## **Supporting information**

## DESignSolvents: An Open Platform for Search and Prediction of the Physicochemical Properties of Deep Eutectic Solvents

Valeria Odegova<sup>a</sup>, Anastasia Lavrinenko<sup>a</sup>, Timur Rakhmanov<sup>a,b</sup>, George Sysuev<sup>a</sup>, Andrei Dmitrenko<sup>a</sup>, Vladimir Vinogradov<sup>\*a</sup>

<sup>a</sup> International Institute "Solution Chemistry of Advanced Materials and Technologies", ITMO University, 191002 Saint Petersburg, Russia.

<sup>b</sup> Advanced Engineering School, Almetyevsk, Russia.

### Methods

Descriptors for machine learning were calculated using RDKit (2022.9.3) <sup>1</sup> based on the collected database. The full list of descriptors for each DES property, including their meanings and formulae, are given in Table S1 and Table S2. Several classical machine learning models were employed for prediction tasks, including decision tree, random forest, gradient boosting, k-nearest neighbors, support vector machines, and multilayer perceptron. We used open-source implementations from scikit-learn (0.0.post1) <sup>2</sup>, xgboost (1.7.3) <sup>3</sup> and catboost (1.1.1) <sup>4</sup> Python (3.9.11) packages. For hyperparameter optimization, we performed a grid search and Bayesian search from scikit-learn (0.0.post1) <sup>2</sup> and scikit-optimize (0.9.0) <sup>5</sup> Python packages, respectively. The best model was selected by evaluating R<sup>2</sup> and RMSE metrics (Table S3). Hyperparameters delivering the best performance for each prediction task are given in Table S4. The importance of descriptors was assessed using Shapley values calculated with the SHAP (0.41.0) <sup>6</sup> Python package. Finally, the DESignSolvents web service was developed with the following technological stack: Python (Django (4.1.7)), Javascript () (Vanilla, Plotly, noUiSlider), Docker <sup>7</sup>, Nginx <sup>8</sup>, Ofelia, and PostgreSQL <sup>9</sup>.



Figure. S1. Distribution for the temperature at which the density and viscosity were measured



Figure. S2. Distribution for melting point, density and viscosity depending on the type of DES

 Table S1. Descriptors for melting temperature prediction

Descriptor	Formula / A source	Meaning
Molar fraction of component	Experimental data	Molar fraction of individual component #1 (since only binary DES are used)
Type of DES	Experimental data	Type of DES (I, II, III, IV, V, IL)
Melting temperatures	Experimental data	Melting temperatures of individual components (in kelvin)
Molecular weight	$MW = MW1 * x_1 + MW2 * x_2 + MW * x_3,$ where $MW$ is the total molecular weight, $MW_m$ is the molecular weight in individual component m, and $x_m$ is the molar fraction of the individual component m.	The molecular weight a molecule calculated for the entire molecule, taking into account the molar fractions of individual components <sup>1</sup>
Number of hydrogen bond donors	$HD = HD1 * x_1 + HD2 * x_2 + HD3 * x_3,$ where $HD$ is the total number of H-bond donors, $HD_m$ is the number of H-bond donors in individual component m, and $x_m$ is the molar fraction of the individual component m.	Number of hydrogen bond donors calculated for the entire molecule, taking into account the molar fractions of individual components <sup>1</sup>
Number of different groups	$N_{Groups} = x_1 * n_1 + x_2 * n_2 + x_3 * n_3,$ where $N_{Groups}$ is the total number of different groups, $x_m$ is the molar fraction of the individual component, and $n_m$ is the number of different groups in the individual component.	The number of elements in DES, taking into account the mole fractions of individual components <sup>1</sup> Number of aliphatic carboxylic acids Number of aromatic carboxylic acid Number of aromatic nitrogens Number of aromatic hydroxyl groups Number of Tertiary amines Number of Secondary amines Number of amides
Number of aromatic rings	$AROM = AROM1 * x_1 + AROM2 * x_2 + AROM3 * x_3,$ where <i>AROM</i> is total number of aromatic rings, <i>AROM_m</i> is the number of aromatic rings in individual component m, and <sup>x</sup> <sub>m</sub> is the molar fraction of the individual component m.	
Chemical toxicity evaluation	$ALERTS = ALERTS1 * x_1 + ALERTS2 * x_2 + ALERTS3 * x_3,$ where <i>ALERTS</i> is the total chemical toxicity evaluation, <i>ALERTS_m</i> is the chemical toxicity evaluation in individual component m, and <sup>x_m</sup> is the molar fraction of the individual component m.	Chemical toxicity evaluation of DES, taking into account the mole fractions of individual components <sup>1</sup>
Number of heavy atoms	$HM = HM1 * x_1 + HM2 * x_2 + HM3 * x_3,$ where $HM$ is the total number of heavy atoms, $HM_m$ is the number of heavy atoms in individual component m, and $x_m$ is the molar fraction of the individual component m.	Number of heavy atoms in DES, taking into account the mole fractions of individual components <sup>1</sup>

