Towards comprehensive coverage of chemical space: Quantum mechanical properties of 836k constitutional and conformational closed shell neutral isomers consisting of HCNOFSiPSClBr

Danish Khan, 1,2 Anouar Benali, 3 Scott Y. H. Kim, 1,2 Guido Falk von Rudorff, 4,5 and O. Anatole von Lilienfeld $^{6,2,1,7,8,9,10,\,*}$

¹Chemical Physics Theory Group, Department of Chemistry,
University of Toronto, St. George Campus, Toronto, ON, Canada

²Vector Institute for Artificial Intelligence, Toronto, ON, M5S 1M1, Canada

³Computational Science Division, Argonne National Laboratory, Argonne, Illinois 60439, United States

⁴Institute of Chemistry, University of Kassel, 34109 Kassel, Germany

⁵Center for Interdisciplinary Nanostructure Science and Technology (CINSaT), 34132 Kassel, Germany

⁶Department of Materials Science and Engineering,
University of Toronto, St. George Campus, Toronto, ON, Canada

⁷ML Group, Technische Universität Berlin and Institute for the Foundations of Learning and Data, 10587 Berlin, Germany

⁸Berlin Institute for the Foundations of Learning and Data, 10587 Berlin, Germany

⁹Department of Physics, University of Toronto, St. George Campus, Toronto, ON, Canada

¹⁰Acceleration Consortium, University of Toronto, Toronto, ON

ABSTRACT

The Vector-QM24 (VQM24) dataset attempts to more comprehensively cover all possible neutral closed shell small organic and inorganic molecules and their conformers at state of the art level of theory. We have used density functional theory ($\omega B97X-D3/cc-pVDZ$) to optimize 577k conformational isomers corresponding to 258k constitutional isomers. Isomers included contain up to five heavy atoms (non-hydrogen) consisting of p-block elements C, N, O, F, Si, P, S, Cl, Br. Single point diffusion quantum Monte Carlo (DMC@PBE0(ccECP/cc-pVQZ)) energies are reported for the subset of the lowest conformers of 10,793 molecules with up to 4 heavy atoms. This dataset has been systematically generated by considering all combinatorially possible stoichiometries, and graphs (according to Lewis rules as implemented in the SURGE package), along with all stable conformers identified by GFN2-xTB. Apart from graphs, geometries, rotational constants, and vibrational normal modes, VQM24 includes internal, atomization, electron-electron repulsion, exchange correlation, dispersion, vibrational frequency, Gibbs free, enthalpy, ZPV, molecular orbital energies; as well as entropy, and heat capacities. Electronic properties include multipole moments (dipole, quadrupole, octupole, hexadecapole), electrostatic potentials at nuclei (alchemical potential), Mulliken charges, and molecular wavefunctions. VQM24 represents a highly accurate and unbiased dataset of molecules, ideal for testing and training transferable, scalable, and generative ML models of real quantum systems.

I. BACKGROUND AND SUMMARY

High quality quantum mechanical datasets of molecular properties are a primary requirement for developing approximate physics and statistics based models to enhance the navigation of chemical compound space (CCS). Numerous datasets focusing on distinct, chemically relevant subspaces have paved the way for systematic and quantitative exploration of CCS. Unfortunately, and due to the combinatorial scaling of number of possible stable compounds with size and composition, they are typically incomplete and consequently introduce considerable bias in machine learning (ML) models trained and assessed on them. The effectiveness of ML models relies on the completeness and the accuracy of relevant reference data. Numerous datasets reported so

far include Refs. 1-14 Most of the quantum mechanical (QM) datasets such as QM7¹, QM9³, ANI¹⁰, QMugs⁶, PubChemQC¹⁴ are subsets of other databases such as GDB^{15,16}, ChEMBL¹⁷, PubChem¹⁸ while datasets like GEOM⁷, MultiXC-QM9⁸, ANI-1x¹¹, QM7-X¹² correspond to extensions. While DFT has been key to the development of highly accurate and efficient ML models over the past decade¹⁹, there is still a lack of datasets exhaustively covering specific regions of chemical space at higher QM levels. Note that even for the simplest stoichiometries and smallest subsets of graphs, exhaustive lists of quantum properties are lacking. The primary reason behind this is the combinatorially intractable nature of the problem the number of atoms and of unique chemical elements grows. Nevertheless, it is important to systematically fill these gaps since limited chemical diversity has been shown to lead to less generalizable ML models, as was pointed out recently²⁰. Furthermore, exhaustive coverage of chemical spaces spanned by the smallest systems is key for ML models

^{*} anatole.vonlilienfeld@utoronto.ca

since locality can be exploited to achieve scalable models of extensive properties in many cases^{21–24}.

