

Are Log P Calculators Accurate? Benchmarking on 96 000 Compounds

Igor V. Tetko,¹ Gennadiy I. Poda,² Claude Ostermann,³ Raimund Mannhold^{4,*}

1-Helmholtz Zentrum München, Neuherberg, Germany; 2-Pfizer Global R & D, Chesterfield, USA
3-Nycomed GmbH, Konstanz, Germany; 4-Heinrich-Heine-Universität, Düsseldorf, Germany

You can download this poster, preprint of the article and the ALOGPS program at <http://www.vcllab.org>

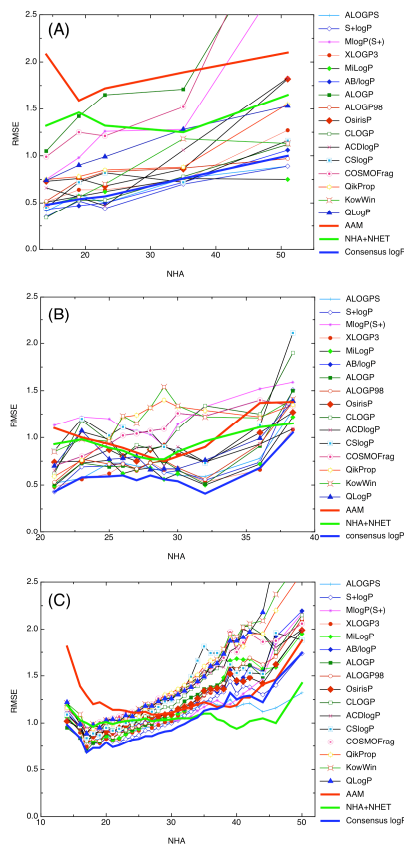
Introduction

Lipophilicity, quantified as log P, is an important parameter daily monitored by medicinal chemists in drug discovery.

The application of log P prediction software tools deserves a comprehensive validity check on large, chemically diverse databases. Several evaluations of calculated vs. experimental log P values appeared in the literature. Most of them suffer from small size, limited chemical diversity, and/or inclusion of only a few log P calculation methods.

Here we review the state-of-the-art in development of log P prediction approaches including substructure- and property-based methods. Then, we compare the predictive power of representative methods for one public (N = 266) and two in-house datasets from Nycomed (N = 882) and Pfizer (N = 95 809). 30 methods were tested for the public and 18 for the industrial datasets.¹

*Mannhold R, Poda G, Ostermann C, Tetko I (2008) Calculations of Molecular Lipophilicity: Natural-like Set and Comparison of Log P Methods on More Than 96,000 Compounds. *J Pharm Sci*, in press.



Method performance as a function of the NHA for the public (A), Nycomed (B) and Pfizer (C) dataset. Each point on the graph shows the RMSE of the analyzed method for molecules with an NHA indicated on the x-axis (minimally 50 molecules were used per point for the public and the Nycomed and 500 molecules for the Pfizer dataset). The red bold line corresponds to the AAM. The blue bold line represents Consensus log P. The green bold line corresponds to the simple two-descriptor model (NC + NHET). Models with errors larger than those calculated with AAM fail to provide predictive model for molecules with a given NHA.

Log P programs used in benchmarking (Programs used for in-house datasets are given in bold)

