surface. The expansion coefficients of the electronic wavefunction are also propagated, and are used to make stochastic changes in the active surface. The coefficient vector \mathbf{c} , expressed in a basis denoted by “rep”, evolves as [126]

$$\frac{d}{dt} \mathbf{c}^{\text{rep}} = -\frac{i}{\hbar} [\mathbf{H}^{\text{rep}} - i\mathbf{T}^{\text{rep}}] \mathbf{c}^{\text{rep}}. \quad (\text{S14})$$

The Hamiltonian matrix is $\mathbf{H}^{\text{rep}} = \langle \psi_n^{\text{rep}} | H | \psi_m^{\text{rep}} \rangle$, and the coupling term is

$$(\mathbf{T}^{\text{rep}})_{nm} = \left\langle \psi_n^{\text{rep}}(\vec{R}(t)) \left| \frac{\partial}{\partial t} \right| \psi_m^{\text{rep}}(\vec{R}(t)) \right\rangle = \vec{v} \cdot \vec{k}_{nm}^{\text{rep}}, \quad (\text{S15})$$

where \vec{v} is the classical velocity of the nuclei. In principle the nuclei can be propagated on a PES in any electronic basis, and hops can be performed between states in that basis. For example, nuclei could be propagated on diabatic PESs and hops could occur between diabatic states. However, it is well-established that surface hopping in the adiabatic basis gives the best results [97]. Hence the nuclei should be propagated in the adiabatic basis, and hops should be decided using \mathbf{c}^{ad} .

While the adiabatic basis should be used for hopping, \mathbf{c} can be *propagated* in the diabatic basis, and then transformed into the adiabatic basis for all subsequent manipulations (e.g. deciding on hops, adding decoherence, etc.). Indeed, using Eq. (S14) in the adiabatic basis can require small time steps for so-called trivial crossings, where derivative NACVs are quite narrow and large [127]. Since global diabatic states are almost never available in *ab initio* simulations, a local diabatization method is often used to propagate \mathbf{c} [88, 127]. This method is

very stable and can be used with a fairly large time step. In this work we have a *global* diabatic basis that can replace local diabatization. Following Refs. [126, 127], we then propagate \mathbf{c} as follows:

$$\mathbf{c}^{\text{ad}}(t + \Delta t) = \mathbf{P}^{\text{ad}}(t, t + \Delta t) \mathbf{c}^{\text{ad}}(t), \quad (\text{S16})$$

$$\mathbf{P}^{\text{ad}}(t, t + \Delta t) = \mathbf{U}^\dagger(t + \Delta t) \mathbf{P}^{\text{d}}(t, t + \Delta t) \mathbf{U}(t) \quad (\text{S17})$$

$$\mathbf{P}^{\text{d}}(t, t + \Delta t) = \prod_{k=1}^K \exp[-i\mathbf{H}_k^{\text{d}}\Delta t/K] \quad (\text{S18})$$

$$\mathbf{H}_k^{\text{d}} = \mathbf{H}^{\text{d}}(t) + \frac{k}{K} [\mathbf{H}^{\text{d}}(t + \Delta t) - \mathbf{H}^{\text{d}}(t)]. \quad (\text{S19})$$

Here $\mathbf{P}^{\text{ad}}(t, t + \Delta t)$ is the propagator of \mathbf{c}^{ad} from t to $t + \Delta t$ (Eq. (S16)), and Δt is the timestep. The adiabatic propagator is a transformation of the diabatic propagator into the adiabatic basis (Eq. (S17)), using the transformation matrix \mathbf{U} (Eq. (2)). The diabatic propagator is a matrix product of K sub-propagators (Eq. (S18)), where K is the number of electronic substeps for each nuclear step. “exp” denotes a matrix exponential, not element-wise exponentiation. The k^{th} sub-propagator uses a linear interpolation between $\mathbf{H}^{\text{d}}(t)$ and $\mathbf{H}^{\text{d}}(t + \Delta t)$ to approximate $\mathbf{H}_d(t + k\Delta t/K)$ (Eq. (S19)). The diabatic Hamiltonian is produced by the neural network model. The probability of hopping from state n to m is [126]

$$p_{n \rightarrow m} = \left(1 - \frac{|c_n^{\text{ad}}(t + \Delta t)|^2}{|c_n^{\text{ad}}(t)|^2}\right) \frac{\text{Re} \left[c_m^{\text{ad}}(t + \Delta t) (P_{mn}^{\text{ad}})^* (c_n^{\text{ad}}(t))^* \right]}{|c_n^{\text{ad}}(t)|^2 - \text{Re} \left[c_n^{\text{ad}}(t + \Delta t) (P_{nn}^{\text{ad}})^* (c_n^{\text{ad}}(t))^* \right]}, \quad (\text{S20})$$

where $\text{Re}(x)$ is the real part of x . Note that while our approach follows that of SHARC [126], we have written our own code and made it publicly available [128]. Our repository also contains code for surface-hopping with the ZN method.

In this work we initialized $\mathbf{c}^{\text{ad}}(t = 0) = [0, 1, 0]$, used Eqs. (S16)-(S19) to generate $\mathbf{c}^{\text{ad}}(t + \Delta t)$, and computed the hopping probability with Eq. (S20). \mathbf{c}^{ad} was expressed in a three-state basis because our model used three diabatic states, and hence \mathbf{U} had three dimensions. Since the model was only trained on the first and second adiabatic energies, we set $p = 0$ for all transitions to the third state. Hops to state m were performed if [126]

$$\sum_{i=1}^{m-1} p_{n \rightarrow i} < r \leq \sum_{i=1}^{m-1} p_{n \rightarrow i} + p_{n \rightarrow m}, \quad (\text{S21})$$

where $0 \leq r \leq 1$ is a random number. The momentum was rescaled after each hop to conserve the total energy. The momentum was multiplied by a constant factor, rather than being re-scaled in the direction of the NACV, to avoid the overhead of a NACV calculation (see below) [88]. If the factor was complex—a so-called “frustrated” hop—then no transition was made [88].

We also tested propagation in the adiabatic basis. In this case we used the NACV to evaluate \mathbf{T} , and constructed P^{ad} from the interpolation of $\mathbf{H}^{\text{ad}} - i\mathbf{T}$ [126]. Both Eq. (S20) and Tully’s original hopping expression [57] were used. The momentum was re-scaled in the direction of the NACV [115]. In all cases the results were very similar to those of the diabatic basis. Diabatic propagation was ultimately chosen because of its reported stability [127], and because of its computational efficiency. In particular, the diabatic propagation requires only diabatic energies and one adiabatic gradient. It therefore uses only one forward and one backward pass through the neural network. By contrast, constructing the NACVs requires gradients for each diabatic element (Eq. (1)). For three diabatic states, this corresponds to one forward pass and six backward passes. While shared convolution layers and caching mean that N gradients do not take $N \times$ as long as one gradient, we found that they still added significant time. Hence diabatic propagation was chosen as the most efficient method.

The decoherence correction of Ref. [129] was used to counter the over-coherence of Tully’s original method. A sign correction was also used to remove random sign changes in the eigenvectors of \mathbf{H}_d . The sign correction A_m for the m^{th} eigenvector \mathbf{v}_m was computed as

$$A_m = \text{sgn}(S_{n'm}) \quad (\text{S22})$$

$$S_{nm} = \mathbf{v}_n^\dagger(t) \mathbf{v}_m(t + \Delta t) = [\mathbf{U}^\dagger(t) \mathbf{U}(t + \Delta t)]_{nm} \quad (\text{S23})$$

$$n' = \text{argmax}_n \{|S_{nm}|\}. \quad (\text{S24})$$

The transformation matrix and NACVs were corrected through

$$U_{nm}(t + \Delta t) \rightarrow A_m U_{nm}(t + \Delta t) \quad (\text{S25})$$

$$\vec{k}_{nm}(t + \Delta t) \rightarrow A_n A_m \vec{k}_{nm}(t + \Delta t) \quad (\text{S26})$$

$$\vec{g}_{nm}(t + \Delta t) \rightarrow A_n A_m \vec{g}_{nm}(t + \Delta t). \quad (\text{S27})$$

Note also that $S_{nm} \approx \langle \psi_n^{\text{ad}}(t) | \psi_m^{\text{ad}}(t + \Delta t) \rangle$:

$$\begin{aligned} \langle \psi_n^{\text{ad}}(t) | \psi_m^{\text{ad}}(t + \Delta t) \rangle &= \sum_{ij} U_{in}^*(t) U_{jm}(t + \Delta t) \langle \psi_i^{\text{d}}(t) | \psi_j^{\text{d}}(t + \Delta t) \rangle \\ &\approx \sum_{ij} U_{ni}^\dagger(t) U_{jm}(t + \Delta t) \langle \psi_i^{\text{d}}(t) | [(1 - (\vec{v}\Delta t) \cdot \nabla_R) | \psi_j^{\text{d}}(t)] \rangle \\ &= \sum_{ij} U_{ni}^\dagger(t) U_{jm}(t + \Delta t) \langle \psi_n^{\text{d}}(t) | \psi_m^{\text{d}}(t) \rangle \\ &= [\mathbf{U}^\dagger(t) \mathbf{U}(t + \Delta t)]_{nm} = S_{nm}. \end{aligned} \quad (\text{S28})$$

Here we used the fact that there is no derivative coupling between diabatic states, and that the diabatic states at a given time are orthonormal (Eq. (S35)).

All simulations were performed with a time step of 0.5 fs and $K = 25$ substeps for electronic propagation. Unlike in ZN dynamics, hops were not restricted to gaps under 0.5 eV. We found that restricting hops improved the ZN results but hurt the FS results. For example, with no maximum gap, most neural ZN trajectories starting from *trans*-azobenzene hopped at ~ 1.5 eV. This did not match the results of Ref. [72], which used *ab initio* ZN with the same level of DFT theory. After adding the restriction, most hops occurred at gaps under 0.1 eV, in agreement with Ref. [72]. By contrast, most FS hops occurred under 0.1 eV even without a maximum restriction, though some still occurred at large gaps. The lifetime and yield without a maximum also better matched previous calculations.

XIII. Active learning

New geometries were selected for QC calculations using two different criteria. In the first three active learning loops, new geometries were chosen based on the prediction variance of two neural networks. Our aim was to select geometries with fairly uncertain predictions. We did not want geometries with extremely high uncertainty, as these usually corresponded to broken graphs that were outside the target learning space for the model. We therefore used a log-normal target distribution for the uncertainty. Target uncertainties were randomly sampled from this distribution, and geometries with the closest variance to the targets were selected. The log-normal probability of obtaining a sample x is given by $P(x) = \exp[-(\ln(x/s) - \mu)^2 / 2\sigma^2] / (x\sigma\sqrt{2\pi})$, where s , σ and μ are positive numbers. Under this distribution, both completely certain and completely uncertain predictions have zero probability, with a peak for predictions with medium uncertainty. The decay of the probability at large uncertainty is slow, such that highly uncertain geometries can still be selected with reasonable probability. We set $\mu = 0$, $\sigma = 1$, and $s = 7$, giving a target distribution with a mode of 2.5 kcal/mol(\AA) and a mean of 11.5 kcal/mol(\AA). Committee variances were computed for both forces and energies for each electronic state, and the largest variance was compared with the target variance from the distribution. For each species we selected 16 geometries from ground state MD and 33 from surface hopping. Only geometries from the 164 literature species were chosen. This led to approximately 8,000 new data points in each active learning cycle. Other, simpler methods of geometry selection are certainly possible; for example, random selection is also known to work quite well [37, 38]. Our approach successfully increased the model quality throughout active learning, but we did not thoroughly compare it to other methods. Such a comparison may be of interest in the future.

In the next two loops we used only azobenzene, with the goal of densely sampling the CI region. Half the geometries were selected randomly from the excited state dynamics, and half were selected based on the gap.

Because the model overestimated the gap in uncertain regions, we did not want to simply select geometries with low predicted gap. Rather, in each trajectory we identified avoided crossing geometries, and assigned equal sampling probability to all geometries before and after that crossing, up to a maximum predicted gap of 1.7 eV. An avoided crossing geometry was defined as having a gap that was lower than in the previous and subsequent time step. In this way we aimed to identify regions that could have small gaps, even if the model did not predict them to be so.

XIV. Validation

A. *Ab initio* NAMD

Ab initio simulations were used for unsubstituted *cis* and *trans* azobenzene. Starting configurations and velocities were generated with *ab initio* MD using Q-Chem to drive the dynamics. Ten MD simulations were performed for each species. Each simulation was initiated with a different geometry produced by ground-state neural network MD. The simulations were run for 10 ps each using a time step of 0.5 fs, with the the BP86 functional [130, 131] and 6-31G* basis [64]. Initial velocities were sampled from a thermal distribution at 300 K, with rotation and translation projected out. The nuclei were propagated with the Velocity Verlet algorithm. A Nosé-Hoover chain of length 5, timescale $\tau = 25$ fs, and temperature $T = 300$ K was used as a thermostat. The Fock extrapolation order was set to six, and twelve Fock matrices from previous steps were used in the extrapolation.

167 and 178 excited state trajectories were generated for *trans* and *cis* azobenzene, respectively. Each trajectory was initialized with a different set of coordinates and velocities, randomly sampled from the ten ground state simulations after 1 ps of equilibration. Spin-flip TDDFT was used with the BHHLYP functional [65] and 6-31G* basis. DIIS with geometric direct minimization (GDM) [132] was used to improve SCF convergence. We found this to be critical: with the usual DIIS algorithm, the SCF cycle failed to converge for 36% of trajectories. This usually occurred within the first 200 fs.

The dynamics were run with in-house scripts, which can be found at <https://github.com/learningmatter-mit/NeuralForceField>. We propagated the elements of the electronic wavefunction in the adiabatic basis, and corrected the sign of the force NACV by minimizing its change from the previous step. The momentum was re-scaled in the direction of the NACV after a hop [115]. All other parameters were unchanged from the DANN-NAMD simulations. *Cis* trajectories were propagated for 200 fs, and *trans* trajectories for a maximum of 1.5 ps. For the former, we extended trajectories that had hopped within the last 50 fs, or not hopped at all, until at least 50 fs had passed since they hopped. For the latter, we stopped a trajectory if it had lasted at least 500 fs and had hopped more than 300 fs earlier.

B. Transfer learning

To validate the yield results for substituted compounds, we performed DANN-NAMD for the top candidates using a set of highly accurate models. Each model was fine-tuned for a single species only, which allowed it to achieve high accuracy on that one molecule. This transfer learning strategy was used in place of *ab initio* NAMD because the latter would be prohibitively slow for all but the smallest molecules. For example, consider molecule **169**, which has only 54 atoms. We found that a single gradient or NACV calculation for this species took approximately 50 minutes with 8 CPU cores. Since *trans* derivatives must be simulated for at least 1 ps, and since the time step must be no larger than 0.5 fs, we would need to perform 2,000 QC calculations for each trajectory. Assuming parallel calculation of the NACVs and gradients at each step, an *ab initio* simulation would take 70 days. By contrast, the fine-tuning approach allows us to perform highly accurate simulations for tens to hundreds of species.

Each new model was refined from the original network using QC data from a single species. The initial training geometries were sampled from DANN-NAMD simulations with the original DANN model. A committee of three

lr	lr _{min}	patience	factor	batch size	ρ_{E_0}	ρ_{E_1}	$\rho_{\Delta E_{01}}$	ρ_{f_0}	ρ_{f_1}	ρ_{ref}	ρ_{nacv}
10^{-4}	10^{-5}	50	0.5	1	0.3	0.3	0.5	1	1	0	1

Table S7. Training parameters used for transfer learning. ‘‘Patience’’ refers to the number of epochs without an improvement in validation loss before the learning rate is reduced. ‘‘Factor’’ is the amount by which the learning rate is lowered.

DANN models, each trained on the entire dataset, was used to select the initial geometries (see below). Three fine-tuned models were then trained. Each was initialized with a different random seed and trained with different data splits. 90% of the data was used for training and 10% for validation. The first model was subsequently used for DANN-NAMD. Geometries were selected from these simulations using the newly-trained committee, and each geometry received QC calculations. This data was added to the training set, each committee model was re-trained, and the cycle was repeated as in Fig. 2(b). After each cycle we evaluated the model accuracy using the geometries generated by DANN-NAMD.

Initially we selected only 50 geometries in each round of active learning. Once we narrowed our focus to two species, we sampled 500 geometries in each round. The geometries were sampled according to the following four strategies. One third was sampled randomly. One third was chosen by the prediction variance in the excited state forces. Those with the highest variance were selected. We used only the uncertainty in the forces, and not the energies, because the former is a better indicator of trajectory instability [81]. One sixth was chosen by proximity to a CI (Eq. (S13) with $\Delta E = 0.2$ eV). The final sixth was chosen to sample excited-state barriers. In particular, we sampled geometries with probability

$$p \propto \exp[(E - E_{\text{min}})/(k_{\text{B}}T)]. \quad (\text{S29})$$

Here E is the excited state energy, E_{min} is the minimum excited state energy in the trajectory, and $T = 300$ K. We only sampled configurations from 50 fs or later in the simulations. This was done to avoid the initial high-energy geometries encountered before relaxation. Equation (S29) is inversely proportional to the room-temperature Boltzmann probability, and thus assigns the highest weight to the highest-energy configurations.