Here, we tackle this task by reporting VECTOR-QM24 (VQM24), a diverse and comprehensive dataset of - small organic and inorganic molecules calculated at the ω B97X-D3/cc-pVDZ level of theory²⁵⁻²⁸. This includes 258.242 unique constitutional isomers and 577,705 conformers of varying stoichiometries. More specifically, this dataset has been generated by first evaluating all possible Lewis structures (according to SURGE²⁹) for molecules consisting of neutral closed shell combinations of up to five heavy atoms drawn from C, N, O, F, Si, P, S, Cl, Br with their most frequent valencies. Thereafter, conformational isomers were generated for all graphs, using GFN2xTB³⁰, and subsequently relaxed using density functional theory ($\omega B97X-D3/cc-pVDZ$). For all converged molecules, we provide the corresponding optimized structures, an extensive list of thermal properties along with vibrational modes and frequencies, and electronic properties and wavefunctions. Goind beyond DFT, we also report diffusion quantum Monte Carlo (DMC) energies for the smaller sub-set of 10,793 lowest lying conformers of molecules composed of only up to 4 heavy atoms. To the best of our knowledge, this constitutes the largest quantum Monte Carlo (QMC) dataset reported yet. The molecules included in this dataset also constitute an overlapping as well as complementary set of the atoms-inmolecules (amons)²¹ dictionary that models GDB and $ZINC^{31}$.

II. METHODS

A. Structure generation

Element	С	N	О	F	Si	Р	S	Cl	Br
Valencies	4	3,5	2	1	4	3,5	2,4,6	1	1

TABLE I. Chemical elements and their corresponding valencies included in the VQM24 dataset.

In order to generate the structures for VQM24, all combinatorially possible sum formulas from molecules with up to five heavy atoms, drawn from the list of the following chemical elements: C, N₃, N₅, O, F, Si, P₃, P₅, S₂, S₄, S₆, Cl, Br (lower index syntax is in line with SURGE²⁹ and correspond to valencies, as also listed in Table 1). This was done by generating all possible combinations for up to 5 selected elements using the native Python package itertools. Since there are 13 possible heavy elements, the number of combinations containing n-heavy atoms is equal to $\frac{(12+n)!}{12!n!} = 13$, 91, 455, 1820, and 6188 for n=1, 2, 3, 4, and 5, respectively. The possible number of hydrogen atoms, $n_{\rm H}$, for any chosen combination of heavy atoms is then given by the following integer partitioning

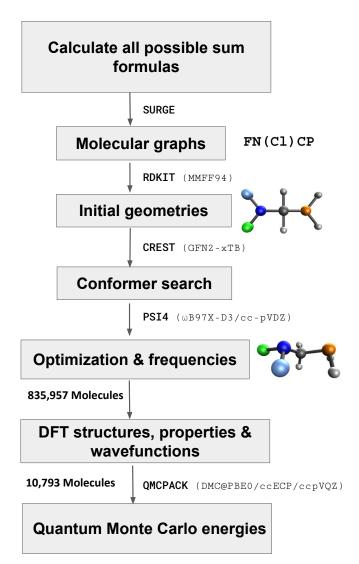


FIG. 1. Workflow used to generate the VQM24 dataset. All possible stoichiometries were first calculated by choosing all combinations of up to 5 heavy atoms (non-Hydrogen) and saturating them with hydrogens to satisfy the valencies. Heavy atoms included along with their valencies are reported in Table 1. For each stoichiometry all possible graphs as identified by the SURGE²⁹ package were evaluated. RDkit³² was then used to generate initial geometries which were initially optimized at the GFN2-xTB³⁰ level of theory followed by a conformer search using CREST³³. All conformers identified at the xTB level of theory were then optimized with DFT (ω B97X-D3/cc-pVDZ²⁶⁻²⁸) using PSI4³⁴ followed by frequency calculations to identify the saddle point orders. For a smaller subset of the most stable conformers with upto 4 heavy atoms, subsequent Diffusion quantum Monte Carlo (DMC) calculations were performed using QMCPACK^{35,36} along with nodal surfaces obtained at the PBE0/ccECP/cc-pVQZ³⁷⁻⁴⁰ level of theory using PySCF.

problem:

$$n_{\rm H} = v - 2n_1 - 4n_2 - 6n_3$$

$$0 \le n_1 \le \frac{v}{2}, \ 0 \le n_2 \le \frac{v}{4}, \ 0 \le n_3 \le \frac{v}{6}$$

$$(1)$$

where n_1 , n_2 , n_3 are integers respectively denoting the number of single, double and triple bonds in the system, while v denotes the total valency of the system (sum over atomic valencies of each element are displayed in Table 1). Note that this procedure also accounts for stoichiometries consistent with ring closures.

