

# 1 Data-Driven Strategies for Accelerated Materials Design

2 *Published as part of the Accounts of Chemical Research special issue “Data Science Meets Chemistry”.*

3 Robert Pollice, Gabriel dos Passos Gomes, Matteo Aldeghi, Riley J. Hickman, Mario Krenn,  
4 Cyrille Lavigne, Michael Lindner-D’Addario, AkshatKumar Nigam, Cher Tian Ser, Zhenpeng Yao,  
5 and Alán Aspuru-Guzik\*



Cite This: <https://dx.doi.org/10.1021/acs.accounts.0c00785>



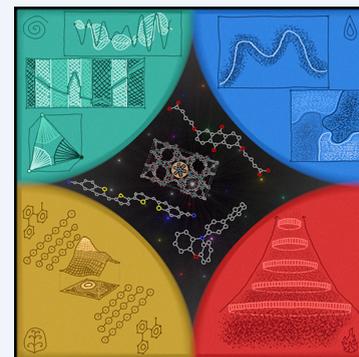
Read Online

ACCESS |

Metrics & More

Article Recommendations

6 **CONSPECTUS:** The ongoing revolution of the natural sciences by the advent of machine  
7 learning and artificial intelligence sparked significant interest in the material science community  
8 in recent years. The intrinsically high dimensionality of the space of realizable materials makes  
9 traditional approaches ineffective for large-scale explorations. Modern data science and  
10 machine learning tools developed for increasingly complicated problems are an attractive  
11 alternative. An imminent climate catastrophe calls for a clean energy transformation by  
12 overhauling current technologies within only several years of possible action available. Tackling  
13 this crisis requires the development of new materials at an unprecedented pace and scale. For  
14 example, organic photovoltaics have the potential to replace existing silicon-based materials to a  
15 large extent and open up new fields of application. In recent years, organic light-emitting diodes  
16 have emerged as state-of-the-art technology for digital screens and portable devices and are  
17 enabling new applications with flexible displays. Reticular frameworks allow the atom-precise  
18 synthesis of nanomaterials and promise to revolutionize the field by the potential to realize  
19 multifunctional nanoparticles with applications from gas storage, gas separation, and electrochemical energy storage to  
20 nanomedicine. In the recent decade, significant advances in all these fields have been facilitated by the comprehensive application  
21 of simulation and machine learning for property prediction, property optimization, and chemical space exploration enabled by  
22 considerable advances in computing power and algorithmic efficiency.



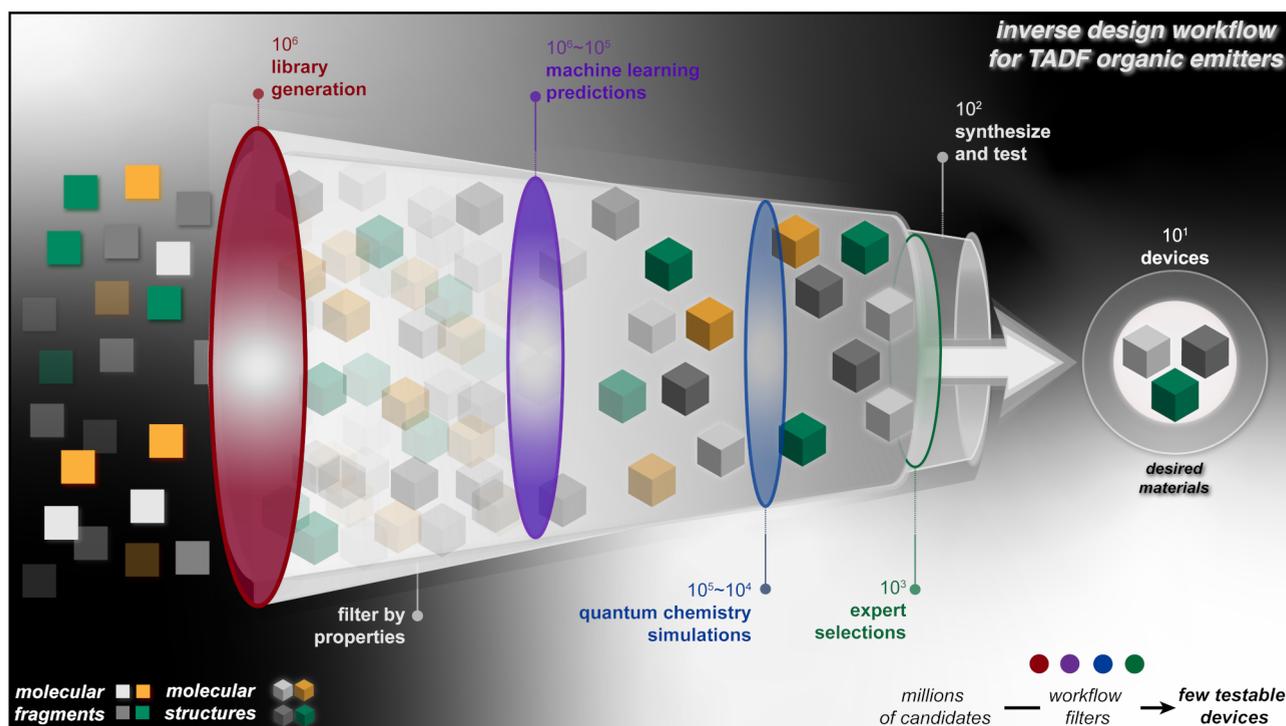
23 In this Account, we review the most recent contributions of our group in this thriving field of machine learning for material science.  
24 We start with a summary of the most important material classes our group has been involved in, focusing on small molecules as  
25 organic electronic materials and crystalline materials. Specifically, we highlight the data-driven approaches we employed to speed up  
26 discovery and derive material design strategies. Subsequently, our focus lies on the data-driven methodologies our group has  
27 developed and employed, elaborating on high-throughput virtual screening, inverse molecular design, Bayesian optimization, and  
28 supervised learning. We discuss the general ideas, their working principles, and their use cases with examples of successful  
29 implementations in data-driven material discovery and design efforts. Furthermore, we elaborate on potential pitfalls and remaining  
30 challenges of these methods. Finally, we provide a brief outlook for the field as we foresee increasing adaptation and implementation  
31 of large scale data-driven approaches in material discovery and design campaigns.

## 32 KEY REFERENCES

33 • Gómez-Bombarelli, R.; Aguilera-Iparraguirre, J.; Hirzel,  
34 T. D.; Duvenaud, D.; Maclaurin, D.; Blood-Forsythe, M.  
35 A.; Chae, H. S.; Einzinger, M.; Ha, D.-G.; Wu, T.;  
36 Markopoulos, G.; Jeon, S.; Kang, H.; Miyazaki, H.;  
37 Numata, M.; Kim, S.; Huang, W.; Hong, S. I.; Baldo, M.;  
38 Adams, R. P.; Aspuru-Guzik, A. Design of Efficient  
39 Molecular Organic Light-Emitting Diodes by a High-  
40 Throughput Virtual Screening and Experimental Ap-  
41 proach. *Nat. Mater.* **2016**, *15* (10), 1120–1127.<sup>1</sup>  
42 *Realization of an integrated inverse design workflow from*  
43 *high-throughput virtual screening to device testing for organic*  
44 *light-emitting diode materials.*

• Yao, Z.; Sanchez-Lengeling, B.; Bobbitt, N. S.; Bucior, B. 45  
J.; Kumar, S. G. H.; Collins, S. P.; Burns, T.; Woo, T. K.; 46  
Farha, O.; Snurr, R. Q.; Aspuru-Guzik, A. Inverse Design 47  
of Nanoporous Crystalline Reticular Materials with Deep 48  
Generative Models. *Nat. Mach. Intell.* **2021**, *3*, 76.<sup>2</sup> *An* 49  
*automated nanoporous materials discovery platform powered* 50  
*by a supramolecular variational autoencoder was built and* 51

Received: December 7, 2020



**Figure 1.** Inverse design workflow for thermally activated delayed fluorescence organic emitters from selecting fragments to device integration and testing.

52 demonstrated for the efficient exploration of the near infinite  
53 reticular chemical space and inverse design of reticular  
54 materials with desired functions like gas separation.

- 55 • Nigam, A.; Friederich, P.; Krenn, M.; Aspuru-Guzik, A.  
56 Augmenting Genetic Algorithms with Deep Neural  
57 Networks for Exploring the Chemical Space. In *International  
58 Conference on Learning Representations*; 2020.<sup>3</sup> *The  
59 proposal of a genetic algorithm enhanced by a neural network  
60 for inverse molecular design that can avoid convergence and  
61 bias molecule generation based on existing data sets.*
- 62 • Häse, F.; Roch, L. M.; Kreisbeck, C.; Aspuru-Guzik, A.  
63 Phoenix: A Bayesian Optimizer for Chemistry. *ACS Cent.  
64 Sci.* **20184** (9), 1134–1145.<sup>4</sup> *A probabilistic global  
65 optimization algorithm based on Bayesian kernel density  
66 estimation for the efficient parallel search of optimal  
67 experimental conditions.*

## 68 ■ INTRODUCTION

69 The tremendous rise of data science and machine learning (ML)  
70 in the last decades led to the suggestion that it constitutes the  
71 fourth pillar of science.<sup>5</sup> While data has always been at the heart  
72 of research, current hardware enables its utilization at an  
73 unprecedented scale.<sup>5</sup> Accordingly, our group, the Matter Lab,  
74 has been using ML extensively to accelerate the discovery of new  
75 materials, especially for clean energy technologies to combat  
76 climate catastrophe and enable innovative technologies.

77 In this Account, we define discovery as observing a previously  
78 unknown natural phenomenon or object,<sup>6,7</sup> and design as  
79 rationally devising an object based on a particular plan.<sup>8</sup>  
80 Typically, discovery precedes and inspires materials design, as  
81 design requires at least minimal knowledge of the necessary  
82 features. Therefore, large scale discovery helps to speed up the  
83 establishment of material design principles, *i.e.*, heuristics to  
84 realize particular designs, because they enable identifying

patterns in known matter with desired properties. In turn, 85  
successful design catalyzes the realization of new materials by 86  
restricting the search space to only the most promising regions in 87  
subsequent campaigns. 88

89 Herein, we review our work on organic electronic materials, 89  
crystalline materials, and data-driven methodologies for 90  
materials discovery and design, particularly high-throughput 91  
virtual screening, supervised learning, inverse molecular design, 92  
and Bayesian optimization. Moreover, we formulate general 93  
strategies for data-driven materials design our lab has adopted 94  
over the years and show how to implement them using ML. 95  
Finally, investigating these approaches critically, we propose 96  
typical use cases and highlight unsolved challenges. 97

## ■ APPLICATIONS

### Organic Electronic Materials

99 One of our research foci has been organic electronic materials.<sup>9</sup> 100  
Compared to silicon-based electronics, they offer several 101  
advantages, including low cost, low density, high mechanical 102  
flexibility and toughness, low energy consumption, and easy 103  
processability. Further, chemical derivatization is well-established 104  
making the accessible candidate space vast. 105

106 Accordingly, solar cells have experienced a remarkable surge 106  
because of the vast energy available from the sun and increasing 107  
efforts against a climate catastrophe. Organic photovoltaics<sup>10</sup> 108  
(OPVs) could replace commercial silicon-based devices if their 109  
power conversion efficiencies (PCEs) surpassed 10% and their 110  
lifetimes exceeded several thousands of hours. Notably, state-of- 111  
the-art OPVs reach 18% PCE in laboratory devices.<sup>11</sup> The 112  
Harvard Clean Energy Project (CEP) was initiated to find 113  
photoactive organic materials with high efficiencies.<sup>12</sup> Starting 114  
from 26 building blocks, selected based on expert knowledge to 115  
maximize performance and synthesizability,<sup>13</sup> 10<sup>7</sup> potential 116  
donors were generated. They were evaluated using high- 117

throughput virtual screening (HTVS, *vide infra*) via increasingly expensive property predictions. First, the library was assessed using linear descriptor models constructed from experimental data. Subsequently, electronic structure calculations were performed, and PCEs were estimated using the Scharber model with a fullerene as acceptor.<sup>14</sup> That way, about 1000 candidates with estimated PCEs of 11% and higher were identified.

