

## Supplementary Information

### Distribution-preserved sampling (DPS) for smarter machine learning assisted ultra-large-scale virtual screening

Alexander Trachtenberg, Alexander Spelkov, and Barak Akabayov\*

Department of Chemistry and Data Science Research Center, Ben-Gurion University of the Negev, Beer-Sheva 8410501, Israel.

\*e-mail:akabayov@bgu.ac.il

#### Brief description of chemical descriptors in RDKit

(<https://github.com/milanimarcel/Descriptors-from-RDKit/blob/master/README.md>):

- **MaxEStateIndex**: Returns a tuple of EState indices for the molecule, Reference: Hall, Mohny and Kier. *JCICS 31* 76-81 (1991)
- **MinEStateIndex**: Returns a tuple of EState indices for the molecule, Reference: Hall, Mohny and Kier. *JCICS 31* 76-81 (1991)
- **MaxAbsEStateIndex**: Returns a tuple of EState indices for the molecule, Reference: Hall, Mohny and Kier. *JCICS 31* 76-81 (1991)
- **MinAbsEStateIndex**: Returns a tuple of EState indices for the molecule, Reference: Hall, Mohny and Kier. *JCICS 31* 76-81 (1991)
- **QED**: QED stands for quantitative amount of similarity of drugs and the concept was first introduced by Richard Bickerton and his colleagues. The empirical logic of the QED measure reflects the underlying distribution of molecular properties, including molecular weight, logP, topological polar surface area, number of hydrogen scavenger donors and acceptors, number of aromatic rings and rotating screens and the presence of unwanted chemicals.
- **SPS**: An empirical measure of molecular three-dimensional complexity and spatial arrangement, useful for assessing scaffold diversity
- **MolWt**: The average molecular weight of the molecule
- **HeavyAtomMolWt**: The average molecular weight of the molecule without the hydrogens
- **ExactMolWt**: The exact molecular weight of the molecule
- **NumValenceElectrons**: The number of valence electrons of the molecule
- **NumRadicalElectrons**: The number of radical electrons of the molecule (says nothing about spin state)
- **MaxPartialCharge**: A partial charge is a non-integer charge value when measured in elementary charge units. Partial charge is more commonly called net atomic charge. It is represented by the Greek lowercase letter  $\delta$ , namely, as  $\delta^-$  or  $\delta^+$ . Partial charges are created due to the asymmetric distribution of electrons in chemical bonds.
- **MinPartialCharge**: Min partial Charge
- **MaxAbsPartialCharge**: Returns molecular charge descriptors
- **MinAbsPartialCharge**: Returns molecular charge descriptors
- **FpDensityMorgan1**: Morgan fingerprint density
- **FpDensityMorgan2**: Morgan fingerprint density
- **FpDensityMorgan3**: Morgan fingerprint density
- **BCUT2D\_MWHI / BCUT2D\_MWLOW**: High and low eigenvalues of a weighted topological matrix derived using atomic molecular weights, capturing extremes in molecular-weight-based topology

