# Supporting Information

# Accelerated Diradical Character Assessment in Large Datasets of Polybenzenoid Hydrocarbons Using xTB Fractional Occupation

Alexandra Wahab[a] and Renana Gershoni-Poranne*[b]

[a]*Laboratory for Organic Chemistry, Department of Chemistry and Applied Biosciences, ETH Zürich, 8093 Zürich, Switzerland*

[b]*Schulich Faculty of Chemistry, Technion – Israel Institute of Technology, Haifa 32000, Israel*

e-mail: rporanne@technion.ac.il

# Contents

# Abbreviations and Acronyms

cc-PBH      *cata*-condensed polybenzenoid hydrocarbon

DFT         density functional theory

FOD         fractional occupation number weighted electron density

FT-DFT      finite-temperature density functional theory

HF          Hartree–Fock

PBH         polybenzenoid hydrocarbon

pc-PBH      *peri*-condensed polybenzenoid hydrocarbon

RKS         restricted Kohn–Sham

UKS         unrestricted Kohn–Sham

xTB         extended tight-binding

# S1 Computational Methods

In this section we provide templates for all input files used in the course of this work. These include benchmarking calculations as well as subsequent calculation of the diradical metrics. All semiempirical calculations were performed at the GFN2-xTB level with xtb[S1] version 6.2. All DFT and HF calculations were performed with ORCA 5.0.3.[S2,S3]

## S1.1 Geometry optimization

The following input templates were used to optimize the geometries of the molecules studied in this work. We used the CAM-B3LYP[S4–S8] functional and the def2-SVP[S9] basis set for optimization, with Grimme's D3[S10] dispersion correction and the Becke-Johnson damping scheme.[S11,S12]

Molecules with low/negligible diradical character (Groups A and C) were calculated with restricted Kohn–Sham (RKS), while molecules with marked diradical character (Groups B and D) were calculated with unrestricted Kohn–Sham (UKS). As detailed in the main text, this determination was performed using different parameters. For the *cata*-condensed polybenzenoid hydrocarbon (cc-PBH) dataset, the threshold was defined by the longest linear stretch in the molecule, such that molecules with $n_{\mathrm{LL}} \geq 6$ were considered to be predominantly diradical. For the *peri*-condensed-PBH (pc-PBH) dataset, we used the $N_{\mathrm{FOD}}^{\mathrm{xTB}}$ threshold we had established to classify the data, such that molecules with $N_{\mathrm{FOD}}^{\mathrm{xTB}} \geq 1.2$ were considered to be predominantly diradical.

The following is the ORCA input file for Groups A and C.

```
# functional and basis set
! cam-b3lyp def2-svp

# accuracy, approximations, and dispersion corrections
! tightscf rijcosX def2/j d3bj

# type of calculation
! opt

# output control
! miniprint

%base"cam-b3lyp_def2-svp_opt_0"

# type of input; charge; multiplicity; input
*xyzfile 0 1 filename.xyz
```

The following is the ORCA input file for Groups B and D.

```
# functional and basis set
! cam-b3lyp def2-svp

# accuracy, approximations, and dispersion corrections
! tightscf rijcosX def2/j d3bj

# type of calculation
! opt uks

# output control
! miniprint

%base"cam-b3lyp_def2-svp_opt_0"

# type of input; charge; multiplicity; input
*xyzfile 0 1 filename.xyz
```

## S1.2 The fractional occupation number weighted electron density (FOD) method

### S1.2.1 With extended tight-binding (xTB)

We calculated the $N_{\text{FOD}}^{\text{xTB}}$ values at a temperature of 5,000 K, as suggested in the xTB documentation.[S13–S15] To obtain the $N_{\text{FOD}}^{\text{xTB}}$, the following command line was used, where cam-b3lyp_def2-svp_opt_0.xyz is the DFT-optimized geometry (see input above in section S1.1).

```
xtb cam-b3lyp_def2-svp_opt_0.xyz --fod --etemp 5000
```

### S1.2.2 With density functional theory (DFT)

To obtain $N_{\text{FOD}}^{\text{DFT}}$ values we performed finite-temperature density functional theory (FT-DFT) calculations (see input below) with the B3LYP[S4–S7] functional the def2-SVP basis set as used for the geometry optimization. The calculations were run at a temperature of 9,000 K as suggested by Grimme and coworkers,[S16] who derived the following empirical formula to obtain the optimal temperature:

$$T_{\text{el}} = 20\,000\,\text{K} \times a_x + 5000\,\text{K}$$

where $a_x$ is the percentage of the included Fock exchange of the functional ($a_x = 20\%$ for B3LYP).

Unlike for the geometry optimizations, we did not use the CAM-B3LYP functional as the Fock exchange, $a_x$, scales with the range.

```
# functional and basis set
! b3lyp def2-svp

# accuracy
! tightscf

%scf
SmearTemp 9000
end


# type of input; charge; multiplicity; input
*xyzfile 0 1 cam-b3lyp_def2-svp_opt_0.xyz
```

## S1.3   Yamaguchi $y$ value

The $y$ values were obtained with a broken symmetry calculation at the HF-level with the same basis set used for geometry optimization, def2-SVP.

```
# method and basis set
! uhf def2-svp

# accuracy
! tightscf

# type of calculation
! uno

%scf brokensym 1,1 end

# type of input; charge; multiplicity; input
*xyzfile 0 1 cam-b3lyp_def2-svp_opt_0.xyz
```

In the main text, we show the following equation to calculate the Yamaguchi $y_i$ value:[S17,S18]

$$y_i = 1 - \frac{4(n_{HONO-i} - n_{LUNO+i})}{4 + (n_{HONO-i} - n_{LUNO+i})^2} \tag{S1}$$

In practice, we used a different formulation of $y_i$, which is expressed in terms of the corresponding orbital overlap $T_i$:

$$y_i = 1 - \frac{2T_i}{1 + T_i^2} \tag{S2}$$

where

$$T_i = \frac{n_{HONO-i} - n_{LUNO+i}}{2} \tag{S3}$$

because it was shown that $T_i$ can be directly calculated from the natural orbital occupation numbers[S17,S18] and is obtained from the ORCA output file straightforwardly.

