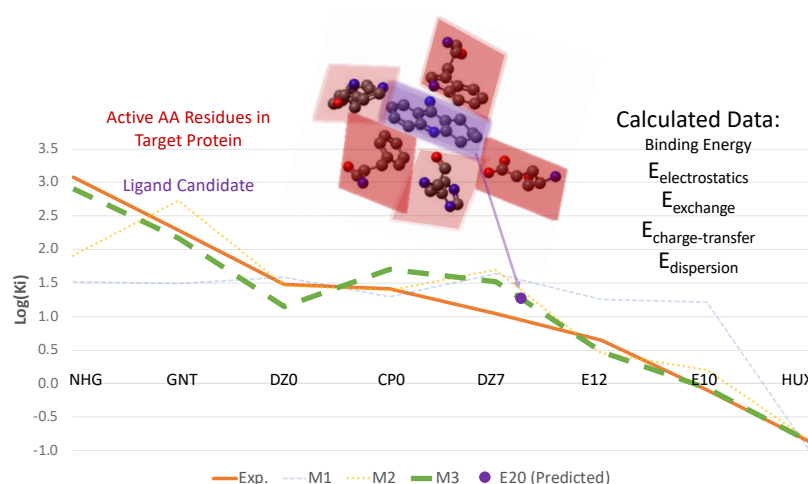# Toward Simple, Predictive Understanding of Protein-Ligand Interactions: Electronic Structure Calculations Join Forces with the Chemist's Intuition

*by Nitai Sylvetsky*[*,§]

[§] Department of Organic Chemistry, Weizmann Institute of Science, 76100 Reḥovot, Israel

KEYWORDS Protein-ligand • noncovalent interactions • localized coupled cluster • local energy decomposition • chemical intuition

ABSTRACT Contemporary efforts for modeling protein-ligand interactions entail a painful tradeoff – as reliable information on *both* noncovalent binding factors and the dynamic behavior of a protein-ligand complex is often beyond practical limits. In the following paper, we demonstrate that information drawn exclusively from *static* molecular structures can be used for the semi-quantitative prediction of experimentally-measured binding affinities for protein-ligand complexes. In the particular case considered here, inhibition constants ($K_i$) were calculated for eight different ligands of torpedo californica acetylcholinesterase using a simple, multiple-linear-regression-based model. The latter, incorporating five informative and independent variables – drawn from QM cluster, DLPNO-CCSD(T) calculations and LED analyses on the eight complexes – is shown to recover no-less-than 96% of the sum of squares for measured $K_i$ values, and used to predict the inhibition potential for yet another ligand (E20, for which no $K_i$ values are available in the literature). This thus challenges the widespread assumption that "static pictures" are inadequate for predicting reactivity properties of flexible and dynamic protein-ligand systems.

**Introduction**

Protein-ligand (PL) interactions have drawn great amounts of scientific attention throughout the last century (see Refs.[1–4] for a few recent textbooks and reviews). Besides being examined for playing crucial roles in a variety of essential biochemical processes, such interactions are often focused on in many drug design studies – revolving around finding inhibitors for proteins such as enzymes and neuroreceptors for the purpose of invoking a desirable biological response.[5–8] Due to such considerations, many researchers from a broad spectrum of scientific disciplines (consisting of computational biologists and biochemists as well as theorists from chemistry and physics) have attempted to provide some general theoretical/computational modeling schemes for predicting biochemically-relevant PL binding events.[9–13]

It has been well-established that PL systems are greatly influenced by noncovalent interactions (NCIs).[14–19] The latter, resulting from subtle electronic effects, are very small in magnitude and cannot virtually be measured by experimental means. Thus, *ab initio* electronic structure methods constitute a precious (and almost exclusive) source of information on biochemically-relevant NCIs – which, in turn, often used for the parametrization and calibration of more approximate computational modeling techniques (such as DFT functionals and molecular mechanics force fields).[20–23] Ideally, one could use such nonempirical electronic structure methods for running molecular dynamics (MD) simulations on realistic PL systems; in such scenario, the information drawn from such simulations would include an adequate description of biochemically-significant NCIs, and it can thus be expected to offer desirable predictive power (which is, after all, the main goal of any theoretical model). However, electronic structure calculations are notorious for their steep computational cost scaling with the system's size – which generally precludes using them for MD simulations on realistically-sized biochemical systems (excluding a few recent approximate approaches, each entailing different methodological challenges; see, for instance, Refs.[24,25]).

For this reason, and since *some* description of NCIs relevant for PL binding is clearly crucial for predictive purposes,[26,27] electronic structure calculations are mostly combined with additional computational techniques used for describing the dynamic, continuous relationship between PL pairs that leads to biochemically-significant (active-site) binding. In this manner, electronic structure calculations are performed on static structures, which are assumed to represent crucial events in the PL binding process (see Ref.[28] for a recent, comprehensive

review). It is generally assumed, for instance, that the actual biochemically-significant binding event – taking place in the protein's active site – must incorporate some description of noncovalent binding factors. Thus, one common piece of information on PL interactions provided by electronic structure methods corresponds to the PL binding energy – calculated as the energetic difference between the bound PL structure and its underlying protein and ligand structures found at infinite separation (Eq. 1):

$$\Delta E_{bind} = E_{PL} - (E_P + E_L) \quad [1]$$

Were P and L stand for protein and ligand, respectively (in their complex-structure geometry). It should be pointed out that the relationship between such calculated energetics and realistic PL systems is quite unclear (as said, PL binding is a continuous, dynamic process; representing it using such "binary" means – i.e., bound complex vs. free structures – clearly ignores this fact); still, quite a few authors have employed such quantities as bits and pieces of information in more-general predictive theoretical/computational schemes – where additional such pieces, obtained using different techniques (e.g., classical MD trajectories), are also used.[29–34] Needless to say, such multi-method efforts require an appropriate multi-method-expertise from the researcher, and entail lots of (perhaps undetectable) sources of error and technical difficulties – as demonstrated in Figure 1.

When interested in predictive modeling of PL systems, we are therefore faced with a painful dilemma (Figure 2). An appropriate description on biochemically-relevant NCIs is, on the one hand, required; the dynamic relationship between PL pairs cannot, on the other hand, be ignored; holding on to one source of information and letting go of the other would make our inquiry simple and elegant, but wrong and unreliable; trying to hold on to both complicates things further, as reasonable interfaces between different kinds of information must be established – giving rise to many corresponding sources of error that cannot necessarily be assessed.

**Conventional QC-based Approach:**

**Starting Point**

Motivation: finding potential inhibitors for, e.g., some enzyme of interest EZ

**Stage 1. Establishing a Pool of Candidates**

Identify a candidate inhibitor α – assumed to "behave" similarly to known ligands. Sources of information: chemical intuition/docking techniques/QSAR methods...

**Stage 2. Candidate Verification**

Confirm/refute the existence of a stable α–EZ complex:

*2.1 Ideal (often unrealistic) way:* Scan the potential energy hypersurface for such complex using QC methods → locate local/global minima, representing bound structures

*2.2 Approximate (practical) way:* Use classical MD simulations for sampling *potentially-bound* structures (e.g., where α is within binding distance to *active* amino acid residues) → investigate further using QC or QM/MM methods.

In case a stable α–EZ complex *has been found*

**Stage 3. Evaluating the Candidate's Inhibition Potential**

Compare α's *binding affinity* ($K_i$) with that of known ligands ($L_n$). But how can $K_i$ be *assessed*?

*3.1 $\Delta E_{bind}$ (equation 1 in main text) as a heuristic measure:* Let X and Y be the binding energies for the $L_1$–EZ and α–EZ complexes, respectively. If X<Y → $K_i$ of α is larger, and *vice versa*.

Binding affinities often reflect *much* more information than calculated binding energies! (Lock-key model: many dynamic factors should also be taken into account)

*3.2 Classical-heuristic approach:* Run classical MD simulations for both ligands → come up with *ad-hoc* measures for $K_i$ based on resulting trajectories

No explicit treatment of quantum mechanical electronic behavior, so no way of knowing whether binding forces are appropriately described!