Additionally, statistical analysis of the top-performing molecules revealed design principles for photoactive donors identifying building blocks more likely to exhibit high performance. Notably, the screening efforts led to the experimental characterization of an organic crystal with one of the highest reported hole mobilities reported at the time.<sup>15</sup> Subsequently, extending the CEP to nonfullerene acceptors, over 51 000 candidates were generated based on 107 expertly chosen fragments.<sup>16</sup> More sophisticated property calibration with Gaussian processes and a modified Scharber model improved PCE predictions with a well-studied electron donor. Overall, 838 molecules with predicted PCEs of 8% or larger were found. Moreover, statistical analysis of the candidate structures was performed with respect to both Morgan fingerprints and the building blocks, establishing a general architecture for nonfullerene acceptors.

Similarly, organic light-emitting diodes<sup>17</sup> (OLEDs) have found wide adoption in small displays, are becoming prevalent in screens and lighting applications, and are entering the market in flexible displays. Thermally activated delayed fluorescence (TADF) emitters have become the main OLED class because of their high quantum efficiency, operational stability, and low cost. Their essential property is a small energy gap between the first excited singlet and triplet states so that energetically favored but nonemissive, triplet excitons can be upconverted to emissive singlet excitons. Based on knowledge about the TADF mechanism, our group carried out HTVS of emitters covering 10<sup>6</sup> candidates (Figure 1).<sup>1</sup> Key methodology included efficient quantum chemistry, calibrated against experiment via supervised learning (*vide infra*). Linear regression and neural networks were used for property predictions across the entire space.

Exploration was performed iteratively using a neural network to predict the most promising candidates, which were then simulated, minimizing evaluations. Not only were known emitters rediscovered, but new structures were also uncovered. Additionally, the systematic exploration exposed both established property trade-offs and unknown property limits. Moreover, the best leads were evaluated by human experts concerning synthesizability and novelty. Consequently, the most promising molecules after both computer and human-based evaluations were synthesized and incorporated into devices leading to high external quantum efficiencies of over 20%. This study serves as a prototype for the entire data-driven discovery pipeline from defining the candidate space to device integration.

Finally, renewable energy like wind and solar is intermittent, requiring large storage capacities to meet consumer demands. Redox-flow batteries (RFBs) resolve that by separating energy from power, enabling large grids to store immense amounts of energy scalable to varying demand loads.<sup>18</sup> Organic RFBs<sup>19</sup> (ORFBs) represent a sensible advancement, as redox-active organic electrolytes are tunable and cheaper than inorganic alternatives.<sup>20</sup> To identify ideal organic electrolytes, our group performed HTVS of quinones, which are well-known for their single-electron redox pairs.<sup>21</sup> The screening spanned 1710

single- and double-electron redox pairs to validate existing studies and find new redox couples.

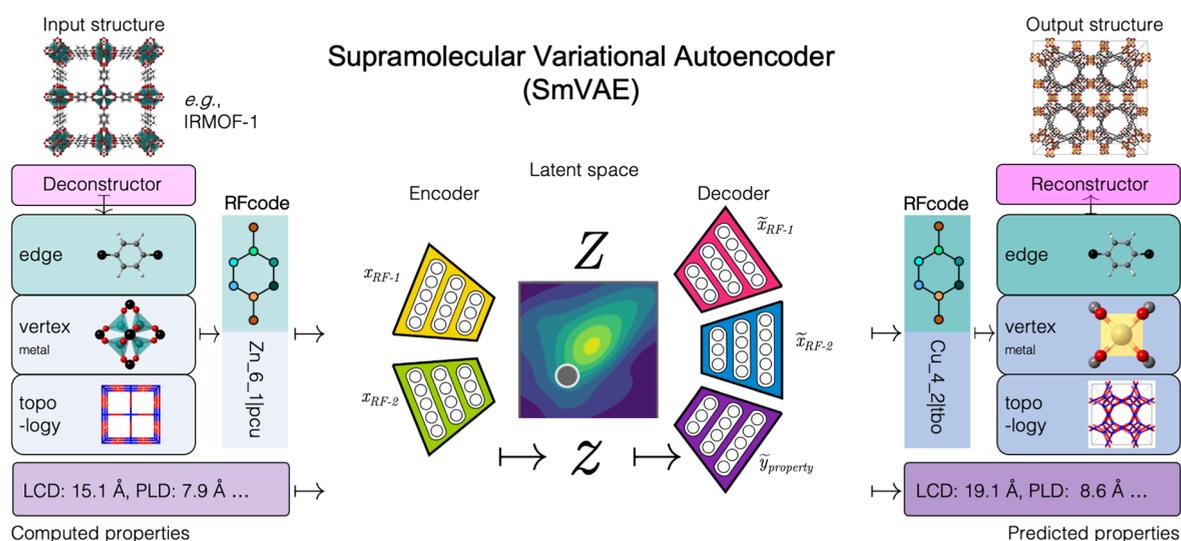
The results indicated that quinone-exclusive electrolytes were promising aqueous ORFBs and revealed that functionalizations near the carbonyl groups largely affected redox potential and those away largely affected solubility. Subsequently, several experimental studies verified these predictions.<sup>22,23</sup> However, decomposition was found to deteriorate battery capacity irreversibly.<sup>24</sup> Hence, our group performed combined computational and experimental studies on the decomposition of quinones in aqueous environments.<sup>18</sup> HTVS was performed for over 140 000 redox pairs, including decomposition product analysis. The results identified a trade-off between redox potential, with a maximum near 0.95 V, close to experimental results at 0.85 V,<sup>25</sup> and stability. These results provide roadmaps for future studies, which are ongoing in our group, as the trade-off suggests that electrolyte stability must be considered.

### Crystalline Materials

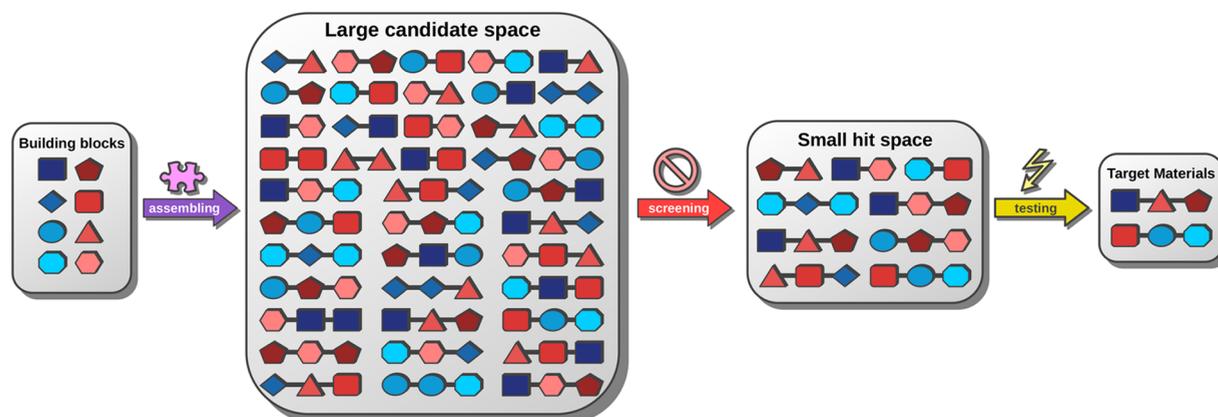
Crystalline energy storage materials with high energy density at low cost are cornerstones of renewable energy applications. For instance, multivalent calcium ion batteries<sup>26</sup> (CIBs) improve upon monovalent lithium-ion counterparts through increased capacities and higher material abundance while maintaining comparable operating voltages.<sup>27</sup> However, the development of CIBs is hindered by the failure of traditional graphite and calcium metal anodes due to the intercalation difficulty and the lack of efficient electrolytes. Recently, a high voltage (4.45 V) CIB cell using tin as the anode was reported to achieve a remarkable cyclability (over 300 cycles).<sup>28</sup>

Importantly, designing CIB anodes with improved performance requires a thorough exploration of the alloying space as calcium mixes with many elements. Hence, our group constructed a workflow to discover novel multivalent CIBs.<sup>29</sup> First, the tin electrochemical calcination reaction was investigated computationally and the reaction driving force as a function of calcium content was simulated. This exploration allowed the identification of threshold voltages governing the calcination limits. Consequently, a four-step screening strategy was adopted to look for high-performance CIB anodes. First, 357 metal–calcium binary and ternary compounds were identified from the Inorganic Crystal Structure Database (ICSD)<sup>30</sup> and further filtered to 115 candidates with existing decalciated metal/metalloid or binary intermetallic compounds. The calcination voltage profiles were calculated, and two threshold calcination voltages were defined, one stricter, based on the tin–calcium system, and the other more relaxed to account for potential differences in the driving force requirements. For each threshold, the maximum capacities, output voltages, volume expansions, and energy densities of the respective material were determined. Finally, metal–calcium systems with higher energy density than tin–calcium were identified, in which metalloids (Si, As, Sb, Ge), post-transition metals (Al, Pb, Cu, Cd, CdCu<sub>2</sub>, Ga, Bi, In, Tl, Hg), and noble metals (Ag, Pt, Pd, Au) showed promise as alloying candidates for CIB anodes and calls for further experimental validations.

Additionally, reticular frameworks<sup>31</sup> (RFs), which include metal–organic frameworks (MOFs), are crystalline porous materials with high internal surface area and high stability and can be used for gas storage, gas separation, and electrochemical energy storage. They are constructed via self-assembly of molecular building blocks and exhibit a near-infinite combinatorial space, complicating their systematic exploration.



**Figure 2.** Automated reticular framework (RF) discovery platform using the supramolecular variational autoencoder (SmVAE). We construct the intermediate representation, RFcode, using unique, decomposed nets as a tuple of edges, vertices, and topologies. We consider the edges as SMILES, while vertices and topologies are categorical variables from known structures. SmVAE is a multicomponent variational autoencoder encoding and decoding each part of the RFcode separately ( $x_{\text{edge}} \rightarrow \tilde{x}_{\text{edge}}$ ,  $x_{\text{RFcom}} \rightarrow \tilde{x}_{\text{RFcom}}$ ). Structures are converted into/back from RFcode using the deconstructor/reconstructor, then transferred into continuous vectors ( $z$ ). To organize the latent space based on properties, we add a supervised model to predict properties ( $\tilde{y}_{\text{property}}$ ) based on labeled data ( $y$ ). Data from ref 2.