Equation S1 is obtained by plugging equation S3 into equation S2.
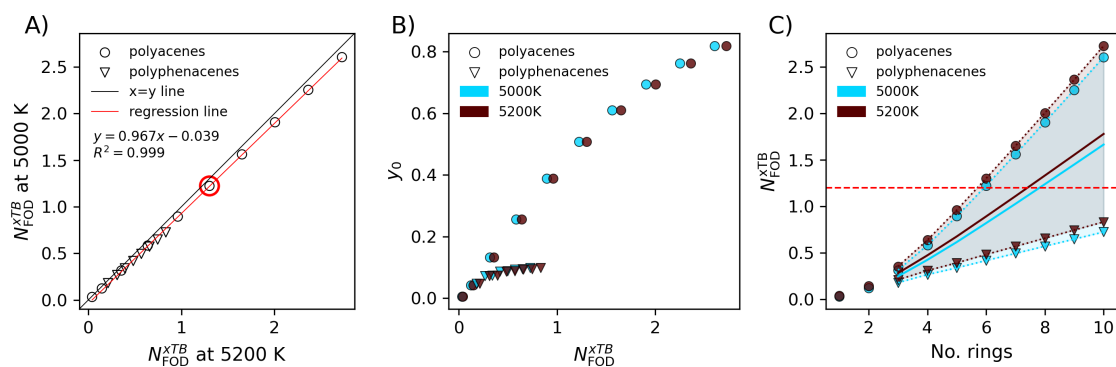
# S2 The effect of $T_{\text{el}}$

In this work, we performed all $N_{\text{FOD}}^{\text{xTB}}$ calculations using $T_{\text{el}} = 5{,}000$ K, as recommended by Grimme and coworkers.[S14] Very shortly before finalizing the current manuscript, a new report by Lischka and coworkers recommended using $T_{\text{el}} = 5{,}200$ K.[S19]

To investigate the effect of changing $T_{\text{el}}$, we calculated $N_{\text{FOD}}^{\text{xTB}}$ at the suggested optimal temperature $T_{\text{el}} = 5{,}200$ K for our benchmarking set of 18 molecules, which comprises all polyacenes and polyphenacenes with $n_{rings} \leq 10$ (benzene and naphthalene are included in the polyacene series). The results of these calculations are detailed in Table S2.

**Table S2.** $y$ and $N_{\text{FOD}}^{\text{xTB}}$ values, obtained at $T_{\text{el}} = 5{,}000$ K and 5,200 K, for the benchmarking molecules, listed by $n_{\text{rings}}$.

| $n_{\text{rings}}$ | Polyacenes | | | Polyphenacenes | | |
|---|---|---|---|---|---|---|
| | $N_{\text{FOD}}^{\text{xTB}}$ @ 5,000 K | $N_{\text{FOD}}^{\text{xTB}}$ @ 5,200 K | $y$ | $N_{\text{FOD}}^{\text{xTB}}$ @ 5,000 K | $N_{\text{FOD}}^{\text{xTB}}$ @ 5,200 K | $y$ |
| 1 | 0.03 | 0.04 | 1% | - | - | - |
| 2 | 0.12 | 0.15 | 4% | - | - | - |
| 3 | 0.32 | 0.36 | 13% | 0.18 | 0.21 | 5% |
| 4 | 0.58 | 0.64 | 26% | 0.27 | 0.31 | 7% |
| 5 | 0.89 | 0.96 | 39% | 0.34 | 0.39 | 7% |
| 6 | 1.22 | 1.30 | 51% | 0.42 | 0.48 | 9% |
| 7 | 1.56 | 1.65 | 61% | 0.50 | 0.57 | 9% |
| 8 | 1.91 | 2.01 | 69% | 0.57 | 0.66 | 9% |
| 9 | 2.25 | 2.36 | 76% | 0.65 | 0.75 | 9% |
| 10 | 2.60 | 2.72 | 82% | 0.73 | 0.83 | 10% |

Figure S1 shows the comparison between the values obtained at $T_{\text{el}} = 5{,}200$ K versus the values obtained at $T_{\text{el}} = 5{,}000$ K. As can be seen in Panel A, an excellent linear correlation is obtained, with a slope close to 1. The values calculated using $T_{\text{el}} = 5{,}200$ K are slightly higher, and this gap increases slightly as the diradical character increases. This means that, if one were to use $T_{\text{el}} = 5{,}200$ K, the threshold values would all be slightly raised. Nevertheless, because the difference between the two is systematic, the overall trends would remain the same and it is unlikely that specific molecules would be assigned different character (closed- or open-shell). This is also apparent from Panel B , which compares the $N_{\text{FOD}}^{\text{xTB}}$ values obtained at the two temperatures with the $y$ values. It can be seen from Panel C that the result of raising the temperature is a slight upward shift of the range of $N_{\text{FOD}}^{\text{xTB}}$ values.
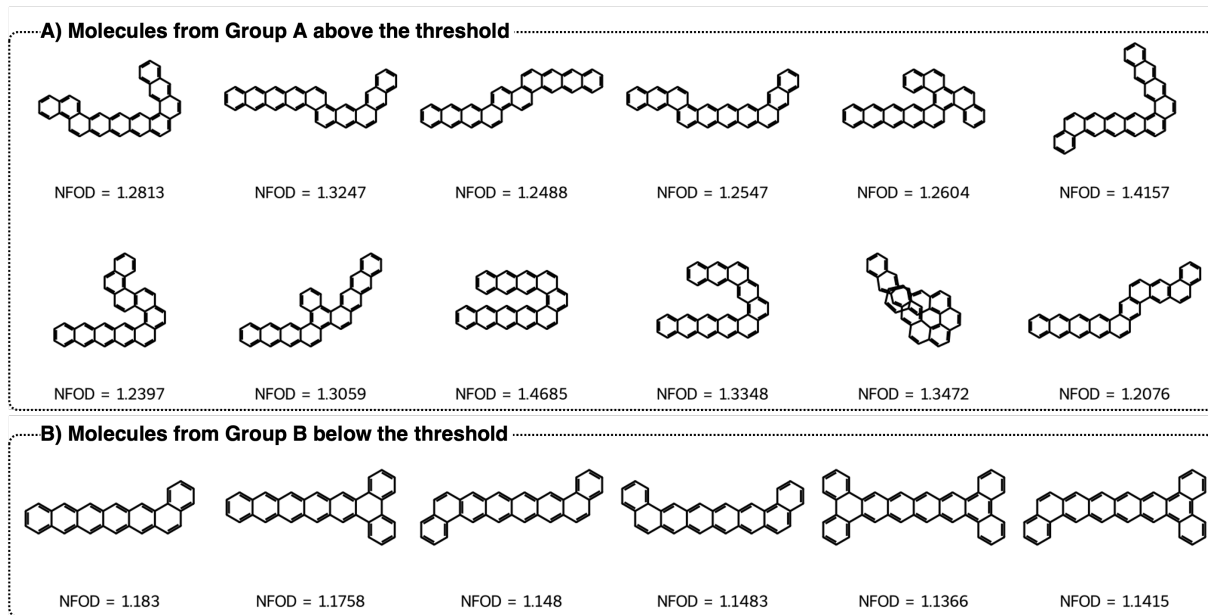
**Figure S1.** The effect of $T_{el}$ on the $N_{FOD}^{xTB}$ values for the benchmarking set. A) Scatter plot of $N_{FOD}^{xTB}$ values calculated with $T_{el} = 5{,}000$ K versus $T_{el} = 5{,}200$ K. B) Scatter plot of $N_{FOD}^{xTB}$ at the two temperatures versus $y_0$. C) $N_{FOD}^{xTB}$ values at the two temperatures arranged according to $n_{rings}$. The colored areas show the ranges of the $N_{FOD}^{xTB}$ values (min to max). The bold lines denote the respective mean values. In all plots, the values obtained at $T_{el} = 5{,}000$ K are shown in blue and the values obtained at $T_{el} = 5{,}200$ K are shown in magenta. Polyacenes are shown as circles and polyphenacenes are shown as triangles.

