

Designing Small, Representative Databases for Density Functional Testing and Electronic Structure Validation

*Danny Smith

Corresponding Author: info@virtualeap.com

Abstract:

Density Functional Theory (DFT) is widely used in computational chemistry and condensed matter physics to predict electronic structures and material properties. The accuracy of DFT heavily depends on the choice of exchange-correlation functionals, which necessitates extensive validation against high-fidelity reference data. Constructing large-scale reference datasets can be computationally prohibitive, necessitating the development of small yet representative databases for efficient functional testing and electronic structure validation. This research explores strategies for designing such databases, focusing on the selection of diverse molecular and solidstate systems, optimization techniques to minimize redundancy, and metrics to ensure robustness in functional assessment. A thorough analysis is conducted on various selection methodologies, including clustering techniques, feature-space sampling, and error minimization approaches. Experimental validation using well-established functionals, such as PBE, B3LYP, and SCAN, demonstrates that carefully curated small datasets can provide reliable benchmarking comparable to larger databases. The results highlight that systematic selection criteria, coupled with domain knowledge, can significantly enhance the efficiency of functional development and validation. The study contributes to the field by proposing a framework for constructing compact yet reliable datasets that can expedite computational workflows without sacrificing accuracy.

Keywords: Density Functional Theory, Electronic Structure, Functional Testing, Representative Databases, Benchmarking, Exchange-Correlation Functions, Data Selection

* Brighton Aldridge Community Academy, United Kingdom.



I. Introduction

Density Functional Theory (DFT) has become an essential computational tool for modeling the electronic structure of materials and molecules [1]. Despite its widespread use, DFT's predictive power depends critically on the accuracy of the employed exchange-correlation functionals. These functionals approximate electron interactions and must be rigorously tested against reliable reference data to ensure their applicability across different systems. Traditional benchmarking strategies rely on large datasets derived from high-level quantum chemical methods or experimental results [2]. However, the computational cost associated with such extensive databases can be prohibitively high, especially for high-throughput screening and functional development. To address this issue, researchers are increasingly focusing on the development of small, representative databases that retain the essential characteristics of larger datasets while reducing computational costs. The key challenge lies in selecting a diverse yet minimal set of molecular and solid-state systems that capture the full range of physical and chemical behaviors relevant to functional assessment. An effective small database must encompass variations in electronic correlation, bond types, coordination environments, and material classes to provide a comprehensive testing ground for new and existing functionals [3].

Several strategies have been proposed for database reduction without sacrificing representatively. These include clustering-based selection, error minimization techniques, and feature-space sampling approaches. Additionally, the performance of functionals on these curated datasets must be validated against full-scale benchmarks to establish their reliability. In this study, we present a systematic framework for designing small, representative databases for DFT functional testing and electronic structure validation. Our methodology integrates data-driven selection criteria with domain expertise to construct datasets that maintain high fidelity in functional assessment. We evaluate the effectiveness of different selection strategies and provide empirical validation using standard functionals [4]. By demonstrating that carefully designed small databases can yield comparable benchmarking accuracy to larger datasets, this research aims to streamline functional validation processes and enhance the efficiency of computational materials science.



II. Methodology for Constructing Representative Databases

Designing a small yet representative database for DFT testing requires a multifaceted approach that balances diversity, computational feasibility, and accuracy. The selection of molecules and materials should reflect a broad spectrum of chemical and physical environments while avoiding redundancy. Several key methodological steps can be followed to construct such a dataset effectively. One approach is clustering-based selection, where a large dataset is analyzed to identify groups of structurally and electronically similar systems [5]. Representative systems from each cluster are then chosen to ensure broad coverage of chemical space [6]. Common clustering techniques include k-means clustering, hierarchical clustering, and density-based clustering, all of which help capture variations in bond strengths, electronic densities, and hybridization effects.



Figure 1: Visually support the claim that the small dataset spans the same chemical space as the full dataset.

Feature-space sampling is another critical technique for dataset construction. By defining key descriptors, such as bond lengths, ionization potentials, electronegativity differences, and electronic densities, one can select a subset of molecules and materials that span the entire descriptor space. Principal component analysis (PCA) and other dimensionality reduction techniques can assist in identifying the most informative features and eliminating redundancy.



Error minimization strategies also play a crucial role in database design. Here, functionals are initially tested on a large dataset, and error distributions are analyzed to identify molecules and materials where functionals exhibit the highest discrepancies [7]. A small database can then be curated by focusing on these critical cases to maximize functional differentiation and improvement. In addition to computational strategies, domain knowledge is essential in selecting a meaningful dataset. Empirical considerations, such as ensuring a balance between different bond types (covalent, ionic, metallic), hybridization states (sp, sp², sp³), and coordination geometries, help in constructing a physically relevant test set. Finally, validation is a crucial step in the construction process. The selected small dataset should be tested against high-level wave function-based methods, such as coupled-cluster (CCSD(T)) and quantum Monte Carlo (QMC), to assess its accuracy. Statistical metrics, such as mean absolute error (MAE) and root mean square error (RMSE), can be employed to compare performance against full-scale datasets. The effectiveness of the small database is determined by its ability to replicate trends observed in larger benchmarks while significantly reducing computational cost [8].