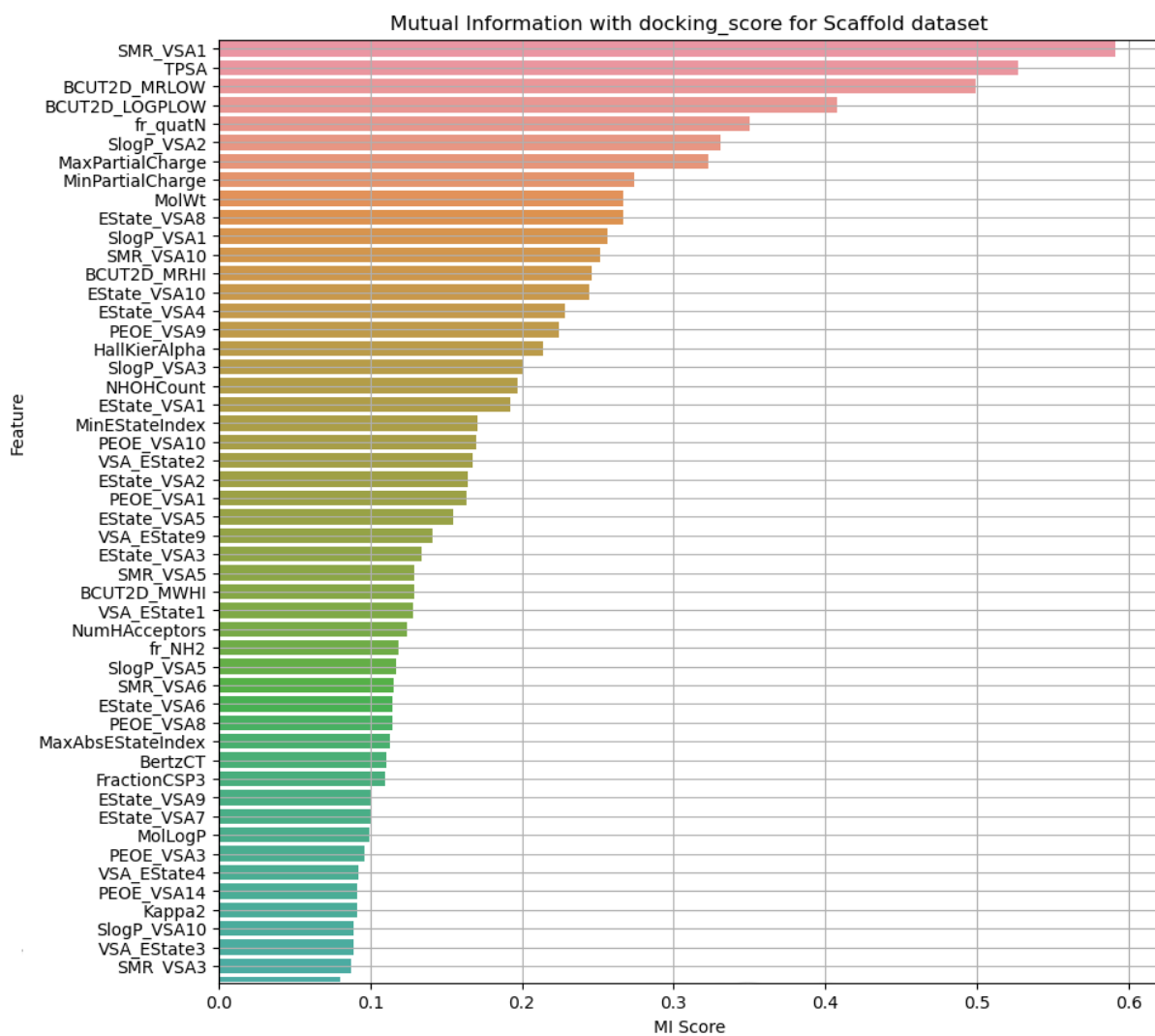
- **BCUT2D\_CHGHI / BCUT2D\_CHGLO**: High and low eigenvalues based on partial atomic charges, representing charge distribution extremes
- **BCUT2D\_LOGPHI / BCUT2D\_LOGPLOW**: Eigenvalue descriptors using LogP contributions, capturing hydrophobicity extremes across the molecule
- **BCUT2D\_MRHI / BCUT2D\_MRLOW**: High and low eigenvalues derived from atomic molar refractivity values, reflecting polarizability extremes
- **AvgIpc**: Average Inner Polar Surface Area - An average per-atom inner polar surface measure, indicative of polarity distribution across the molecule
- **BalabanJ**: Balaban's J value for a molecule, Chem. Phys. Lett. 89:399-404 (1982).
- **BertzCT**: "A topological index meant to quantify "complexity" of molecules. J. Am. Chem. Soc. 103:3599-601 (1981)."
- **Chi0**: From equations (1),(9) and (10) of Rev. Comp. Chem. vol 2, 367-422, (1991)
- **Chi0n**: Similar to Hall Kier Chi0v, but uses nVal instead of valence This makes a big difference after we get out of the first row. Rev. Comput. Chem. 2:367-422 (1991).
- **Chi0v**: From equations (5),(9) and (10) of Rev. Comp. Chem. vol 2, 367-422, (1991)
- **Chi1**: From equations (1),(11) and (12) of Rev. Comp. Chem. vol 2, 367-422, (1991)
- **Chi1n**: Similar to Hall Kier Chi1v, but uses nVal instead of valence. Rev. Comput. Chem. 2:367-422 (1991).
- **Chi1v**: From equations (5),(11) and (12) of Rev. Comp. Chem. vol 2, 367-422, (1991)
- **Chi2n**: Similar to Hall Kier Chi2v, but uses nVal instead of valence This makes a big difference after we get out of the first row. Rev. Comput. Chem. 2:367-422 (1991).
- **Chi2v**: From equations (5),(15) and (16) of Rev. Comp. Chem. vol 2, 367-422, (1991)