# S3    Examples of mistaken classification when using $n_{\mathrm{LL}}$

In this work, the chemical space of all *cata*-condensed PBHs with $n_{rings} \leq 10$ was divided into Groups A and B based on a structural parameter: the longest linear stretch, $n_{\mathrm{LL}}$. Group A contained all molecules with $n_{\mathrm{LL}} < 6$ and Group B contained all molecules with $n_{\mathrm{LL}} \geq 6$. This division was made based on the large body of literature showing that diradical character becomes non-negligible for polyacenes starting from hexacene.

As we note in the main text, this division apparently does a rather good job of classifying closed- versus open-shell systems. However, there are some cases that are wrongly characterized when using this structural parameters, as becomes evident when the $y$ and $N_{\mathrm{FOD}}$ diagnostics are calculated. This can be seen in Figure 4A of the main text: some Group A molecules cross the threshold value, and some Group B molecules appear below the threshold.

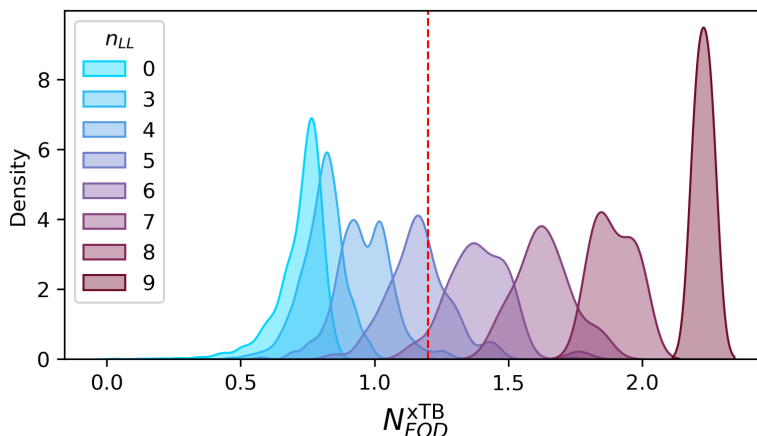In Figure S2 we present representative examples of molecules that are mistakenly classified when using the $n_{\mathrm{LL}}$ parameter.



**Figure S2.** Representative examples of molecules wrongly classified when using $n_{\mathrm{LL}}$ as the basis for classification. A) Group A molecules that are open-shell. B) Group B molecules that are closed-shell.

It appears that the first case (molecules in Group A that were wrongly classified as closed-shell) includes molecules containing two linear stretches of at least 4 rings. For such structures, even though $n_{\mathrm{LL}} < 6$, the existence of two relatively long linear stretches increases the diradical character. The second case (molecules in Group B that were wrongly classified as open-shell) includes molecules that have $n_{\mathrm{LL}} = 6$ with additional angular or trigonal annulations at one or two edges of the linear stretch. This indicates that the diradical character associated with the hexacene moiety is tempered by addition of more rings in any type of annulation that does not extend the linear stretch.

# S4 The effect of the longest linear stretch

As shown in the main text and in our previous work,[S20–S22] the longest linear stretch (for which $n_{\mathrm{LL}}$ represents the number of rings annulated linearly) is the structural feature most closely associated with the extent of diradical character (and small HOMO-LUMO gaps). Therefore, it serves as a simple and convenient feature for estimating the diradical character of cc-PBHs.

Figure S3 below shows the KDE distributions of the $N_{\mathrm{FOD}}^{\mathrm{xTB}}$ values, grouped by $n_{\mathrm{LL}}$. This figure nicely shows upward trend, whereby increasing $n_{\mathrm{LL}}$ corresponds to an increase in diradical character cc-PBHs.



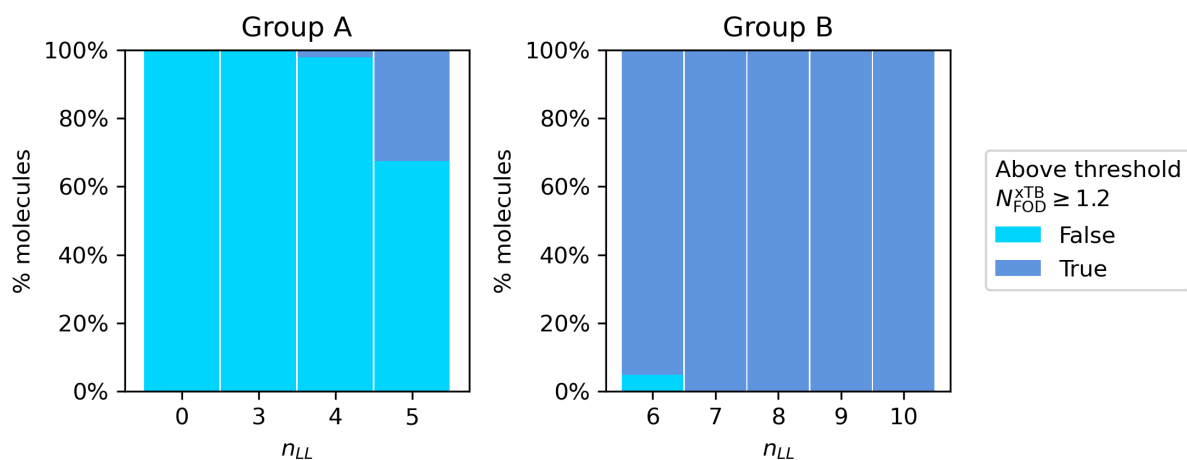**Figure S3.** KDE plots of the $N_{\mathrm{FOD}}^{\mathrm{xTB}}$ of the cc-PBH data grouped by longest linear stretch ($n_{\mathrm{LL}}$). The set threshold of 1.2 is indicated by a vertical red line.

