

1 **Identifying structure-absorption relationships and predicting absorption strength of non-**  
2 **fullerene acceptors for organic photovoltaics**

3  
4 *Jun Yan,<sup>a,#</sup> Xabier Rodríguez-Martínez,<sup>\*,b,c,#</sup> Drew Pearce,<sup>a</sup> Hana Douglas,<sup>a</sup> Danai Bili,<sup>a</sup>*  
5 *Mohammed Azzouzi,<sup>a</sup> Flurin Eisner,<sup>a</sup> Alise Virbule,<sup>a</sup> Elham Rezasoltani,<sup>a</sup> Valentina Belova,<sup>c</sup>*  
6 *Bernhard Döring,<sup>c</sup> Sheridan Few,<sup>a,f</sup> Anna A. Szumska,<sup>a</sup> Xueyan Hou,<sup>a</sup> Guichuan Zhang,<sup>d</sup> Hin-*  
7 *Lap Yip,<sup>d,e</sup> Mariano Campoy-Quiles<sup>\*,c</sup> and Jenny Nelson<sup>\*,a</sup>*

8  
9 # J.Y. and X.R.-M. contributed equally to this work.

10  
11 <sup>a</sup> Department of Physics, Imperial College London, SW7 2AZ, London, United Kingdom

12 Email: [jenny.nelson@imperial.ac.uk](mailto:jenny.nelson@imperial.ac.uk)

13  
14 <sup>b</sup> Electronic and Photonic Materials (EFM), Department of Physics, Chemistry and Biology  
15 (IFM), Linköping University, Linköping, SE 581 83 Sweden

16 Email: [xabier.rodriguez.martinez@liu.se](mailto:xabier.rodriguez.martinez@liu.se)

17  
18 <sup>c</sup> Instituto de Ciencia de Materiales de Barcelona, ICMAB-CSIC, Campus UAB, Bellaterra  
19 08193, Spain

20 Email: [mcampoy@icmab.es](mailto:mcampoy@icmab.es)

21  
22 <sup>d</sup> Institute of Polymer Optoelectronic Materials and Devices, State Key Laboratory of  
23 Luminescent Materials and Devices, South China University of Technology, Guangzhou  
24 510640, P. R. China

25  
26 <sup>e</sup> Department of Materials Science and Engineering, City University of Hong Kong, Tat Chee  
27 Avenue, Kowloon, Hong Kong

28  
29 <sup>f</sup> Sustainability Research Institute, School of Earth and Environment, University of Leeds,  
30 Leeds, LS2 9JT

33 **Keywords:** small-molecule acceptors, organic solar cells, absorption coefficient, machine-  
34 learning, density functional theory

35 **Abstract**

36 Non-fullerene acceptors (NFAs) are excellent light harvesters, yet the origin of such high  
37 optical extinction is not well understood. In this work, we investigate the absorption strength of  
38 NFAs by building a database of time-dependent density functional theory (TDDFT)  
39 calculations of  $\sim 500$   $\pi$ -conjugated molecules. The calculations are first validated by comparison  
40 with experimental measurements on solution and solid state using common fullerene and non-  
41 fullerene acceptors. We find that the molar extinction coefficient ( $\epsilon_{d,max}$ ) shows reasonable  
42 agreement between calculation in vacuum and experiment for molecules in solution,  
43 highlighting the effectiveness of TDDFT for predicting optical properties of organic  $\pi$ -  
44 conjugated molecules. We then perform a statistical analysis based on molecular descriptors to  
45 identify which features are important in defining the absorption strength. This allows us to  
46 identify structural features that are correlated with high absorption strength in NFAs and could  
47 be used to guide molecular design: highly absorbing NFAs should possess a planar, linear, and  
48 fully conjugated molecular backbone with highly polarisable heteroatoms. We then exploit a  
49 random decision forest to draw predictions for  $\epsilon_{d,max}$  using a computational framework based  
50 on extended tight-binding Hamiltonians, which shows reasonable predicting accuracy with  
51 lower computational cost than TDDFT. This work provides a general understanding of the  
52 relationship between molecular structure and absorption strength in  $\pi$ -conjugated organic  
53 molecules, including NFAs, while introducing predictive machine-learning models of low  
54 computational cost.

55

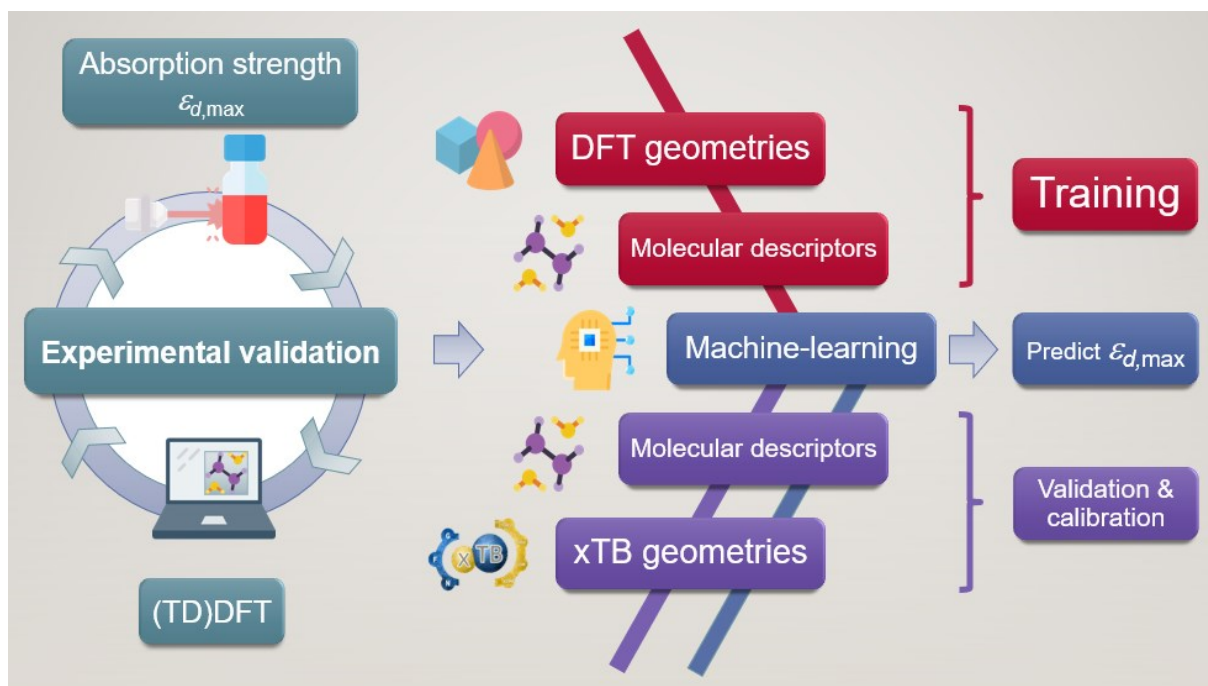
56 **Broader context**

57 The synthetic versatility of organic  $\pi$ -conjugated semiconductors converts them onto the ideal  
58 candidates for rational molecular design based on high-throughput screening techniques.  
59 Significant advances had been made by trial and error with new but increasingly diverse  
60 moieties and materials, primarily using non-fullerene acceptors (NFAs). These have raised the  
61 efficiency of organic photovoltaics (OPVs) above 19% in single junctions, to a large extent  
62 owing to their high absorption strength. However, the reasons for that remain elusive, thus  
63 preventing the molecular tailoring of NFAs with further enhanced light harvesting capabilities  
64 that enable breakthrough OPV efficiencies in the years to come. Here we exploit time-  
65 dependent density functional theory (TDDFT) calculations performed on NFA molecules and  
66  $\pi$ -conjugated oligomers to investigate what drives their absorption strength higher. The  
67 statistical analysis of thousands of molecular descriptors reveals that molecular linearity,  
68 planarity, polarizability, and number of  $\pi$ -conjugated carbon atoms correlate strongly with the  
69 absorption strength, hence forming a structure-absorption strength relationship that is further  
70 exploited to introduce design rules for highly absorbing NFAs. We identify frequent moieties  
71 (i.e. molecular fragments) and combinations thereof to drive absorption strength higher in novel  
72 NFAs. To speed up the screening of NFA molecular candidates at lower computational cost,  
73 we propose exploiting a state-of-the-art machine-learning (ML) model in combination with  
74 extended tight-binding Hamiltonians to predict the absorption strength of  $\pi$ -conjugated organic  
75 molecules. This work contributes to an improved understanding of the absorption strength of  
76  $\pi$ -conjugated organic molecules in general while suggesting ways the OPV community to  
77 design highly absorbing NFAs that maximize the light harvesting capabilities of materials for  
78 solar energy conversion.

79

80 **TOC:**

81 We combine experiments with density functional theory calculations, statistical analysis, and  
82 machine-learning to reveal the structure-absorption strength relationship and predict the  
83 absorption strength in organic non-fullerene acceptors.



84

85

## 86 1. Introduction

87 Organic photovoltaic (OPV) energy conversion is a promising option among next generation  
88 renewable and sustainable energy technologies for a low-carbon energy future.<sup>1-3</sup> OPV has  
89 shown promising potential for various applications, such as indoor photovoltaics (PV),<sup>4-6</sup> semi-  
90 transparent solar windows,<sup>7,8</sup> PV greenhouses,<sup>9</sup> and off-grid power supply.<sup>10</sup> Recent OPV  
91 devices based on non-fullerene acceptors (NFAs) have demonstrated certified power  
92 conversion efficiencies (PCEs) exceeding 19% in a single junction configuration,<sup>11</sup> much closer  
93 to the efficiencies observed in inorganic semiconductor PV technologies such as crystalline  
94 silicon and perovskite solar cells, and far higher than values thought attainable in OPV when  
95 using fullerene derivatives as the electron-acceptors.<sup>12</sup> The startling progress led by NFAs can  
96 be attributed to various advantages over fullerene derivatives, such as band-gap tunability, sharp  
97 absorption onset, high emission, high absorption, and low energy losses.<sup>13-15</sup> Among these  
98 advantages, the absorption strength of state-of-the-art NFAs is particularly outstanding, as  
99 exemplified in **Figure 1c** (a detailed list of chemical names and nomenclatures is provided in  
100 **Supplementary Note 1**).<sup>16</sup> For instance, Y6 shows a maximum extinction coefficient ( $\kappa_{max}$ )  
101 over 1.5 in the visible part of the electromagnetic spectrum, as compared to less than 0.75 for  
102 fullerene derivatives (PC61BM and PC71BM). High extinction coefficient increases the chance  
103 of high quantum efficiency and photogenerated current density, and makes it possible to  
104 fabricate highly absorbing OPV films with just a few tens of nanometre-thick photoactive layers.  
105 In comparison with workhorse fullerene acceptors, OPV devices based on highly absorbing  
106 NFAs could be made comparably thinner than the former, which exponentially raises the output  
107 power per weight (i.e. the specific weight in  $W\ g^{-1}$ ) of OPV devices<sup>17</sup> and might be an effective  
108 route toward lower production costs (as less material could employed to achieve an equivalent  
109 PCE) and even increase device thermal stability<sup>18</sup>. Moreover, through detailed balance between  
110 photon absorption and emission,<sup>19,20</sup> high absorption strength in principle should lead to high  
111 emission from the NFAs, while strong NFA emission is believed to be a key reason for NFA-  
112 based OPVs to possess low nonradiative voltage losses.<sup>21-25</sup> Despite the clear advantage of  
113 strong photo-absorption of NFAs over fullerene derivatives, the phenomenon has attracted  
114 much less attention than other properties of NFAs.<sup>21-23,25-28</sup> Conceptually, symmetry rules (i.e.,  
115 Laporte rule) can explain the qualitative difference between NFAs and fullerene derivatives in  
116 terms of absorption strength, yet such rules cannot predict differences in absorption strength  
117 among structures for which the lowest transitions are symmetry allowed. The features  
118 empirically and theoretically proposed<sup>29,30</sup> to lead to strong absorption in  $\pi$ -conjugated  
119 polymers are molecular stiffness, linearity, extended  $\pi$ -conjugation and large molecular size. It

120 is therefore of interest to establish whether the same (or other) molecular features are  
121 quantitatively associated or not with increased absorption strength in NFAs, while seeking  
122 molecular design rules to drive absorption and performance higher in new molecules.

123 Excited state calculations based on quantum chemistry methods, such as time-dependent  
124 density functional theory (TDDFT),<sup>8,29,31–33</sup> Hartree-Fock method,<sup>34</sup> ab initio Monte Carlo  
125 method,<sup>35</sup> second order Møller-Plesset theory (MP2),<sup>36</sup> and coupled cluster method,<sup>37</sup> have  
126 been applied to predict the electronic and optical properties of molecules. Among them, TDDFT  
127 is the most widely applied method for excited state calculations, and has shown reasonable  
128 accuracy in calculating and predicting the trends in absorption strength of organic  
129 molecules,<sup>29,31</sup> as also demonstrated in this work. However, the rapid scaling of computation  
130 time with molecular size has been the real obstacle limiting the applicability of TDDFT for  
131 excited state calculations on molecules with hundreds of atoms. Given the size and diverse  
132 structure of modern NFAs, faster and more efficient methods are therefore needed to establish  
133 the relationship between excited-state and molecular properties in NFAs.

134 The emergence of artificial intelligence (AI) has made it possible to study quantitative  
135 structure–property relationships (QSPRs) in molecules with massively improved computational  
136 efficiency. As the most popular branch of AI, machine-learning (ML) has attracted much  
137 attention in materials science over the last decade, and has been widely applied for material  
138 property prediction and material discovery.<sup>38–41</sup> Recently, ML has also gained popularity in  
139 OPV scenarios,<sup>42–52</sup> yet existing ML studies related to OPVs have been primarily focused either  
140 on the energetics<sup>42,43,53–56</sup> or directly on PCE,<sup>42,48,57–64</sup> with little attention paid to the absorption  
141 strength of the photoactive materials.<sup>65,66</sup> Moreover, there are no ML studies explicitly focused  
142 on the absorption strength of NFAs beyond the identification of moieties of frequent appearance  
143 in highly absorbing molecules.<sup>42</sup> However, QSPR and ML models have been successfully  
144 applied to investigate the absorption strength of fluorophores or dyes typically employed in  
145 bioimaging, showing encouraging results.<sup>30,67,68</sup> Therefore, it is appealing to apply ML methods  
146 in combination with QSPR models to investigate the origin of the large absorption strength in  
147 state-of-the-art NFAs.

148 Here, we present an experimental, TDDFT, QSPR, statistical and ML study of the absorption  
149 strength of NFAs to identify the key chemical and structural features that lead to high optical  
150 absorption in state-of-the-art NFAs. We exploit a database of nearly 500 unique organic  
151 molecules (or 3500 calculations) generated using DFT and TDDFT over several years. We  
152 obtain good quantitative agreement between TDDFT calculations of absorption strength and

153 experimental values for state-of-the-art NFAs and fullerenes, which supports the use of TDDFT  
154 results for further statistical and QSPR modelling. Accordingly, we extract molecular  
155 information from the DFT-optimized geometries by computing nearly 6000 molecular  
156 descriptors and first looking for correlations with the absorption strength. The strongest  
157 correlations are found between experimentally measured maximum molar extinction coefficient  
158 ( $\epsilon_{d,max}$ ) and two main molecular descriptors from calculations:  $\lambda_{l,p}$  and  $C2SP2$ , which describe  
159 the size of the molecule in the direction of maximal atomic polarizability, and the number of  
160  $sp^2$  hybridized carbon atoms that are bound to two other carbons (C2), respectively. These  
161 quantities can be related to a few key material features leading to high absorption strength:  
162 linearity, planarity, and extension of the  $\pi$ -conjugation in the form of fused and closed-ring  
163 moieties, in good agreement with previous ML reports on fluorophores and dyes.<sup>30</sup> We further  
164 identify several moieties and paired combinations thereof that are frequently found in highly  
165 absorbing NFAs, corresponding to thieno[3,2-b]thiophene (TT), thiophene (T), 2-(5,6-difluoro-  
166 3-oxo-2,3-dihydro-1H-inden-1-ylidene)malononitrile (2FIC), 2-(3-oxo-2,3-dihydro-1H-inden-  
167 1-ylidene)malononitrile (IC) and indaceno[1,2-b:5,6-b']dithiophene (IDT). These form a  
168 catalogue of molecular design rules to further enhance the absorption strength of organic  $\pi$ -  
169 conjugated molecules, such as next-generation NFAs. We then train and test an ensemble  
170 learning method, namely a random decision forest (RF), to predict  $\epsilon_{d,max}$  and provide further  
171 information about the most important features in the modelling of absorption strength in organic  
172  $\pi$ -conjugated molecules. Finally, we explore the possibility to predict  $\epsilon_{d,max}$  while using a  
173 cheaper molecular geometry optimization method based on semiempirical extended tight-  
174 binding (xTB) Hamiltonians instead of the expensive DFT approach. We do so by training a  
175 RF with our TDDFT database and proving its predictive properties in terms of  $\epsilon_{d,max}$  when  
176 interpolated using xTB-optimized geometries. This approach shows application potential in  
177 high-throughput screening studies in combination with generative molecular models.

## 178 **2. Results and discussion**

### 179 **2.1. Experimental validation of calculated absorption strength using TDDFT**

180 Quantifying how well the TDDFT derived excited state properties agree with the experimental  
181 measurements in terms of absorption strength is of utmost importance to validate our theoretical  
182 calculations and support further conclusions extracted thereof. Accordingly, we first evaluate  
183 the agreement between TDDFT calculations and experimental data in terms of the absorption  
184 strength. We compare the absorption strength of a broad catalogue (~10 molecules) of NFA

185 molecules and widely studied fullerene derivatives (PC61BM and PC71BM, with their  
186 molecular structures shown in **Figure 1a**) as obtained from TDDFT calculations, with a variety  
187 of optical measurements in both solution and solid state. For the most representative NFAs  
188 examined, we verify that their frontier molecular orbital energy levels as retrieved from TDDFT  
189 calculations are properly aligned, relative to those of a set of common polymer donors, for the  
190 NFAs to act as electron acceptor in a bulk heterojunction blend with those donors (**Figure S1**).  
191 The measured refractive index ( $n$ ) and extinction coefficient ( $\kappa$ ) of those molecules in thin film  
192 obtained using our variable-angle spectroscopic ellipsometry (VASE) measurements are shown  
193 in **Figure 1b** and c. Solution state data shown in **Figure 1d** and e are collected from a variety  
194 of literature references as detailed in the Supporting database.

195 As a metric for absorption strength, we initially consider several candidates such as the  
196 oscillator strength ( $f_{osc}$ ), the absorption coefficient ( $\alpha$ ) or the imaginary part of the dielectric  
197 function ( $\epsilon_2$ ). In this work, we eventually focus on the maximum molar extinction coefficient  
198 ( $\epsilon_{d,max}$ ,  $M^{-1} cm^{-1}$ ) of NFAs as it shows the best agreement between experimental and theoretical  
199 data, as we demonstrate below.  $\epsilon_{d,max}$  constitutes a typical experimental measurement in  
200 solution that can also be accessed from myriad literature references. Note that the usual  
201 calculations based on single molecules using TDDFT cannot account for solid state effects as  
202 they are performed for isolated molecules in vacuum or surrounded by an isotropic medium  
203 (such as a solvent using the polarizable-continuum-solvent-model, PCM, **Figure S2**). The  
204 derivation of the theoretical  $\epsilon_d$  is provided in the Methods section, which results in a  
205 mathematical expression for  $\epsilon_{d,max}$  as

$$206 \quad \epsilon_{d,max} = 10 \log_{10}(e) N_A \frac{2\pi e\hbar}{3\epsilon_0 m_0 n_r c} f_{osc,max} \frac{1}{\sigma\sqrt{2\pi}}, \quad (\text{Equation 1})$$

207 Where  $N_A$  is the Avogadro constant,  $e$  the elementary charge,  $\hbar$  the reduced Planck constant,  
208  $\epsilon_0$  the vacuum permittivity,  $m_0$  the electron mass,  $n_r$  is the refractive index in solution (assumed  
209 to be 1.3 of a common organic solvent throughout this study), and  $c$  the speed of light.  $f_{osc,max}$   
210 is the oscillator strength of the strongest transition among the calculated states within the  
211 visible-IR part of the spectrum, and  $E_{max}$  is the energy of that transition. The brightest  
212 transition is very often the lowest-energy transition in commonly used  $\pi$ -conjugated  
213 molecules.<sup>29</sup> We note here that the delta function in Eq. (S13) is replaced with a gaussian  
214 distribution function with a peak intensity of  $\frac{1}{\sigma\sqrt{2\pi}}$ , where  $\sigma$  is the gaussian width and assumed  
215 to be 0.1 eV for a common organic pi-conjugated molecule.



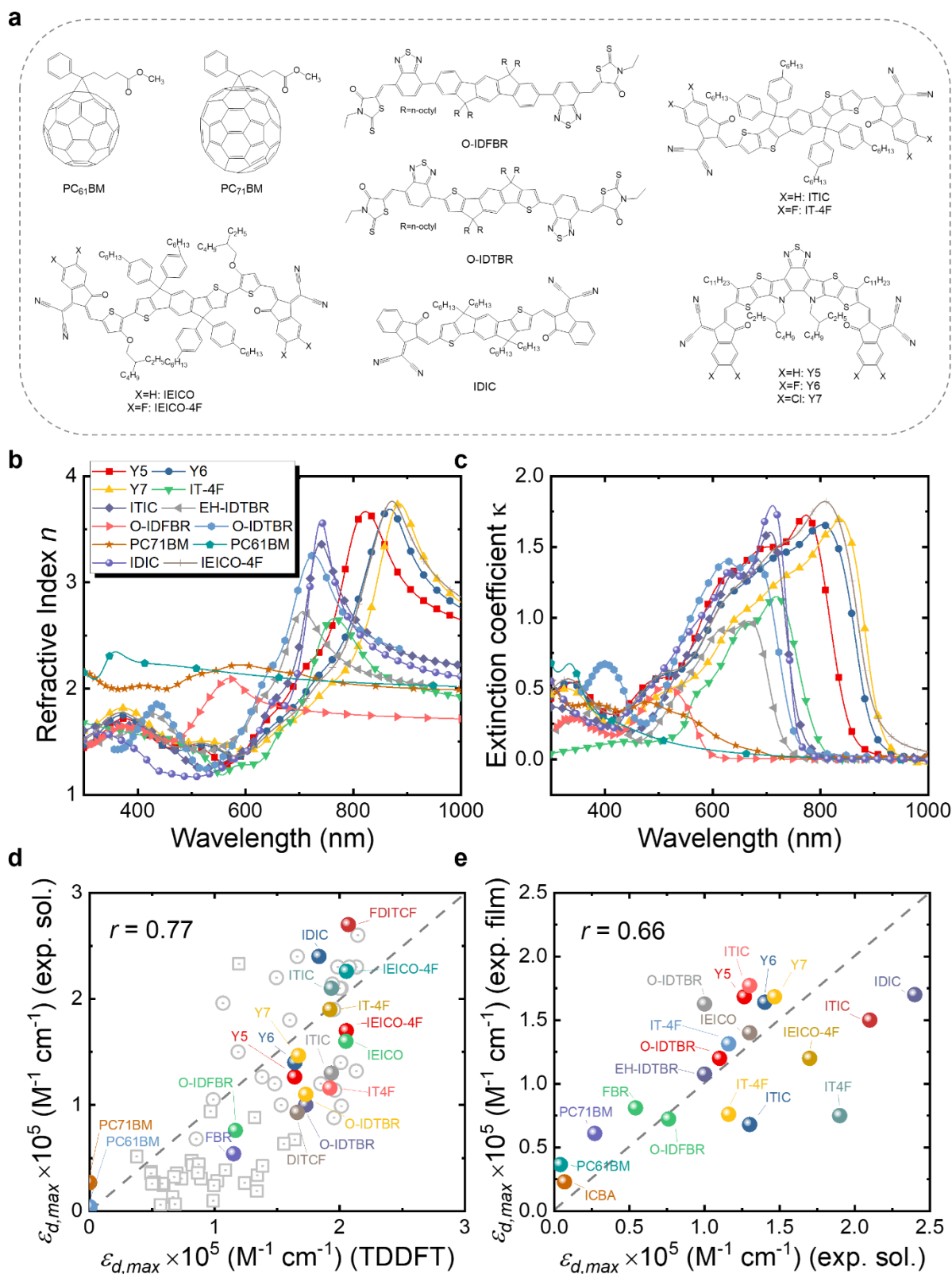
216 The experimental  $\varepsilon_{d,max}$  from solution can be obtained using the optical density (OD)  
217 measurements performed using UV-visible spectroscopy, via

218 
$$\varepsilon_{d,max} = \frac{OD_{max}}{\rho d}, \quad (\text{Equation 2})$$

219 where  $OD_{max}$  is the maximum optical density,  $\rho$  is the molar concentration (M), and  $d$  the light  
220 path length of the cuvette (cm). Similarly, the experimental  $\varepsilon_{d,max}$  from film can be estimated  
221 assuming a mass concentration  $\rho_M$  in the film of  $1000 \text{ g L}^{-1}$  (as a typical value for conjugated  
222 polymers and small molecules),<sup>29</sup> either from the maximum absorption coefficient  $\alpha_{cm,max}$   
223 ( $\text{cm}^{-1}$ ) or extinction coefficient ( $\kappa_{max}$ ) (**Figure 1c**), via

224 
$$\varepsilon_{d,max} = \log_{10}(e) \alpha_{cm,max} \frac{M_w}{\rho_M} = \log_{10}(e) \frac{4\pi\kappa_{max} M_w}{\lambda_{max} \rho_M}, \quad (\text{Equation 3})$$

225 where  $M_w$  is the molecular weight in  $\text{g mol}^{-1}$ , and  $\lambda_{max}$  the wavelength at  $\kappa_{max}$  in centimetre.



