

Finding the Right Bricks for Molecular Lego: A Data Mining Approach to Organic Semiconductor Design

Christian Kunkel, Christoph Schober, Johannes T. Margraf, Karsten Reuter, and
Harald Oberhofer*

*Chair for Theoretical Chemistry and Catalysis Research Center, Technische Universität
München, Lichtenbergstraße 4, D-85747 Garching, Germany*

E-mail: harald.oberhofer@tum.de

Abstract

Improving charge carrier mobilities in organic semiconductors is a challenging task that has hitherto primarily been tackled by empirical structural tuning of promising core compounds. Knowledge-based methods can greatly accelerate such local exploration, while a systematic analysis of large chemical databases can point towards promising design strategies. Here, we demonstrate such data mining by clustering an in-house database of > 64.000 organic molecular crystals for which two charge-transport descriptors, the electronic coupling and the reorganization energy, have been calculated from first principles. The clustering is performed according to the Bemis-Murcko scaffolds of the constituting molecules and according to the sidegroups with which these molecular backbones are functionalized. In both cases, we obtain statistically significant structure-property relationships with certain scaffolds (sidegroups) consistently leading to favorable charge-transport properties. Functionalizing promising scaffolds

with favorable sidegroups results in engineered molecular crystals for which we indeed compute improved charge-transport properties.

1 Introduction

In the last few years, organic semiconductors have received considerable attention due to their comparatively small ecological and economical footprint, as well as due to their great versatility promising a wide spectrum of novel material properties.¹⁻⁵ Their potential uses include organic field effect transistors (OFETs),⁶ photovoltaics (OPVs),⁷ light emitting diodes (OLEDs)⁸ or even sensors.⁹ A major limitation to the commercial application of this promising class of materials has been the failure to reproducibly yield high charge-carrier mobilities as a prerequisite for the electronic performance.^{10,11} Thanks to intense research efforts,¹²⁻¹⁵ a number of materials now exist with mobilities exceeding even that of amorphous silicon.¹⁶ Yet, the sheer endless materials spaces constituted by assembling building blocks of functionalized organic molecules raise the suspicion that these few materials are nothing but the top of an iceberg.

So far, experimental discovery of new materials commonly follows a cycle of synthesis, device manufacture, performance evaluation and subsequent molecular tuning, based on chemical intuition, experience and—to a certain degree—trial and error.^{5,6,17} Theoretical investigation is then often used to understand and exploit the key charge-transport mechanisms.¹⁸⁻²² Although highly successful, such efforts are limited by the pace of iterative improvement, while only locally exploring the chemical space around a compound family. As a more recent development, a number of groups have therefore started to conduct *in silico* materials discovery calculations,²³⁻²⁸ to guide synthetic studies to promising targets. In these screening approaches, readily available or easily calculable properties of the system—so called descriptors—which are correlated with critical observables of the system (such as the carrier mobility) are determined for a larger set of candidate materials. These then

yield a ranking of the screened candidate materials based on their suitability as an organic semiconductor.

In other fields of science, such as heterogeneous catalysis,^{29–31} superhard materials,³² or drug discovery,³³ such screening studies are already followed up with modern large-scale data mining or even machine learning approaches to extract more general design criteria.³⁴ Corresponding approaches undertaken for organic semiconductors (OSCs) have so far been limited to a much smaller scale.^{35,36} In this work, we now take this step and combine the statistical tools of modern data mining approaches with the molecular analysis methods known as cheminformatics to uncover guiding principles for the design of organic semiconducting materials.

1.1 Descriptor-Based Data Mining of OSCs

Two of the most prominent charge-transport descriptors for organic semiconductors are the network of electronic couplings $|H_{ab}|$ between microscopic sites involved in the charge transport and the reorganization energy λ , which measures the cost of accommodating a new charge state after the carrier has moved to the next site. While $|H_{ab}|$ and λ are most closely associated with the small polaron hopping mechanism of carrier transport known from Marcus theory,^{37,38} both also significantly influence the carrier mobility in other transport regimes.³⁹ In evaluating λ , usually only the dominant local, intra-molecular reorganization energy is considered.⁴⁰ In contrast to $|H_{ab}|$, the reorganization energy can then be estimated cheaply as a single site property, while at the same time decisively—in Marcus-like models even exponentially—determining the carrier mobility.^{18,41}

A number of screening studies thus initially focused on the optimization of λ through modification of molecular sidegroups in, for example, dinaphtho-thienothiophene (DNTT),⁴² triarylamine-based donor molecules for solar cells,⁴³ or high-performance polymer semiconductors.⁴⁴ Taken together, the studies impressively demonstrated how *in silico* exploration of material spaces can efficiently identify high-performance organic semiconductors. While

similar approaches are now also starting to surface in experimental studies,⁴⁵ most recent *in silico* design studies start to optimize λ in newly generated derivatives. The design problem is thereby mainly tackled by construction and computational screening of highly focused molecular libraries based on a number of different strategies. These range from the generation of polymer semiconductors from preselected building blocks^{41,46} to the optimization of λ through heteroatom replacement,^{47,48} introduction of sidegroups^{41,48–55} or ring fusion.^{56,57}

Modern data-driven approaches instead pursue a complementary strategy. Rather than screening a focused library for potential candidate materials, in particular data mining techniques unveil correlations in large data sets and thereby provide leads for a rational optimization beyond the originally considered design space.^{58–60} Size and especially diversity of the underlying data source are thus key for the generality and transferability of the identified correlations. A suitable representation of the data critically determines the efficiency with which these correlations can be extracted, whether by machine learning or more traditional fitting techniques. In several fields, development and application of such knowledge-based approaches has already led to a greatly accelerated materials design.^{61–63}

Here, we show, that a data-mining rooted in the chemical understanding of the structure—often termed chemical intuition—can indeed yield general design principles for organic semiconductors. To this end, we concentrate on the class of crystalline molecular organics and analyze an in-house dataset containing > 64.000 experimentally characterized mono-molecular crystals. This dataset was originally assembled to computationally screen for hitherto unknown organic crystals with favorable charge-transport properties.²⁵ As further described below, it contains an unprecedented molecular diversity and a wealth of promising semiconductors as identified on the basis of computed high electronic couplings $|H_{ab}|$ and low reorganization energies λ . In this work, we now perform a two-step data mining procedure on this dataset. Using the concept of Bemis-Murcko scaffolds⁶⁴ that has already been successfully used in pharmaceutical research, we establish a useful data representation that uncovers statistically significant performance differences between the individual molecular

building blocks of the crystals. Clustering the data first according to their Bemis-Murcko scaffolds, i.e. the molecular backbones, we find clear structure-property relationships, with a range of scaffolds consistently leading to favorable charge-transport properties, despite the positional disorder brought about by differing sidegroups and anchor points on them. In a second step, we then cluster the data according to the sidegroups that are attached to the molecular scaffolds. Again, we obtain clear relationships with a range of sidegroups generally lowering the reorganization energy, regardless of the scaffold they are attached to. These findings suggest the combination of certain, high-coupling scaffolds with reorganization energy-reducing sidegroups as a promising design criterion for optimized charge transport in organic molecular crystals. We validate the transferability and accuracy of this criterion by a corresponding sidegroup-engineering of molecular testsets and indeed find significant improvements in the charge-transport properties.

2 Methods

2.1 Dataset

Our study is based on the “64k-dataset”,²⁵ containing information on 64.725 organic crystals composed of 61.770 unique molecular structures, extracted from the Cambridge Structural Database⁶⁵ and screened for favorable charge-transport related descriptors λ , $|H_{ab}|$. All crystals are mono-molecular, i.e. they contain only a single type of molecule. For a cheminformatics based analysis, *xyz*-coordinates of the molecular structures were extracted from the crystal and converted to canonical smiles strings⁶⁶ using Openbabel.⁶⁷ All subsequent molecular analysis was carried out using the RDKit.⁶⁸ Marvin was used for drawing chemical structures.⁶⁹ Pymatgen was used to analyze the occurrence of molecular point groups and crystal space groups.⁷⁰