For the calculated sum formulas, molecular graphs were generated using SURGE. ²⁹ The graphs were then converted to geometries with MMFF94 as implemented in RDKit³², which were then optimized initially using the GFN2-xTB³⁰ semi-empirical method. Following this, a conformer search was conducted with Crest³³, and all conformers were added to the dataset. This workflow resulted in $\sim 1.1 \mathrm{M}$ geometries.

B. DFT optimization and calculations

Subsequently, we optimized all geometries with DFT using the $\omega B97X-D3/cc-pVDZ$ level of theory²⁶⁻²⁸ and the Gaussian Tight convergence criteria, as implemented in PSI4³⁴. Density fitting was employed in all calculations using the cc-pVDZ-JKFIT⁴¹ auxiliary basis set. All these calculations were conducted using the PSI4 software package (version 1.7).³⁴ The optimization was performed in three passes. In the first pass the default settings in PSI4 for geometry optimizations were used (DIIS method⁴² for SCF and RFO for geometry optimization in redundant internal coordinates⁴³) with a maximum of 100 optimization steps. The molecules that did not converge entered the 2nd pass in which 2nd order SCF convergence method employing the full Newton step was used along with ultrafine Lebedev-Treutler^{44,45} exchange-correlation integration grid (590 spherical, 99 radial points) and a maximum of 100 geometry optimization steps. In the third pass we further employed full Hessian evaluations at the initial step and every 20th goemetry optimization step afterwards with a maximum of 50 optimization steps, in addition to the 2nd pass settings. and we switched to Cartesian only. Those molecules that did not converge after all three passes were left unconverged, and are denoted as such. Following this procedure, we obtained a grand-total of 835,947 converged molecules, and 262,542 that did not converge. The relaxation of converged systems was subsequently followed by vibrational frequency calculations at the same level of theory to identify the saddle point orders of the geometries. In total we found 784,875 molecules to have converged to a local minimum and 51,072 to saddle points. All molecules have been included in the dataset with the minimum geometries and saddle points stored separately.

C. Diffusion Monte Carlo

Quantum Monte Carlo (QMC) techniques are methods that stochastically solve the many-body Schrödinger equation. By explicitly including many-body electronic interactions, these methods achieve mathematical rigor and, in principle, can resolve the Schrödinger equation exactly. However, practical applications require some approximations to maintain computational feasibility, although most of these are controlled and can be rigorously extrapolated at a computational cost. With the proliferation of high-performance computers reaching hundreds of petaflops and the recent deployment of exascale machines, (Summit at Oak Ridge National Laboratory and Aurora at Argonne National Laboratory), QMC methods are poised to significantly take advantage of this computational power; by utilizing stochastic numerical sampling, where samples are evaluated independently, QMC methods achieve embarrassingly parallel processing, enhancing their efficiency for high-performance computing (HPC).

While many variations exist, recent years have seen significant theoretical, algorithmic, and computational advances, particularly in Diffusion Monte Carlo (DMC) Using a projector or Green's function based approach, DMC solves the Schrödinger equation in an imaginary time $\tau = it$. This ensures that any initial state $|\psi\rangle$, not orthogonal to the ground state $|\phi_0\rangle$, will converge to the ground state in a long time limit. During this process, components corresponding to excited states diminish exponentially, ultimately yielding the true ground state.

$$\lim_{\tau \to \infty} \Psi(\mathbf{R}, \tau) = c_0 e^{-\epsilon_0 \tau} \phi_0(\mathbf{R})$$
 (2)

The introduction of a constant energy offset, $E_T = \epsilon_0$, stabilizes the long-time behavior of the system and keeps it finite. The imaginary time Schrödinger equation then resembles a diffusion equation given by:

$$-\frac{\delta\Psi(\mathbf{R},\tau)}{\delta\tau} = \left[\sum_{i=1}^{N} -\frac{1}{2}\nabla_{i}^{2}\Psi(\mathbf{R},\tau)\right] + \left[V(\mathbf{R}) - E_{T}\right]\Psi(\mathbf{R},\tau)$$
(3)