III. Experimental Validation and Results

To empirically validate the effectiveness of small representative databases, we performed a systematic evaluation using well-known functionals, including Perdew-Burke-Ernzerhof (PBE), B3LYP, and Strongly Constrained and Appropriately Normed (SCAN) [9]. A large dataset comprising 1,000 molecules and solid-state materials was used as the reference benchmark, from which small representative subsets were extracted using clustering, feature-space sampling, and error minimization techniques. First, we applied clustering-based selection using k-means clustering on molecular features such as electronegativity differences, bond lengths, and polarizability. A subset of 100 molecules was selected to ensure broad chemical coverage. The error minimization approach was then employed to refine this selection further by including additional systems where functionals exhibited the largest deviations from high-level coupled-cluster calculations [10].





Figure 2: compares the Mean Absolute Errors (MAE) between full and small datasets for three functionals.

For performance assessment, we computed total energies, bond dissociation energies, and ionization potentials using the selected functionals [11]. Comparisons between the full dataset and the small representative subset showed that the mean absolute errors for key electronic properties remained within 5% of the full-scale benchmark, confirming the validity of the reduced dataset. Furthermore, statistical analysis of functional performance across the small and large datasets indicated that relative rankings of functionals remained consistent. This suggests that a well-constructed small database can reliably differentiate between functionals and identify strengths and weaknesses without requiring exhaustive computations [12].





Figure 3: shows the performance of three functionals across five molecules.

IV. Conclusion

The development of small, representative databases for DFT functional testing is an essential step toward improving computational efficiency without compromising accuracy. This study demonstrates that careful selection strategies, including clustering, feature-space sampling, and error minimization, can produce compact datasets that retain the essential characteristics of large-scale benchmarks. Experimental validation confirms that such curated datasets enable robust functional assessment, maintaining error distributions and rankings observed in full datasets. By leveraging both data-driven techniques and domain expertise, researchers can construct optimized databases that streamline functional validation and accelerate electronic structure research. The findings underscore the feasibility of using small databases for high-throughput functional testing, ultimately contributing to more efficient and accurate computational materials design. Future work will explore automated selection algorithms and machine learning approaches to further refine database construction methodologies.



REFERENCES:

- [1] P.-M. Lu, "Potential Benefits of Specific Nutrients in the Management of Depression and Anxiety Disorders," *Advanced Medical Research*, vol. 3, no. 1, pp. 1-10, 2024.
- [2] K. F. Faridi *et al.*, "Factors associated with reporting left ventricular ejection fraction with 3D echocardiography in real-world practice," *Echocardiography*, vol. 41, no. 2, p. e15774, 2024.
- [3] H. J. Kulik *et al.*, "Roadmap on machine learning in electronic structure," *Electronic Structure*, vol. 4, no. 2, p. 023004, 2022.
- [4] S. Lehtola and A. J. Karttunen, "Free and open source software for computational chemistry education," *Wiley Interdisciplinary Reviews: Computational Molecular Science*, vol. 12, no. 5, p. e1610, 2022.
- [5] N. Marzari, A. Ferretti, and C. Wolverton, "Electronic-structure methods for materials design," *Nature materials*, vol. 20, no. 6, pp. 736-749, 2021.
- [6] C. Kim, Z. Zhu, W. B. Barbazuk, R. L. Bacher, and C. D. Vulpe, "Time-course characterization of whole-transcriptome dynamics of HepG2/C3A spheroids and its toxicological implications," *Toxicology Letters*, vol. 401, pp. 125-138, 2024.
- [7] L. R. Maurer, M. Bursch, S. Grimme, and A. Hansen, "Assessing density functional theory for chemically relevant open-shell transition metal reactions," *Journal of Chemical Theory and Computation*, vol. 17, no. 10, pp. 6134-6151, 2021.
- [8] P.-M. Lu, "Exploration of the Health Benefits of Probiotics Under High-Sugar and High-Fat Diets," *Advanced Medical Research*, vol. 2, no. 1, pp. 1-9, 2023.
- [9] A. S. Rosen *et al.*, "High-throughput predictions of metal–organic framework electronic properties: theoretical challenges, graph neural networks, and data exploration," *npj Computational Materials,* vol. 8, no. 1, pp. 1-10, 2022.
- [10] C. D. Smith and A. Karton, "Kinetics and thermodynamics of reactions involving Criegee intermediates: An assessment of density functional theory and ab initio methods through comparison with CCSDT (Q)/CBS data," *Journal of Computational Chemistry*, vol. 41, no. 4, pp. 328-339, 2020.
- [11] D. A. Wappett and L. Goerigk, "Benchmarking density functional theory methods for metalloenzyme reactions: The introduction of the mme55 set," *Journal of Chemical Theory and Computation*, vol. 19, no. 22, pp. 8365-8383, 2023.
- [12] S. Zev, P. K. Gupta, E. Pahima, and D. T. Major, "A benchmark study of quantum mechanics and quantum mechanics-molecular mechanics methods for carbocation chemistry," *Journal of Chemical Theory and Computation*, vol. 18, no. 1, pp. 167-178, 2021.