2.2 Electronic Structure Calculations

The computational settings and methodology were described in detail in the preceding work that established the “64k-dataset”.²⁵ Thus, we here only briefly summarize them for completeness. Charge-transport related descriptors λ and $|H_{ab}|$ were computed using first-principles Density-Functional Theory (DFT) at the GGA (BLYP^{71,72}) level of theory. All calculations were done using the FHI-aims package^{73,74} employing its well-established basis-sets and integration grid settings, detailed below in this section. Our work on the “64k-dataset” showed that even semi-local DFT was able to recover trends among different crystals for both $|H_{ab}|$ ⁷⁵ and λ .²⁵ For each crystal in the dataset, the nearest-neighbor molecular dimers were first extracted by a purpose-built python program. Electronic couplings $|H_{ab}|$ for each dimer were then calculated using the fragment molecular orbital approach (FO-DFT)⁷⁶ in the $H^{2n-1}@D + A$ variant.⁷⁵ The electronic wave functions were expanded in an extended “tier 1” basis set of FHI-aims using light integration settings. The highest electronic coupling among all dimers was extracted from this data and constitutes the $|H_{ab}|$ entering the dataset.²⁵ Keeping in mind that the “64k-dataset” consists solely of crystals composed of a single species of molecule each, the internal reorganization energy was evaluated for each molecule based on the 4-point scheme,⁷⁷ using a non-periodic QM/MM embedding setup.²⁵ In this ONIOM-based scheme,⁷⁸ a shell of nearest neighboring molecules mimics the steric hindrance due to a realistic solid-state environment during geometry optimization of a central molecule. The shell of nearest neighbours was thereby extracted based on the minimal distance found between the atoms of any two molecules in the crystal. All molecules that approach the central molecule to within this distance times 1.5 were included in the neighbor shell. Then, keeping the shell of nearest neighbors fixed, the local structure optimisation of the central molecule was performed using the Broyden-Fletcher-Goldfarb-Shanno (BFGS) implemented in the Atomic Simulation Environment⁷⁹ using a convergence criterion $f_{\max} < 0.05$ eV/Å. The MM parts of these calculations were based on the LAMMPS package⁸⁰ with the “Universal Force Field” (UFF)⁸¹ and the QEq charge equilibration scheme⁸² in

order to determine interaction parameters and fixed partial charges, respectively. To ensure maximum accuracy of our main descriptor, we re-calculated all reorganization energies λ in the present database using default "tight" settings and the respective "tier 2" basis-sets of FHI-aims.

2.3 Statistical Methods

We used the two-tailed Mann-Whitney-U test⁸³ to assess whether two samples are drawn from the same distribution (null hypothesis) or not (alternative hypothesis). Note that the Mann-Whitney-U test is non-parametric and therefore avoids assumptions on the underlying distributions. To account for the great number of statistical tests we carried out, we corrected the computed p -values for each descriptor using the false discovery rate correction of Benjamini-Hochberg,⁸⁴ as implemented in the statsmodels package.⁸⁵ These corrected p -values are then named q -values. Results were marked statistically significant, if the q -value was found to be below 0.05, i.e. if the probability that the statistical testing wrongly rejects the null hypothesis is smaller than 5%. Similar methods and criteria are commonly used in Gene Set Enrichment Analysis⁸⁶ and have also been applied to cheminformatics^{87,88} for the identification of active compound series.

3 Results

3.1 Molecular Scaffolds

Compared to earlier, focused studies on OSCs, the dataset used in this study covers a wide range of molecular diversity to ensure the generality of our data mining results. This is reflected by 83 crystal space groups occurring at least 3 times in the dataset, while the mono-molecular building blocks are correspondingly distributed over more than 20 molecular point groups. The molecules comprise H, B, C, N, O, F, Si, P, S, Cl, Se, Te, Br, I, and As species and therewith cover the entire range of elements commonly found in OSCs. As further

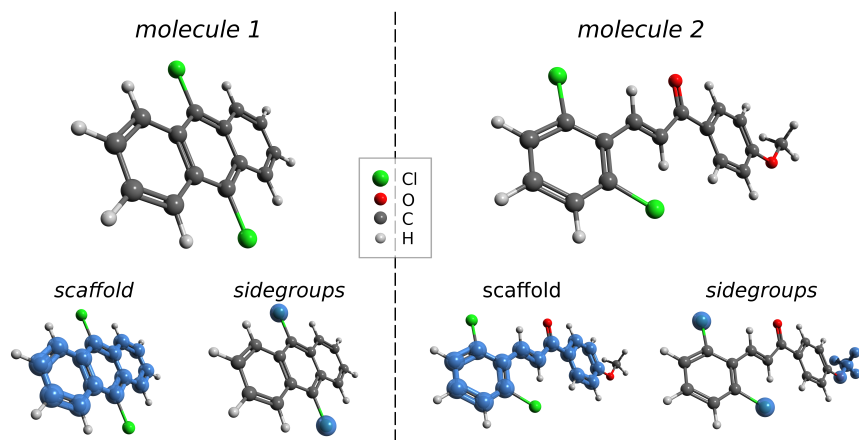


Figure 1: Illustration of the Bemis-Murcko partitioning scheme. For two example molecules, the decomposition into scaffold and sidegroups is illustrated by highlighting the scaffold in light blue on the left and by highlighting the sidegroups on the right.

specified below, a variety of molecular backbones is contained in the database, featuring single to multiple (hetero)cycles. These backbones are functionalized with 4229 different sidegroups, 828 of which occur in more than three different crystals.

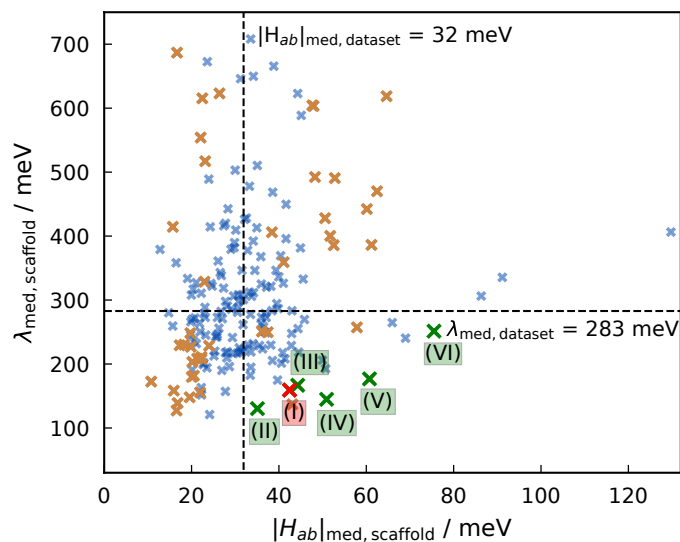
In order to extract correlations from such a diverse set of data, a suitable representation is required that collapses its dimensionality to a tractable size. In contrast to more abstract substructure or statistical learning based approaches,⁸⁹ we here pursue a chemically intuitive concept that partitions the molecular building block of the organic crystal into ring systems, linkers and sidegroup atoms. The Bemis-Murcko (BM) scaffold⁶⁴ is thereby specifically defined as the molecular core comprised of connected ring systems and their linkers, resulting after removal of all sidegroup atoms. Figure 1 provides two illustrative examples of this decomposition scheme. While BM scaffolds have so far mainly been used in medicinal chemistry to study properties of drug-like molecules,^{90–92} the concept is also well-suited for our purpose. It mainly extracts π -conjugated ring systems which to date almost exclusively form the basis of molecular electronics.^{93,94} The scaffold also constitutes the largest part of the molecule which is likely most decisive in terms of electronic structure, shape, conformational flexibility⁸⁷ and possibly geometric arrangement in the crystal. Consequently, we disregard some structures in the “64k-dataset” where very extensive sidegroups constitute more than

50 % of the non-hydrogen atoms.

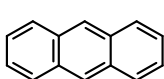
Using this BM representation allows us to cluster the dataset entries based on their respective scaffolds. In order to obtain representative sample sizes, we required each such cluster to contain at least 15 crystals. Although this rigorous filtering discards many clusters due to too small cluster sizes, it leads to a well-analyzable subset of the “64k-dataset” that contains 7.569 molecular crystals distributed over 195 scaffold clusters. These molecular crystals are made of 6.936 unique molecules, since 633 crystals appear as polymorphic structures or had entered the Cambridge Structural Database more than once, e.g. coming from different reported experimental sources. Nevertheless, each scaffold cluster comprises at least 10 unique molecules and we provide a list of BM cluster scaffolds, the detailed distributions of cluster sizes, as well as $|H_{ab}|$ and λ values in this subset in the Supporting Information (SI) Figs. S2-4. For the remainder of the study we will focus on this data subset and will consistently refer to it as our database.

In Fig. 2 we first analyze the median over $|H_{ab}|$ and λ for each scaffold cluster and contrast it with the corresponding medians over the entire database. This provides a chemically intuitive assessment of the performance and tuning potential of each particular scaffold. Within the scope of structures in each scaffold cluster, this averaging step accounts for changes in electronic and crystal structure, brought about by differing anchor points of the varying sidegroups, thus to a certain extent accounting for the so called positional disorder⁹⁵ in each scaffold.