226

227 **Figure 1.** a) Molecular structures of typical organic acceptors, including PC61BM, PC71BM,  
 228 O-IDFBR, O-IDTBR, ITIC, IT-4F, IDIC, IEICO, IEICO-4F, Y5, Y6, and Y7. b) Refractive  
 229 index and c) extinction coefficient of a larger set of typical organic acceptor thin films measured  
 230 using VASE. d) Experimental  $\epsilon_{d,max}$  in solution versus calculated  $\epsilon_{d,max}$  in vacuum using

231 TDDFT of a set of ~80  $\pi$ -conjugated molecules. e) Estimated experimental  $\epsilon_{d,max}$  in film (solid  
232 state) versus that in solution using Eq. (3). Panel (d) contains a subset of well-known NFA  
233 molecules that are highlighted in colour. All TDDFT results in panel (d) were performed using  
234 the functional B3LYP and basis set 6-311+G(d,p), except for the ones (grey squares) taken from  
235 Ref. <sup>69</sup> that are based on the LRC-wPBEh functional and 6-311+G(d) basis set. We also note  
236 here that the side chains of molecules are replaced by H atoms or methyl groups in the  
237 calculations as they are computationally expensive and do not contribute to the  $\pi$ -conjugation,  
238 hence electronic transitions.<sup>29</sup> The experimental data of  $\epsilon_{d,max}$  in film are converted from  
239 maximum values of extinction coefficients shown in panel (c) using Eq. (3), while solution data  
240 are collected from literature, noting that different values may be present for the same material  
241 as retrieved from different sources. Grey dashed lines indicate the perfect match between x and  
242 y axis. The data required for generating panels (d) and (e) in this figure are presented in the  
243 Supplementary Database.

244 **Figure 1d** presents the results of the comparison between experimental  $\epsilon_{d,max}$  in solution and  
245 theoretical  $\epsilon_{d,max}$  calculated from single molecules using TDDFT in vacuum. A brief  
246 discussion of the solvent effect on the absorption strength and the reasons why we choose  
247 vacuum medium are provided in **Figure S2**. Despite the scattering of data points, we observe  
248 the occurrence of a monotonic relationship between solution and calculated  $\epsilon_{d,max}$  with a  
249 Pearson correlation coefficient ( $r$ ) of 0.77. Interestingly, such correlation is no longer observed  
250 when quantifying the absorption strength in terms of  $\alpha_{max}$  neither when adding further data  
251 points from literature on  $\pi$ -conjugated fluorophores to our statistical analysis (**Figure S3a**,  $r =$   
252 0.30), which is believed to be caused by the differences in molecular weight; in that case, only  
253  $\epsilon_{d,max}$  is found to follow a monotonic trend (**Figure S3b**). Some of the material assumptions  
254 on refractive index and density required to obtain  $\alpha_{max}$  values might be responsible for the  
255 observed mismatch. It is worth noting that, expectedly, the correlation between solid state (film)  
256 and solution ( $r = 0.66$ , **Figure 1e**) or calculated  $\epsilon_{d,max}$  ( $r = 0.61$ , **Figure S4**) is not as good  
257 as that from solution data versus calculated  $\epsilon_{d,max}$  ( $r = 0.77$ , **Figure 1d**, neither for  $\alpha_{max}$  as  
258 shown in **Figure S5**). Such discrepancy is attributed to solid-state effects such as the  
259 aggregation effects,<sup>15</sup> intermolecular orientation,<sup>70,71</sup> and side chain interactions,<sup>72</sup> which are  
260 not considered in single molecule excited state calculations.<sup>29</sup> The observed trend that a highly  
261 absorbing material in solution will produce highly absorbing films is, nonetheless, generally  
262 valid and thus solution data is relevant for devices. Since the NFAs analysed here have a rather  
263 similar number of  $\pi$ -electrons ( $n_\pi$ ), the corresponding  $\epsilon_{d,max}$  per  $\pi$ -electron (**Figure S6**) shows

264 a similar trend as that in **Figure 1d**, **Figure 1e**, and **Figure S4**. Despite the simplicity of single  
265 molecule excited state calculations, these data show that using TDDFT calculations of the  
266 excited state to deliver  $\epsilon_{d,max}$  can provide a reasonably good approximation to experimental  
267 measurements. Moreover, dealing with TDDFT calculations gives us room to correlate key  
268 molecular properties, such as molecular size and shape (aspect ratio), linearity, planarity,  
269 grafted side chain positions, or functional groups, to the absorption strength using molecular  
270 descriptors. These observations provide a foundation from molecular structures to identify the  
271 origin and further extend the high optical extinction of NFAs through chemical design rules, as  
272 we show in the upcoming sections.

## 273 **2.2. Statistical analysis of the TDDFT absorption strength dataset**

274 The experimental validation of the TDDFT calculations in NFAs supports the use of such  
275 results to build an extended database of optimized molecular geometries and excited state  
276 properties. The dataset is built by collecting thousands of molecular geometries generated over  
277 the last years in our group, making up a total of 3515 calculations on small molecules and  
278 oligomers. The distribution of number of atoms in a molecule is shown in **Figure S7** with a  
279 majority lying between 50 and 100 atoms. This database is sufficiently diverse to allow us to  
280 detect correlations and chemical/structural design rules that could explain and/or further  
281 enhance optical absorption in conjugated small molecules.

### 282 **2.2.1. Correlation analysis of molecular descriptors**

283 In the simplest statistical analysis of our TDDFT database, we look for correlations of the  
284 absorption strength with respect to a catalogue of molecular descriptors. First, as described in  
285 **Supplementary Note 2**, we filter the pristine TDDFT database by identifying duplicate  
286 molecules (in terms of molecular weight) and selecting the lowest energy conformer (i.e.,  
287 optimized geometries in the ground state) among them. As a result, the curated TDDFT  
288 database employed in this work consists of 479  $\pi$ -conjugated small molecules and oligomers  
289 with a distribution of moieties shown in **Figure S9**.

290 Then, we introduce several target features related with absorption strength, starting from the  
291 maximum oscillator strength of any calculated transition ( $f_{max}$ ); the maximum oscillator  
292 strength of any transition in the visible electromagnetic window (herein constrained between  
293 300-1200 nm or 1-4 eV for its relevance in solar energy harvesting applications) ( $f_{max,vis}$ ); and  
294 the sum of oscillator strengths of all transitions in the visible window,  $f_{sum,vis}$ . These three

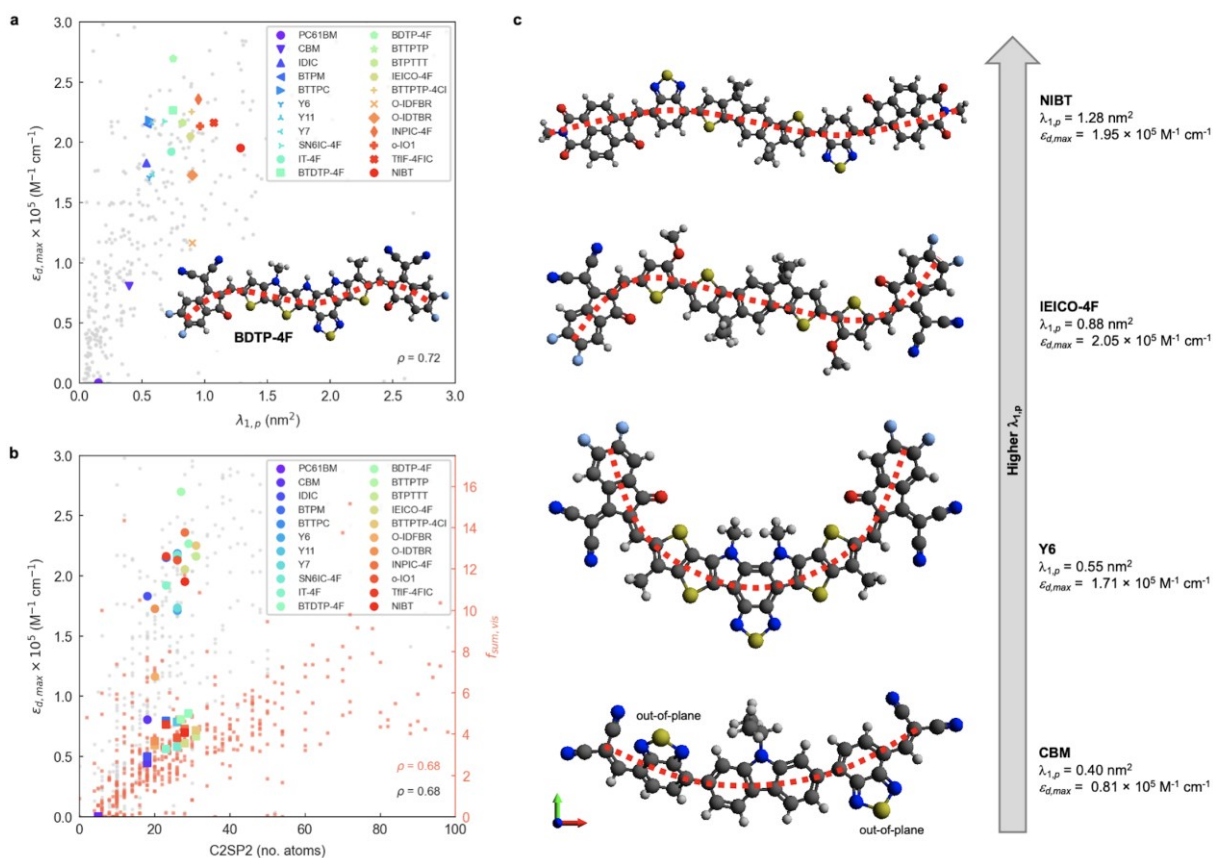
295 features are also evaluated per  $n_\pi$  for the molecule, i.e.,  $f_{max}/n_\pi$ ,  $f_{max,vis}/n_\pi$  and  $f_{sum,vis}/n_\pi$ .  
 296 We then consider the maximum absorption coefficient ( $\alpha_{max}$ ) obtained using Eq. (1) and Eq.  
 297 (3); the maximum of the imaginary part of the dielectric function ( $\varepsilon_{2,max}$ );<sup>29</sup> and  $\varepsilon_{d,max}$ . Finally,  
 298 we compute the spectral overlap between the OD ( $d\alpha(E)$ , where  $d$  is set to a typical film  
 299 thickness value of 100 nm and  $\alpha(E)$  derives from the Gaussian-broadened spectrum of  $f$  in the  
 300 visible spectral range taking a standard deviation of 0.1 eV) and the AM1.5G solar photon flux

301 spectrum ( $\Phi_{AM1.5G}$ ), namely  $f_{overlap} = \frac{\int_{1\text{ eV}}^{4\text{ eV}} \Phi_{AM1.5G}(E)d\alpha(E)dE}{\int_{1\text{ eV}}^{4\text{ eV}} \Phi_{AM1.5G}(E)dE}$ .

302 These features, together with their corresponding histograms (**Figure S14**) in terms of  
 303 Spearman's rank correlation coefficients ( $\rho$ ), are explained in more detail in **Supplementary**  
 304 **Note 2**. Molecular descriptors are calculated using up to four different open-source packages<sup>73-</sup>  
 305 <sup>76</sup> (**Supplementary Note 2**) to generate a (curated) collection of 3239 entries (including 40  
 306 electronic descriptors derived from the TDDFT calculations, namely the energy of the  
 307 molecular orbitals ranging from HOMO-19 to LUMO+19). Then, we scan for statistical  
 308 correlations between those descriptors and all target features introduced above, from which we  
 309 consider as highly correlated descriptors those showing  $\rho \geq 0.7$  as threshold. However, since  
 310 some descriptors are calculated in groups or families where weighting factors are varied among  
 311 atomic masses, van der Waals volumes, electronegativities, ionization potentials or  
 312 polarizabilities, we usually encounter sets of multicollinear descriptors that show very similar  
 313 trends with respect to the target feature. Accordingly, to drop redundant (collinear) descriptors  
 314 we classify them into clusters to select the most representative candidate of each bundle (i.e.,  
 315 cluster). This serves us to simplify the identification of characteristic and well-correlated  
 316 descriptors families. The clustering algorithm applied to analyse multicollinear descriptors  
 317 based on  $\rho$  and  $r$  values is further described in **Supplementary Note 3**.

318 After running the clusterization of descriptors on all target features, we identify strong  
 319 correlations with molecular descriptors for  $f_{max}$ ,  $f_{sum,vis}$  and  $\varepsilon_{d,max}$  (i.e., implicitly  $f_{max,vis}$ ).  
 320 For the remaining target variables ( $f_{overlap}$ ,  $\alpha_{max}$ ,  $\varepsilon_{2,max}$ ,  $f_{max}/n_\pi$ ,  $f_{max,vis}/n_\pi$  and  $f_{sum,vis}/$   
 321  $n_\pi$ ), we do not identify molecular descriptors with  $\rho$  above the threshold value (0.7) and they  
 322 are generally below 0.6 units, see **Figure S14**. The lack of correlation for  $f_{overlap}$  could be  
 323 justified by the existence of a gas-to-solid shift in the corresponding absorption spectrum, which  
 324 prevents proper matching of the Gaussian-broadened absorption features with the solar photon  
 325 flux. Regarding  $\alpha_{max}$  and  $\varepsilon_{2,max}$ , the estimation of these values from TDDFT calculations  
 326 requires taking generalized assumptions on several materials properties (such as density or

327 refractive index) that might be enough to disturb the underlying trends in our heterogeneous  
328 material database. For the quantities normalised by the number of pi electrons, i.e.  $f_{max}/n_{\pi}$ ,  
329  $f_{max,vis}/n_{\pi}$  and  $f_{sum,vis}/n_{\pi}$ , the weak correlation is expected since normalization tends to  
330 deviate from linear correlations depending on the straightness of the molecule.<sup>29</sup> Due to the  
331 strong correlation between size of the molecule and oscillator strength as discussed below based  
332 on C2SP2, the normalised quantity is believed to be a secondary factor, therefore not clear  
333 correlations are observed. In the successful correlation cases (i.e.  $f_{max}$ ,  $f_{sum,vis}$  and  $\epsilon_{d,max}$ ) and  
334 with the given thresholds of 0.7 units for  $\rho$  and  $r$ , we identify a single feature cluster lead by  
335 the  $\lambda_{1,p}$  descriptor in the case of  $f_{max}$  and  $\epsilon_{d,max}$  (**Figure 2a**). For  $f_{sum,vis}$ , a threshold  $\rho$  of 0.68  
336 reveals C2SP2 as a rather descriptive molecular feature (**Figure 2b**). Interestingly, C2SP2 is  
337 also found in the main cluster represented by  $\lambda_{1,p}$  in  $f_{max}$  and  $\epsilon_{d,max}$ , and we could not identify  
338 any strong correlations between the absorption strength (in any of its proposed metrics) and  
339 electronic descriptors (from HOMO-19 to LUMO+19 energy levels). Note that  $\epsilon_{d,max}$  values  
340 in excess of  $2.5 \times 10^5 \text{ M}^{-1} \text{ cm}^{-1}$  in **Figure 2a** and b are mostly attributed to artificially straight  
341 conjugated oligomers with >10 monomers contained in our database, for which the straightness,  
342 hence high  $\epsilon_{d,max}$ , are unlikely to be maintained in the experimental solid state scenario. In fact,  
343 only the exemplary and asymmetric NFA known as BDTP-4F (inset of **Figure 2a**)<sup>77,78</sup> surpasses  
344 that threshold with a record  $\epsilon_{d,max}$  in our NFA dataset ( $2.7 \times 10^5 \text{ M}^{-1} \text{ cm}^{-1}$ , and  $2.4 \times 10^5 \text{ M}^{-1}$   
345  $\text{cm}^{-1}$  measured in  $\text{CHCl}_3$  solution).<sup>77</sup>



346

347 **Figure 2.** (a) Correlation between  $\epsilon_{d,max}$ , as calculated from TDDFT, and  $\lambda_{1,\rho}$  as obtained in  
 348 the database of 479 molecules. The DFT-optimized geometry of BDTP-4F is shown in the inset.  
 349 (b) Correlation between  $\epsilon_{d,max}$  (and  $f_{sum,vis}$  in the secondary axis) and C2SP2 in that same  
 350 database. (c) DFT-optimized geometries of archetypal NFAs ordered by increased values of  
 351  $\lambda_{1,\rho}$  from bottom to top (CBM < Y6 < IEICO-4F < NIBT). Dotted red lines tentatively indicate  
 352 the overall curvature of the main conjugated backbone of the molecule.  $\lambda_{1,\rho}$  and C2SP2 describe  
 353 the size of the molecule in the direction of maximal atomic polarizability, and the number of  
 354 doubly bound carbon atoms ( $\text{sp}^2$  hybridized) bound to two other carbons (C2), respectively.

355  $\lambda_{1,\rho}$  is part of a bundle of three-dimensional molecular size and shape descriptors known as  
 356 weighted holistic invariant molecular (WHIM) descriptors.<sup>79–81</sup> These can be interpreted as a  
 357 generalized search for the principal axes with respect to a defined atomic property.<sup>82</sup> In this  
 358 particular case,  $\lambda_{1,\rho}$  is obtained by performing a principal component analysis (PCA) on the  
 359 centred atomic coordinates of the molecule using a covariance matrix ( $s_{jk}$ ) that is weighted by  
 360 the atomic polarizabilities ( $p_i$ ):

361 
$$s_{jk} = \frac{\sum_{i=1}^A p_i (q_{ij} - \bar{q}_j)(q_{ik} - \bar{q}_k)}{\sum_{i=1}^A p_i}, \quad (\text{Equation 4})$$

362 where  $s_{jk}$  is the weighted covariance between the  $j$ th and  $k$ th atomic coordinates;  $A$  is the total  
363 number of atoms;  $p_i$  is the (tabulated) polarizability of the  $i$ th atom;  $q_{ij}$  and  $q_{ik}$  represent the  
364  $j$ th and  $k$ th coordinate of the  $i$ th atom ( $j, k = x, y, z$ ), respectively; and  $\bar{q}$  is their average  
365 value.<sup>82</sup> After diagonalization of the polarizability-weighted covariance matrix, the first  
366 eigenvalue ( $\lambda_{1,p}$ ) quantifies the size of the molecule in the direction of maximal polarizability  
367 variance. Interestingly, the third eigenvalue ( $\lambda_{3,p}$ ) approaches zero in planar molecules as a  
368 result of absence of variance in the out-of-plane ( $z$ ) direction.<sup>82</sup> On the other hand, C2SP2,  
369 which is not in the WHIM group, accounts for the number of doubly bound carbon atoms ( $sp^2$   
370 hybridized, SP2) bound to two other carbons (C2), thus constituting a two-dimensional  
371 descriptor of fast computation. The correlation between C2SP2 and absorption strength can be  
372 relatively easier to understand, as C2SP2 to some extent represents the size of the conjugated  
373 molecule. Enlarging the size of the molecule increases the total number of  $\pi$ -electrons, which  
374 controls the total oscillator strength following the Thomas-Reiche-Kuhn rule. For the molecules  
375 that are extended along one direction, such as linear oligomers, increasing the size should  
376 enhance the oscillator strength of the first transition,<sup>29</sup> i.e. the dominant one.

377 To further interpret these two magnitudes ( $\lambda_{1,p}$  and C2SP2) as the main correlated descriptors  
378 with  $\varepsilon_{d,max}$  and  $f_{sum,vis}$ , we inspect the DFT-optimized geometries of archetypal NFAs  
379 (**Figure 2c**). The observed trend suggests that optical extinction monotonically increases with  
380  $\lambda_{1,p}$  (**Figure 2a**) in molecules having most of their polarizable atoms arranged along a main axis,  
381 i.e., linear molecules. While CBM shows large torsion angles mainly affecting the 2,1,3-  
382 benzothiadiazole (BT) moieties (thus making the molecule non-planar and increasing  $\lambda_{3,p}$ , see  
383 **Figure S15**), Y6 shows a characteristic curved geometry that limits its  $\varepsilon_{d,max}$  despite showing  
384 improved planarity. The NFA with the highest  $\lambda_{1,p}$  (NIBT) shows both linearity and planarity,  
385 with most of the more polarizable atoms (mainly C and S) lying along the principal polarizable  
386 axis of the molecule. Thus, in terms of molecular geometry, the absorption strength of NFAs  
387 could be further enhanced by distributing most of the atomic polarizability along a main axis  
388 while keeping good planarity and minimizing curvature. However,  $\lambda_{1,p}$  is not the sole molecular  
389 descriptor governing absorption strength, as BDTP-4F shows ca. 40% lower  $\lambda_{1,p}$  ( $0.75 \text{ nm}^2$ ) yet  
390 ca. 40% higher  $\varepsilon_{d,max}$  than NIBT (**Figure 2a**), which suggests that the molecular symmetry of  
391 NFAs could be another important factor affecting  $\varepsilon_{d,max}$ . Our preliminary investigations on  
392 this issue indicate that molecular asymmetry, as quantified by the WHIM symmetry index  $G_u$ ,  
393 might drive absorption strength higher (**Figure S16a**), yet we require a larger NFA database  
394 including more asymmetric molecules to further explore such an observation. Also, we



395 acknowledge that this observation might be biased by the systematic omission of side chains in  
396 the TDDFT calculations. By comparing  $\lambda_{1,p}$  in a selection of small molecule acceptors  
397 geometrically optimized with and without side chains (**Figure S17a**), we observe that in most  
398 cases the addition of side chains either decreases  $\lambda_{1,p}$  slightly or keeps it invariant. Still, the  
399 positive correlation of  $\lambda_{1,p}$  with respect to  $\epsilon_{d,max}$  is maintained (**Figure S17b**). Furthermore,  
400 the presence of naphthalene imide derivatives in the molecular structure of NIBT could be  
401 hindering further increase of the absorption strength with  $\lambda_{1,p}$ , as suggested by our statistical  
402 analysis of frequent moieties in the selection of good light harvesters (presented in the next  
403 section). On the other hand, an increase of  $n_\pi$  in the molecule in the form of closed-ring  
404 conjugated moieties will systematically increase C2SP2 and accordingly  $f_{sum,vis}$ . These  
405 findings support the previously known design rules in terms of molecular linearity and  $\pi$ -  
406 conjugation enabling large oscillator strength in organic small molecules and polymers, and are  
407 consistent with a recent study on chromophores.<sup>30</sup> In particular, trans- conjugated polymer  
408 stereoisomers are known to possess higher optical extinction due to their increased straightness  
409 and persistence length,<sup>29</sup> which agrees with our observations on exemplary curved (Y6) and  
410 more linear (NIBT) NFAs.

411 The energy of the first optical transition ( $E_1$ ) is also of practical importance in light harvesters  
412 such as NFAs as the lower energy part of the solar spectrum, down to  $\sim 1$  eV, contains a higher  
413 photon flux density. Our results show the number of heteroatoms in the molecule as the most  
414 correlated feature with  $E_1$  ( $\rho = -0.72$ , **Figure S18a**) while forming a single feature cluster, yet  
415 neither  $\lambda_{1,p}$  nor C2SP2 show strong correlations with  $E_1$ . This fact prevents the introduction of  
416 molecular design rules targeted at  $E_1$  using  $\lambda_{1,p}$  or C2SP2. However, we acknowledge a negative  
417 correlation between  $E_1$  and  $f_{osc,max}$  among common NFAs that suggests further room for  
418 absorption strength increase as  $E_1$  is reduced (**Figure S19**).

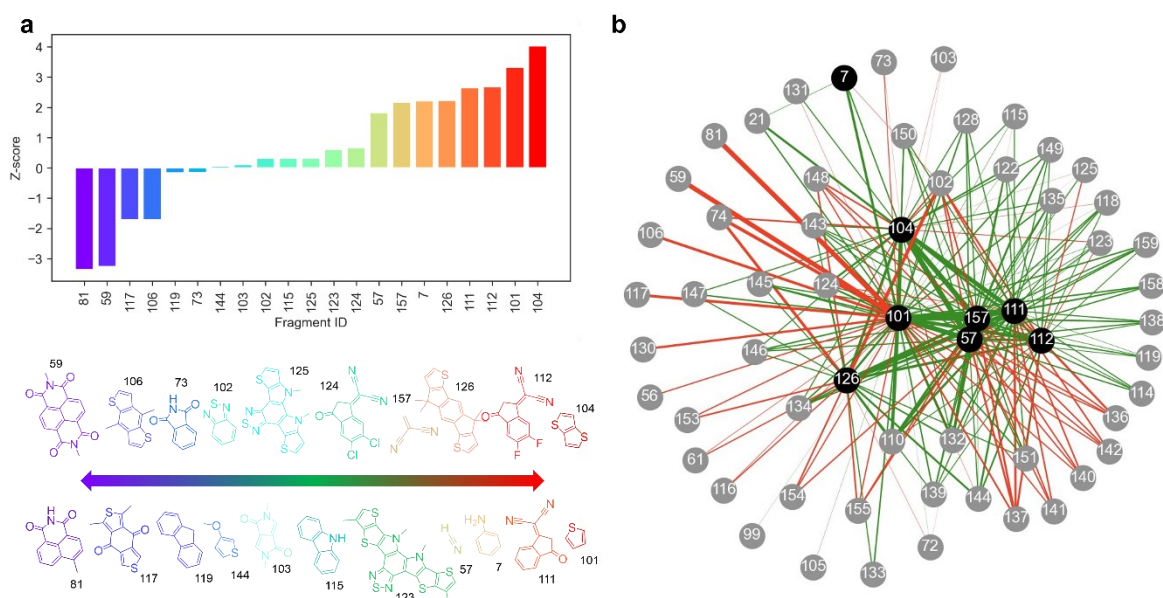
### 419 2.2.2. Chemical insights into highly absorbing molecules

420 Beyond molecular descriptors, we investigate the relationship between the choice of moieties  
421 and absorption strength to provide further material design rules for highly absorbing conjugated  
422 small molecules. Our objective is to identify overrepresented moieties in the subset of high-  
423 absorbing molecules (which we arbitrarily define as those having  $f_{osc,max} > 2.5$ , thus setting a  
424 population of size  $p$ ) with respect to the entire molecular dataset (population of size  $P$ ).  
425 Accordingly, we identify the molecular motifs present in the molecules by comparing their  
426 structures (as derived from SMILES notation) with those of a previously built database of

427 moieties (also SMILES-based). This database of moieties was partly inherited from a previous  
428 work<sup>42</sup> and extended with further motifs present in our particular dataset (see **Supplementary**  
429 **Note 4** and the spreadsheet included as Supplementary database). Afterwards, we consider that  
430 a discrete hypergeometric distribution is adequate to model our molecular dataset and the  
431 fragments found therein<sup>42</sup> to calculate the corresponding Z-scores as  $Z = (k - \bar{k})/\sigma_k$ , where  $k$   
432 is the number of high-absorbing molecules containing certain moiety;  $\bar{k}$  is its expected value,  
433 defined as  $pK/P$  where  $K$  corresponds to the number of molecules in the entire dataset  
434 containing that same moiety; and  $\sigma_k = \sqrt{pK(P - K)(P - p)/(P^2(P - 1))}$  is the standard  
435 deviation of the hypergeometric distribution. Z-scores will indicate (in units of  $\sigma_k$ ) which  
436 moieties are overrepresented or underrepresented in the subset of high-performing molecules  
437 with respect to the expected values when looking at the entire dataset. Our results (**Figure 3**)  
438 suggest that thieno[3,2-b]thiophene (TT), thiophene (T), 2-(5,6-difluoro-3-oxo-2,3-dihydro-  
439 1H-inden-1-ylidene)malononitrile (2FIC), 2-(3-oxo-2,3-dihydro-1H-inden-1-  
440 ylidene)malononitrile (IC), indaceno[1,2-b:5,6-b']dithiophene (IDT), 2-methylene  
441 malononitrile, cyanide, and aniline are particularly frequent in highly absorbing molecules.  
442 Interestingly, four of those molecular fragments (TT, T, 2FIC and IC) are contained in the  
443 chemical structure of the workhorse NFA Y6 (**Figure 2c**). Contrarily, naphthalene imide  
444 derivatives, as typically encountered in n-type small molecules and conjugated polymers;  
445 4H,8H-benzo[1,2-c:4,5-c']dithiophene-4,8-dione and benzo[1,2-b:4,5-b']dithiophene  
446 fragments are mostly underrepresented in the selection of high-performing light harvesters.