The first term captures the diffusion of particles, while the second term is a branching term dependent on the potential capturing the change in the density of these particles. The potential $V(\mathbf{R})$ in Coulombic systems is unbounded, which may cause the rate term $(V(\mathbf{R}) - E_T)$ to diverge. This could lead to considerable fluctuations in particle density and cause substantial statistical errors. Additionally, the equation doesn't account for the fermionic nature of electrons, which requires antisymmetry when particles are exchanged. This requirement introduces nodes in the fermionic wavefunction; if not constrained, would lead to a bosonic solution. This issue is addressed by the fixed-node (FN) approximation 46 . This approximation constrains the wavefunction to maintain the nodal structure of a trial wavefunction, thereby in-

troducing the fixed-node error as the sole source of error in DMC when the reference wavefunction is not exact. The accuracy of DMC thus heavily relies on the quality of the nodes in the trial wavefunction. By introducing a guiding or trial function, $\Psi_G(\mathbf{R})$, that closely approximates the ground state, the following transformation is applied:

$$f(\mathbf{R}, \tau) = \Psi_G(\mathbf{R}) \Psi(\mathbf{R}, \tau), \qquad (4)$$

which modifies equation (3) to:

$$-\frac{\delta f(\mathbf{R}, \tau)}{\delta \tau} = \left[\sum_{i=1}^{N} -\frac{1}{2} \nabla_{i}^{2} f(\mathbf{R}, \tau) \right] - \nabla \cdot \left[\frac{\nabla \Psi(\mathbf{R})}{\Psi(\mathbf{R})} f(\mathbf{R}, \tau) \right] + (E_{L}(\mathbf{R}) - E_{T}) f(\mathbf{R}, \tau),$$
(5)

 E_T is referred to as a "trial energy" and is used to keep the solution normalized over long-time scales, with $E_L(\mathbf{R})$ representing the local energy at position \mathbf{R} . The final term in Eq. (5) is a critical branching term that eliminates any 'walker' crossing a node (where the wavefunction changes sign) and duplicates any walker that reduces the system's energy, bringing it closer to the ground state. This mechanism is often described as the birth and death process in stochastic simulations. The accuracy of DMC hinges largely on the quality of the nodal surface defined by the underneath trial wavefunction. However, it is important to keep in mind that DMC is variational, meaning that the solutions we obtain are always an upper bound to the exact solution.⁴⁷ This allows for the opportunity of testing with various guiding functions to identify the one that minimizes energy. For instance, Bing et al.⁴⁸ demonstrated that using DFT with 3 different exchange-correlation (XC) functionals as guiding functions yielded consistent results within statistical errors, and this for more than 1000 molecule from the QM5 dataset. Additionally, hybrid functionals can be employed to fine-tune the percentage of exchange, optimizing the energy further⁴⁹. More complex trial wavefunctions, such as multi-Slater determinants generated from selected Configuration Interaction (sCI) $^{50-53}$ or an orbital optimization paired with a variational Monte Carlo in the presence of a Jastrow factor⁵⁴, can also be utilized to improve accuracy. These approaches improve the nodal surface and therefore lower the FN-error associated to DMC, but often at a larger computational cost.

Computational details

All 10,793 constitutional isomers (most stable conformer for each) containing up to 4 heavy atoms in VQM24 were selected for DMC calculations. The total energy calculations were performed using the QMCPACK

code^{35,36}. For efficient sampling and to reduce statistical fluctuations, we utilized a Slater-Jastrow type trial wavefunction for all DMC energy evaluations⁵⁵:

$$\Psi_T(\vec{R}) = \exp\left[\sum_i J_i(\vec{R})\right] \sum_k^M C_k D_k^{\uparrow}(\varphi) D_k^{\downarrow}(\varphi), \quad (6)$$

where $D_k^{\downarrow}(\varphi)$ denotes a Slater determinant composed of single-particle orbitals $\varphi_i = \sum_l^{N_b} C_l^i \Phi_l$, in this study, constructed using PBE0^{39,40} Kohn-Sham (KS) orbitals as implemented in the PySCF code⁵⁶. Similarly to the study by Bing et al.⁴⁸, using different functionals does not lead to very different nodal surfaces and results in energy differences of less than 1kcal/mol in DMC, for the small subset of molecules tested. To enhance efficiency and minimize fluctuations in regions close to ionic cores, ccECP pseudopotentials were applied to substitute core electrons.^{37,38} These pseudopotentials, optimized for precise many-body methods like DMC, address non-local effects using the determinant-localization approximation and the t-moves strategy (DLTM).^{57,58} DMC evolving in real space shows in general minimal basis set size dependency, as documented in prior studies^[53and59]. However, when the the basisset is chosen to be too small, the quality of the trial wavefunction can be severly affected degrading the quality of the nodal surface. Given that the cost of evaluating larger basis function is marginal in QMC, we ran all our simulations with the cc-pVQZ basis set, tailored for ccECP.