**Figure 3.** High-throughput virtual screening starts from a large space of candidates (e.g., generated combinatorially, as illustrated). Using virtual screening, most candidates are eliminated, such that fewer (more expensive and time-consuming) experimental tests can be performed.

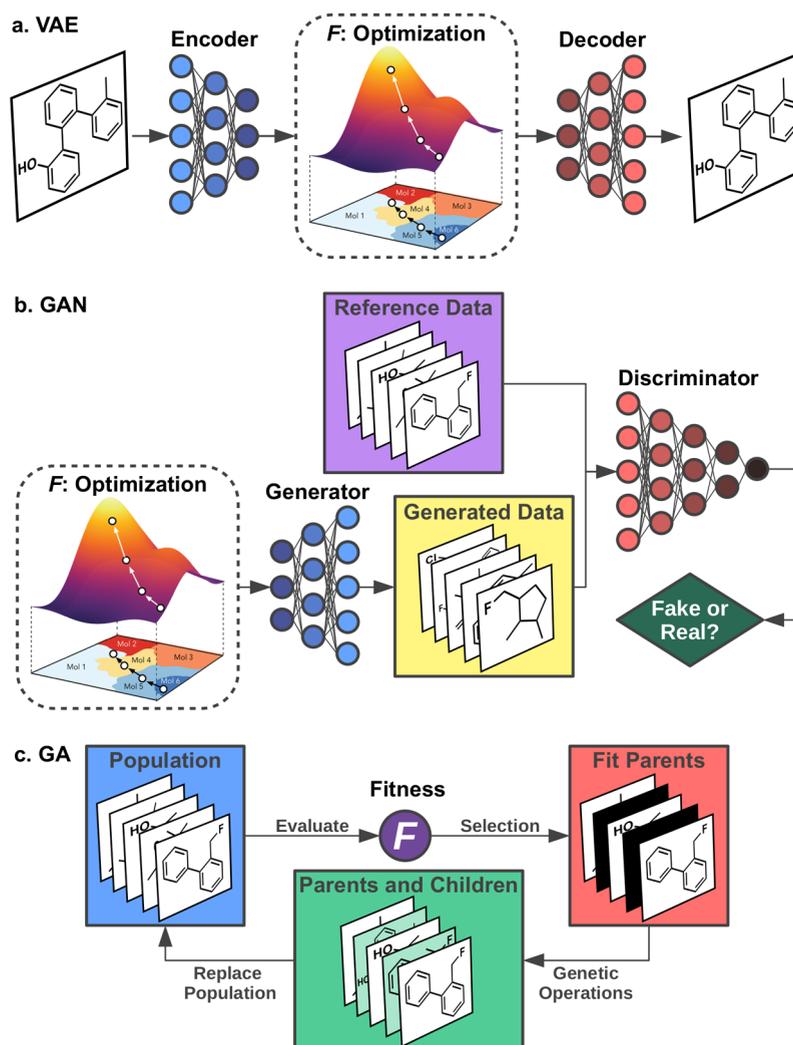
242 Recently, our group developed an invertible and efficient RF  
 243 representation (Figure 2).<sup>2,32</sup> MOF fragments were extracted  
 244 from the computation-ready, experimental (CoRE) MOF  
 245 database<sup>33</sup> and augmented randomly with common functional  
 246 groups. Furthermore, we added sets of multiconnected metal or  
 247 organic nodes and sets of known MOF topologies generating a  
 248 data set with around  $2 \times 10^6$  MOF structures. Moreover,  
 249 property simulations were performed for a random subset of  
 250 about 40 000 MOF structures. The supramolecular variational  
 251 autoencoder (SmVAE) with a MOF structure encoder-decoder,  
 252 property prediction model, and framework generation algorithm  
 253 was constructed with these structures (Figure 2), which can  
 254 locate high performing MOFs through property optimization in  
 255 the latent space. We demonstrated its capabilities for automatic  
 256 design by proposing top candidates for gas separation adsorbent  
 257 materials. We believe that the MOFs discovered are highly  
 258 competitive against the best-performing MOFs/zeolites ever  
 259 reported. Currently, their performance was validated using  
 260 computational methods. Nevertheless, experimental verification

is under way. Furthermore, the as-built platform can be applied  
 to various supramolecular systems (e.g., covalent-organic  
 frameworks, coordination polymers, etc.) and applications  
 (e.g., batteries, catalysis, drug delivery).

## METHODOLOGY

### High-Throughput Virtual Screening (HTVS)

Virtual screening<sup>34</sup> denotes a selection process of candidate  
 materials. Chemicals, either generated on-the-fly or from  
 databases, are subject to simulations that estimate application-  
 specific properties. Candidates failing computational tests are  
 rejected, with the proviso that predicted performance is likely  
 translatable to experimental performance. Thus, HTVS is a  
 technique that reduces large candidate spaces to a manageable  
 set of promising materials (Figure 3). In our search for new  
 TADF emitters (*vide supra*),<sup>1</sup> the candidate space was narrowed  
 down by 5 orders of magnitude via HTVS. Importantly, HTVS  
 on large chemical spaces is inverse molecular design (*vide infra*)  
 because, rather than designing structures directly, the computa-



**Figure 4.** Inverse molecular design based on desired properties ( $F$ ), with variational autoencoders (VAEs, a), generative adversarial networks (GANs, b), and genetic algorithms (GAs, c). Adapted with permission from ref 44. Copyright 2018 American Chemical Society.

279 tional tests and the candidate space are designed, which leads to  
 280 the final hits based on the predicted properties.<sup>35</sup> Moreover, it  
 281 can provide the basis for both generative and supervised models  
 282 (*vide infra*), as they all rely on validated data.

283 Accordingly, HTVS is a powerful accelerator because  
 284 computer simulation can be significantly less expensive than  
 285 the respective experiments.<sup>34</sup> The continuing growth in  
 286 computational power, which will soon reach the exascale, has  
 287 made virtual screening highly scalable as it is embarrassingly  
 288 parallel. Although HTVS is at least almost 20 years old,<sup>36</sup> it only  
 289 recently started transforming materials science by advances in  
 290 the accuracy and efficiency of density functional theory  
 291 (DFT).<sup>37</sup> Besides computational cost, the main appeal of DFT  
 292 was the possibility to tailor functional parameters to reproduce  
 293 experiments, which increased its predictive power significantly.

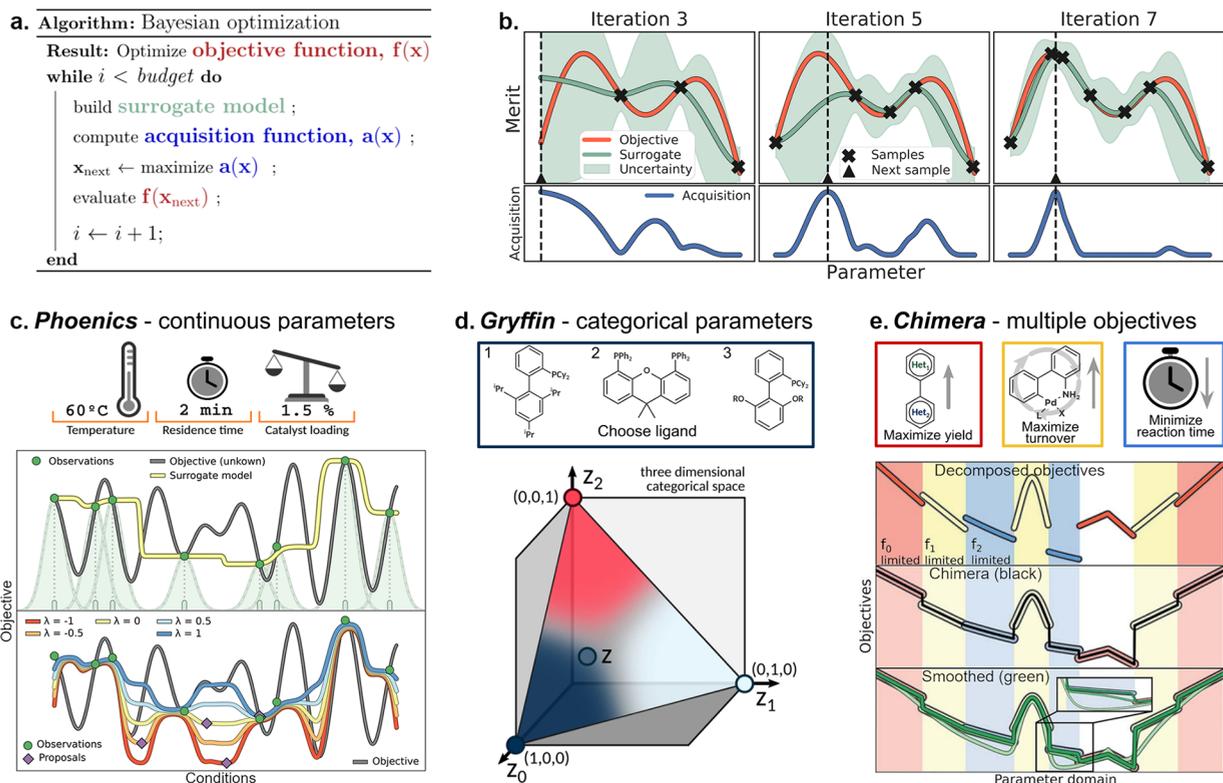
294 For instance, linear response time-dependent DFT (TD-  
 295 DFT) is accurate and computationally inexpensive for excited  
 296 state properties. More importantly, it is robust, can be used in a  
 297 black-box manner, and is readily deployed in simulations of tens  
 298 of thousands of molecules with minimal failure rates.<sup>14</sup>  
 299 However, one pernicious failure mode of TD-DFT is the  
 300 description of excited states with significant double-excitation  
 character, which is, *inter alia*, important in describing molecules

with inverted singlet–triplet gaps,<sup>38,39</sup> such as the INVEST  
 emitters recently described by our group.<sup>40</sup> Nevertheless, as  
 computing power is increasing, more sophisticated *ab initio*  
 approaches can be used in HTVS, allowing one to tackle ever  
 more complicated problems and new material classes.

Yet, the impact of HTVS has been hampered by the difficulty  
 in scaling the experimental confirmation of candidates,<sup>1</sup> as  
 simulations feasible for high-throughput are still largely  
 qualitative for condensed-phase properties.<sup>41</sup> A loose screen  
 that accounts for computational inaccuracies minimizes false  
 negatives, but the high cost of experimental validation means  
 that almost all candidates must be rejected. The accuracy of  
 computational screening can be maximized by implementing  
 self-correcting filters such as checking whether simulations  
 showed proper convergence catching false positives early on in  
 the workflow. Nevertheless, ultimately, improvements in the  
 experimental throughput are essential, calling for self-driving  
 laboratories and closed-loop experimentation.<sup>42,43</sup>

### AI-Powered Inverse Molecular Design

Inverse molecular design<sup>35</sup> starts at the desired properties and  
 explores the chemical space to identify molecules optimizing  
 them. Recently, various ML techniques have been employed to  
 improve inverse molecular design, motivated by advances both



**Figure 5.** (a) General pseudocode for Bayesian optimization. (b) Visualization of Bayesian optimization of an objective function (red curve) using Gaussian processes. (c) Examples of continuous-valued parameters compatible with *Phoenix*, along with a sample surrogate model and acquisition functions generated by the algorithm. Adapted with permission from ref 4. Copyright 2018 American Chemical Society. (d) Depiction of the representation of a categorical variable in *Gryffin* with three options (e.g., three ligands) on a simplex.<sup>51</sup> (e) Example of a multiobjective optimization problem for a chemical reaction, along with the construction of *Chimera* (bottom panel) from three 1-dimensional objective functions. Reproduced with permission from ref 52. Copyright 2018 Royal Society of Chemistry.