To provide a reference of a well-performing scaffold, we highlight in Fig. 2 the anthracene scaffold cluster. With anthracene a well-known base material for organic electronics,⁹⁶ this reference cluster is predominantly composed of high charge-carrier mobility crystals with relatively low λ and high $|H_{ab}|$. Compared to this reference cluster, there is a wide spread of the median descriptors among the 195 scaffold clusters. In particular, there are many scaffolds with exceedingly poor median charge-transport characteristics (very high $\lambda_{\text{med,scaffold}}$ and/or very low $|H_{ab}|_{\text{med,scaffold}}$). On the other hand, there are also several scaffold types



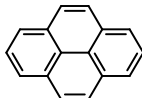
(I) *anthracene*



$$\lambda_{\text{med,scaffold}} = 159$$

$$|H_{ab}|_{\text{med,scaffold}} = 42$$

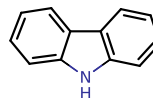
(II) *pyrene*



$$\lambda_{\text{med,scaffold}} = 130$$

$$|H_{ab}|_{\text{med,scaffold}} = 35$$

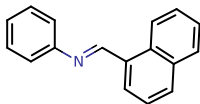
(III) *carbazole*



$$\lambda_{\text{med,scaffold}} = 167$$

$$|H_{ab}|_{\text{med,scaffold}} = 44$$

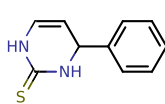
(IV)



$$\lambda_{\text{med,scaffold}} = 145$$

$$|H_{ab}|_{\text{med,scaffold}} = 51$$

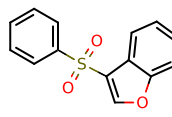
(V)



$$\lambda_{\text{med,scaffold}} = 177$$

$$|H_{ab}|_{\text{med,scaffold}} = 61$$

(VI)



$$\lambda_{\text{med,scaffold}} = 251$$

$$|H_{ab}|_{\text{med,scaffold}} = 76$$

Figure 2: Top: Scatterplot displaying the median $|H_{ab}|_{\text{med,scaffold}}$ and $\lambda_{\text{med,scaffold}}$ for each of the 195 scaffold clusters. Orange crosses mark clusters with a $|H_{ab}|$ and λ distribution significantly differing from each of the respective background distributions, and blue crosses mark clusters with no such significant difference. The medians of the two descriptors over the entire dataset are marked with dashed lines. The anthracene cluster is highlighted in red as a reference cluster containing well-known high charge-carrier mobility crystals. Bottom: Examples of BM-scaffolds further discussed in the text. The datapoints for these scaffolds are labeled and depicted in green color in the top panel.

that exhibit at least as favorable characteristics as the anthracene reference cluster. Note, that in particular improvements in the reorganization energy, i.e. a decrease, can thereby also readily compensate for a deterioration in the electronic coupling (or vice versa). Within a Marcus model, a 50% improvement (deterioration) of λ could for instance tolerate a decrease (increase) of $|H_{ab}|$ by up to 18 meV (26 meV) to still yield the same hopping rates, and by implication mobilities, as anthracene. Essentially all of the scaffolds highlighted in Fig. 2 fall in approximately the same range of mobility. Considering the immense range of $\lambda_{\text{med,scaffold}}$ values included in our dataset and comparing it to the narrow range of available $|H_{ab}|_{\text{med,scaffold}}$ values, we again find that the reorganization energy is an ideal first level target for molecular optimization. The large spread of $\lambda_{\text{med,scaffold}}$ seen in Fig. 2 over all scaffolds indicates thereby the vastly different degrees of inner-sphere reorganization caused by the different molecular cores. A spread along $|H_{ab}|_{\text{med,scaffold}}$, on the other hand, can be caused by a number of factors, such as varying capabilities of scaffolds to stabilize their most favorable packing motives, as well as intrinsic differences of the involved frontier orbital geometries.³⁹

In order to assess the statistical significance of this observed spread in the median descriptors we employ non-parametric statistical testing based on the two-tailed Mann-Whitney-U test.⁸³ As further detailed in the methods section, this test determines non-parametrically the probability p of two sets of values to belong to the same distribution. By correcting detection bias due to the relatively large number of tests being carried out, this yields the q probability and we judge from a value $q < 5\%$ that the distribution of descriptor values in a given scaffold cluster differs significantly from the distribution in the whole database. With respect to the λ descriptor, 105 clusters, i.e. more than half of all considered scaffolds, fulfill this restrictive criterion. With respect to the $|H_{ab}|$ descriptor, this number is a bit lower, but with 69 still extends over more than a third of the 195 scaffolds (the full data for all scaffolds is given in Table S1 in the SI). A total of 44 clusters exhibits statistically different distributions in both descriptors and the corresponding data is color-coded in Fig. 2. The statistical analysis therewith confirms our initial notion that BM scaffolds are a suitable data

representation that allows to differentiate molecular crystals in terms of their charge-transfer properties. Favorable scaffolds with a general tendency to high electronic couplings and low reorganization energies can thus be identified from the available data. Smaller such scaffolds featuring one or two aromatic rings could then be used as promising building blocks, while larger scaffolds can serve as a starting point for further functionalization. In Fig. 2 we specifically highlight a few such identified smaller scaffolds that we further discuss in Section III below.

3.2 Molecular Sidegroups

Having identified suitable molecular scaffolds, we next turn our attention to the sidegroups. The introduction of appropriate sidegroups has been shown to improve carrier mobilities by influencing the stable crystal structure^{20,45,97} or by suppressing vibrationally induced disorder.^{98,99} Here, we focus specifically on the influence on the reorganization energy. Previous work has confirmed this critical influence in comparative studies attaching different sidegroups to a given single scaffold.^{41,48-54} Building on this empirical concept, we now use our extensive database to investigate whether there are molecular sidegroups that generally lower the reorganization energy, independent of the scaffold they are attached to. We achieve this with the same techniques of statistical analysis as before for the scaffolds, i.e. we first identify all sidegroups in the database and then cluster the data. This time, one cluster corresponds to all molecules (of any scaffold type) that contain a given sidegroup. This does not differentiate between specific attachment positions to the scaffold or the coexistence of multiple sidegroups in one molecule, and we return to this point below. To ensure meaningful statistics, we only consider sidegroup clusters comprising at least 10 different molecules, with the notable exception of the $-\text{Se}-\text{CH}_3$ sidegroup. This sidegroup occurs only in six molecular crystals in the database, but is nevertheless considered because of its exceptional properties (*vide infra*). Altogether this yields 70 sidegroup clusters containing a total of 6.162 unique molecules.

For each molecule in a given sidegroup cluster we determine its relative reorganization energy $\Delta\lambda$ as the difference between the reorganization energy λ of the molecule and the median reorganization energy $\lambda_{\text{med,scaffold}}$ over the cluster of molecules with this particular scaffold, i.e. the $\lambda_{\text{med,scaffold}}$ evaluated in the last section. Thus, a negative value of $\Delta\lambda$ indicates that the attachment of the sidegroup to the scaffold leads to compounds of lower reorganization energy as compared to the attachment of other sidegroups to the same scaffold. Note that due to the diversity of the database it is possible that certain sidegroups appear multiple times on a single scaffold. In order not to over-emphasize such cases, we divide each $\Delta\lambda$ by the number of occurrences. Taking the median of all $\Delta\lambda$ within one sidegroup cluster, we finally arrive at $\Delta\lambda_{\text{med,sidegroup}}$, which evaluates how the attachment of the sidegroup generally influences the reorganization energy over different scaffold types.

For many sidegroups, the obtained $\Delta\lambda_{\text{med,sidegroup}}$ differ significantly from zero. This indicates that these sidegroups indeed either consistently worsen or consistently improve the reorganization energy, regardless of the scaffold they are attached to. Mann-Whitney-U testing⁸³ is again employed to verify that these results are not due to chance or a biased sample selection. For 34 of the 70 sidegroup clusters this statistical testing confirms that their relative reorganization energy distributions differ significantly from the distribution in the entire database ($q < 5\%$). In particular for the $-\text{Se}-\text{CH}_3$ functional group a very low q -value of 1% confirms the statistical relevance despite the smaller sample size.

Figure 3 compiles the $\Delta\lambda$ distributions of these 34 sidegroup clusters. Groups appearing at the left side of Fig. 3 tend to reliably decrease λ , while those on the right tend to increase the reorganization energy. As apparent, a $\Delta\lambda$ distribution within a sidegroup cluster differing significantly from the distribution over the entire database does not necessarily mean that the median $\Delta\lambda_{\text{med,sidegroup}}$ has to differ markedly from zero. For the corresponding sidegroups in the middle of Fig. 3, the effect of the sidegroup seems more subtle, sometimes increasing, sometimes decreasing the reorganization energy. We speculate that in these cases the (presently unresolved) specific attachment position or the coexistence of other sidegroups

could play a decisive role.⁵¹ More importantly, there are, however, 12 sidegroups which on average decrease the reorganization energy by more than 25 meV, see Fig. 3. These sidegroups are obviously good candidates for λ -tuning approaches and will be further discussed in Section III below.