Name	Provider	URL/E-mail
ABLogP v. 2.0	Pharma Algorithms, Lithuania/Canada	http://www.ablogp.com
ASSOLV, LSER	Pharma Algorithms, Lithuania/Canada	http://www.ablogp.com
ACDlogP v. 11	Advanced Chemistry Development, USA	http://www.acdchem.com
ALOGP (DragonX 1.4)	Talete Srl, Milano, Italy	http://www.talete.it
ALOGPS	Acetyls Software Inc., USA	http://www.alogps.com
ALOGP v. 2.1	Virtual Computational Chemistry Laboratory, Germany	http://www.vcllab.org
CLUP	University of Geneva, Switzerland	clup@unige.ch
ALOGP v. 4.3 (v. 5.0)	BioByte Inc., USA	http://www.biobyte.com
COSMOfrag v. 2.3	COSMOlogic GmbH & Co. KG, Germany	http://www.cosmologic.de
CSlogP	Chemical Ltd, USA	http://www.chemical.com
GBLogP	Max Tootov, USA	gg@maxtootov.com
HINT	EdoSoft, LLC, USA	http://www.edosoft.com
KowWin v. 1.67	Synapse Inc., USA	http://www.synapse.com
LSER UFL	Heinhold Center for Environ. Research UFZ, Germany	http://www.ufz.de
MLOGP v. 2.2	Molinspiration Cheminformatics, Slovak Republic	http://www.molinspiration.com
MLOGP (DragonX 1.4)	Talete Srl, Milano, Italy	http://www.talete.it
MLOGP(S+), ADMET 2.3	Simulations Plus, Inc., USA	http://www.simulations-plus.com
MolLog	MolSoft LLC, USA	http://www.molsoft.com
NC-NHET	Virtual Computational Chemistry Laboratory, Germany	http://www.vcllab.org
OsisP	Actelion, Switzerland	http://www.actelion.com
QikProp v. 3.0	Schrodinger LLC, USA	http://www.schrodinger.com
QLOGP	University of Miami, USA	http://www.miami.edu
Quantlog	Quantum Pharmaceuticals, Russia	http://www.quantum-pharm.com
S-logP, ADMET 2.3	Simulations Plus, Inc., USA	http://www.simulations-plus.com
SLIPPER-2002	Institute of Physiol. Active Compounds, Russia	http://icad.msk.su
SPARC	Upstream Solutions, Switzerland	http://www.upstream.ch
TILOGP	Institute of Physiol. Active Compounds, Russia	http://www.ti-chem.edu
VEGA	University of Milan, Italy	http://www.vega.unimi.it
VLQOP	TOPKAT, Acetyls Software Inc., USA	http://www.acdchem.com
XLOGP3	Test of Physical Chemistry, Peking University, China	http://www.xlogp3.com
XLOGP3	Institute of Organic Chemistry, Shanghai, China	http://www.xlogp3.com

Conclusions

Most methods produced reliable results for the public dataset, but for the in-house datasets only a few were superior to AAM. Among the best methods, ALOGPS, S-logP, XLOGP3, OsisP, ALOP, and ALOP98 were consistent for both in-house datasets. Surprisingly, these methods include the classical approaches ALOP and MLOGP. Best performances for the in-house datasets were calculated with Consensus log P.

Only Consensus log P calculated an RMSE=1.00 for the Pfizer set. Despite low average accuracy of prediction, the confidence of prediction and the analysis of molecular size seems to distinguish reliable from non-reliable predictions. E.g., molecules with NHA = 18-20 or m molecules with a sm all StdDev had an RMSE=0.75 for Consensus log P. All models poorly predicted molecules with NHA > 30-35. As a base-line estimation of log P with an RMSE = 1.00.

When prediction accuracy is low, measuring log P is highly recommended. Such new values allow to improve prediction power via a user-training option without training as provided by e.g. ALOP3, SLIPPER, KowWin, ABLogP, ACDlogP, XLOGP3, and good local models for improved predictions within homologous series can be derived.²

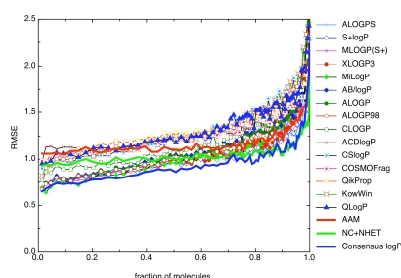
In summary, despite log P calculation is commonly viewed as simple, our analysis showed low prediction accuracy for most of the existing calculation methods. Inaccurate log P prediction can mislead selection of chemical series to follow up from virtual screening, HTS triage and new analog design since it may lead to discarding potentially promising structures due to "poor" physico-chemical properties. Studying the confidence of log P prediction allows to distinguish reliable from non-reliable predictions and, thus, helps to avoid this problem.

2. Tetko, I.V.; Bounoue, P.; Meyers, H.W.; Rehner, D.C.; Poda, G.I. Can we estimate the accuracy of ADME-Tox predictions? *Drug Discov. Today* 2006, 11, 115-161, 700-707.

Reliable versus non-reliable predictions

Can we a priori distinguish reliable from non-reliable log P predictions? We estimated the prediction accuracy of the tested methods by their standard deviations (StdDev) compared to Consensus logP; accuracy gradually decreased with StdDev. For example, 10% of molecules with most confident predictions (i.e., StdDev < 0.3) had log P values predicted with an average RMSE of 0.73 and 0.75 for Consensus logP and ALOP3, respectively. While, there was a shift in the distribution of these 10% of molecules towards smaller NHA, this set also included molecules with different molecular sizes. Thus, analysing the confidence of predictions allowed accurate predictions for molecules outside the range of optimal size, i.e. 16-NHA<21.