Nonetheless, we note that the KDE distributions overlap significantly. In particular, we observe that both the $n_{\mathrm{LL}} = 5$ group and $n_{\mathrm{LL}} = 6$ groups have distributions on either side of the $N_{\mathrm{FOD}}^{\mathrm{xTB}} = 1.2$ threshold value (vertical dashed red line). Thus, using this structural feature to distinguish between closed- and open-shell molecules is not ideal.
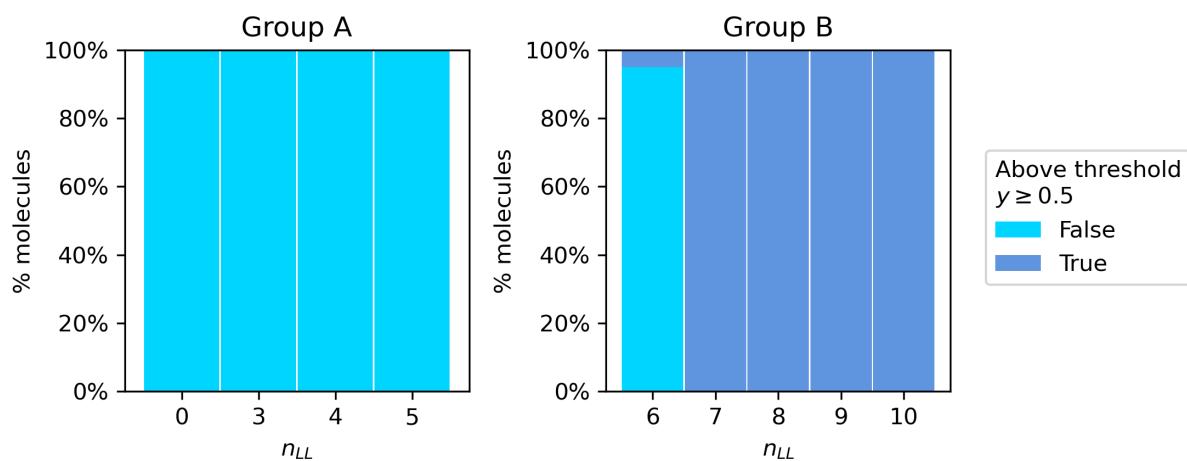
This is further confirmed by the next two figures, in which we plot the percentage of molecules (in each $n_{\mathrm{LL}}$ group) which are identified as being closed-shell (light blue) or open-shell (dark blue) based on the $N_{\mathrm{FOD}}^{\mathrm{xTB}}$ threshold (Figure S4) and the $y$ value (Figure S5).

We perform this analysis for the two datasets, Group A ($n_{\mathrm{LL}} < 6$) and Group B ($n_{\mathrm{LL}} \geq 6$), which are described in the main text. As seen in Figure S4, our $N_{\mathrm{FOD}}^{\mathrm{xTB}}$ threshold categorizes some molecules from Group A as having significant diradical character, and, *vice versa*, in Group B it identifies some molecules as closed-shell.

When the $y = 50\%$ value is used as the threshold (Figure S5), all the molecules in Group A are found to be closed-shell. However, a larger number of molecules in Group B are identified as closed-shell (all of them with $n_{\mathrm{LL}} = 6$). Overall, the $N_{\mathrm{FOD}}^{\mathrm{xTB}}$ threshold marks more molecules as open-shell than the $y$ value suggests. These observations align with our analysis in the main text, where we show that the $N_{\mathrm{FOD}}^{\mathrm{xTB}} = 1.2$ threshold is too conservative, i.e., it results in more false positives than false negatives.

**Figure S4.** Percentages of molecules in each $n_{\mathrm{LL}}$ group categorized as below (light blue) or above (dark blue) the threshold ($N_{\mathrm{FOD}}^{\mathrm{xTB}} = 1.2$). Group A contains all molecules with $n_{\mathrm{LL}} < 6$ and group B contains molecules with $n_{\mathrm{LL}} \geq 6$.



**Figure S5.** Percentages of molecules in each $n_{\mathrm{LL}}$ group categorized as below (light blue) or above (dark blue) the threshold ($y = 0.5$). Group A contains all molecules with $n_{\mathrm{LL}} < 6$ and group B contains molecules with $n_{\mathrm{LL}} \geq 6$.
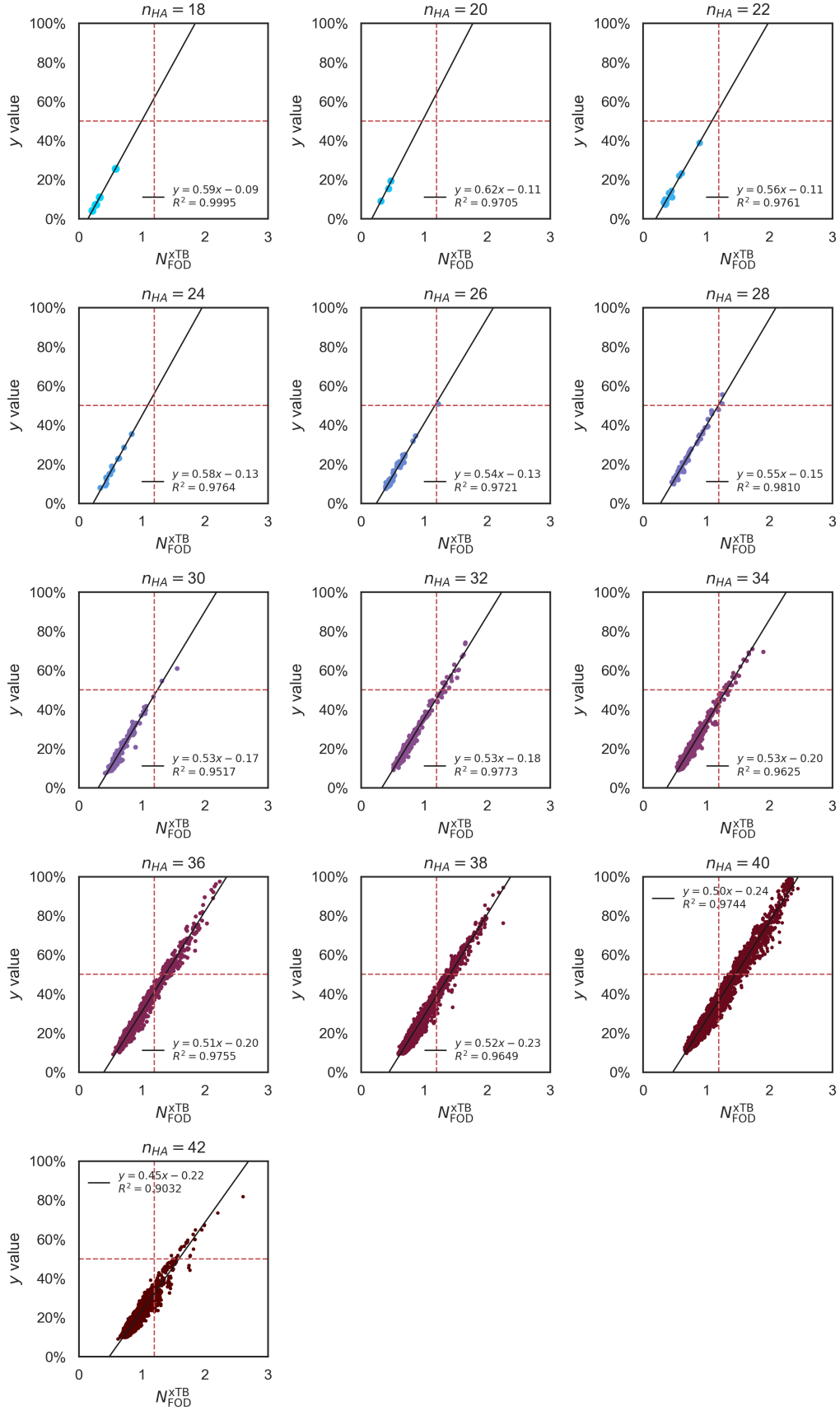
# S5 The effect of size

In the main text, we demonstrate the effect of size on the $N_{\text{FOD}}^{\text{xTB}}$ threshold value and, therefore, conclude that different thresholds should be identified, based on the molecular size. In our analysis, we partition the dataset according to the number of heavy (i.e., non-hydrogen) atoms, $n_{\text{HA}}$, as a measure of size. Table S3 below details the number of molecules in each of the $n_{\text{HA}}$-groups.