3.3 Sidegroup Engineering

Our data analysis to this point has identified suitable scaffolds that yield generally high $|H_{ab}|$ and low λ values, as well as suitable sidegroups that consistently lower λ . This suggests an appealing design principle to arrive at new molecular crystals with optimized charge-transfer properties: Select a favorable scaffold with high $|H_{ab}|$ and further optimize the reorganization energy by addition of favorable sidegroups or by exchange of unfavorable with favorable sidegroups. Unfortunately, this neglects the known sensitive influence of sidegroups on the crystal structure.^{43,100} In contrast to λ , which is predominantly determined by the intramolecular structure, $|H_{ab}|$ describes the electronic coupling between molecules in the crystal and is thus critically dependent on the crystal structure. While sidegroup tuning could improve the reorganization energy, it could also change the crystal structure and then potentially worsen $|H_{ab}|$.

To address this aspect we use matched molecular pair (MMP) analysis^{101,102} to identify those molecules in the database, where the resulting molecule after addition or exchange of a single sidegroup is still contained in the database. This way, the crystal structure of both the initial and the modified molecule is available, enabling us to explicitly analyze the influence of the sidegroup engineering on the electronic coupling $|H_{ab}|$. Specifically, we thereby focus on the exchange of an unfavorable sidegroup ($\Delta\lambda_{\text{med,sidegroup}} > 0$ meV) with a sidegroup identified as very favorable above ($\Delta\lambda_{\text{med,sidegroup}} < -25$ meV), cf. Fig. 3. This generates a testset S1 containing 206 optimized molecules. As a control series, we also generated a testset S2 containing 124 molecules, where a favorable sidegroup ($\Delta\lambda_{\text{med,sidegroup}} < 0$ meV) is replaced with a highly unfavorable one ($\Delta\lambda_{\text{med,sidegroup}} > 25$ meV). In analogy, testset

S3 with 161 molecules (S4 with 91 molecules) contains molecules where a highly favorable (unfavorable) sidegroup has been added to the molecular scaffold. The size discrepancy between the four testsets is due to the different availability of experimental crystal structures. A full description of the testsets and the ensuing analysis is provided in the SI.

Fully supporting our sidegroup analysis above, we obtain marked improvements of the reorganization energy (up to changes of a factor of 5) for the overwhelming majority of modified crystals in the two “optimized” testsets (87% of all crystals in S1, 77% in S3), see Fig. 4. Similarly, the “worsened” control sets S2 and S4 show predominantly a large deterioration of the λ (in 80% of cases in S2, 73% in S4). As expected, the sidegroup engineering also influences the electronic coupling. However, the changes there are mostly unsystematic. In the “optimized” testsets S1 and S3 we find a roughly equal number of cases where $|H_{ab}|$ is improved to those where the coupling deteriorates. Interestingly, in the “worsened” testsets S2 and S4, we still obtain a slight improvement of $|H_{ab}|$ in 58 % and 59 % of all cases, respectively. This is most likely attributed to selection bias, as the same correlation of large couplings with regrettably large reorganization energies is already present in the molecules of “64k-dataset”.²⁵ Notwithstanding, in particular for the relevant larger coupling values, the induced unsystematic changes are comparatively small. As the reorganization energy furthermore enters e.g. the Marcus model exponentially, this suggests that these changes in $|H_{ab}|$ should not break the overall trend. Indeed, when we compute the Marcus transfer rates k_{Marcus} on the basis of the obtained couplings and reorganization energies,¹⁸ we still find increased rates in 71% (61%) of the molecular crystals in the “optimized” testset S1 (S3), and decreased rates in 67% (59%) of the molecular crystals in the “worsened” control testsets S2 (S4). This suggests that sidegroup engineering is (at least to a certain extent) robust with respect to induced changes in the crystal structure and is therefore a viable approach to optimize the charge-transport properties of molecular crystals.

4 Discussion

In this work we explored the possibility of optimizing organic crystals for charge transport using a data-driven approach. Mining a large database of two major charge-transport descriptors, the reorganization energy λ and the electronic coupling $|H_{ab}|$, for close to 7.000 mono-molecular crystals, we transferred the concept of Bemis-Murcko scaffolds from pharmaceutical research to cluster the data according to molecular scaffolds and the attached sidegroups. For both types of clusters we obtain statistically confirmed structure-property relationships, i.e. we identify scaffolds or sidegroups that are generally favorable for charge transport.

A number of promising scaffolds with generally low λ and high $|H_{ab}|$ regardless of their functionalization are identified in the lower right quadrant of Fig. 2. Among them are scaffolds contained in well-known OSCs: anthracene⁹⁶ (Fig. 2 I), pyrene¹⁰³ (Fig. 2 II), or carbazole²¹ (Fig. 2 III). Based on our strict statistical significance criterion, all three exhibit consistently improved charge-transport parameters compared to the database background. The only exception is $|H_{ab}|_{\text{med,scaffold}}$ of the pyrene scaffold II, which shows only little statistical significance due to the limited occurrence of this scaffold in our database.

Next to these more common aromatic systems, there are also a few promising linked ring systems like scaffolds IV and V shown in Fig. 2, which one would not immediately associate with the OSC context. They do, however, exhibit much higher median electronic couplings and comparable reorganization energies as the anthracene reference cluster. Multiple other scaffolds are particularly appealing with regard to only one of the two charge-transport descriptors. In particular, scaffolds like Fig. 2 VI exhibit extraordinarily high electronic couplings suitable for organic electronics, but feature on average unsuitably high reorganization energies. These scaffolds are naturally of specific interest for a targeted optimization through functionalization with favorable sidegroups.

Our data-driven analysis of sidegroups paired with tests of statistical significance allows for a broader perspective on existing empirical strategies of sidegroup tuning. Among the

34 identified sidegroups that significantly alter the reorganization energies of their respective scaffolds are sidegroups like $-\text{OH}$, $-\text{O}-\text{CH}_3$, $-\text{NH}_2$, $-\text{F}$, $-\text{Cl}$, $-\text{CF}_3$, or $-\text{C}\equiv\text{N}$ that had been considered in preceding comparative studies.^{41,48-55,104} In parts, these earlier studies found trends that match the general structure-property relationships established in this work, for example for $-\text{C}\equiv\text{N}$, $-\text{F}$, $-\text{CF}_3$, and $-\text{NH}_2$ substitutions on thiophenes, furans, and pyrroles.⁴¹ Interestingly, with the exception of the cyano group, these often studied sidegroups not always seem to have a favorable effect on reorganization energies, with some ($-\text{CF}_3$, and $-\text{NH}_2$) even showing—on average—the opposite trend. The here confirmed favorable influence of the cyano group has been explained in the literature⁵⁰ by the introduction of a local nonbonding character in the frontier orbitals. In contrast, the common and much studied class of alkyl sidegroups is, with the exception of $-\text{CH}(\text{CH}_3)_2$, not represented among the 34 sidegroups with significant influence. On the basis of our present database we can thus not confirm any systematic improvement or deterioration of the reorganization induced by this class of sidegroups. Instead and most intriguingly, many of the here identified 12 sidegroups that lead to a marked improvement of λ have to the best of our knowledge only rarely been discussed in the OSC context. This holds for $-\text{Se}-\text{CH}_3$, $-\text{N}(\text{CH}_2-\text{CH}_3)_2$, $-\text{SH}$,⁵⁴ $-\text{I}$, $-\text{CH}_2-\text{NO}_2$, $-\text{S}-\text{CH}_3$,¹⁰⁴ $-\text{C}\equiv\text{CH}$,⁵⁴ $-\text{CH}=\text{O}$,⁵⁴ or $-\text{N}(\text{CH}_3)_2$,^{55,104} where we added the few references to works addressing these sidegroups we could find at all. A detailed analysis of the reason for the consistent λ reduction induced by these sidegroups is presented in the SI. It confirms that each of these sidegroups yields first of all an expected additional contribution to the reorganization energy.⁵¹ However, this contribution is more than compensated by a reduction of the scaffold reorganization energy. Especially the best performing $-\text{Se}-\text{CH}_3$ sidegroup achieves this by efficiently delocalizing the charge throughout the scaffold and the sidegroup.

The clear structure-property relationships for both scaffolds and sidegroups suggest the possibility of a targeted molecular design by combining favorable scaffolds with reorganization energy-lowering sidegroups. This effect is expected to lead to improved charge transport

rates, with the reorganization energy entering e.g. a Marcus-like³⁸ model exponentially, while $|H_{ab}|$ only enters quadratically.¹⁸ Our data analysis indicates that such a design strategy is robust (at least to a certain extent) to a change in the crystal structure that might be induced by such functionalization. Accordingly, we would presently predict the addition for instance of most favorable $-\text{Se}-\text{CH}_3$ or $-\text{N}(\text{CH}_2-\text{CH}_3)_2$ sidegroups to any of the favorable scaffolds highlighted in Fig. 2 to yield highly promising OSC candidates for experimental synthesis. Much of the uncertainty in these predictions could be removed by explicitly validating that the electronic couplings remain indeed favorable in the crystal structure adopted by the corresponding compound. Unfortunately, computational crystal structure prediction is still a highly challenging task.^{105,106} In fact, the limitations in rapidly and reliably determining the crystal structure for a larger number of promising candidates represent in our view currently the major bottleneck to a large-scale data-driven *in silico* discovery of improved OSCs. This prevents a targeted extension of the database through additional first-principles calculations, and correspondingly restricts the data mining to compounds for which experimental crystal structures are available in the Cambridge Structural Database.⁶⁵ The concomitant relative scarcity of data also prohibits refined data clustering that notably would distinguish different anchor points of sidegroups to a scaffold or the coexistence of multiple sidegroups. Efficient crystal structure prediction and the concerted buildup of experimental and computational structural databases are therefore key to unleash the full potential of modern data mining or machine learning approaches for organic semiconductors, and thereby navigate the vast and largely unexplored design spaces of this intriguing class of molecular materials.