Trajectories may be clinically dependent on initial conditions → drawing *causal-deterministic* conclusions will be difficult

*3.3 Integrative-heuristic approach:* Combine 3.1 (or alternative QC data) with 3.2 → Predict $K_i$ based on hybrid classical-MD + QC information

Exhaustive multi-method study – requires lots of technical expertise, not easy to manage!
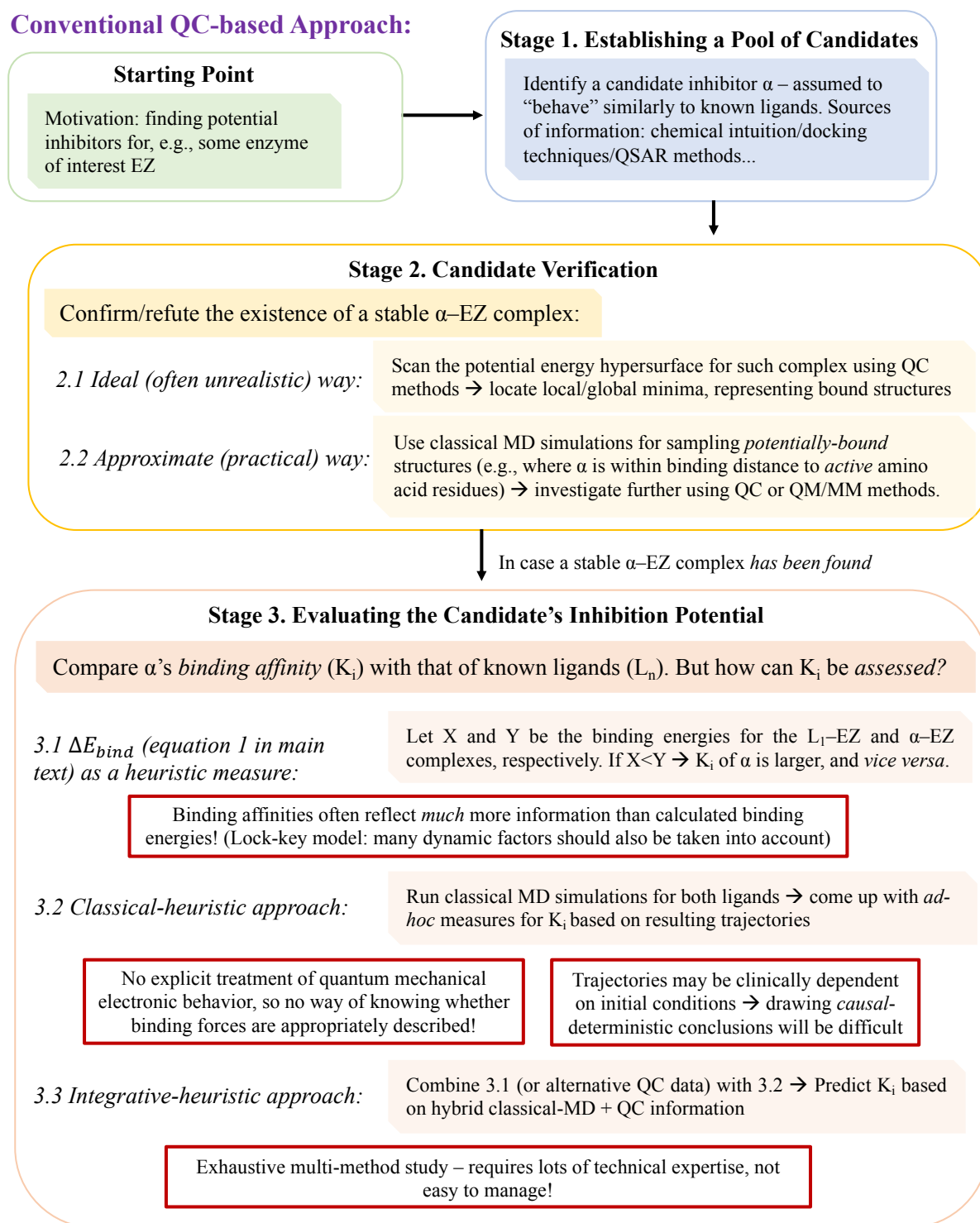
Figure 1. A hypothetical, "conventional" molecular-modeling-based ligand identification process, employing quantum chemistry methods. Compare with Figure 5, which illustrates our approach as proposed in the present paper (Acronyms: QC = quantum chemistry, MD = molecular dynamics, QM/MM = hybrid quantum chemistry – classical molecular mechanics methods. Some crucial problems threatening the process' success are outlined in red).

The main aim of this paper is demonstrating that this dilemma may often be avoided – by using electronic structure calculations on static molecular structures that *also provide some important information on the dynamic nature of PL binding processes*. In such manner, it should be possible to avoid using MD simulations altogether and still arrive at valuable predictive models – which may guide future experiments and drug discovery studies. Being mainly interested in utilizing the *information* offered by electronic-structure methods, and not in specific state-of the-art data analysis and modeling techniques, we will limit our discussion to a very simple predictive model type – based solely on multiple linear regression (MLR). The latter, incorporating independent variables from the aforementioned electronic structure calculations, will be used for a semi-quantitative prediction of experimentally-measured inhibition constants, or $K_i$ values (which are ubiquitously used as a practical measure for binding affinities, and compared across different ligands as a relative, realistic biochemical reactivity potential with respect to a specific target protein).[35–40]

Our assumptions, in this context, may be summarized as follows:
 (a) Active-site binding corresponds to a critical event in the inhibition process; that is, inhibition cannot occur in the absence of such event.
 (b) A combination of independent energetic components derived from a sufficiently-accurate description (which includes noncovalent binding factors) of this binding event is characteristic to a given ligand's isomeric structure and chemical composition. That is, a significant change in the latter would result in qualitatively-different such components.
 (c) Individual local-energy-decomposition (LED; see computational methods section) contributions exhibit well-defined intermolecular distance dependence;[41] they therefore incorporate some dynamic information on NCIs taking place in the active site. (indeed, the latter NCIs result, *inter alia*, from the ligand's electronic properties; thus, they may also reflect additional, potential PL NCIs – taking place *before* active site binding.)
 (d) For quality-control purposes, calculated quantities should *not* implicitly include information from molecular structures or events that are (even slightly) orthogonal to active-site binding. (interaction energies, which employ *optimized* structures for each of the interacting monomers in *vacuum*, do include such implicit information – as opposed to the inter-fragment binding energies used below.)

It should be stressed that the very fundamental principles on which our model lie may simply be traced back to *chemical intuition* – as so many predictive tools, incorporating static

molecular structures as a source of information, are still extensively-used by the general chemistry community for the purpose of studying realistic, dynamic molecular systems. It may seem, in fact, that explaining dynamic processes by means of static molecular structures is a general feature that defines chemistry as a scientific discipline. The interested reader may browse through an account of representative such chemical explanations offered in *Appendix A: Static Solutions to Dynamic Problems*.

## *The PL modeling dilemma:*



**OR**

*Appropriate description of noncovalent interactions*

*Dynamic Behavior*

*Which one should be chosen…?*

Figure 2. An illustration of the main challenge hindering predictive modeling of PL systems (see text).

## **Computational Methods**

All geometries used in this work were obtained in the following manner:

1. Nine crystal structures of torpedo californica acetylcholinesterase (Tc AChE), each containing a different bound ligand (a.k.a inhibitor) in its active site, were drawn from the PDB website (see corresponding research papers in Refs.[35–40,42];).

2. "Potentially-active" amino acid residues, defined to be found within 3 Å from any atom in the ligand structure (thus being capable of significantly-interacting with the latter; see, for instance, section 2.2 in Ref.[43]) were selected *via* 'CONTACT' calculations included in the CCP4 suite.[44]

3. Residues found in the preceding stage for each crystal structure were then simply taken alongside the corresponding bound ligand to create the final active-site + ligand geometries used throughout this paper.

Electronic structure calculations were then performed exclusively on the resulting structures, and thus correspond to "QM cluster" calculations – according to the taxonomy used in Ref.[28] All calculated data considered in this paper were obtained using DLPNO-CCSD(T) calculations and subsequent LED analyses included in the ORCA 4.2 program package;[41,45] "NormalPNO" settings,[46] as well as the def2-SVP basis set,[47] were used for all of the latter. Thus, all data were drawn from LED outputs in the following manner:

- DLPNO-CCSD(T)/SVP inter-fragment binding energies were drawn from the "Sum of INTER-fragment total energies" entry, found in the "INTER- vs INTRA-FRAGMENT TOTAL ENERGIES (Eh)" section in the LED outputs. As a sanity check, we verified that binding energies derived from subtracting the sum of "Intra-fragment total energies" from the "total energy" for a given PL complex (both found in the same section in LED output) produce identical energetic values – as shown in the ESI. Note that different definitions for "binding energies" can be found in the literature (some actually correspond to the "interaction energies" mentioned in the introduction); in our case, the term simply corresponds to the difference in total energies between the super-system and its underlying protein and ligand fragments [which satisfies assumption (d) in the introduction].

- Energetic contributions corresponding to LED components arising from electrostatics, exchange and dispersion were extracted from the "FINAL SUMMARY DLPNO-CCSD ENERGY DECOMPOSITION (Eh)" section in the LED outputs. Charge transfer contributions were drawn from the preceding "DECOMPOSITION OF CCSD STRONG PAIRS INTO DOUBLE EXCITATION TYPES (Eh)". Note that for our purposes [see assumption (c) in the introduction], we were interested in grouping different energetic contributions according to their intermolecular distance dependence; thus, we chose to consider the *sum* of "Charge Transfer 1 to 2" and "Charge Transfer 2 to 1" as the total charge transfer contribution to the binding energy (denoted by $E_{ct}$). Similarly, our account for dispersion corresponds to the sum of the "Dispersion (strong pairs)" and "Dispersion (weak pairs)" contributions found in the LED output.

