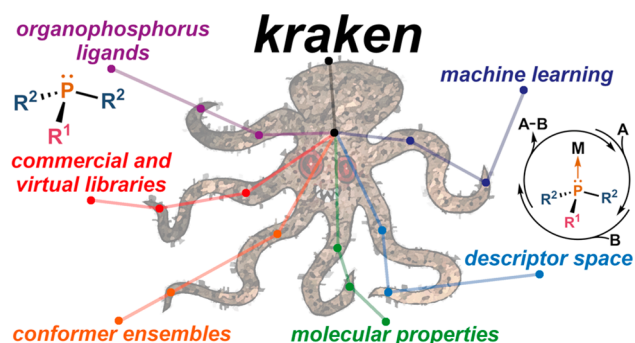


# A Comprehensive Discovery Platform for Organophosphorus Ligands for Catalysis

Tobias Gensch,<sup>\*</sup> Gabriel dos Passos Gomes, Pascal Friederich, Ellyn Peters, Théophile Gaudin, Robert Pollice, Kjell Jorner, AkshatKumar Nigam, Michael Lindner-D'Addario, Matthew S. Sigman,<sup>\*</sup> and Alán Aspuru-Guzik<sup>\*</sup>

**ABSTRACT:** The design of molecular catalysts typically involves reconciling multiple conflicting property requirements, largely relying on human intuition and local structural searches. However, the vast number of potential catalysts requires pruning of the candidate space by efficient property prediction with quantitative structure–property relationships. Data-driven workflows embedded in a library of potential catalysts can be used to build predictive models for catalyst performance and serve as a blueprint for novel catalyst designs. Herein we introduce *kraken*, a discovery platform covering monodentate organophosphorus(III) ligands providing comprehensive physicochemical descriptors based on representative conformer ensembles. Using quantum-mechanical methods, we calculated descriptors for 1558 ligands, including commercially available examples, and trained machine learning models to predict properties of over 300000 new ligands. We demonstrate the application of *kraken* to systematically explore the property space of organophosphorus ligands and how existing data sets in catalysis can be used to accelerate ligand selection during reaction optimization.



## INTRODUCTION

Ligand engineering on the basis of mechanistic hypotheses has been a primary driver of reaction discovery and optimization in catalysis. An emerging and complementary approach applies data-driven methods to molecular design by capturing multidimensional property relationships that directly influence performance.<sup>1–3</sup> The success of such data-driven approaches relies on the availability of powerful molecular representations<sup>4–6</sup> that can be used in a wide range of machine learning (ML) methods.<sup>7–12</sup> Organophosphorous(III) ligands are among the most widely used ligands in homogeneous catalysis. In this study, we establish a comprehensive workflow to study these ubiquitous compounds that can be further extended to other ligand classes. The platform that we developed can be employed for inverse design of novel homogeneous catalysts inspired by past work in the context of both molecular and materials discovery. For example, the Materials Project,<sup>13</sup> OQMD,<sup>14</sup> and AFLOW<sup>15</sup> are tools for exploring the inorganic compound space that include databases, computer scripts for feature extraction, and ML toolkits. Additionally, the Harvard Clean Energy Project<sup>16</sup> has similar goals in the space of organic photovoltaics. Moreover, in the case of heterogeneous catalysis, the Catalysis-Hub<sup>17</sup> contains computed heterogeneous reaction energies with the associated barriers, and the Open Catalyst Project<sup>18</sup> provides density functional theory

(DFT) geometry relaxations for material surfaces with adsorbates. These illustrative examples provide the foundation for how our teams approached the development of a workflow for the chemical space of organophosphorus ligands.

In this context, Tolman introduced experimentally measured descriptors now termed the Tolman electronic parameter (TEP)<sup>19</sup> and the Tolman cone angle<sup>20</sup> to quantify and rationalize phosphorus ligand properties over 50 years ago.<sup>21</sup> These molecular descriptors allowed mapping of the phosphine property space and provided a tool to understand systematic trends in reactivity and stability by using linear free energy relationships and substituent additivity approaches.<sup>22,23</sup> With the emergence of quantum-chemical methods, interest in computed properties of phosphines arose.<sup>24–26</sup> Building on this, the ligand knowledge bases (LKB) developed by Fey and co-workers<sup>27–30</sup> marked an impressive milestone in the mapping of ligand spaces. The LKB-P consists of computed properties of 366 monodentate organophosphorus ligands in

typical coordination environments.<sup>29</sup> The profound impact of ligand conformations on properties<sup>31–34</sup> and catalytic activities<sup>35,36</sup> has been recognized frequently; however, the systematic quantification of ligand flexibility is still underdeveloped. Thus, inspired by the previous approaches to the mapping of ligand space, we aimed at devising a workflow that encompasses a wide range of steric and electronic properties of catalytically relevant ligands including descriptors for their flexibility to enhance the capabilities of data-driven catalyst design.<sup>37</sup>

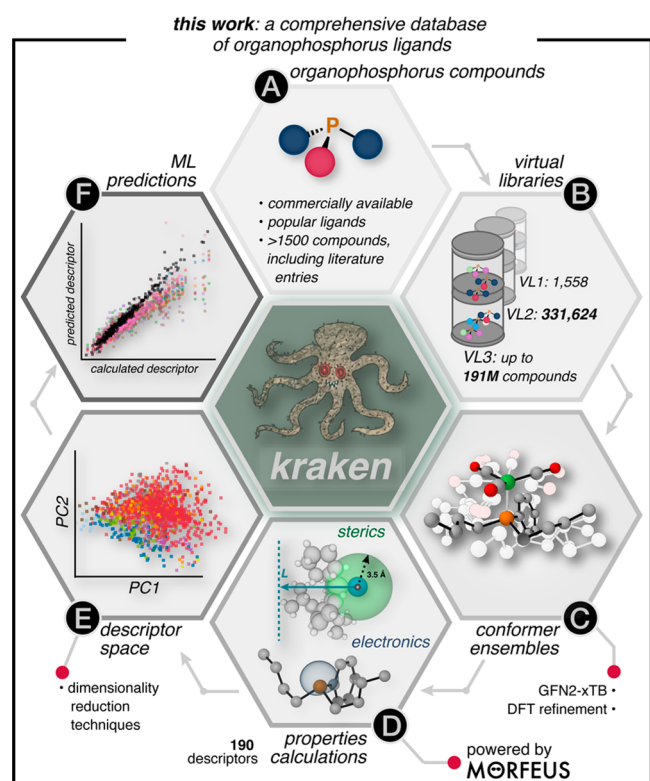
Herein we present *kraken*, an extensive virtual open-access library covering monodentate organophosphorus(III) ligands targeted at facilitating the design and optimization of catalytic processes (Figure 1). To account for conformational flexibility,

a general-purpose physicochemical descriptor set is derived from representative conformer ensembles of both the uncoordinated ligands and a model complex, which we hypothesized provides access to the essential features describing the multitude of intermolecular interactions involved in catalytically relevant steps. Additionally, we demonstrate the application of *kraken* to explore the property space of organophosphorus(III) ligands at a massive scale using increasingly sophisticated models for property estimation of arbitrary organophosphorus(III) compounds. Finally, we showcase the use of *kraken* for inverse catalyst design by constructing multiple linear free energy relationships and other regression models based on experimental data and using it to predict the performance of the entire ligand database, providing the best candidates to be tested in subsequent experiments.

## RESULTS AND DISCUSSION

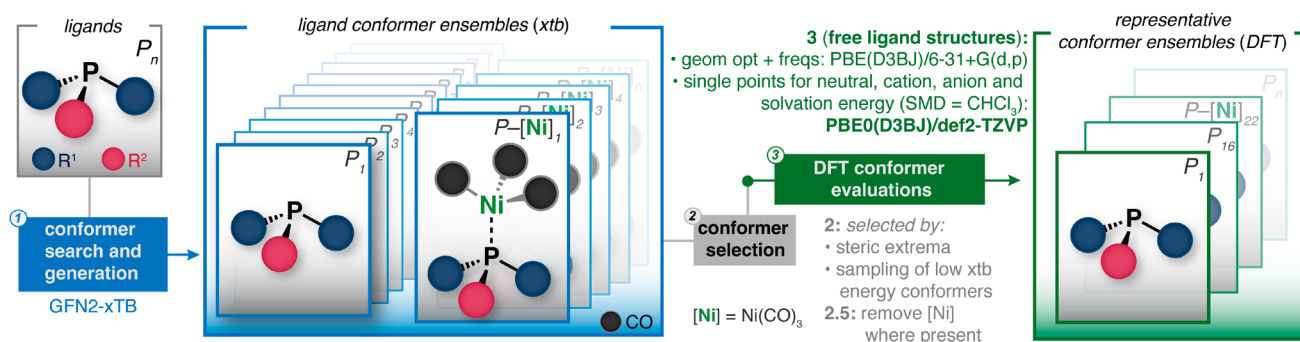
**Library Scope.** A central goal was to comprehensively map the chemical space of monodentate organophosphorus(III) ligands, focusing in particular on structures relevant to applications in catalysis and its use for data-driven ligand optimization campaigns. We initially selected phosphines that were commercially available and prevalent in the organo-(transition)metal chemistry literature. In anticipation of the ML property prediction goals, we surveyed the scientific literature and systematically added ligands with less prevalent substituents based on the core structures found. This was followed by a curation step to avoid structures with additional N-, P-, or S-containing donor sites or acidic moieties in an orientation that might result in additional coordination modes to a metal (e.g., bidentate ligands). Overall, the library contains ligands with various phosphorus–element bonds encompassing H, B, C, N, O, F, Si, and S next to phosphorus in arbitrary combinations. Thus, besides phosphines, other important ligand classes such as phosphoramidites, phosphites, and phosphinamines are also included. In its current state, the library (cf. Figure 1) is constructed on full DFT calculations for 1558 compounds and their conformers, at least 400 of which are commercially available, and it includes the 200 most-cited phosphorus ligands in the literature (virtual library level 1 = VL1). The library also includes 194 compounds with phosphorus in a cyclic structure with the remaining noncyclic structures containing a total of 576 unique substituents. Several representative structures can be found in Figures 4E, 8, and 9, and the full library can be accessed via an interactive web application at <https://kraken.cs.toronto.edu>.

