

Supporting information to: Machine-learning to predict anharmonic frequencies: a study of models and transferability

Basis set information of the ML dataset

Table S1: Clusters used for ML dataset and corresponding basis set information (-r = ring; -c = chain; -l = linear; -rb = boat ring).

| Cluster | basis set (def2-) | Cluster | basis set (def2-) |
|-----------------------------------|-------------------|-----------------------------------|-------------------|
| HF | TZVP | (HBr) ₂ | TZVP |
| (HF) ₂ | TZVP | (HBr) _{3-r} | SVPD |
| (HF) _{3-r} | TZVP | (HBr) _{3-c} | SVPD |
| (HF) _{3-c} ^a | SVPD | (HBr) _{4-r} | SVP |
| (HF) _{3-cl} ^b | TZVP | CH ₃ F | TZVP |
| (HF) _{4-r} | TZVP | (CH ₃ F) ₂ | SVPD |
| (HF) _{4-rb} | TZVP | (CH ₃ F) ₂ | SVP |
| (HF) _{5-r} | TZVP | CH ₃ Cl | TZVP |
| HCl | TZVP | (CH ₃ Cl) ₂ | SVP |
| (HCl) ₂ | TZVP | CH ₃ Br | TZVP |
| (HCl) _{3-r} | TZVP | (CH ₃ Br) ₂ | TZVP |
| (HCl) _{3-c} | SVPD | C ₂ H ₅ F | TZVP |
| (HCl) _{4-r} | SVPD | C ₂ H ₅ Cl | TZVP |
| (HCl) _{4-rb} | SVPD | C ₂ H ₅ Br | TZVP |
| HBr | TZVP | | |

S1 Quantitative analysis of the ML model

Convergence of the GBR model

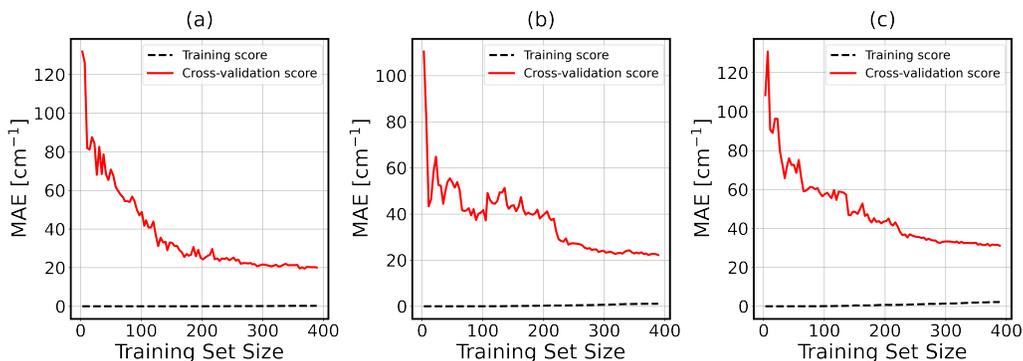


Figure S1: Learning curves of the gradient boosting regression leave-one-out cross validation for (a) Diagonal, (b) VSCF and (c) VSCF-PT2 frequencies

Convergence of the MLR model

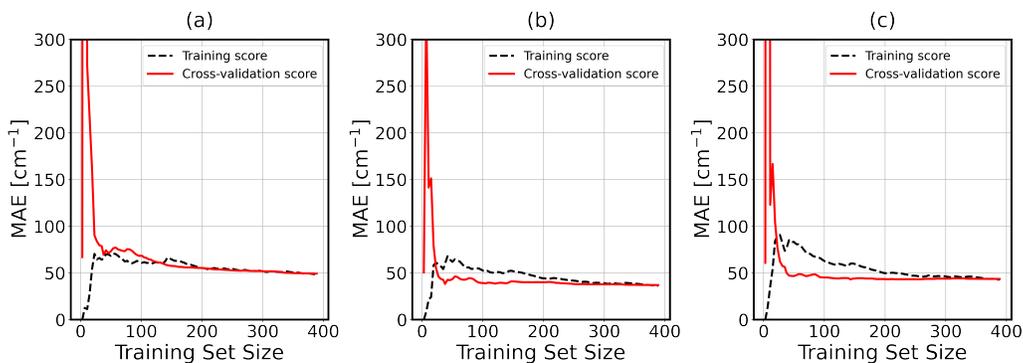


Figure S2: Learning curves of the linear regression leave-one-out cross validation for (a) Diagonal, (b) VSCF and (c) VSCF-PT2 frequencies