**Table S3.** Size of the entire dataset when grouped by number of heavy atoms *nha*.

| $\mathbf{n_{HA}}$ | Dataset size |
|---|---|
| 6 | 1 |
| 10 | 1 |
| 14 | 2 |
| 16 | 1 |
| 18 | 5 |
| 20 | 3 |
| 22 | 14 |
| 24 | 14 |
| 26 | 46 |
| 28 | 69 |
| 30 | 183 |
| 32 | 340 |
| 34 | 811 |
| 36 | 1627 |
| 38 | 3618 |
| 40 | 6704 |
| 42 | 6436 |

To obtain the size-aware thresholds, we performed a linear regression on the $N_{\text{FOD}}^{\text{xTB}}$ against $y$ scatter, for each $n_{\text{HA}}$-group. From the obtained fitting equation, we identified the $N_{\text{FOD}}^{\text{xTB}}$ value corresponding to $y = 0.5$. The regressions of each subset are plotted in Figure S6. In all cases, $R^2 \geq 0.95$. The groups with $n_{\text{HA}} < 18$ were not fit, as they have too few points.
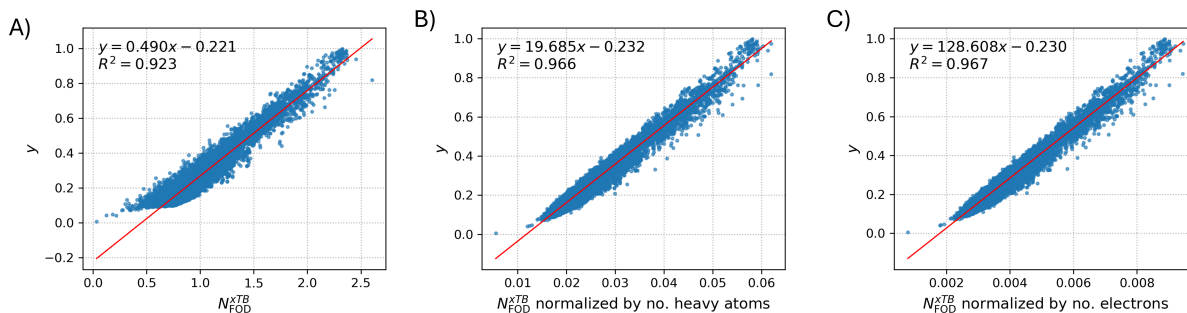
**Figure S6.** Scatter plots and fitting equations of $N_{\mathrm{FOD}}^{\mathrm{xTB}}$ vs $y$ for each $n_{\mathrm{HA}}$-subset starting from $n_{\mathrm{HA}} = 18$.

# S6 Comparison between normalization by $n_e$ versus $n_{HA}$

Due to the size-extensivity of the $N_{FOD}$ metric, we use the number of heavy atoms, $n_{HA}$, to separate the chemical space into groups based on their molecular size. As we note in the main text, Grimme and coworkers previously suggested a similar but different approach: to use the number of electrons, $n_e$, to normalize the $N_{FOD}$ value. In this section, we address two questions. The first is: how does normalizing the data impact the agreement with $y$ and the ability to identify open-shell molecules, as compared using 'raw' $N_{FOD}$ data? In this context, we also explore whether there is a difference when using $n_{HA}$ or $n_e$ for normalization. The second question is: does the size-aware approach improve when using $n_e$ as the structural parameter rather than $n_{HA}$?
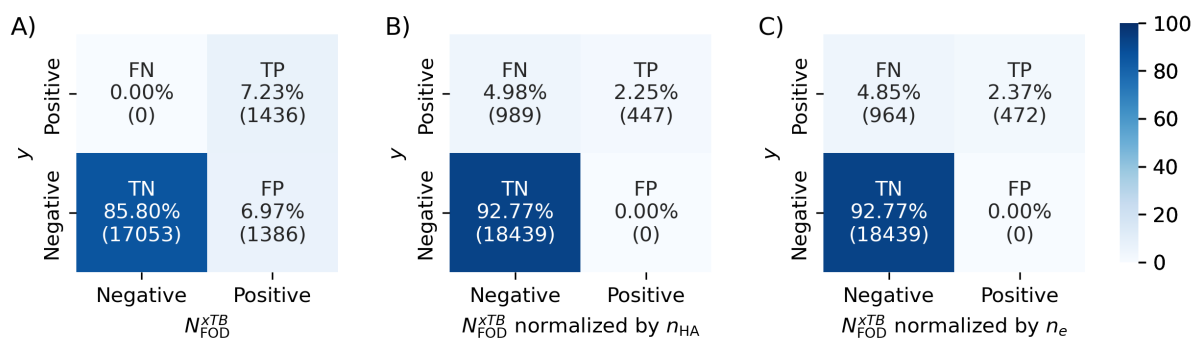
To answer the first question, we first calculated two types of normalized $N_{FOD}$ values for the entire dataset by dividing the $N_{FOD}$ value by $n_e$ and by $n_{HA}$, respectively. We then plotted the normalized values against the respective $y$ values, for each molecule in the dataset. For comparison, we present these results alongside a similar scatter plot of the raw (un-normalized) data in Figure S7.



**Figure S7.** Scatter plots of A) $N_{FOD}^{xTB}$, B) $N_{FOD}^{xTB}$ normalized by $n_{HA}$, and C) $N_{FOD}^{xTB}$ normalized by $n_e$ against $y$.

It is apparent that normalization does improve the correlation to $y$, and that both normalization schemes afford the same improvement. However, the improvement in correlation does not necessarily lead to an improvement in classification. To investigate this, we once again used hexacene as the threshold for identifying open-shell character. The normalized $N_{FOD}^{xTB}$ of hexacene (0.04616 and 0.00698, for $n_{HA}$ and $n_e$, respectively) was applied to the respectively normalized data. The performance of the different thresholds is presented in Figure S8.

As can be seen in the confusion matrices presented in Figure S8, both normalization approaches perform essentially the same, finding a TP rate of approximately 2% and a TN rate of approximately 93%. In both cases there is a decrease in FP identifications, from approximately 6% to 0%, concurrent with an increase in FN identifications from 0% to approximately 5%. In other words, the normalized thresholds cause the majority of diradical molecules to be mistakenly identified as closed-shell, which is opposite to the result of the original threshold (Figure S8A), with which many closed-shell molecules were mistakenly identified as diradical. This indicates that while the 1.2 threshold is too low, the normalized thresholds are too high. To conclude our answer to the first question: normalizing the $N_{FOD}^{xTB}$ values improves their correlation to $y$, but does not meaningfully

**Figure S8.** Confusion matrices for evaluating the classification based on A) original threshold of $N_{\mathrm{FOD}}^{\mathrm{xTB}} = 1.2000$ , B) $n_{\mathrm{HA}}$-normalized threshold of $N_{\mathrm{FOD}}^{\mathrm{xTB}} = 0.04616$, and C) $n_{\mathrm{e}}$-normalized threshold of $N_{\mathrm{FOD}}^{\mathrm{xTB}} = 0.00698$.

improve the ability to identify a threshold for classification of the molecules into closed- and open-shell. The normalized threshold are slightly better than the original threshold of $N_{\mathrm{FOD}}^{\mathrm{xTB}} = 1.2$, but considerably worse than the size-aware approach (which had a cumulative false identification rate of $< 1\%$, see Figure S9A).