**Conformer Ensembles.** One key challenge when defining the *kraken* computational workflow is the representation of the conformational space of each ligand, the conformer energies, and the corresponding contribution to the ligand properties. This is particularly relevant for steric properties that vary significantly with conformation, whereas electronic properties are generally less sensitive.<sup>38</sup> While no individual model system (i.e., free ligand or specific reference complexes) can fully reflect the conformational space accessible to a ligand in any given complex, there are certain limits for attainable geometries and properties. Importantly, investigating these ranges and limits was used to probe the behavior of ligands in catalytic systems and predict their catalytic performance. For example, the buried volume, i.e., the fraction of the volume of a sphere, which is placed at the metal center, occupied by ligand atoms,<sup>39</sup> of a trialkylphosphine could be very large if all chains



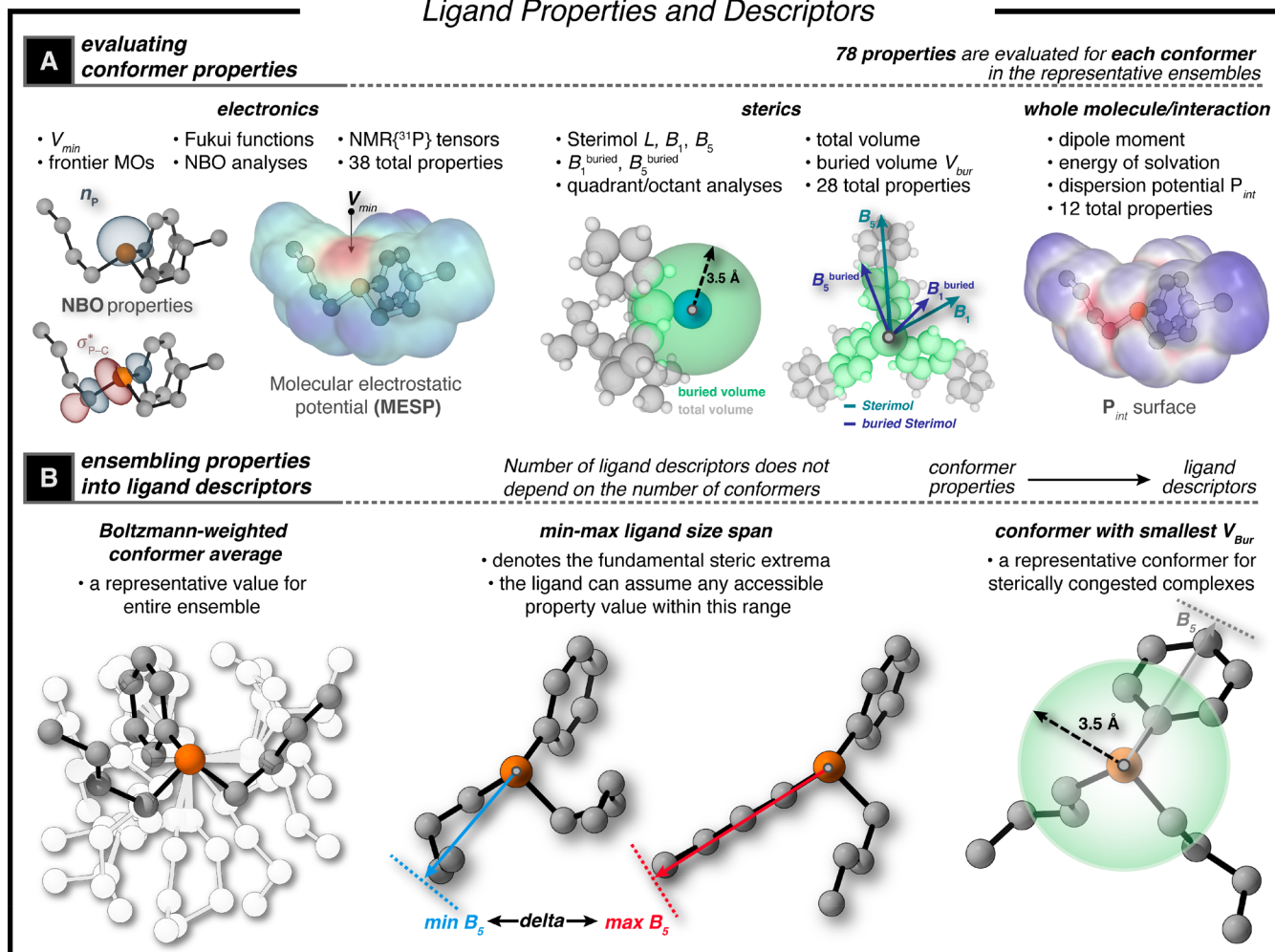
**Figure 1.** *kraken*: a comprehensive database of organophosphorus ligands. (A) A set of 1558 organophosphorus compounds is gathered, including literature and commercial sources. (B) Virtual libraries (VL) are built from the substituents of the initial P(III) set. The first level (Virtual Library 1, VL1) contains the initial set; VL2 results from a combinatorial approach with either all substituents equal or two different substituents per ligand (576 total unique fragments), yielding 331776 compounds; VL3 is a virtual library where all combinations are possible, i.e., all three substituents can be different, with over 191 million entries. (C) Conformer ensembles are generated for each of the P(III) molecules in VL1, at the GFN2-xTB level of theory. Each conformer is reoptimized by using DFT, with a total of 21437 conformers evaluated (average of 13.8 conformers per ligand). (D) 78 physical–organic properties are captured for every calculated conformer; Boltzmann averages, min–max steric extrema, and other representative conformers are curated for a total of 190 descriptors per ligand. (E) Chemical property spaces are defined and visualized by using dimensionality reduction techniques. (F) ML models are built to simulate a virtual property library for 331776 compounds in VL2. VL3 is deployed by querying the ML models on demand.

## generating ligand conformer ensembles



**Figure 2.** Computational workflows used to build *kraken*. Free and  $Ni(CO)_3$ -complexed ligands from VL1 are subjected to a conformer search with CREST at the GFN2-xTB level. Ligand conformer ensembles are subjected to a conformer selection. DFT is used for geometry optimization and single points of the selected conformers as free ligands.  $P_n$  and  $P-[Ni]_n$  refer to individual conformers obtained from the conformer search of the free ligand and the  $LNi(CO)_3$  complex, respectively.

## Ligand Properties and Descriptors



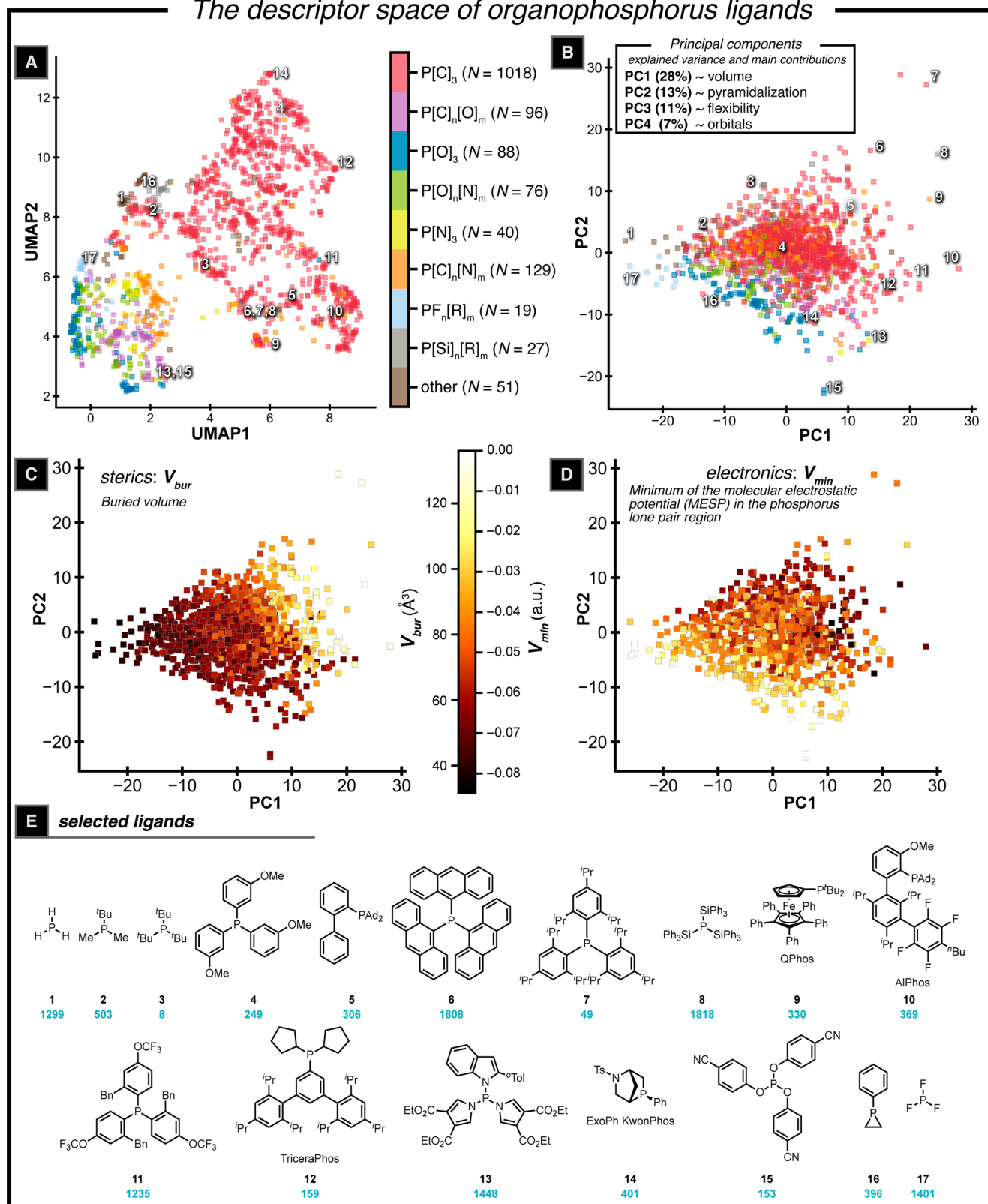
**Figure 3.** (A) Illustrations of some properties computed for each conformer. (B) Ensembling conformer properties to generate ligand descriptors. Note that absolute buried volume in  $\text{\AA}^3$  is used in this library instead of the more common percent buried volume  $\%V_{bur}$  ( $\%V_{bur} = V_{bur}/1.8$ ) to retain comparability with the total volume.

are folded toward the phosphorus lone pair, but it could never be smaller than when all chains are folded away. Thus, the smallest and largest attainable property values within the thermally accessible conformers of each ligand is defined as the

representative range, irrespective of the exact complex environment. Notably, the correct range can only be derived from a (sufficiently) complete conformer ensemble. To allow the workflow to operate at large scale and reasonable cost, we