- **Chi3n**: Similar to Hall Kier Chi3v, but uses nVal instead of valence This makes a big difference after we get out of the first row. Rev. Comput. Chem. 2:367-422 (1991).
- **Chi3v**: From equations (5),(15) and (16) of Rev. Comp. Chem. vol 2, 367-422, (1991)
- **Chi4n**: Similar to Hall Kier Chi4v, but uses nVal instead of valence. This makes a big difference after we get out of the first row. Rev. Comput. Chem. 2:367-422 (1991).
- **Chi4v**: From equations (5),(15) and (16) of Rev. Comp. Chem. vol 2, 367-422, (1991)
- **HallKierAlpha**: The Hall-Kier alpha value for a molecule. Rev. Comput. Chem. 2:367-422 (1991).
- **Ipc**: the information content of the coefficients of the characteristic polynomial of the adjacency matrix of a hydrogen-suppressed graph of a molecule.
- **Kappa1**: Kappa indices are calculated relative to the least branched (linear) and most branched (star) compounds with the same number of atoms as the molecule being investigated.
- **Kappa2**: Hall-Kier Kappa2 value
- **Kappa3**: Hall-Kier Kappa3 value
- **LabuteASA**: Labute's Approximate Surface Area (ASA from MOE)
- **PEOE\_VSA1**: MOE (Molecular Operating Environment) Charge VSA Descriptor 1 ( $-\infty < x < -0.30$ )
- **PEOE\_VSA10**: MOE Charge VSA Descriptor
- **PEOE\_VSA11**: MOE Charge VSA Descriptor 11 ( $0.15 \leq x < 0.20$ )
- **PEOE\_VSA12**: MOE Charge VSA Descriptor 12 ( $0.20 \leq x < 0.25$ )
- **PEOE\_VSA13**: MOE Charge VSA Descriptor 13 ( $0.25 \leq x < 0.30$ )
- **PEOE\_VSA14**: MOE Charge VSA Descriptor 14 ( $0.30 \leq x < \infty$ )
- **PEOE\_VSA2**: MOE Charge VSA Descriptor 2 ( $-0.30 \leq x < -0.25$ )
- **PEOE\_VSA3**: MOE Charge VSA Descriptor 3 ( $-0.25 \leq x < -0.20$ )
- **PEOE\_VSA4**: MOE Charge VSA Descriptor 4 ( $-0.20 \leq x < -0.15$ )