MLR predictions

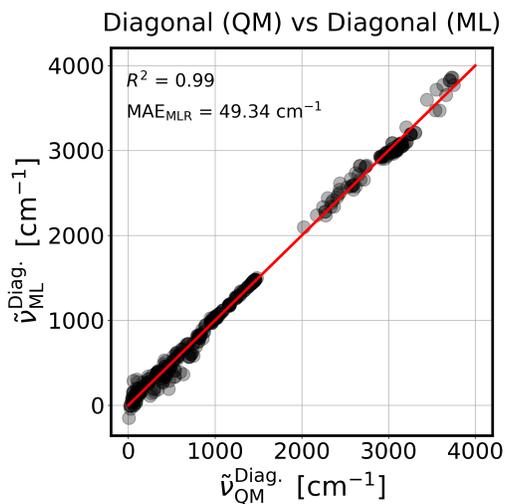


Figure S3: QM calculated Diagonal frequencies ($\tilde{\nu}_{QM}^{Diag.}$) vs. ML predicted Diagonal frequencies ($\tilde{\nu}_{ML}^{Diag.}$). $\tilde{\nu}_{ML}^{Diag.}$ are LOOCV linear regression predictions. Red line is ideal diagonal

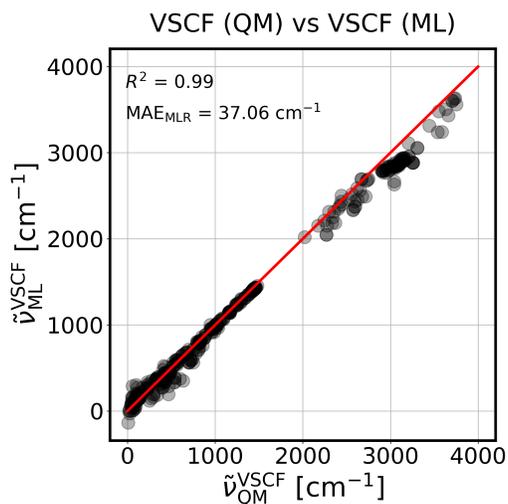


Figure S4: Same as Fig. S3 for VSCF frequencies. $\tilde{\nu}_{ML}^{VSCF}$ are LOOCV linear regression predictions. Red line is ideal diagonal

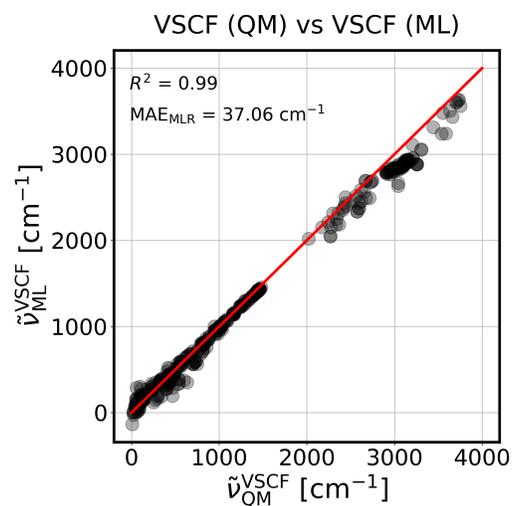


Figure S5: Same as Fig S3 for VSCF-PT2 frequencies. $\tilde{\nu}_{ML}^{\text{VSCF-PT2}}$ are LOOCV linear regression predictions. Red line is ideal diagonal