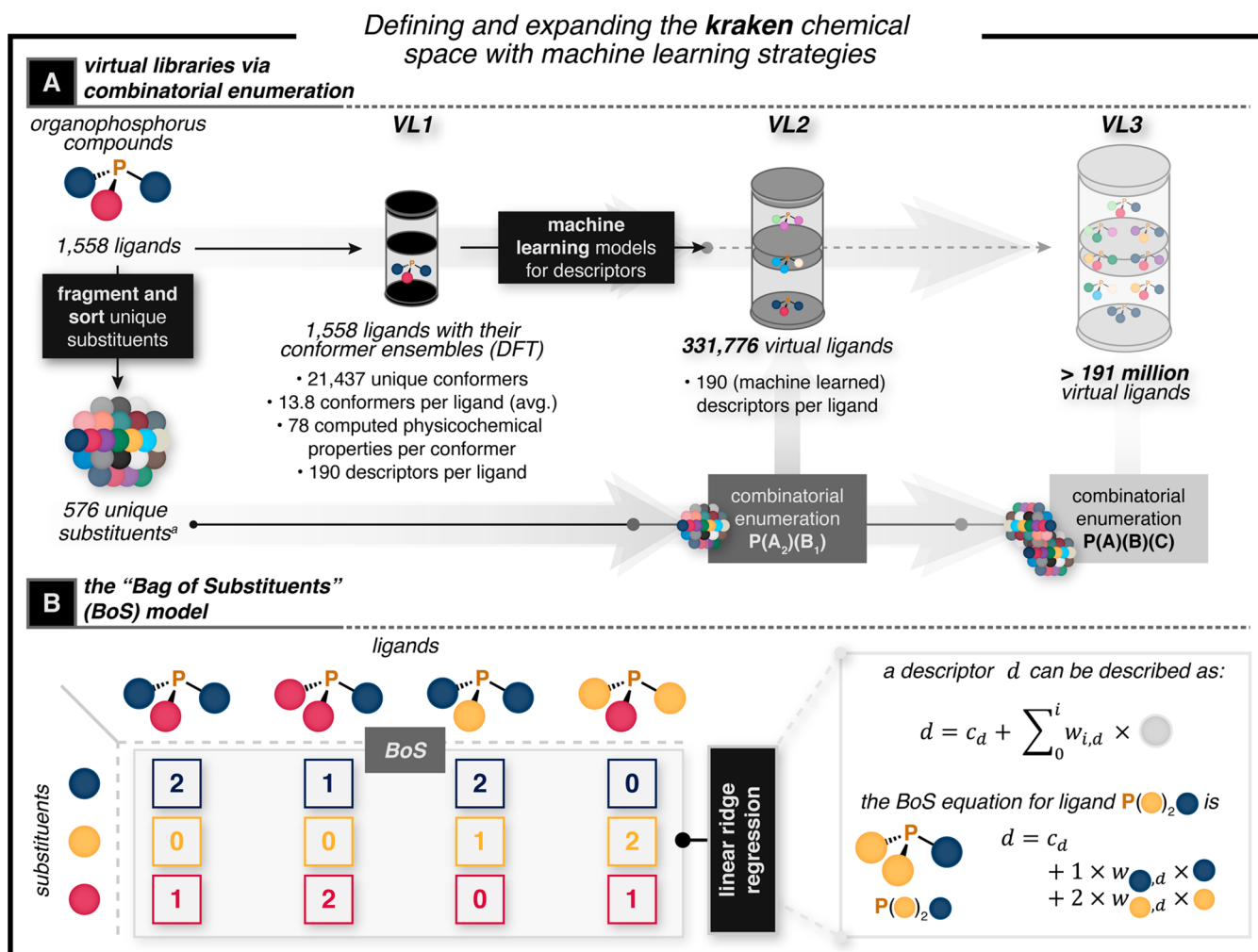


## The descriptor space of organophosphorus ligands



**Figure 4.** Property space visualizations of monodentate organophosphorus(III) ligands using UMAP and PCA. (A) Dimensionality reduction of the descriptor space with UMAP to two dimensions. (B) Dimensionality reduction of the descriptor space with PCA and projection of the corresponding results onto the two largest principal components, PC1 and PC2. The inset describes the types of descriptors with the highest loadings for each of the first four principal components along with the respective relative explained variance. (C) PC1–PC2 projection color-coded by Boltzmann-averaged  $V_{bur}$ . (D) PC1–PC2 projection color-coded by Boltzmann-averaged  $V_{min}$ . (E) Representative compounds that are highlighted in panels A and B (black numbers). The numbering scheme also used on the web app is shown in blue.



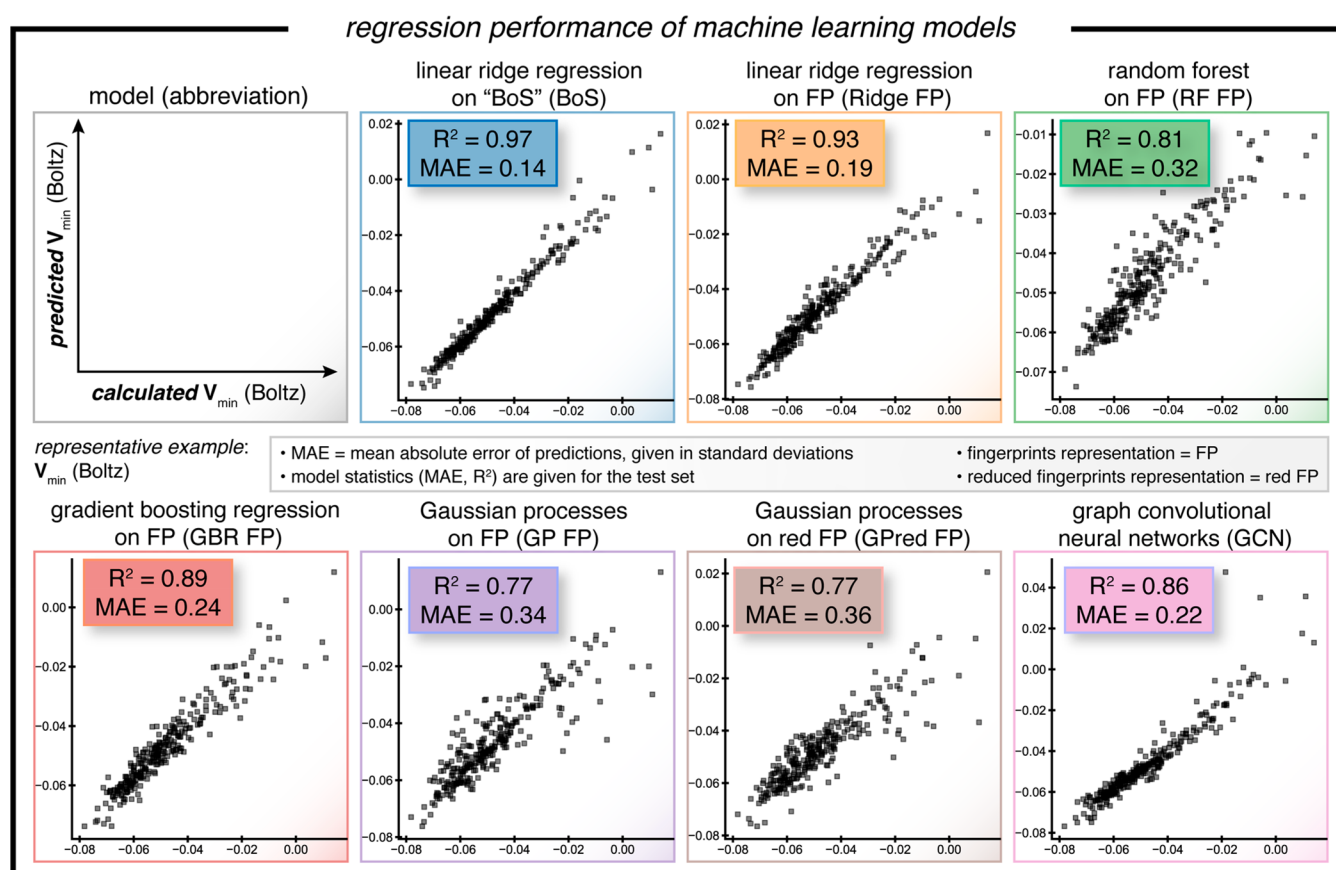


**Figure 5.** Defining and expanding the *kraken* chemical space with machine learning strategies. (A) Construction of the virtual libraries VL2 and VL3, respectively, from virtual library VL1, which comprises 1558 unique ligands. “The number of unique ligands excludes ligands in which the phosphorus atom is within a ring. (B) Illustration of the “Bag of Substituents” model to predict ligand descriptors based on substituent increments.  $d$ : descriptor;  $c_d$ : constant;  $w_{i,d}$ : substituent weight per descriptor;  $i$  in the sum: total number of occurrences for a given substituent in a ligand.

applied GFN2-xTB,<sup>40,41</sup> a semiempirical tight-binding method developed to deliver excellent molecular geometries at the fraction of the cost of DFT, together with the workflows implemented in CREST to generate conformer ensembles (Figure 2).<sup>42,43</sup> Because of the sensitivity of steric properties to structural changes, we used these ensembles to select the structures with extreme values for at least one steric descriptor to be evaluated using DFT. Importantly, the conformational space of each ligand was assessed in two reference states: free ligand and coordinated to  $Ni(CO)_3$ . Generally, conformations in the free ligand tend to occupy more space around the phosphorus lone pair and, hence, free ligands appear more sterically demanding than complexed ones. Both situations are important for describing catalytic processes including potential unwanted side reactions like ligand dissociation. For consistent results, the ligand conformations from both reference states are then optimized as free ligands using DFT. To distinguish between a ligand and its individual conformers, we ascribe *properties* to the individual conformers of a ligand and *descriptors* to a ligand. A total of 78 properties are evaluated for each conformer at the DFT level (for representative examples, see Figure 3A). Figure 3B illustrates the five descriptor variants (Boltzmann-weighted average, min, max,

delta, properties of conformer with smallest  $V_{bur}$ ) that are derived from those properties. All five descriptor variants are then used for properties that were found to vary strongly across conformers (mostly steric properties) whereas only the Boltzmann-weighted average is used for those properties that are less sensitive to conformation (mostly electronic properties) to avoid overly redundant descriptors, resulting in a total of 190 descriptors at the DFT level. For complete details, see section 2 of the [Supporting Information](#).