- **PEOE\_VSA5:** MOE Charge VSA Descriptor 5 ( $-0.15 \leq x < -0.10$ )
- **PEOE\_VSA6:** MOE Charge VSA Descriptor 6 ( $-0.10 \leq x < -0.05$ )
- **PEOE\_VSA7:** MOE Charge VSA Descriptor 7 ( $-0.05 \leq x < 0.00$ )
- **PEOE\_VSA8:** MOE Charge VSA Descriptor 8 ( $0.00 \leq x < 0.05$ )
- **PEOE\_VSA9:** MOE Charge VSA Descriptor 9 ( $0.05 \leq x < 0.10$ )
- **SMR\_VSA1:** MOE MR VSA Descriptor 1 ( $-\infty < x < 1.29$ )
- **SMR\_VSA10:** MOE MR VSA Descriptor 10
- **SMR\_VSA2:** MOE MR VSA Descriptor 2 ( $1.29 \leq x < 1.82$ )
- **SMR\_VSA3:** MOE MR VSA Descriptor 3 ( $1.82 \leq x < 2.24$ )
- **SMR\_VSA4:** MOE MR VSA Descriptor 4 ( $2.24 \leq x < 2.45$ )
- **SMR\_VSA5:** MOE MR VSA Descriptor 5 ( $2.45 \leq x < 2.75$ )
- **SMR\_VSA6:** MOE MR VSA Descriptor 6 ( $2.75 \leq x < 3.05$ )
- **SMR\_VSA7:** MOE MR VSA Descriptor 7 ( $3.05 \leq x < 3.63$ )
- **SMR\_VSA8:** MOE MR VSA Descriptor 8 ( $3.63 \leq x < 3.80$ )
- **SMR\_VSA9:** MOE MR VSA Descriptor 9 ( $3.80 \leq x < 4.00$ )
- **SlogP\_VSA1:** MOE logP VSA Descriptor 1 ( $-\infty < x < -0.40$ )
- **SlogP\_VSA10:** MOE logP VSA Descriptor 10 ( $0.40 \leq x < 0.50$ )
- **SlogP\_VSA11:** MOE logP VSA Descriptor 11 ( $0.50 \leq x < 0.60$ )
- **SlogP\_VSA12:** MOE logP VSA Descriptor 12 ( $0.60 \leq x < \infty$ )
- **SlogP\_VSA2:** MOE logP VSA Descriptor 2 ( $-0.40 \leq x < -0.20$ )
- **SlogP\_VSA3:** MOE logP VSA Descriptor 3 ( $-0.20 \leq x < 0.00$ )
- **SlogP\_VSA4:** MOE logP VSA Descriptor 4 ( $0.00 \leq x < 0.10$ )
- **SlogP\_VSA5:** MOE logP VSA Descriptor 5 ( $0.10 \leq x < 0.15$ )
- **SlogP\_VSA6:** MOE logP VSA Descriptor 6 ( $0.15 \leq x < 0.20$ )
- **SlogP\_VSA7:** MOE logP VSA Descriptor 7 ( $0.20 \leq x < 0.25$ )
- **SlogP\_VSA8:** MOE logP VSA Descriptor 8 ( $0.25 \leq x < 0.30$ )
- **SlogP\_VSA9:** MOE logP VSA Descriptor 9 ( $0.30 \leq x < 0.40$ )
- **TPSA:** The polar surface area (PSA) or topological polar surface area (TPSA) of a molecule is defined as the surface sum over all polar atoms or molecules, primarily oxygen and nitrogen, also including their attached hydrogen atoms. PSA is a commonly used medicinal chemistry metric for the optimization of a drug's ability to permeate cells.
- **EState\_VSA1:** EState VSA Descriptor 1 ( $-\infty < x < -0.39$ )
- **EState\_VSA10:** EState VSA Descriptor 10 ( $9.17 \leq x < 15.00$ )
- **EState\_VSA11:** EState VSA Descriptor 11 ( $15.00 \leq x < \infty$ )
- **EState\_VSA2:** EState VSA Descriptor 2 ( $-0.39 \leq x < 0.29$ )
- **EState\_VSA3:** EState VSA Descriptor 3 ( $0.29 \leq x < 0.72$ )
- **EState\_VSA4:** EState VSA Descriptor 4 ( $0.72 \leq x < 1.17$ )
- **EState\_VSA5:** EState VSA Descriptor 5 ( $1.17 \leq x < 1.54$ )
- **EState\_VSA6:** EState VSA Descriptor 6 ( $1.54 \leq x < 1.81$ )
- **EState\_VSA7:** EState VSA Descriptor 7 ( $1.81 \leq x < 2.05$ )
- **EState\_VSA8:** EState VSA Descriptor 8 ( $2.05 \leq x < 4.69$ )
- **EState\_VSA9:** EState VSA Descriptor 9 ( $4.69 \leq x < 9.17$ )
- **VSA\_EState1:** VSA EState Descriptor 1 ( $-\infty < x < 4.78$ )
- **VSA\_EState10:** VSA EState Descriptor 10 ( $11.00 \leq x < \infty$ )
- **VSA\_EState2:** VSA EState Descriptor 2 ( $4.78 \leq x < 5.00$ )
- **VSA\_EState3:** VSA EState Descriptor 3 ( $5.00 \leq x < 5.41$ )
- **VSA\_EState4:** VSA EState Descriptor 4 ( $5.41 \leq x < 5.74$ )
- **VSA\_EState5:** VSA EState Descriptor 5 ( $5.74 \leq x < 6.00$ )