The Jastrow function includes terms for one-body (electron-ion), two-body (electron-electron), and threebody (electron-electron-ion) interactions. The oneand two-body interactions were defined using spline functions⁶⁰, and the three-body interactions were modeled using polynomials. 61 Specifically, the study utilized 16 parameters for each atom type in one-body terms with a cutoff of 8 Bohr, and 20 parameters per spin-channel for two-body terms with a cutoff of 10 Bohr. The three-body terms incorporated 26 parameters each, with a 5 Bohr cutoff. These parameters in the Jastrow factor were individually optimized for each molecular geometry using a linear optimization method developed by Umrigar et al.⁶². For all DMC simulations, we utilized a timestep of 0.001 a.u., eliminating the necessity for timestep extrapolation. The simulations involved 1500 blocks of 40 imaginary time steps each, with only the 40th step considered for calculating the standard deviation. We used 16,000 walkers to reduce autocorrelation and prevent population bias, achieving average error bars of 0.4 mHa across approximately 2.3 billion samples.

Given the large number of molecules, we used the NEXUS Workflow package⁶³ to generate input files, manage, and monitor jobs across different stages of the calculations. This allowed for a "black box" and fully automated computational campaign. The trial wavefunction generation was conducted on the Argonne

LCRC system, Improv, using a single node composed of 2x AMD EPYC 7713 with 64 cores at 2GHz. Each molecule, on average, required 45 seconds of computing, amounting to a total of 134 node hours for the whole set. The subsequent DMC calculations required 20 nodes per molecule on the Argonne Polaris HPC, using the AMD EPYC 7543P CPU with 64 cores at 2.8GHz. Each molecule took approximately 15 minutes to achieve a sub kcal/mol error bar, totaling around $\tilde{5}4,000$ node hours.

III. DATA RECORDS

The DFT dataset is reported as separate .npz files for, conformational minima, constitutional minima and saddle point structures. Each property is reported in separate arrays with the ordering of the molecules across every array being the same. The keys for accessing each property from the DFT .npz files are tabulated in Table 2. DMC data is similarly reported in a sep-

Property	Unit	Key
Stoichiometry	-	compounds
Atomic Numbers	_	atoms
Cartesian coordinates (XYZ)	Å	coordinates
SMILES	-	graphs
InCHI strings	-	inchi
Total energies	Ha	Etot
Internal energies	Ha	UO
Atomization energies		Eatomization
Electron-electron energies		Eee
Exchange correlation energies	Ha	Exc
Dispersion energy	Ha	Edisp
HOMO-LUMO gap	Ha	gap
Dipole moments	a.u.	dipole
Quadrupole moments	a.u.	quadrupole
Octupole moments	a.u.	octupole
Hexadecapole moments	a.u.	hexadecapole
Rotational constants	MHz	rots
Vibrational eigen modes	Å	vibmodes
Vibrational frequencies	$ \mathrm{cm}^{-1} $	freqs
Free energy (H)	Ha	G 1
Internal (Thermal) energy (H)	Ha	U298
Enthalpy (H)	Ha	H
Zero point vibrational energy (H)	Ha	zpves
Entropy (H)	$\frac{\text{cal}}{\text{mol}_{\cdot}\text{K}}$	S
Heat capacities (H)	l cal	Cv, Cp
Electrostatic potentials at nuclei	mol K	Vesp
Mulliken charges	a.u.	Qmulliken
MO energies (molden files)	На	-
Wavefunctions (molden files)	110	
waveruncoms (morden mes)	_	

TABLE II. DFT properties for all the 835'947 converged molecules with up to 5 heavy atoms, along with the corresponding keys to access them from the reported .npz files in the VQM24 dataset. (H) indicates thermodynamic properties calculated via the Harmonic approximation.

arate .npz file with the corresponding keys recorded in Table 3. All of the data is publicly available at https://doi.org/10.5281/zenodo.11164951

Property	Unit	Key	
Stoichiometry	-	compounds	
Atomic Numbers	-	atoms	
Cartesian coordinates (XYZ)	Å	coordinates	
SMILES	_	graphs	
InCHI strings	_	inchi	
Total energy	На	Etot	
Error bar	На	std	

TABLE III. DMC properties for 10'793 molecules with up to 4 heavy atoms, along with the corresponding keys to access them from the reported .npz files in the VQM24 dataset.