Transfer learning has previously been used to account for solvent effects [37] and to reach higher levels of QC theory [37, 133] with only a small portion of the training set. Typically the majority of the network weights are frozen during re-training. This leaves modifiable parameters in only the final few layers. However, we found that the best results were obtained without any parameter freezing. We therefore allowed all parameters to be modified during re-training. Further, we found it best to start with the normal learning rate rather than a reduced one. We also did not use a reference loss, as its effect on the original DANN results was minor, and possibly even harmful. The full set of training parameters can be found in Table S7. All trajectory parameters were unchanged from the screening phase. For final inference we performed DANN-NAMD for 5 ps with 2,000 trajectories, using 500 ps of ground state MD.

The accuracy of the fine-tuned networks is shown in Tables S8 and S9. Statistics are shown for 500 geometries that were sampled from DANN-NAMD for each molecule using the final trained networks. The model errors are well below 1 kcal/mol(\AA) for geometries sampled by all methods other than uncertainty. The error on the uncertainty-selected configurations is under 2 kcal/mol(\AA) in all cases. Since these geometries are specifically chosen to have the highest error, and since their errors are still rather small, we can be confident in the accuracy of the models. We note, however, that the average error for the uncertain geometries depends on how many trajectories are run. Running more trajectories means sampling more configurations, and hence finding more geometries with high uncertainty. As mentioned above, we ran 100 trajectories for 5 ps each for both screening and transfer learning (2,000 trajectories were used for the final predictions). We saved frames every 15 fs, leading to 33,333 geometries in total. Since we picked the 167 most uncertain frames, our sample is roughly the same as choosing the top 0.5% of all geometries with the highest error.

One reason that we did not use transfer learning for the unsubstituted compounds was severe spin contamination. This is an artifact of unrestricted SF-TDDFT, and is a well-documented problem in NAMD for *trans* azobenzene [72]. After the first round of active learning, we found that the force error for the unsubstituted

Sampled by	Metric	E_0	E_1	ΔE_{01}	\vec{F}_0	\vec{F}_1	\vec{g}_{01}
ZN	MAE (\downarrow)	0.48	0.42	0.44	0.83	0.82	0.52
	R^2 (\uparrow)	1.00	1.00	0.96	0.99	0.99	0.86
Barrier	MAE (\downarrow)	0.62	0.93	0.86	0.67	0.82	0.66
	R^2 (\uparrow)	0.99	0.99	0.99	1.00	0.99	0.88
Random	MAE (\downarrow)	0.74	0.80	0.86	0.71	0.77	0.89
	R^2 (\uparrow)	0.99	0.99	0.98	1.00	1.00	0.88
Uncertainty	MAE (\downarrow)	0.74	1.26	1.09	0.83	1.36	1.04
	R^2 (\uparrow)	0.99	0.99	0.99	0.99	0.97	0.73

Table S8. Test set statistics for the final TL model of species **165**. Results are divided by the method used to select samples. “ZN” uses Zhu-Nakamura gap-based sampling [Eq. (S13)] and “barrier” uses Eq. (S29). 1,244 geometries were used for fine-tuning.

Sampled by	Metric	E_0	E_1	ΔE_{01}	\vec{F}_0	\vec{F}_1	\vec{g}_{01}
ZN	MAE (\downarrow)	0.67	0.54	0.39	0.87	0.84	0.64
	R^2 (\uparrow)	0.98	0.98	0.99	0.99	0.99	0.82
Barrier	MAE (\downarrow)	0.89	0.74	0.41	0.61	0.57	0.51
	R^2 (\uparrow)	0.99	0.95	1.00	1.00	1.00	0.98
Random	MAE (\downarrow)	0.58	0.66	0.51	0.67	0.70	0.78
	R^2 (\uparrow)	0.99	1.00	1.00	1.00	1.00	0.96
Uncertainty	MAE (\downarrow)	0.72	1.66	1.90	0.83	1.53	1.01
	R^2 (\uparrow)	0.99	0.93	0.97	0.99	0.95	0.92

Table S9. As in Table S8, but for species **169**. 2,445 geometries were used for fine-tuning.

models increased to 8 kcal/mol/Å. This error did not drop significantly with more data, even though the new geometries were fairly close to the Franck-Condon region. This issue did not occur with any of the derivatives. We found that the error was correlated with $\langle S^2 \rangle$, and reasoned that our method of singlet selection (Sec. IID) was likely choosing triplet excited states. This highlights the issues inherent in SF-TDDFT, and reinforces the need for low-cost, spin-complete alternatives [28–30].

Lastly, we note that both the *ab initio* and transfer-learned *cis* quantum yields were noticeably higher than in other SF-TDDFT studies [72]. This is likely because we used MD to initiate the trajectories, rather than normal-mode or Wigner sampling based on the harmonic approximation. Our experiments showed that normal-mode sampling led to decreased *cis* yields, closer to those reported in Ref. [72]. Unlike the *trans* isomer, *cis* azobenzene is somewhat flexible, with significant torsions occurring during ground-state MD. This indicates that the harmonic approximation should be avoided when possible. Indeed, using MD ground state sampling together with FS dynamics for *cis* azobenzene, we obtained a yield of $60 \pm 4\%$; this is in excellent agreement with experimental values in non-polar solution, which are close to 55% on average [73]. It is in much better agreement than the value of 34% obtained with FS surface hopping in Ref. [72]. Since we used the same electronic structure method and the same surface hopping approach, we can be confident that the difference is mainly due to MD sampling.

C. Conical intersection pathways

We labeled the CI pathway of each *trans* trajectory according to its last hopping geometry. Each geometry was compared to two reference planar CIs and two reference rotational CIs. The trajectory was labeled by the reference CI that was the closest to the hopping geometry. That is, if a trajectory hopped at a geometry closer to one of the rotational CIs, it was labeled a rotational trajectory, and similarly for the planar CI.

The reference CIs were chosen from the set of all hopping geometries in the trajectories. The two rotational CIs were those with dihedral angles closest to 90° and 270° , respectively, and $\max(\alpha_{\text{CNN}}, \alpha_{\text{NNC}})$ closest to 138° . The two planar CIs were those with dihedral angles closest to 186° and 174° , respectively, and both $\max(\alpha_{\text{CNN}}$ and $\alpha_{\text{NNC}})$ closest to 148° . The angles are those of the optimized CI geometries in Ref. [7]. For the derivative **169**, we further optimized each of the reference structures to minimize the gap and hence obtain a true CI. For each trajectory we computed the RMSD between the CNNC atoms of the hopping geometry and the CNNC atoms of the reference CIs.

We additionally enforced that a hopping geometry could only be considered planar if $|180 - \theta| \leq \theta_0$, where θ is the CNNC dihedral angle and θ_0 is a cutoff value. We chose $\theta_0 = 75^\circ$ for azobenzene and 65° for **169**. Without this constraint, many of the hopping geometries labeled “planar” in **169** actually led to isomerization. This made a dihedral constraint necessary. On the other hand, when we only labeled all geometries with $|\theta - 180| \leq 45^\circ$ as planar [7], we found through visual inspection that many non-reactive CI geometries were mislabeled as rotational. Using the minimal distance to a reference CI, together with a modest dihedral constraint, led to the most qualitatively reasonable results. We confirmed that none of the geometries labeled as “planar CI” with our metric led to switching, which further reinforced the soundness of our approach.

XV. Figure details

A. Computational speed-up

Here we describe how the speeds of ML and QC calculations were computed. For QC, we first computed the run time of one gradient calculation on a single geometry, denoted t_{calc} . The node time was then calculated as $t_{\text{node}} = t_{\text{calc}}/n_{\text{calc}}$, with $n_{\text{calc}} = (\text{cores per node})/(\text{cores per job})$. We assumed 40 cores per node. All QC jobs were performed with 8 cores, and so n_{calc} was equal to 5.

For ML we performed a batched calculation on ten copies of each geometry, and computed one gradient. The node time was computed as $t_{\text{node}} = t_{\text{calc}}/n_{\text{calc}}$, with $n_{\text{calc}} = 10 \cdot (\text{total memory per gpu})/(\text{memory used in calculation}) \cdot (\text{GPUs per node})$. We set (total memory per gpu) = 32 GB and (GPUs per node) = 2. We used a script that performed one batched calculation on ten copies of a random geometry, and re-ran it 7,000 times. For each iteration we used the `nvidia-smi` command with the PID of the current job to access the GPU memory. We re-ran the script many times, rather than running many calculations in one script, because `nvidia-smi` does not account for all freed memory until a job is finished. That is, after one calculation is finished, `nvidia-smi` shows that GPU memory is still being occupied by the job. This occurs even after all local variables are deleted. Hence running multiple calculations in one script would yield an overestimated GPU memory.

In the above calculation, we implicitly assumed that multiplying the batch size by x would lead to an x -fold increase in memory. In practice we have found that the memory is increased by less than that. This means that the ML speedup in Fig. 3 is a conservative estimate.

B. Diabatic energies

Since three diabatic states were used in the DANN model, and since all three were coupled near the CI, it would be incorrect to plot only two states in Fig. 3(b). We therefore applied a fixed rotation matrix, $\mathbf{U} \neq \mathbf{U}(\vec{R})$, to generate a new diabatic Hamiltonian with only two coupled states. The new diabatic Hamiltonian was given

by $\mathbf{H}'_d(\vec{R}) = \mathbf{U}^\dagger \mathbf{H}_d(\vec{R}) \mathbf{U}$. Note that any position-independent rotation matrix can be brought outside the gradient in Eq. (1), and so \mathbf{H}'_d is still diabatic. The rotation matrix was chosen to diagonalize \mathbf{H}_d at the CI, so that $\mathbf{U}^\dagger \mathbf{H}_d(\vec{R}_{\text{CI}}) \mathbf{U} = \text{diag}(\{E\})$. This led to $d'_{00} = d'_{11}$ and $d'_{01} = 0$ at the CI. Hence the lowest two rotated states were the most important contributors to E_0 and E_1 , and were therefore used in Fig. 3(b).

XVI. Intensive and extensive quantities

The DANN model predicts each property by summing over atomic contributions, just as in the PaiNN model. This guarantees size-extensivity. However, as explained below, the off-diagonal elements of \mathbf{H}_d should be intensive, in the sense that atoms not involved in the excitation should not contribute to these values. Nonetheless, we found that summation for these terms gave better results than averaging. This was true even when using a learnable weighted average. Atom-wise summation can still generate accurate d_{nm} , because the readout network can simply map unimportant atoms to zero.

To demonstrate extensivity and intensivity, consider two uncoupled subsystems, A and B . The total clamped nucleus Hamiltonian is $H(\vec{r}, \vec{R}) = H^A(\vec{r}, \vec{R}) + H^B(\vec{r}, \vec{R})$. Let the excitation of interest be in subsystem A . In this case the adiabatic states involve only excitations in subsystem A , so that the n^{th} diabatic wave function is written as

$$\psi_{d,n}(\vec{r}; \vec{R}) = \psi_{\text{ad},0}^B(\vec{r}; \vec{R}) \sum_k U_{nk} \psi_{\text{ad},k}^A(\vec{r}; \vec{R}). \quad (\text{S30})$$

The diabatic state is a direct product of the ground state wave function of system B , $\psi_{\text{ad},0}^B(\vec{r}; \vec{R})$, and a rotation of the adiabatic states of system A , $\psi_{\text{ad},k}^A(\vec{r}; \vec{R})$. The matrix elements of \mathbf{H}_d are then given by

$$\begin{aligned} (\mathbf{H}_d)_{nm} &= \langle \psi_{d,n} | H(\vec{r}, \vec{R}) | \psi_{d,m} \rangle \\ &= \sum_{kl} U_{nk}^* U_{ml} \langle \psi_{\text{ad},0}^B \psi_{\text{ad},k}^A | H_A(\vec{r}, \vec{R}) + H_B(\vec{r}, \vec{R}) | \psi_{\text{ad},0}^B \psi_{\text{ad},l}^A \rangle \\ &= \sum_{kl} U_{nk}^* U_{ml} (E_{k,A} + E_B) \delta_{kl} \\ &= \sum_k U_{in}^* U_{mk} E_{k,A} + E_B (\mathbf{U} \mathbf{U}^\dagger)_{mn} \\ &= \mathbf{H}_{d,nm}^A + E_B \delta_{nm}. \end{aligned} \quad (\text{S31})$$

Hence the diagonal elements of \mathbf{H}_d each gain the adiabatic energy of B , while the off-diagonal elements remain the same. To understand this result physically, consider that adding a scalar multiplied by the identity yields eigenvalues that are each shifted by the scalar. This means that the eigenvalues of \mathbf{H}_d are each shifted by E_B . Therefore the excitation energy is unchanged, which is the expected behavior upon adding an uncoupled system. Note also that extensivity is sometimes described as doubling the energy when the system size is doubled. However, this definition is too narrow, as it applies to only one adiabatic energy at a time. For example, if subsystem B is a copy of subsystem A , then $E_B = E_{0,A}$, meaning that $E_0 = E_{0,A} + E_B = 2E_{0,A}$, and so the ground state energy is indeed doubled. However, $E_1 = E_{1,A} + E_B = 2E_{1,A} - \text{gap}$, where $\text{gap} = E_{1,A} - E_{0,A}$. Therefore, the excited state energy is not doubled. A more general definition is that the energy of the second subsystem is added to each adiabatic energy.

This analysis shows that the off-diagonal elements of \mathbf{H}_d are intensive. The adiabatic energy gap is also intensive. Intensive here means that adding a subsystem that does not participate in the excitation does not modify the quantity. Such properties can naturally be modeled with an attention mechanism [134–137] that learns the importance of each atom to the excitation. For example, the excitation energy can be modeled as an attention-weighted sum over atomwise quantities. However, as discussed in Supplementary Sec. IX, this approach led to worse performance than simply predicting extensive energies for each state.

XVII. Proof of diabaticity

Here we prove Eq. (1) in the main text. Using bra-ket notation for the wave functions, we define the diabatic states as a linear combination of adiabatic states,

$$|\psi_{d,n}\rangle = \sum_m |\psi_{ad,m}\rangle V_{nm}^*, \quad (\text{S32})$$

where \mathbf{V} is a unitary matrix. Multiplying each side by $V_{nn'}$ and summing over n yields

$$|\psi_{ad,n}\rangle = \sum_m |\psi_{d,m}\rangle V_{mn}, \quad (\text{S33})$$

where we have renamed the dummy indices, $n \rightarrow m$ and $n' \rightarrow n$. We have also used the fact that $(\mathbf{V}^\dagger \mathbf{V})_{nm} = \delta_{nm}$ for any unitary matrix, where δ_{nm} is the Kronecker delta. Given the connection between diabatic and adiabatic states, we can relate \mathbf{V} to the derivative coupling \vec{k} :

$$\begin{aligned} \vec{k}_{nm} &\equiv \langle \psi_{ad,n} | \nabla_R | \psi_{ad,m} \rangle \\ &= \sum_{ij} V_{jn}^* (\nabla_R V_{im} \langle \psi_{d,j} | \psi_{d,i} \rangle + V_{im} \langle \psi_{d,j} | \nabla_R | \psi_{d,i} \rangle) \\ &= \sum_i V_{in}^* \nabla_R V_{im}. \end{aligned} \quad (\text{S34})$$

We have used the fact that, by definition, the derivative coupling between any two diabatic states is zero. We have also used the orthonormality of the diabatic states:

$$\langle \psi_{d,j} | \psi_{d,i} \rangle = \sum_{nm} V_{jn} V_{im}^* \langle \psi_{ad,n} | \psi_{ad,m} \rangle = \sum_{nm} V_{jn} V_{im}^* \delta_{nm} = (\mathbf{V} \mathbf{V}^\dagger)_{ij} = \delta_{ij}, \quad (\text{S35})$$

which follows from the orthonormality of the adiabatic wave functions. Meanwhile, by construction, the Hamiltonian produced by the model has eigenvalues equal to the adiabatic energies:

$$\begin{aligned} (\mathbf{H}_d)_{nm} &= (\mathbf{U} \text{diag}(\{E\}) \mathbf{U}^\dagger)_{nm} \\ &= \sum_{ij} U_{ni} (E_i \delta_{ij}) U_{mj}^* \\ &= \left(\sum_i U_{ni} \langle \psi_{ad,i} | \right) \hat{H} \left(\sum_j U_{mj}^* | \psi_{ad,j} \rangle \right), \end{aligned} \quad (\text{S36})$$

where \mathbf{U} is the unitary matrix that diagonalizes \mathbf{H}_d , \hat{H} is the Hamiltonian operator, and we have used the relation $\langle \psi_{ad,j} | \hat{H} | \psi_{ad,i} \rangle = E_i \langle \psi_{ad,j} | \psi_{ad,i} \rangle = E_i \delta_{ij}$. Comparing Eqs. (S36) and (S32), we see that \mathbf{H}_d is the representation of \hat{H} in a different electronic basis. In particular, if we choose $\mathbf{U} = \mathbf{V}$, such that \mathbf{U} satisfies Eq. (S34), then \mathbf{H}_d is the representation of \hat{H} in the diabatic basis.

