# Eugen Hruška, Ph.D.

🌐 Hruska-Lab.github.io    ⓘD 0000-0001-5679-8419    in eugen-hruska
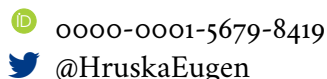
📍 Charles University    🐦 @HruskaEugen    ✉ eugen.hruska@faf.cuni.cz

## Publications

fractional h-index 3 (OpenAlex), h-index 6 (WoS), citations 107 (WoS)

1   Chen, X., Sun, Y., **Hruska**, E., Dixit, V., Yang, J., He, Y., Wang, Y., & Liu, F. (2024). Explainable machine learning identification of superconductivity from single-particle spectral functions. *Preprint.* 🔗 https://arxiv.org/abs/2406.04445

Abstract: The traditional method of identifying symmetry-breaking phase transitions through the emergence of a single-particle gap encounters significant challenges in quantum materials with strong fluctuations. To address this, we have developed a data-driven approach using a domain-adversarial neural network trained on simulated spectra of cuprates. This model compensates for the scarcity of experimental data – a significant barrier to the wide deployment of machine learning in physical research – by leveraging the abundance of theoretically simulated data. When applied to unlabeled experimental spectra, our model successfully distinguishes the true superconducting states from gapped fluctuating states, without the need for fine temperature sampling across the transition. Further, the explanation of our machine learning model reveals the crucial role of the Fermi-surface spectral intensity even in gapped states. It paves the way for robust and direct spectroscopic identification of fluctuating orders, particularly in low-dimensional, strongly correlated materials.

2   Suwała, D., & **Hruska**, E. (2024). The wins and failures of current docking methods tested on the flexible active site of cytochromes p450. *Preprint.*
🔗 https://doi.org/10.26434/chemrxiv-2024-05299

Abstract: In this work, we benchmark 4 selected open source docking engines for use in the cytochrome P450 protein family. The key enzymes family of phase I metabolism is characterized by a wide variety of accepted substrates due to flexible active site. This work is a benchmark study which aims to evaluate the capabilities of current rigid and induced-fit docking methods for prediction of correct heme-ligand orientation. To asses it, we use two unique distances to heme iron and a SuCOS score to quantify reconstruction of orientation and chemical features. We selected three rigid protein docking engines: GNINA, AutoDock VINA, GalaxyDock2 HEME and a flexible docking model, RosettaFold-All-Atoms to test them on a dataset of 128 CYP-binding ligands. We report mean absolute error for RosettaFold-All-Atom on key distance, to the atom closest to heme iron in experimental reference structure, 3 times lower than AutoDock VINA engine in the same simulation. Our results indicate that induced fit method is a significant improvement over rigid methods for flexible active site, but still offer limited predictivity. During crossdocking, RosettaFold-All-Atoms was able to recreate over a quarter of distances up to 20 percent difference from experiment. Further analysis indicates a low overlap in the distribution of ligand chemical features, based on a SuCOS score, which suggests a space for further improvement.

3   Chen, X., Li, P., **Hruska**, E., & Liu, F. (2023). Δ-machine learning for quantum chemistry prediction of solution-phase molecular properties at the ground and excited states. *Phys. Chem. Chem. Phys.*, *25*(19), 13417–13428. 🔗 https://doi.org/10.1039/D3CP00506B

Abstract: Due to the limitation of solvent models, quantum chemistry calculation of solution-phase molecular properties often deviates from experimental measurements. Recently, $\Delta$-machine learning ($\Delta$-ML) was shown to be a promising approach to correcting errors in the quantum chemistry calculation of solvated molecules. However, this approach's applicability to different molecular properties and its performance in various cases are still unknown. In this work, we tested the performance of $\Delta$-ML in correcting redox potential and absorption energy calculations using four types of input descriptors and various ML methods. We sought to understand the dependence of $\Delta$-ML performance on the property to predict the quantum chemistry method, the data set distribution/size, the type of input feature, and the feature selection techniques. We found that $\Delta$-ML can effectively correct the errors in redox potentials calculated using density functional theory (DFT) and absorption energies calculated by time-dependent DFT. For both properties, the $\Delta$-ML-corrected results showed less sensitivity to the DFT functional choice than the raw results. The optimal input descriptor depends on the property, regardless of the specific ML method used. The solvent–solute descriptor (SS) is the best for redox potential, whereas the combined molecular fingerprint (cFP) is the best for absorption energy. A detailed analysis of the feature space and the physical foundation of different descriptors well explained these observations. Feature selection did not further improve the $\Delta$-ML performance. Finally, we analyzed the limitation of our $\Delta$-ML solvent effect approach in data sets with molecules of varying degrees of electronic structure errors.