ACKNOWLEDGMENTS

We acknowledge the support of the Natural Sciences and Engineering Research Council of Canada (NSERC). [funding reference number RGPIN-2023-04853]. Cette recherche a été financée par le Conseil de recherches en sciences naturelles et en génie du Canada (CRSNG), [numéro de référence RGPIN-2023-04853]. This research was undertaken thanks in part to funding provided to the University of Toronto's Acceleration Consortium from the Canada First Research Excellence Fund, grant number: CFREF-2022-00042. O.A.v.L. has received support as the Ed Clark Chair of Advanced Materials and as a Canada CIFAR AI Chair. O.A.v.L. has received funding from the European Research Council (ERC) under the European Union's Horizon 2020 research and innovation programme (grant agreement No. 772834). The authors are grateful to Compute Canada and the Acceleration Consortium for computational resources. A.B was funded by the U.S. Department of Energy, Office of Science, Basic Energy Sciences, Materials Sciences and Engineering Division, as part of the Computational Materials Sciences Program and Center for Predictive Simulation of Functional Materials. DMC calculations used an award of computer time provided by the Innovative and Novel Computational Impact on Theory and Experiment (INCITE) program. This research used resources of the Argonne Leadership Computing Facility, which is a DOE Office of Science User Facility supported under contract DE-AC02-06CH11357. We also gratefully acknowledge the computing resources provided on IMPROV, a highperformance computing cluster operated by the Laboratory Computing Resource Center (LCRC) at Argonne National Laboratory.

^[1] M. Rupp, A. Tkatchenko, K.-R. Müller, and O. A. von Lilienfeld, Phys. Rev. Lett. **108**, 058301 (2012).

^[2] G. Montavon, M. Rupp, V. Gobre, A. Vazquez-Mayagoitia, K. Hansen, A. Tkatchenko, K.-R. Müller,

- and O. A. von Lilienfeld, New Journal of Physics 15, 095003 (2013).
- [3] R. Ramakrishnan, P. O. Dral, M. Rupp, and O. A. von Lilienfeld, Scientific Data 1 (2014).
- [4] S. Chmiela, A. Tkatchenko, H. E. Sauceda, I. Poltavsky, K. T. Schütt, and K.-R. Müller, Science Advances 3 (2017), 10.1126/sciadv.1603015.
- [5] A. S. Christensen and O. A. von Lilienfeld, Machine Learning: Science and Technology 1, 045018 (2020).
- [6] C. Isert, K. Atz, J. Jiménez-Luna, and G. Schneider, Scientific Data 9, 273 (2022).
- [7] S. Axelrod and R. Gómez-Bombarelli, Scientific Data 9 (2022), 10.1038/s41597-022-01288-4.
- [8] S. Nandi, T. Vegge, and A. Bhowmik, Scientific Data 10 (2023), 10.1038/s41597-023-02690-2.
- [9] J. S. Smith, O. Isayev, and A. E. Roitberg, Scientific Data 4 (2017), 10.1038/sdata.2017.193.
- [10] J. S. Smith, O. Isayev, and A. E. Roitberg, Chem. Sci. 8, 3192 (2017).
- [11] J. S. Smith, R. Zubatyuk, B. Nebgen, N. Lubbers, K. Barros, A. E. Roitberg, O. Isayev, and S. Tretiak, Scientific Data 7 (2020), 10.1038/s41597-020-0473-z.
- [12] J. Hoja, L. M. Sandonas, B. G. Ernst, A. Vazquez-Mayagoitia, R. A. DiStasio, and A. Tkatchenko, Scientific Data 8 (2021), 10.1038/s41597-021-00812-2.
- [13] X. Qu, A. Jain, N. N. Rajput, L. Cheng, Y. Zhang, S. P. Ong, M. Brafman, E. Maginn, L. A. Curtiss, and K. A. Persson, Comput. Mater. Sci. 103, 56 (2015).
- [14] M. Nakata and T. Shimazaki, Journal of Chemical Information and Modeling 57, 1300 (2017), pMID: 28481528, https://doi.org/10.1021/acs.jcim.7b00083.
- [15] L. C. Blum and J.-L. Reymond, Journal of the American Chemical Society 131, 8732 (2009).
- [16] L. Ruddigkeit, R. Van Deursen, L. C. Blum, and J.-L. Reymond, Journal of chemical information and modeling 52, 2864 (2012).
- [17] A. Gaulton, L. J. Bellis, A. P. Bento, J. Chambers, M. Davies, A. Hersey, Y. Light, S. McGlinchey, D. Michalovich, B. Al-Lazikani, and J. P. Overington, Nucleic Acids Research 40, D1100 (2011).
- [18] S. Kim, P. A. Thiessen, E. E. Bolton, J. Chen, G. Fu, A. Gindulyte, L. Han, J. He, S. He, B. A. Shoemaker, J. Wang, B. Yu, J. Zhang, and S. H. Bryant, Nucleic Acids Research 44, D1202 (2015).
- [19] B. Huang, G. F. von Rudorff, and O. A. von Lilienfeld, Science 381, 170 (2023), https://www.science.org/doi/pdf/10.1126/science.abn3445.
- [20] M. Glavatskikh, J. Leguy, G. Hunault, T. Cauchy, and B. Da Mota, Journal of Cheminformatics 11 (2019), 10.1186/s13321-019-0391-2.
- [21] B. Huang and O. A. von Lilienfeld, Nature Chemistry 12, 945 (2020).
- [22] O. T. Unke, M. Stöhr, S. Ganscha, T. Unterthiner, H. Maennel, S. Kashubin, D. Ahlin, M. Gastegger, L. Medrano Sandonas, J. T. Berryman, et al., Science Advances 10, eadn4397 (2024).
- [23] D. Khan, A. J. A. Price, M. L. Ach, and O. A. von Lilienfeld, "Adaptive hybrid density functionals," (2024), arXiv:2402.14793 [physics.chem-ph].
- [24] D. Khan, M. L. Ach, and O. A. von Lilienfeld, "Adaptive atomic basis sets," (2024), arXiv:2404.16942 [physics.chem-ph].
- [25] J.-D. Chai and M. Head-Gordon, Phys. Chem. Chem. Phys. 10, 6615 (2008).