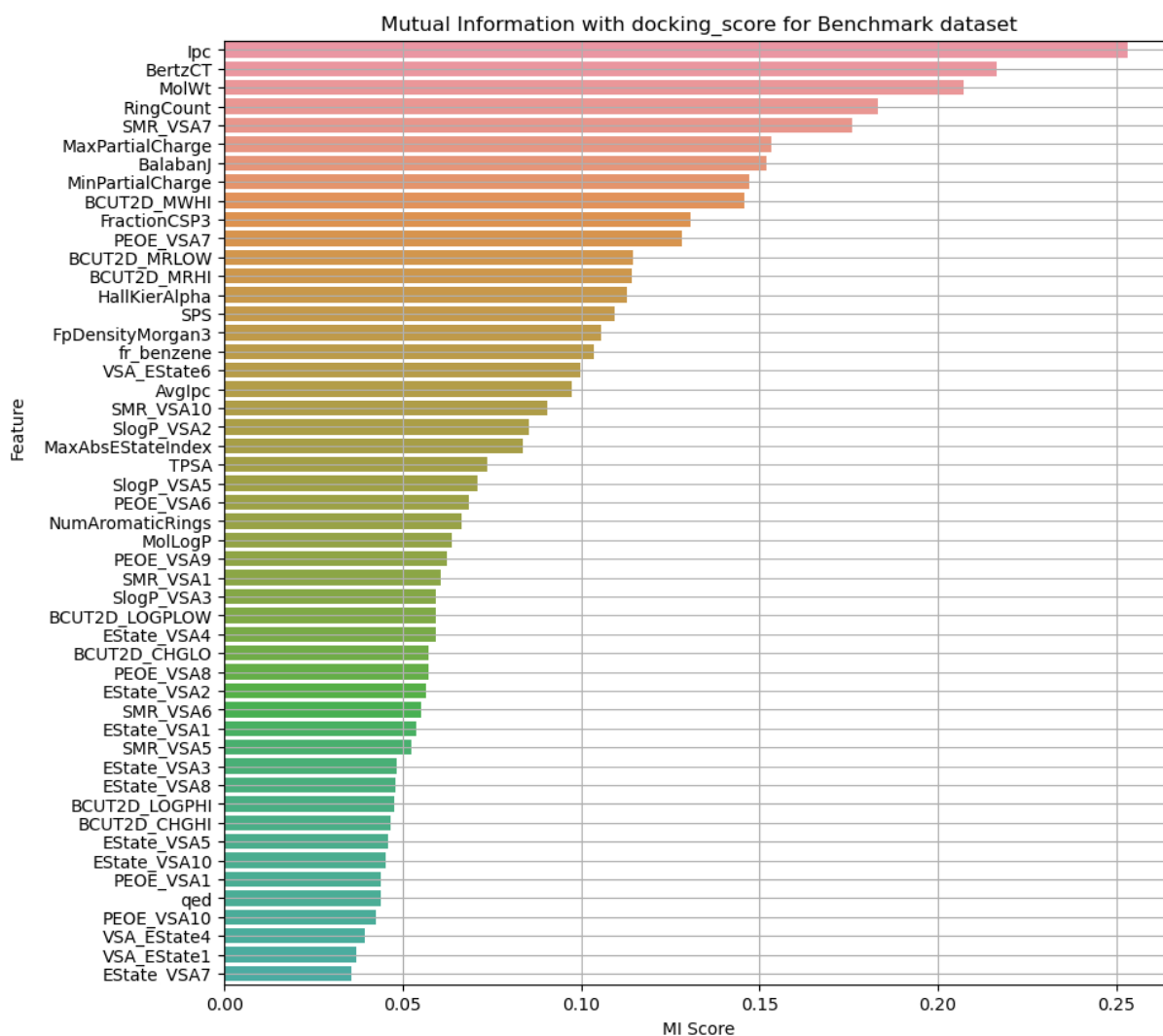
- **VSA\_EState6:** VSA EState Descriptor 6 (  $6.00 \leq x < 6.07$  )
- **VSA\_EState7:** VSA EState Descriptor 7 (  $6.07 \leq x < 6.45$  )
- **VSA\_EState8:** VSA EState Descriptor 8 (  $6.45 \leq x < 7.00$  )
- **VSA\_EState9:** VSA EState Descriptor 9 (  $7.00 \leq x < 11.00$  )
- **FractionCSP3:** The fraction of C atoms that are SP3 hybridized.
- **HeavyAtomCount:** Number of heavy atoms of a molecule.
- **NHOHCount:** Number of -NH and -OH groups
- **NOCCount:** Number of nitrogen and oxygen atoms
- **NumAliphaticCarbocycles:** returns the number of aliphatic (containing at least one non-aromatic bond) carbocycles for a molecule
- **NumAliphaticHeterocycles:** returns the number of aliphatic (containing at least one non-aromatic bond) heterocycles for a molecule
- **NumAliphaticRings:** returns the number of aliphatic (containing at least one non-aromatic bond) rings for a molecule
- **NumAmideBonds:** The total count of amide bonds in the molecule, suggesting potential hydrogen-bonding capacity
- **NumAromaticCarbocycles:** The number of aromatic carbocycles for a molecule
- **NumAromaticHeterocycles:** The number of aromatic rings containing heteroatoms (e.g., N, O, S), indicative of key pharmacophores
- **NumAromaticRings:** The number of aromatic rings for a molecule
- **NumAtomStereoCenters:** Count of stereogenic atoms with defined chirality centers
- **NumBridgeheadAtoms:** Number of atoms at bridgehead positions in polycyclic structures, reflecting topological complexity
- **NumHAcceptors:** Number of Hydrogen Bond Acceptors
- **NumHDonors:** Number of Hydrogen Bond Donors
- **NumHeteroatoms:** Number of Heteroatoms
- **NumHeterocycles:** Total number of ring systems containing at least one heteroatom
- **NumRotatableBonds:** Number of Rotatable Bonds
- **NumSaturatedCarbocycles:** returns the number of saturated carbocycles for a molecule
- **NumSaturatedHeterocycles:** returns the number of saturated heterocycles for a molecule
- **NumSaturatedRings:** Number of Saturated Rings
- **NumSpiroAtoms:** Count of atoms that serve as shared nodes in spiro ring systems
- **NumUnspecifiedAtomStereoCenters:** Count of stereocenters present in the molecule whose chiral configuration is unspecified
- **RingCount:** Number of All Rings
- **Phi:** Flexibility Index / Kier Flexibility Index – A metric of overall molecular flexibility computed according to topological and valence properties
- **MolLogP:** octanol/water partition coefficient (Wildman-Crippen LogP value)
- **MolMR:** Molar Refractivity(Wildman-Crippen MR value)
- **fr\_Al\_COO:** Number of aliphatic carboxylic acids
- **fr\_Al\_OH:** Number of aliphatic hydroxyl groups
- **fr\_Al\_OH\_noTert:** Number of aliphatic hydroxyl groups excluding tert-OH
- **fr\_ArN:**Number of N functional groups attached to aromatics
- **fr\_Ar\_COO:** Number of Aromatic carboxylic acids
- **fr\_Ar\_N:** Number of aromatic nitrogens
- **fr\_Ar\_NH:** Number of aromatic amines
- **fr\_Ar\_OH:** Number of aromatic hydroxyl groups
- **fr\_COO:** Number of carboxylic acids