S2 Comparison of VSCF-calculated, AIMD and experimental^a spectra of CH₃F

AIMD computational details

Velocity-Verlet algorithm as implemented in NWChem was used for AIMD simulations. Electronic potential was calculated using density functional theory with def2-TZVP basis set and BLYP exchange correlation functional. Grimme's dispersion correction with Becke-Johnson damping D3(BJ) were employed to account for Van der Waals interactions. The nuclear time step was set to 0.2419 fs and in total 5000 steps were simulated. NVE ensemble, i.e. no thermostat with the temperature of 298.15 K, to generate initial velocities from Maxwell-Boltzmann distribution, were utilized. Corresponding IR spectrum was generated using NWChem postprocessing tool `qmd_analysis`.

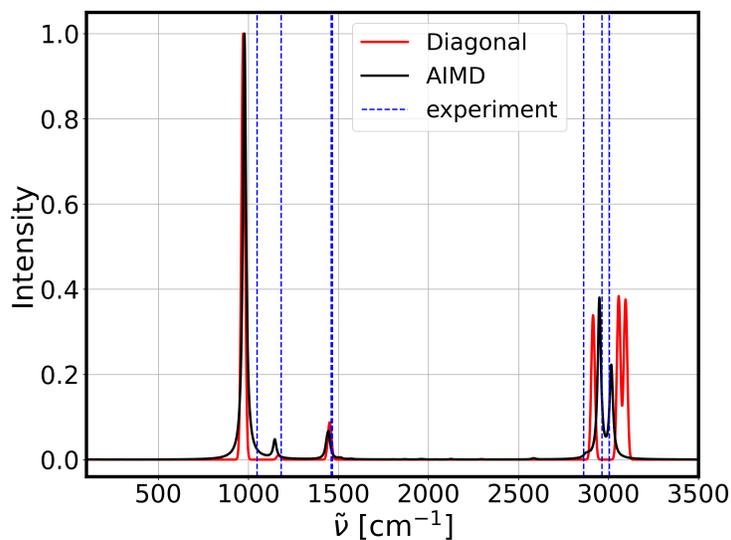


Figure S6: Diagonal and AIMD spectra of CH₃F; blue dashed vertical lines are experimental intensity peak locations

^aJ. Chem. Phys. 76, 809–816 (1982); J. Chem. Phys. 66, 970–975 (1977)

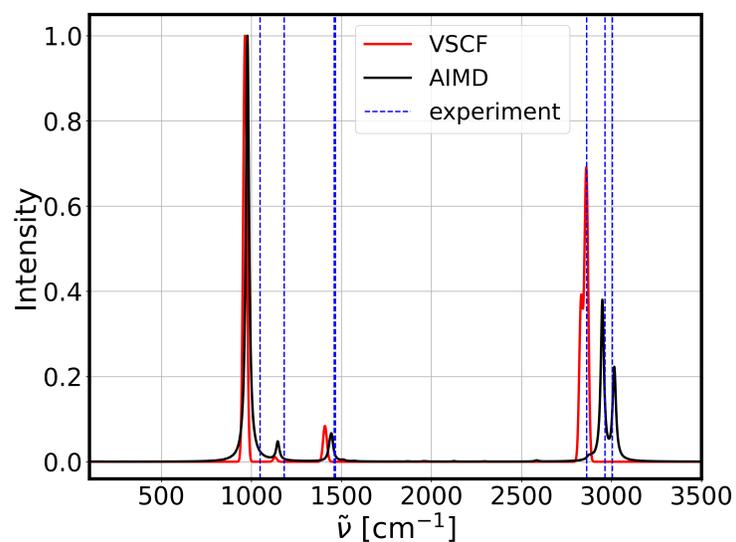


Figure S7: VSCF and AIMD spectra of CH_3F ; blue dashed vertical lines are experimental intensity peak locations