Note that whereas the nonempirical DLPNO-CCSD(T) method and LED approach are used for generating the data considered in this paper – other methods (such as those based on a

perturbation theory formalism) may generally be used for similar purposes.[48–51] It should also be mentioned that the above basis set and PNO domains may rightfully be considered inadequate for quantitatively-accurate electronic structure calculations (resulting in energetics found within 1 kcal/mol from a reliable reference level) of noncovalent interactions in *vacuum*.[46] That being said, it should be stressed that accurate calculation of NCI energetics should *not* be recognized as one of the main goals of the current paper. Instead, we will focus on using the very basic *information* derived from LED calculations for predictive purposes. The *semi-quantitative* manifestation of this purpose, as we shall show below, is independent of such extreme quantitative accuracy.

Multiple linear regression was carried out using the "Analysis Toolpak" add-in for Microsoft Excel 2018 (Macintosh version); 95% confidence intervals were consistently employed for all resulting models. For reproduction purposes, all relevant geometries and ORCA input files used for this paper are provided in the ESI, alongside a dedicated spreadsheet containing our raw and calculated data.

We would like to suggest that our methodological choices and considerations may be of particular interest on their own (and not just as means for achieving the main, stated goal of this paper); the interested reader may browse through associated methodological discussions, questions and answers – all provided in *Appendix B: Methodological Meditations*.


## Results and Discussion

Experimentally-measured inhibition constants [expressed as $\log(K_i)$], and the calculated data for all ligands under consideration (i.e., inter-fragment binding energies and corresponding LED contributions, all given in kcal/mol), are provided in Table 1.

Table.1 Experimentally-measured log($K_i$) values, and calculated data for eight Tc AChE ligands – obtained at the levels of theory specified in the methods section. Similar data calculated for yet another ligand (E20), for which no experimental $K_i$ value was measured, are given in the last row. All energetic components are in kcal/mol.

| PDB ID | Ligand | Log($K_i$) (Refs.[35–40]) | Binding Energy | $E_{elstat}$ | $E_{exch}$ | $E_{ct}$ | $E_{disp}$ |
|--------|--------|---------------------------|----------------|--------------|------------|----------|------------|
| 3ZV7 | NHG | 3.079 | 81.767 | 50.900 | 13.564 | 8.025 | 14.784 |
| 1W6R | GNT | 2.279 | 98.819 | 52.101 | 20.069 | 9.356 | 23.725 |
| 5NAU | DZ0 | 1.475 | 49.015 | 28.012 | 8.192 | 4.639 | 11.611 |
| 1U65 | CP0 | 1.415 | 201.247 | 112.476 | 36.497 | 17.648 | 46.225 |
| 5NAP | DZ7 | 1.046 | 19.606 | 12.764 | 2.953 | 2.989 | 3.577 |
| 1H23 | E12 | 0.653 | 220.404 | 138.560 | 34.812 | 19.335 | 41.187 |
| 1H22 | E10 | -0.097 | 243.558 | 154.083 | 38.869 | 23.252 | 44.367 |
| 1E66 | HUX | -0.886 | 1439.501 | 1250.114 | 117.549 | 60.216 | 42.091 |
| 1EVE | E20 | - | 85.203 | 50.170 | 14.598 | 8.658 | 18.184 |

Clearly, drawing predictive conclusions from this data using nothing but the naked eye is no easy task – as there is no clear one-to-one correspondence between any of the calculated variables and the experimentally-measured quantities. Let us therefore resort to an MLR-based picture – from which, as will be shown below, some interesting conclusions may be drawn.

For illustration purposes, a plot of experimental $K_i$ values is given in Figure 3.A; the predictive value provided by a simple, MLR-based model (to be denoted M1 from this point onwards) – employing calculated binding energies as a single predictive variable – is accordingly demonstrated in Figure 3.B. Both residual errors and regression statistics (Table 2) testify that binding energies simply do not possess enough information for reproducing the general trend created by experimentally-measured $K_i$ values – as M1 recovers only 50% of the sum of squares (SSQ) for the latter. In other words, the variation in $K_i$ values is not trivially explained by means of the corresponding binding energies. In addition, residual errors as large as ~0.5 log($K_i$) – having clear implications on the model's predictive value – can be observed for five of the predicted inhibition constants. These findings clearly fit our expectations regarding the possibility of reducing binding affinities to calculated binding energies – as pointed out in the introduction (see also Figure 1).
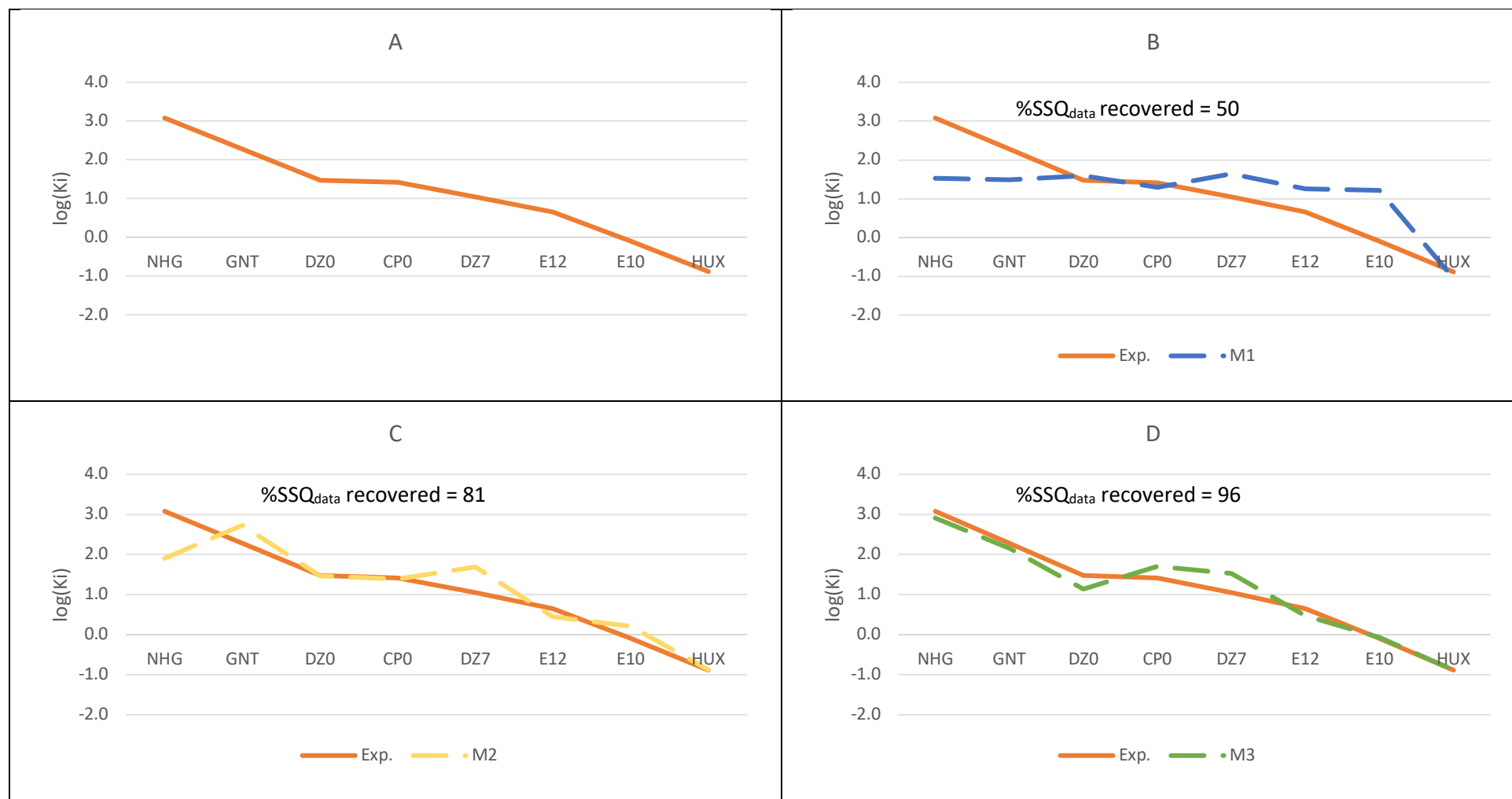
Figure 3. (A) Experimentally-measured $K_i$ values (nM) for eight different ligands, taken from Refs.[35–40]. (B) Multiple-linear-regression model [M1] based on inter-fragment binding energies calculated for the above ligands and the corresponding active amino-acid residues in the Tc AChE active site (C) A similar model [M2] based on specific noncovalent interactions calculated for the same systems using the LED approach (D) Our best such model, employing both calculated binding energies and specific noncovalent interactions used in (A-B). Clearly, M3 is the most robust model considered – as indicated by statistical parameters (see Table 2) as well as by the similarity between the resulting predicted curve and that of (A).