 Table S2. Descriptors for density and viscosity prediction

Descriptor	Formula / A source	Meaning
Molar fractions	Experimental data	Molar fractions of individual components in DES
	Experimental data	Type of DES (LILLIII, IV, V, Ternary)
Temperature	Experimental data	The temperature at which the density or viscosity
remperature	Experimental data	measurement was carried out
VdWVolume	$VvdW = \sum all atom \ contributions - 5.92 * NB - 14.7 * RA - 3.8 * RNA,$	Van der Waals volumes of individual components
	where $NB$ is the number of bonds. $RA$ is the number of aromatic rings, and $RNA$ is the number of nonaromatic	calculated using Van der Waals radii obtained from
	rings.	RDKit <sup>1,10</sup>
	$V_{ball} = \frac{4}{3} * \pi * R^3$ , To calculate the contributions of atoms, we use the ball volume formula:	
	where <i>R</i> is Van der Waals radius of atom.	
	To calculate the number of bonds, we use the following formula:	
	NB = N - 1 + RA + RNA,	
	where <sup>N</sup> is the total number of atoms.	
NumHeteroatoms	$NHet = NHet1 * x_1 + NHet2 * x_2 + NHet3 * x_3,$	The number of heteroatoms in a molecule calculated
	where $^{NHet}$ is the total number of heteroatoms, $^{NHet}_m$ is the number of heteroatoms in individual component	for the entire molecule, taking into account the molar
	m, and $\frac{x_m}{m}$ is the molar fraction of the individual component m.	fractions of individual components <sup>1</sup>
RingCount	$Ring C = Ring C1 * x_1 + Ring C2 * x_2 + Ring C3 * x_3,$	The number of rings in a molecule calculated for the
	where $RingC$ is the total number of rings, $RingC_m$ is the number of rings in individual component m, and $x_m$ is	entire molecule, taking into account the molar fractions
	the molar fraction of the individual component m.	of individual components <sup>1</sup>
InertialShapeFactor	$pm_1 * pm_3$	Inertial shape factor (it is related to the moment of
	where $pm_n$ is the principal moment of inertia <i>n</i> of DES calculated by the following formula: $pm_n = PMI_1 + x_1 + PMI_2 + x_2 + PMI_3 + x_3$	inertia of the molecule and the characteristics of its shape) of DES taking into account the molar fractions of
	where $PMI_m$ is the principal moment of inertian of DES of individual component m and $x_m$ is the molar	individual components <sup>1,11</sup>
	fraction of the individual component m	
SnherocityIndex	$\frac{1}{3pherocity index - 3 + \frac{1}{m_1 + m_2 + m_2}},$	Spherocity Index (related to the shape of molecules
opheroentymaex		and aspects that can be associated with sphericity or
	where $pm_n$ is the principal moment of inertia <i>n</i> of DES calculated by the following formula: $pm_n = PMI_{-}1_n * x_1 + PMI_{-}2_n * x_2 + PMI_{-}3_n * x_3$	elongation) of DES taking into account the molar
	where $PMI_m_n$ is the principal moment of inertia n of DES of individual component m and $x_m$ is the molar	fractions of individual components <sup>1,11</sup>
	fraction of the individual component m.	
Mass fraction of metal	$\frac{MetalMassFraction}{x_1 * mw_1 + x_2 * mw_2 + x_3 * mw_3},$	Mass fraction of metal in DES, taking into account the
	х х	mole fractions of individual components
	where "m is the molar fraction of the individual component," m is the number of metals in the individual $MW(Mg)$ is the number of metals in the individual	
N	component, <i>in the molar mass of metal, and <math>n</math> is the molar mass of the individual component.</i> $N_{Elaw} = x_1 * n_1 + x_2 * n_2 + x_2 * n_2.$	
Number of elements	$\frac{1}{2} \frac{1}{2} \frac{1}$	ine number of elements in DES taking into account the
	where $m$ is the molar fraction of the individual component and $m$ is the number of elements in the individual	mole tractions of individual components <sup>1</sup>
	component m.	Density: Li, C, N, O, F, Na, Wig, Ai, P, S, Ci, K, Ca, Cr, Min,
		Fe, Co, Ni, Cu, Zn, Br