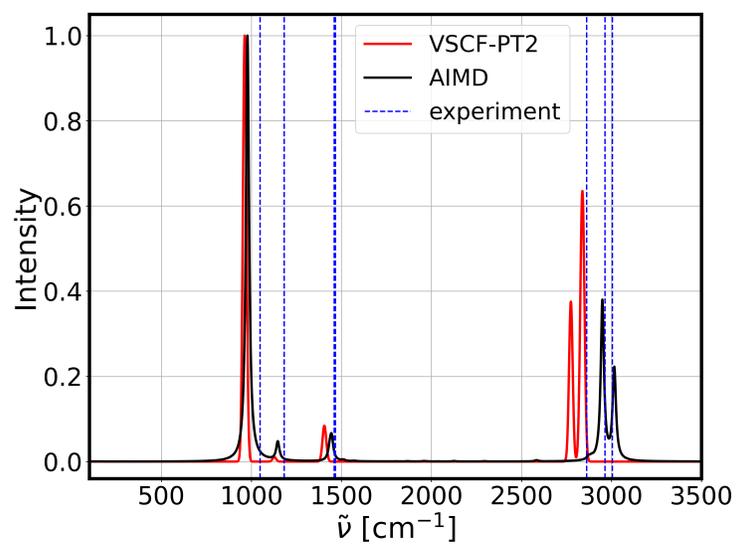


Figure S8: VSCF-PT2 and AIMD spectra of CH_3F ; blue dashed vertical lines are experimental intensity peak locations

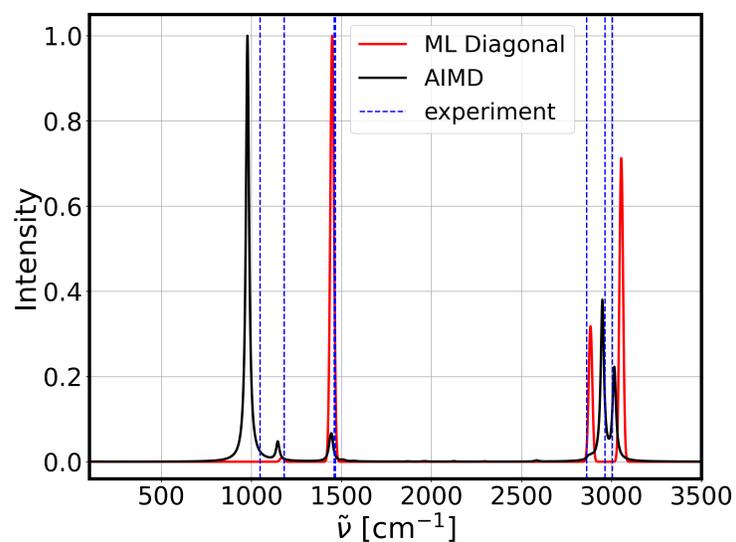


Figure S9: ML predicted Diagonal and AIMD spectra of CH_3F ; blue dashed vertical lines are experimental intensity peak locations

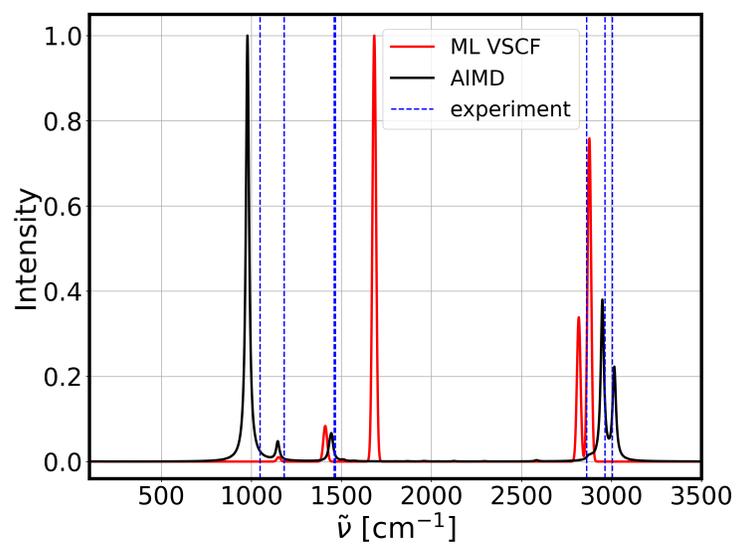


Figure S10: ML predicted VSCF and AIMD spectra of CH_3F ; blue dashed vertical lines are experimental intensity peak locations