**4**  **Hruska**, E., Gale, A., Huang, X., & Liu, F. (2022). AutoSolvate : A Toolkit for Automating Quantum Chemistry Design and Discovery of Solvated Molecules. *J. Chem. Phys., 156*(12). 🔗 https://doi.org/10.1063/5.0084833

Abstract: The availability of large, high-quality data sets is crucial for artificial intelligence design and discovery in chemistry. Despite the essential roles of solvents in chemistry, the rapid computational data set generation of solution-phase molecular properties at the quantum mechanical level of theory was previously hampered by the complicated simulation procedure. Software toolkits that can automate the procedure to set up high-throughput explicit-solvent quantum chemistry (QC) calculations for arbitrary solutes and solvents in an open-source framework are still lacking. We developed AutoSolvate, an open-source toolkit to streamline the workflow for QC calculation of explicitly solvated molecules. It automates the solvated-structure generation, force field fitting, configuration sampling, and the final extraction of microsolvated cluster structures that QC packages can readily use to predict molecular properties of interest. AutoSolvate is available through both a command line interface and a graphical user interface, making it accessible to the broader scientific community. To improve the quality of the initial structures generated by AutoSolvate, we investigated the dependence of solute-solvent closeness on solute/solvent identities and trained a machine learning model to predict the closeness and guide initial structure generation. Finally, we tested the capability of AutoSolvate for rapid data set curation by calculating the outer-sphere reorganization energy of a large data set of 166 redox couples, which demonstrated the promise of the AutoSolvate package for chemical discovery efforts.

**5**  **Hruska**, E., Gale, A., & Liu, F. (2022). Bridging the experiment-calculation divide: Machine learning corrections to redox potential calculations in implicit and explicit solvent models. *J. Chem. Theory Comput.* 🔗 https://doi.org/10.1021/acs.jctc.1c01040

Abstract: Prediction of redox potentials is essential for catalysis and energy storage. Although density functional theory (DFT) calculations have enabled rapid redox potential predictions for numerous compounds, prominent errors persist compared to experimental measurements. In this work, we develop machine learning (ML) models to reduce the errors of redox potential calculations in both implicit and explicit solvent models. Training and testing of the ML correction models are based on the diverse ROP313 dataset with experimental redox potentials measured for organic and

organometallic compounds in a variety of solvents. For the implicit solvent approach, our ML models can reduce both the systematic bias and the number of outliers. ML corrected redox potentials also demonstrate less sensitivity to DFT functional choice. For the explicit solvent approach, we significantly reduce the computational costs by embedding the microsolvated cluster in implicit bulk solvent, obtaining converged redox potential results with a smaller solvation shell. This combined implicit-explicit solvent model, together with GPU-accelerated quantum chemistry methods, enabled rapid generation of a large dataset of explicit-solvent-calculated redox potentials for 165 organic compounds, allowing detailed investigation of the error sources in explicit solvent redox potential calculations.

6   **Hruska**, E., & Liu, F. (2022). *Quantum chemistry in the age of machine learning, chapter 6: Machine learning: An overview*. Elsevier.

Abstract: In this chapter, we give an overview of machine learning in the context of chemistry applications. Four types of machine learning are introduced together with corresponding quantum chemistry applications: supervised learning, unsupervised learning, semi-supervised learning, and reinforcement learning. We also present some basic concepts, techniques, and terminologies related to practical machine learning applications. Finally, we demonstrate the machine learning concepts with case studies focused on predicting the multi-reference character of molecules.

7   **Hruska**, E., Zhao, L., & Liu, F. (2022). Ground truth explanation dataset for chemical property prediction on molecular graphs. *Preprint*.
    🔗 https://doi.org/10.26434/chemrxiv-2022-96slq-v2

Abstract: Interpretation of chemistry on an atomic scale improves with explainable artificial intelligence (XAI). The parts of the molecule with the most significant influence on the chemical property of interest can be visualized with atomwise and bondwise attributions. Nonetheless, the attributions from different XAI methods regularly disagree substantially, causing uncertainty about which explanation is correct. To determine a ground truth for attributions, we define chemical operations which avoid alchemical steps or approximations and allow extracting one attribution per atom or bond from existing datasets of chemical properties. This general procedure allows for generating large datasets of ground truth attributions. The approach allowed us to create a ground truth explanation dataset with more than 5 million data points for the HOMO-LUMO gap chemical property. This open-source dataset of atomistic ground truth explanations may serve as a reference for XAI approaches.

8   Gale, A., **Hruska**, E., & Liu, F. (2021). Quantum chemistry for molecules at extreme pressure on graphical processing units: Implementation of extreme-pressure polarizable continuum model. *J. Chem. Phys*, *154*, 244103. 🔗 https://doi.org/10.1063/5.0056480