447 We further study the existing correlation between pairs of moieties to understand in which way  
448 the different molecular fragments should (or should not) be combined to retrieve highly-  
449 absorbing molecules. Our analysis starts by creating molecular subsets determined by the  
450 presence of a given moiety, which acts as source node (coloured in black) in the network graph  
451 shown in **Figure 3b**. Within that subset, we identify the high-absorbing molecules ( $f_{osc,max} >$   
452 2.5) and compute the Z-scores of their moieties (child nodes, coloured in grey in **Figure 3b**)  
453 with respect to the molecules of the entire molecular subset. As per the network shown in  
454 **Figure 3b**, the absolute Z-scores will determine the width of the edges connecting the nodes  
455 (moieties) and its sign the colour of the edge (green for positive Z-score [overrepresentation]  
456 and red for negative Z-score [underrepresentation]). Therefore, green and thick edges connect  
457 pairs of molecules that are more frequently found in high-absorbing molecules whereas thick  
458 and red edges indicate combinations of moieties that lead to less absorbing molecules. In this  
459 analysis, we set up 8 different source nodes corresponding to the most overrepresented moieties

460 observed in **Figure 3a**. As a result, **Figure 3b** can be interpreted as a catalogue of design rules  
461 relating pairs of moieties with high oscillator strength in  $\pi$ -conjugated small molecules.



462  
463 **Figure 3.** (a) Z-scores obtained from the discrete hypergeometric distribution of moieties in the  
464 highly-absorbing molecules ( $f_{osc,max} > 2.5$ ) with respect to the entire molecular dataset, for  
465 moieties activated at least 10 times. The corresponding structures of identified moieties are  
466 shown. (b) Network graph of Z-scores relating pairs of moieties. Source nodes are coloured in  
467 black whereas child nodes are coloured in grey. The colour of the edges corresponds to the sign  
468 of the Z-score (green for positive, red for negative). The width of the edges scales with the  
469 absolute value of the Z-score.

### 470 2.3. Machine-learning modelling of the absorption strength

471 Besides providing useful chemical insights from a material design perspective, molecular  
472 descriptors can be exploited to feed regression models and draw predictions on certain target  
473 features, forming the so-called quantitative structure-property relationship (QSPR) and  
474 quantitative structure-activity relationship (QSAR) models.<sup>79,80,82,83</sup> In the present study, we  
475 train and test several ML models fed with molecular and electronic descriptors obtained from  
476 TDDFT calculations to predict the value of  $\epsilon_{d,max}$  in conjugated small molecules and  
477 oligomers. Finally, we propose exploiting such ML model (trained with TDDFT data) to predict  
478  $\epsilon_{d,max}$  in molecules optimized using a semi-empirical quantum chemistry method, i.e. xTB.<sup>84</sup>  
479 This renders possible thanks to the geometrical similarity of the TDDFT and xTB ground state  
480 conformers, which lead to similar (geometrical) descriptors values; and the calibration of their

481 corresponding energy levels, as per the required inputs of the ML model herein employed.  
482 Therefore, further molecular candidates beyond the pristine dataset could be geometrically  
483 optimized using solely xTB Hamiltonians and their absorption strength predicted using such  
484 ML model. This approach effectively bypasses the use of TDDFT calculations when screening  
485 the absorption strength of novel molecules, which results in less demanding computations and  
486 higher throughput. The present ML workflow will open the possibility to accelerate the  
487 screening of high-performing molecular candidates with low-to-moderate computational  
488 requirements (further discussed in **Supplementary Note 5**).

### 489 2.3.1. Modelling $\varepsilon_{d,max}$ with random decision forests

490 From the analysis of descriptors shown in Section 2.2.1, we identified two main feature clusters  
491 represented by  $\lambda_{1,p}$  and C2SP2. We tentatively consider these two descriptors as independent  
492 variables in baseline models (such as 1-nearest neighbour and linear regression) targeted to  
493  $\varepsilon_{d,max}$ . For the model training and testing, we split our pristine dataset onto two subsets, namely  
494 the training set (gathering 70% of the data, randomly selected) and the testing set (gathering the  
495 remaining 30% of the data). Such baseline models are picked according to a recently introduced  
496 catalogue of good practices in the ML field,<sup>85</sup> to demonstrate the requirement of more advanced  
497 regressors (namely ML) in successful data modelling. The models are scored and quantitatively  
498 compared based on workhorse fitting metrics, such as their coefficient of determination ( $R^2$ );  
499 their adjusted coefficient of determination ( $R_{adj}^2$ , which adds penalties as the number of  
500 parameters increases, see **Supplementary Note 2**); and their Pearson correlation coefficient ( $r$ ),  
501 as retrieved in the training (fitting) and test sets. The inherent mathematical simplicity of the  
502 baseline models results in poor fitting scorings (**Figure S20** and **Table S1**) yet they suggest that  
503 feature selection procedures could end up in higher-performing models.

504 Accordingly, we deploy a state-of-the-art ML method, namely a RF, to aid in both aspects:  
505 feature selection and building of  $\varepsilon_{d,max}$  models of higher accuracy. RFs constitute one of the  
506 simplest and most widely applied ML methods in molecular screening and data mining  
507 studies.<sup>30,43,46,86</sup> They are particularly appealing for their straightforward implementation  
508 through open-source Python libraries such as Scikit-Learn,<sup>87</sup> and also for their inherent  
509 robustness against overfitting and fast optimization. RFs are formed by an ensemble of decision  
510 trees (estimators) that are executed in parallel and independently from each other. Decision  
511 trees serve to classify data by starting from a single root node that is subsequently divided into  
512 child nodes, the latter being chosen randomly among the input features. At every node splitting

513 step (i.e., decision making), the algorithm selects the pathway that minimizes the mean square  
514 error (MSE). Eventually, when every tree reaches its maximum extension (which is set  
515 arbitrarily via model hyperparameters), the predictions of all trees are averaged (ensembled),  
516 hence constituting the final predicted value of the RF. At this stage, myriad cross-validation  
517 (CV) techniques exist to evaluate the quality of the model and help in the tuning of  
518 hyperparameters. CV methods can estimate the ML model performance, evaluate potential  
519 over- or underfitting, and quantify how accurate the model is on drawing predictions on unseen  
520 data. In this work, we adopt two common cross-validation schemes, namely a repeated holdout  
521 CV; and a leave-one-out cross-validation (LOOCV). On the one hand, in a repeated holdout  
522 CV the pristine dataset is randomly split onto two distinct subsets, namely the training (here  
523 gathering 70% of the data) and testing (the remaining fraction of data, i.e. 30%) subsets. The  
524 model is trained and tested on the respective subsets, and the corresponding statistical metrics  
525 ( $R^2$ ,  $r$ , MSE, etc.) annotated. Eventually, the process is repeated  $k$  times (10-fold in this work),  
526 and all metrics are averaged to evaluate the ML model performance (its CV score). On the other  
527 hand, in a LOOCV the holdout process is taken to the extreme as the testing subset consists of  
528 a single data point while the remaining data is used in the training step. The process runs  
529 recursively for all data, thus eventually all data points are used for training and testing in the  
530 LOOCV protocol. Yet being computationally expensive, a LOOCV results in a more accurate  
531 estimate of model performance.

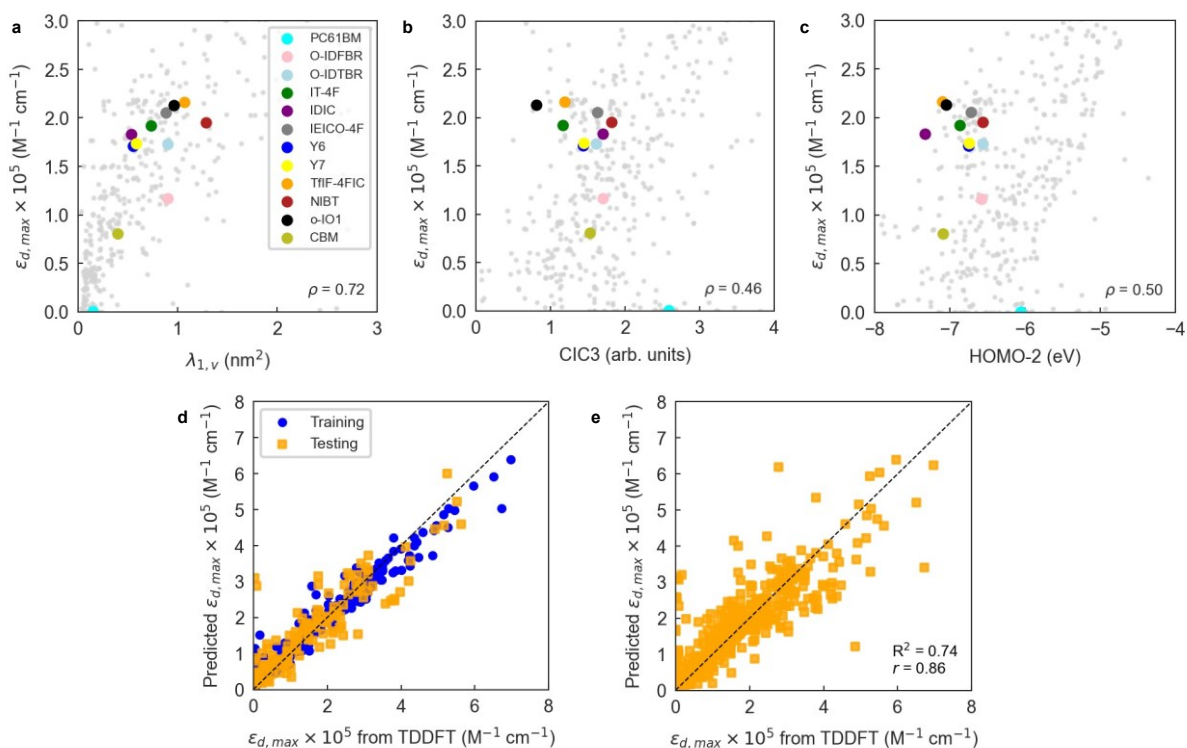
532 **Table S1** includes the performance of an out-of-the-box RF model trained and cross-validated  
533 using 300 trees (estimators). Exemplary comparisons between the two previous baseline models  
534 (1-nearest neighbor and linear regression) and the out-of-the-box RF model are found in **Figure**  
535 **S20**. The RF models indicate that scoring functions ( $R^2$ ,  $r$ ) well above 0.6-0.8 are feasible upon  
536 careful feature selection and further optimization of the RF regressor. Feature selection in RFs  
537 is usually performed by filtering variables based on their feature importance, which is a metric  
538 that accounts for how much a feature decreases the weighted variance in the node splitting steps  
539 of the decision trees. This property enables feature ranking to then apply myriad algorithms to  
540 filter out the least important variables as seen by the RF regressor. In this work, we perform a  
541 recursive feature elimination (RFE) procedure to the initial library of 3239 descriptors as  
542 described in **Supplementary Note 2**. In a RFE protocol, a significant fraction of the initial  
543 population of features is dropped in successive training steps of the RF ensemble. Features are  
544 dropped based on their corresponding feature importance until reaching an arbitrarily low  
545 number of input variables, hence simplifying the original model. Our RFE analysis shows that  
546 a threshold average  $R^2$  of 0.70 is achieved using a 12-variable model ( $R^2 = 0.70 \pm 0.05$ ,  $r =$

547  $0.84 \pm 0.03$ ), which outperforms the RF model presented earlier while including a drastic  
548 reduction in the number of variables (from 3239 to 12). The sweet spot in model accuracy and  
549 number of degrees of freedom is found for the 10-variable model, which shows the maximum  
550 average  $R_{adj}^2$  ( $0.67 \pm 0.06$ ).

551 Notably, a threshold  $R^2$  of 0.60 is already achieved training a 3-parameter RF model ( $R^2 =$   
552  $0.63 \pm 0.06$ ,  $R_{adj}^2 = 0.62 \pm 0.06$ ,  $r = 0.80 \pm 0.03$ ), which is particularly appealing given its  
553 simplicity. The resulting three-variable model includes one three-dimensional descriptor ( $\lambda_{1,v}$   
554 or WHIM\_45, as computed by the RDKit library, **Figure 4a**), one two-dimensional descriptor  
555 (CIC3, as computed by PaDEL software, **Figure 4b**) and one electronic descriptor, in this case  
556 the energy level of the second molecular orbital below the frontier HOMO (HOMO-2, **Figure**  
557 **4c**).  $\lambda_{1,v}$  refers to the first eigenvalue of the covariance matrix weighted by the atomic van der  
558 Waals volumes; thus,  $\lambda_{1,v}$  is included in the multicollinear feature cluster represented by  $\lambda_{1,p}$   
559 that we previously and statistically identified, showing nearly perfect correlation ( $r = 0.99$ ) with  
560  $\lambda_{1,p}$ . Accordingly,  $\lambda_{1,v}$  can be exchanged by  $\lambda_{1,p}$  without loss of performance in the RF model.  
561 This finding confirms that the linearity of the molecule (either quantified in terms of  
562 polarizabilities or van der Waals volumes) plays a key role in determining its absorption  
563 strength in the form of  $\epsilon_{d,max}$ . On the other hand, CIC3 is a graph-based, third-order  
564 neighbourhood symmetry index<sup>82</sup> which lacks a straightforward interpretation due to its  
565 mathematical complexity. We observe, however, that it linearly scales as  $\log_2 A$ , with A being  
566 the total number of vertices (atoms) in the graph (molecule)<sup>82</sup> thus likely reflecting the size of  
567  $\pi$ -conjugation as per the characteristics of our dataset. The interpretation of HOMO-2 as an  
568 important descriptor is more challenging, and it is not possible to substitute it by a different  
569 descriptor without a noticeable drop in the model performance (excepting HOMO-1, which  
570 shows  $r = 0.96$ ).

571 Interestingly, electronic descriptors (in particular) are required for the RF models to achieve  
572 their highest potential and scoring despite we have not observed strong correlations in our  
573 earlier statistical analysis. To probe it, we have performed the same RFE protocol yet skipping  
574 the set of electronic descriptors among the input features. Our results show that the top  
575 performing RF models (selecting 29 variables and getting  $R^2 = 0.58 \pm 0.06$ ,  $R_{adj}^2 =$   
576  $0.48 \pm 0.07$ ,  $r = 0.78 \pm 0.04$ ; or selecting 9 variables to obtain  $R_{adj}^2 = 0.52 \pm 0.06$ , see  
577 **Figure S13**) are yet behind the scorings recorded when the electronic descriptors are included  
578 in the list of features. Note that the performance without electronic descriptors is lower than the

579 3-parameter model that includes HOMO-2 as descriptor, highlighting its positive effect on the  
580 performance of the RF regressor.



581  
582 **Figure 4.** Correlation plots for  $\epsilon_{d,max}$  and the three most important descriptors retrieved by the  
583 RF model: (a)  $\lambda_{1,v}$ ; (b) CIC3; and (c) HOMO-2. (d) Holdout cross-validation run of a RF  
584 ensemble to predict  $\epsilon_{d,max}$ . 70% of the data is randomly selected for training and the remaining  
585 fraction is used for testing; the process is repeated 10 times and the statistical metrics averaged.  
586 The RF model is trained with three molecular descriptors ( $\lambda_{1,v}$ ; CIC3; and HOMO-2) and a  
587 Morgan fingerprint vector of 64 bits. (e) Leave-one-out cross-validation (LOOCV) of that same  
588 RF model using the optimized hyperparameter of 1200 estimators.

589 Molecular fingerprints have also been extensively exploited as input vectorial descriptors in  
590 statistical and ML models focused on feature prediction.<sup>42,88–90</sup> Molecular fingerprints are  
591 usually represented as bit activation vectors of arbitrary length and degree of complexity,  
592 representing the absence or presence of certain molecular (bonding) pattern, moiety, functional  
593 group, or atom. In this work, we exploit the RDKit library to generate moiety fingerprints,  
594 MACCS keys, Morgan fingerprints, path-based or topological fingerprints, E-state fingerprints,  
595 and Coulomb vectors. These fingerprints are quickly computed and serve to complement and  
596 improve the learning process of the ML models employed herein.

597 To better analyse the influence of the different fingerprint vectors in improving the RF scoring,  
598 we trained and cross-validated the 3-parameter RF model previously found in combination with  
599 all fingerprint vectors generated. The results shown in **Table S2** indicate that by adding a  
600 Morgan fingerprint vector of 64 bits to the initial set of input features the model performance  
601 can be substantially improved:  $R^2$  increases by 10% (relative), and  $r$  by another (relative) 5%  
602 (see **Figure 4d**). Therefore, Morgan fingerprints are particularly suitable to fine-tune the  
603 training and prediction accuracy of  $\epsilon_{d,max}$  in RF models although lacking of a straightforward  
604 physical interpretation. Additional refinement of the RF hyperparameters results in further  
605 improved models. We performed this optimization through a randomized search (in 350  
606 iterations) of the hyperparameters controlling the number of estimators in the RF, the minimum  
607 number of samples per leaf node and the minimum number of samples required to split an  
608 internal node, which constitute the main adjustable hyperparameters of the RF algorithm. These  
609 results are shown in **Table S3**, together with the scoring obtained in a rigorous LOOCV of the  
610 optimized RF model (**Figure 4e**). As an alternative ensemble of decision trees, we have also  
611 tested and optimized an Extra Trees (ET) regressor in Scikit-Learn. Its performance is, however,  
612 very close to that attained in the workhorse RF regressor (**Table S3** and **Figure S21**).

### 613 2.3.2. Bypassing TDDFT calculations through machine-learning and extended tight-binding

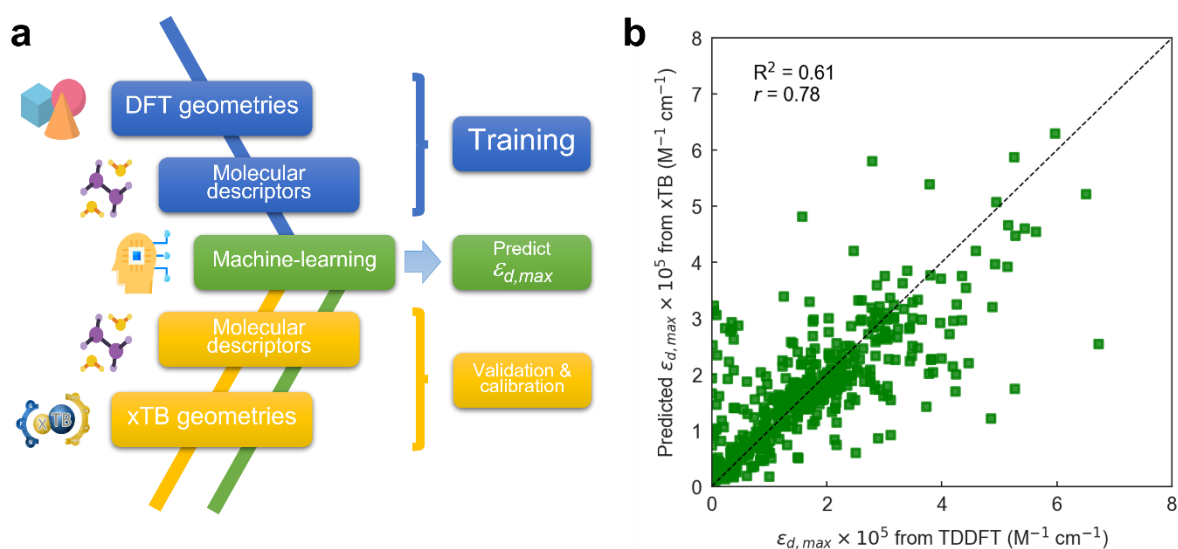
614 xTB Hamiltonians have recently emerged as semi-empirical and low computational cost  
615 quantum chemistry methods.<sup>84</sup> These have a remarkable potential in molecular screening when  
616 implemented in multilevel workflows where xTB is exploited first to identify plausible  
617 candidates using a minimal fraction of computational resources, to then leave room for higher-  
618 level DFT methods in selected candidates.<sup>84</sup> In this work, we propose exploiting a ML model  
619 trained with DFT data to predict  $\epsilon_{d,max}$  in molecular geometries optimized using xTB (**Figure**  
620 **5a**). This is expected to enable faster molecular screening and geometrical optimization steps,  
621 as both being entirely run using xTB Hamiltonians; followed by absorption strength ( $\epsilon_{d,max}$ )  
622 prediction in a TDDFT-trained RF model. Notably, our estimations show that the geometrical  
623 optimization step using GFN2-xTB is ca. 3000 times faster than using DFT with a hybrid  
624 functional (B3LYP/6-311+G(d,p)), as discussed in **Supplementary Note 5** and **Table S4**.

625 Nevertheless, the dissimilarity between xTB- and DFT-optimized molecular geometries might  
626 have a direct impact on the value of the (three-dimensional) molecular descriptors, and hence  
627 on the final accuracy of the interpolated ML model if some of those are included. Accordingly,  
628 we have first quantitatively compared both sets of molecular (non-electronic) descriptors by



629 computing  $r$  in all of them and found that the median of their distributions is very close to unity  
630 in all cases (**Figure S22**). Based on this finding, we proceed by training the RF model with  
631 TDDFT-derived descriptors and exploring how well the model interpolates when fed with xTB-  
632 derived descriptors. **Figure S23a** shows a leave-one-out interpolation of a RF model trained  
633 using TDDFT data and interpolated on GFN2-xTB-optimized molecules, descriptors and  
634 energy levels.<sup>84,91,92</sup> In this kind of model validation, all TDDFT data is used in the training step  
635 excepting that for a single molecule, for which we retrieve its corresponding xTB-optimized  
636 geometry and descriptors as the sole interpolation (testing) dataset; this procedure is  
637 subsequently repeated for all molecules. Thus, the model performance is assessed by comparing  
638 the actual TDDFT-derived  $\epsilon_{d,max}$  of the molecules (x-axis in **Figure 5b**) with that predicted by  
639 a RF model trained with TDDFT data and interpolated using xTB-derived descriptors (y-axis  
640 in **Figure 5b**). This is useful to evaluate whether such RF model fed with TDDFT data could  
641 be exploited to predict  $\epsilon_{d,max}$  in unseen molecules that are geometrically optimized through  
642 xTB Hamiltonians.

643 Our first model takes as inputs the three molecular descriptors found previously to be the most  
644 important features in the RF model together with their corresponding (64-bit) Morgan  
645 fingerprints. The scoring of the LOOCV in this preliminary model ( $R^2 = 0.53, r = 0.74$ ) is  
646 limited due to the existence of a mismatch between the absolute energy levels retrieved by either  
647 DFT (B3LYP) or GFN2-xTB methods (**Figure S23b**). Thus, the RF model trained on TDDFT  
648 data needs proper calibration of the energy levels obtained through GFN2-xTB, which we  
649 perform using either a linear regression, a support vector regressor (SVR) or an additional RF  
650 model (**Figure S23c**). By applying such calibration on the HOMO-2 energy levels, we obtain  
651 the champion RF model ( $R^2 = 0.61, r = 0.78$ ) shown in **Figure 5b** using three molecular  
652 descriptors and a 64-bit Morgan fingerprint vector. Hence, **Figure 5b** shows that molecular  
653 databases of xTB-optimized geometries could be exploited in combination with TDDFT-  
654 trained ML models to predict the absorption strength ( $\epsilon_{d,max}$ ) at significantly lower  
655 computational cost and with reasonable accuracy. The statistical analysis and ML modelling  
656 framework introduced here is thus expected to show large potential in the high-throughput  
657 screening of highly absorbing molecular candidates in combination with generative models  
658 (autoencoders and neural networks) as part of future work in the group.



659

660 **Figure 5.** (a) ML workflow used in this work to draw  $\epsilon_{d,max}$  predictions. A RF model is trained  
 661 on TDDFT data and interpolated (validated) on xTB geometries, including also their  
 662 corresponding molecular descriptors. To improve the accuracy of the model, energy levels  
 663 obtained using the GFN2–xTB Hamiltonian require calibration with TDDFT values (Figure  
 664 S23). (b) Leave-one-out interpolation of the resulting RF model using three input molecular  
 665 descriptors (including calibrated energy levels) and a 64-bit Morgan fingerprint vector.