325 on the algorithmic (powerful ML libraries) and the hardware  
 326 sides (GPU improvements for large neural networks).  
 327 Importantly, inverse molecular design approaches can be  
 328 separated roughly into two classes: model-based ML algorithms  
 329 and evolutionary techniques.

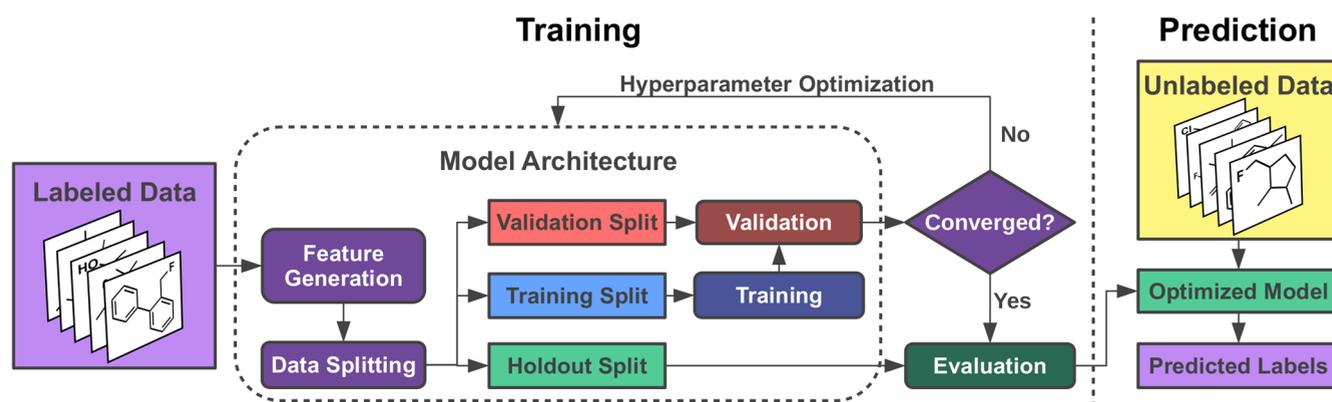
330 Model-based ML algorithms for inverse design models use  
 331 neural networks to learn patterns in molecular structures from  
 332 existing data. After training, these models suggest new molecules  
 333 covering important chemical features from the data set. Several  
 334 methodologies exist. Herein we will discuss variational  
 335 autoencoders (VAEs) and generative adversarial networks  
 336 (GANs) because our group, to the best of our knowledge, was  
 337 the first to apply these tools in chemistry. VAEs (Figure 4a) are  
 338 capable of forming continuous (latent) spaces from discrete  
 339 representations. They are trained to minimize the combined  
 340 losses of latent space smoothness and input reconstruction  
 341 enabling gradient-based optimization in the latent space. For  
 342 inverse design, the latent space of VAEs is coupled with a  
 343 property estimation model using supervised learning (*vide*  
 344 *infra*).<sup>44</sup> Consequently, the latent space is arranged based on the  
 345 property values allowing for a direct search of desired materials.  
 346 GANs (Figure 4b) are generative models with joint training of  
 347 two competing networks, a generator, and a discriminator. The  
 348 generator produces examples from a high dimensional (often  
 349 Gaussian) space, attempting to fool the discriminator, which  
 350 tries to distinguish generated samples from reference structures.  
 351 For molecules, our group proposed a sequential GAN  
 352 (ORGAN), where the model is trained using reinforcement

learning.<sup>45</sup> Desired molecular properties are used as a reward for  
 generating good structures.

Notably, both VAEs and GANs are trained in a supervised  
 way. Hence, they rely on existing data and mimic their  
 distribution. Thus, they are limited in the exploration of the  
 chemical space as compared to evolutionary techniques such as  
 genetic algorithms (GAs, cf. Figure 4c). As its name implies, GAs  
 are inspired by natural evolution. An initial population seeds the  
 algorithm, each member being evaluated. The top-performing  
 members proceed to the next iteration, the worst members are  
 removed or replaced by better offspring. For inverse molecular  
 design, the fitness function corresponds to the determination of  
 desired molecular properties.

In contrast to deep learning-based models, GAs are not biased  
 by user-defined data sets. Therefore, they are superior in  
 unbiased explorations.<sup>3</sup> Recently, we have shown that GAs  
 augmented with neural networks to estimate the similarity of a  
 molecule with a given data set can explore specific structural  
 classes without the large data requirements of GANs and VAEs.  
 Additionally, neural network-based learning was used to detect  
 and avoid local minima trapping the GA to amplify exploration  
 by avoiding convergence.<sup>3</sup> Notably, this shows that ML-based  
 inverse design techniques can be effectively combined with  
 evolutionary algorithms.

Importantly, in all these approaches, molecular representation  
 plays a crucial role. Molecular graphs are used for computational  
 efficiency, as they avoid conformations. Simplified Molecular  
 Input Line Entry System (SMILES)<sup>46</sup> strings are commonly  
 used as a flat encoding of molecular graphs. However, they have



**Figure 6.** Workflow for supervised learning of molecular properties. A known (labeled) data set is used to optimize a model, which is subsequently used to estimate molecular properties for an unknown (unlabeled) data set.

382 a complex structure making a large fraction of molecules  
 383 decoded from arbitrary SMILES invalid. This problem was  
 384 solved recently by our group in a fundamental way by replacing  
 385 SMILES with SELFIES (Self-Referencing Embedded Strings),<sup>47</sup>  
 386 which is available on GitHub.<sup>48</sup> SELFIES is a 100% valid  
 387 molecular string representation suitable as input for any inverse-  
 388 design algorithm that outperformed alternative approaches in  
 389 many benchmarks, such as validity and diversity of generated  
 390 molecules, molecular density in the latent space of VAEs, or  
 391 molecular optimization tasks with GAs.<sup>3</sup>

### 392 Bayesian Optimization

393 Several tasks across chemistry can be framed as optimization  
 394 problems, where controllable parameters optimizing a desired  
 395 objective are sought. For materials, such optimizations are  
 396 challenging, as they are typically high-dimensional, nonconvex,  
 397 and subject to noise and the objectives are expensive to evaluate.  
 398 Suitable optimization strategies ought to be sample-efficient,  
 399 global, and noise-tolerant. That is, they need to identify optimal  
 400 parameter choices with as few measurements as possible, be able  
 401 to escape local minima, and mitigate the detrimental effect of  
 402 noise. A plethora of experiment planning strategies for  
 403 optimization are currently available,<sup>49</sup> from traditional design  
 404 of experiment to evolutionary and heuristic approaches. Among  
 405 these, Bayesian optimization<sup>50</sup> (BO) has emerged as the strategy  
 406 that best meets these requirements.

407 BO is an experiment planning algorithm that, in contrast to  
 408 most other approaches, uses an ML model to learn from  
 409 previous observations before suggesting the next iteration  
 410 (Figure 5a).<sup>50</sup> In its most widely adopted form, BO employs  
 411 techniques such as Gaussian processes to build a *surrogate* model  
 412 that captures the features of the underlying objective function.  
 413 Based on this surrogate, an *acquisition* function is defined, which  
 414 determines the strategy used to propose new experiments  
 415 (Figure 5b). Just like BO formulations using different ML  
 416 models exist, various acquisition functions have been developed.  
 417 Due to the use of an ML model, BO is sample-efficient. It is also  
 418 noise-tolerant, as these models explicitly account for it. Finally,  
 419 BO is a global approach that balances the *exploitation* of the best  
 420 local optima identified with the *exploration* of unprobed areas of  
 421 parameter space.

422 Typical BO approaches are inherently sequential and require  
 423 heavy computations for each iteration. Therefore, BO can be  
 424 unduly expensive when used in conjunction with high-  
 425 throughput evaluations. Thus, our group has developed *Phoenix*  
 426 (Figure 5c), a linear-scaling BO approach that supports parallel

427 experiments.<sup>4</sup> *Phoenix* employs Bayesian neural networks 427  
 (BNNs) to build a kernel density estimate of the objective 428  
 function, and its acquisition function allows to select batches of 429  
 evaluations to be run in parallel. Importantly, *Phoenix* is suitable 430  
 for the optimization of continuous parameters, such as 431  
 temperature and concentration. To also optimize categorical 432  
 parameters, such as the choice of solvent, we developed *Gryffin* 433  
 (Figure 5d), which uses categorical kernel densities that can be 434  
 relaxed to continuous ones.<sup>51</sup> In addition, *Gryffin* allows for 435  
 expert knowledge, in the form of descriptors for each categorical 436  
 choice, to be provided to improve the optimization efficiency. 437  
 Often, multiple competing objectives are present in materials 438  
 science. *Chimera* (Figure 5e) is a general-purpose approach to 439  
 multiobjective optimization.<sup>52</sup> It allows defining a hierarchy of 440  
 objective preferences, which are combined into a single function 441  
 to be optimized with any algorithm of choice. 442

443 Importantly, all the aforementioned algorithms can be 443  
 combined with automated laboratories to enable autonomous 444  
 experimentation.<sup>42</sup> These self-driving platforms are able to 445  
 execute closed-loop workflows for the self-optimization of 446  
 materials and processes. However, this requires robust software 447  
 connections between automated hardware and experiment 448  
 planning methods. *ChemOS* is a flexible, modular, open source 449  
 and portable *Python* package that provides this interface between 450  
 experiment planning and automated experiments.<sup>53,54</sup> Accord- 451  
 ingly, in our laboratory, we have deployed *ChemOS*, together 452  
 with *Phoenix*, *Gryffin*, and *Chimera*, for the autonomous 453  
 optimization of manufacturing processes of thin-film materi- 454  
 als,<sup>55</sup> multicomponent polymer OPV blends,<sup>56</sup> and reaction 455  
 conditions of stereoselective Suzuki coupling.<sup>57</sup> 456

### 457 Supervised Learning

458 The costs associated with property measurement, from both 458  
 experiments and simulations, are a major obstacle to the 459  
 widespread expansion of HTVS, optimization, and inverse 460  
 design. All of these techniques require some form of data 461  
 acquisition, *i.e.*, simulations, measurements, or data mining. 462  
 However, adapting experimental design to suit the needs of 463  
 automated protocols is challenging, despite self-driving 464  
 approaches likely being overall cost-effective. The promise of 465  
 accurate and practically free inference of new results from 466  
 existing data via supervised learning is a major driver of the 467  
 ongoing ML revolution in the physical sciences.<sup>58</sup> 468

469 Supervised learning requires a data set of features and labels.<sup>59</sup>  
 470 For molecular property prediction, this data set contains 470  
 molecules in a specific representation (features) and their 471

472 corresponding properties (labels). First, the data set is split into  
473 three, training, validation and holdout sets. The model is trained  
474 stepwise on the training set, usually by gradient descent or  
475 related algorithms. In general, hyperparameters, *i.e.*, choice of  
476 features, training set, and model architecture, influence  
477 predictive performance. These hyperparameters are optimized  
478 by maximizing prediction accuracy on the validation set.  
479 Eventually, model performance is evaluated via prediction  
480 accuracy for the holdout set, and the final model can be used to  
481 predict properties for unlabeled molecules. The entire workflow  
482 is illustrated in Figure 6. Our group developed several model  
483 architectures for supervised learning of molecular properties,  
484 most notably graph convolutional neural networks.<sup>60,61</sup>

485 Importantly, supervised learning has been used successfully  
486 for materials discovery. For example, our group used the CEP  
487 data set for property prediction.<sup>62</sup> After training on more than  
488 200 000 molecules, a neural network predicted the result of DFT  
489 calculations consistently at a fraction of the computational  
490 expense. Additionally, our group applied this approach to reduce  
491 the number of simulations in HTVS significantly, with training  
492 on a set of similar size.<sup>1</sup> Moreover, our group also used Gaussian  
493 process regression to calibrate for systematic errors in DFT.<sup>16</sup>  
494 Crucially, in these studies, ML algorithms, representations,  
495 acquisition of training data, and validation procedures for  
496 models were tightly integrated with an understanding of the  
497 problem space, as opposed to sole reliance on existing data from  
498 various sources. We believe these considerations are key when it  
499 comes to the practical application of ML in chemistry.