A similar MLR-based model (M2), based solely on calculated LED components, clearly represents a substantial improvement: it recovers 81% of the SSQ for the experimentally-measured $K_i$ values (Table 2; Figure 3.C). Additionally, residual errors of qualitative significance are fewer in number compared to M1 – and correspond to only three predictions. The distribution of errors is generally narrower than that of M1 (as also indicated by $SSQ_{residue}$ for each of the models). Such improvements suggest that information representing *particular* NCIs taking place in ligand binding may be used to better predict experimental results – that is, compared to information exclusively drawn from binding energies. That being said, such outcome may partly be attributed to the fact that more informative variables are fitted to approximate the experimental $K_i$ curve (the fitting process, however, cannot *exclusively* be held responsible for our models' predictive capabilities, as demonstrated in *Appendix B*).

Table 2. Sum of squares of the data and the residual errors for models M1-3 (left). Particular residual errors (or eM[n], n=1-3) for the corresponding predicted log(Ki) values are also provided (right). The heatmap should be read as follows: red marks positive errors, blue corresponds to negative ones. The latter's relative magnitude is indicated by color intensity.

| | M1 | M2 | M3 | PDB ID | Ligand | eM1 | eM2 | eM3 |
|---|---|---|---|---|---|---|---|---|
| $N_{parameters}$ | 1 | 4 | 5 | 3ZV7 | NHG | 1.555 | 1.172 | 0.169 |
| $N_{data}$ | | 8 | | 1W6R | GNT | 0.787 | -0.450 | 0.118 |
| $SSQ_{data}$ | | 1.11E+01 | | 5NAU | DZ0 | -0.111 | 0.013 | 0.332 |
| $SSQ_{Residue}$ | 5.54E+00 | 2.13E+00 | 4.94E-01 | 1U65 | CP0 | 0.118 | 0.022 | -0.289 |
| %Residue | 49.8% | 19.1% | 4.4% | 5NAP | DZ7 | -0.596 | -0.647 | -0.476 |
| %$SSQ_{data}$ recovered | 50% | 81% | 96% | 1H23 | E12 | -0.608 | 0.202 | 0.175 |
| | | | | 1H22 | E10 | -1.314 | -0.302 | -0.025 |
| | | | | 1E66 | HUX | 0.168 | -0.010 | -0.005 |

Let us now consider a third model (M3) – incorporating *both* binding energies and LED data employed in M1 and M2, respectively. As shown in Table 2 and Figure 3.D, this model clearly has striking predictive power – as it recovers no-less-than 96% of the SSQ for the experimentally-measured $K_i$ values. Residual errors are also much smaller, and their distribution narrower, compared to M1-2: the single-largest error amounts to 0.476 log($K_i$), which almost matches the very *smallest*, qualitatively-significant errors found for the previous models. What this all means is that the totality of information corresponding to *both* overall binding strength and specific NCIs for each of the ligands may be used for reproducing the experimentally-measured $K_i$ curve in a semi-quantitative manner.

Indeed, some of the particular binding affinities predicted by M3 are, to some extent, overestimated/underestimated (note DZ0/CP0/DZ7). Nevertheless, despite not having *ideal* predictive value – this simple model may clearly provide useful *comparative* estimates, corresponding to realistic inhibition potentials for different ligands.

As a basic demonstration, let us use the data calculated for yet another ligand, E20 (taken from the 1EVE structure, see Ref.[42]) – which is also given in Table 1. Measured $K_i$ values for this ligand are, to the best of our knowledge, unavailable in current literature. We can therefore trivially use M3 for providing an estimate for its inhibition potential, as shown in Figure 4. It turns out that E20 has a predicted $\log(K_i)$ value of ~1.3 – which is similar in magnitude to those measured for, e.g., CP0. We may accordingly predict that E20 may constitute an effective competitive inhibitor for ligands such as DZ7 and E12 [having $\log(K_i)$ values of 1.046 and 0.653, respectively], and an ineffective one for NHG (3.079). Thus, despite its simplistic architecture, M3 may be seen as valuable predictive tool – as it may guide experimentalists in *practical attempts* for using different ligands for various biochemical purposes.
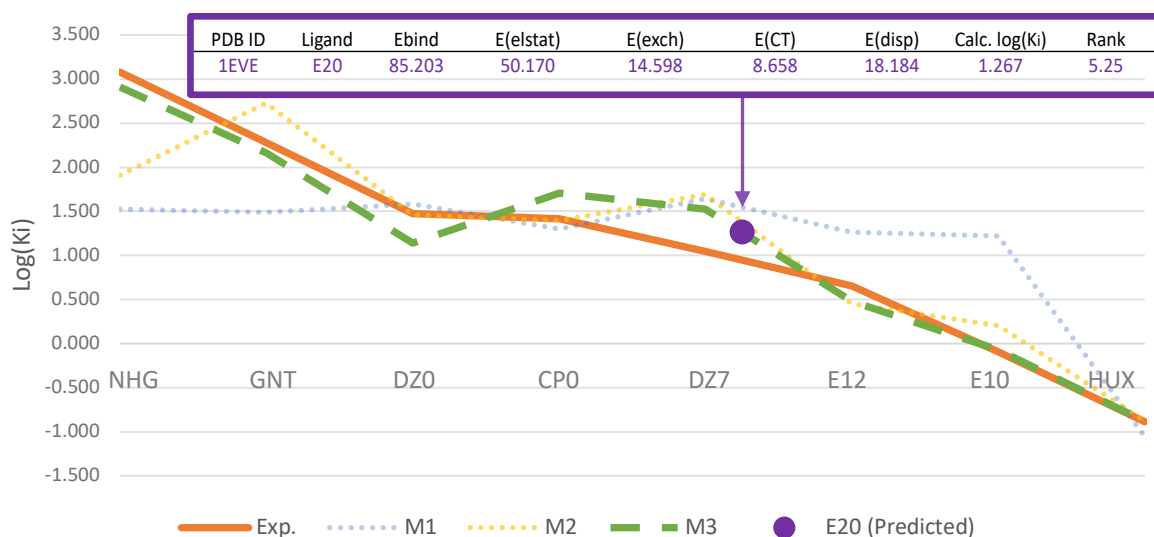


| PDB ID | Ligand | Ebind | E(elstat) | E(exch) | E(CT) | E(disp) | Calc. log($K_i$) | Rank |
|--------|--------|--------|-----------|---------|--------|---------|-----------------|------|
| 1EVE | E20 | 85.203 | 50.170 | 14.598 | 8.658 | 18.184 | 1.267 | 5.25 |

Figure 4. Predicted binding affinity ($K_i$, nM) for the E20 ligand, according to M3 (see Table 2). All calculated energetic contributions are given in kcal/mol. Based on the latter model, E20 can be expected to have the 5th-highest affinity among all nine ligands considered in our work.

Obviously, one should not expect such simple application of our approach to be appropriate for all possible PL systems (some particular cases might prove to be problematic, as demonstrated in *Appendix B, Q3*); still, more elaborate applications – incorporating additional

sources of information – may also be considered if need be, and we hope to explore such options in future projects. Still, and as mentioned in the introduction, it is rather striking that a simple MLR-based model, incorporating information from *static* molecular structures, may be used to explain complex biochemical phenomena – often said to have infinite degrees of freedom. In this context, a few words should be said about the scientific *knowledge* gained by the above results should be added. For the sake of the current discussion, let us follow the classic text by Sanders[52] – which presented knowledge as resulting from the purposeful use of information in an appropriate, well-defined context. Considering the above discussion, a take-home message can be summarized as follows: *static* quantum molecular information may, in principle, be used to provide predictive explanations for *dynamic* protein-ligand processes. This statement clearly has substantial implications on contemporary chemistry knowledge – and we hope it will be of service in future scientific efforts concerning such interactions. As mentioned in the introduction, the general idea which underlies our current approach (illustrated in Figure 5) is, by no means, new. Quite a few great chemists have attempted to conduct similar arguments (see *Appendix A*), but they seem to have lacked the appropriate technical means needed for establishing solid, data-based conclusions. Luckily, state-of-the-art quantum chemistry methods have provided us with the required missing pieces of information – and thus bringing the corresponding puzzle of scientific knowledge to its completion has finally become possible.

It is, perhaps, an appropriate time to remind the reader that we do *not*, on a general principle, recommend using simplistic MLR-based models such as above for practical investigations of PL systems. What we *did* mean to demonstrate is that some particular static molecular structures, described using nonempirical electronic structure methods, clearly possess enough information for modeling biochemically-significant PL interactions. We can only hope that this basic insight will, eventually, be combined with more elaborate and robust modeling techniques. As a side note, we would like to mention that our above results, methodological considerations and assumptions may be of interest for several additional reasons (which had not been discussed in preceding sections): [a] physical meaning of LED contributions is different than that of "realistic" NCIs – which do not necessarily exhibit a well-defined dependence on the intermolecular distance; in addition, the relationship between such calculated components and the total binding energy is nontrivial; [b] using and validating MLR models for confirming the very *informativeness* of predictive variables is a fundamentally different task than establishing statistically-robust models for practical

applications – although the two may easily be confused. Thus, we hereby encourage the reader to browse through *Appendix B*, where such matters are discussed in appropriate length.