Abstract: Pressure plays essential roles in chemistry by altering structures and controlling chemical reactions. The extreme-pressure polarizable continuum model (XP-PCM) is an emerging method with an efficient quantum mechanical description of small- and medium-sized molecules at high pressure (on the order of GPa). However, its application to large molecular systems was previously hampered by a CPU computation bottleneck: the Pauli repulsion potential unique to XP-PCM requires the evaluation of a large number of electric field integrals, resulting in significant computational overhead compared to the gas-phase or standard-pressure polarizable continuum model calculations. Here, we exploit advances in graphical processing units (GPUs) to accelerate the XP-PCM-integral evaluations. This enables high-pressure quantum chemistry simulation of proteins that used to be computationally intractable. We benchmarked the performance using 18 small proteins in aqueous solutions. Using a single GPU, our method evaluates the XP-PCM free energy of a protein with over 500 atoms and 4000

basis functions within half an hour. The time taken by the XP-PCM-integral evaluation is typically 1% of the time taken for a gas-phase density functional theory (DFT) on the same system. The overall XP-PCM calculations require less computational effort than that for their gas-phase counterpart due to the improved convergence of self-consistent field iterations. Therefore, the description of the high-pressure effects with our GPU-accelerated XP-PCM is feasible for any molecule tractable for gas-phase DFT calculation. We have also validated the accuracy of our method on small molecules whose properties under high pressure are known from experiments or previous theoretical studies.

9   **Hruska**, **E.**, Balasubramanian, V., Lee, H., Jha, S., & Clementi, C. (2020). Extensible and scalable adaptive sampling on supercomputers. *J. Chem. Theory Comput.*
    🔗 https://doi.org/10.1021/acs.jctc.0c00991

Abstract: The accurate sampling of protein dynamics is an ongoing challenge despite the utilization of high-performance computer (HPC) systems. Utilizing only "brute force"molecular dynamics (MD) simulations requires an unacceptably long time to solution. Adaptive sampling methods allow a more effective sampling of protein dynamics than standard MD simulations. Depending on the restarting strategy, the speed up can be more than 1 order of magnitude. One challenge limiting the utilization of adaptive sampling by domain experts is the relatively high complexity of efficiently running adaptive sampling on HPC systems. We discuss how the ExTASY framework can set up new adaptive sampling strategies and reliably execute resulting workflows at scale on HPC platforms. Here, the folding dynamics of four proteins are predicted with no a priori information.

10  **Hruska**, **E.**, Abella, J. R., Nüske, F., Kavraki, L. E., & Clementi, C. (2018). Quantitative comparison of adaptive sampling methods for protein dynamics. *J. Chem. Phys.*, *149*(24), 244119. 🔗 https://doi.org/10.1063/1.5053582

Abstract: Adaptive sampling methods, often used in combination with Markov state models, are becoming increasingly popular for speeding up rare events in simulation such as molecular dynamics (MD) without biasing the system dynamics. Several adaptive sampling strategies have been proposed, but it is not clear which methods perform better for different physical systems. In this work, we present a systematic evaluation of selected adaptive sampling strategies on a wide selection of fast folding proteins. The adaptive sampling strategies were emulated using models constructed on already existing MD trajectories. We provide theoretical limits for the sampling speed-up and compare the performance of different strategies with and without using some a priori knowledge of the system. The results show that for different goals, different adaptive sampling strategies are optimal. In order to sample slow dynamical processes such as protein folding without a priori knowledge of the system, a strategy based on the identification of a set of metastable regions is consistently the most efficient, while a strategy based on the identification of microstates performs better if the goal is to explore newer regions of the conformational space. Interestingly, the maximum speed-up achievable for the adaptive sampling of slow processes increases for proteins with longer folding times, encouraging the application of these methods for the characterization of slower processes, beyond the fast-folding proteins considered here.

11  Balasubramanian, V., Bethune, I., Shkurti, A., Breitmoser, E., **Hruska**, **E.**, Clementi, C., Laughton, C., & Jha, S. (2016). Extasy: Scalable and flexible coupling of md simulations and advanced sampling techniques, 361–370.
    🔗 https://doi.org/10.1109/eScience.2016.7870921

Abstract: For many macromolecular systems the accurate sampling of the relevant regions on the potential energy surface cannot be obtained by a single, long Molecular Dynamics (MD) trajectory. New approaches are required to promote more efficient sampling. We present the design and

implementation of the Extensible Toolkit for Advanced Sampling and analYsis (ExTASY) for building and executing advanced sampling workflows on HPC systems. ExTASY provides Python based "templated scripts" that interface to an interoperable and high-performance pilot-based run time system, which abstracts the complexity of managing multiple simulations. ExTASY supports the use of existing highly-optimised parallel MD code and their coupling to analysis tools based upon collective coordinates which do not require a priori knowledge of the system to bias. We describe two workflows which both couple large "ensembles" of relatively short MD simulations with analysis tools to automatically analyse the generated trajectories and identify molecular conformational structures that will be used on-the-fly as new starting points for further "simulation-analysis" iterations. One of the workflows leverages the Locally Scaled Diffusion Maps technique; the other makes use of Complementary Coordinates techniques to enhance sampling and generate start-points for the next generation of MD simulations. We show that the ExTASY tools have been deployed on a range of HPC systems including ARCHER (Cray CX30), Blue Waters (Cray XE6/XK7), and Stampede (Linux cluster), and that good strong scaling can be obtained up to 1000s of MD simulations, independent of the size of each simulation. We discuss how ExTASY can be easily extended or modified by end-users to build their own workflows, and ongoing work to improve the usability and robustness of ExTASY.