500 Moreover, fruitful applications of supervised learning in  
501 materials science start from well-defined scientific goals. In  
502 contrast, the excitement brought upon by ML has generated  
503 many studies that focus on learning performance rather than  
504 scientific objectives. Generally, this is based on the (debatable  
505 and often unsupported) idea that performance metrics on one  
506 data set are transferable to other data sets or related problems.  
507 However, ML algorithms are highly parametrized and thus can  
508 readily overfit.<sup>63</sup> Indeed, the model choice can itself become a  
509 form of overfitting, especially when done on performance  
510 considerations alone.<sup>64</sup> Moreover, training data bias can  
511 contaminate predictions<sup>65</sup> but accounting for these biases  
512 appropriately is problem-specific. Furthermore, many studies  
513 are focused on error estimates obtained from statistical measures  
514 such as cross-validation. Although validation error can be a  
515 useful guide to the true prediction error on new data, it is not a  
516 replacement for it<sup>66</sup> and is often too optimistic.<sup>67</sup> In many ways,  
517 these issues arise when focus on the scientific goals is lost, as  
518 ultimately the best test of supervised learning is whether it solves  
519 problems.

## 520 ■ CONCLUSION AND OUTLOOK

521 In this Account, we have reviewed data-driven approaches our  
522 group has employed for the design of materials, especially for  
523 clean energy applications, in the past decade. One of the first  
524 large scale campaigns our group embarked on was the CEP,  
525 where we implemented supervised learning together with HTVS  
526 using quantum chemistry simulations to investigate 10<sup>7</sup>  
527 potential donor molecules for organic solar cells and devised  
528 design principles by statistical analysis of structure–function  
529 relationships.<sup>12</sup> In the subsequent years, we refined these ML  
530 strategies and expanded our efforts toward other important  
531 materials such as OLEDs, OFRBs, multivalent CIBs, and RFs. In  
532 all these projects, data-driven workflows were key to speed up  
533 both the discovery and the design of new materials.

However, we believe that the full potential of data-driven  
strategies is yet to be unleashed. For instance, many properties  
are currently not investigated in HTVS because of their  
prohibitive computational cost. One such property is molecular  
stability with respect to common decomposition pathways. The  
associated problem is the huge dimensionality of potential  
reactions molecules can undergo, which greatly exceeds the  
chemical compound space in complexity. Recently, our group  
developed a method for the automatic discovery of chemical  
reactions based on the selection of reactive internal coordinates  
such as weak chemical bonds.<sup>68</sup> We believe this approach,  
together with empirical rules or heuristics for selecting reactive  
internal coordinates, could be used for HTVS of reactivity and  
stability of materials, and research in that direction is ongoing.  
Other properties too prohibitive for HTVS include the influence  
of explicit solvation on spectroscopic properties and the direct  
simulation of amorphous solid-state structures and properties.  
The main challenge therein is the large number of particles and  
degrees of freedom in the model systems and the associated  
multitude of interactions.

Furthermore, some of the methodologies we developed have  
only been tested on benchmark problems but are yet to be  
employed in real applications. Particularly, the genetic algorithm  
augmented with neural networks using SELFIES as molecular  
representation<sup>47</sup> our group proposed recently has outperformed  
most alternative generative models in benchmarks. However, it  
has yet to be implemented for designing functional materials,  
and we are actively working on that.<sup>3</sup> Finally, one of the most  
critical challenges of ML is model interpretability. Typically,  
supervised learning approaches are employed in a black box  
fashion without gaining insight into what the model actually  
learned. However, our group has shown recently that regression  
methods such as gradient boosting, when trained on molecular  
graph features, can be used to reveal important chemical  
moieties influencing the properties.<sup>69,70</sup> The trained model can  
be interpreted by human experts and rationalizing the feature  
importance can lead to new scientific understanding. We believe  
that similar approaches have the potential to change the way  
science is carried out in the near future.

However, the bottleneck of materials design campaigns is  
experimental synthesis and characterization, usually by a large  
margin.<sup>71</sup> Any material, no matter how good its (predicted)  
performance, needs to be synthesized for it to be used in real life.  
In particular for clean energy applications, material syntheses  
need to be performed on a huge scale requiring reliable, safe and  
green chemical processes. Accordingly, the continuing speed-up  
in computer power providing unprecedented prediction  
capabilities needs to be paralleled by increased experimental  
throughput. Accelerating materials design ultimately requires  
close integration of computer simulation, ML and experimenta-  
tion in self-driving platforms, which our group termed Materials  
Acceleration Platforms (MAPs).<sup>43</sup>

One essential feature of MAPs is a closed-loop materials  
discovery workflow incorporating experimentation, computa-  
tion, and human intuition. Online characterization techniques in  
conjunction with automated robotic synthesis<sup>72–74</sup> are central  
enabling technologies in these platforms. Making and measuring  
molecules on-demand in a feedback loop with self-correcting  
computational screening and ML is key to finding true “needle-  
in-a-haystack” materials. Currently, our group is implementing  
such a MAP for the realization of innovative materials making  
use of robust cross coupling chemistry, parallel robotic synthesis,  
and in-line characterization of spectroscopic properties coupled

597 with computer simulation and ML. Details of this implementa-  
598 tion will be described in an upcoming Account our group is  
599 working on in due course. Accordingly, the data-driven methods  
600 described above are a stepping stone to accelerate materials  
601 design. However, to realize their true potential, they need to  
602 percolate into experimental systems, and we are looking forward  
603 to witnessing applications of these methods in closed-loop  
604 experimental material design campaigns in the near future.

## 605 ■ AUTHOR INFORMATION

### 606 Corresponding Author

607 **Alán Aspuru-Guzik** – *Chemical Physics Theory Group,*  
608 *Department of Chemistry and Department of Computer*  
609 *Science, University of Toronto, Toronto, Ontario M5S 3H6,*  
610 *Canada; Vector Institute for Artificial Intelligence, Toronto,*  
611 *Ontario MSG 1M1, Canada; Lebovic Fellow, Canadian*  
612 *Institute for Advanced Research (CIFAR), Toronto, Ontario*  
613 *MSG, Canada; [orcid.org/0000-0002-8277-4434](https://orcid.org/0000-0002-8277-4434);*  
614 *Email: [aspuru@utoronto.ca](mailto:aspuru@utoronto.ca)*

### 615 Authors

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

620 **Gabriel dos Passos Gomes** – *Chemical Physics Theory Group,*  
621 *Department of Chemistry and Department of Computer*  
622 *Science, University of Toronto, Toronto, Ontario M5S 3H6,*  
623 *Canada; [orcid.org/0000-0002-8235-5969](https://orcid.org/0000-0002-8235-5969)*

624 **Matteo Aldeghi** – *Chemical Physics Theory Group, Department*  
625 *of Chemistry and Department of Computer Science, University*  
626 *of Toronto, Toronto, Ontario M5S 3H6, Canada; Vector*  
627 *Institute for Artificial Intelligence, Toronto, Ontario MSG*  
628 *1M1, Canada; [orcid.org/0000-0003-0019-8806](https://orcid.org/0000-0003-0019-8806)*

629 **Riley J. Hickman** – *Chemical Physics Theory Group,*  
630 *Department of Chemistry and Department of Computer*  
631 *Science, University of Toronto, Toronto, Ontario M5S 3H6,*  
632 *Canada*

633 **Mario Krenn** – *Chemical Physics Theory Group, Department of*  
634 *Chemistry and Department of Computer Science, University of*  
635 *Toronto, Toronto, Ontario M5S 3H6, Canada; Vector*  
636 *Institute for Artificial Intelligence, Toronto, Ontario MSG*  
637 *1M1, Canada*

638 **Cyrille Lavigne** – *Chemical Physics Theory Group, Department*  
639 *of Chemistry and Department of Computer Science, University*  
640 *of Toronto, Toronto, Ontario M5S 3H6, Canada;*  
641 *[orcid.org/0000-0003-2778-1866](https://orcid.org/0000-0003-2778-1866)*

642 **Michael Lindner-D'Addario** – *Chemical Physics Theory*  
643 *Group, Department of Chemistry and Department of*  
644 *Computer Science, University of Toronto, Toronto, Ontario*  
645 *M5S 3H6, Canada*

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

650 **Cher Tian Ser** – *Chemical Physics Theory Group, Department*  
651 *of Chemistry and Department of Computer Science, University*  
652 *of Toronto, Toronto, Ontario M5S 3H6, Canada*

653 **Zhenpeng Yao** – *Chemical Physics Theory Group, Department*  
654 *of Chemistry and Department of Computer Science, University*  
655 *of Toronto, Toronto, Ontario M5S 3H6, Canada;*  
656 *[orcid.org/0000-0001-8286-8257](https://orcid.org/0000-0001-8286-8257)*

Complete contact information is available at: 557  
<https://pubs.acs.org/10.1021/acs.accounts.0c00785> 558

### 559 Author Contributions 560

M.A., R.J.H., M.K., C.L., M.L.-D., A.K.N., C.T.S., and Z.Y. 561  
contributed equally to this work. R.P. and G.P.G. conceived the 562  
general outline and structure of this manuscript, and all authors 563  
contributed toward refining the structure. The manuscript was 564  
written through contributions of all authors. All authors have 565  
approved the final version of the manuscript.