<span style="color:purple">**Our Energy-Decomposition-Analysis-Based Approach:**</span>

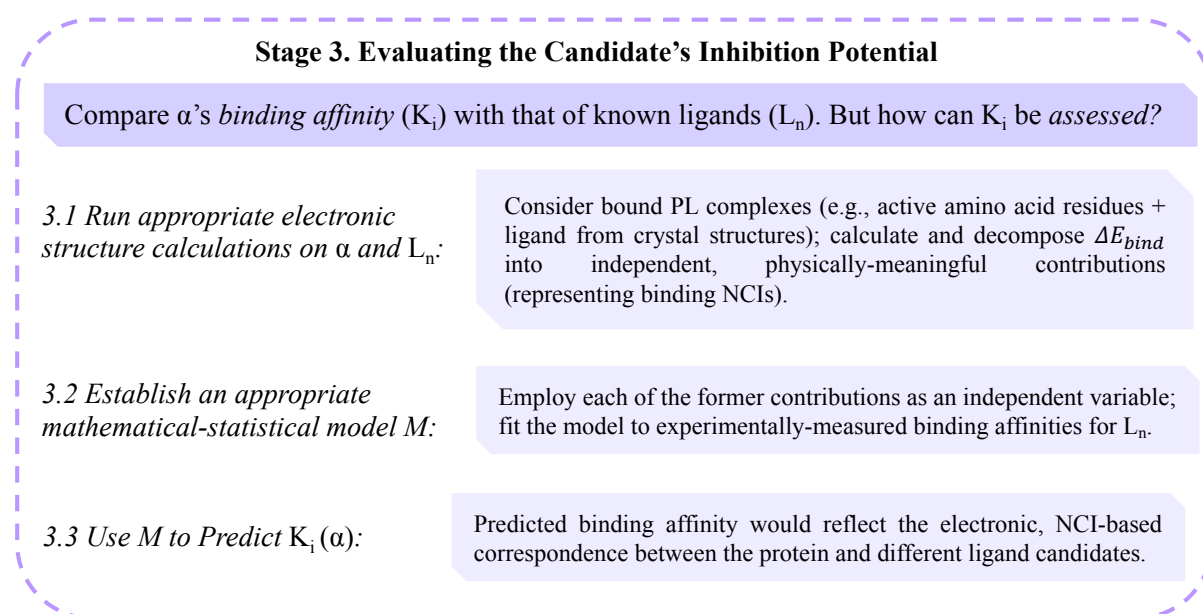- Follow original starting point and Stages 1-2 (Figure 1)

**Stage 3. Evaluating the Candidate's Inhibition Potential**

Compare α's *binding affinity* ($K_i$) with that of known ligands ($L_n$). But how can $K_i$ be *assessed?*

*3.1 Run appropriate electronic structure calculations on* α *and* $L_n$*:*

Consider bound PL complexes (e.g., active amino acid residues + ligand from crystal structures); calculate and decompose $\Delta E_{bind}$ into independent, physically-meaningful contributions (representing binding NCIs).

*3.2 Establish an appropriate mathematical-statistical model M:*

Employ each of the former contributions as an independent variable; fit the model to experimentally-measured binding affinities for $L_n$.

*3.3 Use M to Predict* $K_i$ (α)*:*

Predicted binding affinity would reflect the electronic, NCI-based correspondence between the protein and different ligand candidates.

Figure 5. An "alternate ending" to the process presented in Figure 1, making use of our own energy-decomposition-analysis-based approach as outlined above (Acronyms: PL = protein-ligand, NCI = noncovalent interactions).

## **Summary and Conclusions**

Based on our above investigation of the Tc AChE enzyme and associated ligands/inhibitors, the following conclusions may concisely be summarized:

- We have seen that informative, *static* molecular structures – corresponding to bound protein-ligand complexes – *can* be used for a semi-quantitative prediction of the corresponding, experimentally-measured $K_i$ values. Such successful prediction is by no means trivial, due to the fact that binding affinities are assumed to result from a large variety of *dynamic* factors affecting the *continuous* PL binding process.

- Multiple-linear-regression-based models incorporating *either* inter-fragment binding energies *or* LED components calculated for the bound PL structures do *not* possess desirable predictive power – as they cover only 50% and 81% of the sum of squares

for the experimental $K_i$ values, respectively. In addition, large residual errors (having clear qualitative significance) are observed for both models.

- In contrast, a model employing *both* binding energies and LED components *does* offer surprisingly-useful predictive capabilities, covering no less than 96% of the sum of squares for the experimentally-measured values. Additionally, the residual error distribution is significantly narrower than those of the above alternatives; the *largest* such error, in this case, amounts to the *smallest* significant ones found for former models.

- The above model may easily be used for predicting $K_i$ values which have not yet been experimentally-measured. In such manner, the realistic inhibition potential for the E20 ligand was assessed to be greater than those of DZ7 and E12, but smaller than that of NHG.

- The statistical significance of calculated binding energies and LED components cannot be exclusively attributed to the number of independent parameters and corresponding fitting coefficients used in each model (*Appendix B, Q2*). Thus, our calculated data clearly has *inherent* predictive value.

- Despite the fact that LED components do *not* represent physically-realistic noncovalent interactions (arising from subtle, *dynamic* electronic effects), they do incorporate highly-valuable information on the latter (*Appendix B, Q1*). Such information may be combined with additional data (in our case, calculated binding energies) for the purpose of predicting realistic chemical quantities.

Our above conclusions may also be intuitively-understood using the classic "lock-key" analogy[53] for PL binding: overall active-site binding energetics may be considered to provide some information on a given keyhole's "size", while PL complex-specific NCIs (represented by specific LED contributions) incorporate information on its corresponding "shape". Whereas an entire lock's mechanism cannot simply be *inferred* from its keyhole's properties – focusing on the latter may often suffice for practical predictive purposes.
(Interestingly enough, the "lock-key" analogy had later been "replaced" by an alternative, "glove-hand" version[54] – which seemed to demonstrate some important aspects of PL binding processes which were "ignored" by its predecessor. That, however, clearly entailed simplifying/neglecting other such aspects. In any case, analogies of this sort are merely used for facilitating intuitive understanding and should not be taken too literally.)

As a final remark, we would like to express our hopes and great anticipation for additional efforts concentrated on supplying predictive scientific explanations based on chemical intuition. The latter, which may be seen as one of the most prominent achievements of modern science, has apparently not been fully utilized by means of currently-available scientific methods and techniques. We can only hope to contribute to such efforts and witness the science of chemistry as it reaches new and exciting heights.

## **Acknowledgments**

## **Supporting information**

All geometries and ORCA input files used in this work are provided in the ESI, alongside a dedicated spreadsheet containing our raw and calculated data.

## References

[1] G.U. Nienhaus, *Protein-Ligand Interactions* (Humana Press, New Jersey, 2005).

[2] H. Gohlke, editor , *Protein-Ligand Interactions* (Wiley-VCH Verlag GmbH & Co. KGaA, Weinheim, Germany, 2012).

[3] M.A. Williams, in *Methods Mol. Biol.* (Humana Press Inc., 2013), pp. 3–34.

[4] X. Du, Y. Li, Y.-L. Xia, S.-M. Ai, J. Liang, P. Sang, X.-L. Ji, and S.-Q. Liu, Int. J. Mol. Sci. **17**, 144 (2016).

[5] G. Klebe, Drug Discov. Today **11**, 580 (2006).

[6] A.R. Leach, B.K. Shoichet, and C.E. Peishoff, J. Med. Chem. **49**, 5851 (2006).

[7] M. Congreve, G. Chessari, D. Tisi, and A.J. Woodhead, J. Med. Chem. **51**, 3661 (2008).

[8] D.J. Livingstone, *Drug Design Strategies* (Royal Society of Chemistry, Cambridge, 2011).

[9] T.J.A. Ewing, S. Makino, A.G. Skillman, and I.D. Kuntz, J. Comput. Aided. Mol. Des. **15**, 411 (2001).

[10] H. Gohlke, M. Hendlich, and G. Klebe, J. Mol. Biol. **295**, 337 (2000).

[11] M.K. Gilson and H.-X. Zhou, Annu. Rev. Biophys. Biomol. Struct. **36**, 21 (2007).

[12] S.F. Sousa, A.J.M. Ribeiro, J.T.S. Coimbra, R.P.P. Neves, S.A. Martins, N.S.H.N. Moorthy, P.A. Fernandes, and M.J. Ramos, Curr. Med. Chem. **20**, 2296 (2013).

[13] M. Ragoza, J. Hochuli, E. Idrobo, J. Sunseri, and D.R. Koes, J. Chem. Inf. Model. **57**, 942 (2017).

[14] E.A. Meyer, R.K. Castellano, and F. Diederich, Angew. Chemie Int. Ed. **42**, 1210 (2003).

[15] D.H. Williams, E. Stephens, D.P. O'Brien, and M. Zhou, Angew. Chemie Int. Ed. **43**, 6596 (2004).