	Viscosity: Li, C, N, O, F, Mg, Al, P, S, Cl, K, Cr, Mn, Fe, Co,
	Cu, Zn, Br

Table S3. Metrics for evaluating the accuracy of machine learning models

Metric	Meaning	Formula	
Coefficient of	A metric that shows the proportion of the		
Determination (R <sup>2</sup> )	variance in the response variable of a	$\hat{\mathbf{y}}_{i+1}$	
	regression model that can be explained by	where $\mathcal{I}$ is the true value of the parameter, $\mathcal{I}$ is the predicted value of the parameter, and $\mathcal{I}$ is the average value of the	
	the predictor variables <sup>12</sup> .	variable.	
Root Mean Square Error	A metric that tells how far apart the	$\sqrt{\sum_{i=1}^{n}}$ <i>n</i>	
(RMSE)	predicted values are from the observed	$\hat{\mathbf{y}}_{i+1}$	
	values in a dataset, on average <sup>12</sup> .	where $\frac{1}{2}$ is the true value of the parameter and $\frac{1}{2}$ is the predicted value of the parameter.	
Average absolute relative	A metric used to measure the accuracy of a	$n \sum_{i \in \mathcal{Y}_i}  y_i $	
deviation (AARD)	model by calculating the average	$\hat{y}_{i}$	
	percentage difference between predicted	where $n$ is the total number of data points or observations, $\mathcal{F}_i$ is the true value of the parameter, and $\mathcal{F}_i$ is the predicted value of the parameter.	
	and actual values <sup>13</sup> .		

Table S4. Optimal hyperparameters for the best models

Property	ML model	Hyperparameters
Melting point	Cat Boosting Regression	iterations: 123 learning_rate: 0.061 depth: 6
Density	Cat Boosting Regression	iterations: 600 learning_rate: 0.050 depth: 4
Viscosity	Cat Boosting Regression	iterations: 600 learning_rate: 0.050 depth: 6

MELTING TEMPERATURE













Figure. S3. Performance of predictive machine learning models

#### MELTING TEMPERATURE



Figure. S4. Metrics for different machine learning models for melting temperature



Figure. S5. Metrics for different machine learning models for density



Figure. S6. Metrics for different machine learning models for viscosity



Figure. S7. RMSE for different types of deep eutectic solvents

# **SI References**

- 1 RDKit: Open-source cheminformatics; http://www.rdkit.org .
- 2 F. Pedregosa FABIANPEDREGOSA, V. Michel, O. Grisel OLIVIERGRISEL, M. Blondel, P. Prettenhofer, R. Weiss, J. Vanderplas, D. Cournapeau, F. Pedregosa, G. Varoquaux, A. Gramfort, B. Thirion, O. Grisel, V. Dubourg, A. Passos, M. Brucher, M. Perrot andÉdouardand, andÉdouard Duchesnay and Fré. Duchesnay EDOUARDDUCHESNAY, *Scikit-learn: Machine Learning in Python Gaël Varoquaux Bertrand Thirion Vincent Dubourg Alexandre Passos PEDREGOSA, VAROQUAUX, GRAMFORT ET AL. Matthieu Perrot*, 2011, vol. 12.
- 3 T. Chen and C. Guestrin, in *Proceedings of the ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, Association for Computing Machinery, 2016, vol. 13-

17-August-2016, pp. 785–794.

- 4 L. Prokhorenkova, G. Gusev, A. Vorobev, A. V. Dorogush and A. Gulin, *CatBoost: unbiased boosting with categorical features*.
- 5 T. Head, M. Kumar, H. Nahrstaedt and G. Louppe, *scikit-optimize/scikit-optimize (v0.9.0)*. *Zenodo*.
- 6 S. M. Lundberg and S.-I. Lee, *A Unified Approach to Interpreting Model Predictions*.
- 7 D. Merkel, *Linux Journal*.
- 8 W. Reese, *Linux Journal*.
- 9 *PostgreSQL: The World's Most Advanced Open Source Relational Database*, 2019.
- 10 Y. H. Zhao, M. H. Abraham and A. M. Zissimos, *Journal of Organic Chemistry*, 2003, **68**, 7368–7373.
- 11 Peter. Comba and T. W. Hambley, *Molecular modeling of inorganic compounds*, Wiley-VCH, 2001.
- 12 D. Chicco Corresp, M. J. Warrens, G. Jurman and D. Chicco, *The coefficient of determination R-squared is more informative than SMAPE, MAE, MAPE, MSE, and RMSE in regression analysis evaluation analysis evaluation* 4, 2021.
- 13 W. Su, Y. Hwang, Y. shao, S. Deng, L. Zhao, X. Nie and Y. Zhang, *Energy*, 2019, **166**, 414–425.