### 566 Notes

The authors declare the following competing financial 567  
interest(s): A.A.-G. is co-founder and Chief Visionary Officer 568  
of Kebotix, Inc. 569

### 570 Biographies

**Robert Pollice** is an SNSF postdoctoral fellow at the University of 571  
Toronto. 572

**Gabriel dos Passos Gomes** is an NSERC Banting postdoctoral fellow 573  
at the University of Toronto. 574

**Matteo Aldeghi** is a postdoctoral fellow at the Vector Institute for 575  
Artificial Intelligence and the University of Toronto. 576

**Riley J. Hickman** is a PhD student at the University of Toronto. 577

**Mario Krenn** is an Erwin Schrödinger postdoctoral fellow at the 578  
University of Toronto and the Vector Institute for Artificial 579  
Intelligence. 580

**Cyrille Lavigne** is a postdoctoral fellow at the University of Toronto. 581

**Michael Lindner-D'Addario** is a PhD student at the University of 582  
Toronto. 583

**AkshatKumar Nigam** is a researcher at the University of Toronto. 584

**Cher Tian Ser** is a PhD student at the University of Toronto. 585

**Zhenpeng Yao** is a postdoctoral fellow at the University of Toronto. 586

**Alán Aspuru-Guzik** is a Professor of Chemistry and Computer Science 587  
at the University of Toronto, a Canada 150 Research Chair in 588  
Theoretical Chemistry, a Canada CIFAR AI Chair at the Vector 589  
Institute, a CIFAR Lebovic Fellow in the Biologically Inspired Solar 590  
Energy program, and a Google Industrial Research Chair in Quantum 591  
Computing. 592

## 593 ■ ACKNOWLEDGMENTS

We thank all our co-workers and collaborators who contributed 594  
to the projects highlighted in this account. R.P. acknowledges 595  
funding through a Postdoc.Mobility fellowship by the Swiss 596  
National Science Foundation (SNSF, Project No. 191127). 597  
G.P.G. gratefully acknowledges the Natural Sciences and 598  
Engineering Research Council of Canada (NSERC) for the 599  
Banting Postdoctoral Fellowship. R.J.H. gratefully acknowledges 600  
NSERC for provision of the Postgraduate Scholarships-Doctoral 601  
Program (PGSD3-534584-2019). M.K. acknowledges support 602  
from the Austrian Science Fund (FWF) through the Erwin 603  
Schrödinger fellowship No. J4309. M.L.-D. gratefully acknowl- 604  
edges the Fonds de Recherche Quebec Nature et Technologies 605  
(FRQNT) for the B1X Master's Scholarship. M.L.-D. also 606  
acknowledges support from the Queen Elizabeth II Graduate 607  
Scholarship in Science and Technology (QEII-GSST). We 608  
acknowledge the Defense Advanced Research Projects Agency 609  
(DARPA) under the Accelerated Molecular Discovery Program 610  
under Cooperative Agreement No. HR00111920027 dated 611

712 August 1, 2019. The content of the information presented in this  
713 work does not necessarily reflect the position or the policy of the  
714 Government. A.A.-G. thanks Anders G. Frøseth for his generous  
715 support. A.A.-G. also acknowledges the generous support of  
716 Natural Resources Canada and the Canada 150 Research Chairs  
717 program. We also acknowledge the Department of Navy award  
718 (N00014-19-1-2134) issued by the Office of Naval Research.  
719 The United States Government has a royalty-free license  
720 throughout the world in all copyrightable material contained  
721 herein. Any opinions, findings, and conclusions or recommen-  
722 dations expressed in this material are those of the authors and do  
723 not necessarily reflect the views of the Office of Naval Research.

## 724 ■ REFERENCES

- 725 (1) Gómez-Bombarelli, R.; Aguilera-Iparraguirre, J.; Hirzel, T. D.;  
726 Duvenaud, D.; Maclaurin, D.; Blood-Forsythe, M. A.; Chae, H. S.;  
727 Einzinger, M.; Ha, D.-G.; Wu, T.; Markopoulos, G.; Jeon, S.; Kang, H.;  
728 Miyazaki, H.; Numata, M.; Kim, S.; Huang, W.; Hong, S. I.; Baldo, M.;  
729 Adams, R. P.; Aspuru-Guzik, A. Design of Efficient Molecular Organic  
730 Light-Emitting Diodes by a High-Throughput Virtual Screening and  
731 Experimental Approach. *Nat. Mater.* **2016**, *15* (10), 1120–1127.
- 732 (2) Yao, Z.; Sanchez-Lengeling, B.; Bobbitt, N. S.; Bucior, B. J.;  
733 Kumar, S. G. H.; Collins, S. P.; Burns, T.; Woo, T. K.; Farha, O.; Snurr,  
734 R. Q.; Aspuru-Guzik, A. Inverse Design of Nanoporous Crystalline  
735 Reticular Materials with Deep Generative Models. *Nat. Mach. Intell.*  
736 **2021**, *3*, 76.
- 737 (3) Nigam, A.; Friederich, P.; Krenn, M.; Aspuru-Guzik, A.  
738 Augmenting Genetic Algorithms with Deep Neural Networks for  
739 Exploring the Chemical Space. In *International Conference on Learning*  
740 *Representations*; 2020.
- 741 (4) Häse, F.; Roch, L. M.; Kreisbeck, C.; Aspuru-Guzik, A. Phoenix:  
742 A Bayesian Optimizer for Chemistry. *ACS Cent. Sci.* **2018**, *4* (9), 1134–  
743 1145.
- 744 (5) Hey, T.; Tansley, S.; Tolle, K. *The Fourth Paradigm: Data-Intensive*  
745 *Scientific Discovery*; Microsoft Research: Redmond, WA, 2009.
- 746 (6) Schindler, S. Scientific Discovery: That-Whats and What-Thats.  
747 *Ergo, an Open Access Journal of Philosophy* **2015**, *2*, 123–148.
- 748 (7) Kuhn, T. S. Historical Structure of Scientific Discovery. *Science*  
749 **1962**, *136* (3518), 760–764.
- 750 (8) March, S. T.; Smith, G. F. Design and Natural Science Research on  
751 Information Technology. *Decision Support Systems* **1995**, *15* (4), 251–  
752 266.
- 753 (9) Ostroverkhova, O. Organic Optoelectronic Materials: Mecha-  
754 nisms and Applications. *Chem. Rev.* **2016**, *116* (22), 13279–13412.
- 755 (10) Hedley, G. J.; Ruseckas, A.; Samuel, I. D. W. Light Harvesting for  
756 Organic Photovoltaics. *Chem. Rev.* **2017**, *117* (2), 796–837.
- 757 (11) Liu, Q.; Jiang, Y.; Jin, K.; Qin, J.; Xu, J.; Li, W.; Xiong, J.; Liu, J.;  
758 Xiao, Z.; Sun, K.; Yang, S.; Zhang, X.; Ding, L. 18% Efficiency Organic  
759 Solar Cells. *Science Bulletin* **2020**, *65* (4), 272–275.
- 760 (12) Hachmann, J.; Olivares-Amaya, R.; Atahan-Evrenk, S.; Amador-  
761 Bedolla, C.; Sánchez-Carrera, R. S.; Gold-Parker, A.; Vogt, L.;  
762 Brockway, A. M.; Aspuru-Guzik, A. The Harvard Clean Energy Project:  
763 Large-Scale Computational Screening and Design of Organic Photo-  
764 voltaics on the World Community Grid. *J. Phys. Chem. Lett.* **2011**, *2*  
765 (17), 2241–2251.
- 766 (13) Olivares-Amaya, R.; Amador-Bedolla, C.; Hachmann, J.; Atahan-  
767 Evrenk, S.; Sánchez-Carrera, R. S.; Vogt, L.; Aspuru-Guzik, A.  
768 Accelerated Computational Discovery of High-Performance Materials  
769 for Organic Photovoltaics by Means of Cheminformatics. *Energy*  
770 *Environ. Sci.* **2011**, *4* (12), 4849–4861.
- 771 (14) Hachmann, J.; Olivares-Amaya, R.; Jinich, A.; Appleton, A. L.;  
772 Blood-Forsythe, M. A.; Seress, L. R.; Román-Salgado, C.; Trepte, K.;  
773 Atahan-Evrenk, S.; Er, S.; Shrestha, S.; Mondal, R.; Sokolov, A.; Bao, Z.;  
774 Aspuru-Guzik, A. Lead Candidates for High-Performance Organic  
775 Photovoltaics from High-Throughput Quantum Chemistry – the  
776 Harvard Clean Energy Project. *Energy Environ. Sci.* **2014**, *7* (2), 698–  
777 704.
- (15) Sokolov, A. N.; Atahan-Evrenk, S.; Mondal, R.; Akkerman, H. B.;  
778 Sánchez-Carrera, R. S.; Granados-Focil, S.; Schrier, J.; Mannsfeld, S. C.  
779 B.; Zoombelt, A. P.; Bao, Z.; Aspuru-Guzik, A. From Computational  
780 Discovery to Experimental Characterization of a High Hole Mobility  
781 Organic Crystal. *Nat. Commun.* **2011**, *2* (1), 437.
- (16) Lopez, S. A.; Sanchez-Lengeling, B.; de Goes Soares, J.; Aspuru-  
783 Guzik, A. Design Principles and Top Non-Fullerene Acceptor  
784 Candidates for Organic Photovoltaics. *Joule* **2017**, *1* (4), 857–870.
- (17) Zou, S.-J.; Shen, Y.; Xie, F.-M.; Chen, J.-D.; Li, Y.-Q.; Tang, J.-X.  
786 Recent Advances in Organic Light-Emitting Diodes: Toward Smart  
787 Lighting and Displays. *Mater. Chem. Front.* **2020**, *4* (3), 788–820.
- (18) Tabor, D. P.; Gomez-Bombarelli, R.; Tong, L.; Gordon, R. G.;  
789 Aziz, M. J.; Aspuru-Guzik, A. Mapping the Frontiers of Quinone  
790 Stability in Aqueous Media: Implications for Organic Aqueous Redox  
791 Flow Batteries. *J. Mater. Chem. A* **2019**, *7* (20), 12833–12841.
- (19) Luo, J.; Hu, B.; Zhao, Y.; Liu, T. L. Status and Prospects  
793 of Organic Redox Flow Batteries toward Sustainable Energy Storage.  
794 *ACS Energy Lett.* **2019**, *4* (9), 2220–2240.
- (20) Lin, K.; Gómez-Bombarelli, R.; Beh, E. S.; Tong, L.; Chen, Q.;  
796 Valle, A.; Aspuru-Guzik, A.; Aziz, M. J.; Gordon, R. G. A Redox-Flow  
797 Battery with an Alloxazine-Based Organic Electrolyte. *Nature Energy*  
798 **2016**, *1* (9), 1–8.
- (21) Er, S.; Suh, C.; Marshak, M. P.; Aspuru-Guzik, A. Computational  
800 Design of Molecules for an All-Quinone Redox Flow Battery. *Chemical*  
801 *Science* **2015**, *6* (2), 885–893.
- (22) Yang, Z.; Tong, L.; Tabor, D. P.; Beh, E. S.; Goulet, M.-A.; De  
803 Porcellinis, D.; Aspuru-Guzik, A.; Gordon, R. G.; Aziz, M. J. Alkaline  
804 Benzoquinone Aqueous Flow Battery for Large-Scale Storage of  
805 Electrical Energy. *Adv. Energy Mater.* **2018**, *8* (8), 1702056.
- (23) Kwabi, D. G.; Lin, K.; Ji, Y.; Kerr, E. F.; Goulet, M.-A.; De  
807 Porcellinis, D.; Tabor, D. P.; Pollack, D. A.; Aspuru-Guzik, A.; Gordon,  
808 R. G.; Aziz, M. J. Alkaline Quinone Flow Battery with Long Lifetime at  
809 PH 12. *Joule* **2018**, *2* (9), 1894–1906.
- (24) Goulet, M.-A.; Tong, L.; Pollack, D. A.; Tabor, D. P.; Odom, S.  
811 A.; Aspuru-Guzik, A.; Kwan, E. E.; Gordon, R. G.; Aziz, M. J. Extending  
812 the Lifetime of Organic Flow Batteries via Redox State Management. *J.*  
813 *Am. Chem. Soc.* **2019**, *141* (20), 8014–8019.
- (25) Hooper-Burkhardt, L.; Krishnamoorthy, S.; Yang, B.; Murali, A.;  
815 Nirmalchandar, A.; Prakash, G. K. S.; Narayanan, S. R. A New Michael-  
816 Reaction-Resistant Benzoquinone for Aqueous Organic Redox Flow  
817 Batteries. *J. Electrochem. Soc.* **2017**, *164* (4), A600.
- (26) Arroyo-de Dompablo, M. E.; Ponrouch, A.; Johansson, P.;  
819 Palacín, M. R. Achievements, Challenges, and Prospects of Calcium  
820 Batteries. *Chem. Rev.* **2020**, *120* (14), 6331–6357.
- (27) Ponrouch, A.; Frontera, C.; Bardé, F.; Palacín, M. R. Towards a  
822 Calcium-Based Rechargeable Battery. *Nat. Mater.* **2016**, *15* (2), 169–  
823 172.
- (28) Wang, M.; Jiang, C.; Zhang, S.; Song, X.; Tang, Y.; Cheng, H.-M.  
825 Reversible Calcium Alloying Enables a Practical Room-Temperature  
826 Rechargeable Calcium-Ion Battery with a High Discharge Voltage. *Nat.*  
827 *Chem.* **2018**, *10* (6), 667–672.
- (29) Yao, Z.; Hegde, V. I.; Aspuru-Guzik, A.; Wolverton, C. Discovery  
829 of Calcium-Metal Alloy Anodes for Reversible Ca-Ion Batteries. *Adv.*  
830 *Energy Mater.* **2019**, *9* (9), 1802994.
- (30) Belsky, A.; Hellenbrandt, M.; Karen, V. L.; Luksch, P. New  
832 Developments in the Inorganic Crystal Structure Database (ICSD):  
833 Accessibility in Support of Materials Research and Design. *Acta*  
834 *Crystallogr., Sect. B: Struct. Sci.* **2002**, *58* (3–1), 364–369.
- (31) Lyu, H.; Ji, Z.; Wuttke, S.; Yaghi, O. M. Digital Reticular  
836 Chemistry. *Chem.* **2020**, *6* (9), 2219–2241.
- (32) Bucior, B. J.; Rosen, A. S.; Haranczyk, M.; Yao, Z.; Ziebel, M. E.;  
837 Farha, O. K.; Hupp, J. T.; Siepmann, J. I.; Aspuru-Guzik, A.; Snurr, R. Q.  
839 Identification Schemes for Metal–Organic Frameworks To Enable  
840 Rapid Search and Cheminformatics Analysis. *Cryst. Growth Des.* **2019**,  
841 *19* (11), 6682–6697.
- (33) Chung, Y. G.; Camp, J.; Haranczyk, M.; Sikora, B. J.; Bury, W.;  
843 Krungleviciute, V.; Yildirim, T.; Farha, O. K.; Sholl, D. S.; Snurr, R. Q.  
844 Computation-Ready, Experimental Metal–Organic Frameworks: A 845