**Chemical Space Analysis.** With this data set, we set out to map the associated property space, understand the corresponding property limits, and unveil uncharted regions potentially inspiring forays toward new unique ligand classes. The traditional analysis of phosphine properties uses Tolman’s steric and electronic map, with the TEP on the abscissa and the Tolman cone angle on the ordinate.<sup>21</sup> This simple yet powerful visualization technique has helped chemists to survey available ligands rapidly and select structures with appropriate steric and electronic properties for specific applications. A more sophisticated version of Tolman’s map has been introduced by Fey and co-workers using LKB<sup>27,29</sup> (see above) by reducing multiple descriptors to fewer dimensions via principal component analysis (PCA). Inspired by this work, we applied



**Figure 6.** Regression performance of machine learning models. Illustrative performance of all seven types of ML models from this study for the prediction of  $V_{\min}$  (Boltz). BoS = Bag of Substituents; FP = fingerprint representation: circular fingerprints, radius = 2, folded to 1024 dimensions; red FP = reduced fingerprints representation: 100 most important fingerprint dimensions based on the feature importance of the GBR FP model. For additional details on the ML models see the [Supporting Information](#).

the recently developed Uniform Manifold Approximation and Projection<sup>44</sup> (UMAP, [Figure 4A](#)) as well as PCA ([Figure 4B](#)) to our entire database of DFT-computed ligands with all computed descriptors. These dimensionality reduction representations are available to interrogate on the interactive web application (<https://kraken.cs.toronto.edu>).

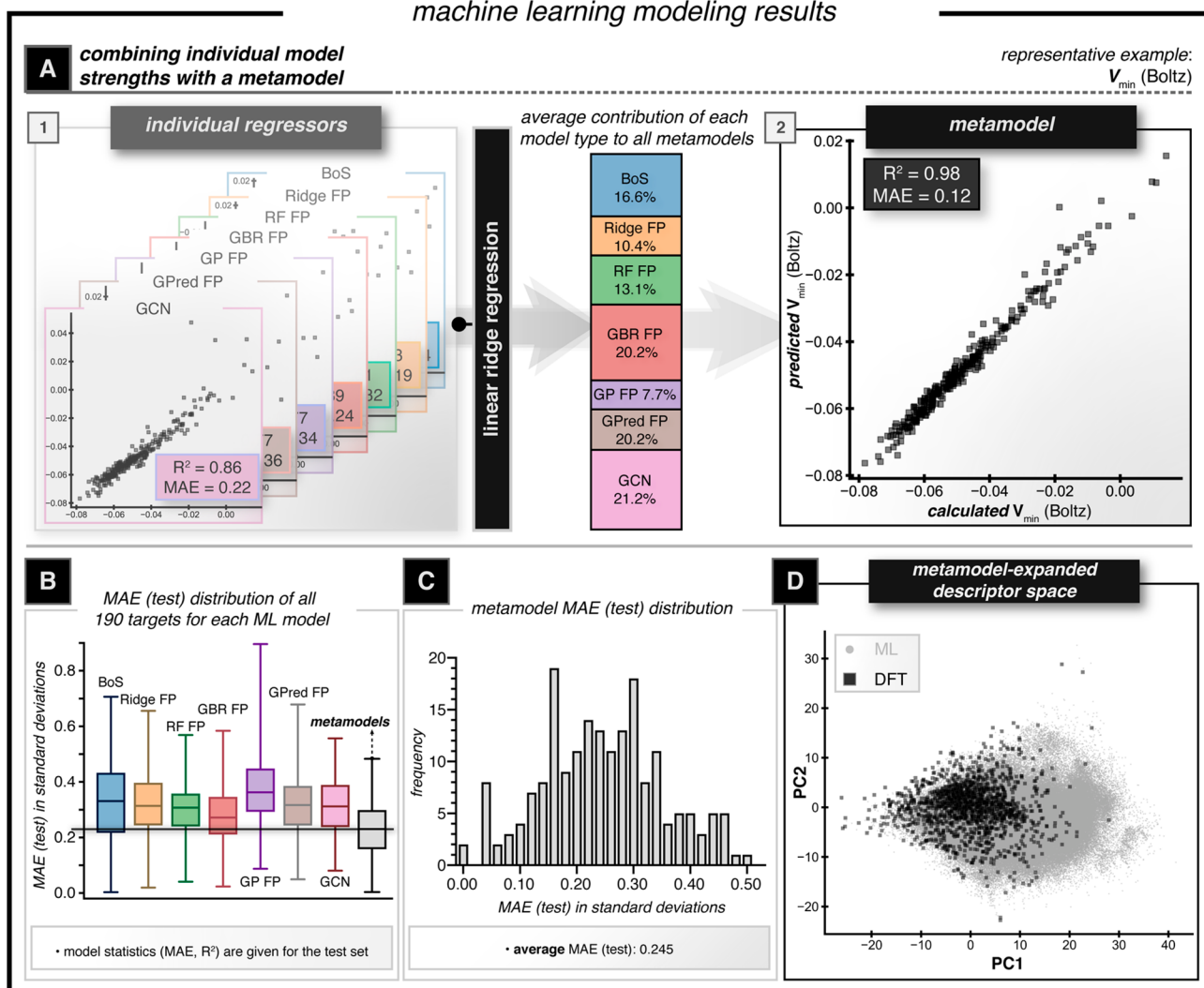
Nonlinear dimensionality reduction techniques can be employed to cluster compounds with a similar distribution of properties and for segregating distinct ligand classes from each other. For this purpose, we applied the UMAP technique as it preserves both local and global structure in the data and is computationally efficient.<sup>44</sup> The corresponding result is shown in [Figure 4A](#) by using the elements bonded directly to phosphorus for color-coding to illustrate the major phosphorus ligand classes. It is immediately obvious that the various ligand classes are well separated, demonstrating the superior ability for data classification of UMAP. This suggested that our descriptor set contains relevant information to differentiate chemically distinct ligand types. Notably, UMAP essentially segregates the database into two important ligand superclasses, i.e., phosphorus bound to relatively electropositive elements like carbon and silicon and phosphorus bound to at least one relatively electronegative element like oxygen or nitrogen, with some overlaps between these two. Importantly, this aligns well with the binding affinities of these ligands as the atom type bound to phosphorus affects this property most severely.

The principal components obtained from PCA define a linearly uncorrelated descriptor set condensing the information

contained in the database to as few dimensions as possible, while approximately preserving distance information in the descriptor space. This preservation of distances allows us to interpolate linearly between points in the descriptor space and, hence, understand the properties of unexplored regions as well. Accordingly, the resulting first two principal components (PC) were used to visualize the property space as depicted in [Figure 4B](#); illustrations with PC1–4 are found in section 3 of the [Supporting Information](#). Again, by coloring the data points with respect to the corresponding elements attached to the phosphorus atom, we can explore the relationships between common ligand classes, such as a smooth transition from phosphines (red) to phosphites (blue) via the intermediate phosphinites and phosphonites (purple) in the lower left of the chemical space.

Furthermore, not only can various ligand classes be distinguished, but the resulting principal components can be analyzed with respect to the properties they are encoding by investigating the most important descriptor loadings. PC1 generally represents total volume and PC2 pyramidalization. Evaluating the next most heavily weighted principal components, PC3 is mainly determined by flexibility descriptors related to the inclusion of conformer ensemble property information and PC4 contains general orbital descriptors (a more detailed analysis can be found in section 3 of the [Supporting Information](#)). Importantly, the added information from the computationally derived properties incorporates both depth and precision to compound

## machine learning modeling results



**Figure 7.** Machine learning modeling results. (A) Stacked linear ridge regression of the seven models was used to create a metamodel for each descriptor, shown with  $V_{\min}$ (Boltz) as an example. The model contributions shown in the middle are averages from all 190 metamodels. (B) Comparison of the mean absolute errors (MAEs) of the seven initial model classes and the metamodels across all descriptors. (C) Distribution of the MAEs of the metamodels across all descriptors (same data as the gray boxplot in part B). (D) Expansion of the descriptor space from VL1 to VL2 with the metamodels as illustrated by PCA with VL2 being projected onto the first two principal components obtained from VL1.

representation as compared to Tolman's mapping. Nevertheless, since the PCs combine various descriptors simultaneously, they represent a more integrated representation of the ligand space. To provide a more intuitive illustration of the PCA property mapping, the individual data points on the PCA plots were colored with respect to the buried volume<sup>39</sup> ( $V_{\text{bur}}$ , Figure 4C) and the minimum molecular electrostatic potential (MESP) in the phosphorus lone pair region,<sup>45</sup> which is correlated to the experimentally determined TEP ( $V_{\min}$ , Figure 4D). Notably, these plots demonstrate that PC1 generally trends with  $V_{\text{bur}}$  and PC2 trends with  $V_{\min}$ , even though it is not strongly collinear.

It is envisioned that these property maps can be used intuitively by chemists that may not be experts in data science. Specifically, when the basic requirements in terms of sterics and electronics are known from previous experiments, the rational selection of the best ligand types that meet various process needs such as cost, environmental, and/or performance goals should be straightforward, similar to how the

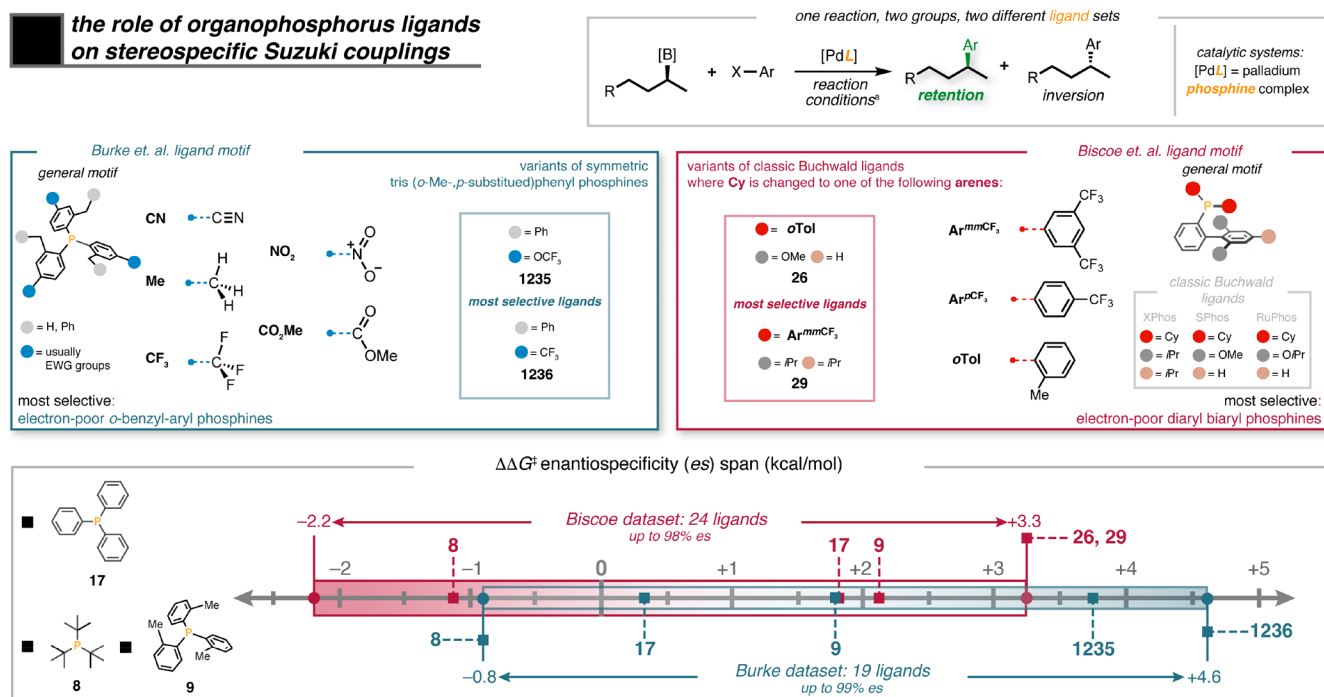
Solvent Selection Tool is applied by process chemists to locate the best solvent for a given reaction.<sup>46</sup>