[16] H.-J. Schneider, Angew. Chemie Int. Ed. **48**, 3924 (2009).

[17] L.M. Salonen, M. Ellermann, and F. Diederich, Angew. Chemie Int. Ed. **50**, 4808 (2011).

[18] A.S. Mahadevi and G.N. Sastry, Chem. Rev. **113**, 2100 (2013).

[19] P. Politzer, J.S. Murray, and T. Clark, Phys. Chem. Chem. Phys. **15**, 11178 (2013).

[20] J. Řezáč and P. Hobza, Chem. Rev. **116**, 5038 (2016).

[21] P. Hobza, Acc. Chem. Res. **45**, 663 (2012).

[22] L.A. Burns, M.S. Marshall, and C.D. Sherrill, J. Chem. Theory Comput. **10**, 49 (2014).

[23] M.J. Gillan, D. Alfè, P.J. Bygrave, C.R. Taylor, and F.R. Manby, J. Chem. Phys. **139**, 114101 (2013).

[24] M.S. Gordon, D.G. Fedorov, S.R. Pruitt, and L. V. Slipchenko, Chem. Rev. **112**, 632 (2012).

[25] J. Liu, T. Zhu, X. Wang, X. He, and J.Z.H. Zhang, J. Chem. Theory Comput. **11**, 5897 (2015).

[26] Y. Gao, X. Lu, L.L. Duan, J.Z.H. Zhang, and Y. Mei, J. Phys. Chem. B **116**, 549 (2012).

[27] C. Ji and Y. Mei, Acc. Chem. Res. **47**, 2795 (2014).

[28] U. Ryde and P. Söderhjelm, Chem. Rev. **116**, 5520 (2016).

[29] K. Roos, J. Viklund, J. Meuller, K. Kaspersson, and M. Svensson, J. Chem. Inf. Model. **54**, 818 (2014).

[30] P. Saparpakorn, M. Kobayashi, S. Hannongbua, and H. Nakai, Int. J. Quantum Chem. **113**, 510 (2013).

[31] J. Heimdal and U. Ryde, Phys. Chem. Chem. Phys. **14**, 12592 (2012).

[32] G. König, P.S. Hudson, S. Boresch, and H.L. Woodcock, J. Chem. Theory Comput. **10**, 1406 (2014).

[33] C.J. Woods, K.E. Shaw, and A.J. Mulholland, J. Phys. Chem. B **119**, 997 (2015).

[34] M.A. Olsson, P. Söderhjelm, and U. Ryde, J. Comput. Chem. **37**, 1589 (2016).

[35] C. Bartolucci, J. Stojan, Q. Yu, N.H. Greig, and D. Lamba, Biochem. J. **444**, 269 (2012).

[36] H.M. Greenblatt, C. Guillou, D. Guénard, A. Argaman, S. Botti, B. Badet, C. Thal, I. Silman, and J.L. Sussman, J. Am. Chem. Soc. **126**, 15405 (2004).

[37] R. Caliandro, A. Pesaresi, L. Cariati, A. Procopio, M. Oliverio, and D. Lamba, J. Enzyme Inhib. Med. Chem. **33**, 794 (2018).

[38] M. Harel, J.L. Hyatt, B. Brumshtein, C.L. Morton, K.J.P. Yoon, R.M. Wadkins, I. Silman, J.L. Sussman, and P.M. Potter, Mol. Pharmacol. **67**, 1874 (2005).

[39] D.M. Wong, H.M. Greenblatt, H. Dvir, P.R. Carlier, Y.-F. Han, Y.-P. Pang, I. Silman, and J.L. Sussman, J. Am. Chem. Soc. **125**, 363 (2003).

[40] H. Dvir, D.M. Wong, M. Harel, X. Barril, M. Orozco, F.J. Luque, D. Muñoz-Torrero, P. Camps, T.L. Rosenberry, I. Silman, and J.L. Sussman, Biochemistry **41**, 2970 (2002).

[41] W.B. Schneider, G. Bistoni, M. Sparta, M. Saitow, C. Riplinger, A.A. Auer, and F. Neese, J. Chem. Theory Comput. **12**, 4778 (2016).

[42] G. Kryger, I. Silman, and J.L. Sussman, Structure **7**, 297 (1999).

[43] H.F. Lodish, A. Berk, S.L. Zipursky, P. Matsudaira, D. Baltimore, and J. Darnell, *Molecular Cell Biology*, 4th editio (W.H. Freeman, New York, 2001).

[44] M.D. Winn, C.C. Ballard, K.D. Cowtan, E.J. Dodson, P. Emsley, P.R. Evans, R.M. Keegan, E.B. Krissinel, A.G.W. Leslie, A. McCoy, S.J. McNicholas, G.N. Murshudov, N.S. Pannu, E.A. Potterton, H.R. Powell, R.J. Read, A. Vagin, and K.S. Wilson, Acta Crystallogr. Sect. D Biol. Crystallogr. **67**, 235 (2011).

[45] M. Schütz, J. Chem. Phys. **113**, 9986 (2000).

[46] D.G. Liakos, M. Sparta, M.K. Kesharwani, J.M.L. Martin, and F. Neese, J. Chem. Theory Comput. **11**, 1525 (2015).

[47] F. Weigend and R. Ahlrichs, Phys. Chem. Chem. Phys. **7**, 3297 (2005).

[48] B. Jeziorski, R. Moszynski, and K. Szalewicz, Chem. Rev. **94**, 1887 (1994).

[49] H.L. Williams and C.F. Chabalowski, J. Phys. Chem. A **105**, 646 (2001).

[50] A.J. Misquitta and K. Szalewicz, Chem. Phys. Lett. **357**, 301 (2002).

[51] T.M. Parker, L.A. Burns, R.M. Parrish, A.G. Ryno, and C.D. Sherrill, J. Chem. Phys. **140**, 094106 (2014).

[52] J. Sanders, in *2016 SAI Comput. Conf.* (IEEE, 2016), pp. 223–228.

[53] E. Fischer, Berichte Der Dtsch. Chem. Gesellschaft **27**, 2985 (1894).

[54] D.E. Koshland, Angew. Chemie Int. Ed. English **33**, 2375 (1995).

[55] P. Pechukas, in *Dyn. Mol. Collisions* (Springer US, Boston, MA, 1976), pp. 269–322.

[56] B.G. Miller and R. Wolfenden, Annu. Rev. Biochem. **71**, 847 (2002).

[57] H. Eyring, J. Chem. Phys. **3**, 107 (1935).

[58] E. Runge and E.K.U. Gross, Phys. Rev. Lett. **52**, 997 (1984).

[59] P.K. Chattaraj and S. Sengupta, J. Phys. Chem. **100**, 16126 (1996).

[60] L. Pauling, Am. Sci. **36**, 50 (1948).

[61] E.A. Carter, G. Ciccotti, J.T. Hynes, and R. Kapral, Chem. Phys. Lett. **156**, 472 (1989).

[62] W. Swegat, J. Schlitter, P. Krüger, and A. Wollmer, Biophys. J. **84**, 1493 (2003).

[63] K.N. Houk, A.G. Leach, S.P. Kim, and X. Zhang, Angew. Chemie Int. Ed. **42**, 4872 (2003).

[64] P. Tiwary, V. Limongelli, M. Salvalaglio, and M. Parrinello, Proc. Natl. Acad. Sci. **112**, E386 (2015).

[65] T.J. Giese and D.M. York, Int. J. Quantum Chem. **98**, 388 (2004).

[66] J. Cioslowski and P.R. Surján, J. Mol. Struct. THEOCHEM **255**, 9 (1992).

[67] R.D. Retherford and M.K. Choe, *Statistical Models for Causal Analysis* (John Wiley & Sons, Inc., Hoboken, NJ, USA, 1993).

[68] B. Efron and T. Hastie, *Computer Age Statistical Inference* (Cambridge University Press, 2016).

[69] P.G. De Benedetti and F. Fanelli, Drug Discov. Today **15**, 859 (2010).

*Appendix A: Static Solutions to Dynamic Problems*

The current paper presents an attempt for drawing information on dynamic PL binding processes from informative, "static pictures". The latter are assumed to represent critical events, which *must* occur as part of biochemically-significant PL binding (see introduction). It might be important to try and understand why one should, in fact, even expect such "static pictures" to be adequate for describing any dynamic process in the first place.

A realistic chemical sample contains an immense ensemble of molecules, found in various conformations and energetic states. Thus, drawing predictive conclusions from discrete, static molecular structures seems physically-unintuitive. That being said, gathering information describing *all* molecules within the ensemble would be unrealistic. However, what if we could find a way of reducing the behavior of this complicated, realistic ensemble to a few "building blocks", or crucial factors (or pieces of information), represented by a few *specific* static molecular structures? Based on some intuitive assumptions, it should be possible to do so and – in that sense – make large-scale, complicated dynamic phenomena predictively-understandable by relatively simple means.[55,56]

Classical Eyring transition state theory (TST) has, in fact, provided chemists with this very possibility (i.e., reducing dynamic problems to static ones).[57] It allowed, for the first time, the prediction of reaction rates – which reflect dynamic information on "collisions" between molecular species – by means of reduction to static, well-defined molecular structures, corresponding to *critical events* along the reaction path (i.e., the forming of a reactive transition state). Indeed, the theory has been shown to offer great predictive power, despite lacking strict physical rigor in its original form (a fact which did not seem to particularly concern Eyring himself; later, more physically-rigorous versions of the theory are reviewed in Ref.[55]).