The model could be trained directly with Eq. (S34). However, the equation is numerically ill-posed because \vec{k}_{nm} diverges at conical intersections. It is preferable to work instead with the force coupling, \vec{g}_{nm} . We now show that Eq. (1), which is defined in terms of \vec{g}_{nm} , holds only if (S34) is satisfied. The left-hand side of Eq. (1) can be written as

$$\begin{aligned} &\sum_{ij} U_{in}^* \nabla_R (\mathbf{U} \text{diag}(\{E\}) \mathbf{U}^\dagger)_{ij} U_{jm} \\ &= \sum_{ijk} U_{in}^* [(\nabla_R E_k) U_{ik} U_{jk}^* + E_k (\nabla_R U_{ik}) U_{jk}^* + E_k U_{ik} (\nabla_R U_{jk}^*)] U_{jm}. \end{aligned} \quad (\text{S37})$$

The first term is

$$\sum_k (\mathbf{U}^\dagger \mathbf{U})_{nk} (\nabla_R E_k) (\mathbf{U}^\dagger \mathbf{U})_{km} = (\nabla_R E_n) \delta_{nm}, \quad (\text{S38})$$

where we have used the fact that $(\mathbf{U}^\dagger \mathbf{U})_{nm} = \delta_{nm}$. Substituting Eq. (S34) with $\mathbf{V} = \mathbf{U}$, the second term is

$$\sum_k \vec{k}_{nk} (\mathbf{U}^\dagger \mathbf{U})_{km} E_k = \vec{k}_{nm} E_m. \quad (\text{S39})$$

Performing the same substitution for the third term gives

$$\sum_k \vec{k}_{mk}^* (\mathbf{U}^\dagger \mathbf{U})_{nk} E_k = -\vec{k}_{nm} E_n, \quad (\text{S40})$$

where we have used the anti-Hermitian property $\vec{k}_{mn}^* = -\vec{k}_{nm}$. Adding the three terms gives

$$(\mathbf{U}^\dagger (\nabla_R H_d) \mathbf{U})_{nm} = \begin{cases} \nabla_R E_n, & \text{if } n = m, \\ (E_m - E_n) \vec{k}_{nm}, & \text{if } n \neq m. \end{cases} \quad (\text{S41})$$

Noting that $\vec{f}_n = -\nabla_R E_n$ and $\vec{g}_{nm} = (E_m - E_n) \vec{k}_{nm}$ gives Eq. (1) in the main text. Hence Eq. (1) can only hold if Eq. (S34) is true. Therefore, enforcing Eq. (1) ensures that there is no derivative coupling between any pair of diabatic states. Note that in the main text we used three diabatic states, and trained only on energies and couplings between the first two adiabatic states. This still means that all three diabatic states are properly diabatic: if Eq. (1) is satisfied for even one pair of adiabatic states, then the derivative coupling must be zero between *all* pairs of diabatic states.

XVIII. Training species

Here we provide the motifs and substituents used for combinatorial molecule generation, together with a list of the literature species used for dense configurational sampling. Species with a net charge were excluded from training but are given here for completeness. In some cases only the *cis* or *trans* isomer was actually investigated experimentally, but in all cases we reference the publication for both isomers.

Many species in both the training and test set had experimental S_1 yields in non-polar solution. We did not put all such molecules in the test set, because this would mean losing hundreds of thousands of training geometries. Instead we used the 40 species with the fewest QC calculations.

Table S10. Motifs used for combinatorial species generation. Examples of literature species for each motif are also given. The species numbers are those in Tables S12 and S13.

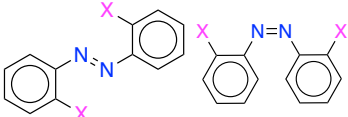
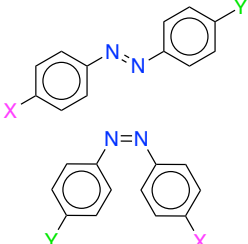
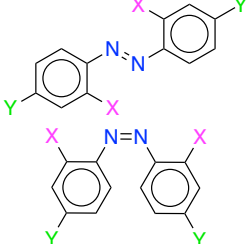
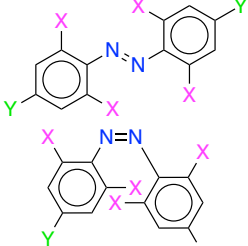
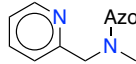
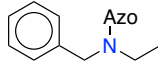
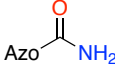
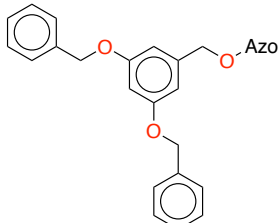
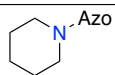
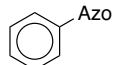
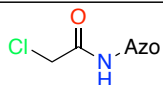
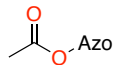
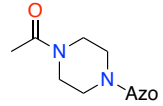
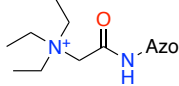
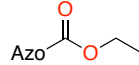
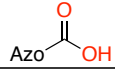
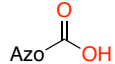
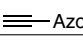
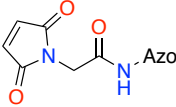

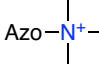
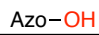
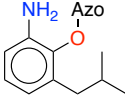
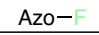
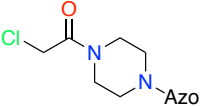

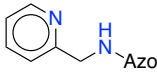
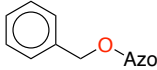
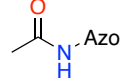
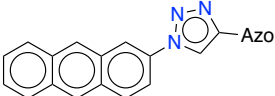
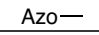
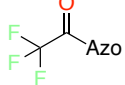
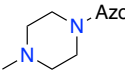
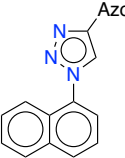
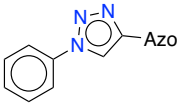
Graphs	Literature species
	1, 2, 25, 26, 63, 64, 125, 126
	19, 20, 27, 28, 37, 38, 39, 40, 61, 62, 69, 70, 73, 74, 75, 76, 83, 84, 97, 98, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 127, 128, 129, 130, 147, 148, 155, 156, 157, 158, 163, 164
	67, 68, 87, 88, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108
	9, 10, 13, 14, 15, 16, 17, 18, 21, 22, 95, 96, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 161, 162

Table S11. Substituents used for combinatorial species generation. "Azo" denotes the azobenzene attachment site.

SMILES	Graph
[Azo]N(C)Cc1cccn1	
[Azo]N	Azo-NH ₂
[Azo]N(CC)Cc1ccccc1	
[Azo]C(N)=O	
[Azo]OCc1cc(OCc2ccccc2)cc(OCc2ccccc2)c1	
[Azo]N1CCCCC1	
[Azo]c1ccccc1	
[Azo]NC(=O)CCl	
[Azo]OC(C)=O	
[Azo]N1CCN(C(C)=O)CC1	
[Azo]NC(=O)C[N+](CC)(CC)CC	
[Azo]Cl	Azo-Cl
[Azo]CO	Azo-CH ₂ -OH
[Azo]C(=O)OCC	
[Azo]CC	
[Azo]C(=O)O	
[Azo]C#C	

<chem>[Azo]NC(=O)CN1C(=O)C=CC1=O</chem>	
<chem>[Azo][N+](=O)[O-]</chem>	
<chem>[Azo][N+](C)(C)C</chem>	
<chem>[Azo]O</chem>	
<chem>[Azo]Oc1c(N)cccc1CC(C)C</chem>	
<chem>[Azo]F</chem>	
<chem>[Azo]N1CCN(C(=O)CCl)CC1</chem>	
<chem>[Azo]Oc1c(N)cccc1C(C)C</chem>	
<chem>[Azo]NCc1ccccn1</chem>	
<chem>[Azo]OCc1cccc1</chem>	
<chem>[Azo]NC(C)=O</chem>	
<chem>[Azo]c1cn(-c2ccc3cc4ccccc4cc3c2)nn1</chem>	
<chem>[Azo]C</chem>	
<chem>[Azo]C(=O)C(F)(F)F</chem>	
<chem>[Azo]N1CCN(C)CC1</chem>	
<chem>[Azo]c1cn(-c2cccc3ccccc23)nn1</chem>	
<chem>[Azo]c1cn(-c2ccccc2)nn1</chem>	

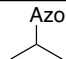
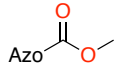
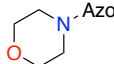
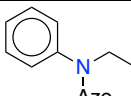
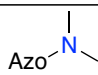
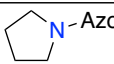
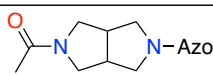

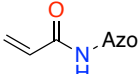
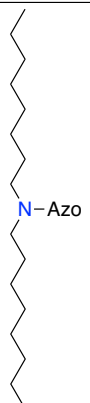
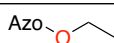
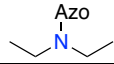
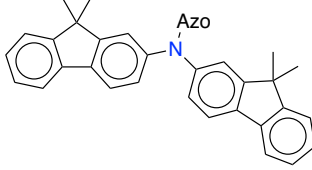
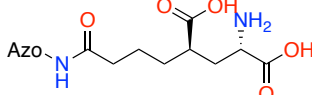
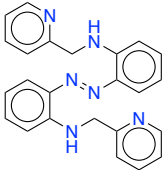
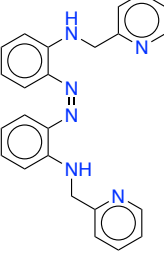
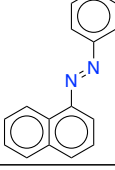
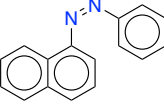
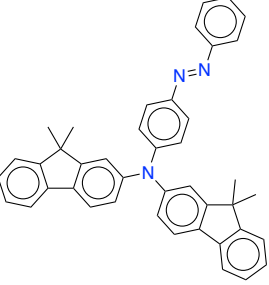
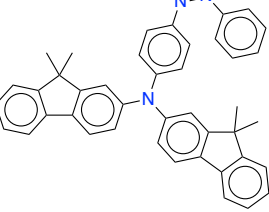
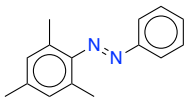
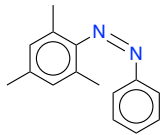
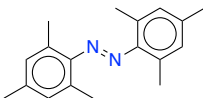
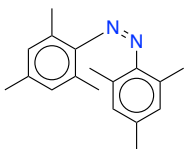
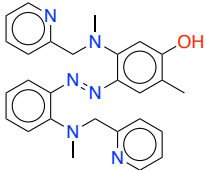
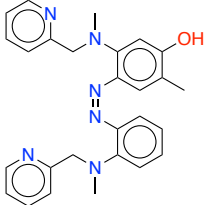
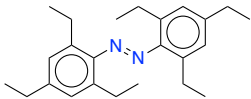
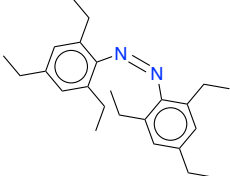
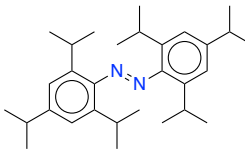
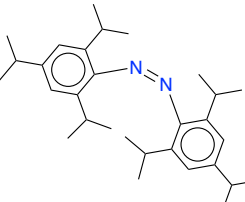
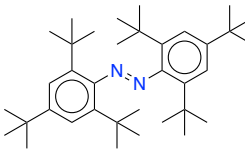
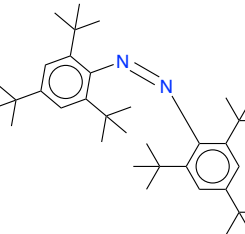
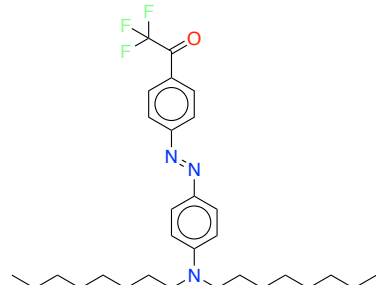
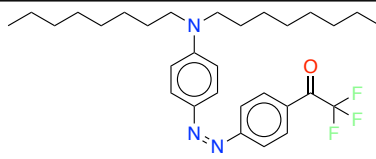
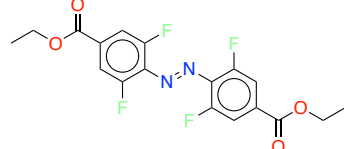
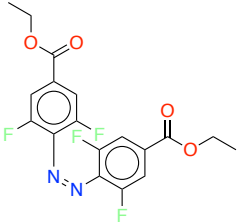
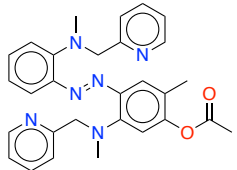
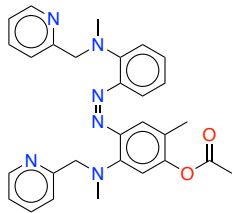
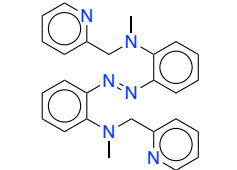
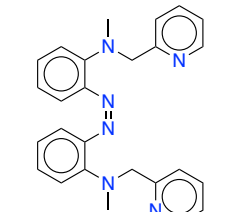
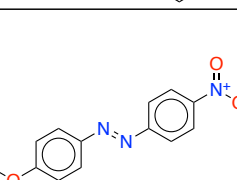
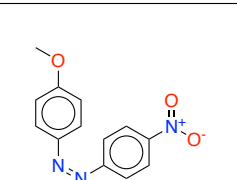
<chem>[Azo]C(C)C</chem>	
<chem>[Azo]C(=O)OC</chem>	
<chem>[Azo]N1CCOCC1</chem>	
<chem>[Azo]N(CC)c1ccccc1</chem>	
<chem>[Azo]N(C)C</chem>	
<chem>[Azo]N1CCCC1</chem>	
<chem>[Azo]N1CC2CN(C(C)=O)CC2C1</chem>	
<chem>[Azo]C(C)(C)C</chem>	
<chem>[Azo]NC(=O)C=C</chem>	
<chem>[Azo]N(CCCCCCCC)CCCCCCCC</chem>	
<chem>[Azo]OCC</chem>	
<chem>[Azo]N(CC)CC</chem>	
<chem>[Azo]N(c1ccc2c(c1)C(C)(C)c1ccccc1-2)c1ccc2c(c1)C(C)(C)c1ccccc1-2</chem>	
<chem>[Azo]NC(=O)CCC[C@H](C[C@H](N)C(=O)O)C(=O)O</chem>	

Table S12. Test set species and literature quantum yields. For papers that reported the yield at different wavelengths, we chose the wavelength closest to the $n - \pi^*$ (S_1) absorption maximum. Quantum yields computed with FS surface hopping using the diabatic model are also shown.