- **fr\_COO2**: Number of carboxylic acids
- **fr\_C\_O**: Number of carbonyl O
- **fr\_C\_O\_noCOO**: Number of carbonyl O, excluding COOH
- **fr\_C\_S**: Number of thiocarbonyl
- **fr\_HOCCN**: Number of C(OH)CCN-Ctert-alkyl or C(OH)CCNcyclic
- **fr\_Imine**: Number of Imines
- **fr\_NH0**: Number of Tertiary amines
- **fr\_NH1**: Number of Secondary amines
- **fr\_NH2**: Number of Primary amines
- **fr\_N\_O**: Number of hydroxylamine groups
- **fr\_Ndealkylation1**: Number of XCCNR groups
- **fr\_Ndealkylation2**: Number of tert-alicyclic amines (no heteroatoms, not quinine-like bridged N)
- **fr\_Nhpyrrole**: Number of H-pyrrole nitrogens
- **fr\_SH**: Number of thiol groups
- **fr\_aldehyde**: Number of aldehydes
- **fr\_alkyl\_carbamate**: Number of alkyl carbamates (subject to hydrolysis)
- **fr\_alkyl\_halide**: Number of alkyl halides
- **fr\_allylic\_oxid**: Number of allylic oxidation sites excluding steroid dienone
- **fr\_amide**: Number of amides
- **fr\_amidine**: Number of amidine groups
- **fr\_aniline**: Number of anilines
- **fr\_aryl\_methyl**: Number of aryl methyl sites for hydroxylation
- **fr\_azide**: Number of azide groups
- **fr\_azo**: Number of azo groups
- **fr\_barbitur**: Number of barbiturate groups
- **fr\_benzene**: Number of benzene rings
- **fr\_benzodiazepine**: Number of benzodiazepines with no additional fused rings
- **fr\_bicyclic**: Number of Bicyclic groups
- **fr\_diazo**: Number of diazo groups
- **fr\_dihydropyridine**: Number of dihydropyridines
- **fr\_epoxide**: Number of epoxide rings
- **fr\_ester**: Number of esters
- **fr\_ether**: Number of ether oxygens (including phenoxy)
- **fr\_furan**: Number of furan rings
- **fr\_guanido**: Number of guanidine groups
- **fr\_halogen**: Number of halogens
- **fr\_hdrzine**: Number of hydrazine groups
- **fr\_hdrzone**: Number of hydrazone groups
- **fr\_imidazole**: Number of imidazole rings
- **fr\_imide**: Number of imide groups
- **fr\_isocyan**: Number of isocyanates
- **fr\_isothiocyan**: Number of isothiocyanates
- **fr\_ketone**: Number of ketones
- **fr\_ketone\_Topless**: Number of ketones excluding diaryl, a,b-unsat. dienones, heteroatom on Calpha
- **fr\_lactam**: Number of beta lactams
- **fr\_lactone**: Number of cyclic esters (lactones)