**Expanding the Space with Machine Learning.** While we achieved a substantial coverage of the organophosphorus ligand space using quantum-chemical simulations, 1558 compounds merely constitute a fraction of the conceivable space of this ligand class. Our computational workflow is too resource-intensive to probe all possible compounds of interest and explore the sparsely covered territory more comprehensively (see Figures 4 and 5A). Hence, to complement the simulations described above, we investigated several complementary ML methods to expand the compound space in our library significantly and provide descriptor estimates for >300000 molecules.

Inspired by the Benson group-increment theory<sup>23,47</sup> in thermochemistry and the demonstration of substituent additivity for the TEP by Tolman,<sup>19</sup> we tested if descriptors can be expressed as the sum of constant contributions from each substituent at phosphorus. To accomplish this, we represented each ligand as a matrix of all unique substituents



## the role of organophosphorus ligands on stereospecific Suzuki couplings



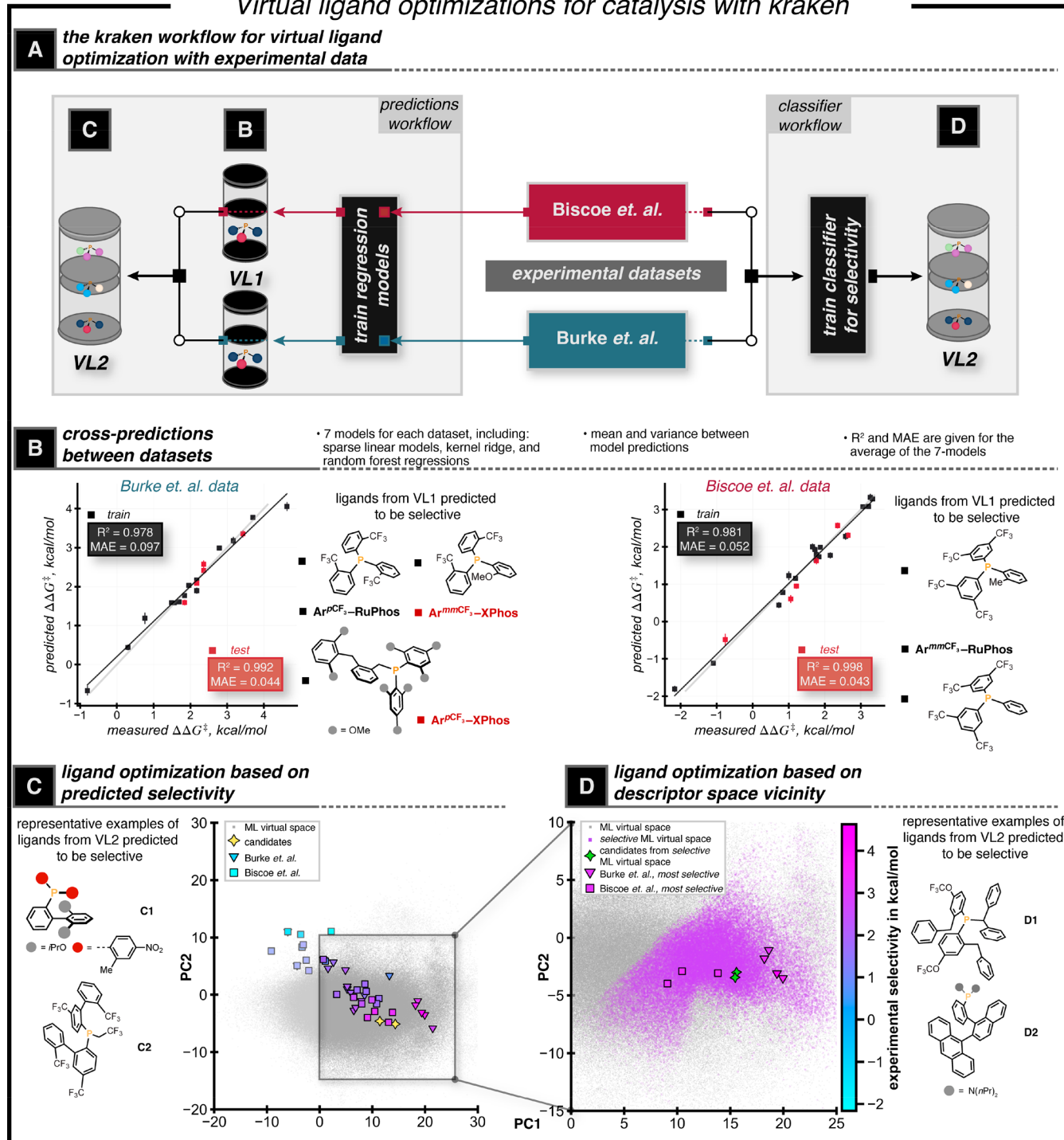
**Figure 8.** Phosphines in enantiospecific Pd-catalyzed  $sp^3$ - $sp^2$  cross-coupling reactions of alkylboronic acids and aryl halides as an application case study. “Conditions Biscoe et al.: [B] =  $BF_3K$  (R), R = Ph, X-Ar = 1-Cl-4- $CO_2Et$ - $C_6H_4$ , [PdL] = G3 Buchwald precatalyst (10 mol %), base =  $K_2CO_3$  (3 equiv), solvent = toluene:H<sub>2</sub>O (2:1), T = 100 °C, t = 24 h.<sup>53</sup> Conditions Burke et al.: [B] = B(OH)<sub>2</sub> (S), R = H, X-Ar = 1-Br-4-Ph- $C_6H_4$ , [PdL] =  $Pd_2dba_3$  (5 mol %) + 10 mol % L, base =  $Ag_2O$  (3 equiv), solvent = dioxane, T = 85 °C, t = 24 h.<sup>54</sup>

bound to the central phosphorus atom containing the number of each substituent present in a particular compound (Figure 5B), which we term “Bag of Substituents” (BoS). For instance,  $PMe_2tBu$  would be encoded by the features “Me” and “tBu”, with a value of 2 in the former column, a value of 1 in the latter, and zeros in all other feature columns (576 in total). Linear regression of each descriptor individually was used with the BoS encoding to assess the additivity hypothesis. The coefficients of determination are a measure of how well the additivity assumptions hold for a descriptor and the trained weights correspond to the group increments. It should be noted that this model is inherently incapable of extrapolating to unseen substituents. As a consequence, all substituents needed to be included at least once in the training data, and when possible, we enforced it to contain at least two occurrences. Apart from this constraint, we used a random 60:20:20 train-validation-test split. Good prediction quality was observed for a number of descriptors (58 properties with  $R^2_{test} \geq 0.80$ ). As expected,  $V_{min}^{(Boltz)}$  ( $R^2 = 0.97$ , cf. Figure 6) and  $V_{bur}^{(Boltz)}$  ( $R^2 = 0.95$ ) are well predicted. Interestingly, several descriptors that may not be expected *a priori* to be “additive” are also predicted with good accuracy, such as the Boltzmann-averaged NBO partial charge at the phosphorus atom ( $R^2_{test} > 0.99$ ).

While the BoS encoding strategy is relatively effective, some descriptors are not well predicted (45 properties with  $R^2 < 0.50$ ) as is expected when substituent interactions or conformational effects are present that this simple model cannot incorporate. Thus, we used molecular fingerprints and graphs as more generalizable features to expand our predictive capacities. With those representations, we also applied other model types such as random forest (RF),<sup>7</sup> gradient boosting regressions (GBR),<sup>48,49</sup> Gaussian processes (GP),<sup>49,50</sup> and graph convolutional neural networks<sup>51</sup> (GCN; see Figure 6 for

the performance on one representative descriptor; more details on the ML models are in the Supporting Information). Each of the models was found to be accurately predictive for a subset of descriptors. However, as none of the approaches were consistently the best for all the descriptors considered, we generated one metamodel for each descriptor. This was accomplished by ensembling all the models linearly to maximize the overall prediction quality. The performance of the metamodel predictors is illustrated in Figure 7A with  $V_{min}$  as an example and in Figure 7B,C for all targets. We then applied the metamodels to the >300000 compounds arising from unary and binary combinations (i.e., general structures  $PA_3$  and  $PA_2B$ ) of all unique substituents present in our original library (VL1) to create an extensive virtual library (VL2) with estimated descriptor values. This chemical space can be visualized in a new PCA plot revealing the virtual space now available (see Figure 7D). Compared to the PCA plot of VL1 (cf. Figure 4B), the plot of VL2 appears more continuous in the descriptor ranges covered and extrapolates considerably into underexplored chemical space, thereby encompassing many new structures that one might want to explore in future applications in a single lookup table. Intriguingly, this concept can readily be applied to all ternary substituent combinations of type PABC to obtain an even larger virtual library (VL3) with ca. 191 million entries that holds additional potential for processes involving P-chiral compounds. However, we only envision on-demand queries at this stage because hosting the entire data set is impractical due to its size, especially considering the lower practical utility associated with the much more difficult synthesis required of phosphorus compounds bearing three instead of two distinct substituents. While the metamodels generally perform well, the predictions should still only be treated as estimates with limited accuracy, especially in the extrapolated parts of the chemical space. The

## Virtual ligand optimizations for catalysis with kraken



**Figure 9.** Using *kraken* for virtual ligand optimizations in asymmetric catalysis by using the data shown in Figure 8. (A) General workflow for the case study. (B) Statistical modeling of experimental results to predict how data from one reported reaction could inform ligand choice in the other through a virtual screen of VL1 for ligands that are predicted to result in high selectivity for the stereoretentive cross-coupling. (C) Combining the statistical models for both reactions to evaluate the entirety of VL2 for new selective ligands. (D) Exploring the PCA descriptor space to determine ligands with novel structures in the high-selectivity regime.

estimates could for example suffice to obtain ligand suggestions from the desired location in chemical space, which can then be subjected to the computational workflow and obtain the actual DFT-level descriptors.