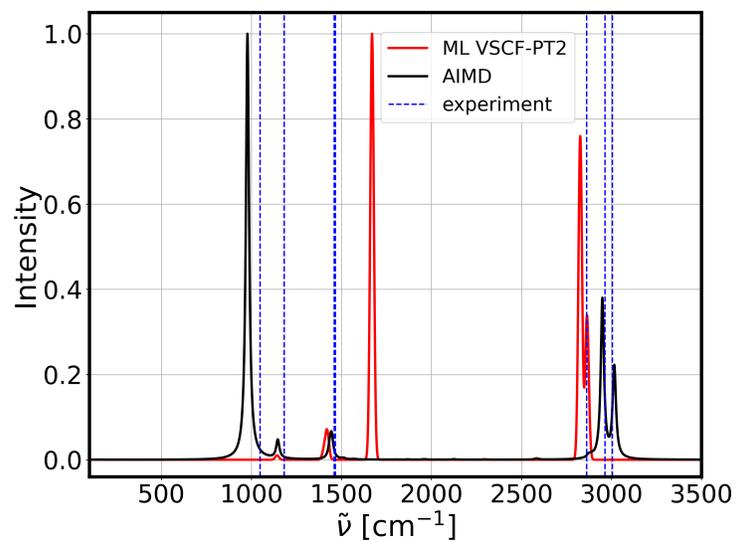


Figure S11: ML predicted VSCF-PT2 and AIMD spectra of CH₃F; blue dashed vertical lines are experimental intensity peak locations

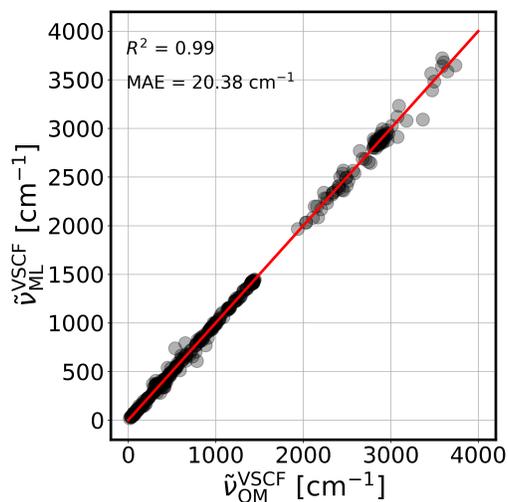


Figure S12: QM calculated vs. ML predicted VSCF frequencies. Here “Diagonal - VSCF” shift is used as target and Diagonal frequencies are included to the original harmonic based descriptor set.

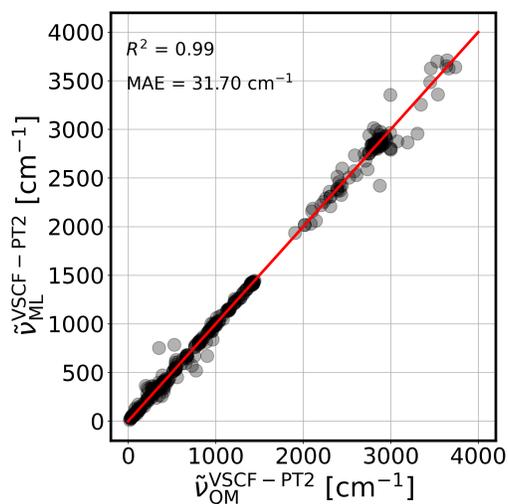


Figure S13: QM calculated vs. ML predicted VSCF-PT2 frequencies. Here “Diagonal - VSCF-PT2” shift is used as target and Diagonal frequencies are included to the original harmonic based descriptor set.