- **fr\_methoxy**: Number of methoxy groups (-OCH<sub>3</sub>)
- **fr\_morpholine**: Number of morpholine rings
- **fr\_nitrile**: Number of nitriles
- **fr\_nitro**: Number of nitro groups
- **fr\_nitro\_ arom**: Number of nitro benzene ring substituents
- **fr\_nitro\_ arom\_ nonortho**: Number of non-ortho nitro benzene ring substituents
- **fr\_nitroso**: Number of nitroso groups, excluding NO<sub>2</sub>
- **fr\_oxazole**: Number of oxazole rings
- **fr\_oxime**: Number of oxime groups
- **fr\_para\_hydroxylation**: Number of *para*-hydroxylation sites
- **fr\_phenol**: Number of phenols
- **fr\_phenol\_noOrthoHbond**: Number of phenolic -OH groups, excluding *ortho* intramolecular H bond substituents
- **fr\_phos\_acid**: Number of phosphoric acid groups
- **fr\_phos\_ester**: Number of phosphoric ester groups
- **fr\_piperdine**: Number of piperdine rings
- **fr\_piperzine**: Number of piperzine rings
- **fr\_priamide**: Number of primary amides
- **fr\_prisulfonamd**: Number of primary sulfonamides
- **fr\_pyridine**: Number of pyridine rings
- **fr\_quatN**: Number of quarternary nitrogens
- **fr\_sulfide**: Number of thioether
- **fr\_sulfonamd**: Number of sulfonamides
- **fr\_sulfone**: Number of sulfone groups
- **fr\_term\_acetylene**: Number of terminal acetylenes
- **fr\_tetrazole**: Number of tetrazole rings
- **fr\_thiazole**: Number of thiazole rings
- **fr\_thiocyan**: Number of thiocyanates
- **fr\_thiophene**: Number of thiophene rings
- **fr\_unbrch\_alkane**: Number of unbranched alkanes of at least 4 members (excludes halogenated alkanes)
- **fr\_urea**: Number of urea groups