**Inverse Ligand Design.** Finally, we aimed to demonstrate the immediate practical applicability of *kraken* to a typical problem common in reaction development and ligand

design.<sup>52</sup> Specifically, we wanted to utilize the ML-predicted database to identify viable alternative ligands for a selective catalytic reaction. To do this, we revisited two independent studies by Biscoe and Burke, respectively,<sup>53,54</sup> that reported enantiospecific Pd-catalyzed  $\text{sp}^3\text{-sp}^2$  cross-coupling reactions of stereodefined alkylboronic acid derivatives with aryl halides. The two studies identified unique ligands that successfully

achieve high levels of stereoretention (Figure 8). In the Biscoe study, the ligand discovery was guided by using predictions from statistical modeling that electron-poor Buchwald-type<sup>55,56</sup> biaryl phosphine ligands were the best performers. The Burke study also discovered that electron-poor ligands were required, but a different core structure, one based on *o*-tolyl phosphines, was found to be required for highly selectivity in this reaction. Intuitively, the best ligands from either study are structurally unique, and a practicing organic chemist would not necessarily think to substitute one with the other. In addition, the reaction conditions, while distinct, are similar enough to expect qualitatively comparable selectivity of the ligands under each condition.

On the basis of these findings, we hypothesized that *kraken*'s descriptors applied to an original data set could be used to predict similar ligand structures found to be optimal in the complementary reaction (Figure 9A). As a first step, several statistical models of each data set were constructed (details in the Supporting Information) by correlating experimental results to the ligand descriptors, which were included in VL1 (Figure 9B). Unique models were averaged to provide robust predictions of which ligands would provide high selectivity.<sup>57</sup> Gratifyingly, trained on the results reported by the Burke group, our predictions identify the exact ligands reported by the Biscoe group as most selective. Similarly, regressing the Biscoe data set and virtually screening VL1 revealed untested electron-poor triaryl phosphines, in particular, Buchwald and *o*-tolyl derivatives, as most selective. This suggests that these two reactions likely proceed via similar mechanistic pathways in the stereo-determining events.

After this successful validation of the interconnectivity of the two reactions, we combined the two data sets to enhance the robustness of the predictions while exploring the entire virtual search space of VL2. This is visualized in the PCA plot depicted in Figure 9C wherein the black-framed points represent the experimental data from the two studies, atop the ML library in gray. We were then interested in comparing two distinct approaches to suggest novel ligands in a large search space. First, we applied the averaged regression models that were trained on the experimental results to the entire VL2 to obtain selectivity predictions and robustness estimates. The ligand predictions were then curated by filtering structures through descriptor limits reported for this process (small ligands)<sup>53</sup> and ligands that presumably would not form a metal complex (very large ligands). As a result, we obtained ~100 ligands that are predicted to provide selective stereoretentive cross-coupling. Many of these are bulky and electron-poor Buchwald-type ligands, represented by the two structures in Figure 9D. Notably, ligand **D1** merges structural elements into a hybrid of both the Biscoe and Burke ligand designs.

While this approach likely provides relatively safe predictions with structural similarity to the best experimental ligands, we envisioned an explorative strategy providing more structural diversity by analyzing the relative positions of ligands in the descriptor space. We classified the "more selective" and "less selective" regions in this space by proximity to the nearest experimental data point in the first four principal components and ranked the resulting >30000 structures by minimizing the distance to the most selective experimental ligands. This explorative classification method suggests unexplored ligands that upon inspection have some structural familiarity to both Burke's and Biscoe's ligand designs, which is highly encouraging. This strategy would be especially effective when

a researcher has relatively sparse data early in an optimization campaign as the local neighborhoods of the active space could be rapidly explored. We also envision this process will be valuable in iterative ligand searches, especially when commercial ligands only provide modest performance.

## CONCLUSIONS AND OUTLOOK

We have developed *kraken*, which covers 300000 monodentate organophosphorus(III) ligands with 190 property descriptors including an extensive description of their conformer dependence, mapping essentially the complete space of conceivable structures that could be used in organo(transition)metal reactions. We demonstrate its application in visualizing the associated property space, predicting properties of molecules not subjected to our full quantum-chemical workflow, and applying the corresponding results to inverse catalyst design.

*Kraken* is accessible as a web application (<https://kraken.cs.toronto.edu>). Computed data are available at the semi-empirical QM, DFT, and ML levels of theory. For 1558 organophosphorus compounds, there are both semiempirical QM and DFT data comprising 190 computed descriptors and properties as well as the coordinates information for the associated conformers. The ML data consist of 331776 entries obtained by generating all organophosphorus ligands with up to two distinct substituents combinatorially and training the models on the DFT data set (see above). Lastly, around 191 million distinct organophosphorus compounds can be queried to generate the ML property predictions on the fly.

Overall, we believe that the property maps generated by common dimensionality reduction techniques included in the *kraken* platform can be a valuable aid in the understanding of the space of organophosphorus ligands. We envision that it will enable synthetic chemists to perform computer-assisted interactive ligand exploration and provide new insights into relevant properties to solve a given problem. The *kraken* tool may enable informed catalyst design based on organophosphorus ligands, facilitate the optimization of reaction process parameters, inspire new ligand choices, and promote the synthesis of new organophosphorus compounds. The database and tools reported herein are currently being applied to enhance reaction optimization<sup>58–61</sup> and mechanistic workflows.<sup>62</sup> The open-source nature of our codes, as well as the open database, is designed to be extended by others, and we welcome further contributions by the community.



## AUTHOR INFORMATION

### Corresponding Authors

**Tobias Gensch** – Department of Chemistry, University of Utah, Salt Lake City, Utah 84112, United States; Department of Chemistry, TU Berlin, 10623 Berlin, Germany; Present Address: Department of Chemistry, Carnegie Mellon University, 4400 Fifth Avenue, Pittsburgh, Pennsylvania 15213, United States; and Department of Chemical Engineering, Carnegie Mellon University, 5000 Forbes Avenue, Pittsburgh, Pennsylvania 15213, United States; [orcid.org/0000-0002-1937-0285](https://orcid.org/0000-0002-1937-0285); Email: [tobias.gensch@tu-berlin.de](mailto:tobias.gensch@tu-berlin.de)

**Matthew S. Sigman** – Department of Chemistry, University of Utah, Salt Lake City, Utah 84112, United States; [orcid.org/0000-0002-5746-8830](https://orcid.org/0000-0002-5746-8830); Email: [sigman@chem.utah.edu](mailto:sigman@chem.utah.edu)

**Alán Aspuru-Guzik** – Chemical Physics Theory Group, Department of Chemistry, University of Toronto, Toronto, Ontario M5S 3H6, Canada; Department of Computer Science, University of Toronto, Toronto, Ontario M5T 3A1, Canada; Vector Institute for Artificial Intelligence, Toronto, Ontario M5G 1M1, Canada; Lebovic Fellow, Canadian Institute for Advanced Research (CIFAR), Toronto, Ontario M5G, Canada; [orcid.org/0000-0002-8277-4434](https://orcid.org/0000-0002-8277-4434); Email: [alan@aspuru.com](mailto:alan@aspuru.com)

### Authors

**Gabriel dos Passos Gomes** – Chemical Physics Theory Group, Department of Chemistry, University of Toronto, Toronto, Ontario M5S 3H6, Canada; Department of Computer Science, University of Toronto, Toronto, Ontario M5T 3A1, Canada; Vector Institute for Artificial Intelligence, Toronto, Ontario M5G 1M1, Canada; [orcid.org/0000-0002-8235-5969](https://orcid.org/0000-0002-8235-5969)

**Pascal Friederich** – Chemical Physics Theory Group, Department of Chemistry, University of Toronto, Toronto, Ontario M5S 3H6, Canada; Department of Computer Science, University of Toronto, Toronto, Ontario M5T 3A1, Canada; Institute of Nanotechnology, Karlsruhe Institute of Technology, 76344 Eggenstein-Leopoldshafen, Germany; [orcid.org/0000-0003-4465-1465](https://orcid.org/0000-0003-4465-1465)

**Ellyn Peters** – Department of Chemistry, University of Utah, Salt Lake City, Utah 84112, United States

**Théophile Gaudin** – Department of Computer Science, University of Toronto, Toronto, Ontario M5T 3A1, Canada; IBM Research Zurich, 8803 Rüschlikon, Switzerland

**Robert Pollice** – Chemical Physics Theory Group, Department of Chemistry, University of Toronto, Toronto, Ontario M5S 3H6, Canada; Department of Computer Science, University of Toronto, Toronto, Ontario M5T 3A1, Canada; [orcid.org/0000-0001-8836-6266](https://orcid.org/0000-0001-8836-6266)

**Kjell Jorner** – Chemical Physics Theory Group, Department of Chemistry, University of Toronto, Toronto, Ontario M5S 3H6, Canada; Department of Computer Science, University of Toronto, Toronto, Ontario M5T 3A1, Canada; Early Chemical Development, Pharmaceutical Sciences, R&D, AstraZeneca, Macclesfield K10 2NA, United Kingdom

**AkshatKumar Nigam** – Chemical Physics Theory Group, Department of Chemistry, University of Toronto, Toronto, Ontario M5S 3H6, Canada; Department of Computer Science, University of Toronto, Toronto, Ontario M5T 3A1, Canada; [orcid.org/0000-0002-5152-2082](https://orcid.org/0000-0002-5152-2082)

**Michael Lindner-D'Addario** – Chemical Physics Theory Group, Department of Chemistry, University of Toronto, Toronto, Ontario M5S 3H6, Canada; Department of Computer Science, University of Toronto, Toronto, Ontario M5T 3A1, Canada

### Author Contributions

T.Ge., G.P.G., and P.F. contributed equally to this work.

### Notes

The authors

declare the following competing financial interest(s): A.A.-G. is a co-founder and the Chief Visionary Officer at Kebotix Inc.