- 846 Tool To Enable High-Throughput Screening of Nanoporous Crystals. *Chem. Mater.* **2014**, *26* (21), 6185–6192.
- 848 (34) Pyzer-Knapp, E. O.; Suh, C.; Gómez-Bombarelli, R.; Aguilera-  
849 Iparraguirre, J.; Aspuru-Guzik, A. What Is High-Throughput Virtual  
850 Screening? A Perspective from Organic Materials Discovery. *Annu. Rev.*  
851 *Mater. Res.* **2015**, *45* (1), 195–216.
- 852 (35) Sanchez-Lengeling, B.; Aspuru-Guzik, A. Inverse Molecular  
853 Design Using Machine Learning: Generative Models for Matter  
854 Engineering. *Science* **2018**, *361* (6400), 360–365.
- 855 (36) Schapira, M.; Raaka, B. M.; Das, S.; Fan, L.; Totrov, M.; Zhou, Z.;  
856 Wilson, S. R.; Abagyan, R.; Samuels, H. H. Discovery of Diverse  
857 Thyroid Hormone Receptor Antagonists by High-Throughput  
858 Docking. *Proc. Natl. Acad. Sci. U. S. A.* **2003**, *100* (12), 7354–7359.
- 859 (37) Ceder, G.; Chiang, Y.-M.; Sadoway, D. R.; Aydinol, M. K.; Jang,  
860 Y.-I.; Huang, B. Identification of Cathode Materials for Lithium  
861 Batteries Guided by First-Principles Calculations. *Nature* **1998**, *392*  
862 (6677), 694–696.
- 863 (38) de Silva, P. Inverted Singlet–Triplet Gaps and Their Relevance  
864 to Thermally Activated Delayed Fluorescence. *J. Phys. Chem. Lett.* **2019**,  
865 *10* (18), 5674–5679.
- 866 (39) Ehrmaier, J.; Rabe, E. J.; Pristash, S. R.; Corp, K. L.; Schlenker, C.  
867 W.; Sobolewski, A. L.; Domcke, W. Singlet–Triplet Inversion in  
868 Heptazine and in Polymeric Carbon Nitrides. *J. Phys. Chem. A* **2019**,  
869 *123* (38), 8099–8108.
- 870 (40) Pollice, R.; Friederich, P.; Lavigne, C.; dos Passos Gomes, G.;  
871 Aspuru-Guzik, A. Organic Molecules with Inverted Gaps between First  
872 Excited Singlet and Triplet States and Appreciable Fluorescence Rates.  
873 *ChemRxiv*, October 29, 2020, ver. 1. DOI: 10.26434/chem-  
874 rxiv.13087319.v1.
- 875 (41) Chen, J.; Chan, B.; Shao, Y.; Ho, J. How Accurate Are  
876 Approximate Quantum Chemical Methods at Modelling Solute–  
877 Solvent Interactions in Solvated Clusters? *Phys. Chem. Chem. Phys.*  
878 **2020**, *22* (7), 3855–3866.
- 879 (42) Häse, F.; Roch, L. M.; Aspuru-Guzik, A. Next-Generation  
880 Experimentation with Self-Driving Laboratories. *TRECHEM* **2019**, *1*  
881 (3), 282–291.
- 882 (43) Flores-Leonar, M. M.; Mejía-Mendoza, L. M.; Aguilar-Granda,  
883 A.; Sanchez-Lengeling, B.; Tribukait, H.; Amador-Bedolla, C.; Aspuru-  
884 Guzik, A. Materials Acceleration Platforms: On the Way to  
885 Autonomous Experimentation. *Current Opinion in Green and*  
886 *Sustainable Chemistry* **2020**, *25*, 100370.
- 887 (44) Gómez-Bombarelli, R.; Wei, J. N.; Duvenaud, D.; Hernández-  
888 Lobato, J. M.; Sánchez-Lengeling, B.; Sheberla, D.; Aguilera-  
889 Iparraguirre, J.; Hirzel, T. D.; Adams, R. P.; Aspuru-Guzik, A.  
890 Automatic Chemical Design Using a Data-Driven Continuous  
891 Representation of Molecules. *ACS Cent. Sci.* **2018**, *4* (2), 268–276.
- 892 (45) Guimaraes, G. L.; Sanchez-Lengeling, B.; Outeiral, C.; Farias, P.  
893 L. C.; Aspuru-Guzik, A. Objective-Reinforced Generative Adversarial  
894 Networks (ORGAN) for Sequence Generation Models. *arXiv (Machine*  
895 *Learning)*, February 7, 2018, 1705.10843, ver. 3.
- 896 (46) Weininger, D. SMILES, a Chemical Language and Information  
897 System. 1. Introduction to Methodology and Encoding Rules. *J. Chem.*  
898 *Inf. Model.* **1988**, *28* (1), 31–36.
- 899 (47) Krenn, M.; Hase, F.; Nigam, A.; Friederich, P.; Aspuru-Guzik, A.  
900 Self-Referencing Embedded Strings (SELFIES): A 100% Robust  
901 Molecular String Representation. *Mach. Learn.: Sci. Technol.* **2020**, *1*,  
902 045024.
- 903 (48) Aspuru-Guzik-Group/Selfies: [aspuru-guzik-group/selfies](https://github.com/aspuru-guzik-group/selfies). *Gi-*  
904 *tHub*, 2020. <https://github.com/aspuru-guzik-group/selfies>.
- 905 (49) Häse, F.; Aldeghi, M.; Hickman, R. J.; Roch, L. M.; Christensen,  
906 M.; Liles, E.; Hein, J. E.; Aspuru-Guzik, A. Olympos: A Benchmarking  
907 Framework for Noisy Optimization and Experiment Planning. *arXiv*  
908 *(Machine Learning)*, October 8, 2020, 2010.04153, ver. 1.
- 909 (50) Shahriari, B.; Swersky, K.; Wang, Z.; Adams, R. P.; de Freitas, N.  
910 Taking the Human Out of the Loop: A Review of Bayesian  
911 Optimization. *Proc. IEEE* **2016**, *104* (1), 148–175.
- 912 (51) Häse, F.; Roch, L. M.; Aspuru-Guzik, A. Gryffin: An Algorithm  
913 for Bayesian Optimization for Categorical Variables Informed by  
Physical Intuition with Applications to Chemistry. *arXiv (Machine*  
914 *Learning)*, March 26, 2020, 2003.12127, ver 1. 915
- (52) Häse, F.; Roch, L. M.; Aspuru-Guzik, A. Chimera: Enabling  
916 Hierarchy Based Multi-Objective Optimization for Self-Driving  
917 Laboratories. *Chem. Sci.* **2018**, *9* (39), 7642–7655. 918
- (53) Roch, L. M.; Häse, F.; Kreisbeck, C.; Tamayo-Mendoza, T.;  
919 Yunker, L. P. E.; Hein, J. E.; Aspuru-Guzik, A. ChemOS: Orchestrating  
920 Autonomous Experimentation. *Sci. Rob.* **2018**, *3* (19), eaat5559. 921
- (54) Roch, L. M.; Häse, F.; Kreisbeck, C.; Tamayo-Mendoza, T.;  
922 Yunker, L. P. E.; Hein, J. E.; Aspuru-Guzik, A. ChemOS: An  
923 Orchestration Software to Democratize Autonomous Discovery. 924  
*PLoS One* **2020**, *15* (4), e0229862. 925
- (55) MacLeod, B. P.; Parlane, F. G. L.; Morrissey, T. D.; Häse, F.;  
926 Roch, L. M.; Dettelbach, K. E.; Moreira, R.; Yunker, L. P. E.; Rooney,  
927 M. B.; Deeth, J. R.; Lai, V.; Ng, G. J.; Situ, H.; Zhang, R. H.; Elliott, M.  
928 S.; Haley, T. H.; Dvorak, D. J.; Aspuru-Guzik, A.; Hein, J. E.;  
929 Berlinguette, C. P. Self-Driving Laboratory for Accelerated Discovery of  
930 Thin-Film Materials. *Science Advances* **2020**, *6* (20), eaaz8867. 931
- (56) Langner, S.; Häse, F.; Perea, J. D.; Stubhan, T.; Hauch, J.; Roch,  
932 L. M.; Heumueller, T.; Aspuru-Guzik, A.; Brabec, C. J. Beyond Ternary  
933 OPV: High-Throughput Experimentation and Self-Driving Laborato-  
934 rories Optimize Multicomponent Systems. *Adv. Mater.* **2020**, *32* (14),  
935 2070110. 936
- (57) Christensen, M.; Yunker, L. P. E.; Adedeji, F.; Häse, F.; Roch, L.  
937 M.; Gensch, T.; dos Passos Gomes, G.; Zepel, T.; Sigman, M. S.;  
938 Aspuru-Guzik, A. Data-Science Driven Autonomous Process Opti-  
939 mization. *ChemRxiv* November 2, 2020, ver 1. DOI: 10.26434/  
940 chemrxiv.13146404.v1. 941
- (58) Butler, K. T.; Davies, D. W.; Cartwright, H.; Isayev, O.; Walsh, A.  
942 Machine Learning for Molecular and Materials Science. *Nature* **2018**,  
943 *559* (7715), 547–555. 944
- (59) Goodfellow, I.; Bengio, Y.; Courville, A. *Deep Learning*; MIT  
945 Press, 2016. 946
- (60) Duvenaud, D. K.; Maclaurin, D.; Iparraguirre, J.; Bombarell, R.;  
947 Hirzel, T.; Aspuru-Guzik, A.; Adams, R. P. Convolutional Networks on  
948 Graphs for Learning Molecular Fingerprints. In *Advances in Neural*  
949 *Information Processing Systems* 28; Cortes, C., Lawrence, N. D., Lee, D.  
950 D., Sugiyama, M., Garnett, R., Eds.; Curran Associates, Inc., 2015; pp  
951 2224–2232. 952
- (61) Flam-Shepherd, D.; Wu, T.; Friederich, P.; Aspuru-Guzik, A.  
953 Neural Message Passing on High Order Paths. *arXiv (Machine*  
954 *Learning)*, February 24, 2020, 2002.10413, ver. 1. 955
- (62) Pyzer-Knapp, E. O.; Li, K.; Aspuru-Guzik, A. Learning from the  
956 Harvard Clean Energy Project: The Use of Neural Networks to  
957 Accelerate Materials Discovery. *Adv. Funct. Mater.* **2015**, *25* (41),  
958 6495–6502. 959
- (63) Roelofs, R.; Shankar, V.; Recht, B.; Fridovich-Keil, S.; Hardt, M.;  
960 Miller, J.; Schmidt, L. A Meta-Analysis of Overfitting in Machine  
961 Learning. In *Advances in Neural Information Processing Systems* 32;  
962 Wallach, H., Larochelle, H., Beygelzimer, A., Alché-Buc, F., Fox, E.,  
963 Garnett, R., Eds.; Curran Associates, Inc., 2019; pp 9179–9189. 964
- (64) Cawley, G. C.; Talbot, N. L. C. On Over-Fitting in Model  
965 Selection and Subsequent Selection Bias in Performance Evaluation. *J.*  
966 *Mach. Learn. Res.* **2010**, *11* (70), 2079–2107. 967
- (65) Ambrose, C.; McLachlan, G. J. Selection Bias in Gene Extraction  
968 on the Basis of Microarray Gene-Expression Data. *Proc. Natl. Acad. Sci.*  
969 *U. S. A.* **2002**, *99* (10), 6562–6566. 970
- (66) Dupuy, A.; Simon, R. M. Critical Review of Published Microarray  
971 Studies for Cancer Outcome and Guidelines on Statistical Analysis and  
972 Reporting. *J. Natl. Cancer Inst.* **2007**, *99* (2), 147–157. 973
- (67) Shi, L.; Campbell, G.; Jones, W. D.; Campagne, F.; Wen, Z.;  
974 Walker, S. J.; Su, Z.; Chu, T.-M.; Goodsaid, F. M.; Pusztai, L.;  
975 Shaughnessy, J. D.; Oberthuer, A.; Thomas, R. S.; Paules, R. S.; Fielden,  
976 M.; Barlogie, B.; Chen, W.; Du, P.; Fischer, M.; Furlanello, C.; Gallas, B.  
977 D.; Ge, X.; Megherbi, D. B.; Symmans, W. F.; Wang, M. D.; Zhang, J.;  
978 Bitter, H.; Brors, B.; Bushel, P. R.; Bylesjo, M.; Chen, M.; Cheng, J.;  
979 Cheng, J.; Chou, J.; Davison, T. S.; Delorenzi, M.; Deng, Y.;  
980 Devanarayan, V.; Dix, D. J.; Dopazo, J.; Dorff, K. C.; Elloumi, F.;  
981 Fan, J.; Fan, S.; Fan, X.; Fang, H.; Gonzaludo, N.; Hess, K. R.; Hong, H.;  
982