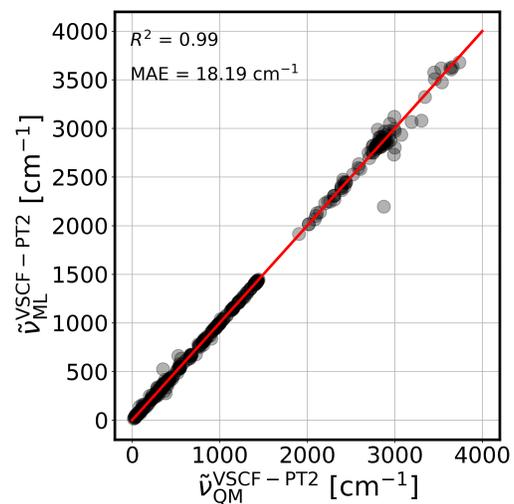


Figure S14: QM calculated vs. ML predicted VSCF-PT2 frequencies. Here “VSCF - VSCF-PT2” shift is used as target and VSCF frequencies are included to the original harmonic based descriptor set.

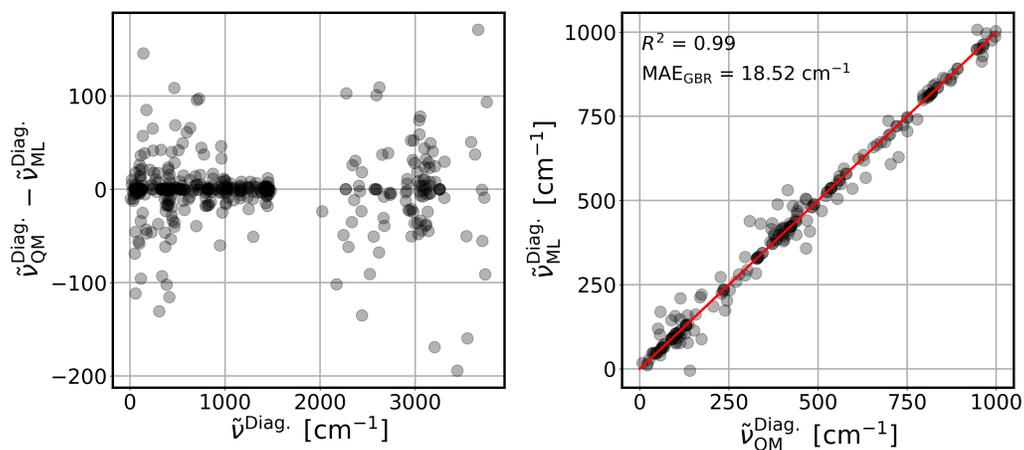


Figure S15: Left panel: difference between ML predicted and QM calculated Diagonal frequencies vs. Diagonal frequencies; Right panel: predictions of the Diagonal frequencies below 1000 cm^{-1}