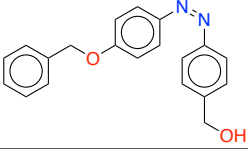
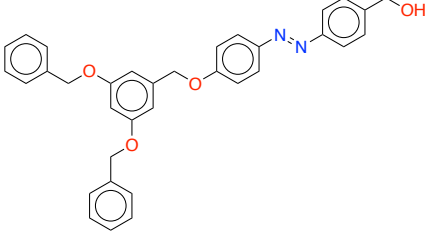
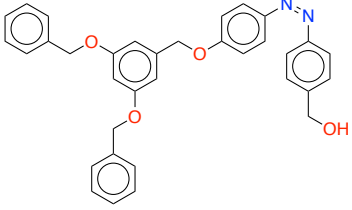
#	SMILES	Graph	Ref.	(S_1 yield, solvent)
1	<chem>c1ccc(CNc2ccccc2/N=N/c2ccccc2NCc2cccn2)nc1</chem>		[73]	Exp: (0^2 , many) Calc: 0 ± 0
2	<chem>c1ccc(CNc2ccccc2/N=N/c2ccccc2NCc2cccn2)nc1</chem>		[73]	Exp: None Calc: 0.51 ± 0.02
3	<chem>c1ccc(/N=N/c2cccc3ccccc23)cc1</chem>		[101]	Exp: (0.25, methyl cyclopentane) Calc: 0.43 ± 0.02
4	<chem>c1ccc(/N=N/c2cccc3ccccc23)cc1</chem>		[101]	Exp: (0.49, methyl cyclopentane) Calc: 0.57 ± 0.02
5	<chem>CC1(C)c2ccccc2-c2ccc(N(c3ccc(/N=N/c4ccccc4)cc3)c3ccc4c(c3)C(C)(C)c3ccccc3-4)cc21</chem>		[138]	Exp: (0.29, toluene) Calc: 0.19 ± 0.05
6	<chem>CC1(C)c2ccccc2-c2ccc(N(c3ccc(/N=N/c4ccccc4)cc3)c3ccc4c(c3)C(C)(C)c3ccccc3-4)cc21</chem>		[138]	Exp: (0.31, toluene) Calc: 0.55 ± 0.02

7	<chem>Cc1cc(C)c(/N=N/c2ccccc2)c(C)c1</chem>		[139]	Exp: (0.16, methylcyclohexane and isohehexane) Calc: 0.34 ± 0.03
8	<chem>Cc1cc(C)c(/N=N\c2ccccc2)c(C)c1</chem>		[139]	Exp: (0.38, methylcyclohexane and isohehexane) Calc: 0.60 ± 0.02
9	<chem>Cc1cc(C)c(/N=N/c2c(C)cc(C)cc2C)c(C)c1</chem>		[102, 139]	Exp: (0.16-0.22, methylcyclohexane and isohehexane) [139], (0.24, n-hexane) [102] Calc: 0.41 ± 0.02
10	<chem>Cc1cc(C)c(/N=N\c2c(C)cc(C)cc2C)c(C)c1</chem>		[102, 139]	Exp: (0.44, methylcyclohexane and isohehexane) [139], (0.5, n-hexane) [102] Calc: 0.59 ± 0.02
11	<chem>Cc1cc(/N=N/c2ccccc2N(C)Cc2ccccc2)c(N(C)Cc2ccccc2)cc1O</chem>		[74]	Exp: (0 ² , many) Calc: 0.33 ± 0.04
12	<chem>Cc1cc(/N=N\c2ccccc2N(C)Cc2ccccc2)c(N(C)Cc2ccccc2)cc1O</chem>		[74]	Exp: None Calc: 0.53 ± 0.02
13	<chem>CCc1cc(CC)c(/N=N/c2c(CC)cc(CC)cc2CC)c(CC)c1</chem>		[102]	Exp: (0.25, n-hexane) Calc: 0.32 ± 0.02
14	<chem>CCc1cc(CC)c(/N=N\c2c(CC)cc(CC)cc2CC)c(CC)c1</chem>		[102]	Exp: (0.5, n-hexane) Calc: 0.49 ± 0.02

15	<chem>CC(C)c1cc(C(C)C)c(/N=N/c2c(C(C)C)cc(C(C)C)cc2C(C)C)c(C(C)C)c1</chem>		[102]	Exp: (0.19, n-hexane) Calc: 0.27 ± 0.02
16	<chem>CC(C)c1cc(C(C)C)c(/N=N/c2c(C(C)C)cc(C(C)C)cc2C(C)C)c(C(C)C)c1</chem>		[102]	Exp: (0.55, n-hexane) Calc: 0.53 ± 0.02
17	<chem>CC(C)(C)c1cc(C(C)(C)C)c(/N=N/c2c(C(C)(C)C)cc(C(C)(C)C)cc2C(C)(C)C)c(C(C)(C)C)c1</chem>		[102]	Exp: (0, n-hexane) Calc: 0.01 ± 0.00
18	<chem>CC(C)(C)c1cc(C(C)(C)C)c(/N=N/c2c(C(C)(C)C)cc(C(C)(C)C)cc2C(C)(C)C)c(C(C)(C)C)c1</chem>		[102]	Exp: None Calc: 0.78 ± 0.02
19	<chem>CCCCCCCCN(CCCCCCCC)c1ccc(/N=N/c2ccc(C(=O)C(F)(F)F)cc2)cc1</chem>		[103]	Exp: (0.16 ³ , toluene) Calc: 0.08 ± 0.04
20	<chem>CCCCCCCCN(CCCCCCCC)c1ccc(/N=N/c2ccc(C(=O)C(F)(F)F)cc2)cc1</chem>		[103]	Exp: (0.7-1 ³ , toluene) Calc: 0.66 ± 0.02
21	<chem>CCOC(=O)c1cc(F)c(/N=N/c2c(F)cc(C(=O)OCC)cc2F)c(F)c1</chem>		[104]	Exp: (0.1, hexane) Calc: 0.08 ± 0.02

22	<chem>CCOC(=O)c1cc(F)c(/N=N\c2c(F)cc(C(=O)OCC)cc2F)c(F)c1</chem>		[104]	Exp: (0.53, hexane) Calc: 0.44 ± 0.02
23	<chem>CC(=O)Oc1cc(N(C)Cc2cccn2)c(/N=N/c2ccccc2N(C)Cc2cccn2)cc1C</chem>		[74]	Exp: (0.17, C6D6) Calc: 0.30 ± 0.04
24	<chem>CC(=O)Oc1cc(N(C)Cc2cccn2)c(/N=N/c2ccccc2N(C)Cc2cccn2)cc1C</chem>		[74]	Exp: None Calc: 0.51 ± 0.02
25	<chem>CN(Cc1cccn1)c1ccccc1/N=N/c1ccccc1N(C)Cc1cccn1</chem>		[73]	Exp: (0.19, CDCl3) Calc: 0.37 ± 0.03
26	<chem>CN(Cc1cccn1)c1ccccc1/N=N/c1ccccc1N(C)Cc1cccn1</chem>		[73]	Exp: None Calc: 0.48 ± 0.02
27	<chem>COc1ccc(/N=N/c2ccc([N+](=O)[O-])cc2)cc1</chem>		[139]	Exp: (0.17, methylcyclohexane and isohexane) Calc: 0.05 ± 0.02
28	<chem>COc1ccc(/N=N/c2ccc([N+](=O)[O-])cc2)cc1</chem>		[139]	Exp: (0.55, methylcyclohexane and isohexane) Calc: 0.46 ± 0.02

29	<chem>COc1cc(N(C)Cc2cccn2)c(/N=N/c2ccccc2N(C)Cc2cccn2)cc1C</chem>		[74]	Exp: (0 ² , many) Calc: 0.35 ± 0.03
30	<chem>COc1cc(N(C)Cc2cccn2)c(/N=N\c2ccccc2N(C)Cc2cccn2)cc1C</chem>		[74]	Exp: None Calc: 0.54 ± 0.02
31	<chem>Fc1cccc(F)c1/N=N/c1ccccc1</chem>		[104]	Exp: (0.32, hexane) Calc: 0.36 ± 0.02
32	<chem>Fc1cccc(F)c1/N=N\c1ccccc1</chem>		[104]	Exp: (0.55, hexane) Calc: 0.53 ± 0.02
33	<chem>Fc1cccc(F)c1/N=N/c1c(F)cc(-c2ccccc2)cc1F</chem>		[105]	Exp: (0.15 ± 0.015, n-hexane) Calc: 0.41 ± 0.02
34	<chem>Fc1cccc(F)c1/N=N\c1c(F)cc(-c2ccccc2)cc1F</chem>		[105]	Exp: (0.28 ± 0.028, n-hexane) Calc: 0.51 ± 0.02
35	<chem>Nc1ccccc1/N=N/c1ccccc1NCc1cccn1</chem>		[73]	Exp: (0 ² , many) Calc: 0.04 ± 0.04
36	<chem>Nc1ccccc1/N=N\c1ccccc1NCc1cccn1</chem>		[73]	Exp: None Calc: 0.54 ± 0.02
37	<chem>OCc1ccc(/N=N/c2ccc(OCc3ccccc3)cc2)cc1</chem>		[140]	Exp: (0.4, dichloromethane) Calc: 0.13 ± 0.03

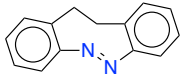
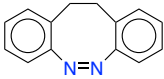
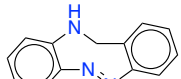
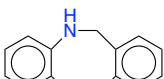
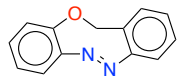
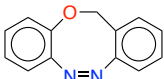
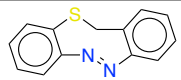
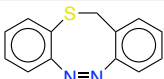
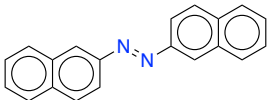
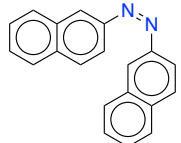
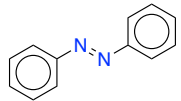
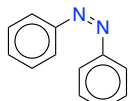
38	<chem>OCc1ccc(/N=N\c2ccc(OCc3ccccc3)cc2)cc1</chem>		[140]	Exp: (0.61, dichloromethane) Calc: 0.53 ± 0.02
39	<chem>OCc1ccc(/N=N/c2ccc(OCc3cc(OCc4ccccc4)cc(OCc4ccccc4)c3)cc2)cc1</chem>		[140]	Exp: (0.36, dichloromethane) Calc: 0.13 ± 0.03
40	<chem>OCc1ccc(/N=N\c2ccc(OCc3cc(OCc4ccccc4)cc(OCc4ccccc4)c3)cc2)cc1</chem>		[140]	Exp: (0.64, dichloromethane) Calc: 0.54 ± 0.02

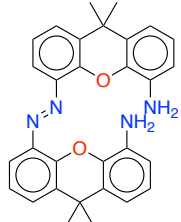
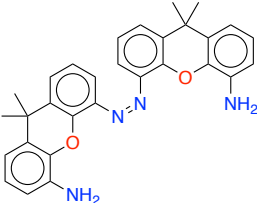
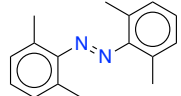
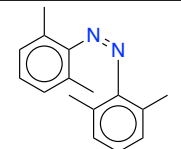
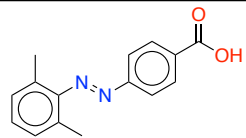
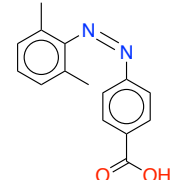
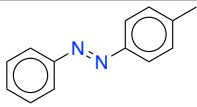
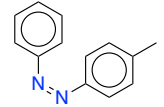
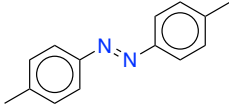
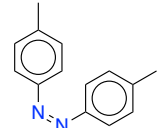
² In Refs. [73, 74], the non-reactive azobenzene derivatives were irradiated from 300 to 600 nm. This range covered both $S_0 \rightarrow S_1$ and $S_0 \rightarrow S_2$ excitation. Only small changes in absorbance were observed, indicating minimal *trans*→*cis* isomerization. However, the S_2 transition was highly red-shifted, and thus had significant overlap with the S_1 tran-

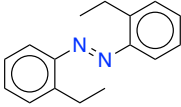
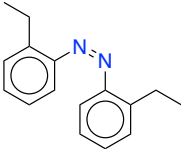
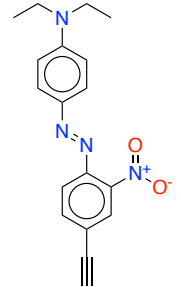
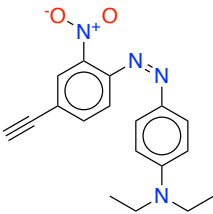
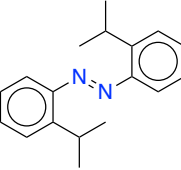
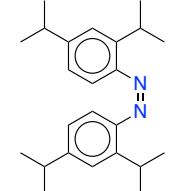
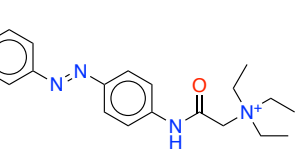
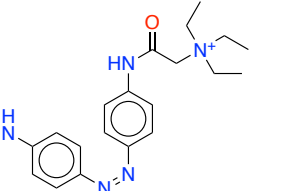
sition. The S_2 transition also had a much higher oscillator strength. Hence the small absorbance changes could have been due to the S_1 transition, and so the S_1 yield may not have been precisely 0. The S_1 yield should therefore be interpreted as “small”, rather than exactly zero.

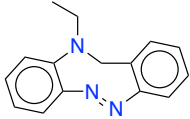
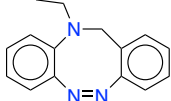
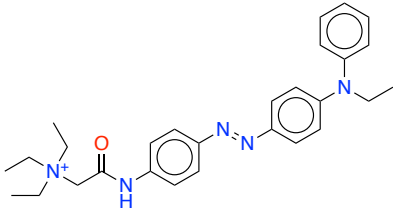
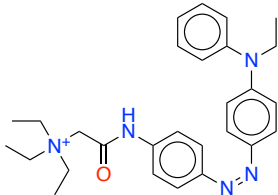
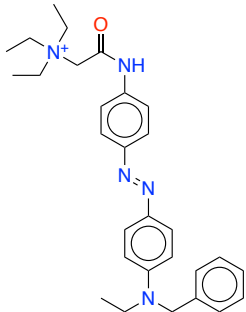
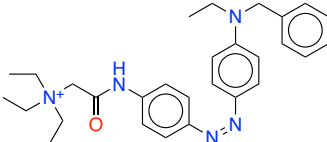
³ Averaged over excitations at 313, 436, and 546 nm.

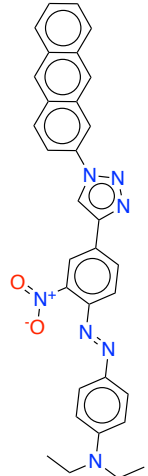
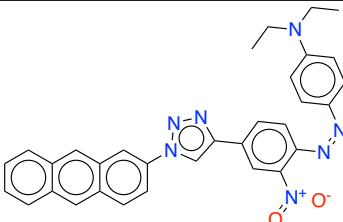
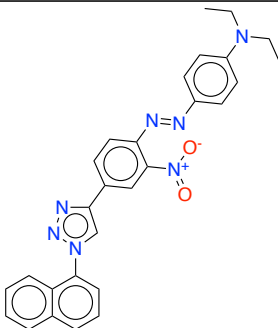
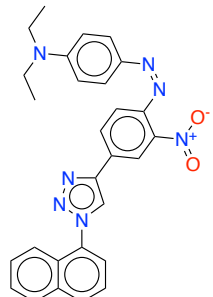
Table S13. Training and validation set literature species used for dense configurational sampling.

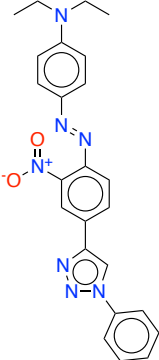
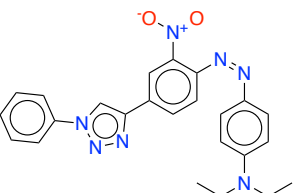
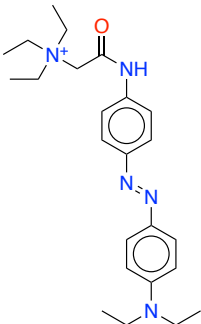
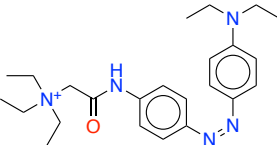
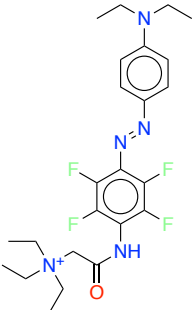
#	SMILES	Graph	Ref.
41	<chem>c1ccc2c(c1)CCc1ccccc1/N=N/2</chem>		[141]
42	<chem>c1ccc2c(c1)CCc1ccccc1/N=N\2</chem>		[141]
43	<chem>c1ccc2c(c1)CNc1ccccc1/N=N/2</chem>		[142]
44	<chem>c1ccc2c(c1)CNc1ccccc1/N=N\2</chem>		[142]
45	<chem>c1ccc2c(c1)COc1ccccc1/N=N/2</chem>		[142]
46	<chem>c1ccc2c(c1)COc1ccccc1/N=N\2</chem>		[142]
47	<chem>c1ccc2c(c1)CSc1ccccc1/N=N/2</chem>		[142]
48	<chem>c1ccc2c(c1)CSc1ccccc1/N=N\2</chem>		[142]
49	<chem>c1ccc2cc(/N=N/c3ccc4ccccc4c3)ccc2c1</chem>		[101]
50	<chem>c1ccc2cc(/N=N\c3ccc4ccccc4c3)ccc2c1</chem>		[101]
51	<chem>c1ccc(/N=N/c2ccccc2)cc1</chem>		[102, 105, 106, 108, 110, 139, 143]
52	<chem>c1ccc(/N=N\c2ccccc2)cc1</chem>		[102, 105, 106, 108, 110, 139, 143]