### 666 3. Conclusion

667 We have demonstrated that TDDFT calculations agree reasonably well with the experimental  
 668 maximum molar extinction coefficient ( $\epsilon_{d,max}$ ) in solution state by exploiting a database of  
 669 TDDFT-optimized small molecular acceptors (NFAs) and donor oligomers collected over the  
 670 years. This finding supports further analysis of the molecular dataset to identify structure-  
 671 absorption relationships by means of statistical and machine-learning (ML) methods. Through  
 672 the exploration of molecular descriptors, we identify two features that are strongly correlated  
 673 with  $\epsilon_{d,max}$ , namely the linearity and planarity of the molecule in the direction of maximum  
 674 atomic polarizability variance; and the number of  $sp^2$ -hybridized carbon atoms bonded to two  
 675 other carbons included in the molecule. These further suggest design rules that highly absorbing  
 676 organic  $\pi$ -conjugated molecules (such as NFAs) should follow, namely a fully conjugated,  
 677 planar and linear molecular backbone with more polarisable heteroatoms. We further identify  
 678 that moieties such as thieno[3,2-b]thiophene (TT), thiophene (T), 2-(5,6-difluoro-3-oxo-2,3-  
 679 dihydro-1H-inden-1-ylidene)malononitrile (2FIC), 2-(3-oxo-2,3-dihydro-1H-inden-1-  
 680 ylidene)malononitrile (IC) and indaceno[1,2-b:5,6-b']dithiophene (IDT) appear more  
 681 frequently in molecules with the highest absorption strength. Finally, we demonstrate the

682 feasibility of random decision forests (RFs) trained with a few (3) molecular descriptors and  
683 64-bit Morgan fingerprint vectors to predict  $\epsilon_{d,max}$  in molecular geometries optimized by a  
684 computationally less demanding method such as extended tight-binding (xTB). This approach  
685 shows the ability to bypass thorough TDDFT calculations, thus facilitating high-throughput  
686 screening of absorption strength in organic  $\pi$ -conjugated molecules in combination with  
687 generative molecular models.

#### 688 **4. Outlook**

689 This work was motivated by the search for molecular design rules to enable higher PCE in  
690 organic solar cells. Although maximizing light absorption for a given optical band gap is a key  
691 requirement to enable record PCE, many additional physical processes contribute to  
692 photovoltaic performance but are not considered directly in the present work, namely, exciton  
693 diffusion, charge transfer, charge separation, charge transport and charge recombination. To  
694 date, there is no holistic modelling framework nor are there sufficient data to relate these  
695 multiple processes to device performance via chemical structure. However, developments in AI  
696 and ML methods are likely to advance the status of models for multiple property - device  
697 performance relationships in the coming years.

698 Nevertheless, understanding how light harvesting alone can be maximized by smart molecular  
699 design is significant for improving several different aspects of OPV performance. Light  
700 absorption is the primary step towards charge generation and is therefore strongly related to the  
701 macroscopic short-circuit current density of the device. According to the reciprocity relation  
702 between absorption and emission,<sup>20</sup> high absorption should in principle lead to strong emission,  
703 therefore reducing the nonradiative energy losses, and benefitting the open-circuit voltage. In  
704 addition, high absorption allows the fabrication of thin devices, therefore facilitating charge  
705 extraction and enhancing fill factor.<sup>93</sup> Moreover, based on the causality principle, high  
706 absorption strength would lead to higher refractive index, which takes the first interference  
707 maximum of electric field to lower thicknesses, resulting in large light harvesting potential in  
708 thinner devices. Therefore, designing highly absorbing organic  $\pi$ -conjugated molecules has the  
709 potential to enhance different aspects relating to the performance of OPVs in conjunction with  
710 the proposed predictive ML model.

711 A separate aspect for future work is the impact of solid-state molecular interactions on light  
712 absorption. This paper concerns the optical absorption of isolated molecules while applications  
713 normally require thin films of molecules. Although intermolecular interactions can strongly

714 impact the strength as well as the spectrum of thin film absorption,<sup>94</sup> this has been neglected in  
715 the present study due to the lack of a suitable database of computations and the lack of solid  
716 state packing information. In the future, ML approaches could be used to better understand and  
717 predict how solid state interactions affect optical absorption, and thereby improve molecular  
718 design rules. Such advances may be enabled by the growing capability in computational  
719 structure prediction as well as improved understanding of the impact of intermolecular  
720 interactions on excited state properties.

721

722

## 723 5. Experimental and theoretical methods

724 Excited state calculation database and experimental  $\epsilon_{d,\max}$  database: TDDFT results in this  
725 study are based on the functional B3LYP and were performed by present and past group  
726 members in Prof. Jenny Nelson's group at Imperial College London, making up more than 3500  
727 entries (corresponding to 479 unique molecules). The majority of experimental solid state thin  
728 film  $\epsilon_{d,\max}$  values for NFAs shown in **Figure 1a**, b and c were measured using variable-angle  
729 spectroscopic ellipsometry (VASE) for the present study. Neat films were deposited from  
730 solution by either spin- or blade-coating on glass substrates at distinct thicknesses (typically  
731 ranging from 30 to 150 nm). Ellipsometry data were acquired at three to five angles of incidence  
732 ( $55^\circ$ - $75^\circ$ ) using a Sopralab GES-5E rotating polarizer spectroscopic ellipsometer (SEMILAB)  
733 coupled to a charge-coupled device (CCD) detector. Experimental solution  $\epsilon_{d,\max}$  were mostly  
734 collected from literature with a majority of data taken from Ref. <sup>95</sup>, and Y5, Y6, and Y7  
735 measured using UV-visible spectroscopy. The complete database and sources are presented in  
736 the Supplementary Database.

737 Theoretical description of molar extinction coefficient ( $\epsilon_d$ ): To calculate the molar extinction  
738 coefficient  $\epsilon_d$ , let us start with defining the absorption coefficient  $\alpha$  in a quantum picture (we  
739 stay with SI units for the moment). The absorption coefficient for transition from state 1 to state  
740 2 can be defined as <sup>19,96</sup>

$$741 \quad \frac{dI}{dx} = -\alpha_{12}I, \quad (\text{Equation S1})$$

742 where  $I$  is light intensity, determined by the energy density of an electromagnetic wave via

$$743 \quad I = \frac{1}{2}nc\epsilon_0|E_0|^2, \quad (\text{Equation S2})$$

744 where  $n$  is the refractive index,  $\epsilon_0$  vacuum permittivity,  $c$  the speed of light, and  $E_0$  the  
745 amplitude of the electric field. For an electromagnetic wave, the rate of intensity attenuation  $\frac{dI}{dx}$   
746 is equal to the rate of loss of energy density from the field  $-\frac{dU}{dt}$ , and the latter is the product of  
747 transition rate  $\Gamma_{12}$  and transition energy  $\hbar\omega_{12}$ , and we have

$$748 \quad \frac{dI}{dx} = -N\Gamma_{12}\hbar\omega_{12}, \quad (\text{Equation S3})$$

749 where  $N$  is the volume density of molecules and  $\hbar$  the reduced Planck constant. Substituting  
 750 for  $\frac{dI}{dx}$  and  $I$  in the definition of  $\alpha_{12}$  we get

$$751 \quad \alpha_{12} = \frac{2N\hbar\omega_{12}\Gamma_{12}}{nc\epsilon_0|E_0|^2} \quad (\text{Equation S4})$$

752 The transition rate  $\Gamma_{12}$  can be defined by Fermi's Golden Rule and the perturbing Hamiltonian  
 753 given by  $H = d_{12}\mathbf{E}_0$  using dipole approximation, where  $d_{12}$  is the transition dipole moment of  
 754 the transition. Considering randomly oriented transition dipoles relative to the direction of the  
 755 exciting electromagnetic field, we have

$$756 \quad \Gamma_{12} = \frac{2\pi}{3\hbar} d_{12}^2 |E_0|^2 \delta(\hbar\omega - E_2 + E_1) \quad (\text{Equation S5})$$

757 Using  $E_2 - E_1 = \hbar\omega_{12}$ , we get

$$758 \quad \alpha_{12} = \frac{4\pi N\omega_{12}}{3nc\epsilon_0\hbar} d_{12}^2 \delta(\omega - \omega_{12}) \quad (\text{Equation S6})$$

759 From an arbitrary transition from state  $i$  to state  $j$ , we can express above equation using oscillator  
 760 strength of the transition ( $f_{ij}$ ):

$$761 \quad \alpha_{ij} = \frac{2\pi Ne^2}{3\epsilon_0 m_0 nc} f_{ij} \delta(\omega - \omega_{ij}), \quad (\text{Equation S7})$$

762 where  $e$  is the elementary charge, and  $f_{ij} = \frac{2m_0\omega_{12}}{e^2\hbar} d_{ij}^2$ . Integrating over all transitions, we have

$$763 \quad \alpha(\omega) = \frac{2\pi Ne^2}{3\epsilon_0 m_0 nc} \sum_{ij} f_{ij} \delta(\omega - \omega_{ij}) \quad (\text{Equation S8})$$

764 To correlate the absorption coefficient ( $\alpha$ ) with the molar extinction coefficient ( $\epsilon_d$ ), we need  
 765 the definition of optical density ( $OD$ ) and optical depth ( $\alpha d$ ). Light is attenuated by passing  
 766 through a depth  $d$  of material such that

$$767 \quad I(d) = I_0 e^{-\alpha d} = I_0 10^{-OD} \quad (\text{Equation S9})$$

768 And optical density, or called sometimes absorbance is defined as

$$769 \quad OD = \rho \epsilon_d d, \quad (\text{Equation S10})$$

770 where  $\rho$  is concentration in molar (M or mol L<sup>-1</sup>), and  $d$  is sample length in cm. Consequently,  
 771 we have

772 
$$\varepsilon_d = \frac{\log_{10}(e)}{\rho} \alpha_{cm} = \frac{\alpha_{cm}}{2.303\rho}, \quad (\text{Equation S11})$$

773 noting that we now write the absorption coefficient per cm to distinguish from the expression  
 774 for  $\alpha$  above, which we did assuming SI units, hence  $\alpha_{cm} = \frac{\alpha}{100}$ .  $\rho$  is moles of molecules per  
 775  $\text{dm}^3$ . We now have

776 
$$\varepsilon_d(\omega) = 10 \log_{10}(e) N_A \frac{2\pi e^2}{3\varepsilon_0 m_0 n c} \sum_{ij} f_{ij} \delta(\omega - \omega_{ij}) \quad (\text{Equation S12})$$

777 Let us recast this in terms of photon energy  $E$  in eV, i.e.  $E = \frac{\hbar\omega}{e}$ , rather than angular frequency,  
 778 so it is easier to consider the magnitude, and finally we have  $\varepsilon_d$  in the unit of  $\text{M}^{-1} \text{cm}^{-1}$ .

779 
$$\varepsilon_d(E) = 10 \log_{10}(e) N_A \frac{2\pi e \hbar}{3\varepsilon_0 m_0 n c} \sum_{ij} f_{ij} \delta(E - E_{ij}) \quad (\text{Equation S13})$$

780 This allows us to compute the theoretical  $\varepsilon_d$  using the calculated oscillator strength at different  
 781 transitions. And the common method to calculate the oscillator strength is time-dependent-  
 782 density-functional-theory (aka TDDFT).

783 Converting complex refractive index from solid state ellipsometry measurements to  $\varepsilon_d$ : Using  
 784 ellipsometry measurements from film (solid state), we can extract the complex refractive index,  
 785  $\eta$

786 
$$\eta = n + i\kappa \quad (\text{Equation S14})$$

787 Where  $n$  is the refractive index, and  $\kappa$  the extinction coefficient. The absorption coefficient  
 788 ( $\alpha_{cm}$ ) is then determined by

789 
$$\alpha_{cm} = \frac{4\pi\kappa}{\lambda_{cm}} \quad (\text{Equation S15})$$

790 Where  $\lambda_{cm}$  is the wavelength in centimetre. Using Equation S11, and the relationship between  
 791 molar concentration  $\rho$  and mass concentration  $\rho_M$ , i.e.,  $\rho = \frac{\rho_M}{M_w}$ , we have

792 
$$\varepsilon_d = \log_{10}(e) \alpha_{cm} \frac{M_w}{\rho_M} = \log_{10}(e) \frac{4\pi\kappa M_w}{\lambda_{cm} \rho_M} \quad (\text{Equation S16})$$

793 Where  $M_w$  is the molecular weight,  $\text{g mol}^{-1}$ , and  $\rho_M$  has the unit of  $\text{g L}^{-1}$ , and is typically  
 794 assumed to be  $1000 \text{ g L}^{-1}$ .

795 **Author contributions**

796 J.Y. and X.R.-M. contributed equally to this work and drafted the paper. J.Y. performed DFT  
797 and TDDFT calculations, absorption strength analysis, and data collection. X.R.-M. performed  
798 the statistical analysis and machine-learning study. D.P., H.D., D.B., M.A., A.V., S.F., A.A.S.,  
799 and X.H. shared their DFT/TDDFT calculation results. F.E. prepared thin films of NFAs for  
800 VASE measurements. X.R.-M., V.B., and B.D. did VASE measurements. E.R. did UV-vis  
801 measurements of Y5, Y6, and Y7 in solution. G.Z. and H.-L.Y. provided Y5, Y6, and Y7. All  
802 authors gave critical review on this work. J.N. and M.C.-Q. supervised this work.

803 **Conflicts of interest**

804 There are no conflicts to declare.

805 **Acknowledgements**

806 J.N., J.Y., D.P., M.A., F.E., and E.R. thank the European Research Council for support under  
807 the European Union's Horizon 2020 research and innovation program (Grant Agreement No.  
808 742708 and No. 648901). The authors at ICMAB acknowledge financial support from the  
809 Spanish Ministry of Science and Innovation through the Severo Ochoa" Program for Centers  
810 of Excellence in R&D (No. CEX2019-000917-S), and project PGC2018-095411-B-I00. E.R.  
811 is grateful to the Fonds de Recherche du Quebec-Nature et technologies (FRQNT) for a  
812 postdoctoral fellowship and acknowledges financial support from the European Cooperation in  
813 Science and Technology. M.A. thanks the Engineering and Physical Sciences Research Council  
814 (EPSRC) for support via doctoral studentships. F.E. thanks the Engineering and Physical  
815 Sciences Research Council (EPSRC) for support via the Post-Doctoral Prize Fellowship. X.R.-  
816 M. acknowledges Prof. Olle Inganäs and the Knut and Alice Wallenberg Foundation for  
817 funding of his current postdoctoral position. H.-L. Yip thanks the support from Guangdong  
818 Major Project of Basic and Applied Basic Research (2019B030302007). The TOC figure and  
819 Figure 5a in the manuscript include freely available resources from Flaticon.com. J.Y. thank  
820 Xiaodan Ge for her support.



821 **References**

- 822 1 J. Nelson, *Mater. Today*, 2011, **14**, 462–470.
- 823 2 G. Li, R. Zhu and Y. Yang, *Nat. Photonics*, 2012, **6**, 153–161.
- 824 3 A. J. Heeger, *Adv. Mater.*, 2014, **26**, 10–28.
- 825 4 M. Mainville and M. Leclerc, *ACS Energy Lett.*, 2020, **5**, 1186–1197.
- 826 5 H. K. H. Lee, J. Wu, J. Barbé, S. M. Jain, S. Wood, E. M. Speller, Z. Li, F. A. Castro,  
827 J. R. Durrant and W. C. Tsoi, *J. Mater. Chem. A*, 2018, **6**, 5618–5626.
- 828 6 Y. Cui, Y. Wang, J. Bergqvist, H. Yao, Y. Xu, B. Gao, C. Yang, S. Zhang, O. Inganäs,  
829 F. Gao and J. Hou, *Nat. Energy*, 2019, **4**, 768–775.
- 830 7 Y. Li, J. D. Lin, X. Che, Y. Qu, F. Liu, L. S. Liao and S. R. Forrest, *J. Am. Chem. Soc.*,  
831 2017, **139**, 17114–17119.
- 832 8 S. Difley and T. Van Voorhis, *J. Chem. Theory Comput.*, 2011, **7**, 594–601.
- 833 9 C. J. M. Emmott, J. A. Röhr, Mariano Campoy-Quiles, Thomas Kirchartz,  
834 Antonio Urbina, N. J. Ekins-Daukes and Jenny Nelson, *Energy Environ. Sci.*, 2015, **8**,  
835 1317–1328.
- 836 10 C. J. M. Emmott, D. Moia, P. Sandwell, N. Ekins-Daukes, M. Hösel, L. Lukoschek, C.  
837 Amarasinghe, F. C. Krebs and J. Nelson, *Sol. Energy Mater. Sol. Cells*, 2016, **149**,  
838 284–293.
- 839 11 L. Zhu, M. Zhang, J. Xu, C. Li, J. Yan, G. Zhou, W. Zhong, T. Hao, J. Song, X. Xue,  
840 Z. Zhou, R. Zeng, H. Zhu, C.-C. Chen, R. C. I. MacKenzie, Y. Zou, J. Nelson, Y.  
841 Zhang, Y. Sun and F. Liu, *Nat. Mater.*, 2022, 1–8.
- 842 12 J. Zhao, Y. Li, G. Yang, K. Jiang, H. Lin, H. Ade, W. Ma and H. Yan, *Nat. Energy*,  
843 2016, **1**, 15027.
- 844 13 P. Cheng, G. Li, X. Zhan and Y. Yang, *Nat. Photonics*, 2018, **12**, 131–142.
- 845 14 Y. Wang, J. Lee, X. Hou, C. Labanti, J. Yan, E. Mazzolini, A. Parhar, J. Nelson, J. Kim  
846 and Z. Li, *Adv. Energy Mater.*, 2021, **11**, 2003002.

- 847 15 J. Hou, O. Inganäs, R. H. Friend and F. Gao, *Nat. Mater.*, 2018, **17**, 119–128.
- 848 16 J. Yuan, Y. Zhang, L. Zhou, G. Zhang, H.-L. Yip, T.-K. Lau, X. Lu, C. Zhu, H. Peng,  
849 P. A. Johnson, M. Leclerc, Y. Cao, J. Ulanski, Y. Li and Y. Zou, *Joule*, 2019, **3**, 1140–  
850 1151.
- 851 17 M. Kaltenbrunner, M. S. White, E. D. Głowacki, T. Sekitani, T. Someya, N. S.  
852 Sariciftci and S. Bauer, *Nat. Commun.* 2012 31, 2012, **3**, 1–7.
- 853 18 W. Yang, W. Wang, Y. Wang, R. Sun, J. Guo, H. Li, M. Shi, J. Guo, Y. Wu, T. Wang,  
854 G. Lu, C. J. Brabec, Y. Li and J. Min, *Joule*, 2021, **5**, 1209–1230.
- 855 19 J. Nelson, *The Physics of Solar Cells*, Imperial College Press, 2003.
- 856 20 U. Rau, *Phys. Rev. B*, 2007, **76**, 085303.
- 857 21 M. Azzouzi, J. Yan, T. Kirchartz, K. Liu, J. Wang, H. Wu and J. Nelson, *Phys. Rev. X*,  
858 2018, **8**, 031055.
- 859 22 F. D. Eisner, M. Azzouzi, Z. Fei, X. Hou, T. D. Anthopoulos, T. J. S. J. S. Dennis, M.  
860 Heeney and J. Nelson, *J. Am. Chem. Soc.*, 2019, **141**, 6362–6374.
- 861 23 J. Yan, E. Rezasoltani, M. Azzouzi, F. Eisner and J. Nelson, *Nat. Commun.*, 2021, **12**,  
862 3642.
- 863 24 A. Classen, C. L. Chochos, L. Lüer, V. G. Gregoriou, J. Wortmann, A. Osvet, K.  
864 Forberich, I. McCulloch, T. Heumüller and C. J. Brabec, *Nat. Energy*, 2020, **5**, 711–  
865 719.
- 866 25 X.-K. Chen, D. Qian, Y. Wang, T. Kirchartz, W. Tress, H. Yao, J. Yuan, M. Hülsbeck,  
867 M. Zhang, Y. Zou, Y. Sun, Y. Li, J. Hou, O. Inganäs, V. Coropceanu, J.-L. Bredas and  
868 F. Gao, *Nat. Energy*, 2021, **6**, 799–806.
- 869 26 J. Benduhn, K. Tvingstedt, F. Piersimoni, S. Ullbrich, Y. Fan, M. Tropicano, K. A. A.  
870 McGarry, O. Zeika, M. K. K. Riede, C. J. J. Douglas, S. Barlow, S. R. R. Marder, D.  
871 Neher, D. Spoltore and K. Vandewal, *Nat. Energy*, 2017, **2**, 17053.
- 872 27 X.-K. Chen, V. Coropceanu, J.-L. Brédas and J.-L. Brédas, *Nat. Commun.*, 2018, **9**,  
873 5295.

- 874 28 D. Qian, Z. Zheng, H. Yao, W. Tress, T. R. Hopper, S. S. Chen, S. Li, J. Liu, S. S.  
875 Chen, J. Zhang, X.-K. K. Liu, B. Gao, L. Ouyang, Y. Jin, G. Pozina, I. A. Buyanova,  
876 W. M. Chen, O. Inganäs, V. Coropceanu, J.-L. L. Bredas, H. Yan, J. Hou, F. Zhang, A.  
877 A. Bakulin and F. Gao, *Nat. Mater.*, 2018, **17**, 703–709.
- 878 29 M. S. Vezie, S. Few, I. Meager, G. Pieridou, B. Dörling, R. S. Ashraf, A. R. Goñi, H.  
879 Bronstein, I. McCulloch, S. C. Hayes, M. Campoy-Quiles and J. Nelson, *Nat. Mater.*,  
880 2016, **15**, 746–753.
- 881 30 B. Kang, C. Seok and J. Lee, *J. Chem. Inf. Model.*, 2020, **60**, 5984–5994.
- 882 31 S. Few, J. M. Frost, J. Kirkpatrick and J. Nelson, *J. Phys. Chem. C*, 2014, **118**, 8253–  
883 8261.
- 884 32 Y. Yi, V. Coropceanu and J.-L. Brédas, *J. Mater. Chem.*, 2011, **21**, 1479.
- 885 33 T. Liu and A. Troisi, *J. Phys. Chem. C*, 2011, **115**, 2406–2415.
- 886 34 J. C. Slater, *Phys. Rev.*, 1951, **81**, 385.
- 887 35 B. L. Hammond, W. A. Lester and P. J. Reynolds, *Monte Carlo Methods in Ab Initio*  
888 *Quantum Chemistry*, WORLD SCIENTIFIC, 1994, vol. 1.
- 889 36 R. A. Friesner, *Proc. Natl. Acad. Sci.*, 2005, **102**, 6648–6653.
- 890 37 Y. Shao, L. F. Molnar, Y. Jung, J. Kussmann, C. Ochsenfeld, S. T. Brown, A. T. B.  
891 Gilbert, L. V. Slipchenko, S. V. Levchenko, D. P. O’Neill, R. A. DiStasio Jr, R. C.  
892 Lochan, T. Wang, G. J. O. Beran, N. A. Besley, J. M. Herbert, C. Yeh Lin, T. Van  
893 Voorhis, S. Hung Chien, A. Sodt, R. P. Steele, V. A. Rassolov, P. E. Maslen, P. P.  
894 Korambath, R. D. Adamson, B. Austin, J. Baker, E. F. C. Byrd, H. Dachsel, R. J.  
895 Doerksen, A. Dreuw, B. D. Dunietz, A. D. Dutoi, T. R. Furlani, S. R. Gwaltney, A.  
896 Heyden, S. Hirata, C.-P. Hsu, G. Kedziora, R. Z. Khalliulin, P. Klunzinger, A. M. Lee,  
897 M. S. Lee, W. Liang, I. Lotan, N. Nair, B. Peters, E. I. Proynov, P. A. Pieniazek, Y.  
898 Min Rhee, J. Ritchie, E. Rosta, C. David Sherrill, A. C. Simmonett, J. E. Subotnik, H.  
899 Lee Woodcock III, W. Zhang, A. T. Bell, A. K. Chakraborty, D. M. Chipman, F. J.  
900 Keil, A. Warshel, W. J. Hehre, H. F. Schaefer III, J. Kong, A. I. Krylov, P. M. W. Gill,  
901 M. Head-Gordon, Yihan Shao, L. Fusti Molnar, Yousung Jung, Jörg Kussmann,  
902 Christian Ochsenfeld, S. T. Brown, A. T.B. Gilbert, L. V. Slipchenko, S.