The truth is, however, that chemists have been used to solving problems this way long before Eyring has published his theory: it may be argued that predicting chemical properties based on, e.g., relevant Lewis/Kekulé structures is founded on similar logic. Obviously, Kekulé structures of organic molecules cannot be expected to represent how the latter "*look*" in a dynamic environment: the electron density for a given molecule clearly undergoes many changes with respect to time, to the point where it is not even trivial to speak of a well-defined molecular structure that is maintained throughout dynamic processes (see, for instance, discussions in Refs.[58,59]). However, *some* static structures *do* interest us because they somehow provide simple and elegant, practical answers for questions of great

importance. Such static structures are therefore of more interest than others, and their identity depends on particular questions we are interested in answering. In the context of chemical reactivity, for instance, energetic *minima* and *maxima* may indeed be of particular interest. Other static structures may allow us to answer questions of different kinds (think, for instance, of solubility; energetic minima/maxima are often less-important than energetically-invariant structural features, such as ones corresponding to an "*ability*" of forming hydrogen bonds). Thus, chemistry may generally be perceived as the science which *reduces dynamic molecular processes to static, informative molecular structures* for the purpose of gaining desirable predictive power (Figure 1A).
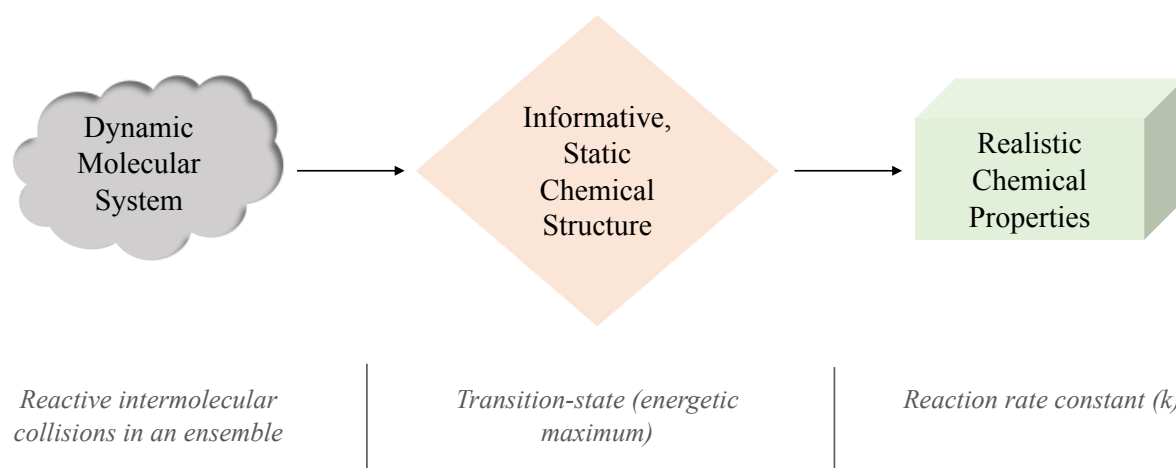


| Dynamic Molecular System | → | Informative, Static Chemical Structure | → | Realistic Chemical Properties |
|---|---|---|---|---|

| *Reactive intermolecular collisions in an ensemble* | *Transition-state (energetic maximum)* | *Reaction rate constant (k)* |
|---|---|---|

Figure A1. A schematic representation of chemical explanations (above); specific examples concerning Eyring's transition-state-theory are also provided (below, in *italics*)

Some thoughts regarding harnessing TST to PL systems have been later expressed by Pauling.[60] Trying to rationalize the catalytic "benefit" of enzymatic reactions, he claimed that the latter increase the relative concentration of reactive transition state species compared to corresponding enzyme-free reactions. Regardless of the practical implications for such claim, the motivation behind it may be understood based on the above discussion of TST: instead of bothering ourselves with a vast amount of information associated with the entire PL complex, we may focus on specific features associated with *a single reactive structure*. Thus, if we wish to quantify a given enzyme's catalytic efficiency, all we need to do is compare the relative concentrations of such species in both enzyme/enzyme-free scenarios. Needless to say, the assumption according to which enzymatic reactions take place *via* single, well-defined reaction paths may indeed be exposed to criticism; however, the same chemical-intuition-driven motivation recognized in the very foundation of TST is again seen to offer

great explanatory value. Its practical utility, however, is somewhat difficult to assess – due to the fact that finding reactive species of interest using current computational techniques is often impractical.

Since Pauling's above claim had emerged, a few related attempts – based on similar argumentation in different practical "packaging" – have also been published.[56,61–64] Still, searching for reactive transition-state structures, telling the "main story" of continuous PL processes, had continued to be the main bottleneck limiting the latter's applicability. That being said, the fact that practical applications of TST to PL systems have met with many difficulties does not mean that a static representation of such systems, having desirable predictive power, cannot be established. Proving that such possibility exists was, in fact, the main motivation for the present paper. Note that our work corresponds, in a sense, to an *inverse* version of TST – since we use *bound* PL complex structures (i.e., energetic *minima*) as a source of information on dynamic processes. Indeed, it turns out that these structures simply provide enough information for answering practical questions of interest (such as ones requiring comparative assessments of binding affinities), and that – from our point of view – plainly stands for yet another useful manifestation of chemical intuition.

*Appendix B: Methodological Meditations*

In order to facilitate the readers of this text, we hereby provide representative answers for important questions that may seem to arise from it. Both questions and answers have resulted, in fact, from actual discussions with colleagues from diverse scientific backgrounds.

Q1. LED vs. Physical Realism

- LED contributions do not represent *realistic* NCIs – which may not exhibit a well-defined intermolecular distance dependence (*true dispersion*, for instance, cannot be said to have the $R^{-6}$ dependence as $E_{disp}$; see Ref.[65] and references therein). Besides, considering that no consistent level of theory was used for obtaining all calculated LED contributions – the latter cannot be summed to reproduce the "original" inter-fragment binding energy values. If that is not enough, geometries drawn from crystal structures do not represent true energetic minima on well-defined potential energy surfaces for the PL complexes at hand. Thus, the data used in M1-3 consist of

approximate quantities that are not rooted in quantum mechanics. How can it then be considered as a reliable source of predictive information?

Indeed, the data used in M1-3 may, perhaps, represent an *approximation* to the "true" NCIs taking place between each ligand and Tc AChE. That being said, and putting aside the fact that simple regression statistics confirm the informativeness of our predictors (see Q2 below), it is important to note that we had no intention of committing our scientific worldview to quantum mechanics alone. To us, electronic structure calculations are simply a tool for establishing general conclusions – aimed to benefit the science of chemistry. Thus, the use of approximate quantities (which were not derived from a fundamental, consistent physical model – embodied in a single, well-defined electronic structure method) does not harm any of the conclusions drawn in this paper (recall that the main question considered here is: is it possible to reduce PL binding to static chemical structures?).

On a different note, it should be emphasized that despite their seemingly-approximate nature, LED contributions *do* incorporate "high-resolution" information on NCIs – as they are derived exclusively from the nonempirical electronic structure calculated for a given system, using a well-defined protocol (in this context, see Cioslowski & Surján's "generalized observability criterion" in Ref.[66]). Thus, and despite the fact they are obtained using different levels of theory, the latter may still be calculated and compared *across different systems* for the purpose of arriving at desirable conclusions – as we have chosen to do in this work.

In this context, what really matters is how the above data are being used, and the purposes for which they are used for. Since we have used total IEs and LED contributions as independent variables in MLR-based models, all we care about is comparing their corresponding values across all PL systems under consideration – as opposed to considering different energetic contributions calculated for a single PL complex. For this reason, the strict physical rigor associated with the LED approach is not of our current concern – contrary to its predictive-informative potential. At last, it should actually be mentioned that the fact that the LED contributions do not fully add up to their corresponding binding energies *does represent a practical advantage*; it allows us to avoid intervariable linear-dependence issues that would prevent us from successfully applying MLR in the first place.

As previously mentioned in the introduction, all geometries considered in this work simply correspond to informative static structures (*hypothetical* energetic minima, perhaps?), that can be used for predicting experimental binding affinities (that is, to offer *useful and practical understanding* of PL systems). Indeed, finding these structures by means of well-

converged, computational geometry optimizations would probably prove to be difficult (for reasons similar to those responsible for the impracticalness of Pauling's idea, as mentioned in *Appendix A*). However, considering the true nature of information and methods used in this paper, it is possible to conclude that *doing so should not even be necessary*: the crystal structures used by us, which result from statistical averaging of many individual complex structures and clearly do not correspond to such well-defined minima – have provided us with the information we are interested in for achieving our goals. Following a similar logic, previously-unknown informative molecular structures may practically be found through, e.g., a combination of semi-converged geometry optimizations and simple chemical intuition. The question whether a specific structure truly corresponds to a *physically-realistic* energetic minimum should not also not be of particular concern – as long as it is capable of providing an informative estimate for a specific PL binding energy, as well as for particular NCIs involved in the binding process.