The database can be accessed and used free of charge via the web application at <https://kraken.cs.toronto.edu>. MORFEUS can be freely accessed at <https://github.com/kjelljorner/morfeus>. The collection of workflow codes and machine learning models used in this project will be available at <https://github.com/aspuru-guzik-group/kraken>.

### ACKNOWLEDGMENTS

T.Ge. thanks the Leopoldina Fellowship Programme of the German National Academy of Sciences Leopoldina (LPDS 2017-18). T.Ge. is funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany's Excellence Strategy (EXC 2008/1-390540038) and by a Liebig Fellowship of the Fonds der Chemischen Industrie. G.P.G. gratefully acknowledges the Natural Sciences and Engineering Research Council of Canada (NSERC) for the Banting Postdoctoral Fellowship. R.P. acknowledges funding through a Postdoc.Mobility fellowship by the Swiss National Science Foundation (SNSF, Project No. 191127). K.J. was a fellow of the AstraZeneca Postdoc Programme (2018–2020). M.L.D. gratefully acknowledges the Fonds de Recherche Québec Nature et Technologies (FRQNT) for the BIX Master's Scholarship and support from the Queen Elizabeth II Graduate Scholarship in Science and Technology (QEII-GSST). The support and resources from the Center for High Performance Computing at the University of Utah are gratefully acknowledged. We acknowledge the Defense Advanced Research Projects Agency (DARPA) under the Accelerated Molecular Discovery Program under Cooperative Agreement No. HR00111920027 dated August 1, 2019. The content of the information presented in this work does not necessarily reflect the position or the policy of the Government. A.A.-G. thanks Anders G. Frøseth for his generous support. A.A.-G. also acknowledges the generous support of Natural Resources Canada and the Canada 150 Research Chairs program. We thank Compute Canada for computational resources. DFT and xtb calculations were performed on the Niagara supercomputer at the SciNet HPC Consortium. SciNet is funded by the Canada Foundation for Innovation; the Government of Ontario; Ontario Research Fund - Research Excellence; and the University of Toronto. Machine learning models were developed and trained on the supercomputer Beluga from École de technologie supérieure, managed by Calcul Québec and Compute Canada. The operation of this supercomputer is funded by the Canada Foundation for Innovation (CFI), the ministère de l'Économie, de la science et de l'innovation du Québec (MESI), and the Fonds de recherche du Québec – Nature et technologies

(FRQ-NT). We are grateful to UofT Matter Lab system administrators Dr. Claire Yu and Chris Crebolder for helping with the deployment of the web app. M.S.S. and E.P. thank the NSF under the CCI Center for Computer Assisted Synthesis (CHE-1925607) for support.

## REFERENCES

- (1) Williams, W. L.; Zeng, L.; Gensch, T.; Sigman, M. S.; Doyle, A. G.; Anslyn, E. V. The Evolution of Data-Driven Modeling in Organic Chemistry. *ACS Cent. Sci.* **2021**, 7 (10), 1622–1637.
- (2) Coley, C. W.; Eyke, N. S.; Jensen, K. F. Autonomous Discovery in the Chemical Sciences Part I: Progress. *Angew. Chem., Int. Ed.* **2020**, 59 (51), 22858–22893.
- (3) Coley, C. W.; Eyke, N. S.; Jensen, K. F. Autonomous Discovery in the Chemical Sciences Part II: Outlook. *Angew. Chem., Int. Ed.* **2020**, 59 (52), 23414–23436.
- (4) Sanchez-Lengeling, B.; Aspuru-Guzik, A. Inverse Molecular Design Using Machine Learning: Generative Models for Matter Engineering. *Science* **2018**, 361 (6400), 360–365.
- (5) David, L.; Thakkar, A.; Mercado, R.; Engkvist, O. Molecular Representations in AI-Driven Drug Discovery: A Review and Practical Guide. *J. Cheminf.* **2020**, 12 (1), 56.
- (6) Goscinski, A.; Fraux, G.; Imbalzano, G.; Ceriotti, M. The Role of Feature Space in Atomistic Learning. *Mach. Learn. Sci. Technol.* **2021**, 2 (2), 025028.
- (7) Ahneman, D. T.; Estrada, J. G.; Lin, S.; Dreher, S. D.; Doyle, A. G. Predicting Reaction Performance in C–N Cross-Coupling Using Machine Learning. *Science* **2018**, 360 (6385), 186–190.
- (8) Zahrt, A. F.; Henle, J. J.; Rose, B. T.; Wang, Y.; Darrow, W. T.; Denmark, S. E. Prediction of Higher-Selectivity Catalysts by Computer-Driven Workflow and Machine Learning. *Science* **2019**, 363 (6424), No. eaau5631.
- (9) St. John, P. C.; Guan, Y.; Kim, Y.; Kim, S.; Paton, R. S. Prediction of Organic Homolytic Bond Dissociation Enthalpies at near Chemical Accuracy with Sub-Second Computational Cost. *Nat. Commun.* **2020**, 11 (1), 1–12.
- (10) Santiago, C. B.; Guo, J.-Y.; Sigman, M. S. Predictive and Mechanistic Multivariate Linear Regression Models for Reaction Development. *Chem. Sci.* **2018**, 9 (9), 2398–2412.
- (11) Wen, M.; Blau, S. M.; Spotte-Smith, E. W. C.; Dwaraknath, S.; Persson, K. A. BonDNet: A Graph Neural Network for the Prediction of Bond Dissociation Energies for Charged Molecules. *Chem. Sci.* **2021**, 12 (5), 1858–1868.
- (12) Pyzer-Knapp, E. O.; Li, K.; Aspuru-Guzik, A. Learning from the Harvard Clean Energy Project: The Use of Neural Networks to Accelerate Materials Discovery. *Adv. Funct. Mater.* **2015**, 25 (41), 6495–6502.
- (13) Jain, A.; Ong, S. P.; Hautier, G.; Chen, W.; Richards, W. D.; Dacek, S.; Cholia, S.; Gunter, D.; Skinner, D.; Ceder, G.; Persson, K. A. Commentary: The Materials Project: A Materials Genome Approach to Accelerating Materials Innovation. *APL Mater.* **2013**, 1 (1), 011002.
- (14) Kirklin, S.; Saal, J. E.; Meredig, B.; Thompson, A.; Doak, J. W.; Aykol, M.; Rühl, S.; Wolverton, C. The Open Quantum Materials Database (OQMD): Assessing the Accuracy of DFT Formation Energies. *npj Comput. Mater.* **2015**, 1 (1), 15010.
- (15) Taylor, R. H.; Rose, F.; Toher, C.; Levy, O.; Yang, K.; Buongiorno Nardelli, M.; Curtarolo, S. A RESTful API for Exchanging Materials Data in the AFLOWLIB. *Comput. Mater. Sci.* **2014**, 93, 178–192.
- (16) Hachmann, J.; Olivares-Amaya, R.; Atahan-Evrenk, S.; Amador-Bedolla, C.; Sánchez-Carrera, R. S.; Gold-Parker, A.; Vogt, L.; Brockway, A. M.; Aspuru-Guzik, A. The Harvard Clean Energy Project: Large-Scale Computational Screening and Design of Organic Photovoltaics on the World Community Grid. *J. Phys. Chem. Lett.* **2011**, 2 (17), 2241–2251.
- (17) Mamun, O.; Winther, K. T.; Boes, J. R.; Bligaard, T. High-Throughput Calculations of Catalytic Properties of Bimetallic Alloy Surfaces. *Sci. Data* **2019**, 6 (1), 1–9.
- (18) Chanussot, L.; Das, A.; Goyal, S.; Lavril, T.; Shuaibi, M.; Riviere, M.; Tran, K.; Heras-Domingo, J.; Ho, C.; Hu, W.; Palizhati, A.; Sriram, A.; Wood, B.; Yoon, J.; Parikh, D.; Zitnick, C. L.; Ulissi, Z. Open Catalyst 2020 (OC20) Dataset and Community Challenges. *ACS Catal.* **2021**, 11 (10), 6059–6072.
- (19) Tolman, C. A. Electron Donor-Acceptor Properties of Phosphorus Ligands. Substituent Additivity. *J. Am. Chem. Soc.* **1970**, 92 (10), 2953–2956.
- (20) Tolman, C. A. Phosphorus Ligand Exchange Equilibria on Zerovalent Nickel. A Dominant Role for Steric Effects. *J. Am. Chem. Soc.* **1970**, 92 (10), 2956–2965.
- (21) Tolman, C. A. Steric Effects of Phosphorus Ligands in Organometallic Chemistry and Homogeneous Catalysis. *Chem. Rev.* **1977**, 77 (3), 313–348.
- (22) Hansch, C.; Leo, A.; Taft, R. W. A Survey of Hammett Substituent Constants and Resonance and Field Parameters. *Chem. Rev.* **1991**, 91 (2), 165–195.
- (23) Benson, S. W.; Cruickshank, F. R.; Golden, D. M.; Haugen, G. R.; O'Neal, H. E.; Rodgers, A. S.; Shaw, R.; Walsh, R. Additivity Rules for the Estimation of Thermochemical Properties. *Chem. Rev.* **1969**, 69 (3), 279–324.
- (24) Bjørsvik, H.-R.; Hansen, U. M.; Carlson, R.; et al. Principal Properties of Monodentate Phosphorus Ligands. Predictive Model for the Carbonyl Absorption Frequencies in Ni(CO)<sub>3</sub>L Complexes. *Acta Chem. Scand.* **1997**, 51, 733–741.
- (25) Perrin, L.; Clot, E.; Eisenstein, O.; Loch, J. A.; Crabtree, R. H. Computed Ligand Electronic Parameters from Quantum Chemistry and Their Relation to Tolman Parameters, Lever Parameters, and Hammett Constants. *Inorg. Chem.* **2001**, 40 (23), 5806–5811.
- (26) Cooney, K. D.; Cundari, T. R.; Hoffman, N. W.; Pittard, K. A.; Temple, M. D.; Zhao, Y. A Priori Assessment of the Stereoelectronic Profile of Phosphines and Phosphites. *J. Am. Chem. Soc.* **2003**, 125 (14), 4318–4324.
- (27) Fey, N.; Tsepis, A. C.; Harris, S. E.; Harvey, J. N.; Orpen, A. G.; Mansson, R. A. Development of a Ligand Knowledge Base, Part 1: Computational Descriptors for Phosphorus Donor Ligands. *Chem. - Eur. J.* **2006**, 12 (1), 291–302.
- (28) Jover, J.; Fey, N.; Harvey, J. N.; Lloyd-Jones, G. C.; Orpen, A. G.; Owen-Smith, G. J. J.; Murray, P.; Hose, D. R. J.; Osborne, R.; Purdie, M. Expansion of the Ligand Knowledge Base for Monodentate P-Donor Ligands (LKB-P). *Organometallics* **2010**, 29 (23), 6245–6258.
- (29) Durand, D. J.; Fey, N. Computational Ligand Descriptors for Catalyst Design. *Chem. Rev.* **2019**, 119 (11), 6561–6594.
- (30) Durand, D. J.; Fey, N. Building a Toolbox for the Analysis and Prediction of Ligand and Catalyst Effects in Organometallic Catalysis. *Acc. Chem. Res.* **2021**, 54 (4), 837–848.
- (31) Fey, N.; Howell, J. A. S.; Lovatt, J. D.; Yates, P. C.; Cunningham, D.; McArdle, P.; Gottlieb, H. E.; Coles, S. J. A Molecular Mechanics Approach to Mapping the Conformational Space of Diaryl and Triarylphosphines. *Dalt. Trans.* **2006**, 44 (46), 5464.
- (32) Suresh, C. H. Molecular Electrostatic Potential Approach to Determining the Steric Effect of Phosphine Ligands in Organometallic Chemistry. *Inorg. Chem.* **2006**, 45 (13), 4982–4986.
- (33) Baber, R. A.; Haddow, M. F.; Middleton, A. J.; Orpen, A. G.; Pringle, P. G.; Haynes, A.; Williams, G. L.; Papp, R. Ligand Stereoelectronic Effects in Complexes of Phospholanes, Phosphinanes, and Phosphapanes and Their Implications for Hydroformylation Catalysis. *Organometallics* **2007**, 26 (3), 713–725.
- (34) Barder, T. E.; Buchwald, S. L. Rationale Behind the Resistance of Dialkylbiarylphosphines toward Oxidation by Molecular Oxygen. *J. Am. Chem. Soc.* **2007**, 129 (16), 5096–5101.
- (35) Crawford, J. M.; Sigman, M. S. Conformational Dynamics in Asymmetric Catalysis: Is Catalyst Flexibility a Design Element? *Synthesis* **2019**, 51 (05), 1021–1036.