- [26] J. Dunning, Thom H., The Journal of Chemical Physics 90, 1007 (1989), https://pubs.aip.org/aip/jcp/articlepdf/90/2/1007/15358102/1007_1_online.pdf.
- [27] J.-D. Chai and M. Head-Gordon, Physical Chemistry Chemical Physics 10, 6615 (2008).
- [28] S. Grimme, J. Antony, S. Ehrlich, and H. Krieg, J. Chem. Phys. 132, 154104 (2010).
- [29] B. D. McKay, M. A. Yirik, and C. Steinbeck, Journal of Cheminformatics 14, 24 (2022).
- [30] C. Bannwarth, S. Ehlert, and S. Grimme, Journal of Chemical Theory and Computation 15, 1652 (2019).
- [31] B. Huang and O. A. von Lilienfeld, "Dictionary of 140k gdb and zinc derived amons," (2020), arXiv:2008.05260 [physics.chem-ph].
- [32] G. Landrum, Release 1, 4 (2013).
- [33] P. Pracht, F. Bohle, and S. Grimme, Physical Chemistry Chemical Physics **22**, 7169 (2020).
- [34] D. G. Smith, L. A. Burns, A. C. Simmonett, R. M. Parrish, M. C. Schieber, R. Galvelis, P. Kraus, H. Kruse, R. Di Remigio, A. Alenaizan, et al., The Journal of chemical physics 152 (2020).
- [35] J. Kim, A. T. Baczewski, T. D. Beaudet, A. Benali, M. C. Bennett, M. A. Berrill, N. S. Blunt, E. J. L. Borda, M. Casula, D. M. Ceperley, et al., Journal of Physics: Condensed Matter 30, 195901 (2018).
- [36] P. R. C. Kent, A. Annaberdiyev, A. Benali, M. C. Bennett, E. J. Landinez Borda, P. Doak, H. Hao, K. D. Jordan, J. T. Krogel, I. Kylänpää, J. Lee, Y. Luo, F. D. Malone, C. A. Melton, L. Mitas, M. A. Morales, E. Neuscamman, F. A. Reboredo, B. Rubenstein, K. Saritas, S. Upadhyay, G. Wang, S. Zhang, and L. Zhao, The Journal of Chemical Physics 152, 174105 (2020), https://doi.org/10.1063/5.0004860.
- [37] A. Annaberdiyev, G. Wang, C. A. Melton, M. C. Bennett, L. Shulenburger, and L. Mitas, J. Chem. Phys. 149, 134108 (2018).
- [38] M. C. Bennett, C. A. Melton, A. Annaberdiyev, G. Wang, L. Shulenburger, and L. Mitas, J. Chem. Phys. 147, 224106 (2017).
- [39] M. Ernzerhof and G. E. Scuseria, J. Chem. Phys. 110, 5029 (1999).
- [40] C. Adamo and V. Barone, J. Chem. Phys. 110, 6158 (1999).
- [41] F. Weigend, Phys. Chem. Chem. Phys. 4, 4285 (2002).
- [42] P. Pulay, Chemical Physics Letters **73**, 393 (1980).
- [43] V. Bakken and T. Helgaker, The Journal of Chemical Physics 117, 9160 (2002), https://pubs.aip.org/aip/jcp/article-pdf/117/20/9160/19243383/9160_1_online.pdf.
- R. The [44] O. Treutler and Ahlrichs, Jour-**Physics** Chemical **102**. 346 nal of (1995),https://pubs.aip.org/aip/jcp/articlepdf/102/1/346/19034811/346_1_online.pdf.
- [45] V. I. Lebedev and D. N. Laikov, Doklady Mathematics 59, 477 (1999).
- [46] J. B. Anderson, The Journal of Chemical Physics 65, 4121 (1976), https://doi.org/10.1063/1.432868.
- [47] D. M. Ceperley and B. J. Alder, The Journal of Chemical Physics 81, 5833 (1984), https://doi.org/10.1063/1.447637.
- [48] B. Huang, O. A. von Lilienfeld, J. T. Krogel, and A. Benali, Journal of Chemical Theory and Computation 19, 1711 (2023), pMID: 36857531, https://doi.org/10.1021/acs.jctc.2c01058.