- 903 V. Levchenko, D. P. O'Neill, R. A. D. Jr, R. C. Lochan, Tao Wang, G. J.O. Beran, N.  
904 A. Besley, J. M. Herbert, C. Y. Lin, T. V. Voorhis, S. H. Chien, Alex Sodt, R.  
905 P. Steele, V. A. Rassolov, P. E. Maslen, P. P. Korambath, R. D. Adamson,  
906 Brian Austin, Jon Baker, E. F. C. Byrd, Holger Dachsel, R. J. Doerksen,  
907 Andreas Dreuw, B. D. Dunietz, A. D. Dutoi, T. R. Furlani, S. R. Gwaltney,  
908 Andreas Heyden, So Hirata, Chao-Ping Hsu, Gary Kedziora, R. Z. Khalliulin,  
909 Phil Klunzinger, A. M. Lee, M. S. Lee, WanZhen Liang, Itay Lotan, Nikhil Nair,  
910 Baron Peters, E. I. Proynov, P. A. Pieniazek, Y. M. Rhee, Jim Ritchie, Edina Rosta, C.  
911 D. Sherrill, A. C. Simmonett, J. E. Subotnik, H. L. W. H. F. S. III, Weimin Zhang, A.  
912 T. Bell, A. K. Chakraborty, D. M. Chipman, F. J. Keil, Arieh Warshel, W. J. Hehre, H.  
913 L. W. H. F. S. III, Jing Kong, A. I. Krylov, P. M. W. Gill and Martin Head-Gordon,  
914 *Phys. Chem. Chem. Phys.*, 2006, **8**, 3172–3191.
- 915 38 P. Raccuglia, K. C. Elbert, P. D. F. Adler, C. Falk, M. B. Wenny, A. Mollo, M. Zeller,  
916 S. A. Friedler, J. Schrier and A. J. Norquist, *Nat. 2016 5337601*, 2016, **533**, 73–76.
- 917 39 K. T. Butler, D. W. Davies, H. Cartwright, O. Isayev and A. Walsh, *Nat. 2018*  
918 *5597715*, 2018, **559**, 547–555.
- 919 40 J. Westermayr and P. Marquetand, *Chem. Rev.*, 2021, **121**, 9873–9926.
- 920 41 F. Häse, L. M. Roch, P. Friederich and A. Aspuru-Guzik, *Nat. Commun.*, 2020, **11**, 1–  
921 11.
- 922 42 S. A. Lopez, B. Sanchez-Lengeling, J. de Goes Soares and A. Aspuru-Guzik, *Joule*,  
923 2017, **1**, 857–870.
- 924 43 H. Sahu, W. Rao, A. Troisi and H. Ma, *Adv. Energy Mater.*, 2018, **8**, 1–9.
- 925 44 W. K. Tatum, D. Torrejon, A. B. Resing, J. W. Onorato and C. K. Luscombe, *Comput.*  
926 *Mater. Sci.*, 2021, **197**, 110599.
- 927 45 N. Majeed, M. Saladina, M. Krompiec, S. Greedy, C. Deibel and R. C. I. MacKenzie,  
928 *Adv. Funct. Mater.*, 2019, **1907259**, 1907259.
- 929 46 X. Rodríguez-Martínez, E. Pascual-San-José, Z. Fei, M. Heeney, R. Guimerà and M.  
930 Campoy-Quiles, *Energy Environ. Sci.*, 2021, **14**, 986–994.
- 931 47 X. Rodríguez-Martínez, E. Pascual-San-José and M. Campoy-Quiles, *Energy Environ.*

- 932 *Sci.*, 2021, **14**, 3301–3322.
- 933 48 K. Kranthiraja and A. Saeki, *Adv. Funct. Mater.*, 2021, **31**, 1–11.
- 934 49 E. O. Pyzer-Knapp, C. Suh, R. Gómez-Bombarelli, J. Aguilera-Iparraguirre and A.  
935 Aspuru-Guzik, *Annu. Rev. Mater. Res.*, 2015, **45**, 195–216.
- 936 50 J. Hachmann, R. Olivares-amaya, S. Atahan-evrenk, C. Amador-bedolla, R. S.  
937 Sanchez-carrera, A. Gold-parker, L. Vogt, A. M. Brockway and A. Aspuru-Guzik, *J.*  
938 *Phys. Chem. Lett.*, 2011, **2**, 2241–2251.
- 939 51 A. Mahmood, J.-L. Wang, Asif Mahmood and Jin-Liang Wang, *Energy Environ. Sci.*,  
940 2021, **14**, 90–105.
- 941 52 P. Malhotra, S. Biswas, F.-C. Chen and G. D. Sharma, *Sol. Energy*, 2021, **228**, 175–  
942 186.
- 943 53 A. Kuzmich, D. Padula, H. Ma and A. Troisi, *Energy Environ. Sci.*, 2017, **10**, 395–401.
- 944 54 I. Y. Kanal, S. G. Owens, J. S. Bechtel and G. R. Hutchison, *J. Phys. Chem. Lett.*,  
945 2013, **4**, 1613–1623.
- 946 55 L. Wilbraham, E. Berardo, L. Turcani, K. E. Jelfs and M. A. Zwijnenburg, *J. Chem.*  
947 *Inf. Model.*, 2018, **58**, 2450–2459.
- 948 56 E. O. Pyzer-Knapp, K. Li and A. Aspuru-Guzik, *Adv. Funct. Mater.*, 2015, **25**, 6495–  
949 6502.
- 950 57 J. Hachmann, R. Olivares-Amaya, A. Jinich, A. L. Appleton, M. A. Blood-Forsythe, L.  
951 R. Seress, C. Román-Salgado, K. Trepte, S. Atahan-Evrenk, S. Er, S. Shrestha, R.  
952 Mondal, A. Sokolov, Z. Bao and A. Aspuru-Guzik, *Energy Environ. Sci.*, 2014, **7**, 698–  
953 704.
- 954 58 N. Bérubé, V. Gosselin, J. Gaudreau and M. Côté, *J. Phys. Chem. C*, 2013, **117**, 7964–  
955 7972.
- 956 59 S. Nagasawa, E. Al-Naamani and A. Saeki, *J. Phys. Chem. Lett.*, 2018, **9**, 2639–2646.
- 957 60 Y. Huang, J. Zhang, E. S. Jiang, Y. Oya, A. Saeki, G. Kikugawa, T. Okabe, T. Okabe,  
958 F. S. Ohuchi and F. S. Ohuchi, *J. Phys. Chem. C*, 2020, **124**, 12871–12882.

- 959 61 W. Sun, M. Li, Y. Li, Z. Wu, Y. Sun, S. Lu, Z. Xiao, B. Zhao and K. Sun, *Adv. Theory*  
960 *Simulations*, 2019, **2**, 1–9.
- 961 62 R. Olivares-Amaya, C. Amador-Bedolla, J. Hachmann, S. Atahan-Evrenk, R. S.  
962 Sánchez-Carrera, L. Vogt and A. Aspuru-Guzik, *Energy Environ. Sci.*, 2011, **4**, 4849–  
963 4861.
- 964 63 H. Sahu, F. Yang, X. Ye, J. Ma, W. Fang and H. Ma, *J. Mater. Chem. A*, 2019, **7**,  
965 17480–17488.
- 966 64 D. Padula, J. D. D. Simpson and A. Troisi, *Mater. Horizons*, 2019, **6**, 343–349.
- 967 65 L. Simine, T. C. Allen and P. J. Rossky, *Proc. Natl. Acad. Sci.*, 2020, **117**, 13945–  
968 13948.
- 969 66 J. F. Joung, M. Han, J. Hwang, M. Jeong, D. H. Choi and S. Park, *JACS Au*, 2021, **1**,  
970 427–438.
- 971 67 S. Ye, W. Hu, X. Li, J. Zhang, K. Zhong, G. Zhang, Y. Luo, S. Mukamel and J. Jiang,  
972 *Proc. Natl. Acad. Sci.*, 2019, **116**, 201821044.
- 973 68 C. Nantasenamat, C. Isarankura-Na-Ayudhya, N. Tansila, T. Naenna and V.  
974 Prachayasittikul, *J. Comput. Chem.*, 2007, **28**, 1275–1289.
- 975 69 E. J. Beard, G. Sivaraman, Á. Vázquez-Mayagoitia, V. Vishwanath and J. M. Cole, *Sci.*  
976 *Data*, 2019, **6**, 307.
- 977 70 S. Varghese and S. Das, *J. Phys. Chem. Lett.*, 2011, **2**, 863–873.
- 978 71 C. L. Donley, J. Zaumseil, J. W. Andreasen, M. M. Nielsen, H. Sirringhaus, R. H.  
979 Friend and J. S. Kim, *J. Am. Chem. Soc.*, 2005, **127**, 12890–12899.
- 980 72 P. J. Brown, D. S. Thomas, A. Köhler, J. S. Wilson, J.-S. S. Kim, C. M. Ramsdale, H.  
981 Sirringhaus and R. H. Friend, *Phys. Rev. B*, 2003, **67**, 064203.
- 982 73 C. W. Yap, *J. Comput. Chem.*, 2011, **32**, 1466–1474.
- 983 74 D.-S. Cao, Q.-S. Xu, Q.-N. Hu and Y.-Z. Liang, *Bioinformatics*, 2013, **29**, 1092–1094.
- 984 75 H. Moriwaki, Y.-S. Tian, N. Kawashita and T. Takagi, *J. Cheminform.*, 2018, **10**, 4.

- 985 76 RDKit: Open-Source Cheminformatics Software.
- 986 77 Z. Luo, R. Ma, Y. Xiao, T. Liu, H. Sun, M. Su, Q. Guo, G. Li, W. Gao, Y. Chen, Y.  
987 Zou, X. Guo, M. Zhang, X. Lu, H. Yan and C. Yang, *Small*, 2020, **16**, 2001942.
- 988 78 M. Y. Mehboob, M. Adnan, R. Hussain and Z. Irshad, *Synth. Met.*, 2021, **277**, 116800.
- 989 79 R. Todeschini, M. Lasagni and E. Marengo, *J. Chemom.*, 1994, **8**, 263–272.
- 990 80 R. Todeschini and P. Gramatica, in *3D QSAR in Drug Design*, Kluwer Academic  
991 Publishers, Dordrecht, pp. 355–380.
- 992 81 R. Todeschini and P. Gramatica, *Quant. Struct. Relationships*, 1997, **16**, 113–119.
- 993 82 R. Todeschini and V. Consonni, *Molecular Descriptors for Chemoinformatics*, Wiley,  
994 Second Edi., 2009, vol. 41.
- 995 83 R. Todeschini and V. Consonni, *Handbook of Molecular Descriptors*, Wiley, 2000.
- 996 84 C. Bannwarth, E. Caldeweyher, S. Ehlert, A. Hansen, P. Pracht, J. Seibert, S. Spicher  
997 and S. Grimme, *WIREs Comput. Mol. Sci.*, , DOI:10.1002/wcms.1493.
- 998 85 N. Artrith, K. T. Butler, F.-X. Coudert, S. Han, O. Isayev, A. Jain and A. Walsh, *Nat.*  
999 *Chem.*, 2021, **13**, 505–508.
- 1000 86 M. Lee, *Adv. Energy Mater.*, 2019, 1900891.
- 1001 87 F. Pedregosa, R. Weiss, M. Brucher, G. Varoquaux, A. Gramfort, V. Michel, B.  
1002 Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas,  
1003 A. Passos, D. Cournapeau, M. Brucher, M. Perrot and É. Duchesnay, *J. Mach. Learn.*  
1004 *Res.*, 2011, **12**, 2825–2830.
- 1005 88 D. Rogers and M. Hahn, *J. Chem. Inf. Model.*, 2010, **50**, 742–754.
- 1006 89 L. H. Hall and L. B. Kier, *J. Chem. Inf. Comput. Sci.*, 1995, **35**, 1039–1045.
- 1007 90 D. C. Elton, Z. Boukouvalas, M. S. Butrico, M. D. Fuge and P. W. Chung, *Sci. Rep.*,  
1008 2018, **8**, 9059.
- 1009 91 C. Bannwarth, S. Ehlert and S. Grimme, *J. Chem. Theory Comput.*, 2019, **15**, 1652–  
1010 1671.

- 1011 92 S. Grimme, C. Bannwarth and P. Shushkov, *J. Chem. Theory Comput.*, 2017, **13**,  
1012 1989–2009.
- 1013 93 F. Deledalle, T. Kirchartz, M. S. Vezie, M. Campoy-Quiles, P. S. Tuladhar, J. Nelson  
1014 and J. R. Durrant, *Phys. Rev. X*, 2015, **5**, 1–13.
- 1015 94 F. C. Spano, *Acc. Chem. Res.*, 2010, **43**, 429–439.
- 1016 95 G. Forti, A. Nitti, P. Osw, G. Bianchi, R. Po and D. Pasini, *Int. J. Mol. Sci.*, 2020, **21**,  
1017 8085.
- 1018 96 L. A. A. Pettersson, L. S. Roman and O. Inganäs, *J. Appl. Phys.*, 1999, **86**, 487–496.
- 1019 97 S. Karuthedath, J. Gorenflot, Y. Firdaus, N. Chaturvedi, C. S. P. De Castro, G. T.  
1020 Harrison, J. I. Khan, A. Markina, A. H. Balawi, T. A. Dela Peña, W. Liu, R.-Z. Liang,  
1021 A. Sharma, S. H. K. Paleti, W. Zhang, Y. Lin, E. Alarousu, D. H. Anjum, P. M.  
1022 Beaujuge, S. De Wolf, I. McCulloch, T. D. Anthopoulos, D. Baran, D. Andrienko and  
1023 F. Laquai, *Nat. Mater.*, 2021, **20**, 378–384.
- 1024 98 Z. Cong, B. Zhao, Z. Chen, W. Wang, H. Wu, J. Liu, J. Wang, L. Wang, W. Ma and C.  
1025 Gao, *ACS Appl. Mater. Interfaces*, 2019, **11**, 16795–16803.
- 1026 99 W. Zhao, S. Li, H. Yao, S. Zhang, Y. Zhang, B. Yang and J. Hou, *J. Am. Chem. Soc.*,  
1027 2017, **139**, 7148–7151.
- 1028 100 W. Wang, B. Zhao, Z. Cong, Y. Xie, H. Wu, Q. Liang, S. Liu, F. Liu, C. Gao, H. Wu  
1029 and Y. Cao, *ACS Energy Lett.*, 2018, **3**, 1499–1507.
- 1030 101 S. Holliday, R. S. Ashraf, A. Wadsworth, D. Baran, S. A. Yousaf, C. B. Nielsen, C. H.  
1031 Tan, S. D. Dimitrov, Z. Shang, N. Gasparini, M. Alamoudi, F. Laquai, C. J. Brabec, A.  
1032 Salleo, J. R. Durrant and I. McCulloch, *Nat. Commun.*, 2016, **7**, 1–11.
- 1033 102 D. Baran, T. Kirchartz, S. Wheeler, S. Dimitrov, M. Abdelsamie, J. Gorman, R. S.  
1034 Ashraf, S. Holliday, A. Wadsworth, N. Gasparini, P. Kaienburg, H. Yan, A. Amassian,  
1035 C. J. Brabec, J. R. Durrant and I. McCulloch, *Energy Environ. Sci.*, 2016, **9**, 3783–  
1036 3793.
- 1037 103 N. A. Cooling, E. F. Barnes, F. Almyahi, K. Feron, M. F. Al-Mudhaffer, A. Al-Ahmad,  
1038 B. Vaughan, T. R. Andersen, M. J. Griffith, A. S. Hart, A. G. Lyons, W. J. Belcher and



- 1039 P. C. Dastoor, *J. Mater. Chem. A*, 2016, **4**, 10274–10281.
- 1040 104 M. M. Wienk, J. M. Kroon, W. J. H. Verhees, J. Knol, J. C. Hummelen, P. A. van Hal  
1041 and R. A. J. Janssen, *Angew. Chemie Int. Ed.*, 2003, **42**, 3371–3375.
- 1042 105 D. Baran, R. S. Ashraf, D. A. Hanifi, M. Abdelsamie, N. Gasparini, J. A. Röhr, S.  
1043 Holliday, A. Wadsworth, S. Lockett, M. Neophytou, C. J. M. Emmott, J. Nelson, C. J.  
1044 Brabec, A. Amassian, A. Salleo, T. Kirchartz, J. R. Durrant and I. McCulloch, *Nat.*  
1045 *Mater.*, 2017, **16**, 363–369.
- 1046 106 M. Li, Y. Liu, W. Ni, F. Liu, H. Feng, Y. Zhang, T. Liu, H. Zhang, X. Wan, B. Kan, Q.  
1047 Zhang, T. P. Russell and Y. Chen, *J. Mater. Chem. A*, 2016, **4**, 10409–10413.
- 1048 107 N. Qiu, H. Zhang, X. Wan, C. Li, X. Ke, H. Feng, B. Kan, H. Zhang, Q. Zhang, Y. Lu  
1049 and Y. Chen, *Adv. Mater.*, 2017, **29**, 1604964.
- 1050 108 C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D.  
1051 Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer,  
1052 M. H. van Kerkwijk, M. Brett, A. Haldane, J. F. del Río, M. Wiebe, P. Peterson, P.  
1053 Gérard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke and T.  
1054 E. Oliphant, *Nature*, 2020, **585**, 357–362.
- 1055 109 P. Virtanen, R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E.  
1056 Burovski, P. Peterson, W. Weckesser, J. Bright, S. J. van der Walt, M. Brett, J. Wilson,  
1057 K. J. Millman, N. Mayorov, A. R. J. Nelson, E. Jones, R. Kern, E. Larson, C. J. Carey,  
1058 Í. Polat, Y. Feng, E. W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I.  
1059 Henriksen, E. A. Quintero, C. R. Harris, A. M. Archibald, A. H. Ribeiro, F. Pedregosa  
1060 and P. van Mulbregt, *Nat. Methods*, 2020, **17**, 261–272.
- 1061 110 N. S. Raju, R. Bilgic, J. E. Edwards and P. F. Fleeer, *Appl. Psychol. Meas.*, 1997, **21**,  
1062 291–305.
- 1063 111 G. Moreau, *J. Chem. Inf. Comput. Sci.*, 1997, **37**, 929–938.
- 1064

1065 Supporting Information

1066 **Identifying structure-absorption relationships and predicting absorption strength of non-**  
1067 **fullerene acceptors for organic photovoltaics**

1068  
1069 *Jun Yan,<sup>a,#</sup> Xabier Rodríguez-Martínez,<sup>\*,b,c,#</sup> Drew Pearce,<sup>a</sup> Hana Douglas,<sup>a</sup> Danai Bili,<sup>a</sup>*  
1070 *Mohammed Azzouzi,<sup>a</sup> Flurin Eisner,<sup>a</sup> Alise Virbule,<sup>a</sup> Elham Rezasoltani,<sup>a</sup> Valentina Belova,<sup>c</sup>*  
1071 *Bernhard Dörfling,<sup>c</sup> Sheridan Few,<sup>a,f</sup> Anna A. Szumska,<sup>a</sup> Xueyan Hou,<sup>a</sup> Guichuan Zhang,<sup>d</sup> Hin-*  
1072 *Lap Yip,<sup>d,e</sup> Mariano Campoy-Quiles <sup>\*c</sup> and Jenny Nelson <sup>\*a</sup>*

1073

1074 <sup>#</sup> J.Y. and X.R.-M. contributed equally to this work.

1075

1076 <sup>a</sup> Department of Physics, Imperial College London, SW7 2AZ, London, United Kingdom

1077 Email: [jenny.nelson@imperial.ac.uk](mailto:jenny.nelson@imperial.ac.uk)

1078

1079 <sup>b</sup> Electronic and Photonic Materials (EFM), Department of Physics, Chemistry and Biology  
1080 (IFM), Linköping University, Linköping, SE 581 83 Sweden

1081 Email: [xabier.rodriguez.martinez@liu.se](mailto:xabier.rodriguez.martinez@liu.se)

1082

1083 <sup>c</sup> Instituto de Ciencia de Materiales de Barcelona, ICMAB-CSIC, Campus UAB, Bellaterra  
1084 08193, Spain

1085 Email: [mcampoy@icmab.es](mailto:mcampoy@icmab.es)

1086

1087 <sup>d</sup> Institute of Polymer Optoelectronic Materials and Devices, State Key Laboratory of  
1088 Luminescent Materials and Devices, South China University of Technology, Guangzhou  
1089 510640, P. R. China

1090

1091 <sup>e</sup> Department of Materials Science and Engineering, City University of Hong Kong, Tat Chee  
1092 Avenue, Kowloon, Hong Kong

1093

1094 <sup>f</sup> Sustainability Research Institute, School of Earth and Environment, University of Leeds,  
1095 Leeds, LS2 9JT

1096

1097

1098 **Supplementary Note 1.** Chemical names and nomenclature of the materials  
1099 highlighted in this work.

1100 **PC61BM:** [6,6]-Phenyl-C<sub>61</sub>-butyric acid methyl ester

1101 **PC71BM:** [6,6]-Phenyl-C<sub>71</sub>-butyric acid methyl ester

1102 **ICBA:** 1',1'',4',4''-Tetrahydro-di[1,4]methanonaphthaleno[1,2:2',3',56,60:2'',3'']<sub>5,6</sub>fullerene-  
1103 C<sub>60</sub>

1104 **Y5:** (2,2'-((2Z,2'Z)-((12,13-bis(2-ethylhexyl)-3,9-diundecyl-12,13-  
1105 dihydro[1,2,5]thiadiazolo[3,4e]thieno[2'',3'':4',5']thieno[2',3':4,5]pyrrolo[3,2-g]  
1106 thieno[2',3':4,5]thieno[3,2-b]-indole-2,10-diyl)bis(methanylylidene))bis(3-oxo-2,3-dihydro-  
1107 1H-indene-2,1-diylidene))dimalononitrile)

1108 **Y6:** 2,2'-((2Z,2'Z)-((12,13-bis(2-ethylhexyl)-3,9-diundecyl-12,13-dihydro-  
1109 [1,2,5]thiadiazolo[3,4-e]thieno[2'',3'':4',5']thieno[2',3':4,5]pyrrolo[3,2-  
1110 g]thieno[2',3':4,5]thieno[3,2-b]indole-2,10-diyl)bis(methanylylidene))bis(5,6-difluoro-3-oxo-  
1111 2,3-dihydro-1H-indene-2,1-diylidene))dimalononitrile)

1112 **Y7:** 2,2'-((2Z,2'Z)-((12,13-bis(2-ethylhexyl)-3,9-diundecyl-12,13-dihydro-  
1113 [1,2,5]thiadiazolo[3,4-e]thieno[2'',3'':4',5']thieno[2',3':4,5]pyrrolo[3,2-  
1114 g]thieno[2',3':4,5]thieno[3,2-b]indole-2,10-diyl)bis(methanylylidene))bis(5,6-dichloro-3-oxo-  
1115 2,3-dihydro-1H-indene-2,1-diylidene))dimalononitrile)

1116 **Y11:** 2,2'-((2Z,2'Z)-((6,12,13-tris(2-ethylhexyl)-3,9-diundecyl-12,13-dihydro-6H-  
1117 thieno[2'',3'':4',5']thieno[2',3':4,5]pyrrolo[3,2-g]thieno[2',3':4,5]thieno[3,2-  
1118 b][1,2,3]triazolo[4,5-e]indole-2,10-diyl)bis(methanylylidene))bis(5,6-difluoro-3-oxo-2,3-  
1119 dihydro-1H-indene-2,1-diylidene))dimalononitrile)

1120 **Y12:** 2,2'-((2Z,2'Z)-((12,13-bis(2-butyloctyl)-3,9-diundecyl-12,13-dihydro-  
1121 [1,2,5]thiadiazolo[3,4-e]thieno[2'',3'':4',5']thieno[2',3':4,5]pyrrolo[3,2-  
1122 g]thieno[2',3':4,5]thieno[3,2-b]indole-2,10-diyl)bis(methanylylidene))bis(5,6-difluoro-3-oxo-  
1123 2,3-dihydro-1H-indene-2,1-diylidene))dimalononitrile)

1124 **O-IDTBR:** (5Z,5'Z)-5,5'-(((4,4,9,9-tetrakis(n-octyl)-4,9-dihydro-s-indaceno[1,2-b:5,6-  
1125 b']dithiophene-2,7-diyl)bis(benzo[c][1,2,5]thiadiazole-7,4-diyl))bis(methaneylylidene))bis(3-  
1126 ethyl-2-thioxothiazolidin-4-one)

- 1127 **O-IDFBR:** (5Z,5'Z)-5,5'-(((6,6,12,12-tetraoctyl-6,12-dihydroindeno[1,2-b]fluorene-2,8-  
1128 diyl)bis(benzo[c][1,2,5]thiadiazole-7,4-diyl))bis(methaneylylidene))bis(3-ethyl-2-  
1129 thioxothiazolidin-4-one)
- 1130 **EH-IDTBR:** (5Z,5'Z)-5,5'-(((4,4,9,9-tetrakis(2-ethylhexyl)-4,9-dihydro-s-indaceno[1,2-b:5,6-  
1131 b']dithiophene-2,7-diyl)bis(benzo[c][1,2,5]thiadiazole-7,4-diyl))bis(methaneylylidene))bis(3-  
1132 ethyl-2-thioxothiazolidin-4-one)
- 1133 **IDIC:** 2,2'-((2Z,2'Z)-((4,4,9,9-tetrahexyl-4,9-dihydro-s-indaceno[1,2-b:5,6-b']dithiophene-  
1134 2,7-diyl)bis(methaneylylidene))bis(3-oxo-2,3-dihydro-1H-indene-2,1-  
1135 diylidene))dimalononitrile
- 1136 **SN6IC-4F:** 2,2'-((2Z,2'Z)-((thieno[3,2-  
1137 b]thieno[2''',3''':4'',5'']pyrrolo[2'',3'':4',5']thieno[2',3':4,5]thieno[2,3- d]pyrrole,4,9-dihydro-4,9-  
1138 di-1-octylnonyl-2,7-diyl)bis(methaneylylidene))bis((5,6-difluoro-3-oxo-2,3 -dihydro-1H-  
1139 indene-2,1-diylidene))dimalononitrile
- 1140 **ITIC:** 2,2'-[[6,6,12,12-tetrakis(4-hexylphenyl)-6,12-dihydrodithieno[2,3-d:2',3'-d']-s-  
1141 indaceno[1,2-b:5,6-b']dithiophene-2,8-diyl]bis[methylidyne(3-oxo-1H-indene-2,1(3H)-  
1142 diylidene)]]bis[propanedinitrile]
- 1143 **ITIC-C<sub>2</sub>C<sub>6</sub>:** 2,2'-[[6,6,12,12-tetrakis(2-ethylhexyl)-6,12-dihydrodithieno[2,3-d:2',3'-d']-s-  
1144 indaceno[1,2-b:5,6-b']dithiophene-2,8-diyl]bis[methylidyne(3-oxo-1H-indene-2,1(3H)-  
1145 diylidene)]]bis[propanedinitrile]
- 1146 **ITIC-C<sub>8</sub>:** 2,2'-[[6,6,12,12-tetrakis(n-octyl)-6,12-dihydrodithieno[2,3-d:2',3'-d']-s-  
1147 indaceno[1,2-b:5,6-b']dithiophene-2,8-diyl]bis[methylidyne(3-oxo-1H-indene-2,1(3H)-  
1148 diylidene)]]bis[propanedinitrile]
- 1149 **IT-4F:** 9-Bis(2-methylene-((3-(1,1-dicyanomethylene)-6,7-difluoro)-indanone))-5,5,11,11-  
1150 tetrakis(4-hexylphenyl)-dithieno[2,3-d:2',3'-d']-s-indaceno[1,2-b:5,6-b']dithiophene
- 1151 **CBM:** 2,2'-(7,7'-(9-(heptadecan-9-yl)-9H-carbazole-2,7-diyl)bis(benzo[c][1,2,5]thiadiazole-  
1152 7,4-diyl))bis(methan-1-yl-1-ylidene)dimalononitrile
- 1153 **FBR:** 5,5'-[(9,9-Dioctyl-9H-fluorene-2,7-diyl)bis(2,1,3-benzothiadiazole-7,4-  
1154 diylmethylidyne)]bis[3-ethyl-2-thioxo-4-thiazolidinone]