Comparison of methods using the total RMSE may lead to wrong conclusions about method performance. For example, ClogP produced an RMSE of 1.29 for the Pfizer set, i.e. had a larger RMSE compared to those of the NC-NHET and AAM models. However, for the same 10% threshold of the most confident predictions ClogP produced an RMSE of 0.87, indicating a better performance than the AAM (RMSE=1.07) and NC-NHET (RMSE=0.95) models. Thus, despite the fact that the NC-NHET equation had a lower RMSE compared to ClogP for the entire set, the latter model provided higher accuracy for the molecules with most confident predictions.



The RMSE of methods for Pfizer data as function of the fraction of molecules sorted along increasing StdDev values. Each point (at least 500 molecules) averages errors of m methods with the same (or very similar) StdDev values. The first 10% of molecules have a StdDev < 0.3.

Benchmarking of methods using the public dataset

Datasets and definitions

30 methods compared in separate analysis of StarList and Non-Star List.

Entire dataset: 266 molecules of a dataset from Alex Audev

StarList subset: 223 molecules also present in BioByte list

Non-Star List subset: 43 molecules outside BioByte list (NCEs)

Arithmetic Average Model (AAM): mean log P, used as model that predicts the same value for all dataset molecules; mean log P of entire dataset = 2.32, R²=0 between predicted and experimental values

Rank II: models with root mean squared errors (RMSE) close to or larger than that of the AAM, i.e. models are non-predictive

Rank I: methods with statistical results identical or close to ABLlogP and ALOP3 (smallest RMSE for Star and Non-Star set)

Rank I: remaining models

Consensus logP: average of predicted log P from all rank I and II methods

Benchmarking Results

Models are ranked according to their predictive performance for both datasets.

Prediction results for the Star set were statistically almost identical for the top-ranked methods ABLogP, S-logP, and ACDlogP.

For the Non-Star set, 17 methods non-significantly differed in performance from top-ranking ALOP3. Its small size (N=43) excludes statistical differentiation between methods. For all tested calculation methods, the prediction performance was on average by 0.35 log units lower for the Non-Star set versus the Star set.

6 methods for the Star set and 8 methods for the Non-Star set showed results statistically similar to the AAM model. Thus, these methods failed to provide accurate ranking of these datasets.

Counting acceptable (A<0.5), disputable (A=0.5-1.0), and unacceptable (A>1.0) predictions allows accurate ranking of predictive power and quick checking of the general applicability of the methods.

Performance of algorithms for the public dataset

Method	Star set (N = 223)					Non-Star set (N = 43)				
	RMSE	rank	% within error range	RMSE	rank	RMSE	rank	% within error range	RMSE	rank
ABLogP	0.41	1	84	12	4	1.00	1	42	23	35
S-logP	0.45	1	79	22	3	0.87	1	40	35	26
ACDlogP	0.50	2	75	17	1	1.00	1	40	35	26
Consensus log P	0.52	1	74	18	9	0.70	1	17	28	26
QLOGP	0.52	1	74	28	5	0.91	1	47	28	26
VLQOP OPS	0.52	1	64	21	7	1.07	1	33	28	26
ALOP3	0.53	1	71	20	8	0.82	1	42	26	26
MLOGP	0.57	1	69	22	9	0.86	1	49	30	21
XLOGP3	0.62	1	60	30	10	0.89	1	47	23	30
KowWin	0.64	1	68	21	11	1.05	1	40	30	30
CSlogP	0.65	1	66	22	12	0.93	1	58	19	23
ALOP3	0.69	1	60	26	13	0.82	1	42	30	33
MLOGP	0.69	1	61	25	14	0.93	1	40	35	26
ALOP3	0.70	1	61	26	13	1.00	1	30	37	33
VEGA	0.71	1	59	26	15	0.94	1	42	35	33
VLQOP	0.72	1	65	22	14	1.13	1	40	28	33
TILOGP	0.74	1	67	16	16	1.26	1	30	37	30
ASSOLV	0.75	1	53	30	17	1.02	1	49	28	23
QikProp	0.77	1	53	30	17	1.24	1	40	26	35
Quantlog	0.80	1	47	30	22	1.17	1	35	40	32
SLIPPER-2002	0.80	1	62	22	15	1.16	1	35	23	42
COSMOfrag	0.84	1	48	26	18	1.03	1	26	35	33
XLOGP2	0.87	1	57	22	20	1.16	1	35	23	42
QLOGP	0.96	1	49	26	20	1.42	1	21	26	53
CLUP	1.04	1	47	27	26	1.24	1	28	30	42
VEGA	1.05	1	41	25	30	1.54	1	33	9	49
TILOGP	1.07	1	44	26	30	1.26	1	35	16	56
MLOGP (Sim)	1.26	1	38	30	33	1.58	1	26	28	47
SPARC	1.36	1	45	22	32	1.70	1	28	25	49
MLOGP(Dragon)	1.52	1	39	26	35	2.45	1	19	30	53
HINT	1.60	1	36	23	41	1.75	1	19	32	47
AAM	1.62	1	22	34	33	3.16	1	23	38	43
VLQOP-MOPS	1.76	1	31	27	37	1.39	1	30	31	43
NC-NHET	1.80	1	34	22	44	1.72	1	30	6	65
GBLogP	1.98	1	32	26	42	1.75	1	19	16	65