To answer the second question, we repeated our size-based analysis, this time separating the molecules into groups according to their $n_{\mathrm{e}}$. We obtained fitting equations and identified the size-aware threshold for each group. We then compared the classification of closed- or open-shell to the $y$ values, and obtained the TP, TN, FP, and FN percentages. The results of this analysis are presented in Figure S9, alongside the same evaluation performed using $n_{\mathrm{HA}}$.



**Figure S9.** Confusion matrices assessing the performance of the different threshold approaches: A) single-threshold using $N_{\mathrm{FOD}}^{\mathrm{xTB}} = 1.2$; B) size-aware threshold based on $n_{\mathrm{HA}}$ and C) size-aware threshold based on $n_{\mathrm{e}}$. TP, FN, FP, and TN percentages and numbers (in parentheses) are noted. The color corresponds to the percentage. Groups containing fewer than 3 molecules were omitted (a total of 9 molecules).

As can be seen in the confusion matrices in Figure S9, there is no difference between using $n_{\mathrm{HA}}$ or $n_{\mathrm{e}}$ to obtain the size-aware thresholds. Both approaches afford similar performance in terms of identifying closed- versus open-shell PBHs.
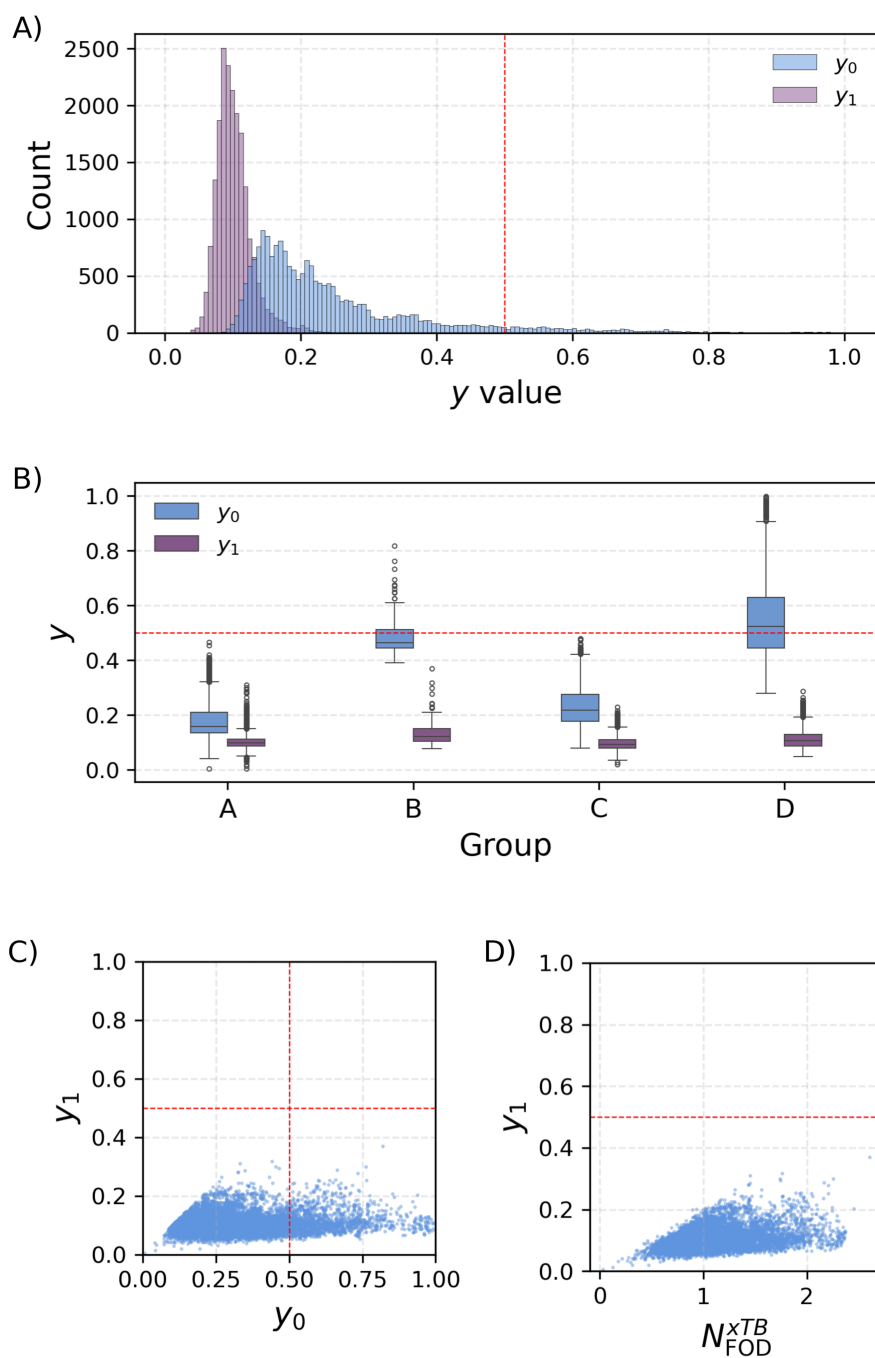
# S7 Investigation of $y_1$ values

In this work, we use the Yamaguchi $y$ value as the 'ground truth' of diradical character. As we mentioned in the main text and in the Computational Details section, there are several Yamaguchi $y$ values, which are given by Equation S1. For a 'perfect' closed-shell system, the occupation of the HONO is 2 and the occupation of the LUNO is 0; accordingly, $y_0 = 0$ and $y_1 = 0$. For a 'perfect' diradical system (i.e., originating from a transition of a single electron from the HONO to the LUNO), the occupation of both the HONO and the LUNO is 1, resulting in $y_0 = 1$ and $y_1 = 0$. For a tetraradical case (i.e., a single electron in each of the orbitals HONO-1, HONO, LUNO, and LUNO+1, respectively), $y_0 = 1$ and $y_1 = 1$.

In this study, we used only the $y_0$ value to evaluate the PBHs under investigation, and comapred it to the $N_{\mathrm{FOD}}$ value. In principle, these two diagnostics do not measure the same thing. The $y_0$ value measures *diradical* character, whereas the $N_{\mathrm{FOD}}$ values measures *open-shell* character in a more general sense (i.e., may include other forms of polyradical character). Hence, one possible concern is that high $N_{\mathrm{FOD}}$ values are actually indicating tetraradical character, rather than diradical.