Supporting Information Available

Additional detailed information on the BM scaffold dataset: Distributions of $|H_{ab}|$, λ over the dataset, as well as a distribution of the number of molecules in the scaffolds and a list of all scaffolds. In addition, an extended boxplot for all sidegroup clusters as well as a depiction

of all identified sidegroups is provided. A series of $-H$, $-S-CH_3$ and $-Se-CH_3$ exchanges on different molecules further illustrates the favorable impact of the latter two groups on gas-phase reorganization energies, while a detailed discussion of the underlying factors is provided for one example.

Acknowledgement

We want to thank Lynne Stecher for very helpful discussions on the statistical methods. We further acknowledge support from the Solar Technologies Go Hybrid initiative of the State of Bavaria and the Leibniz Supercomputing Centre for high-performance computing time at the SuperMUC facility.

References

- (1) Minemawari, H.; Yamada, T.; Matsui, H.; Tsutsumi, J.; Haas, S.; Chiba, R.; Kumai, R.; Hasegawa, T. Inkjet printing of single-crystal films. *Nature* **2011**, *475*, 364.
- (2) Stavriniidou, E.; Gabrielsson, R.; Gomez, E.; Crispin, X.; Nilsson, O.; Simon, D. T.; Berggren, M. Electronic plants. *Sci. Adv.* **2015**, *1*, e1501136.
- (3) Xu, J.; Wang, S.; Wang, G.-J. N.; Zhu, C.; Luo, S.; Jin, L.; Gu, X.; Chen, S.; Feig, V. R.; To, J. W. et al. Highly stretchable polymer semiconductor films through the nanoconfinement effect. *Science* **2017**, *355*, 59–64.
- (4) Nikolka, M.; Nasrallah, I.; Rose, B.; Ravva, M. K.; Broch, K.; Sadhanala, A.; Harkin, D.; Charmet, J.; Hurhangee, M.; Brown, A. et al. High operational and environmental stability of high-mobility conjugated polymer field-effect transistors through the use of molecular additives. *Nat. Mater.* **2017**, *16*, 356.

- (5) Wang, C.; Dong, H.; Jiang, L.; Hu, W. Organic semiconductor crystals. *Chem. Soc. Rev.* **2018**, *47*, 422–500.
- (6) Wang, C.; Dong, H.; Hu, W.; Liu, Y.; Zhu, D. Semiconducting π -Conjugated Systems in Field-Effect Transistors: A Material Odyssey of Organic Electronics. *Chem. Rev.* **2012**, *112*, 2208–2267.
- (7) Lin, Y.; Li, Y.; Zhan, X. Small molecule semiconductors for high-efficiency organic photovoltaics. *Chem. Soc. Rev.* **2012**, *41*, 4245–4272.
- (8) Xu, R.-P.; Li, Y.-Q.; Tang, J.-X. Recent advances in flexible organic light-emitting diodes. *J. Mater. Chem. C* **2016**, *4*, 9116–9142.
- (9) Anikeeva, P.; Koppes, R. A. Restoring the sense of touch. *Science* **2015**, *350*, 274–275.
- (10) Anthony, J. E. Organic electronics: addressing challenges. *Nat. Mater.* **2014**, *13*, 773.
- (11) Yavuz, I. Dichotomy between the band and hopping transport in organic crystals: insights from experiments. *Phys. Chem. Chem. Phys.* **2017**, *19*, 25819–25828.
- (12) Venkateshvaran, D.; Nikolka, M.; Sadhanala, A.; Lemaur, V.; Zelazny, M.; Kepa, M.; Hurhangee, M.; Kronemeijer, A. J.; Pecunia, V.; Nasrallah, I. et al. Approaching disorder-free transport in high-mobility conjugated polymers. *Nature* **2014**, *515*, 384.
- (13) Ostroverkhova, O. Organic optoelectronic materials: mechanisms and applications. *Chem. Rev.* **2016**, *116*, 13279–13412.
- (14) Kang, S. D.; Snyder, G. J. Charge-transport model for conducting polymers. *Nat. Mater.* **2017**, *16*, 252.
- (15) Brédas, J.-L.; Sargent, E. H.; Scholes, G. D. Photovoltaic concepts inspired by coherence effects in photosynthetic systems. *Nat. Mater.* **2017**, *16*, 35.

- (16) Kim, G.; Kang, S.-J.; Dutta, G. K.; Han, Y.-K.; Shin, T. J.; Noh, Y.-Y.; Yang, C. A thienoisindigo-naphthalene polymer with ultrahigh mobility of 14.4 cm²/Vs that substantially exceeds benchmark values for amorphous silicon semiconductors. *J. Am. Chem. Soc.* **2014**, *136*, 9477–9483.
- (17) Mei, J.; Diao, Y.; Appleton, A. L.; Fang, L.; Bao, Z. Integrated Materials Design of Organic Semiconductors for Field-Effect Transistors. *J. Am. Chem. Soc.* **2013**, *135*, 6724–6746.
- (18) Oberhofer, H.; Blumberger, J. Revisiting electronic couplings and incoherent hopping models for electron transport in crystalline C 60 at ambient temperatures. *Phys. Chem. Chem. Phys.* **2012**, *14*, 13846–13852.
- (19) Zhang, C.; Zhu, X. Thieno[3,4-b]thiophene-Based Novel Small-Molecule Optoelectronic Materials. *Acc. Chem. Res.* **2017**, *50*, 1342–1350.
- (20) Sosorev, A. Y. Role of intermolecular charge delocalization and its dimensionality in efficient band-like electron transport in crystalline 2,5-difluoro-7,7,8,8-tetracyanoquinodimethane (F2-TCNQ). *Phys. Chem. Chem. Phys.* **2017**, *19*, 25478–25486.
- (21) Reig, M.; Bagdziunas, G.; Volyniuk, D.; Grazulevicius, J. V.; Velasco, D. Tuning the ambipolar charge transport properties of tricyanovinyl-substituted carbazole-based materials. *Phys. Chem. Chem. Phys.* **2017**, *19*, 6721–6730.
- (22) Li, S.; Ye, L.; Zhao, W.; Yan, H.; Yang, B.; Liu, D.; Li, W.; Ade, H.; Hou, J. A Wide Band Gap Polymer with a Deep Highest Occupied Molecular Orbital Level Enables 14.2% Efficiency in Polymer Solar Cells. *J. Am. Chem. Soc.* **2018**, *140*, 7159–7167.
- (23) Olivares-Amaya, R.; Amador-Bedolla, C.; Hachmann, J.; Atahan-Evrenk, S.; Sanchez-Carrera, R. S.; Vogt, L.; Aspuru-Guzik, A. Accelerated computational discovery of

- high-performance materials for organic photovoltaics by means of cheminformatics. *Energy Environ. Sci.* **2011**, *4*, 4849–4861.
- (24) Akimov, A. V.; Prezhdov, O. V. Large-scale computations in chemistry: a bird’s eye view of a vibrant field. *Chem. Rev.* **2015**, *115*, 5797–5890.
- (25) Schober, C.; Reuter, K.; Oberhofer, H. Virtual Screening for High Carrier Mobility in Organic Semiconductors. *J. Phys. Chem. Lett.* **2016**, *7*, 3973–3977.
- (26) Mercado, R.; Fu, R.-S.; Yakutovich, A. V.; Talirz, L.; Haranczyk, M.; Smit, B. In Silico Design of 2D and 3D Covalent Organic Frameworks for Methane Storage Applications. *Chemistry of Materials* **2018**, *30*, 5069–5086.
- (27) Meyer, B.; Sawatlon, B.; Heinen, S.; von Lilienfeld, O. A.; Corminboeuf, C. Machine learning meets volcano plots: computational discovery of cross-coupling catalysts. *Chem. Sci.* **2018**, *9*, 7069–7077.
- (28) Meredig, B.; Agrawal, A.; Kirklın, S.; Saal, J. E.; Doak, J. W.; Thompson, A.; Zhang, K.; Choudhary, A.; Wolverton, C. Combinatorial screening for new materials in unconstrained composition space with machine learning. *Phys. Rev. B* **2014**, *89*, 094104.
- (29) Nørskov, J. K.; Bligaard, T.; Rossmeisl, J.; Christensen, C. H. Towards the computational design of solid catalysts. *Nat. Chem.* **2009**, *1*, 37.
- (30) Andersen, M.; Medford, A. J.; Nørskov, J. K.; Reuter, K. Scaling-Relation-Based Analysis of Bifunctional Catalysis: The Case for Homogeneous Bimetallic Alloys. *ACS Catalysis* **2017**, *7*, 3960–3967.
- (31) Goldsmith, B. R.; Esterhuizen, J.; Liu, J.-X.; Bartel, C. J.; Sutton, C. Machine learning for heterogeneous catalyst design and discovery. *AIChE Journal* **2018**, *64*, 2311–2323.