Number of molecules in the dataset. RMSE values for the 10 methods with the lowest RMSE values. RMSE values for the 10 methods with the highest RMSE values. RMSE values for the 10 methods with the highest RMSE values. RMSE values for the 10 methods with the highest RMSE values.

Benchmarking of methods using in-house datasets

Datasets

18 methods compared in separate analysis of Pfizer and Nycomed set. Due to the large size of the in-house datasets we included those methods that could be run in batch mode.

Pfizer dataset

95 809 experimental log P values generated at a number of legacy sites. We retained all experimental log P data points and added log D values for species predominantly neutral at pH 7.4. Multiple log PID data points for compounds with identical structures were averaged. Salts were stripped off. The original tautomers and stereoisomers were preserved (as deposited in the proprietary database).

Nycomed dataset

The dataset of 882 compounds was prepared in a similar fashion.

Benchmarking Results

Mean log P values of 2.92 and 3.15 were calculated as AAMs for compounds from the Pfizer and Nycomed datasets, respectively.

We observed dramatic differences in method performance between public and in-house datasets. Based on RMSE, out of 18 methods only 8 (Nycomed) and 9 (Pfizer) produced results significantly better than AAM. Following methods conformed for both sets: ALOP3 and S-logP, XLOGP3, ALOP3, OsisP, MLOGP and ALOP98.

MLOGP(S+) and ABLogP produced good results for the Pfizer dataset, but failed to produce models with RMSE significantly lower than that of AAM for the Nycomed dataset. QLOGP showed the opposite behavior.

In summary, atom-based approaches coupled with non-linear neural network learning techniques exhibited higher prediction accuracy in this validity check than fragment-based approaches. The highest accuracy for both in-house datasets was calculated using Consensus logP.

Performance of algorithms for in-house datasets

Method	Pfizer set (N = 95 809)					Nycomed set (N = 882)				
	RMSE	rank	% within error range	RMSE	rank	RMSE	rank	% within error range	RMSE	rank
Consensus log P	1.02	1	43	29	0.74	0.68	1	41	32	18
ALOP3	1.02	1	44	30	29	1.01	1	51	34	18
S-logP	1.02	1	44	29	27	1.00	1	58	27	15
NC-NHET	1.04	1	36	35	32	1.04	1	58	42	32
MLOGP(S+)	1.05	1	40	29	31	1.05	1	11	32	41
XLOGP3	1.07	1	43	28	29	1.06	1	55	34	12
MLOGP	1.10	27	11	41	30	1.09	1	67	1	26
ABLogP	1.12	24	11	48	29	1.11	1	68	45	28
ALOP3	1.12	1	39	32	1.12	0.72	1	52	33	15
ALOP98	1.12	1	40	28	32	1.10	1	52	31	17
OsisP	1.13	6	39	28	33	1.12	1	65	43	34
AAM	1.14	1	31	22	38	1.16	1	54	42	31
QLOGP	1.23	1	37	28	35	1.21	1	101	46	28
ACDlogP	1.28	1	35	27	38	1.28	1	87	46	23
CSlogP	1.29	20	37	27	36	1.28	1	106	38	33
COSMOfrag	1.30	1088	32	27	40	1.30	1	106	39	40
QikProp	1.32	103	31	26	43	1.32	1	117	27	49
KowWin	1.32	16	33	25	41	1.31	1	120	39	27
QLOGP	1.33	24	34	27	39	1.32	1	100	50	37
XLOGP2	1.80	1	15	67	180	0.94	39	31	29	1
MLOGP(Dragon)	2.03	1	34	24	42	2.03	0.90	31	45	30

*Number of molecules with no tautomers failures due to errors or omissions of programs. All methods predicted all molecules for the Nycomed dataset. RMSE calculated after excluding of 189 tautomers compounds from the Pfizer dataset. Most molecules failed by COSMOfrag are tautomers.