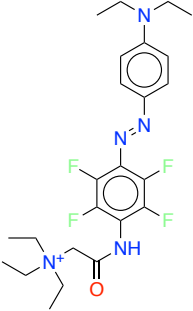
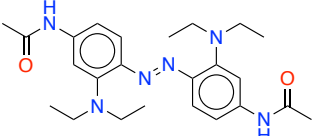
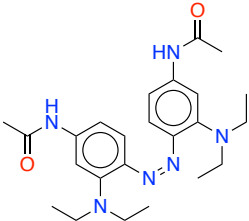
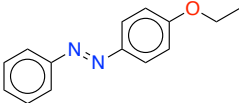
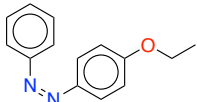
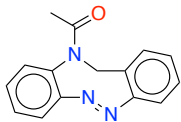
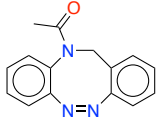
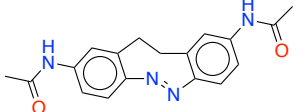
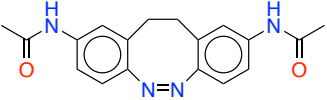
53	<chem>CC1(C)c2cccc(N)c2Oc2c(/N=N/c3cccc4c3Oc3c(N)cccc3C4(C)C)cccc21</chem>		[144]
54	<chem>CC1(C)c2cccc(N)c2Oc2c(/N=N\c3cccc4c3Oc3c(N)cccc3C4(C)C)cccc21</chem>		[144]
55	<chem>Cc1cccc(C)c1/N=N/c1c(C)cccc1C</chem>		[145]
56	<chem>Cc1cccc(C)c1/N=N\c1c(C)cccc1C</chem>		[145]
57	<chem>Cc1cccc(C)c1/N=N/c1ccc(C(=O)O)cc1</chem>		[146]
58	<chem>Cc1cccc(C)c1/N=N\c1ccc(C(=O)O)cc1</chem>		[146]
59	<chem>Cc1ccc(/N=N/c2ccccc2)cc1</chem>		[143]
60	<chem>Cc1ccc(/N=N\c2ccccc2)cc1</chem>		[143]
61	<chem>Cc1ccc(/N=N/c2ccc(C)cc2)cc1</chem>		[143]
62	<chem>Cc1ccc(/N=N\c2ccc(C)cc2)cc1</chem>		[143]

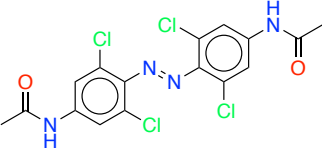
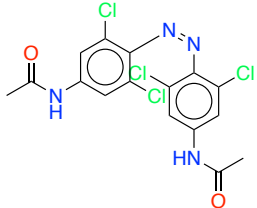
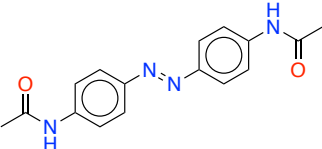
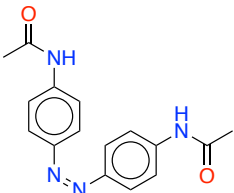
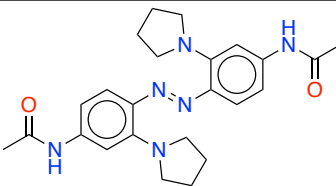
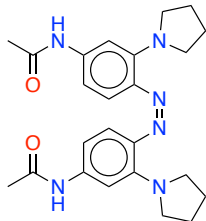
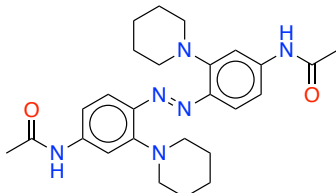
63	<chem>CCc1ccccc1/N=N/c1ccccc1CC</chem>		[145]
64	<chem>CCc1ccccc1/N=N\c1ccccc1CC</chem>		[145]
65	<chem>C#Cc1ccc(/N=N/c2ccc(N(CC)CC)cc2)c([N+](=O)[O-])c1</chem>		[147]
66	<chem>C#Cc1ccc(/N=N\c2ccc(N(CC)CC)cc2)c([N+](=O)[O-])c1</chem>		[147]
67	<chem>CC(C)c1ccc(/N=N/c2ccc(C(C)C)cc2C(C)C)c(C(C)C)c1</chem>		[102]
68	<chem>CC(C)c1ccc(/N=N\c2ccc(C(C)C)cc2C(C)C)c(C(C)C)c1</chem>		[102]
69	<chem>C=CC(=O)Nc1ccc(/N=N/c2ccc(NC(=O)C[N+](CC)(CC)CC)cc2)cc1</chem>		[148]
70	<chem>C=CC(=O)Nc1ccc(/N=N\c2ccc(NC(=O)C[N+](CC)(CC)CC)cc2)cc1</chem>		[148]

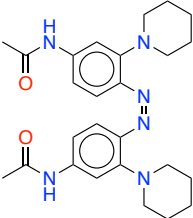
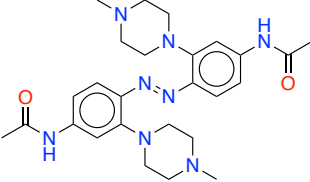
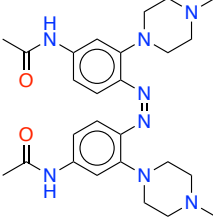
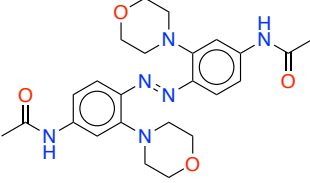
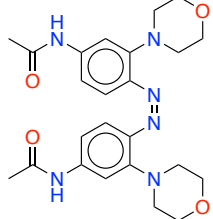
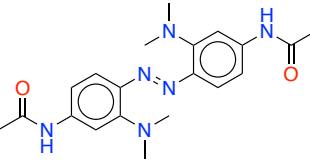
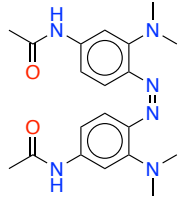
71	<chem>CCN1Cc2ccccc2/N=N/c2ccccc21</chem>		[142]
72	<chem>CCN1Cc2ccccc2/N=N\c2ccccc21</chem>		[142]
73	<chem>CCN(c1ccccc1)c1ccc(/N=N/c2ccc(NC(=O)C[N+](CC)(CC)CC)cc2)cc1</chem>		[148]
74	<chem>CCN(c1ccccc1)c1ccc(/N=N\c2ccc(NC(=O)C[N+](CC)(CC)CC)cc2)cc1</chem>		[148]
75	<chem>CCN(Cc1ccccc1)c1ccc(/N=N/c2ccc(NC(=O)C[N+](CC)(CC)CC)cc2)cc1</chem>		[148]
76	<chem>CCN(Cc1ccccc1)c1ccc(/N=N\c2ccc(NC(=O)C[N+](CC)(CC)CC)cc2)cc1</chem>		[148]

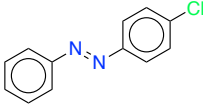
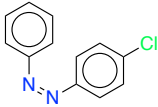
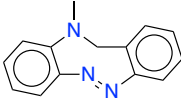
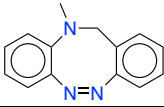
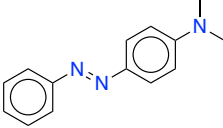
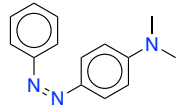
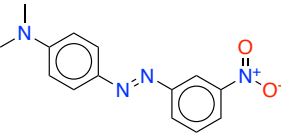
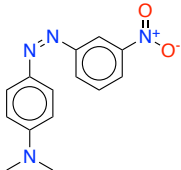
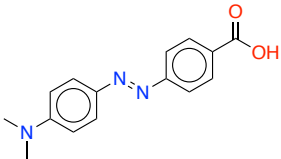
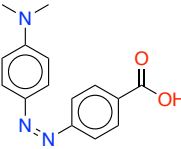
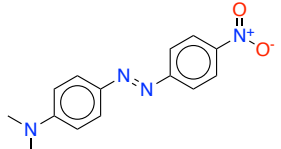
77	<chem>CCN(CC)c1ccc(/N=N/c2ccc(-c3cn(-c4ccc5cc6ccccc6cc5c4)nn3)cc2[N+](=O)[O-])cc1</chem>		[147]
78	<chem>CCN(CC)c1ccc(/N=N\c2ccc(-c3cn(-c4ccc5cc6ccccc6cc5c4)nn3)cc2[N+](=O)[O-])cc1</chem>		[147]
79	<chem>CCN(CC)c1ccc(/N=N/c2ccc(-c3cn(-c4cccc5ccccc45)nn3)cc2[N+](=O)[O-])cc1</chem>		[147]
80	<chem>CCN(CC)c1ccc(/N=N\c2ccc(-c3cn(-c4cccc5ccccc45)nn3)cc2[N+](=O)[O-])cc1</chem>		[147]

81	<chem>CCN(CC)c1ccc(/N=N/c2ccc(-c3cn(-c4ccccc4)nn3)cc2[N+](=O)[O-])cc1</chem>		[147]
82	<chem>CCN(CC)c1ccc(/N=N\c2ccc(-c3cn(-c4ccccc4)nn3)cc2[N+](=O)[O-])cc1</chem>		[147]
83	<chem>CCN(CC)c1ccc(/N=N/c2ccc(NC(=O)C[N+](CC)(CC)CC)cc2)cc1</chem>		[148]
84	<chem>CCN(CC)c1ccc(/N=N\c2ccc(NC(=O)C[N+](CC)(CC)CC)cc2)cc1</chem>		[148]
85	<chem>CCN(CC)c1ccc(/N=N/c2c(F)c(F)c(NC(=O)C[N+](CC)(CC)CC)c(F)c2F)cc1</chem>		[148]

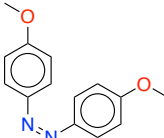
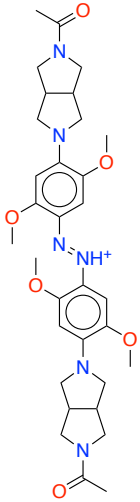
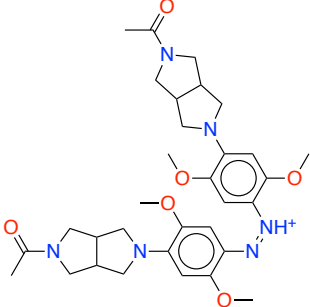
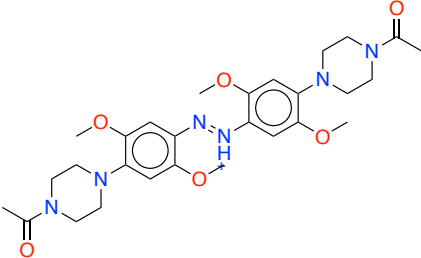
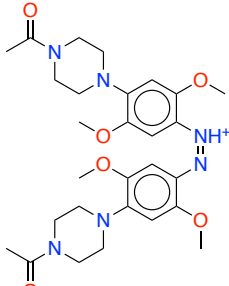
86	<chem>CCN(CC)c1ccc(/N=N\c2c(F)c(F)c(NC(=O)C[N+](CC)(CC)CC)c(F)c2F)cc1</chem>		[148]
87	<chem>CCN(CC)c1ccc(NC(C)=O)ccc1/N=N/c1ccc(NC(C)=O)cc1N(CC)CC</chem>		[149]
88	<chem>CCN(CC)c1ccc(NC(C)=O)ccc1/N=N\c1ccc(NC(C)=O)cc1N(CC)CC</chem>		[149]
89	<chem>CCOc1ccc(/N=N/c2ccccc2)cc1</chem>		[143]
90	<chem>CCOc1ccc(/N=N\c2ccccc2)cc1</chem>		[143]
91	<chem>CC(=O)N1Cc2ccccc2/N=N/c2ccccc21</chem>		[142]
92	<chem>CC(=O)N1Cc2ccccc2/N=N\c2ccccc21</chem>		[142]
93	<chem>CC(=O)Nc1ccc2c(c1)CCc1cc(NC(C)=O)ccc1/N=N/2</chem>		[150]
94	<chem>CC(=O)Nc1ccc2c(c1)CCc1cc(NC(C)=O)ccc1/N=N\2</chem>		[150]

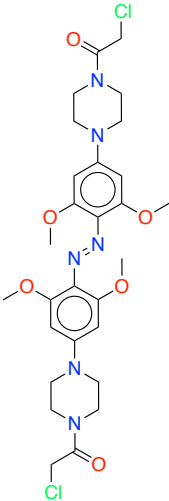
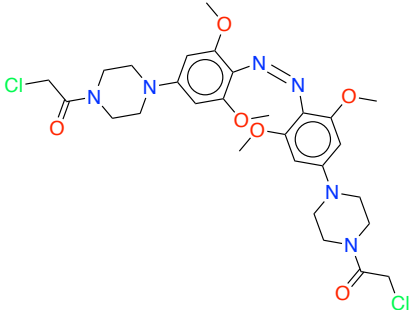
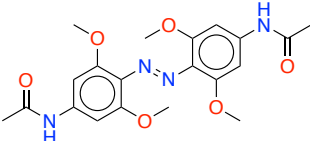
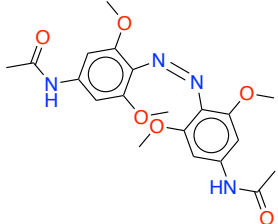
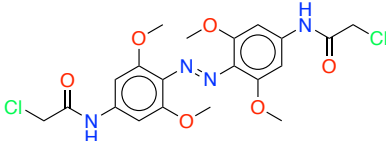
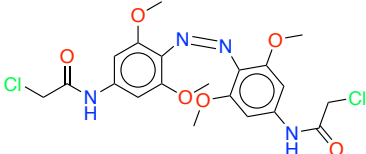
95	<chem>CC(=O)Nc1cc(Cl)c(/N=N/c2c(Cl)cc(NC(C)=O)cc2Cl)c(Cl)c1</chem>		[151]
96	<chem>CC(=O)Nc1cc(Cl)c(/N=N\c2c(Cl)cc(NC(C)=O)cc2Cl)c(Cl)c1</chem>		[151]
97	<chem>CC(=O)Nc1ccc(/N=N/c2ccc(NC(C)=O)cc2)cc1</chem>		[152]
98	<chem>CC(=O)Nc1ccc(/N=N\c2ccc(NC(C)=O)cc2)cc1</chem>		[152]
99	<chem>CC(=O)Nc1ccc(/N=N/c2ccc(NC(C)=O)cc2N2CCCC2)c(N2CCCC2)c1</chem>		[149]
100	<chem>CC(=O)Nc1ccc(/N=N\c2ccc(NC(C)=O)cc2N2CCCC2)c(N2CCCC2)c1</chem>		[149]
101	<chem>CC(=O)Nc1ccc(/N=N/c2ccc(NC(C)=O)cc2N2CCCCC2)c(N2CCCCC2)c1</chem>		[149]

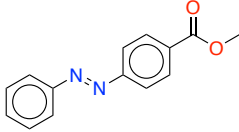
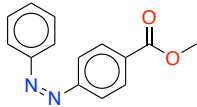
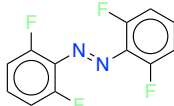
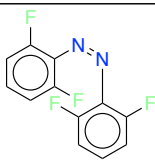
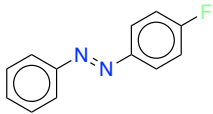
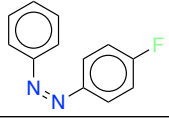
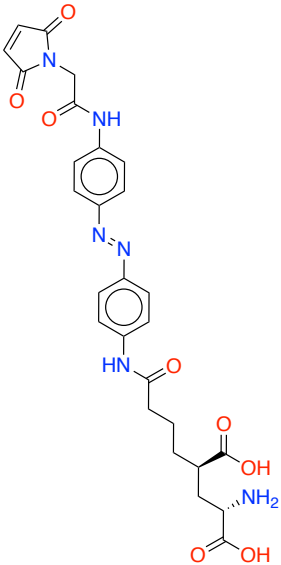
102	<chem>CC(=O)Nc1ccc(/N=N\c2ccc(NC(C)=O)cc2N2CCCC2)c(N2CCCC2)c1</chem>		[149]
103	<chem>CC(=O)Nc1ccc(/N=N/c2ccc(NC(C)=O)cc2N2CCN(C)CC2)c(N2CCN(C)CC2)c1</chem>		[149]
104	<chem>CC(=O)Nc1ccc(/N=N\c2ccc(NC(C)=O)cc2N2CCN(C)CC2)c(N2CCN(C)CC2)c1</chem>		[149]
105	<chem>CC(=O)Nc1ccc(/N=N/c2ccc(NC(C)=O)cc2N2CCOCC2)c(N2CCOCC2)c1</chem>		[149]
106	<chem>CC(=O)Nc1ccc(/N=N\c2ccc(NC(C)=O)cc2N2CCOCC2)c(N2CCOCC2)c1</chem>		[149]
107	<chem>CC(=O)Nc1ccc(/N=N/c2ccc(NC(C)=O)cc2N(C)C)c(N(C)C)c1</chem>		[149]
108	<chem>CC(=O)Nc1ccc(/N=N\c2ccc(NC(C)=O)cc2N(C)C)c(N(C)C)c1</chem>		[149]