- (32) Mansouri Tehrani, A.; Oliynyk, A. O.; Parry, M.; Rizvi, Z.; Couper, S.; Lin, F.; Miyagi, L.; Sparks, T. D.; Brgoch, J. Machine Learning Directed Search for Ultraincompressible, Superhard Materials. *J. Am. Chem. Soc.* **2018**, *140*, 9844–9853, PMID: 30010335.
- (33) Chen, H.; Engkvist, O.; Wang, Y.; Olivecrona, M.; Blaschke, T. The rise of deep learning in drug discovery. *Drug Discov. Today* **2018**, *23*, 1241–1250.
- (34) Zunger, A. Inverse design in search of materials with target functionalities. *Nat. Rev. Chem.* **2018**, *2*, 0121, Perspective.
- (35) Misra, M.; Andrienko, D.; Baumeier, B.; Faulon, J.-L.; von Lilienfeld, O. A. Toward Quantitative Structure–Property Relationships for Charge Transfer Rates of Polycyclic Aromatic Hydrocarbons. *J. Chem. Theory Comput.* **2011**, *7*, 2549–2555.
- (36) Sahu, H.; Rao, W.; Troisi, A.; Ma, H. Toward Predicting Efficiency of Organic Solar Cells via Machine Learning and Improved Descriptors. *Adv. Energy Mater.* **2018**, *8*, 1801032.
- (37) Marcus, R. A. On the Theory of Oxidation-Reduction Reactions Involving Electron Transfer. I. *J. Chem. Phys.* **1956**, *24*, 966–978.
- (38) Marcus, R. A. Electron Transfer Reactions in Chemistry. Theory and Experiment. *Rev. Mod. Phys.* **1993**, *65*, 599–610.
- (39) Oberhofer, H.; Reuter, K.; Blumberger, J. Charge Transport in Molecular Materials: An Assessment of Computational Methods. *Chem. Rev.* **2017**, *117*, 10319–10357.
- (40) McMahon, D. P.; Troisi, A. Evaluation of the External Reorganization Energy of Polyacenes. *J. Phys. Chem. Lett.* **2010**, *1*, 941–946.

- (41) Hutchison, G. R.; Ratner, M. A.; Marks, T. J. Hopping Transport in Conductive Heterocyclic Oligomers: Reorganization Energies and Substituent Effects. *J. Am. Chem. Soc.* **2005**, *127*, 2339–2350.
- (42) Sokolov, A. N.; Atahan-Evrenk, S.; Mondal, R.; Akkerman, H. B.; Sánchez-Carrera, R. S.; Granados-Focil, S.; Schrier, J.; Mannsfeld, S. C. B.; Zoombelt, A. P.; Bao, Z. et al. From computational discovery to experimental characterization of a high hole mobility organic crystal. *Nat. Commun.* **2011**, *2*, 437.
- (43) Alberga, D.; Ciofini, I.; Mangiatordi, G. F.; Pedone, A.; Lattanzi, G.; Roncali, J.; Adamo, C. Effects of Substituents on Transport Properties of Molecular Materials for Organic Solar Cells: A Theoretical Investigation. *Chem. Mater.* **2017**, *29*, 673–681.
- (44) Li, J.; Zhao, Y.; Tan, H. S.; Guo, Y.; Di, C.-A.; Yu, G.; Liu, Y.; Lin, M.; Lim, S. H.; Zhou, Y. et al. A stable solution-processed polymer semiconductor with record high-mobility for printed transistors. *Sci. Rep.* **2012**, *2*, 754.
- (45) Liu, Z.; Zhang, G.; Zhang, D. Modification of Side Chains of Conjugated Molecules and Polymers for Charge Mobility Enhancement and Sensing Functionality. *Acc. Chem. Res.* **2018**, *51*, 1422–1432.
- (46) Moral, M.; Garzón-Ruiz, A.; Castro, M.; Canales-Vázquez, J.; Sancho-García, J. C. Virtual Design in Organic Electronics: Screening of a Large Set of 1,4-Bis(phenylethynyl)benzene Derivatives as Molecular Semiconductors. *J. Phys. Chem. C* **2017**, *121*, 28249–28261.
- (47) Zhu, R.; Duan, Y.-A.; Geng, Y.; Wei, C.-Y.; Chen, X.-Y.; Liao, Y. Theoretical evaluation on the reorganization energy of five-ring-fused benzothiophene derivatives. *Comput. Theor. Chem.* **2016**, *1078*, 16–22.
- (48) Oshi, R.; Abdalla, S.; Springborg, M. Theoretical study on functionalized anthracene

- and tetraceneas starting species to produce promising semiconductor materials. *Comput. Theor. Chem.* **2018**, *1128*, 60–69.
- (49) Oshi, R.; Abdalla, S.; Springborg, M. Study of the influence of functionalization on the reorganization energy of naphthalene using DFT. *Comput. Theor. Chem.* **2017**, *1099*, 209–215.
- (50) Chang, Y.-C.; Chao, I. An Important Key to Design Molecules with Small Internal Reorganization Energy: Strong Nonbonding Character in Frontier Orbitals. *J. Phys. Chem. Lett.* **2010**, *1*, 116–121.
- (51) Geng, H.; Niu, Y.; Peng, Q.; Shuai, Z.; Coropceanu, V.; Brédas, J.-L. Theoretical study of substitution effects on molecular reorganization energy in organic semiconductors. *J. Chem. Phys.* **2011**, *135*, 104703.
- (52) Lee, C.; Sohlberg, K. The effect of substitution on reorganization energy and charge mobility in metal free phthalocyanine. *Chem. Phys.* **2010**, *367*, 7–19.
- (53) Vedova-Brook, N.; Matsunaga, N.; Sohlberg, K. Correlating substituent parameter values to electron transport properties of molecules. *Chem. Phys.* **2004**, *299*, 89–95.
- (54) Zhao, C.; Wang, W.; Ma, Y. Molecular design toward good hole transport materials based on anthra[2,3-c]thiophene: A theoretical investigation. *Comput. Theor. Chem.* **2013**, *1010*, 25–31.
- (55) Pan, J.-H.; Chiu, H.-L.; Chen, L.; Wang, B.-C. Theoretical investigations of triphenylamine derivatives as hole transporting materials in OLEDs: Correlation of the Hammett parameter of the substituent to ionization potential, and reorganization energy level. *Comput. Mater. Sci.* **2006**, *38*, 105–112.
- (56) Chen, W.-C.; Chao, I. Molecular Orbital-Based Design of π -Conjugated Organic Mate-

- rials with Small Internal Reorganization Energy: Generation of Nonbonding Character in Frontier Orbitals. *J. Phys. Chem. C* **2014**, *118*, 20176–20183.
- (57) Kuo, M.-Y.; Liu, C.-C. Molecular Design toward High Hole Mobility Organic Semiconductors: Tetraceno[2,3-c]thiophene Derivatives of Ultrasmall Reorganization Energies. *J. Phys. Chem. C* **2009**, *113*, 16303–16306.
- (58) Cao, B.; Adutwum, L. A.; Oliynyk, A. O.; Lubber, E. J.; Olsen, B. C.; Mar, A.; Buriak, J. M. How To Optimize Materials and Devices via Design of Experiments and Machine Learning: Demonstration Using Organic Photovoltaics. *ACS Nano* **2018**, *12*, 7434–7444.
- (59) Gómez-Bombarelli, R.; Aguilera-Iparraguirre, J.; Hirzel, T. D.; Duvenaud, D.; Maclaurin, D.; Blood-Forsythe, M. A.; Chae, H. S.; Einzinger, M.; Ha, D.-G.; Wu, T. et al. Design of efficient molecular organic light-emitting diodes by a high-throughput virtual screening and experimental approach. *Nat. Mater.* **2016**, *15*, 1120.
- (60) Agrawal, A.; Choudhary, A. Perspective: Materials informatics and big data: Realization of the “fourth paradigm” of science in materials science. *APL Mater.* **2016**, *4*, 053208.
- (61) Butler, K. T.; Davies, D. W.; Cartwright, H.; Isayev, O.; Walsh, A. Machine learning for molecular and materials science. *Nature* **2018**, *559*, 547–555.
- (62) Ferguson, A.; Hachmann, J. Machine learning and data science in materials design: a themed collection. *Mol. Syst. Des. Eng.* **2018**, *3*, 429–430.
- (63) Tabor, D. P.; Roch, L. M.; Saikin, S. K.; Kreisbeck, C.; Sheberla, D.; Montoya, J. H.; Dwaraknath, S.; Aykol, M.; Ortiz, C.; Tribukait, H. et al. Accelerating the discovery of materials for clean energy in the era of smart automation. *Nat. Rev. Mater.* **2018**, *3*, 5–20.