- [49] B. Busemeyer, M. Dagrada, S. Sorella, M. Casula, and L. K. Wagner, Phys. Rev. B 94, 035108 (2016).
- [50] M. A. Morales, J. McMinis, B. K. Clark, J. Kim, and G. E. Scuseria, Journal of Chemical Theory and Computation 8, 2181 (2012), http://dx.doi.org/10.1021/ct3003404.
- [51] M. Caffarel, T. Applencourt, E. Giner, and A. Scemama, The Journal of Chemical Physics 144, 151103 (2016), https://doi.org/10.1063/1.4947093.
- [52] A. Scemama, A. Benali, D. Jacquemin, M. Caffarel, and P.-F. Loos, The Journal of Chemical Physics 149, 034108 (2018), https://doi.org/10.1063/1.5041327.
- [53] F. D. Malone, A. Benali, M. A. Morales, M. Caffarel, P. R. Kent, and L. Shulenburger, "Systematic comparison and cross-validation of fixed-node diffusion monte carlo and phaseless auxiliary-field quantum monte carlo in solids," (2020).
- [54] E. Slootman, I. Poltavsky, R. Shinde, J. Cocomello, S. Moroni, A. Tkatchenko, and C. Filippi, arXiv preprint arXiv:2404.09755 (2024).
- [55] K. E. Schmidt and J. W. Moskowitz, The Journal of Chemical Physics 93, 4172 (1990).
- [56] Q. Sun, X. Zhang, S. Banerjee, P. Bao, M. Barbry, N. S. Blunt, N. A. Bogdanov, G. H. Booth, J. Chen, Z.-H. Cui, J. J. Eriksen, Y. Gao, S. Guo, J. Hermann, M. R. Hermes, K. Koh, P. Koval, S. Lehtola, Z. Li, J. Liu,

- N. Mardirossian, J. D. McClain, M. Motta, B. Mussard, H. Q. Pham, A. Pulkin, W. Purwanto, P. J. Robinson, E. Ronca, E. R. Sayfutyarova, M. Scheurer, H. F. Schurkus, J. E. T. Smith, C. Sun, S.-N. Sun, S. Upadhyay, L. K. Wagner, X. Wang, A. White, J. D. Whitfield, M. J. Williamson, S. Wouters, J. Yang, J. M. Yu, T. Zhu, T. C. Berkelbach, S. Sharma, A. Y. Sokolov, and G. K.-L. Chan, The Journal of Chemical Physics 153, 024109 (2020), https://pubs.aip.org/aip/jcp/article-pdf/doi/10.1063/5.0006074/16722275/024109_1_online.pdf.
- [57] A. Zen, J. G. Brandenburg, A. Michaelides, and D. Alfe, J. Chem. Phys. 151, 134105 (2019).
- [58] M. Casula, S. Moroni, S. Sorella, and C. Filippi, J. Chem. Phys. 132, 154113 (2010).
- [59] M. c. v. Dubecký, J. Chem. Theory Comput. 13, 3626 (2017).
- [60] K. Esler, J. Kim, D. Ceperley, and L. Shulenburger, Computing in Science Engineering 14, 40 (2012).
- [61] N. D. Drummond, M. D. Towler, and R. J. Needs, Phys. Rev. B 70, 235119 (2004).
- [62] C. J. Umrigar, J. Toulouse, C. Filippi, S. Sorella, and R. G. Hennig, Phys. Rev. Lett. 98, 110201 (2007).
- [63] J. T. Krogel, Computer Physics Communications 198, 154 (2016).