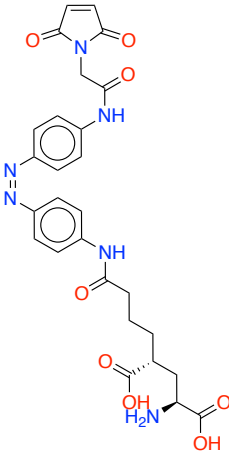
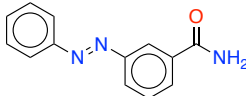
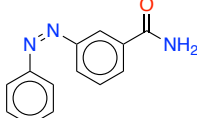
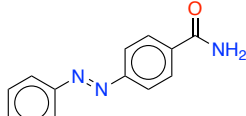
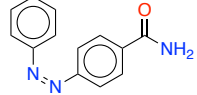
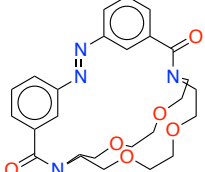
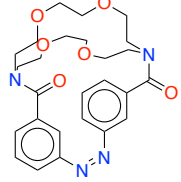
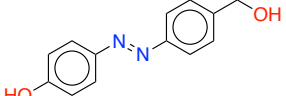
109	<chem>Clc1ccc(/N=N/c2ccccc2)cc1</chem>		[143]
110	<chem>Clc1ccc(/N=N\c2ccccc2)cc1</chem>		[143]
111	<chem>CN1Cc2ccccc2/N=N/c2ccccc21</chem>		[142]
112	<chem>CN1Cc2ccccc2/N=N\c2ccccc21</chem>		[142]
113	<chem>CN(C)c1ccc(/N=N/c2ccccc2)cc1</chem>		[153, 154]
114	<chem>CN(C)c1ccc(/N=N\c2ccccc2)cc1</chem>		[153, 154]
115	<chem>CN(C)c1ccc(/N=N/c2ccc([N+](=O)[O-])c2)cc1</chem>		[145]
116	<chem>CN(C)c1ccc(/N=N\c2ccc([N+](=O)[O-])c2)cc1</chem>		[145]
117	<chem>CN(C)c1ccc(/N=N/c2ccc(C(=O)O)cc2)cc1</chem>		[155]
118	<chem>CN(C)c1ccc(/N=N\c2ccc(C(=O)O)cc2)cc1</chem>		[155]
119	<chem>CN(C)c1ccc(/N=N/c2ccc([N+](=O)[O-])cc2)cc1</chem>		[156]

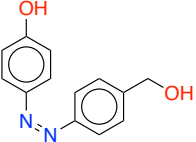
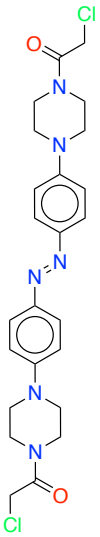
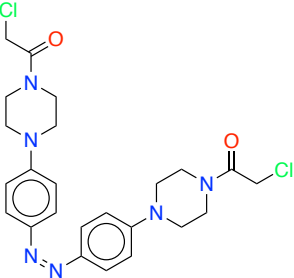
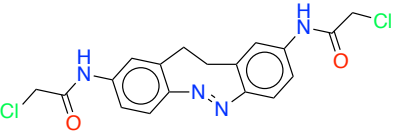
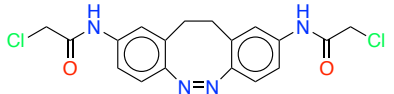
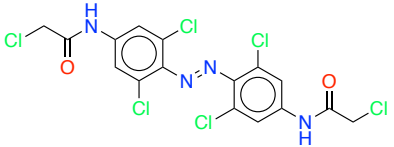
120	<chem>CN(C)c1ccc(/N=N\c2ccc([N+](=O)[O-])cc2)cc1</chem>		[156]
121	<chem>CN(C)c1ccc(/N=N/c2ccc(S(=O)(=O)O)cc2)cc1</chem>		[145]
122	<chem>CN(C)c1ccc(/N=N\c2ccc(S(=O)(=O)O)cc2)cc1</chem>		[145]
123	<chem>C[N+](C)(C)c1ccc(/N=N/c2ccc([N+](C)(C)C)cc2)cc1</chem>		[157]
124	<chem>C[N+](C)(C)c1ccc(/N=N\c2ccc([N+](C)(C)C)cc2)cc1</chem>		[157]
125	<chem>COc1ccccc1/N=N/c1ccccc1OC</chem>		[145, 158]
126	<chem>COc1ccccc1/N=N\c1ccccc1OC</chem>		[145, 158]
127	<chem>COc1ccc(/N=N/c2ccc(CO)cc2)cc1</chem>		[140]
128	<chem>COc1ccc(/N=N\c2ccc(CO)cc2)cc1</chem>		[140]
129	<chem>COc1ccc(/N=N/c2ccc(OC)cc2)cc1</chem>		[145]

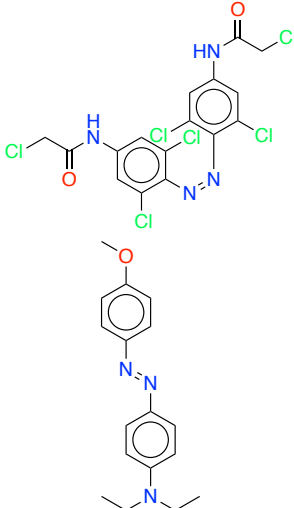
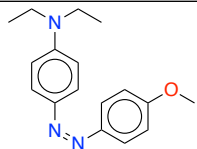
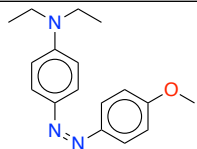
130	<chem>COc1ccc(/N=N\c2ccc(OC)cc2)cc1</chem>		[145]
131	<chem>COc1cc(N2CC3CN(C(C)=O)CC3C2)c(OC)cc1/N=[NH+]/c1cc(OC)c(N2CC3CN(C(C)=O)CC3C2)cc1OC</chem>		[159]
132	<chem>COc1cc(N2CC3CN(C(C)=O)CC3C2)c(OC)cc1/N=[NH+]/c1cc(OC)c(N2CC3CN(C(C)=O)CC3C2)cc1OC</chem>		[159]
133	<chem>COc1cc(N2CCN(C(C)=O)CC2)c(OC)cc1/N=[NH+]/c1cc(OC)c(N2CCN(C(C)=O)CC2)cc1OC</chem>		[159]
134	<chem>COc1cc(N2CCN(C(C)=O)CC2)c(OC)cc1/N=[NH+]/c1cc(OC)c(N2CCN(C(C)=O)CC2)cc1OC</chem>		[159]

135	<chem>COc1cc(N2CCN(C(=O)CCl)CC2)cc(OC)c1/N=N/c1c(OC)cc(N2CCN(C(=O)CCl)CC2)cc1OC</chem>		[160]
136	<chem>COc1cc(N2CCN(C(=O)CCl)CC2)cc(OC)c1/N=N/c1c(OC)cc(N2CCN(C(=O)CCl)CC2)cc1OC</chem>		[160]
137	<chem>COc1cc(NC(C)=O)cc(OC)c1/N=N/c1c(OC)cc(NC(C)=O)cc1OC</chem>		[151]
138	<chem>COc1cc(NC(C)=O)cc(OC)c1/N=N/c1c(OC)cc(NC(C)=O)cc1OC</chem>		[151]
139	<chem>COc1cc(NC(=O)CCl)cc(OC)c1/N=N/c1c(OC)cc(NC(=O)CCl)cc1OC</chem>		[151]
140	<chem>COc1cc(NC(=O)CCl)cc(OC)c1/N=N/c1c(OC)cc(NC(=O)CCl)cc1OC</chem>		[151]

141	<chem>COC(=O)c1ccc(/N=N/c2ccccc2)cc1</chem>		[143]
142	<chem>COC(=O)c1ccc(/N=N\c2ccccc2)cc1</chem>		[143]
143	<chem>Fc1cccc(F)c1/N=N/c1c(F)cccc1F</chem>		[104]
144	<chem>Fc1cccc(F)c1/N=N\c1c(F)cccc1F</chem>		[104]
145	<chem>Fc1ccc(/N=N/c2ccccc2)cc1</chem>		[143]
146	<chem>Fc1ccc(/N=N\c2ccccc2)cc1</chem>		[143]
147	<chem>N[C@@H](C[C@@H](CCCC(=O)Nc1ccc(/N=N/c2ccc(NC(=O)CN3C(=O)C=CC3=O)cc2)cc1)C(=O)O)C(=O)O</chem>		[161]

148	<chem>N[C@@H](C[C@@H](CCCC(=O)Nc1ccc(/N=N\c2ccc(NC(=O)CN3C(=O)C=CC3=O)cc2)cc1)C(=O)O)C(=O)O</chem>		[161]
149	<chem>NC(=O)c1cccc(/N=N/c2ccccc2)c1</chem>		[162]
150	<chem>NC(=O)c1cccc(/N=N\c2ccccc2)c1</chem>		[162]
151	<chem>NC(=O)c1ccc(/N=N/c2ccccc2)cc1</chem>		[162]
152	<chem>NC(=O)c1ccc(/N=N\c2ccccc2)cc1</chem>		[162]
153	<chem>O=C1c2cccc(c2)/N=N/c2cccc(c2)C(=O)N2CCOCCOCCN1CCOCCOCC2</chem>		[106, 163]
154	<chem>O=C1c2cccc(c2)/N=N\c2cccc(c2)C(=O)N2CCOCCOCCN1CCOCCOCC2</chem>		[106, 163]
155	<chem>OCc1ccc(/N=N/c2ccc(O)cc2)cc1</chem>		[140]

156	<chem>OCc1ccc(/N=N\c2ccc(O)cc2)cc1</chem>		[140]
157	<chem>O=C(CCl)N1CCN(c2ccc(/N=N/c3ccc(N4CCN(C(=O)CCl)CC4)cc3)cc2)CC1</chem>		[164]
158	<chem>O=C(CCl)N1CCN(c2ccc(/N=N/c3ccc(N4CCN(C(=O)CCl)CC4)cc3)cc2)CC1</chem>		[164]
159	<chem>O=C(CCl)Nc1ccc2c(c1)CCc1cc(NC(=O)CCl)ccc1/N=N/2</chem>		[150]
160	<chem>O=C(CCl)Nc1ccc2c(c1)CCc1cc(NC(=O)CCl)ccc1/N=N/2</chem>		[150]
161	<chem>O=C(CCl)Nc1cc(Cl)c(/N=N/c2c(Cl)cc(NC(=O)CCl)cc2Cl)c(Cl)c1</chem>		[151]

162	<chem>O=C(CCl)Nc1cc(Cl)c(/N=N\c2c(Cl)cc(NC(=O)CCl)cc2Cl)c(Cl)c1</chem>		[151]
163	<chem>CCN(CC)c1ccc(/N=N/c2ccc(OC)cc2)cc1</chem>		[154]
164	<chem>CCN(CC)c1ccc(/N=N\c2ccc(OC)cc2)cc1</chem>		[154]

-
- [1] Rachel C Evans, Peter Douglas, and Hugh D Burrow. *Applied photochemistry*. Springer, 2013.
- [2] Alexie M Kolpak and Jeffrey C Grossman. Azobenzene-functionalized carbon nanotubes as high-energy density solar thermal fuels. *Nano letters*, 11(8):3156–3162, 2011.
- [3] Sebastian Mai and Leticia González. Molecular photochemistry: Recent developments in theory. *Angewandte Chemie International Edition*, 59(39):16832–16846, 2020.
- [4] Johannes Broichhagen, James Allen Frank, and Dirk Trauner. A roadmap to success in photopharmacology. *Accounts of chemical research*, 48(7):1947–1960, 2015.
- [5] Michael M Lerch, Mickel J Hansen, Gooitzen M van Dam, Wiktor Szymanski, and Ben L Feringa. Emerging targets in photopharmacology. *Angewandte Chemie International Edition*, 55(37):10978–10999, 2016.
- [6] Sebastian Mai and Leticia González. Unconventional two-step spin relaxation dynamics of [Re (CO) 3 (im)(phen)]⁺ in aqueous solution. *Chemical science*, 10(44):10405–10411, 2019.
- [7] Jimmy K Yu, Christoph Bannwarth, Ruibin Liang, Edward G Hohenstein, and Todd J Martínez. Nonadiabatic dynamics simulation of the wavelength-dependent photochemistry of azobenzene excited to the $n\pi^*$ and $\pi\pi^*$ excited states. *Journal of the American Chemical Society*, 142(49):20680–20690, 2020.
- [8] Christoph Bannwarth, Jimmy K Yu, Edward G Hohenstein, and Todd J Martínez. Hole-hole tamm-dancoff-approximated density functional theory: A highly efficient electronic structure method incorporating dynamic and static correlation. *The Journal of Chemical Physics*, 153(2):024110, 2020.
- [9] JohnáC Tully. Mixed quantum–classical dynamics. *Faraday Discussions*, 110:407–419, 1998.
- [10] Dmitrii V Shalashilin. Quantum mechanics with the basis set guided by ehrenfest trajectories: Theory and application to spin-boson model. *The Journal of chemical physics*, 130(24):244101, 2009.
- [11] Michal Ben-Nun, Jason Quenneville, and Todd J Martínez. *Ab initio* multiple spawning: Photochemistry from first principles quantum molecular dynamics. *The Journal of Physical Chemistry A*, 104(22):5161–5175, 2000.
- [12] Toru Shiozaki, Werner Györfy, Paolo Celani, and Hans-Joachim Werner. Communication: Extended multi-state complete active space second-order perturbation theory: Energy and nuclear gradients. *Journal of Chemical Physics*, 135(8):081106–081106, 2011.
- [13] Jeppe Olsen, Björn O Roos, Poul Jørgensen, and Hans Jørgen Aa. Jensen. Determinant based configuration inter-

- action algorithms for complete and restricted configuration interaction spaces. *The Journal of chemical physics*, 89(4):2185–2192, 1988.
- [14] Dongxia Ma, Giovanni Li Manni, and Laura Gagliardi. The generalized active space concept in multiconfigurational self-consistent field methods. *The Journal of chemical physics*, 135(4):044128, 2011.
- [15] Giovanni Li Manni, Francesco Aquilante, and Laura Gagliardi. Strong correlation treated via effective hamiltonians and perturbation theory. *The Journal of chemical physics*, 134(3):034114, 2011.
- [16] Joseph Ivanic. Direct configuration interaction and multiconfigurational self-consistent-field method for multiple active spaces with variable occupations. I. Method. *The Journal of chemical physics*, 119(18):9364–9376, 2003.
- [17] Chenchen Song and Todd J Martínez. Reduced scaling extended multi-state caspt2 (xms-caspt2) using supporting subspaces and tensor hyper-contraction. *The Journal of chemical physics*, 152(23):234113, 2020.
- [18] Chenchen Song, Jeffrey B Neaton, and Todd J Martínez. Reduced scaling formulation of caspt2 analytical gradients using the supporting subspace method. *The Journal of Chemical Physics*, 154(1):014103, 2021.
- [19] Stefan Seritan, Christoph Bannwarth, B. Scott Fales, Edward G. Hohenstein, Sara I. L. Kokkila-Schumacher, Nathan Luehr, James W. Snyder Jr., Chenchen Song, Alexey V. Titov, Ivan S. Ufimtsev, and Todd J. Martínez. TeraChem: Accelerating electronic structure and ab initio molecular dynamics with graphical processing units. *The Journal of chemical physics*, 152(22):224110, 2020.
- [20] Konrad H Marti and Markus Reiher. New electron correlation theories for transition metal chemistry. *Physical Chemistry Chemical Physics*, 13(15):6750–6759, 2011.
- [21] Sandeep Sharma and Garnet Kin-Lic Chan. Spin-adapted density matrix renormalization group algorithms for quantum chemistry. *The Journal of chemical physics*, 136(12):124121, 2012.
- [22] Matthew R Hermes and Laura Gagliardi. Multiconfigurational self-consistent field theory with density matrix embedding: The localized active space self-consistent field method. *Journal of chemical theory and computation*, 15(2):972–986, 2019.
- [23] Christel M Marian, Adrian Heil, and Martin Kleinschmidt. The DFT/MRCI Method. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, 9(2):e1394, 2019.
- [24] Giovanni Li Manni, Rebecca K Carlson, Sijie Luo, Dongxia Ma, Jeppe Olsen, Donald G Truhlar, and Laura Gagliardi. Multiconfiguration pair-density functional theory. *Journal of chemical theory and computation*, 10(9):3669–3680, 2014.
- [25] Laura Gagliardi, Donald G Truhlar, Giovanni Li Manni, Rebecca K Carlson, Chad E Hoyer, and Junwei Lucas Bao. Multiconfiguration pair-density functional theory: A new way to treat strongly correlated systems. *Accounts of chemical research*, 50(1):66–73, 2017.
- [26] Christopher J. Stein and Markus Reiher. Automated selection of active orbital spaces. *Journal of chemical theory and computation*, 12(4):1760, 2016.
- [27] Christopher J. Stein and Markus Reiher. autoCAS: A program for fully automated multiconfigurational calculations. *Journal of computational chemistry*, 40(25):2216, 2019.
- [28] Jimmy K Yu, Christoph Bannwarth, Edward G Hohenstein, and Todd J Martínez. *Ab initio* nonadiabatic molecular dynamics with hole-hole Tamm–Dancoff approximated density functional theory. *Journal of Chemical Theory and Computation*, 16(9):5499–5511, 2020.
- [29] Michael Filatov. Spin-restricted ensemble-referenced kohn–sham method: basic principles and application to strongly correlated ground and excited states of molecules. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, 5(1):146–167, 2015.
- [30] Seunghoon Lee, Michael Filatov, Sangyoub Lee, and Cheol Ho Choi. Eliminating spin-contamination of spin-flip time dependent density functional theory within linear response formalism by the use of zeroth-order mixed-reference (mr) reduced density matrix. *The Journal of chemical physics*, 149(10):104101, 2018.
- [31] Teresa Cusati, Giovanni Granucci, Emilio Martínez-Núñez, Francesca Martini, Maurizio Persico, and Saulo Vázquez. Semiempirical Hamiltonian for simulation of azobenzene photochemistry. *The Journal of Physical Chemistry A*, 116(1):98–110, 2012.
- [32] Mayu Inamori, Takeshi Yoshikawa, Yasuhiro Iwabata, Yoshifumi Nishimura, and Hiromi Nakai. Spin-flip approach within time-dependent density functional tight-binding method: Theory and applications. *Journal of computational chemistry*, 41(16):1538–1548, 2020.
- [33] Marc de Wergifosse, Christoph Bannwarth, and Stefan Grimme. A simplified spin-flip time-dependent density functional theory approach for the electronic excitation spectra of very large diradicals. *The Journal of Physical Chemistry A*, 123(27):5815–5825, 2019.
- [34] Zhuoran Qiao, Matthew Welborn, Animashree Anandkumar, Frederick R Manby, and Thomas F Miller III. OrbNet: Deep learning for quantum chemistry using symmetry-adapted atomic-orbital features. *The Journal of Chemical Physics*, 153(12):124111, 2020.
- [35] Kristof T. Schütt, Oliver T. Unke, and Michael Gastegger. Equivariant message passing for the prediction of tensorial