- (36) Brethomé, A. V.; Fletcher, S. P.; Paton, R. S. Conformational Effects on Physical-Organic Descriptors: The Case of Sterimol Steric Parameters. *ACS Catal.* **2019**, *9* (3), 2313–2323.
- (37) Fey, N.; Harvey, J. N.; Lloyd-Jones, G. C.; Murray, P.; Orpen, A. G.; Osborne, R.; Purdie, M. Computational Descriptors for Chelating P,P- And P,N-Donor Ligands. *Organometallics* **2008**, *27* (7), 1372–1383.
- (38) Wilbraham, L.; Berardo, E.; Turcani, L.; Jelfs, K. E.; Zwijnenburg, M. A. High-Throughput Screening Approach for the Optoelectronic Properties of Conjugated Polymers. *J. Chem. Inf. Model.* **2018**, *58* (12), 2450–2459.
- (39) Hillier, A. C.; Sommer, W. J.; Yong, B. S.; Petersen, J. L.; Cavallo, L.; Nolan, S. P. A Combined Experimental and Theoretical Study Examining the Binding of N-Heterocyclic Carbenes (NHC) to the Cp\*RuCl (Cp\* =  $\eta^5$ -C<sub>5</sub>Me<sub>5</sub>) Moiety: Insight into Stereoelectronic Differences Between Unsaturated and Saturated NHC Ligands. *Organometallics* **2003**, *22* (21), 4322–4326.
- (40) Grimme, S.; Bannwarth, C.; Shushkov, P. A Robust and Accurate Tight-Binding Quantum Chemical Method for Structures, Vibrational Frequencies, and Noncovalent Interactions of Large Molecular Systems Parametrized for All spd-Block Elements (Z = 1–86). *J. Chem. Theory Comput.* **2017**, *13* (5), 1989–2009.
- (41) Bannwarth, C.; Ehlert, S.; Grimme, S. GFN2-xTB - An Accurate and Broadly Parametrized Self-Consistent Tight-Binding Quantum Chemical Method with Multipole Electrostatics and Density-Dependent Dispersion Contributions. *J. Chem. Theory Comput.* **2019**, *15* (3), 1652–1671.
- (42) Grimme, S. Exploration of Chemical Compound, Conformer, and Reaction Space with Meta-Dynamics Simulations Based on Tight-Binding Quantum Chemical Calculations. *J. Chem. Theory Comput.* **2019**, *15* (5), 2847–2862.
- (43) Pracht, P.; Bohle, F.; Grimme, S. Automated Exploration of the Low-Energy Chemical Space with Fast Quantum Chemical Methods. *Phys. Chem. Chem. Phys.* **2020**, *22* (14), 7169–7192.
- (44) McInnes, L.; Healy, J.; Melville, J. UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction. *arXiv* **2018**, 1802.03426; <https://arxiv.org/abs/1802.03426> (accessed 2021-10-31).
- (45) Suresh, C. H.; Koga, N. Quantifying the Electronic Effect of Substituted Phosphine Ligands via Molecular Electrostatic Potential. *Inorg. Chem.* **2002**, *41* (6), 1573–1578.
- (46) Diorazio, L. J.; Hose, D. R. J.; Adlington, N. K. Toward a More Holistic Framework for Solvent Selection. *Org. Process Res. Dev.* **2016**, *20* (4), 760–773.
- (47) Benson, S. W.; Buss, J. H. Additivity Rules for the Estimation of Molecular Properties. Thermodynamic Properties. *J. Chem. Phys.* **1958**, *29* (3), 546–572.
- (48) Friederich, P.; Krenn, M.; Tamblyn, I.; Aspuru-Guzik, A. Scientific Intuition Inspired by Machine Learning-Generated Hypotheses. *Mach. Learn. Sci. Technol.* **2021**, *2* (2), 025027.
- (49) Friederich, P.; Dos Passos Gomes, G.; De Bin, R.; Aspuru-Guzik, A.; Balcells, D. Machine Learning Dihydrogen Activation in the Chemical Space Surrounding Vaska's Complex. *Chem. Sci.* **2020**, *11* (18), 4584–4601.
- (50) Lopez, S. A.; Sanchez-Lengeling, B.; de Goes Soares, J.; Aspuru-Guzik, A. Design Principles and Top Non-Fullerene Acceptor Candidates for Organic Photovoltaics. *Joule* **2017**, *1* (4), 857–870.
- (51) Duvenaud, D.; Maclaurin, D.; Aguilera-Iparraguirre, J.; Gómez-Bombarelli, R.; Hirzel, T.; Aspuru-Guzik, A.; Adams, R. P. Convolutional Networks on Graphs for Learning Molecular Fingerprints. In *Advances in Neural Information Processing Systems*; Cortes, C., Lawrence, N., Lee, D., Sugiyama, M., Garnett, R., Eds.; Curran Associates, Inc.: 2015; Vol. 28, pp 399–411.
- (52) Foscatto, M.; Jensen, V. R. Automated in Silico Design of Homogeneous Catalysts. *ACS Catal.* **2020**, *10* (3), 2354–2377.
- (53) Zhao, S.; Gensch, T.; Murray, B.; Niemeyer, Z. L.; Sigman, M. S.; Biscoe, M. R. Enantiodivergent Pd-Catalyzed C–C Bond Formation Enabled through Ligand Parameterization. *Science* **2018**, *362* (6415), 670–674.
- (54) Lehmann, J. W.; Crouch, I. T.; Blair, D. J.; Trobe, M.; Wang, P.; Li, J.; Burke, M. D. Axial Shielding of Pd(II) Complexes Enables Perfect Stereoretention in Suzuki-Miyaura Cross-Coupling of Csp<sup>3</sup> Boronic Acids. *Nat. Commun.* **2019**, *10* (1), 1263.
- (55) Surry, D. S.; Buchwald, S. L. Dialkylbiarylphosphines in Pd-Catalyzed Amination: A User's Guide. *Chem. Sci.* **2011**, *2* (1), 27–50.
- (56) Ingoglia, B. T.; Wagen, C. C.; Buchwald, S. L. Biarylmonophosphine Ligands in Palladium-Catalyzed C–N Coupling: An Updated User's Guide. *Tetrahedron* **2019**, *75* (32), 4199–4211.
- (57) Brethomé, A. V.; Paton, R. S.; Fletcher, S. P. Retooling Asymmetric Conjugate Additions for Sterically Demanding Substrates with an Iterative Data-Driven Approach. *ACS Catal.* **2019**, *9* (8), 7179–7187.
- (58) Christensen, M.; Yunker, L. P. E.; Adediji, F.; Häse, F.; Roch, L. M.; Gensch, T.; dos Passos Gomes, G.; Zepel, T.; Sigman, M. S.; Aspuru-Guzik, A.; Hein, J. E. Data-Science Driven Autonomous Process Optimization. *Commun. Chem.* **2021**, *4* (1), 112.
- (59) De Jesus Silva, J.; Bartalucci, N.; Jelier, B.; Grosslight, S.; Gensch, T.; Schunemann, C.; Muller, B.; Kamer, P. C. J.; Coperet, C.; Sigman, M. S.; Togni, A. Development and Molecular Understanding of a Pd□catalyzed Cyanation of Aryl Boronic Acids Enabled by High-Throughput Experimentation and Data Analysis. *Helv. Chim. Acta* **2021**.
- (60) Gensch, T.; Smith, S. R.; Colacot, T. J.; Timsina, Y.; Xu, G.; Glasspoole, B. W.; Sigman, M. S. Design and Application of a Screening Set for Monophosphine Ligands in Metal Catalysis. 2021-08-13. *ChemRxiv*; DOI: 10.33774/chemrxiv-2021-fgm7v (accessed 2021-12-07).
- (61) Zell, D.; Kingston, C.; Jermaks, J.; Smith, S. R.; Seeger, N.; Wassmer, J.; Sirois, L. E.; Han, C.; Zhang, H.; Sigman, M. S.; Gosselin, F. Stereoconvergent and -Divergent Synthesis of Tetrasubstituted Alkenes by Nickel-Catalyzed Cross-Couplings. *J. Am. Chem. Soc.* **2021**, *143* (45), 19078–19090.
- (62) Newman-Stonebraker, S. H.; Smith, S. R.; Borowski, J. E.; Peters, E.; Gensch, T.; Johnson, H. C.; Sigman, M. S.; Doyle, A. G. Univariate Classification of Phosphine Ligation State and Reactivity in Cross-Coupling Catalysis. *Science* **2021**, *374* (6565), 301–308.