In practice, however, this appears not to be a substantial concern. The excellent agreement we obtain between the $y_0$ value and the $N_{\mathrm{FOD}}$ values already suggests that they are, in fact, measuring the same diradical behavior. In other words, it suggests that the open-shell character of the PBHs in this study originates from diradical character, with little to no contribution from other sources of polyradical character. We further verified this by calculating the $y_1$ values, which measure the tetraradical character of the molecules, for all molecules in this study.

Figure S10A shows the histograms of $y_0$ and $y_1$ values for all molecules. As can be seen, the $y_1$ values are centered around 0.1 and do not exceed 0.2 for most molecules, which indicates that tetraradical character is negligible for these molecules. A small number of molecules have slightly higher values, but none go beyond 0.3. In Figure S10B, we investigate the behavior of the two values across the four groups (A-D) that we defined in this work. While the $y_0$ value for the Groups A and C are quite low, the values for Groups B and D are substantially higher, which is in agreement with their expected diradical character. In contrast, for $y_1$, there is no significant change in the distribution of values across the four groups. Figures S10C and D show scatter plots of the $y_1$ values against $y_0$ and $N_{\mathrm{FOD}}^{\mathrm{xTB}}$, respectively, in which there appears to be no correlation between an increase in diradical character and an increase in tetraradical character, for both diagnostics..
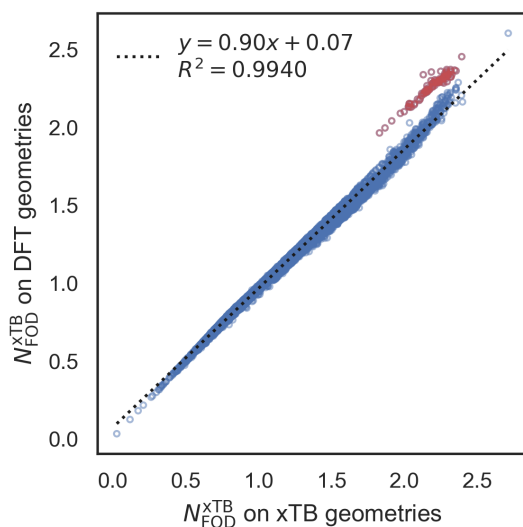
**Figure S10.** A) Histograms of $y_0$ and $y_1$ values for all molecules in this work; B) box-and-whisker plots of the $y_0$ and $y_1$ distributions for the individual groups defined in this work; C) scatter plot of $y_0$ versus $y_1$; D) scatter plot of $y_1$ vs $N_{\mathrm{FOD}}^{\mathrm{xTB}}$. The thresholds of 0.5 (50%) are marked with dashed red lines.

# S8    The effect of geometry on $N_{\text{FOD}}^{\text{xTB}}$

For consistency, we calculated all the descriptors we studied in this work, $N_{\text{FOD}}^{\text{xTB}}$, $N_{\text{FOD}}^{\text{DFT}}$, and $y$, on the same DFT-optimized geometries. However, when working with large datasets, it is common to employ semi-empirical calculations for the geometry optimizations, as well. Therefore, to investigate the consequences of such a choice, we performed a similar analysis as in the main text, with the only difference being that in the second analysis we used xTB-optimized geometries for the $N_{\text{FOD}}^{\text{xTB}}$ calculations.
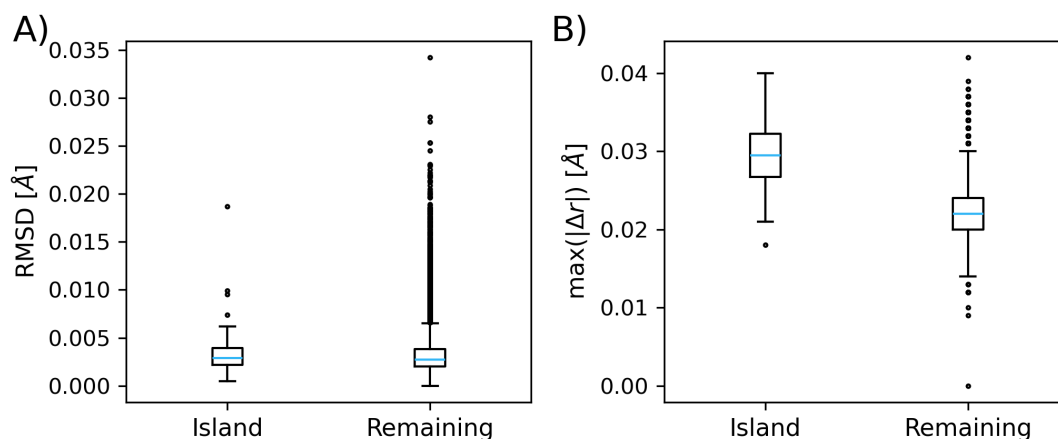
Overall, the geometries are very similar (for a more in depth comparison, we refer the reader to References S20 and S23). Nevertheless, slight changes do appear, in particular in close-to-planar geometries.

When applying the same series of steps to the xTB-optimized geometries, we found that the $N_{\text{FOD}}^{\text{xTB}}$ threshold increases from 1.2 to 1.3, as the values in general shift a bit upwards. Figure S11 shows the $N_{\text{FOD}}^{\text{xTB}}$ values calculated on xTB-geometries versus the ones calculated on DFT-geometries. A very good correlation is observed between the two datasets ($R^2 = 0.9940$). The slope of the fitting equation is 0.9 and the offset is 0.7. This suggests that the differences are systematic, and that the variations in geometry optimization are uniform throughout the structures included in the dataset.



**Figure S11.** Scatter plot of $N_{\text{FOD}}^{\text{xTB}}$ values calculated for xTB- vs DFT-geometries. All molecules in this work are included. The regression line is indicated by the black stripped line. The fitting equation and $R^2$ are written on the plot. The data points colored in red are the structures where the values disagree due to geometry.

Nonetheless, we note that there is an 'island' of 84 molecules (highlighted in red in Figure S11) that display a discrepancy from the rest of the data. It is possible to think of two reasons for this discrepancy. The first possibility is that they have very similar geometries, but that the slight differences cause large shifts in the $N_{\text{FOD}}^{\text{xTB}}$ values. This would suggest a very high sensitivity of the metric to geometric changes. The second possibility is that the geometries are actually sufficiently different to cause these discrepancies, which would suggest that xTB does not optimize them well. Considering the well-behaved correlation

**Figure S12.** Boxplots of: A) RMSD of heavy atoms between xTB- and DFT-optimized geometries; B) max absolute change in C-C bond length between xTB- and DFT-geometries.

for all of the other compounds, it seems that xTB is not overly sensitive to minor geometric shifts, and therefore the second explanation is more likely.

To investigate this, we calculated the root mean squared difference (RMSD) as well as the variance in C-C bond length between the xTB- and DFT-optimized geometries (Figure S12).