- properties and molecular spectra. *arXiv preprint arXiv:2102.03150*, 2021.
- [36] Oliver T Unke, Stefan Chmiela, Michael Gastegger, Kristof T Schütt, Huziel E Sauceda, and Klaus-Robert Müller. SpookyNet: Learning force fields with electronic degrees of freedom and nonlocal effects. *arXiv preprint arXiv:2105.00304*, 2021.
- [37] Shi Jun Ang, Wujie Wang, Daniel Schwalbe-Koda, Simon Axelrod, and Rafael Gómez-Bombarelli. Active learning accelerates ab initio molecular dynamics on reactive energy surfaces. *Chem*, 7(3):738, 2021.
- [38] Wujie Wang, Tzuhsiung Yang, William H. Harris, and Rafael Gómez-Bombarelli. Active learning and neural network potentials accelerate molecular screening of ether-based solvate ionic liquids. *Chemical Communications*, 56(63):8920, 2020.
- [39] Wen-Kai Chen, Xiang-Yang Liu, Wei-Hai Fang, Pavlo O Dral, and Ganglong Cui. Deep learning for nonadiabatic excited-state dynamics. *The journal of physical chemistry letters*, 9(23):6702–6708, 2018.
- [40] Pavlo O Dral, Mario Barbatti, and Walter Thiel. Nonadiabatic excited-state dynamics with machine learning. *The journal of physical chemistry letters*, 9(19):5660–5663, 2018.
- [41] Deping Hu, Yu Xie, Xusong Li, Lingyue Li, and Zhenggang Lan. Inclusion of machine learning kernel ridge regression potential energy surfaces in on-the-fly nonadiabatic molecular dynamics simulation. *The journal of physical chemistry letters*, 9(11):2725–2732, 2018.
- [42] Jingbai Li, Patrick Reiser, Benjamin R Boswell, André Eberhard, Noah Z Burns, Pascal Friederich, and Steven A Lopez. Automatic discovery of photoisomerization mechanisms with nanosecond machine learning photodynamics simulations. *Chemical Science*, 2021.
- [43] Julia Westermayr, Michael Gastegger, Maximilian FSJ Menger, Sebastian Mai, Leticia González, and Philipp Marquetand. Machine learning enables long time scale molecular photodynamics simulations. *Chemical science*, 10(35):8100–8107, 2019.
- [44] Julia Westermayr, Michael Gastegger, and Philipp Marquetand. Combining SchNet and SHARC: The SchNarc machine learning approach for excited-state dynamics. *The journal of physical chemistry letters*, 11(10):3828–3834, 2020.
- [45] Julia Westermayr and Philipp Marquetand. Machine learning for electronically excited states of molecules. *Chemical Reviews*, 2020.
- [46] Dmitry V Makhov, William J Glover, Todd J Martinez, and Dmitrii V Shalashilin. Ab initio multiple cloning algorithm for quantum nonadiabatic molecular dynamics. *The Journal of chemical physics*, 141(5):054110, 2014.
- [47] Chaoyuan Zhu, Shikha Nangia, Ahren W Jasper, and Donald G Truhlar. Coherent switching with decay of mixing: an improved treatment of electronic coherence for non-born–oppenheimer trajectories. *The Journal of chemical physics*, 121(16):7658–7670, 2004.
- [48] GW Richings, Iakov Polyak, KE Spinlove, GA Worth, Irene Burghardt, and Benjamin Lasorne. Quantum dynamics simulations using gaussian wavepackets: the vmcg method. *International Reviews in Physical Chemistry*, 34(2):269–308, 2015.
- [49] Ali Abedi, Neepa T Maitra, and Eberhard KU Gross. Exact factorization of the time-dependent electron-nuclear wave function. *Physical review letters*, 105(12):123002, 2010.
- [50] Ali Abedi, Federica Agostini, and E KU Gross. Mixed quantum-classical dynamics from the exact decomposition of electron-nuclear motion. *EPL (Europhysics Letters)*, 106(3):33001, 2014.
- [51] Seung Kyu Min, Federica Agostini, Ivano Tavernelli, and Eberhard KU Gross. Ab initio nonadiabatic dynamics with coupled trajectories: A rigorous approach to quantum (de) coherence. *The journal of physical chemistry letters*, 8(13):3048–3055, 2017.
- [52] Basile FE Curchod and Federica Agostini. On the dynamics through a conical intersection. *The journal of physical chemistry letters*, 8(4):831–837, 2017.
- [53] Jong-Kwon Ha, In Seong Lee, and Seung Kyu Min. Surface hopping dynamics beyond nonadiabatic couplings for quantum coherence. *The journal of physical chemistry letters*, 9(5):1097–1104, 2018.
- [54] Michael H Beck, Andreas Jäckle, Graham A Worth, and H-D Meyer. The multiconfiguration time-dependent Hartree (MCTDH) method: a highly efficient algorithm for propagating wavepackets. *Physics reports*, 324(1):1–105, 2000.
- [55] Haobin Wang and Michael Thoss. Multilayer formulation of the multiconfiguration time-dependent Hartree theory. *The Journal of chemical physics*, 119(3):1289–1299, 2003.
- [56] Irene Burghardt, H-D Meyer, and LS Cederbaum. Approaches to the approximate treatment of complex molecular systems by the multiconfiguration time-dependent hartree method. *The Journal of chemical physics*, 111(7):2927–2939, 1999.
- [57] John C Tully. Molecular dynamics with electronic transitions. *The Journal of Chemical Physics*, 93(2):1061–1071, 1990.
- [58] Chaoyuan Zhu and Hiroki Nakamura. The two-state linear curve crossing problems revisited. ii. analytical approximations for the stokes constant and scattering matrix: The landau–zener case. *The Journal of chemical physics*, 97

- (11):8497–8514, 1992.
- [59] Chaoyuan Zhu and Hiroki Nakamura. The two-state linear curve crossing problems revisited. iii. analytical approximations for stokes constant and scattering matrix: Nonadiabatic tunneling case. *The Journal of chemical physics*, 98(8):6208–6222, 1993.
- [60] Yinan Shu and Donald G Truhlar. Diabatization by machine intelligence. *Journal of Chemical Theory and Computation*, 16(10):6456–6464, 2020.
- [61] David MG Williams and Wolfgang Eisfeld. Neural network diabatization: A new ansatz for accurate high-dimensional coupled potential energy surfaces. *The Journal of chemical physics*, 149(20):204106, 2018.
- [62] Yafu Guan, Dong H Zhang, Hua Guo, and David R Yarkony. Representation of coupled adiabatic potential energy surfaces using neural network based quasi-diabatic hamiltonians: 1, 2 2 a' states of lifh. *Physical Chemistry Chemical Physics*, 21(26):14205–14213, 2019.
- [63] Yihan Shao, Martin Head-Gordon, and Anna I Krylov. The spin-flip approach within time-dependent density functional theory: Theory and applications to diradicals. *The Journal of chemical physics*, 118(11):4807–4818, 2003.
- [64] Michelle M Francl, William J Pietro, Warren J Hehre, J Stephen Binkley, Mark S Gordon, Douglas J DeFrees, and John A Pople. Self-consistent molecular orbital methods. XXIII. A polarization-type basis set for second-row elements. *The Journal of Chemical Physics*, 77(7):3654–3665, 1982.
- [65] Axel D Becke. A new mixing of Hartree–Fock and local density-functional theories. *The Journal of chemical physics*, 98(2):1372–1377, 1993.
- [66] Benjamin G Levine, Chaehyuk Ko, Jason Quenneville, and Todd J Martínez. Conical intersections and double excitations in time-dependent density functional theory. *Molecular Physics*, 104(5-7):1039–1051, 2006.
- [67] Michael Filatov. Assessment of density functional methods for obtaining geometries at conical intersections in organic molecules. *Journal of chemical theory and computation*, 9(10):4526–4541, 2013.
- [68] C Alden Mead and Donald G Truhlar. Conditions for the definition of a strictly diabatic electronic basis for molecular systems. *The Journal of Chemical Physics*, 77(12):6090–6098, 1982.
- [69] Michael Baer and Robert Engelman. A study of the diabatic electronic representation within the Born-Oppenheimer approximation. *Molecular Physics*, 75(2):293–303, 1992.
- [70] A Toniolo, C Ciminelli, Maurizio Persico, and TJ Martínez. Simulation of the photodynamics of azobenzene on its first excited state: Comparison of full multiple spawning and surface hopping treatments. *The Journal of chemical physics*, 123(23):234308, 2005.
- [71] Horst Köppel, Joachim Gronki, and Susanta Mahapatra. Construction scheme for regularized diabatic states. *The Journal of Chemical Physics*, 115(6):2377, 2001.
- [72] Ling Yue, Yajun Liu, and Chaoyuan Zhu. Performance of TDDFT with and without spin-flip in trajectory surface hopping dynamics: cis \rightleftharpoons trans azobenzene photoisomerization. *Physical Chemistry Chemical Physics*, 20(37):24123–24139, 2018.
- [73] HM Dhammika Bandara, Tracey R Friss, Miriam M Enriquez, William Isley, Christopher Incarvito, Harry A Frank, Jose Gascon, and Shawn C Burdette. Proof for the concerted inversion mechanism in the trans \rightarrow cis isomerization of azobenzene using hydrogen bonding to induce isomer locking. *The Journal of organic chemistry*, 75(14):4817–4827, 2010.
- [74] HM Dhammika Bandara, Shannon Cawley, José A Gascón, and Shawn C Burdette. Short-circuiting azobenzene photoisomerization with electron-donating substituents and reactivating the photochemistry with chemical modification, 2011.
- [75] Rafael Gómez-Bombarelli, Jorge Aguilera-Iparraguirre, Timothy D Hirzel, David Duvenaud, Dougal Maclaurin, Martin A Blood-Forsythe, Hyun Sik Chae, Markus Einzinger, Dong-Gwang Ha, Tony Wu, et al. Design of efficient molecular organic light-emitting diodes by a high-throughput virtual screening and experimental approach. *Nature materials*, 15(10):1120–1127, 2016.
- [76] Willem A Velema, Wiktor Szymanski, and Ben L Feringa. Photopharmacology: beyond proof of principle. *Journal of the American Chemical Society*, 136(6):2178–2191, 2014.
- [77] HM Dhammika Bandara and Shawn C Burdette. Photoisomerization in different classes of azobenzene. *Chemical Society Reviews*, 41(5):1809–1825, 2012.
- [78] AR Dias, ME Minas Da Piedade, JA Martinho Simoes, JA Simoni, C Teixeira, HP Diogo, Yang Meng-Yan, and G Pilcher. Enthalpies of formation of cis-azobenzene and trans-azobenzene. *The Journal of Chemical Thermodynamics*, 24(4):439–447, 1992.
- [79] Le Yu, Chao Xu, Yibo Lei, Chaoyuan Zhu, and Zhenyi Wen. Trajectory-based nonadiabatic molecular dynamics without calculating nonadiabatic coupling in the avoided crossing case: Trans \rightleftharpoons cis photoisomerization in azobenzene. *Physical Chemistry Chemical Physics*, 16(47):25883–25895, 2014.
- [80] Zhuoran Qiao, Feizhi Ding, Matthew Welborn, Peter J. Bygrave, Animashree Anandkumar, Frederick R. Manby, and Thomas F. Miller III. Multi-task learning for electronic structure to predict and explore molecular potential

- energy surfaces. *arXiv preprint arXiv:2011.02680*, 2020.
- [81] Daniel Schwalbe-Koda, Aik Rui Tan, and Rafael Gómez-Bombarelli. Differentiable sampling of molecular geometries with uncertainty-based adversarial attacks. *Nature Communications*, 12(5104), 2021.
- [82] Troy Van Voorhis, Tim Kowalczyk, Benjamin Kaduk, Lee-Ping Wang, Chiao-Lun Cheng, and Qin Wu. The diabatic picture of electron transfer, reaction barriers, and molecular dynamics. *Annual review of physical chemistry*, 61: 149–170, 2010.
- [83] Michael S Schuurman and David R Yarkony. On the vibronic coupling approximation: A generally applicable approach for determining fully quadratic quasidiabatic coupled electronic state Hamiltonians. *The Journal of chemical physics*, 127(9):094104, 2007.
- [84] Yihan Shao, Zhengting Gan, Evgeny Epifanovsky, Andrew TB Gilbert, Michael Wormit, Joerg Kussmann, Adrian W Lange, Andrew Behn, Jia Deng, Xintian Feng, et al. Advances in molecular quantum chemistry contained in the Q-Chem 4 program package. *Molecular Physics*, 113(2):184–215, 2015.
- [85] Shuichi Nosé. A unified formulation of the constant temperature molecular dynamics methods. *The Journal of chemical physics*, 81(1):511–519, 1984.
- [86] William G Hoover. Canonical dynamics: Equilibrium phase-space distributions. *Physical review A*, 31(3):1695, 1985.
- [87] Hanneli R Hudock, Benjamin G Levine, Alexis L Thompson, Helmut Satzger, David Townsend, N Gador, Susanne Ullrich, Albert Stolow, and Todd J Martínez. Ab initio molecular dynamics and time-resolved photoelectron spectroscopy of electronically excited uracil and thymine. *The Journal of Physical Chemistry A*, 111(34):8500–8508, 2007.
- [88] Sebastian Mai, Philipp Marquetand, and Leticia González. Nonadiabatic dynamics: The SHARC approach. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, 8(6):e1370, 2018.
- [89] John M Herbert, Xing Zhang, Adrian F Morrison, and Jie Liu. Beyond time-dependent density functional theory using only single excitations: Methods for computational studies of excited states in complex systems. *Accounts of chemical research*, 49(5):931–941, 2016.
- [90] TJ Martínez, M Ben-Nun, and RD Levine. Molecular collision dynamics on several electronic states. *The Journal of Physical Chemistry A*, 101(36):6389–6402, 1997.
- [91] Eugene S Kryachko and David R Yarkony. Diabatic bases and molecular properties. *International Journal of Quantum Chemistry*, 76(2):235–243, 2000.
- [92] Chad E Hoyer, Kelsey Parker, Laura Gagliardi, and Donald G Truhlar. The DQ and DQ Φ electronic structure diabatization methods: Validation for general applications. *The Journal of chemical physics*, 144(19):194101, 2016.
- [93] Hisao Nakamura and Donald G Truhlar. The direct calculation of diabatic states based on configurational uniformity. *The Journal of Chemical Physics*, 115(22):10353–10372, 2001.
- [94] Felix Plasser, Sandra Gómez, Maximilian FSJ Menger, Sebastian Mai, and Leticia González. Highly efficient surface hopping dynamics using a linear vibronic coupling model. *Physical Chemistry Chemical Physics*, 21(1):57–69, 2019.
- [95] Yafu Guan, Bina Fu, and Dong H Zhang. Construction of diabatic energy surfaces for LiFH with artificial neural networks. *The Journal of chemical physics*, 147(22):224307, 2017.
- [96] Yinan Shu, Zoltan Varga, Antonio Gustavo Sampaio de Oliveira-Filho, and Donald G Truhlar. Permutationally restrained diabatization by machine intelligence. *Journal of Chemical Theory and Computation*, 17(2):1106–1116, 2021.
- [97] Martin Richter, Philipp Marquetand, Jesús González-Vázquez, Ignacio Sola, and Leticia González. SHARC: *ab initio* molecular dynamics with surface hopping in the adiabatic representation including arbitrary couplings. *Journal of chemical theory and computation*, 7(5):1253–1258, 2011.
- [98] Erich Runge and Eberhard KU Gross. Density-functional theory for time-dependent systems. *Physical Review Letters*, 52(12):997, 1984.
- [99] Samer Gozem, Federico Melaccio, Alessio Valentini, Michael Filatov, Miquel Huix-Rotllant, Nicolas Ferré, Luis Manuel Frutos, Celestino Angeli, Anna I Krylov, Alexander A Granovsky, et al. Shape of multireference, equation-of-motion coupled-cluster, and density functional theory potential energy surfaces at a conical intersection. *Journal of chemical theory and computation*, 10(8):3074–3084, 2014.
- [100] Alexander Nikiforov, Jose A Gamez, Walter Thiel, Miquel Huix-Rotllant, and Michael Filatov. Assessment of approximate computational methods for conical intersections and branching plane vectors in organic molecules. *The Journal of chemical physics*, 141(12):124122, 2014.
- [101] Shmuel Malkin and Ernst Fischer. Temperature dependence of photoisomerization. Part II. Quantum yields of cis-trans isomerizations in azo-compounds. *The Journal of Physical Chemistry*, 66(12):2482–2486, 1962.
- [102] Hermann Rau and Shen Yu-Quan. Photoisomerization of sterically hindered azobenzenes. *Journal of Photochemistry and Photobiology A: Chemistry*, 42(2-3):321–327, 1988.
- [103] Gerhard J Mohr and Ulrich-W Grummt. Photochemistry of the amine-sensor dye 4-N, N-dioctylamino-4'-trifluoroacetylazobenzene. *Journal of Photochemistry and Photobiology A: Chemistry*, 163(3):341–345, 2004.