- 1155 **BTMP:** 2,2'-((2Z,2'Z)-((12,13-diisobutyl-3,9-dimethyl-5,7,12,13-tetrahydro-  
1156 [1,2,5]thiadiazolo[3,4-e]thieno[2'',3'':4',5']thieno[2',3':4,5]pyrrolo[3,2-  
1157 g]thieno[2',3':4,5]thieno[3,2-b]indole-2,10-diyl)bis(methaneylylidene))bis(6-methyl-3-oxo-  
1158 2,3-dihydro-1H-indene-2,1-diylidene))dimalononitrile
- 1159 **BTTPC:** 2,2'-((6Z,6'Z)-((12,13-diisobutyl-3,9-dimethyl-5,7,12,13-tetrahydro-  
1160 [1,2,5]thiadiazolo[3,4-e]thieno[2'',3'':4',5']thieno[2',3':4,5]pyrrolo[3,2-  
1161 g]thieno[2',3':4,5]thieno[3,2-b]indole-2,10-diyl)bis(methaneylylidene))bis(5-oxo-5,6-dihydro-  
1162 7H-indeno[5,6-b]thiophene-6,7-diylidene))dimalononitrile
- 1163 **BTDTTP-4F:** 2,2'-((2Z,2'Z)-((3,12-dimethyl-13,14-dihydro-12H-[1,2,5]thiadiazolo[3,4-  
1164 e]thieno[2'',3'':4',5']pyrrolo[2',3':4,5]thieno[3,2-  
1165 b]thieno[2'',3'':4',5']thieno[2',3':4,5]pyrrolo[3,2-g]indole-2,10-  
1166 diyl)bis(methaneylylidene))bis(5,6-difluoro-3-oxo-2,3-dihydro-1H-indene-2,1-  
1167 diylidene))dimalononitrile
- 1168 **BDTP-4F:** 2,2'-(((1Z,1'Z)-(1,11-dimethyl-4,6,6c,10,11,11b,12,13-octahydro-2H-  
1169 [1,2,5]thiadiazolo[3,4-e]thieno[2'',3'':4',5']pyrrolo[2',3':4,5]thieno[3,2-  
1170 b]thieno[2',3':4,5]pyrrolo[3,2-g]indole-2,9(1H)-diylidene)bis(methaneylylidene))bis(5,6-  
1171 difluoro-3-oxo-2,3,3a,6,7,7a-hexahydro-1H-indene-2-yl-1-ylidene))dimalononitrile
- 1172 **BTTPTP-2OYPD:** 2,2'-((2Z,2'Z)-((13,14-diisobutyl-5,7,13,14-tetrahydro-  
1173 [1,2,5]thiadiazolo[3,4-e]thieno[2''',3''':4'',5'']thieno[2'',3'':4',5']thieno[2',3':4,5]pyrrolo[3,2-  
1174 g]thieno[2'',3'':4',5']thieno[2',3':4,5]thieno[3,2-b]indole-2,10-  
1175 diyl)bis(methaneylylidene))bis(3-oxo-2,3-dihydro-1H-indene-2,1-diylidene))dimalononitrile
- 1176 **BTPTTT-2OYPD:** 2,2'-((2Z,2'Z)-((13,14-diisobutyl-5,7,13,14-tetrahydro-  
1177 [1,2,5]thiadiazolo[3,4-e]thieno[2''',3''':4'',5'']thieno[2'',3'':4',5']thieno[2',3':4,5]pyrrolo[3,2-  
1178 g]thieno[2'',3'':4',5']thieno[2',3':4,5]thieno[3,2-b]indole-2,10-  
1179 diyl)bis(methaneylylidene))bis(3-oxo-2,3-dihydro-1H-indene-2,1-diylidene))dimalononitrile
- 1180 **IEICO:** 2,2'-((2Z,2'Z)-((5,5'-(4,4,9,9-tetrakis(4-hexylphenyl)-4,9-dihydros-indaceno[1,2-  
1181 b:5,6-b']dithiophene-2,7-diyl)bis(4-((2-ethylhexyl)oxy)thiophene-5,2-  
1182 diyl))bis(methaneylylidene))bis(3-oxo-2,3-dihydro-1H-indene-2,1-diylidene))dimalononitrile
- 1183 **IEICO-4F:** 2,2'-((2Z,2'Z)-(((4,4,9,9-tetrakis(4-hexylphenyl)-4,9-dihydro-sindaceno[1,2-b:5,6-  
1184 b']dithiophene-2,7-diyl)bis(4-((2-ethylhexyl)oxy)thiophene-5,2-

1185 diyl))bis(methanylylidene))bis(5,6-difluoro-3-oxo-2,3-dihydro-1H-indene-2,1-  
1186 diylidene))dimalononitrile

1187 **BTTPTP-4Cl:** 2,2'-((2Z,2'Z)-((13,14-diisobutyl-5,7,13,14-tetrahydro-[1,2,5]thiadiazolo[3,4-  
1188 e]thieno[2''',3''':4'',5'']thieno[2'',3'':4',5']thieno[2',3':4,5]pyrrolo[3,2-  
1189 g]thieno[2'',3'':4',5']thieno[2',3':4,5]thieno[3,2-b]indole-2,10-  
1190 diyl)bis(methaneylylidene))bis(5,6-dichloro-3-oxo-2,3-dihydro-1H-indene-2,1-  
1191 diylidene))dimalononitrile

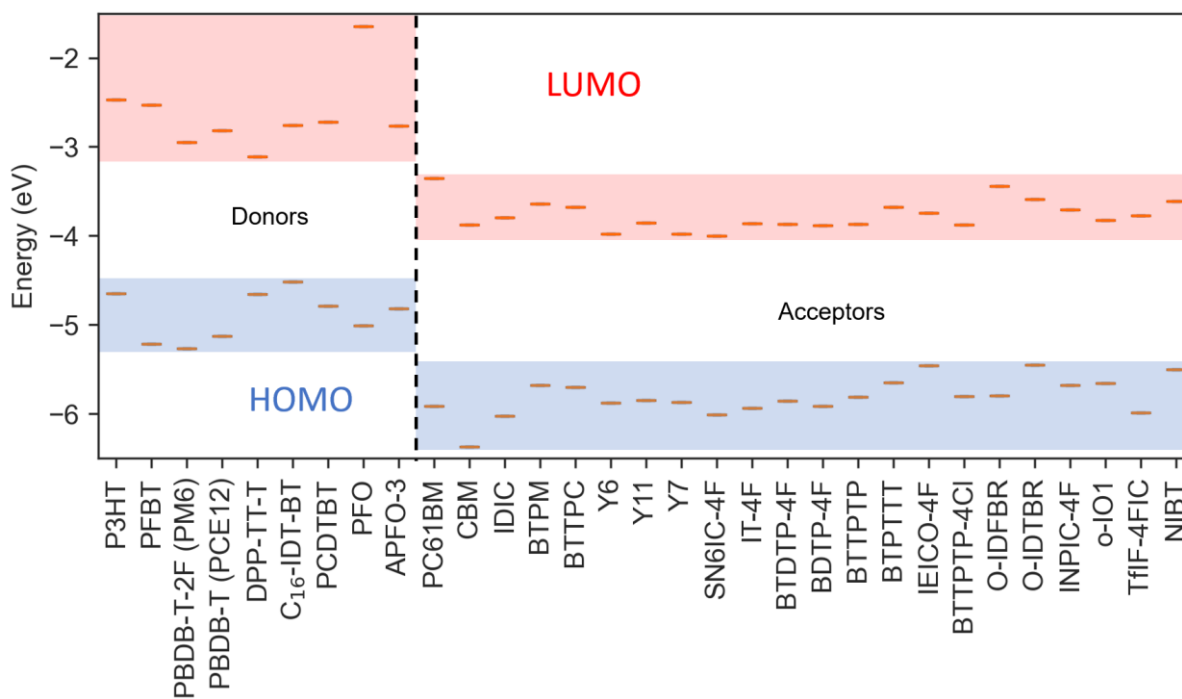
1192 **INPIC-4F:** [(Z)-2-( {24-[(Z)-(1-Dicyanomethylidene-5,6-difluoro-3-oxo-2-  
1193 indanylidene)methyl]-15,15,30,30-tetrakis(p-hexylphenyl)-12,27-dioctyl-5,8,20,23-tetrathia-  
1194 12,27-diazanonacyclo[16.12.0.0<sup>3,16</sup>.0<sup>4,14</sup>.0<sup>6,13</sup>.0<sup>7,11</sup>.0<sup>19,29</sup>.0<sup>21,28</sup>.0<sup>22,26</sup>]}triaconta-  
1195 1(18),2,4(14),6(13),7(11),9,16,19(29),21(28),22(26),24-undecaen-9-yl} methylidene)-5,6-  
1196 difluoro-3-oxo-1-indanylidene]propanedinitrile

1197 **o-IO1:** 2-((Z)-2-((5-(7-(5-(((Z)-1-(dicyanomethylene)-5,6-difluoro-3-oxo-1,3-dihydro-2H-  
1198 inden-2-ylidene)methyl)-3-((2-ethylhexyl)oxy)thiophen-2-yl)-4,4,9,9-tetraoctyl-4,9-dihydro-  
1199 s-indaceno[1,2-b:5,6-b']dithiophen-2-yl)-4-(2-ethylhexyl)thiophen-2-yl)methylene)-5,6-  
1200 difluoro-3-oxo-2,3-dihydro-1H-inden-1-ylidene)malononitrile

1201 **TfIF-4FIC:** [(Z)-2-( {26-[(Z)-(1-Dicyanomethylidene-5,6-difluoro-3-oxo-2-  
1202 indanylidene)methyl]-7,7,16,16,23,23,32,32-octaoctyl-11,27-  
1203 dithianonacyclo[17.13.0.0<sup>3,17</sup>.0<sup>4,15</sup>.0<sup>6,13</sup>.0<sup>8,12</sup>.0<sup>20,31</sup>.0<sup>22,29</sup>.0<sup>24,28</sup>]}dotriaconta-  
1204 1(19),2,4(15),5,8(12),9,13,17,20(31),21,24(28),25,29-tridecaen-10-yl} methylidene)-5,6-  
1205 difluoro-3-oxo-1-indanylidene]propanedinitrile

1206 **NIBT:** (7Z,7'Z)-7,7'-(((4,4,9,9-tetrakis(4-octylphenyl)-4,9-dihydro-s-indaceno[1,2-b:5,6-  
1207 b']dithiophene-2,7-diyl)bis(benzo[c][1,2,5]thiadiazole-7,4-diyl))bis(methaneylylidene))bis(2-  
1208 (2-ethylhexyl)-1H-indeno[6,7,1-def]isoquinoline-1,3,6(2H,7H)-trione)

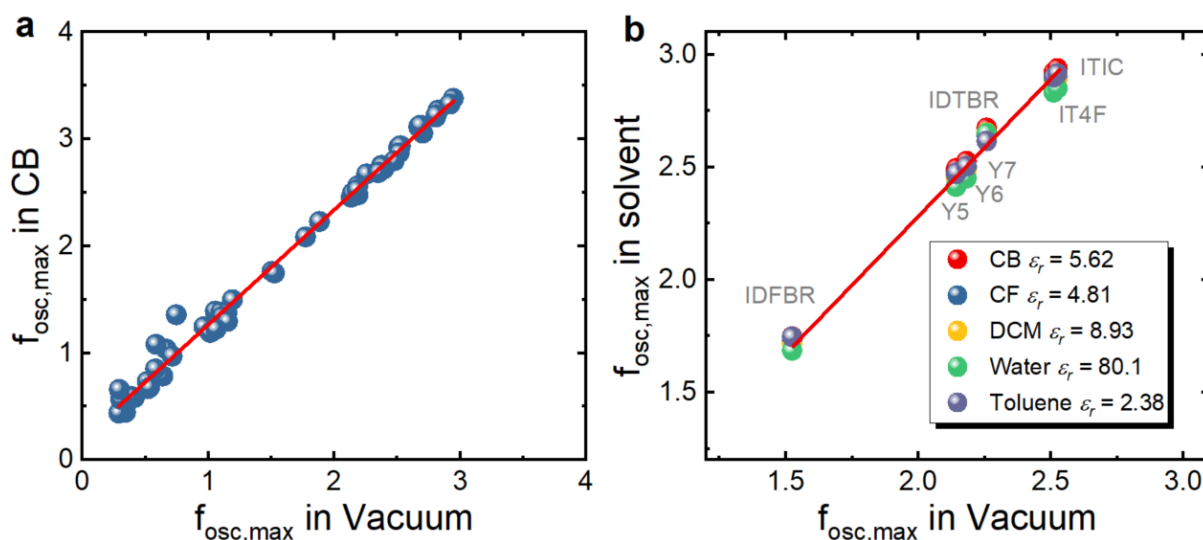
1209



1210

1211 **Figure S1. Highest occupied molecular orbital (HOMO) and lowest unoccupied molecular**  
 1212 **orbital (LUMO) energy levels of representative donor and acceptor molecules as retrieved**  
 1213 **from TDDFT calculations.** Molecules considered as NFAs in this work show proper HOMO  
 1214 and LUMO energy level alignment to act as electron acceptor when in a bulk heterojunction  
 1215 blend with commonly used donors, such as P3HT, PCDTBT, PM6 or PBDB-T.

1216

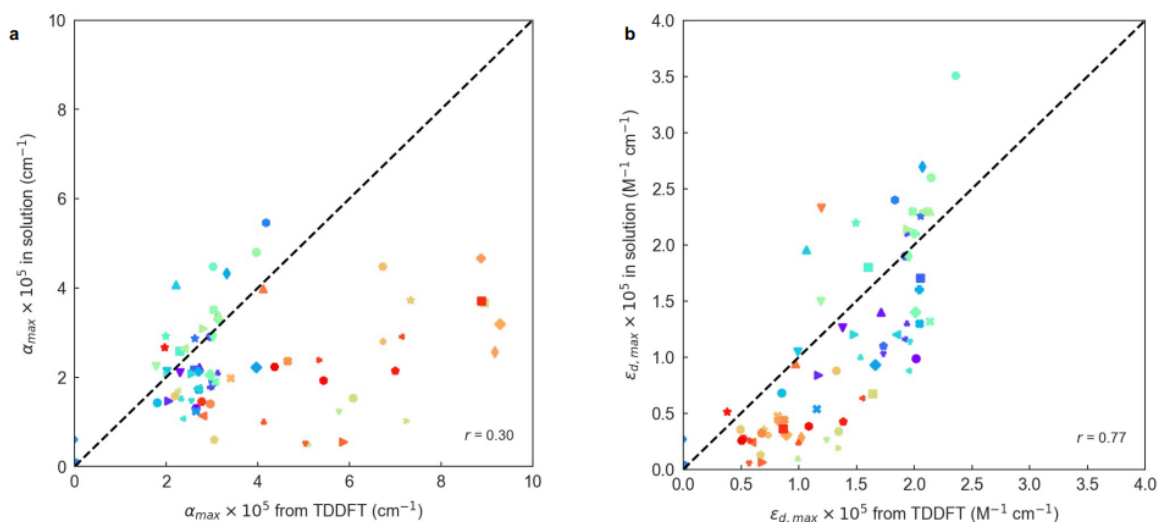


1217

1218 **Figure S2. Solvent effect on the maximum oscillator strength ( $f_{osc,max}$ ) in TDDFT**  
 1219 **calculations. (a)**  $f_{osc,max}$  in Chlorobenzene (CB) versus in vacuum for 56 organic molecules  
 1220 including common NFAs. **(b)**  $f_{osc,max}$  in various organic solvents versus in vacuum for 7  
 1221 different organic molecules, O-IDFBR, O-IDTBR, IT-4F, ITIC, Y5, Y6, and Y7. Noting here  
 1222 that  $\epsilon_{d,max} \propto f_{osc,max}$ . TDDFT was performed under B3LYP/6-311+G(d,p) using Polarizable-  
 1223 continuum-solvent-model (PCM). We can see that the choice of solvent does not affect  $f_{osc,max}$   
 1224 much, and that a good linear correlation between solvent and vacuum  $f_{osc,max}$  is obtained. This  
 1225 tells us that the same correlation between TDDFT and experiments will be maintained based  
 1226 on either vacuum environment or polarized medium, which allows us focus on TDDFT results  
 1227 from vacuum calculations only. This is a great benefit since most of the TDDFT calculations  
 1228 by the present and past group members were done in vacuum, allowing us to have a larger  
 1229 database.

1230



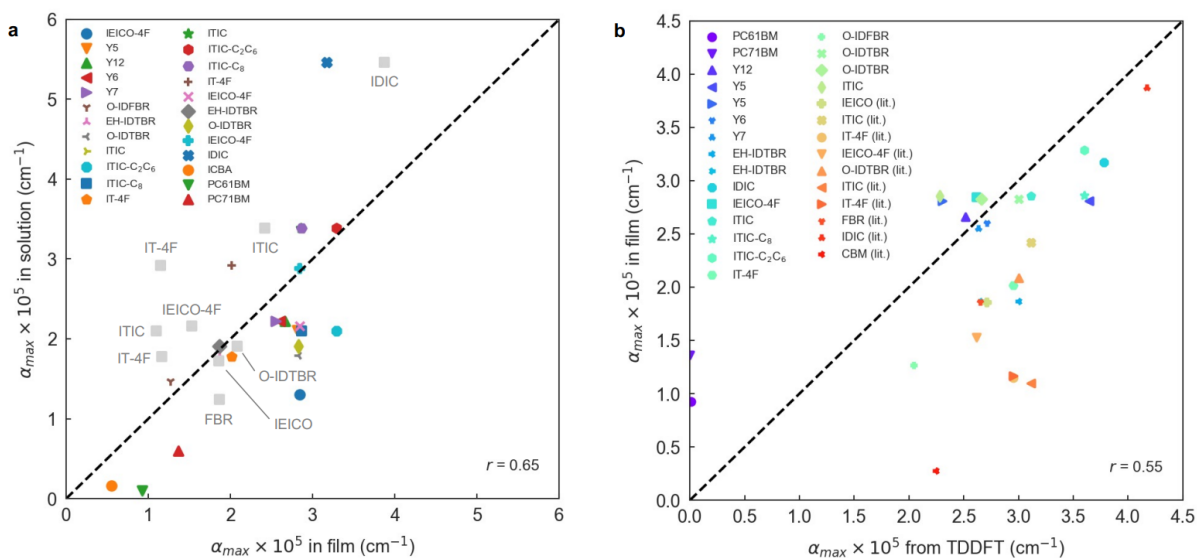


1231

1232 **Figure S3.** (a) Correlation between the maximum of the absorption coefficient ( $\epsilon_{d,max}$ )  
 1233 obtained in solution state with the TDDFT calculated values. (b) Same dataset yet plotted in  
 1234 terms of maximum molar extinction coefficient  $\epsilon_{d,max}$ ). A significant fraction of this dataset  
 1235 was collected from literature.<sup>95,97–107</sup> When required, a refractive index of 1.5 and a solid density  
 1236 of  $1000 \text{ g L}^{-1}$  were considered for all materials.

1237



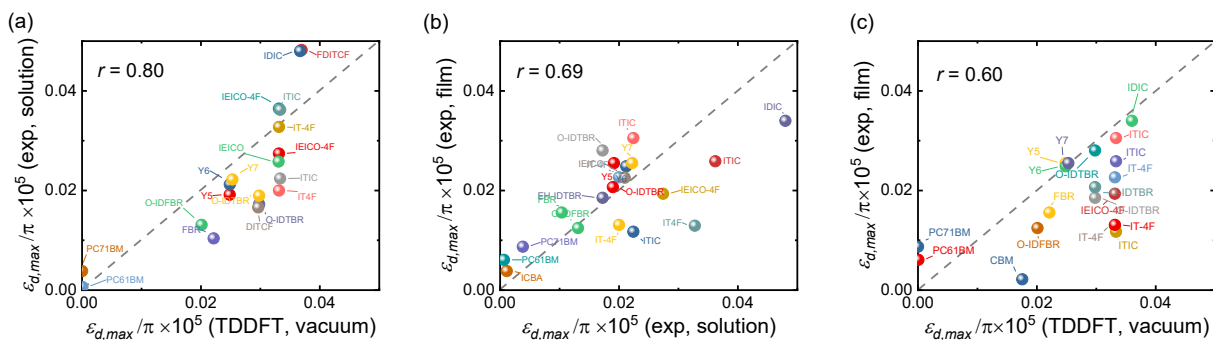


1248

1249 **Figure S5.** (a) Maximum absorption coefficient ( $\alpha_{max}$ ) in solution and film. Grey squares  
 1250 correspond to data obtained from literature.<sup>97,99</sup> (b) Maximum absorption coefficient in film and  
 1251 as obtained in their corresponding TDDFT calculations. A few data points (labelled as lit.)  
 1252 correspond to values extracted from literature.<sup>97,99</sup>

1253

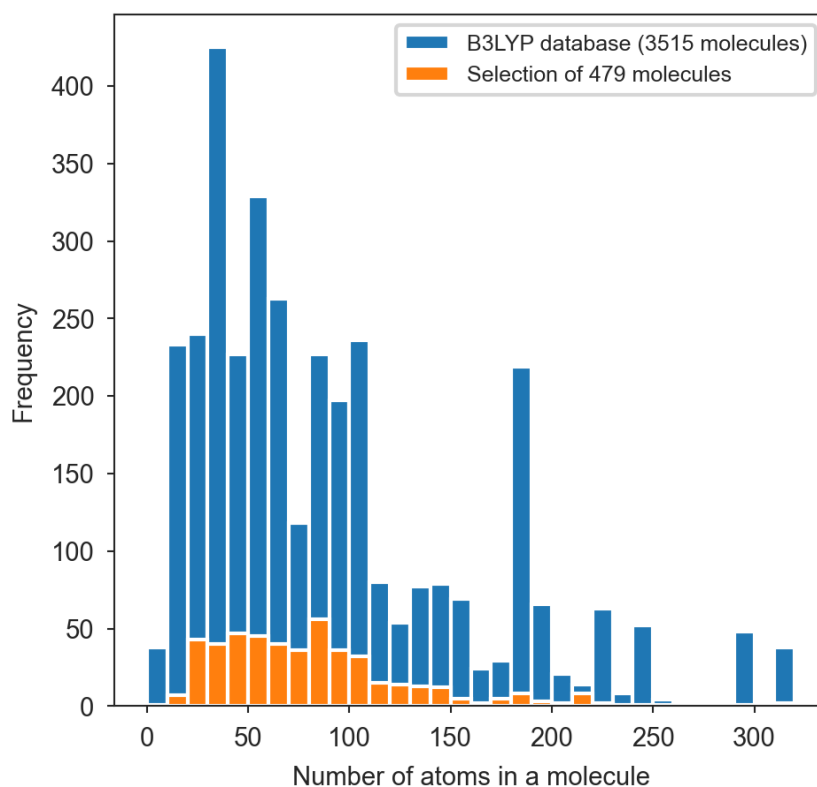
1254



1255

1256 **Figure S6.** Effect of the number of  $\pi$ -electrons on the comparison between experimental  
 1257 (solution and film) maximum molar extinction coefficients and TDDFT results for the NFAs  
 1258 studied. (a) Experimental  $\epsilon_{d,max}/\pi$  in solution versus calculated  $\epsilon_{d,max}/\pi$  using TDDFT, (b)  
 1259 experimental  $\epsilon_{d,max}/\pi$  in film (solid state) versus that in solution; and (c) experimental  
 1260  $\epsilon_{d,max}/\pi$  in solid state film versus calculated  $\epsilon_{d,max}/\pi$  using TDDFT. Unit of  $\epsilon_{d,max}/\pi$  is  $M^{-1}$   
 1261  $cm^{-1}$ .

1262

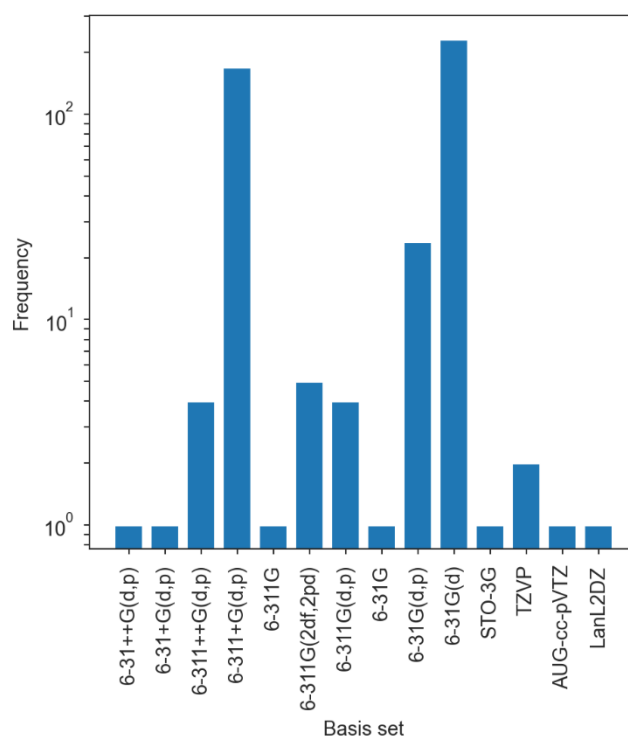


1263  
 1264 **Figure S7.** Histogram of the number of atoms present in the molecules of our TDDFT (B3LYP)  
 1265 dataset. Blue bars correspond to the molecules found originally in the dataset (3515 entries).  
 1266 Orange bars represent the distribution of the number of atoms found in the 479 molecules  
 1267 selected based on lowest energy conformation criteria.

1268

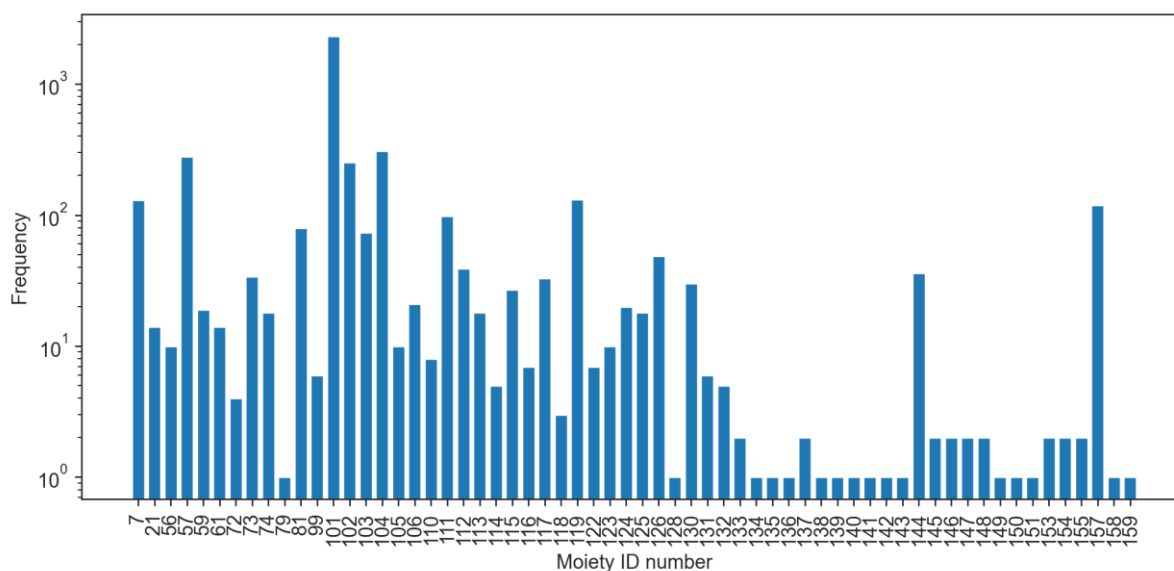
1269 **Supplementary Note 2.** Description of the TDDFT database, statistical and machine-learning  
1270 methods.