Although the RMSDs are similar for the molecules in the general correlation and those in the 'island', it appears that the molecules in the 'island' do have larger differences in the C-C bond lengths. Because the electronic behavior of polycyclic aromatic molecules is strongly related to bond-length alternation (an indication of delocalization), it is reasonable to conclude that such a difference in C-C bond lengths could explain the observed discrepancies $N_{\text{FOD}}^{\text{xTB}}$ values. Thus, it would appear that GFN2-xTB struggles to obtain the correct optimized geometries for this group of molecules, which were identified and found to be zethrene derivatives. Their structures are depicted in Figure S13.

hc_c40h22_0pent_5994    hc_c40h22_0pent_5995    hc_c40h22_0pent_5996    hc_c40h22_0pent_5998    hc_c40h22_0pent_5999    hc_c40h22_0pent_6000

hc_c40h22_0pent_6001    hc_c40h22_0pent_6002    hc_c40h22_0pent_6007    hc_c40h22_0pent_6015    hc_c40h22_0pent_6018    hc_c40h22_0pent_6020

hc_c40h22_0pent_6093    hc_c40h22_0pent_6106    hc_c40h22_0pent_6107    hc_c40h22_0pent_6108    hc_c40h22_0pent_6109    hc_c40h22_0pent_6110

hc_c40h22_0pent_6111    hc_c40h22_0pent_6113    hc_c40h22_0pent_6114    hc_c40h22_0pent_6625    hc_c40h22_0pent_6626    hc_c40h22_0pent_6628

hc_c40h22_0pent_6631    hc_c40h22_0pent_6644    hc_c40h22_0pent_6735    hc_c40h22_0pent_6759    hc_c40h22_0pent_6762    hc_c40h22_0pent_6788

hc_c40h22_0pent_6793    hc_c40h22_0pent_6797    hc_c40h22_0pent_6804    hc_c40h22_0pent_6827    hc_c40h22_0pent_6833    hc_c40h22_0pent_6840

hc_c40h22_0pent_6853    hc_c40h22_0pent_6866    hc_c40h22_0pent_6867    hc_c40h22_0pent_6887    hc_c40h22_0pent_6949    hc_c40h22_0pent_6955

hc_c40h22_0pent_6969    hc_c36h20_0pent_1221    hc_c36h20_0pent_1224    hc_c36h20_0pent_1225    hc_c36h20_0pent_1226    hc_c36h20_0pent_1240

hc_c36h20_0pent_1351    hc_c36h20_0pent_1354    hc_c38h20_0pent_171    hc_c38h20_0pent_850    hc_c38h20_0pent_1151    hc_c38h20_0pent_1866

hc_c40h22_0pent_6631  hc_c40h22_0pent_6644  hc_c40h22_0pent_6735  hc_c40h22_0pent_6759  hc_c40h22_0pent_6762  hc_c40h22_0pent_6788

hc_c40h22_0pent_6793  hc_c40h22_0pent_6797  hc_c40h22_0pent_6804  hc_c40h22_0pent_6827  hc_c40h22_0pent_6833  hc_c40h22_0pent_6840

hc_c40h22_0pent_6853  hc_c40h22_0pent_6866  hc_c40h22_0pent_6867  hc_c40h22_0pent_6887  hc_c40h22_0pent_6949  hc_c40h22_0pent_6955

hc_c40h22_0pent_6969  hc_c36h20_0pent_1221  hc_c36h20_0pent_1224  hc_c36h20_0pent_1225  hc_c36h20_0pent_1226  hc_c36h20_0pent_1240

hc_c36h20_0pent_1351  hc_c36h20_0pent_1354  hc_c38h20_0pent_171  hc_c38h20_0pent_850  hc_c38h20_0pent_1151  hc_c38h20_0pent_1866

**Figure S13.** Molecules that show the biggest deviations between $N_{\mathrm{FOD}}^{\mathrm{xTB}}$ calculated on DFT-geometries versus xTB-ones.

# References

(S1)  Bannwarth, C.; Ehlert, S.; Grimme, S. *J. Chem. Theory Comput.* **2019**, *15*, 1652–1671.

(S2)  Neese, F. *Wiley Interdiscip. Rev. Comput. Mol. Sci.* **2012**, *2*, 73–78.

(S3)  Neese, F. *Wiley Interdiscip. Rev. Comput. Mol.* **2022**, *12*, e1606.

(S4)  Becke, A. D. *J. Chem. Phys* **1993**, *98*, 5648–5652.

(S5)  Lee, C.; Yang, W.; Parr, R. G. *Phys. Rev. B* **1988**, *37*, 785–789.

(S6)  Miehlich, B.; Savin, A.; Stoll, H.; Preuss, H. *Chem. Phys. Lett.* **1989**, *157*, 200–206.

(S7)  Hertwig, R. H.; Koch, W. *Chem. Phys. Lett.* **1997**, *268*, 345–351.

(S8)  Yanai, T.; Tew, D. P.; Handy, N. C. *Chem. Phys. Lett.* **2004**, *393*, 51–57.

(S9)  Weigend, F.; Ahlrichs, R. *Phys. Chem. Chem. Phys.* **2005**, *7*, 3297.

(S10)  Grimme, S.; Antony, J.; Ehrlich, S.; Krieg, H. *J. Chem. Phys* **2010**, *132*, 154104.

(S11)  Johnson, E. R.; Becke, A. D. *J. Chem. Phys* **2005**, *123*, 024101.

(S12)  Grimme, S.; Ehrlich, S.; Goerigk, L. *J Comput Chem* **2011**, *32*, 1456–1465.

(S13)  https://xtb-docs.readthedocs.io/en/latest/properties.html#fractional-occupation-density-fod-calculation (accessed 11/21/2024).

(S14)  Grimme, S.; Hansen, A. *Angew. Chem. Int. Ed.* **2015**, *54*, 12308–12313.

(S15)  Bauer, C. A.; Hansen, A.; Grimme, S. *Chemistry – A European Journal* **2017**, *23*, 6150–6164.

(S16)  Grimme, S. *Angew. Chem. Int. Ed.* **2013**, *52*, 6306–6312.

(S17)  Yamaguchi, K. *Chem. Phys. Lett.* **1975**, *33*, 330–335.

(S18)  Nakano, M. *The Chemical Record* **2017**, *17*, 27–62.

(S19)  Carvalho, J. R.; Nieman, R.; Kertesz, M.; Aquino, A. J. A.; Hansen, A.; Lischka, H. *Theor Chem Acc* **2024**, *143*, 69.

(S20)  Wahab, A.; Pfuderer, L.; Paenurk, E.; Gershoni-Poranne, R. *J. Chem. Inf. Model.* **2022**, *62*, 3704–3713.

(S21)  Fite, S.; Wahab, A.; Paenurk, E.; Gross, Z.; Gershoni-Poranne, R. *J. Phys. Org. Chem.* **2023**, *36*, e4458.

(S22)  Weiss, T.; Wahab, A.; Bronstein, A. M.; Gershoni-Poranne, R. *J. Org. Chem.* **2023**, *88*, 9645–9656.

(S23)  Wahab, A.; Gershoni-Poranne, R. *Phys. Chem. Chem. Phys.* **2024**, *26*, 15344–15357.