- [104] Christopher Knie, Manuel Utecht, Fangli Zhao, Hannes Kulla, Sergey Kovalenko, Albert M Brouwer, Peter Saalfrank, Stefan Hecht, and David Bléger. *ortho*-Fluoroazobenzenes: Visible light switches with very long-lived Z isomers. *Chemistry—A European Journal*, 20(50):16492–16501, 2014.
- [105] J Moreno, M Gerecke, AL Dobryakov, IN Ioffe, AA Granovsky, D Bléger, S Hecht, and SA Kovalenko. Two-photon-induced versus one-photon-induced isomerization dynamics of a bistable azobenzene derivative in solution. *The Journal of Physical Chemistry B*, 119(37):12281–12288, 2015.
- [106] Hermann Rau. Further evidence for rotation in the π , π^* and inversion in the n, π^* photoisomerization of azobenzenes. *Journal of photochemistry*, 26(2-3):221–225, 1984.
- [107] Ernst Fischer. Calculation of photostationary states in systems $A \rightleftharpoons B$ when only A is known. *The Journal of Physical Chemistry*, 71(11):3704–3706, 1967.
- [108] George Zimmerman, Lue-Yung Chow, and Un-Jin Paik. The photochemical isomerization of azobenzene. *Journal of the American Chemical Society*, 80(14):3528–3531, 1958.
- [109] Philipp Marquetand, Martin Richter, Jesús González-Vázquez, Ignacio Sola, and Leticia González. Nonadiabatic ab initio molecular dynamics including spin-orbit coupling and laser fields. *Faraday discussions*, 153:261–273, 2011.
- [110] Pietro Bortolus and Sandra Monti. Cis-trans photoisomerization of azobenzene. Solvent and triplet donors effects. *Journal of Physical Chemistry*, 83(6):648–652, 1979.
- [111] Jacques Ronayette, René Arnaud, Philippe Lebourgeois, and Jacques Lemaire. Isomérisation photochimique de l’azobenzène en solution. i. *Canadian Journal of Chemistry*, 52(10):1848–1857, 1974.
- [112] Pietro Bortolus and Sandra Monti. Cis *rightleftharpoons* trans photoisomerization of azobenzene-cyclodextrin inclusion complexes. *Journal of Physical Chemistry*, 91(19):5046–5050, 1987.
- [113] Giustiniano Tiberio, Luca Muccioli, Roberto Berardi, and Claudio Zannoni. How does the trans-cis photoisomerization of azobenzene take place in organic solvents? *ChemPhysChem*, 11(5):1018, 2010.
- [114] Giovanni Granucci and Maurizio Persico. Excited state dynamics with the direct trajectory surface hopping method: azobenzene and its derivatives as a case study. *Theoretical Chemistry Accounts*, 117(5):1131–1143, 2007.
- [115] Brian R Landry and Joseph E Subotnik. How to recover Marcus theory with fewest switches surface hopping: Add just a touch of decoherence. *The Journal of chemical physics*, 137(22):22A513, 2012.
- [116] Hsing-Ta Chen and David R Reichman. On the accuracy of surface hopping dynamics in condensed phase non-adiabatic problems. *The Journal of Chemical Physics*, 144(9):094104, 2016.
- [117] J Patrick Zobel, Moritz Heindl, Felix Plasser, Sebastian Mai, and Leticia González. Surface hopping dynamics on vibronic coupling models. *Accounts of chemical research*, 54(20):3760–3771, 2021.
- [118] Lea M Ibele, Yorick Lassmann, Todd J Martínez, and Basile FE Curchod. Comparing (stochastic-selection) ab initio multiple spawning with trajectory surface hopping for the photodynamics of cyclopropanone, fulvene, and dithiane. *The Journal of Chemical Physics*, 154(10):104110, 2021.
- [119] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. 2019.
- [120] Johannes Klicpera, Janek Groß, and Stephan Günnemann. Directional message passing for molecular graphs. *arXiv preprint arXiv:2003.03123*, 2020.
- [121] Julia Westermayr and Philipp Marquetand. Deep learning for UV absorption spectra with SchNarc: First steps toward transferability in chemical compound space. *The Journal of Chemical Physics*, 153(15):154112, 2020.
- [122] Mario Barbatti and Kakali Sen. Effects of different initial condition samplings on photodynamics and spectrum of pyrrole. *International Journal of Quantum Chemistry*, 116(10):762–771, 2016.
- [123] Ask Hjorth Larsen, Jens Jørgen Mortensen, Jakob Blomqvist, Ivano E Castelli, Rune Christensen, Marcin Dułak, Jesper Friis, Michael N Groves, Bjørk Hammer, Cory Hargus, Eric D Hermes, Paul C Jennings, Peter Bjerre Jensen, James Kermode, John R Kitchin, Esben Leonhard Kolsbjerg, Joseph Kubal, Kristen Kaasbjerg, Steen Lysgaard, Jón Bergmann Maronsson, Tristan Maxson, Thomas Olsen, Lars Pastewka, Andrew Peterson, Carsten Rostgaard, Jakob Schiøtz, Ole Schütt, Mikkel Strange, Kristian S Thygesen, Tejs Vegge, Lasse Vilhelmsen, Michael Walter, Zhenhua Zeng, and Karsten W Jacobsen. The atomic simulation environment—a Python library for working with atoms. *Journal of Physics: Condensed Matter*, 29(27):273002, 2017.
- [124] Loup Verlet. Computer “experiments” on classical fluids. I. Thermodynamical properties of Lennard-Jones molecules. *Physical review*, 159(1):98, 1967.
- [125] Ling Yue, Le Yu, Chao Xu, Yibo Lei, Yajun Liu, and Chaoyuan Zhu. Benchmark performance of global switching versus local switching for trajectory surface hopping molecular dynamics simulation: Cis \rightleftharpoons trans azobenzene photoisomerization. *ChemPhysChem*, 18(10):1274–1287, 2017.
- [126] Sebastian Mai, Philipp Marquetand, and Leticia González. A general method to describe intersystem crossing

- dynamics in trajectory surface hopping. *International Journal of Quantum Chemistry*, 115(18):1215–1231, 2015.
- [127] Felix Plasser, Giovanni Granucci, Jiri Pittner, Mario Barbatti, Maurizio Persico, and Hans Lischka. Surface hopping dynamics using a locally diabatic formalism: Charge transfer in the ethylene dimer cation and excited state dynamics in the 2-pyridone dimer. *The Journal of chemical physics*, 137(22):22A514, 2012.
- [128] Neural Force Field. <https://github.com/learningmatter-mit/NeuralForceField>.
- [129] Giovanni Granucci and Maurizio Persico. Critical appraisal of the fewest switches algorithm for surface hopping. *The Journal of chemical physics*, 126(13):134114, 2007.
- [130] Axel D Becke. Density-functional exchange-energy approximation with correct asymptotic behavior. *Physical review A*, 38(6):3098, 1988.
- [131] John P Perdew. Density-functional approximation for the correlation energy of the inhomogeneous electron gas. *Physical Review B*, 33(12):8822, 1986.
- [132] Troy Van Voorhis and Martin Head-Gordon. A geometric approach to direct minimization. *Molecular Physics*, 100(11):1713–1721, 2002.
- [133] Justin S Smith, Benjamin T Nebgen, Roman Zubatyuk, Nicholas Lubbers, Christian Devereux, Kipton Barros, Sergei Tretiak, Olexandr Isayev, and Adrian Roitberg. Outsmarting quantum chemistry through transfer learning. 2019.
- [134] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [135] Yoon Kim, Carl Denton, Luong Hoang, and Alexander M Rush. Structured attention networks. *arXiv preprint arXiv:1702.00887*, 2017.
- [136] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017.
- [137] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- [138] Motoki Kurita, Miku Makihara, and Hideyuki Nakano. Photochromic reactions of 4-[bis (9, 9-dimethylfluoren-2-yl) amino] azobenzene in low molecular-mass organogels. *Soft materials*, 12(1):42–46, 2014.
- [139] Dina Gegiou, KA Muszkat, and Ernst Fischer. Temperature dependence of photoisomerization. V. Effect of substituents on the photoisomerization of stilbenes and azobenzenes. *Journal of the American Chemical Society*, 90(15):3907–3918, 1968.
- [140] Paul Sierocki, Huub Maas, Patrick Dragut, Gabriele Richardt, Fritz Vögtle, Luisa De Cola, Fred Brouwer, and Jeffrey I Zink. Photoisomerization of azobenzene derivatives in nanostructured silica. *The Journal of Physical Chemistry B*, 110(48):24390–24398, 2006.
- [141] Ron Siewertsen, Hendrikje Neumann, Bengt Buchheim-Stehn, Rainer Herges, Christian Nather, Falk Renth, and Friedrich Temps. Highly efficient reversible Z-E photoisomerization of a bridged azobenzene with visible light through resolved $S_1(n\pi^*)$ absorption bands. *Journal of the American Chemical Society*, 131(43):15594–15595, 2009.
- [142] Pascal Lentès, Eduard Stadler, Fynn Röhrlich, Arne Brahms, Jens Gröbner, Frank D Sönnichsen, Georg Gescheidt, and Rainer Herges. Nitrogen bridged diazocines: Photochromes switching within the near-infrared region with high quantum yields in organic solvents and in water. *Journal of the American Chemical Society*, 141(34):13592–13600, 2019.
- [143] PP Birnbaum and DWG Style. The photo-isomerization of some azobenzene derivatives. *Transactions of the Faraday Society*, 50:1192–1196, 1954.
- [144] S Anitha Nagamani, Yasuo Norikane, and Nobuyuki Tamaoki. Photoinduced hinge-like molecular motion: studies on xanthene-based cyclic azobenzene dimers. *The Journal of organic chemistry*, 70(23):9304–9313, 2005.
- [145] John Olmsted III, Jerry Lawrence, and Geary G Yee. Photochemical storage potential of azobenzenes. *Solar Energy*, 30(3):271–274, 1983.
- [146] Tao Chen, Atsushi Yamaguchi, Kazumasa Igarashi, Naoya Nakagawa, Hidenori Nishioka, Hiroyuki Asanuma, and Mikio Yamashita. Ultrafast photoisomerization and its single-shot pump pulse efficiency of trans-azobenzene derivative: Compound for photosensitive DNA. *Optics Communications*, 285(6):1206–1211, 2012.
- [147] Alexis Goulet-Hanssens, T Christopher Corkery, Arri Priimagi, and Christopher J Barrett. Effect of head group size on the photoswitching applications of azobenzene Disperse Red 1 analogues. *Journal of Materials Chemistry C*, 2(36):7505–7512, 2014.
- [148] Alexandre Mourot, Michael A Kienzler, Matthew R Banghart, Timm Fehrentz, Florian ME Huber, Marco Stein, Richard H Kramer, and Dirk Trauner. Tuning photochromic ion channel blockers. *ACS chemical neuroscience*, 2(9):536–543, 2011.
- [149] Oleg Sadvovskii, Andrew A Beharry, Fuzhong Zhang, and G Andrew Woolley. Spectral tuning of azobenzene photo-switches for biological applications. *Angewandte Chemie International Edition*, 48(8):1484–1486, 2009.

- [150] Subhas Samanta, Chuanguang Qin, Alan J Lough, and G Andrew Woolley. Bidirectional photocontrol of peptide conformation with a bridged azobenzene derivative. *Angewandte Chemie International Edition*, 51(26):6452–6455, 2012.
- [151] Subhas Samanta, Andrew A Beharry, Oleg Sadovski, Theresa M McCormick, Amirhossein Babalhavaeji, Vince Tropepe, and G Andrew Woolley. Photoswitching azo compounds in vivo with red light. *Journal of the American Chemical Society*, 135(26):9777–9784, 2013.
- [152] Marta Gascón-Moya, Arnau Pejoan, Mercè Izquierdo-Serra, Silvia Pittolo, Gisela Cabré, Jordi Hernando, Ramon Alibés, Pau Gorostiza, and Félix Busqué. An optimized glutamate receptor photoswitch with sensitized azobenzene isomerization. *The Journal of organic chemistry*, 80(20):9915–9925, 2015.
- [153] Hermann Rau, Gerhard Greiner, Guenter Gauglitz, and Herbert Meier. Photochemical quantum yields in the $A (+h\nu) \rightleftharpoons B (+h\nu, \Delta)$ system when only the spectrum of A is known. *Journal of physical chemistry*, 94(17):6523–6524, 1990.
- [154] Angelo Albini, Elisa Fasani, and Silvio Pietra. The photochemistry of azo dyes. Photoisomerisation versus photoreduction from 4-diethylaminoazobenzene and 4-diethylamino-4'-methoxyazobenzene. *Journal of the Chemical Society, Perkin Transactions 2*, (7):1021–1024, 1983.
- [155] Yu Jin Lee, Sung Ik Yang, Dae Seung Kang, and Sang-Woo Joo. Solvent dependent photo-isomerization of 4-dimethylaminoazobenzene carboxylic acid. *Chemical Physics*, 361(3):176–179, 2009.
- [156] FO Koller, C Sobotta, TE Schrader, T Cordes, WJ Schreier, A Sieg, and P Gilch. Slower processes of the ultrafast photo-isomerization of an azobenzene observed by IR spectroscopy. *Chemical Physics*, 341(1-3):258–266, 2007.
- [157] Tetsuro Umemoto, Yuta Ohtani, Takamasa Tsukamoto, Tetsuya Shimada, and Shinsuke Takagi. Pinning effect for photoisomerization of a dicationic azobenzene derivative by anionic sites of the clay surface. *Chemical Communications*, 50(3):314–316, 2014.
- [158] Ruriko Tahara, Tatsuya Morozumi, Hiroshi Nakamura, and Masatsugu Shimomura. Photoisomerization of azobenzocrown ethers. Effect of complexation of alkaline earth metal ions. *The Journal of Physical Chemistry B*, 101(39):7736–7743, 1997.
- [159] M Dong, A Babalhavaeji, MJ Hansen, L Kalman, and GA Woolley. Red, far-red, and near infrared photoswitches based on azonium ions. *Chemical Communications*, 51(65):12981–12984, 2015.
- [160] Subhas Samanta, Amirhossein Babalhavaeji, Ming-xin Dong, and G Andrew Woolley. Photoswitching of ortho-substituted azonium ions by red light in whole blood. *Angewandte Chemie*, 125(52):14377–14380, 2013.
- [161] Elizabeth C Carroll, Shai Berlin, Joshua Levitz, Michael A Kienzler, Zhe Yuan, Dorte Madsen, Delmar S Larsen, and Ehud Y Isacoff. Two-photon brightness of azobenzene photoswitches designed for glutamate receptor optogenetics. *Proceedings of the National Academy of Sciences*, 112(7):E776–E785, 2015.
- [162] Andreas Archut, Fritz Vögtle, Luisa De Cola, Gianluca Camillo Azzellini, Vincenzo Balzani, PS Ramanujam, and Rolf H Berg. Azobenzene-functionalized cascade molecules: photoswitchable supramolecular systems. *Chemistry—A European Journal*, 4(4):699–706, 1998.
- [163] Igor K Lednev, Tian-Qing Ye, Laurence C Abbott, Ronald E Hester, and John N Moore. Photoisomerization of a capped azobenzene in solution probed by ultrafast time-resolved electronic absorption spectroscopy. *The Journal of Physical Chemistry A*, 102(46):9161–9166, 1998.
- [164] Andrew A Beharry, Oleg Sadovski, and G Andrew Woolley. Photo-control of peptide conformation on a timescale of seconds with a conformationally constrained, blue-absorbing, photo-switchable linker. *Organic & biomolecular chemistry*, 6(23):4323–4332, 2008.