- (64) Bemis, G. W.; Murcko, M. A. The Properties of Known Drugs. 1. Molecular Frameworks. *J. Med. Chem.* **1996**, *39*, 2887–2893.
- (65) Allen, F. H. The Cambridge Structural Database: a quarter of a million crystal structures and rising. *Acta Crystallogr. B* **2002**, *58*, 380–388.
- (66) Weininger, D.; Weininger, A.; Weininger, J. L. SMILES. 2. Algorithm for generation of unique SMILES notation. *J. Chem. Inf. Comput. Sci.* **1989**, *29*, 97–101.
- (67) O’Boyle, N. M.; Banck, M.; James, C. A.; Morley, C.; Vandermeersch, T.; Hutchison, G. R. Open Babel: An open chemical toolbox. *J. Cheminformatics* **2011**, *3*, 33.
- (68) Landrum, G. RDKit: Open-source cheminformatics. <http://www.rdkit.org>, 2018; [Online; accessed 07-August-2018].
- (69) ChemAxon, Marvin 17.5.0. <http://www.chemaxon.com>, 2017; <http://www.chemaxon.com>, [Online; accessed 07-August-2018].
- (70) Ong, S. P.; Richards, W. D.; Jain, A.; Hautier, G.; Kocher, M.; Cholia, S.; Gunter, D.; Chevrier, V. L.; Persson, K. A.; Ceder, G. Python Materials Genomics (pymatgen): A robust, open-source python library for materials analysis. *Comput. Mater. Sci.* **2013**, *68*, 314–319.
- (71) Becke, A. D. Density-functional exchange-energy approximation with correct asymptotic behavior. *Phys. Rev. A* **1988**, *38*, 3098–3100.
- (72) Lee, C.; Yang, W.; Parr, R. G. Development of the Colle-Salvetti correlation-energy formula into a functional of the electron density. *Phys. Rev. B* **1988**, *37*, 785–789.
- (73) Blum, V.; Gehrke, R.; Hanke, F.; Havu, P.; Havu, V.; Ren, X.; Reuter, K.; Scheffler, M. Ab initio molecular simulations with numeric atom-centered orbitals. *Comp. Phys. Commun.* **2009**, *180*, 2175–2196.

- (74) Zhang, I. Y.; Ren, X.; Rinke, P.; Blum, V.; Scheffler, M. Numeric atom-centered-orbital basis sets with valence-correlation consistency from H to Ar. *New J. Phys.* **2013**, *15*, 123033.
- (75) Schober, C.; Reuter, K.; Oberhofer, H. Critical analysis of fragment-orbital DFT schemes for the calculation of electronic coupling values. *J. Chem. Phys.* **2016**, *144*, 054103.
- (76) Senthilkumar, K.; Grozema, F. C.; Bickelhaupt, F. M.; Siebbeles, L. D. A. Charge transport in columnar stacked triphenylenes: Effects of conformational fluctuations on charge transfer integrals and site energies. *J. Chem. Phys.* **2003**, *119*, 9809–9817.
- (77) Nelsen, S. F.; Blackstock, S. C.; Kim, Y. Estimation of inner shell Marcus terms for amino nitrogen compounds by molecular orbital calculations. *J. Am. Chem. Soc.* **1987**, *109*, 677–682.
- (78) Svensson, M.; Humbel, S.; Froese, R. D. J.; Matsubara, T.; Sieber, S.; Morokuma, K. ONIOM: A Multilayered Integrated MO + MM Method for Geometry Optimizations and Single Point Energy Predictions. A Test for Diels-Alder Reactions and Pt(P(*t*-Bu)₃)₂ + H₂ Oxidative Addition. *J. Phys. Chem.* **1996**, *100*, 19357–19363.
- (79) Larsen, A. H.; Mortensen, J. J.; Blomqvist, J.; Castelli, I. E.; Christensen, R.; Dułak, M.; Friis, J.; Groves, M. N.; Hammer, B.; Hargus, C. et al. The atomic simulation environment—a Python library for working with atoms. *Journal of Physics: Condensed Matter* **2017**, *29*, 273002.
- (80) Plimpton, S. Fast Parallel Algorithms for Short-range Molecular Dynamics. *J. Comput. Phys.* **1995**, *117*, 1–19.
- (81) Rappe, A. K.; Casewit, C. J.; Colwell, K. S.; Goddard, W. A.; Skiff, W. M. UFF, a full periodic table force field for molecular mechanics and molecular dynamics simulations. *J. Am. Chem. Soc.* **1992**, *114*, 10024–10035.

- (82) Rappe, A. K.; Goddard, W. A. Charge equilibration for molecular dynamics simulations. *J. Phys. Chem.* **1991**, *95*, 3358–3363.
- (83) Mann, H. B.; Whitney, D. R. On a Test of Whether one of Two Random Variables is Stochastically Larger than the Other. *Ann. Math. Statist.* **1947**, *18*, 50–60.
- (84) Benjamini, Y.; Hochberg, Y. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *J. Royal Stat. Soc. B* **1995**, *57*, 289–300.
- (85) Seabold, S.; Perktold, J. Statsmodels: econometric and statistical modeling with Python. Proceedings of the 9th Python in Science Conference. 2010; pp 57–61.
- (86) Subramanian, A.; Tamayo, P.; Mootha, V. K.; Mukherjee, S.; Ebert, B. L.; Gillette, M. A.; Paulovich, A.; Pomeroy, S. L.; Golub, T. R.; Lander, E. S. et al. Gene set enrichment analysis: A knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. USA* **2005**, *102*, 15545–15550.
- (87) Varin, T.; Gubler, H.; Parker, C. N.; Zhang, J.-H.; Raman, P.; Ertl, P.; Schuffenhauer, A. Compound Set Enrichment: A Novel Approach to Analysis of Primary HTS Data. *J. Chem. Inf. Model.* **2010**, *50*, 2067–2078.
- (88) Napolitano, F.; Sirci, F.; Carrella, D.; di Bernardo, D. Drug-set enrichment analysis: a novel tool to investigate drug mode of action. *Bioinformatics* **2016**, *32*, 235–241.
- (89) Shelat, A. A.; Guy, R. K. Scaffold composition and biological relevance of screening libraries. *Nat. Chem. Biol.* **2007**, *3*, 442.
- (90) Hu, Y.; Stumpfe, D.; Bajorath, J. Computational Exploration of Molecular Scaffolds in Medicinal Chemistry. *J. Med. Chem.* **2016**, *59*, 4062–4076.
- (91) Schuffenhauer, A.; Varin, T. Rule-Based Classification of Chemical Structures by Scaffold. *Mol. Inform.* **2011**, *30*, 646–664.

- (92) Shang, J.; Sun, H.; Liu, H.; Chen, F.; Tian, S.; Pan, P.; Li, D.; Kong, D.; Hou, T. Comparative analyses of structural features and scaffold diversity for purchasable compound libraries. *J. Cheminformatics* **2017**, *9*, 25.
- (93) Wang, C.; Dong, H.; Hu, W.; Liu, Y.; Zhu, D. Semiconducting π -Conjugated Systems in Field-Effect Transistors: A Material Odyssey of Organic Electronics. *Chem. Rev.* **2012**, *112*, 2208–2267.
- (94) Jiang, W.; Li, Y.; Wang, Z. Heteroarenes as high performance organic semiconductors. *Chem. Soc. Rev.* **2013**, *42*, 6113–6127.
- (95) Lehnherr, D.; Waterloo, A. R.; Goetz, K. P.; Payne, M. M.; Hampel, F.; Anthony, J. E.; Jurchescu, O. D.; Tykwinski, R. R. Isomerically Pure syn-Anthradithiophenes: Synthesis, Properties, and FET Performance. *Organic Letters* **2012**, *14*, 3660–3663.
- (96) Ando, S.; Nishida, J.-i.; Fujiwara, E.; Tada, H.; Inoue, Y.; Tokito, S.; Yamashita, Y. Novel p- and n-Type Organic Semiconductors with an Anthracene Unit. *Chem. Mater.* **2005**, *17*, 1261–1264.
- (97) Yao, J.; Yu, C.; Liu, Z.; Luo, H.; Yang, Y.; Zhang, G.; Zhang, D. Significant Improvement of Semiconducting Performance of the Diketopyrrolopyrrole–Quaterthiophene Conjugated Polymer through Side-Chain Engineering via Hydrogen-Bonding. *J. Am. Chem. Soc.* **2016**, *138*, 173–185.
- (98) Illig, S.; Eggeman, A. S.; Troisi, A.; Jiang, L.; Warwick, C.; Nikolka, M.; Schweicher, G.; Yeates, S. G.; Henri Geerts, Y.; Anthony, J. E. et al. Reducing dynamic disorder in small-molecule organic semiconductors by suppressing large-amplitude thermal motions. *Nat. Commun.* **2016**, *7*, 10736, Article.
- (99) Chernyshov, I. Y.; Vener, M. V.; Feldman, E. V.; Paraschuk, D. Y.; Sosorev, A. Y. Inhibiting Low-Frequency Vibrations Explains Exceptionally High Electron Mobility in