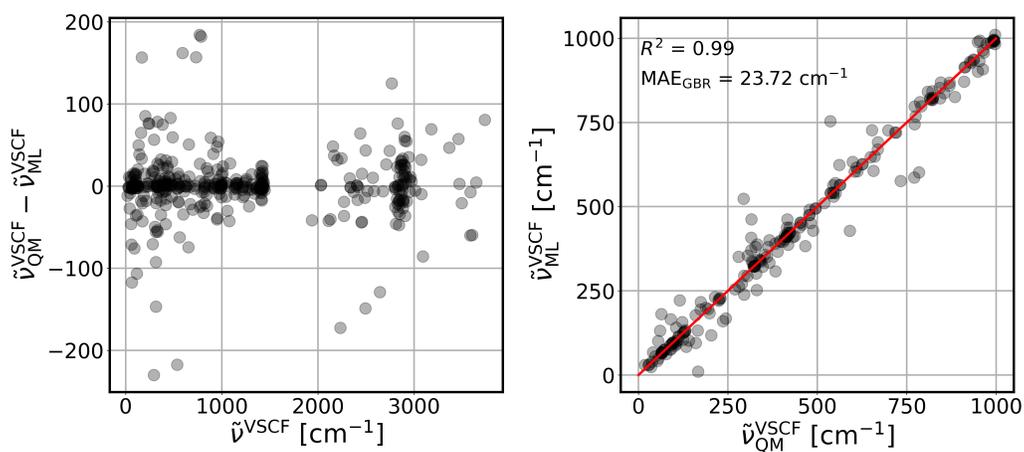


Figure S16: Left panel: difference between ML predicted and QM calculated VSCF frequencies vs. VSCF frequencies; Right panel: predictions of the VSCF frequencies below 1000 cm^{-1}

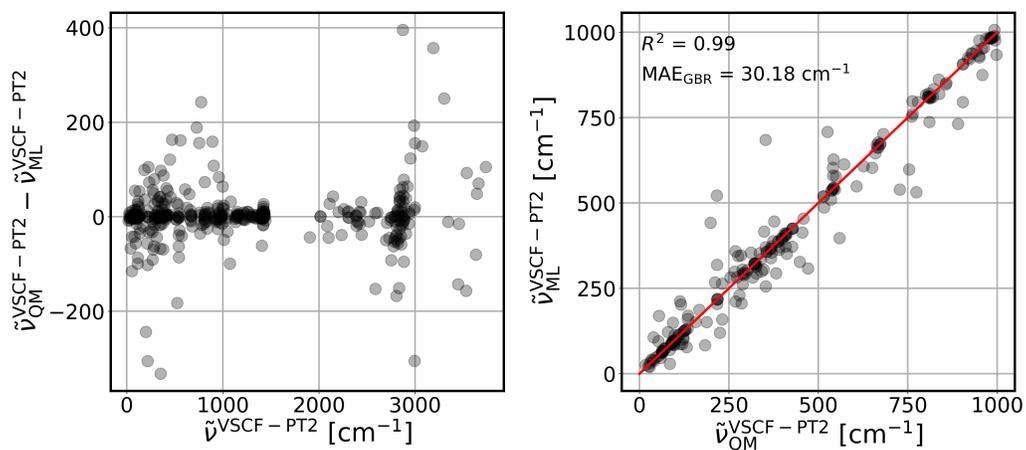


Figure S17: Left panel: difference between ML predicted and QM calculated VSCF-PT2 frequencies vs. VSCF-PT2 frequencies; Right panel: predictions of the VSCF-PT2 frequencies below 1000 cm^{-1}

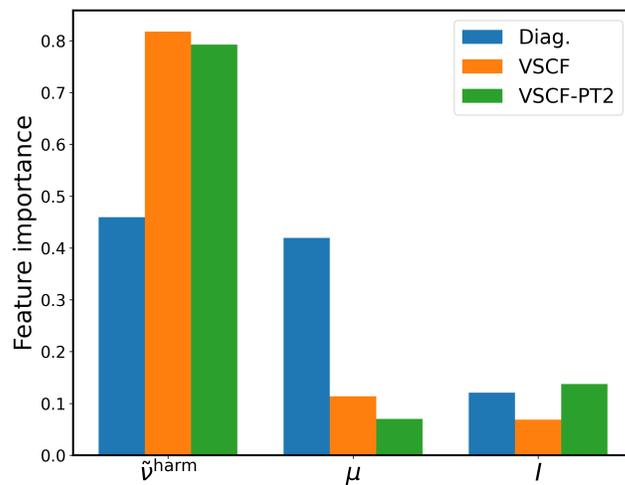


Figure S18: Feature importance of the descriptor set for 3 anharmonic frequency types. Here $\tilde{\nu}^{\text{harm}}$ - harmonic frequency, μ - reduced mass and I - intensity