**Figure S1:** Top 50 RDKit chemical descriptors ranked by mutual information (MI) with docking scores for the **Scaffold**-based dataset.



**Figure S2:** Top 50 RDKit chemical descriptors ranked by mutual information (MI) with docking scores for **Benchmark** dataset.

Dataset	Splitting Method	Subset Fraction	RMSE	R <sup>2</sup>	Training Time [min.]
Scaffold	Random	1%	0.4736	0.8164	1.5220
Scaffold	Preserved	1%	0.4716	0.8179	1.3196
Benchmark	Random	1%	0.7343	0.5055	1.2515
Benchmark	Preserved	1%	0.7292	0.5124	1.0856
Scaffold	Random	5%	0.4421	0.8400	8.1097
Scaffold	Preserved	5%	0.4435	0.8389	7.0252
Benchmark	Random	5%	0.6839	0.5712	7.1124
Benchmark	Preserved	5%	0.6832	0.5720	6.2098
Scaffold	Random	10%	0.4309	0.8480	14.9277
Scaffold	Preserved	10%	0.4316	0.8474	14.1543
Benchmark	Random	10%	0.6626	0.5974	14.8106
Benchmark	Preserved	10%	0.6603	0.6002	13.2891
Scaffold	Random	25%	0.4209	0.8549	38.7931
Scaffold	Preserved	25%	0.4203	0.8553	34.9784
Benchmark	Random	25%	0.6328	0.6326	37.0249
Benchmark	Preserved	25%	0.6331	0.6325	33.9888
Scaffold	Random	50%	0.4173	0.8574	73.8511
Scaffold	Preserved	50%	0.4176	0.8572	71.2890
Benchmark	Random	50%	0.6104	0.6582	76.7545
Benchmark	Preserved	50%	0.6117	0.6571	70.2797
Scaffold	Random	75%	0.4170	0.8575	105.2672
Scaffold	Preserved	75%	0.4179	0.8569	102.3032
Benchmark	Random	75%	0.5974	0.6722	116.5927
Benchmark	Preserved	75%	0.5978	0.6725	114.8751

**Table S1:** RMSE, R<sup>2</sup> and training time in minutes for training a Random Forest model on different subsets of Scaffold and Benchmark datasets using different splitting methods.