1271 The pristine data source of this work consists of a database of 3515 molecules optimized via  
1272 DFT using the B3LYP functional as implemented in Gaussian09 software package. The  
1273 database gathers original calculations performed for this particular study on conjugated small  
1274 molecules as well as others developed in-house during the past years, including diverse  
1275 conjugated small molecules, fullerenes and conjugated (co-)oligomers in distinct conformations  
1276 (i.e., cis-/trans-).<sup>29</sup> Given the variety of input sources, the corresponding data cleaning  
1277 procedure consists of: i) identifying duplicates based on molecular weight; and ii) picking the  
1278 lowest energy molecular conformation among each set of duplicates. The filtering results in a  
1279 final selection of 479 conjugated small molecules and oligomers optimized at the B3LYP level  
1280 of theory. The resulting database gathers a variety of basis sets employed in the geometrical  
1281 optimization step: 48% of the molecules were optimized using the 6-31G(d) set and 36% of  
1282 them using the more computationally-expensive 6-311+G(d,p), see **Figure S8**. Furthermore,  
1283 the chemical heterogeneity of the studied database is leaned toward known molecules and  
1284 moieties of frequent use in high-performing solar energy harvesting applications, see Figure S9  
1285 and Figure S10. Side chains are systematically omitted or substituted by methyl groups in all  
1286 calculations to reduce the computational cost.



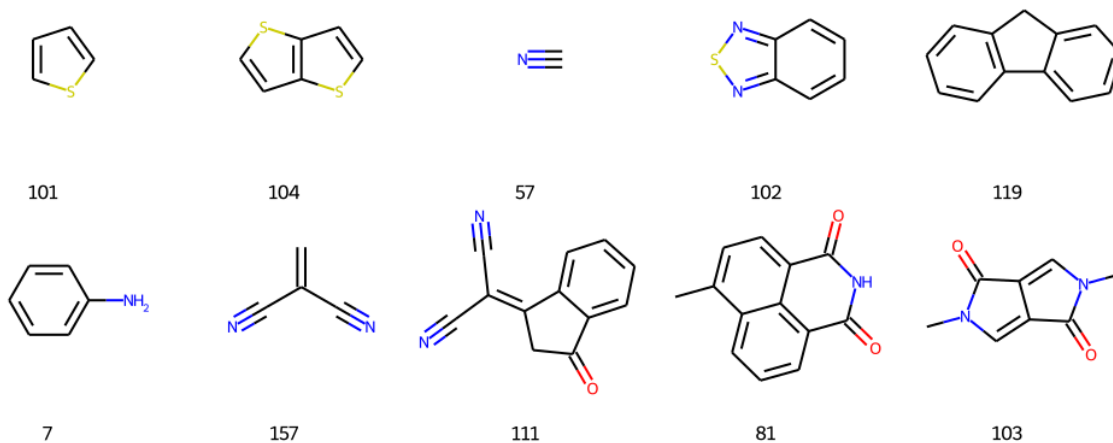
1287

1288 **Figure S8.** Histogram of the basis set employed in the geometrical optimization of the 479  
 1289 molecules found in the curated DFT database.



1290

1291 **Figure S9.** Histogram of the moieties present in the DFT database of 479 conjugated small  
 1292 molecules and oligomers. Moieties labels correspond to the chemical structures shown in  
 1293 Supplementary Note 4.



1294

1295 **Figure S10.** The 10 most frequent moieties together with their corresponding ID number.

1296

1297 Molecular descriptors are computed using four different open-source software packages and  
 1298 Python libraries (such as NumPy),<sup>108</sup> including 1D, 2D and 3D descriptors as retrieved from  
 1299 the corresponding DFT optimized geometries. The software bundle employed includes PaDEL  
 1300 (1874 descriptors),<sup>73</sup> PyChem (1094 descriptors),<sup>74</sup> Mordred (1826 descriptors)<sup>75</sup> and RDKit  
 1301 v2021.03.2. (1039 descriptors).<sup>76</sup> As a result, we obtained an initial set of 5834 descriptors for  
 1302 each molecule, which decreased up to 3239 after curing (i.e. dropping of uninformative or  
 1303 constant descriptors and others containing infinite or NaN values). Note that the same descriptor  
 1304 might be computed by more than one software bundle, yet slight numerical disagreements may  
 1305 arise due to the different computation algorithms. For that reason, we do not filter out redundant  
 1306 descriptors and perform their subsequent statistical analysis using the entire available catalogue.  
 1307 Furthermore, we include electronic features retrieved from the DFT calculations such as the  
 1308 energy levels of the 20 occupied and unoccupied molecular orbitals (HOMOs and LUMOs)  
 1309 closer to the band gap; the electronic band gap energy itself,  $E_{gap}$ ; and the number of  $\pi$  electrons  
 1310 ( $n_{\pi}$ ) in the molecule, which was determined using custom coding based on the RDKit library.  
 1311 The set of molecular fingerprints tested in this work is computed using RDKit and it includes  
 1312 customized coding for the moiety fingerprints and built-in functions for the computation of  
 1313 MACCS keys, Morgan fingerprints,<sup>88</sup> path-based (topological) fingerprints, E-state  
 1314 fingerprints<sup>89</sup> and Coulomb vectors.<sup>90</sup>

1315 The target features in this study focus on the maximum oscillator strength ( $f_{max}$ ) and other  
 1316 derived figures such as the maximum oscillator strength in the visible electromagnetic spectrum  
 1317 ( $f_{max,vis}$ , herein constrained between 300-1200 nm for its relevance in solar energy harvesting



1318 applications); the sum of  $f$  in the visible window,  $f_{max,vis}$ ; the spectral overlap between the  
1319 Gaussian-broadened spectrum of fs in the visible (taking a standard deviation of 0.1 eV) and  
1320 the AM1.5G solar irradiance spectrum,  $f_{overlap}$ ; the maximum absorption coefficient ( $\alpha_{max}$ );  
1321 the maximum of the imaginary part of the dielectric function ( $\epsilon_{2,max}$ ); and the maximum molar  
1322 extinction coefficient,  $f_{max}$ ,  $f_{max,vis}$  and  $f_{sum,vis}$  are also evaluated per number of  $\pi$  electrons  
1323 in the molecule, i.e.  $f_{max}/n_{\pi}$ ,  $f_{max,vis}/n_{\pi}$  and  $f_{sum,vis}/n_{\pi}$ .

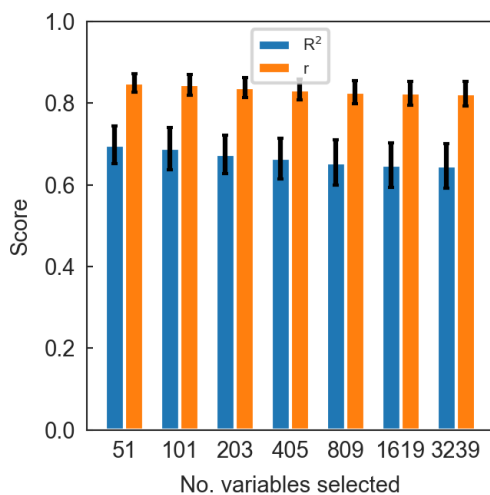
1324 The statistical analysis of descriptors is deployed using the open-source library SciPy<sup>109</sup>  
1325 whereas the machine-learning (ML) models (k-nearest neighbours, linear regression, support  
1326 vector regressor and random forests) are implemented in Scikit-Learn.<sup>87</sup>

1327 Regarding the scoring of the ML models,  $R^2$  ranges from  $-\infty$  to unity, being 1 the best possible  
1328 score and zero an indication of lack of predictive power (as it is always returning the expected  
1329 value of the target function, i.e., its average value);  $R_{adj}^2$  is formulated as<sup>110</sup>

1330 
$$R_{adj}^2 = 1 - (1 - R^2) \frac{n - 1}{n - p - 1},$$

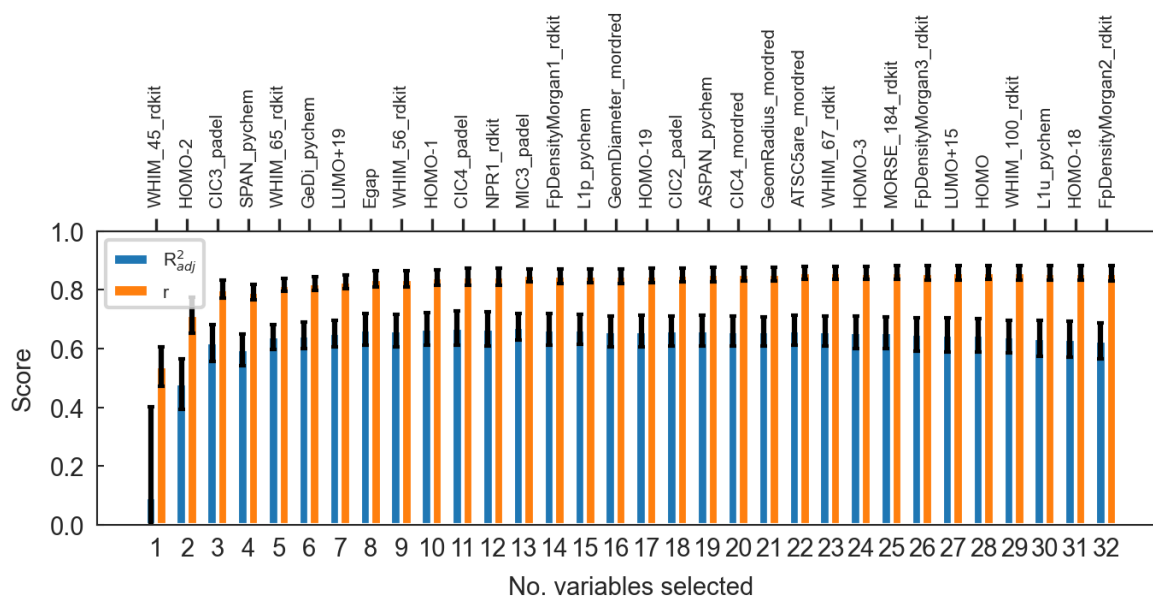
1331 where  $p$  is the number of variables and  $n$  the sample size. Thus,  $R_{adj}^2$  adds penalties if the model  
1332 uses too many variables, which is a useful metric when studying feature selection procedures.  
1333 Test sets comprise 30% of the available data and all models are 10-fold holdout cross-validated  
1334 (unless otherwise stated, using a randomized 70%-30% splitting for the train and test sets,  
1335 respectively).

1336 The recursive feature elimination (RFE) procedure applied in this work starts by decreasing the  
1337 number of input features to 32 (i.e., 1% of the starting descriptor population of 3239 descriptors)  
1338 in six consecutive feature reduction steps, in which after performing successive 10-fold cross-  
1339 validations we drop 50% of the (averaged and least important) descriptors. Rather than  
1340 observing a performance drop, the actual scoring of the RF ensemble improves as the number  
1341 of features is reduced from 3239 ( $R^2 = 0.65 \pm 0.06, r = 0.82 \pm 0.03$ ) to 51 ( $R^2 =$   
1342  $0.70 \pm 0.05, r = 0.85 \pm 0.02$ ) variables in the last RFE iteration (**Figure S11**). After the last  
1343 pruning step (51 variables), we select the 32 most important descriptors and perform a more  
1344 thorough feature selection procedure by successively dropping (one-by-one) the least important  
1345 descriptor (always keeping a 10-fold cross-validation scheme, see **Figure S12**).



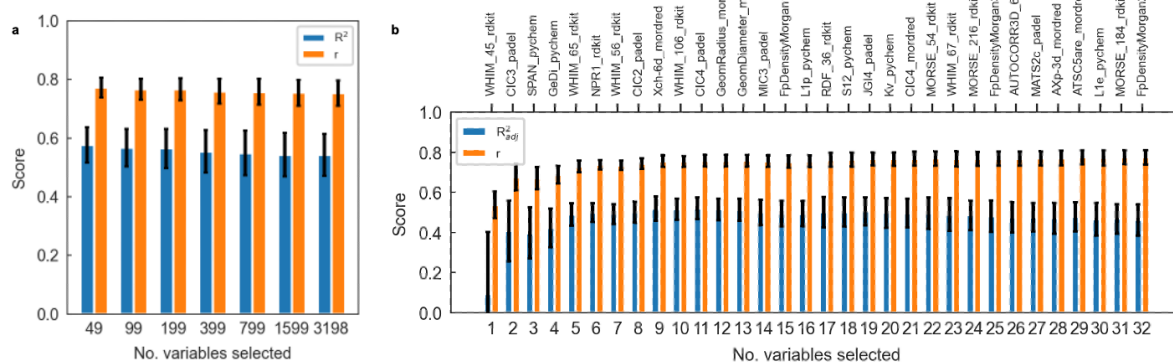
1346

1347 **Figure S11.** Scoring of RF regressors as part of a recursive feature elimination (RFE) loop in  
 1348 which 50% of the least important descriptors are dropped at each step.



1349

1350 **Figure S12.** Scoring of 10-fold cross-validated RF regressors (300 estimators) trained and  
 1351 tested using different amounts of input descriptors as progressively indicated by the RFE  
 1352 algorithm. The top axis indicates, from left to right, the name of the variable that is added to the  
 1353 RF model, thus forming an ordered list of the most important descriptors found by the RF  
 1354 method.



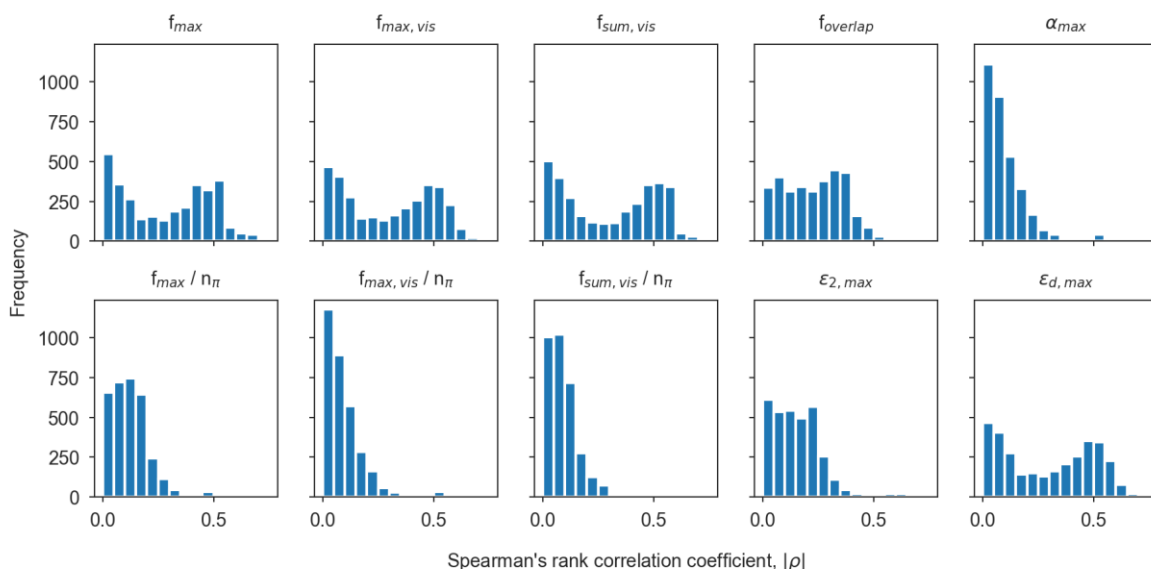
1355

1356 **Figure S13.** Performance of RF regressors trained without including electronic descriptors,  
 1357 using 300 estimators and 10-fold cross-validation. (a) Scoring parameters of cross-validated  
 1358 models as part of the RFE algorithm. (b) Detailed scoring parameters of the last 32 models  
 1359 obtained by RFE.

1360

1361 **Supplementary Note 3.** Clusterization algorithm of multicollinear descriptors.

1362 The clusterization algorithm starts by taking the descriptor with the highest Spearman's rank  
1363 correlation coefficient ( $\rho$ ) and computing the Pearson correlation coefficient ( $r$ ) with respect to  
1364 the remaining elements in the  $\rho$ -ordered list of descriptors with  $\rho \geq 0.7$ . Descriptors from this  
1365 list are dropped if  $r \geq 0.7$  and considered to be in the same cluster; those showing  $r \leq 0.7$  are  
1366 candidates to form a different cluster. The process runs in a recursive-elimination manner until  
1367 naturally leading to a selection of (typically) 1 to 5 descriptor clusters depending on the selected  
1368 thresholds (0.6-0.7). These clusters gather the most statistically relevant and monotonic  
1369 correlation trends with the target feature. Interestingly, by looking at the features stored in each  
1370 of the clusters it is possible to replace some of the descriptors found originally by the algorithm  
1371 by alternative figures of easier interpretation and/or larger physicochemical relevance.

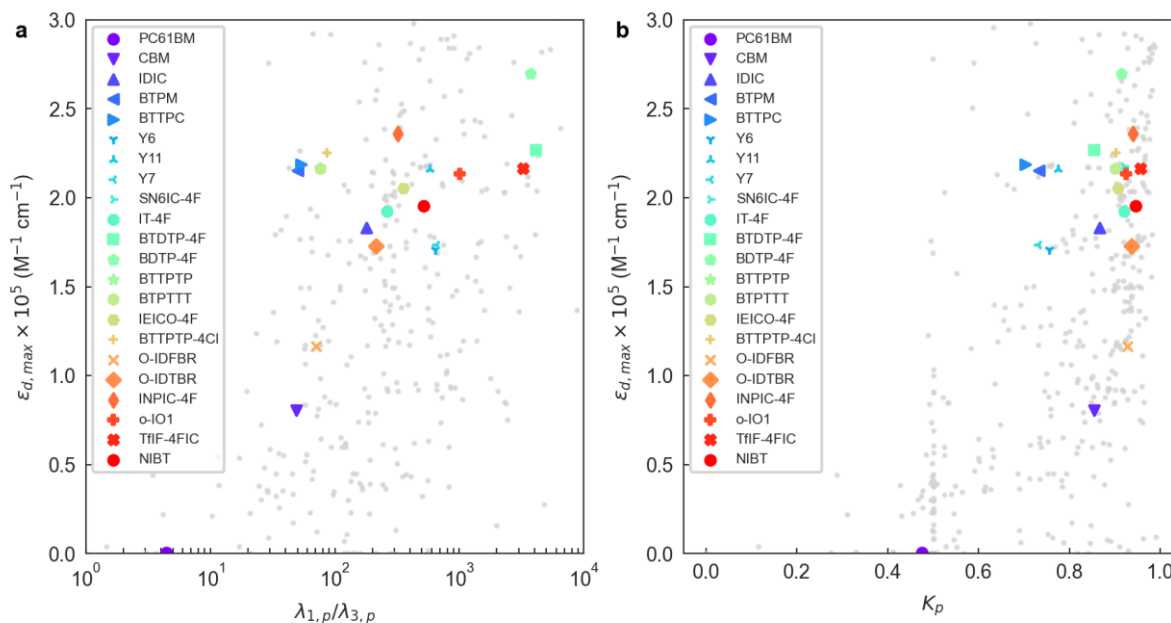


1372

1373 **Figure S14.** Spearman's rank correlation coefficient (in absolute value) histograms for the 3239  
1374 descriptors and the 10 different target features related with optical absorption and oscillator  
1375 strength explored in this work.

1376

1377



1378

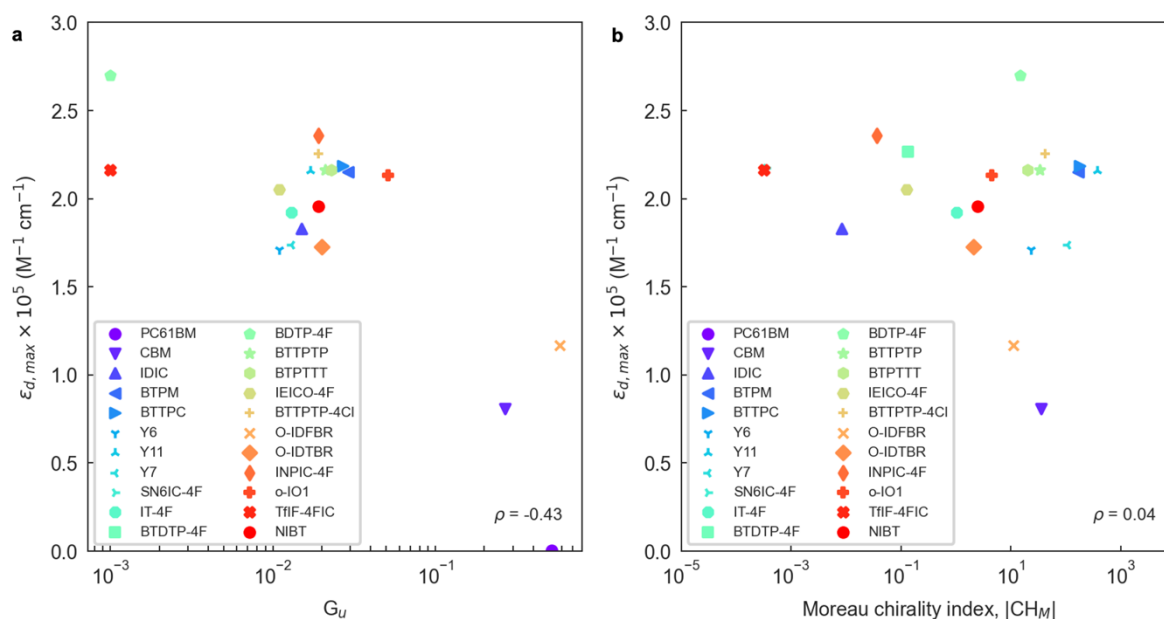
1379 **Figure S15. Influence of the molecular planarity on the maximum molar extinction**  
 1380 **coefficient.** (a) The  $\lambda_{1,p}/\lambda_{3,p}$  ratio correlates positively with  $\epsilon_{d,max}$  as straighter (more linear)  
 1381 molecules show larger  $\lambda_{1,p}$  while enhanced molecular planarity lowers  $\lambda_{3,p}$  values (which  
 1382 approach zero as there is no variance out of the molecular plane). (b) The shape of the molecule  
 1383 quantified with  $K_p$  shows that linear and planar molecules (i.e.,  $K_p$  closer to unity)<sup>82</sup> enable  
 1384 larger  $\epsilon_{d,max}$  values.  $K_p$  is defined as<sup>82</sup>

1385

$$K_p = \frac{\sum_m \left| \frac{\lambda_{m,p}}{\sum_m \lambda_{m,p}} - \frac{1}{3} \right|}{4/3},$$

1386 where  $m = 1,2,3$  and  $0 \leq K_p \leq 1$ .

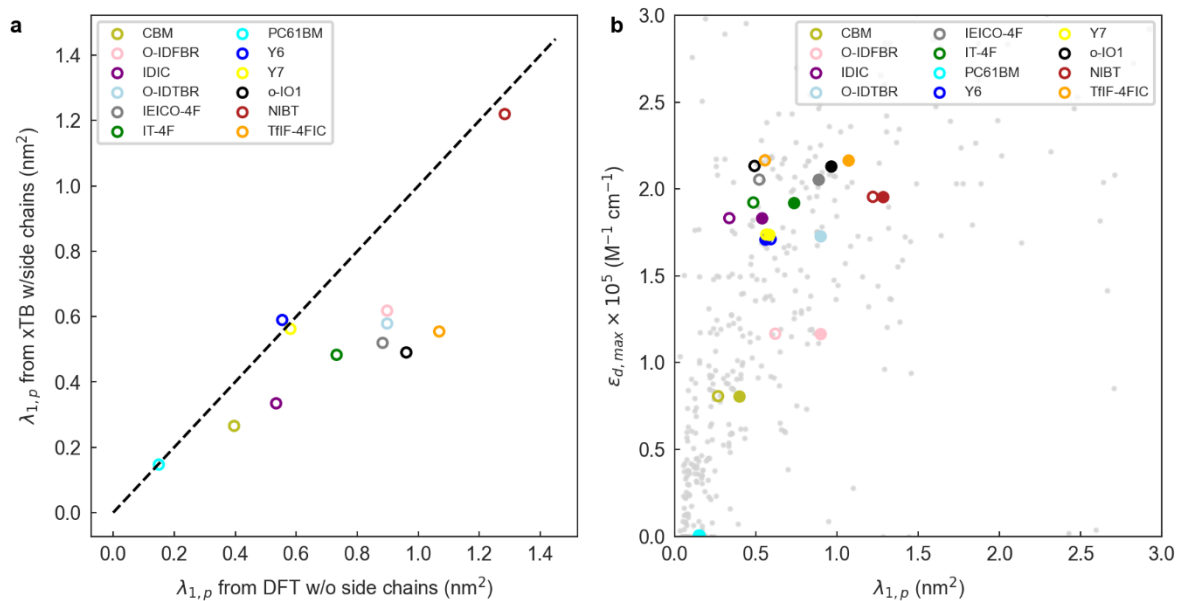
1387



1388

1389 **Figure S16. Influence of molecular symmetry on absorption strength.** (a) Quantification of  
 1390 the total molecular symmetry as per the definition of the WHIM symmetry descriptor  $G_u$   
 1391 (corresponding to the unweighted geometric mean of the directional symmetries,  $G_u =$   
 1392  $\sqrt[3]{\gamma_{1,u} \cdot \gamma_{2,u} \cdot \gamma_{3,u}}$ )<sup>82</sup> shows that as the molecules lose their central symmetry (i.e., lower  $G_u$   
 1393 values),  $\epsilon_{d,max}$  can be further enhanced. (b) Conversely, the Moreau chirality index<sup>111</sup> weighted  
 1394 by atomic coordinates of small molecular absorbers shows poor correlation with  $\epsilon_{d,max}$ .