983 Huan, J.; Irizarry, R. A.; Judson, R.; Juraeva, D.; Lababidi, S.; Lambert,  
984 C. G.; Li, L.; Li, Y.; Li, Z.; Lin, S. M.; Liu, G.; Lobenhofer, E. K.; Luo, J.;  
985 Luo, W.; McCall, M. N.; Nikolsky, Y.; Pennello, G. A.; Perkins, R. G.;  
986 Philip, R.; Popovici, V.; Price, N. D.; Qjan, F.; Scherer, A.; Shi, T.; Shi,  
987 W.; Sung, J.; Thierry-Mieg, D.; Thierry-Mieg, J.; Thodima, V.; Trygg, J.;  
988 Vishnuvajjala, L.; Wang, S. J.; Wu, J.; Wu, Y.; Xie, Q.; Yousef, W. A.;  
989 Zhang, L.; Zhang, X.; Zhong, S.; Zhou, Y.; Zhu, S.; Arasappan, D.; Bao,  
990 W.; Lucas, A. B.; Berthold, F.; Brennan, R. J.; Bunes, A.; Catalano, J. G.;  
991 Chang, C.; Chen, R.; Cheng, Y.; Cui, J.; Czika, W.; Demichelis, F.;  
992 Deng, X.; Dosymbekov, D.; Eils, R.; Feng, Y.; Fostel, J.; Fulmer-  
993 Smentek, S.; Fuscoe, J. C.; Gatto, L.; Ge, W.; Goldstein, D. R.; Guo, L.;  
994 Halbert, D. N.; Han, J.; Harris, S. C.; Hatzis, C.; Herman, D.; Huang, J.;  
995 Jensen, R. V.; Jiang, R.; Johnson, C. D.; Jurman, G.; Kahlert, Y.; Khuder,  
996 S. A.; Kohl, M.; Li, J.; Li, L.; Li, M.; Li, Q.-Z.; Li, S.; Li, Z.; Liu, J.; Liu, Y.;  
997 Liu, Z.; Meng, L.; Madera, M.; Martinez-Murillo, F.; Medina, I.;  
998 Meehan, J.; Miclaus, K.; Moffitt, R. A.; Montaner, D.; Mukherjee, P.;  
999 Mulligan, G. J.; Neville, P.; Nikolskaya, T.; Ning, B.; Page, G. P.; Parker,  
1000 J.; Parry, R. M.; Peng, X.; Peterson, R. L.; Phan, J. H.; Quanz, B.; Ren,  
1001 Y.; Riccadonna, S.; Roter, A. H.; Samuelson, F. W.; Schumacher, M. M.;  
1002 Shambaugh, J. D.; Shi, Q.; Shippy, R.; Si, S.; Smalter, A.; Sotiropoulos, C.;  
1003 Soukup, M.; Staedtler, F.; Steiner, G.; Stokes, T. H.; Sun, Q.; Tan, P.-Y.;  
1004 Tang, R.; Tezak, Z.; Thorn, B.; Tsyganova, M.; Turpaz, Y.; Vega, S. C.;  
1005 Visintainer, R.; von Frese, J.; Wang, C.; Wang, E.; Wang, J.; Wang, W.;  
1006 Westermann, F.; Willey, J. C.; Woods, M.; Wu, S.; Xiao, N.; Xu, J.; Xu,  
1007 L.; Yang, L.; Zeng, X.; Zhang, J.; Zhang, L.; Zhang, M.; Zhao, C.; Puri,  
1008 R. K.; Scherf, U.; Tong, W.; Wolfinger, R. D.; MAQC Consortium. The  
1009 MicroArray Quality Control (MAQC)-II Study of Common Practices  
1010 for the Development and Validation of Microarray-Based Predictive  
1011 Models. *Nat. Biotechnol.* **2010**, *28* (8), 827–838.  
1012 (68) Lavigne, C.; dos Passos Gomes, G.; Pollice, R.; Aspuru-Guzik, A.  
1013 Automatic Discovery of Chemical Reactions Using Imposed Activation.  
1014 *ChemRxiv*, September 29, 2020, ver.1. DOI: [10.26434/chem-](https://doi.org/10.26434/chemrxiv.13008500.v1)  
1015 [rxiv.13008500.v1](https://doi.org/10.26434/chemrxiv.13008500.v1).  
1016 (69) Friederich, P.; dos Passos Gomes, G.; Bin, R. D.; Aspuru-Guzik,  
1017 A.; Balcells, D. Machine Learning Dihydrogen Activation in the  
1018 Chemical Space Surrounding Vaska's Complex. *Chem. Sci.* **2020**, *11*  
1019 (18), 4584–4601.  
1020 (70) Friederich, P.; Krenn, M.; Tamblyn, I.; Aspuru-Guzik, A.  
1021 Scientific Intuition Inspired by Machine Learning Generated  
1022 Hypotheses. *arXiv (Machine Learning)*, December 14, 2020,  
1023 2010.14236, ver. 2.  
1024 (71) Aspuru-Guzik, A.; Lindh, R.; Reiher, M. The Matter Simulation  
1025 (R)Evolution. *ACS Cent. Sci.* **2018**, *4* (2), 144–152.  
1026 (72) Steiner, S.; Wolf, J.; Glatzel, S.; Andreou, A.; Granda, J. M.;  
1027 Keenan, G.; Hinkley, T.; Aragon-Camarasa, G.; Kitson, P. J.; Angelone,  
1028 D.; Cronin, L. Organic Synthesis in a Modular Robotic System Driven  
1029 by a Chemical Programming Language. *Science* **2019**, *363* (6423),  
1030 eaav2211.  
1031 (73) Coley, C. W.; Thomas, D. A.; Lummiss, J. A. M.; Jaworski, J. N.;  
1032 Breen, C. P.; Schultz, V.; Hart, T.; Fishman, J. S.; Rogers, L.; Gao, H.;  
1033 Hicklin, R. W.; Plehiers, P. P.; Byington, J.; Piotti, J. S.; Green, W. H.;  
1034 Hart, A. J.; Jamison, T. F.; Jensen, K. F. A Robotic Platform for Flow  
1035 Synthesis of Organic Compounds Informed by AI Planning. *Science*  
1036 **2019**, *365* (6453), eaax1566.  
1037 (74) Burger, B.; Maffettone, P. M.; Gusev, V. V.; Aitchison, C. M.; Bai,  
1038 Y.; Wang, X.; Li, X.; Alston, B. M.; Li, B.; Clowes, R.; Rankin, N.; Harris,  
1039 B.; Sprick, R. S.; Cooper, A. I. A Mobile Robotic Chemist. *Nature* **2020**,  
1040 *583* (7815), 237–241.