Q2. Causality vs. "Kitchen-Sink Regression"

- In principle, any data may be fitted and used as predictors for reproducing the experimental $K_i$ curve (Figure 3.A); so how do you know that your calculated energetics contain information that has some *causal* relationship with the experimental observations?

Indeed, we have established that binding energies and LED contributions calculated for a set of PL complexes represent enough information for reproducing experimentally-measured $K_i$ values in a semi-quantitative manner; what this really means is that the latter may be reduced, *in the statistical sense*, to the former predictors. This does *not* imply, however, that the binding affinity for a given ligand is determined *exclusively* by the physical factors embodied in our calculated data. Nevertheless, MLR is often used to *corroborate* intuitive causal relationships by showing that a given variable(s) incorporates desirable predictive power; for a review of statistical inference and its relationship to causality (which have been thoroughly discussed in many textbooks), we encourage the reader to browse through, e.g., chapter 2 in Ref.[67]

In many practical applications, fitted regression coefficients are interpreted to reflect causal relationships. That is, in cases where residual errors are small enough, and there seems to be an *intuitive meaning* for the particular regression coefficient calculated for each predictor

(which depends on the researcher's expectations) – the relationship between considered predictors and response variable is said to be causal. In the case considered here, coming up with some intuitive meaning for calculated regression coefficients would be rather difficult (see dedicated spreadsheet in the ESI). That, however, clearly does not mean that our calculated data is uninformative, or that it does not have casual implications on experimentally-measured binding affinities; the contrary, in fact, may easily be confirmed. In order to quantitatively demonstrate the causal-relevance of our predictors, we compared the predictive-performances of models M1-3 with those of three alternatives (A1-3); the latter have been constructed to incorporate a similar number of independent variables (1,4, and 5, respectively), while the specific values "calculated" for each of the ligands considered are *random numbers* found within the range of corresponding quantities used in M1-3. Thus, the single predictor used in A1 spans randomly-generated values between 20–1440 (corresponding to the smallest and largest calculated binding energies used in M1). Similarly, Models A2-3 were established to include random numbers found within the range of the corresponding LED components used in M2-3. It turns out that for a given n (n=1-3), the predictive power of $A_n$ is *substantially lower* than that of $M_n$: A3, for instance, covers just 21% of the SSQ for experimentally-measured $K_i$ values (compared to 96% for M3). Hence, putting chemical intuition aside – the latter findings may indeed be considered as a fair testimony for the *true informative value* of our calculated energetic properties (which is not related to some statistical-pragmatic notion of causality): it means that the predictive potential of M1-3 cannot be *fully-rationalized* by means of fitting regression coefficients to match a number of observed values.

It should be noted that it *might* just happen that some randomly-generated numbers will prove to reflect predictive information after such fitting process (a fact often incorporated into machine learning techniques[68]) – but that would clearly not teach us anything about the causal relevance of our own calculated data to experimentally-measured $K_i$ values. Needless to say, the DLPNO-CCSD(T)/SVP level of theory used above is, by no means, a random number generator: it offers a well-defined path towards obtaining physically-meaningful energetic quantities. Thus, the fact that the latter provide us with desirable predictive power (which cannot trivially be obtained through optimization of regression coefficients, as done for models A1-3) may also suggest *some* causal relationship between our predictors and the considered PL binding affinities.

Q3. Tc AChE and associated ligands compared to other PL systems

- Some of the reasons for the "push" for dynamic modeling of PL interactions are: (a) when binding of the ligand is dynamically accessible – e.g., the binding pocket is hidden in certain conformations (this may be associated with allosteric regulation and other induced conformational changes); (b) the need for including finite temperature effects (such as relevant entropic contributions). To what extent do you think that your approach takes these two factors into account? Would you therefore expect it to be adequate for any possible PL system, or perhaps just for a particular subclass of PL interactions?

First of all, it is worth mentioning that we are *not* claiming that simple models such as M1-3 can be expected to provide useful predictions for just any PL system; while we do believe that our approach (outlined in both main text *and* appendices) should generally be considered for predictive purposes – we certainly suspect that the particular models considered in this work might be inadequate for some PL cases.

Still, some cases which correspond to reason (a) may also, perhaps, be statistically-reduced to specific NCIs between PL pairs (as binding *specificity* to a particular ligand must incorporate some correspondence of this sort). Similarly, reason (b) would indeed push for dynamic modeling as long as considered effects are expected to qualitatively change the very fundamental PL NCI-based correspondence; if that is not the case – there would be no necessary reason to go beyond the static NCI-based picture. Thus, our approach should, in principle, incorporate the above factors – as long as the PL binding process is dominated by the aforementioned NCI-based correspondence. This would clearly not be true for all possible PL binding processes – as may be illustrated using the following simple thought-experiment.

Consider, for instance, a protein that introduces some "random" component into the continuous binding process:

- Case #0: Ongoing conformational changes, for instance, may "hide/reveal" the binding pocket independently from any relationship with a particular ligand. Such random component, however, should not – in principle – affect any prediction of *relative* binding affinities (as it is not ligand-specific and should affect any binding process).

- o Case #1: Taking an additional step forward, it is possible to think of a protein which introduces some random *and* ligand specific such component. In this case, the bound PL structure may not offer enough information for useful predictions – and additional sources of information would then have to be considered.

We would certainly expect models M1-3 to fail in the latter case. Nevertheless, additional variables (obtained using electronic-structure calculations, or other sources of information) may also be used for arriving at useful predictive capabilities – and so our very general approach should not be suspected for having inherent, system-specific shortcomings.

Q4. Relationship between the current work and QSAR approaches

- Your MLR-based model seems similar to ones used in quantitative-structure-activity-relationships (QSAR) studies. Can you explain how your approach differs from previously-proposed QSAR schemes?

Indeed, it is possible to argue that some QSAR models are based upon a "philosophy" similar to ours; that is, chemical predictions are offered based on calculated structural features that are *not* assumed to remain invariant throughout a molecular process of interest. However, approaches of this sort are, to the best of our knowledge, inherently different than the one used in this work:[69]

1. QSAR descriptors (a.k.a independent variables) are not often derived *exclusively* from electronic structure calculations on a bound complex structure. Thus, relationships between all predictors and response variable are not rooted in a consistent, uniform or intuitively-understandable physical picture.
2. QSAR descriptors are not explicitly chosen to reflect 'critical events' in particular molecular processes. In fact, particular QSAR descriptors are often chosen based on the statistical robustness of the resulting model alone, while the added information corresponding to these predictors is not explicitly discussed.

The fact that fitting techniques such as MLR are used for constructing predictive QSAR schemes may certainly suggest some similarities between the latter and the simple models (M1-3) used in our work; however, our emphasis lies on electronic-structures- and chemical-

intuition-based predictor selection. For all it is worth, it should be possible to say that we have chosen to implement our perception of chemistry in a QSAR envelope. The main idea, however, on drawing predictive information on complex molecular systems from static electronic-structure calculations is independent of such particular implementation and should be considered as the main message of our text.

Q5. Implications on Chemical Bonding

- You attempt to statistically-reduce a continuous, *dynamic* biochemical process to information drawn from selected *static* "electron-density snapshots" of isolated molecular species. It is possible to argue that the very notion of a *chemical bond* is also based on somewhat similar statistical reductions. Should your current effort actually be considered as an attempt for generalizing bonding concepts to PL systems?

Despite not being quite aware of this point in the beginning of our investigations, our above reasoning and inferences do resemble predictive arguments involving notions of chemical bonding. In a sense, chemical bonds do not exist in a dynamic molecular system (as any definition of a "bond" would lie on a static electron-density picture). Still, this concept helps us chemists to arrive at astonishingly-useful practical predictions (of, e.g., reactivity, solubility, and the like – depending on the particular questions considered – as mentioned in *Appendix A*). In the current context, we have chosen to consider ligand "binding" as a "critical (instantaneous, therefore static) event" along a dynamic trajectory, that can be used for predictive purposes; since all we are currently interested in is predicting biochemically-significant responses (i.e., enzyme inhibition), "binding" may abstractly be thought of as a set of static pictures which provide us with sufficient predictive power. Such notion of bonding is rather flexible, and may indeed be useful in a variety of (bio)chemical contexts. We surely hope to explore it further in future projects.

Q6. *Ad hominem*

- "I've been in the business for a long while and never witnessed such manuscript getting published. You've got plenty of verbal discussions and no equations to back

them up (especially in the appendices). I'm not sure what it is that you're doing, but it doesn't feel like science to me"

Scientific truths do not depend on anyone's feeling of comfort. If you are interested in criticizing the above arguments - please do so while sticking to logic.