- 2,5-Difluoro-7,7,8,8-tetracyanoquinodimethane (F2-TCNQ) Single Crystals. *J. Phys. Chem. Lett.* **2017**, *8*, 2875–2880.
- (100) Giangreco, I.; Cole, J. C.; Thomas, E. Mining the Cambridge Structural Database for Matched Molecular Crystal Structures: A Systematic Exploration of Isostructurality. *Cryst. Growth Des.* **2017**, *17*, 3192–3203.
- (101) Tyrchan, C.; Evertsson, E. Matched Molecular Pair Analysis in Short: Algorithms, Applications and Limitations. *Comput. Struct. Biotechnol. J.* **2017**, *15*, 86–90.
- (102) Hussain, J.; Rea, C. Computationally Efficient Algorithm to Identify Matched Molecular Pairs (MMPs) in Large Data Sets. *J. Chem. Inf. Model.* **2010**, *50*, 339–348.
- (103) Zöphel, L.; Beckmann, D.; Enkelmann, V.; Chercka, D.; Rieger, R.; Müllen, K. Asymmetric pyrene derivatives for organic field-effect transistors. *Chem. Commun.* **2011**, *47*, 6960–6962.
- (104) Janprapa, N.; Vchirawongkwin, V.; Kritayakornupong, C. Substituent effects on furan-phenylene copolymer for photovoltaic improvement: A density functional study. *Chem. Phys.* **2018**, *510*, 60–69.
- (105) Maddox, J. Crystals from first principles. *Nature* **1988**, *335*, 201.
- (106) Reilly, A. M.; Cooper, R. I.; Adjiman, C. S.; Bhattacharya, S.; Boese, A. D.; Brandenburg, J. G.; Bygrave, P. J.; Bylisma, R.; Campbell, J. E.; Car, R. et al. Report on the Sixth Blind Test of Organic Crystal Structure Prediction Methods. *Acta Crystallogr. B* **2016**, *72*, 439–459.

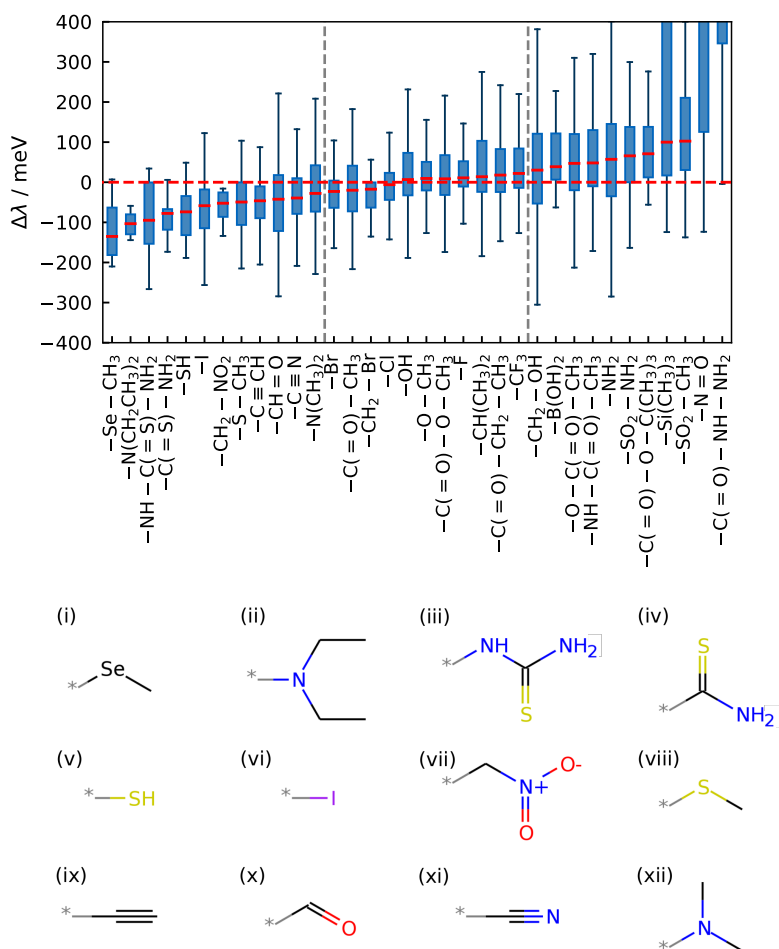


Figure 3: Top: Boxplot showing the distributions of relative reorganization energies $\Delta\lambda$ of 34 sidegroups, see text. They are sorted by their median $\Delta\lambda_{\text{med,sidegroup}}$ which is marked by a solid red line. Blue boxes extend from lower to upper quartiles of the data, while the vertical lines (whiskers) delineate the most extreme points. Gray-dashed vertical lines serve as separators to distinguish sidegroups that generally improve (i.e. lower) the reorganization energy ($\Delta\lambda_{\text{med,sidegroup}} < -25, \text{meV}$, to the left in the plot) and sidegroups that generally worsen (i.e. increase) the reorganization energy ($\Delta\lambda_{\text{med,sidegroup}} > 25, \text{meV}$, to the right in the plot). Due to space constraints, the labels for the different sidegroups use (=O/S) to denote branches, similar as is done in canonical smiles strings.⁶⁶ Bottom: Identified 12 most favorable sidegroups appearing on the left hand side of the boxplot.

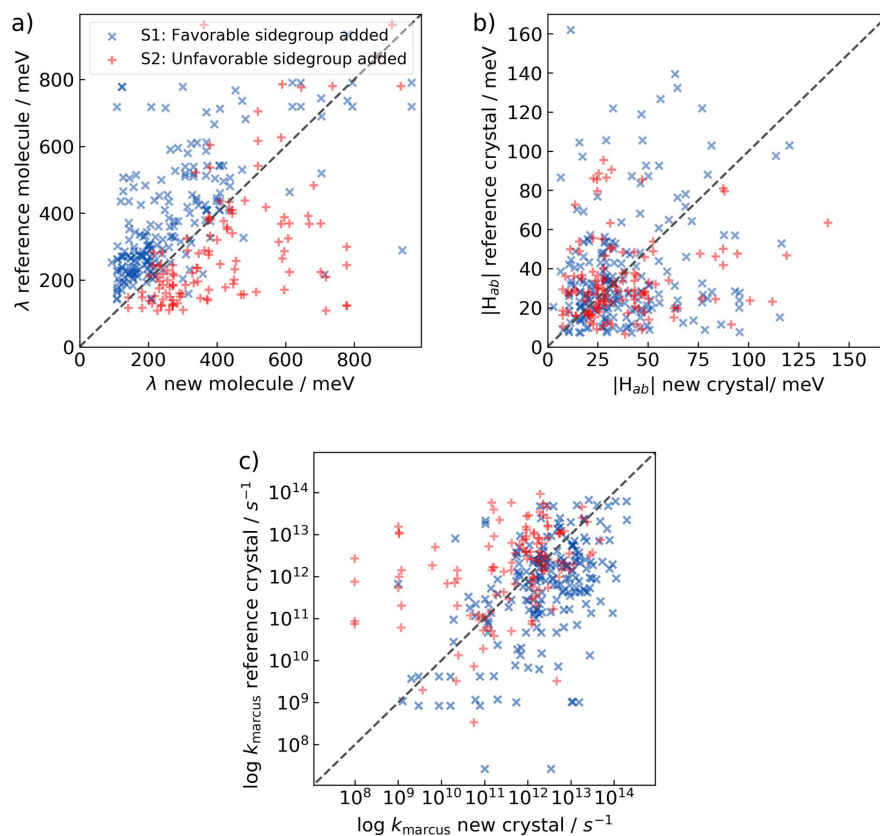


Figure 4: Charge-transfer parameter tuning by sidegroup exchange for the series of molecules contained in the “optimized” testset S1 (blue data points) and the control testset S2 (red data points), which was deliberately made worse. Shown are changes in a) λ , b) $|H_{ab}|$ and c) Marcus type hopping rate k_{Marcus} when going from a reference molecule to a new molecule (see text).