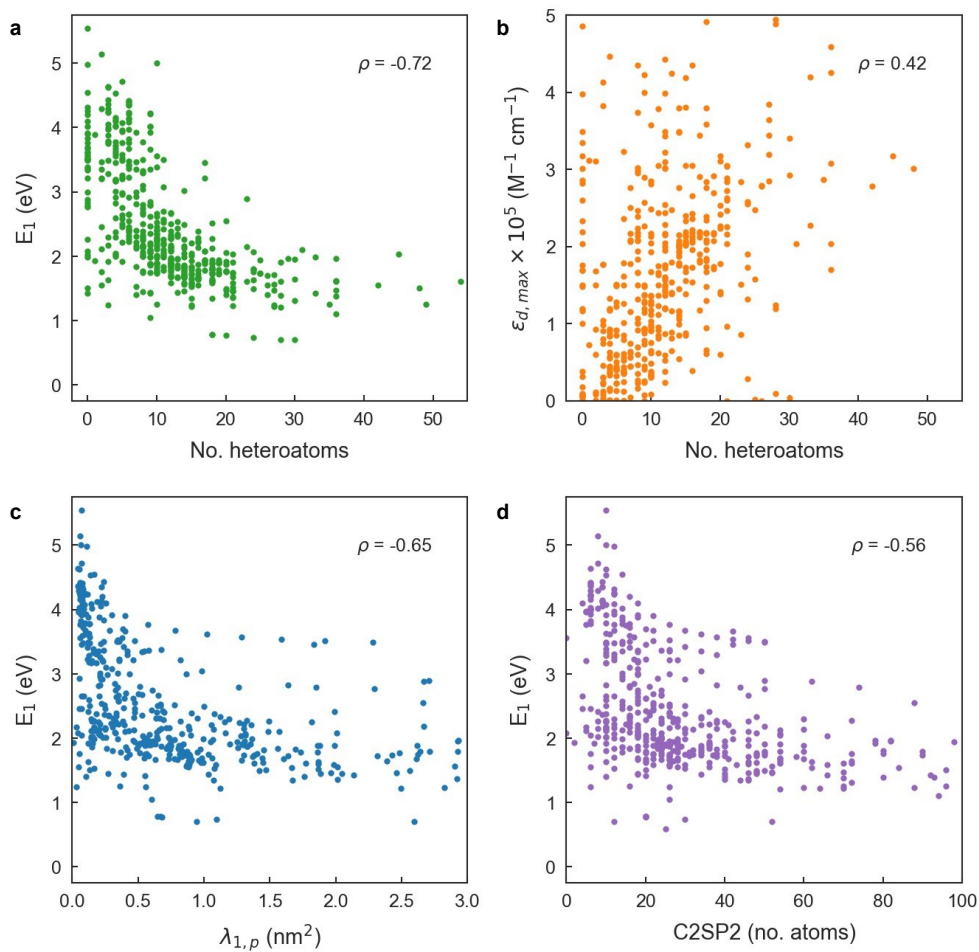
1395



1396

1397 **Figure S17. Influence of the side chains on  $\lambda_{1,p}$  values.** (a) Comparison of  $\lambda_{1,p}$  values for a  
 1398 selection of small molecule acceptors as computed from xTB including side chains (y axis) and  
 1399 DFT-optimized geometries with methyl-substituted side chains (x axis). (b) Maximum molar  
 1400 extinction coefficient ( $\epsilon_{d,max}$ ) as a function of  $\lambda_{1,p}$  for small molecule acceptors optimized with  
 1401 (open circles) and without (filled circles) side chains.

1402

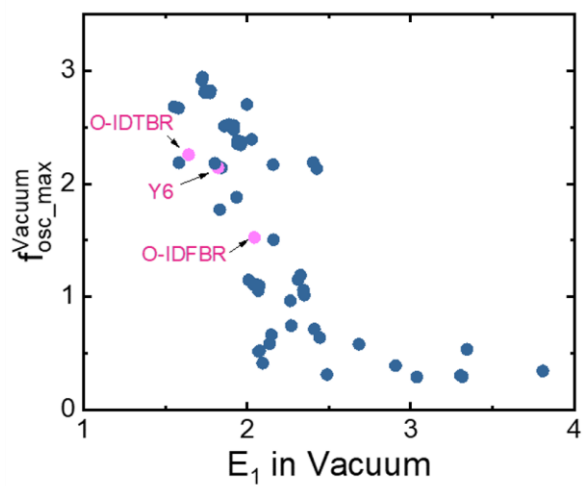


1403

1404 **Figure S18.** (a) Correlation between  $E_1$  and the number of heteroatoms in the molecules. (b)  
 1405 Correlation between the molar extinction coefficient ( $\epsilon_{d,max}$ ) and the number of heteroatoms.  
 1406 (c) Correlation between  $E_1$  and  $\lambda_{1,\rho}$ . (d) Correlation between  $E_1$  and C2SP2. All panels include  
 1407 the corresponding Spearman's rank correlation coefficient ( $\rho$ ).

1408





1409

1410 **Figure S19.** Relationship between the maximum oscillator strength and the energy of the first

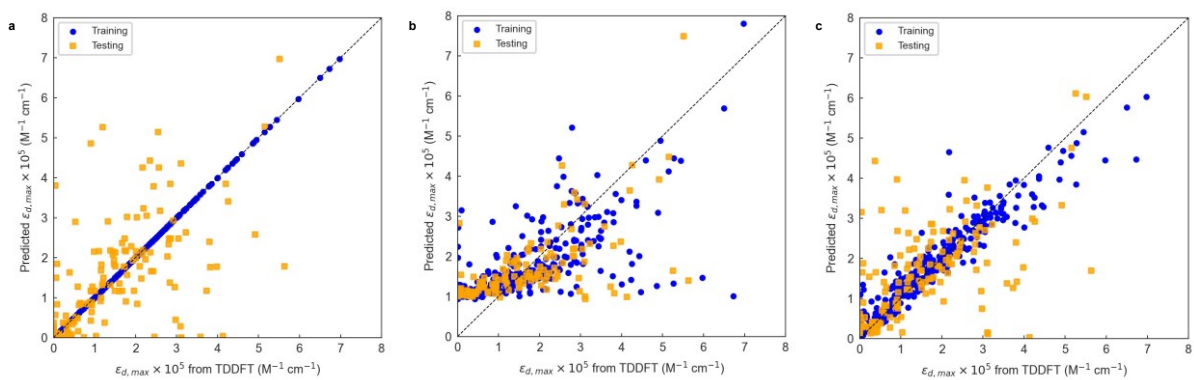
1411 electronic transition in a set of TDDFT-optimized NFAs.

1412

1413 **Table S1.** Statistical performance of a manifold of 10-fold cross-validated baseline models  
1414 using  $\varepsilon_{d,max}$  as target feature.

<b>Model</b>	<b>No. variables</b>	<b>R<sup>2</sup></b>	<b>r</b>
1-nearest neighbour	2	$-0.18 \pm 0.41$	$0.50 \pm 0.05$
	3239	$-0.10 \pm 0.37$	$0.47 \pm 0.09$
Linear regression	2	$0.37 \pm 0.10$	$0.61 \pm 0.09$
Random forest w/300 estimators	2	$0.23 \pm 0.17$	$0.59 \pm 0.06$
	3239	$0.65 \pm 0.06$	$0.82 \pm 0.03$

1415



1416

1417 **Figure S20.** Correlation plots of three (exemplary) baseline models trained and tested on  $\epsilon_{d,max}$   
 1418 using two input descriptors only:  $\lambda_{1,p}$  and C2SP2. (a) 1-nearest neighbour; (b) linear regression;  
 1419 and (c) out-of-the-box RF trained with 300 estimators.

1420

1421 **Table S2.** Performance of RF models trained and 10-fold cross-validated using 300 estimators,  
 1422 3 input molecular descriptors ( $\lambda_{1,v}$ , CIC3 and HOMO-2) and different forms of molecular  
 1423 fingerprint vectors. In the case of Morgan fingerprints, we set the connectivity radius to 4 units,  
 1424 while for topological fingerprints the minimum and maximum path counts are set to 1 and 6  
 1425 units, respectively. Their vector lengths are set to either 64 or 2048 bits to reflect different  
 1426 degrees of model complexity.

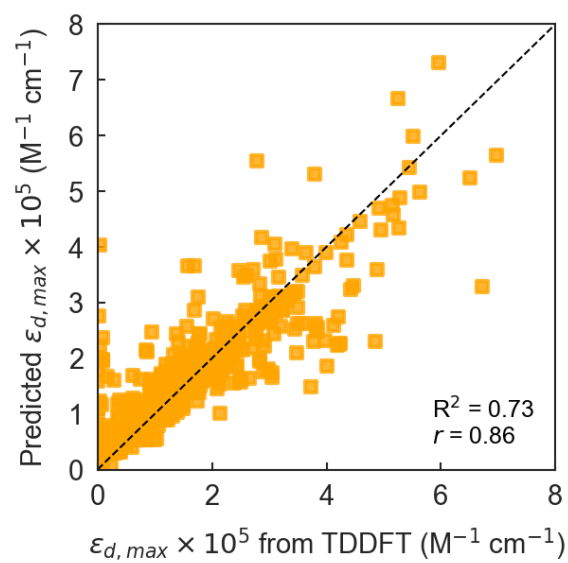
No. molecular descriptors	Fingerprint type	No. bits	Total no. inputs	R <sup>2</sup>	r
3	N/A	N/A	3	0.63 ± 0.06	0.80 ± 0.03
3	Moiety	159	162	0.63 ± 0.06	0.81 ± 0.04
3	MACCS	166	169	0.66 ± 0.04	0.83 ± 0.02
3	Morgan	64	67	0.70 ± 0.05	0.84 ± 0.03
		2048	2051	0.69 ± 0.05	0.84 ± 0.03
3	Topology	64	67	0.68 ± 0.04	0.83 ± 0.02
		2048	2051	0.69 ± 0.05	0.84 ± 0.03
3	E-state	79	82	0.66 ± 0.04	0.82 ± 0.02
3	Coulomb	320	323	0.56 ± 0.07	0.77 ± 0.04

1427

1428 **Table S3.** Scoring of the baseline and hyperparametrically optimized RF and ExtraTrees  
 1429 models, fed with 3 molecular descriptors and a Morgan fingerprint vector of 64 bits.

Model	No. estimators	No. samples per leaf	No. samples to split	Validation	R <sup>2</sup>	r
RF (out-of-the-box)	300	1	2	10-fold CV	0.70 ± 0.05	0.84 ± 0.03
RF (optimized)	1200	1	2	10-fold CV	0.70 ± 0.05	0.85 ± 0.03
RF (optimized)	1200	1	2	LOOCV	0.74	0.86
ExtraTrees (out-of-the-box)	300	1	2	10-fold CV	0.69 ± 0.05	0.85 ± 0.02
ExtraTrees (optimized)	2000	1	2	10-fold CV	0.70 ± 0.04	0.85 ± 0.02
ExtraTrees (optimized)	2000	1	2	LOOCV	0.73	0.86

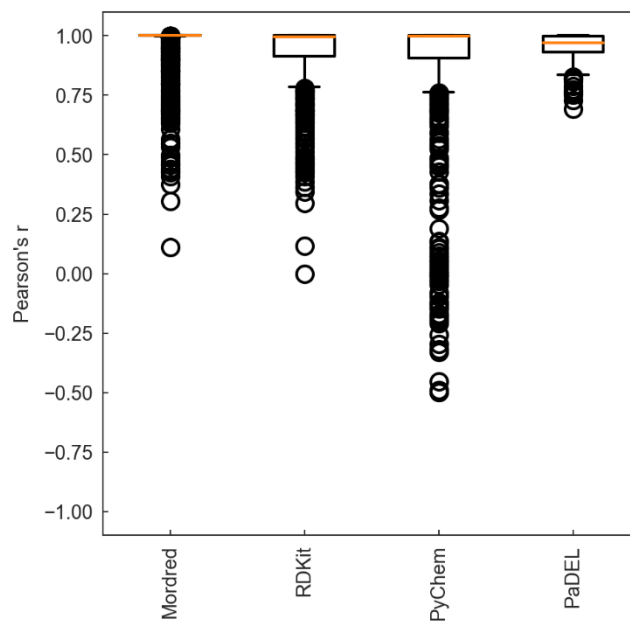
1430



1431

1432 **Figure S21.** LOOCV of the optimized Extra Trees (ET) regressor fed with 3 molecular  
1433 descriptors and a 64-bit vector as Morgan fingerprint.

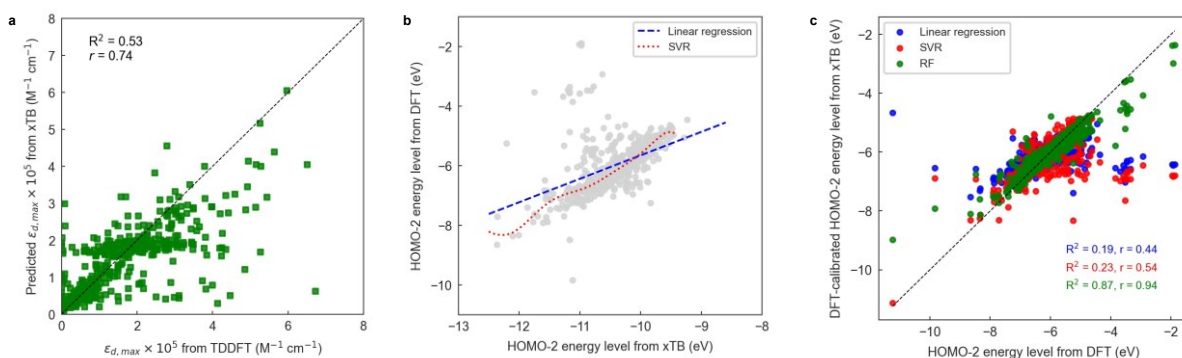
1434



1435

1436 **Figure S22.** Boxplots for the Pearson correlation coefficients between different sets of  
1437 molecular descriptors retrieved from xTB and DFT (B3LYP) optimized geometries.

1438



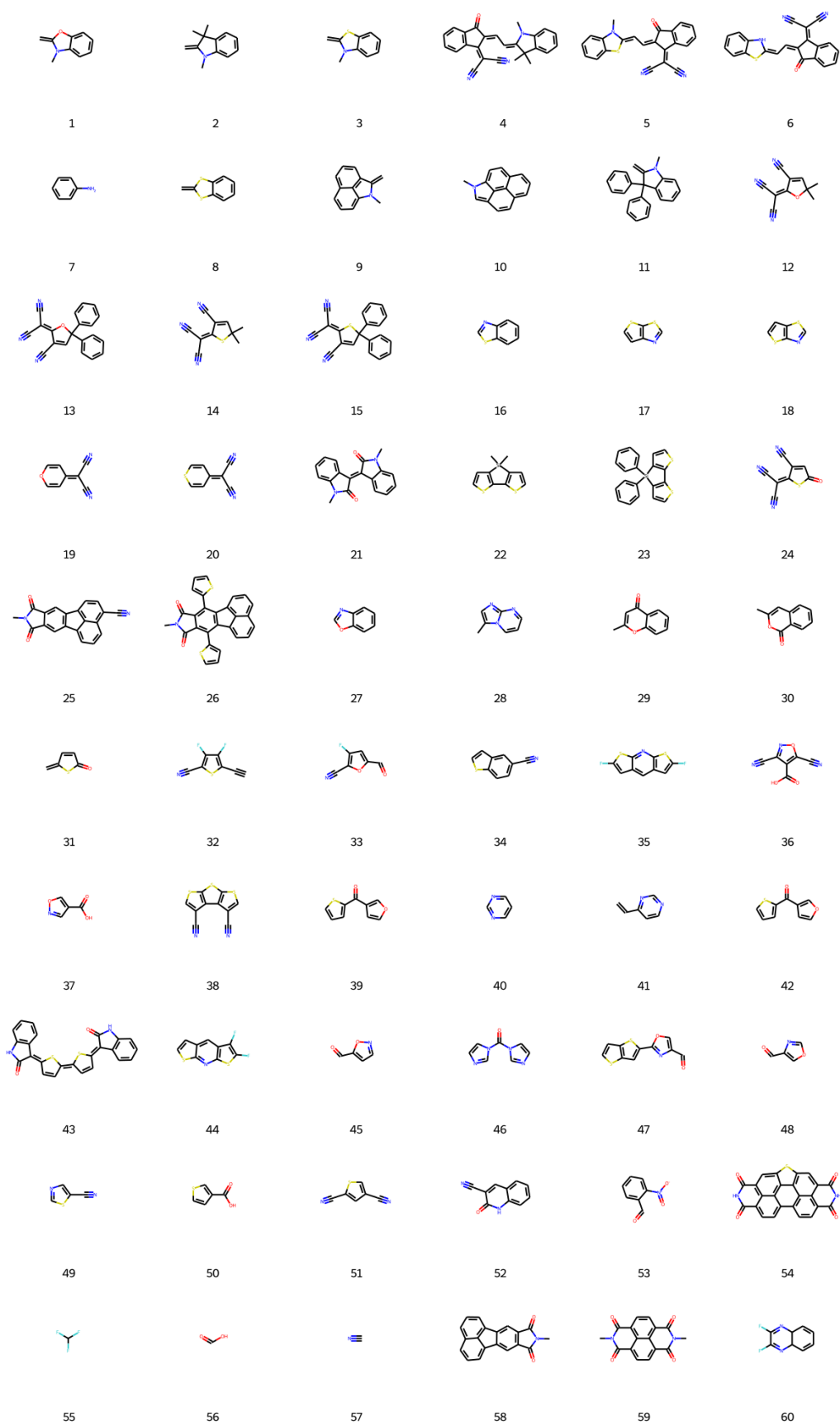
1439

1440 **Figure S23.** (a) Leave-one-out interpolation of a RF model trained using DFT data and tested  
 1441 on xTB-optimized molecules using 3 parameters ( $\lambda_{1,v}$  and CIC3 recalculated from the xTB  
 1442 geometry, and the HOMO-2 energy level as computed in xTB) and a 64-bit vector as Morgan  
 1443 fingerprint. (b) Fitting of linear regression and support vector regressor (SVR) models for  
 1444 calibration of HOMO-2 energy levels as computed in xTB and DFT (B3LYP). The mismatch  
 1445 in the absolute values of HOMO-2 energy levels between DFT and xTB calculations prevents  
 1446 obtaining higher scorings in the RF models depicted in (a). (c) Correlation plot between  
 1447 HOMO-2 energy levels from DFT and the corresponding calibrated values as obtained by linear  
 1448 regression (blue), SVR (red) and RF (green) models. The dashed black line indicates perfect  
 1449 matching between DFT and calibrated values.

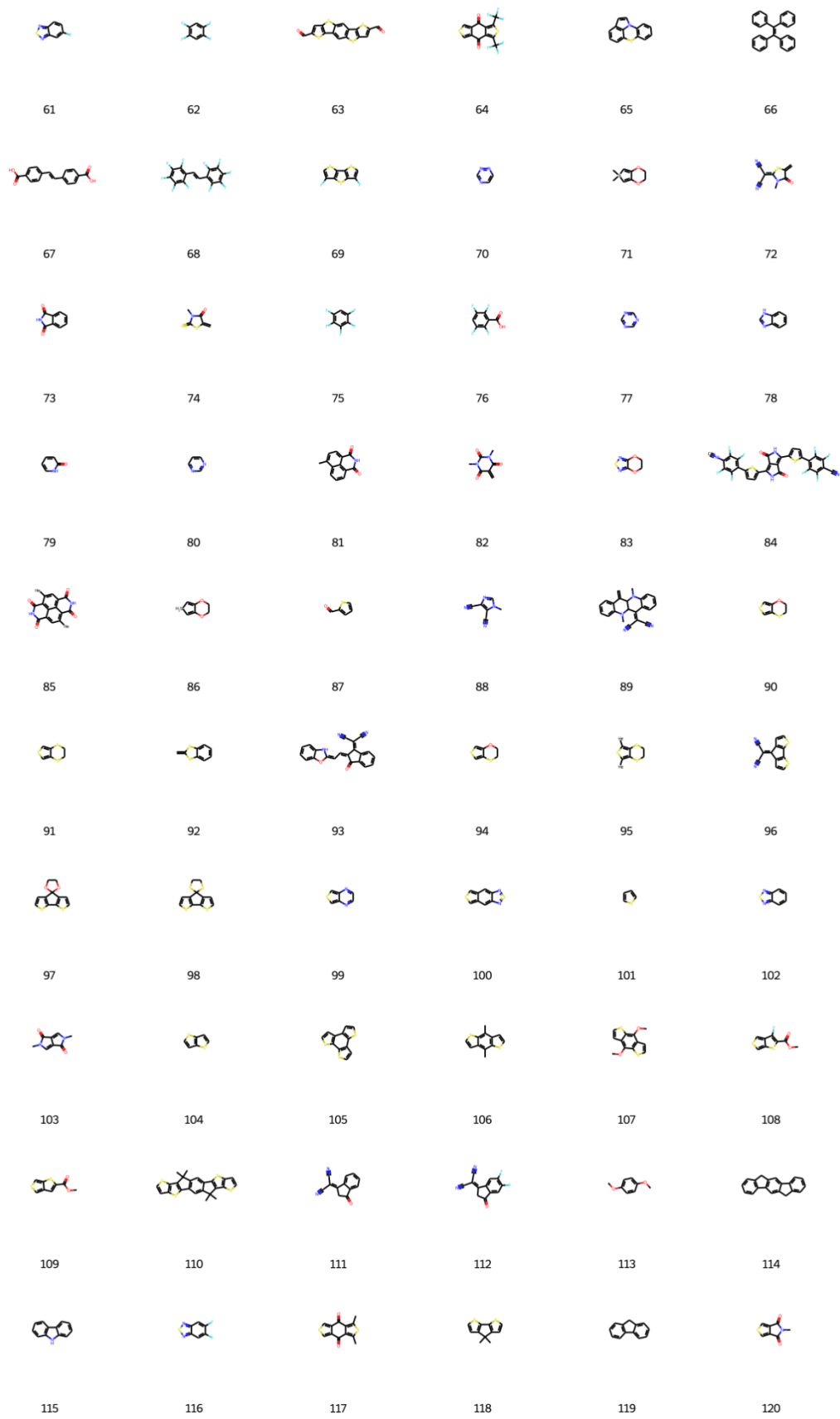
1450



1451 **Supplementary Note 4.** Detailed database of moieties used in this work.

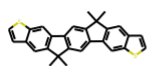


1452



1453

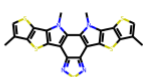
1454



121



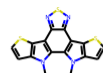
122



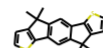
123



124



125



126



127



128



129



130



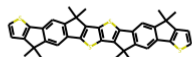
131



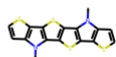
132



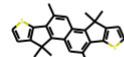
133



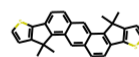
134



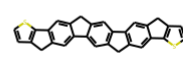
135



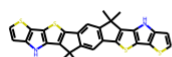
136



137



138



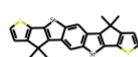
139



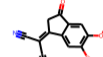
140



141



142



143



144



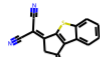
145



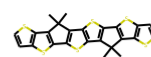
146



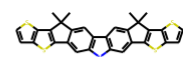
147



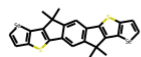
148



149



150



151



152



153



154



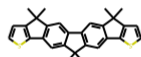
155



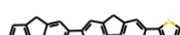
156



157



158



159

1455

1456

1457

1458 **Supplementary Note 5.** Estimation of computation time required to make absorption strength  
1459 predictions using xTB Hamiltonians (in combination with ML models) or rigorous TDDFT.

1460 **Table S4** provides a comparison in terms of the computation time required for the molecular  
1461 geometry optimization step in the TDDFT and xTB approaches. We analysed 194 molecules  
1462 from TDDFT calculations based on B3LYP/6-311+G(d,p) and 475 molecules (**Figure S7**) from  
1463 xTB calculations based on GFN2-xTB to tentatively quantify the difference in computational  
1464 efficiency between both methods. Our analysis suggests that geometry optimization using  
1465 GFN2-xTB is ca. 3000 times faster than TDDFT/B3LYP/6-311+G(d,p), even though GFN2-  
1466 xTB calculations were done on a conventional 12 CPUs laptop as opposed to the 32 CPUs  
1467 dedicated cluster/workstation employed in the TDDFT calculations, thus highlighting the great  
1468 advantage of using xTB over TDDFT.

1469 Furthermore, we have estimated the time consumption for the absorption strength predictions  
1470 using the established ML model in this work. The time required for the ML model training and  
1471 LOOCV steps is below 10 minutes (12 CPUs), whereas the calculation of molecular descriptors  
1472 (>5000 descriptors) for the full data set (479 molecules) takes no less than 180 minutes (12  
1473 CPUs). Hence, for a molecule made up of 100 atoms, the whole absorption strength  
1474 determination (i.e., from geometry optimization to  $\epsilon_{d,max}$  prediction) effectively takes around  
1475 200 minutes using xTB with ML; and 1345 minutes using solely TDDFT. Nevertheless, the  
1476 advantage of the ML approach is more evident as interpolation in the trained model takes less  
1477 than 1 second (per molecule) to compute, which enables at least four orders of magnitude faster  
1478 molecular screening with respect to TDDFT (1345 minutes or 80700 seconds per molecule).

1479 **Table S4.** Computation time required for molecular geometry optimization steps using  
1480 TDDFT/B3LYP/6-311+G(d,p) and xTB/GFN2-xTB.

<b>Approach</b>	DFT/B3LYP/6-311+G(d,p)	xTB/GFN2-xTB
<b>No. of molecules</b>	194	475
<b>No. of atoms</b>	18022	37989
<b>No. of CPUs</b>	32	12
<b>Time elapsed (mins)</b>	206398.355	136
<b>Time elapsed per atom (mins)</b>	11.45257768	0.003579984

1481