

# Journal Name

## ARTICLE TYPE

Cite this: DOI:00.0000/xxxxxxxxxx

# Machine learning-based prediction of fish acute mortality: Implementation, interpretation, and regulatory relevance

Lilian Gasser,<sup>\*a‡</sup> Christoph Schür,<sup>b‡</sup> Fernando Perez-Cruz,<sup>a,c</sup> Kristin Schirmer<sup>b,d,e</sup> and Marco Baity-Jesi<sup>b</sup>

Received Date

Accepted Date

DOI:00.0000/xxxxxxxxxx

- 1 Data
- 2 Methods
- 2.1 Models
- 3 Results and Discussion
  - 3.1 Data quality & variability
  - 3.2 Modeling results
    - 3.2.1 Validation results
    - 3.2.2 Performance on test set
  - 3.3 Including phylogenetic distances
  - 3.4 Explainability
    - 3.4.1 Residuals
    - 3.4.2 Feature importance
    - 3.4.3 Species sensitivity
- References

<sup>a</sup> Swiss Data Science Center (SDSC), Andreasstrasse 5, Zürich, Switzerland. E-mail: [lilian.gasser@sdsc.ethz.ch](mailto:lilian.gasser@sdsc.ethz.ch)

<sup>b</sup> Eawag, Swiss Federal Institute of Aquatic Science and Technology, Dübendorf, Switzerland.

<sup>c</sup> ETH Zürich: Department of Computer Science, Zürich, Switzerland.

<sup>d</sup> ETH Zürich: Department of Environmental Systems Science, Zürich, Switzerland.

<sup>e</sup> EPF Lausanne, School of Architecture, Civil and Environmental Engineering, Lausanne, Switzerland.

‡ These authors contributed equally to this work.

Table 1 Response and feature variables. The column name corresponds to the naming in the input datasets.

Column name	Description	Type of variable
Response variables		
<i>result_conc1_mean</i>	Lethal mass concentration value (in mg/L)	continuous
<i>result_conc1_mean_mol</i>	Lethal molar concentration value (in mol/L)	continuous
<i>result_conc1_mean_log</i>	Lethal mass concentration value after a log10 transformation	continuous
<i>result_conc1_mean_mol_log</i>	Lethal molar concentration value after a log10 transformation	continuous
Experimental features		
<i>result_obs_duration_mean</i>	Observation duration (in hours)	ordinal
<i>test_media_type</i>	Media type	nominal
<i>test_exposure_type</i>	Exposure type	nominal
<i>result_conc1_type</i>	Exposure concentration type	nominal
Chemical properties		
<i>chem_mw</i>	Molecular weight (in g/mol)	continuous
<i>chem_ws</i>	Water solubility (in mg/L)	continuous
<i>chem_mp</i>	Melting point (in °C)	continuous
<i>chem_rdkit_clogp</i>	Octanol-water partition coefficient	continuous
Molecular representations		
<i>chem_MACCS_fp</i>	Collapsed MACCS fingerprint	binary bits
<i>chem_pcp_fp</i>	Collapsed PubChem fingerprint	binary bits
<i>chem_Morgan_fp</i>	Collapsed Morgan fingerprint	binary bits
<i>chem_ToxPrint_fp</i>	Collapsed ToxPrint fingerprint	binary bits
<i>chem_mol2vec[000-299]</i>	300-dimensional mol2vec embedding	continuous
<i>chem_mordred_x</i>	Mordred features	continuous or ordinal
Taxonomic properties		
<i>tax_eco_climate</i>	Ecology, climate zone	nominal
<i>tax_eco_ecozone</i>	Ecology, ecozone	nominal
<i>tax_eco_food</i>	Ecology, food class	nominal
<i>tax_eco_migrate2</i>	Ecology, migratory behavior (binary encoding)	binary
<i>tax_lh_amd</i>	Life history, life span (in d)	continuous
<i>tax_lh_licm</i>	Life history, ultimate body length (in cm)	continuous
<i>tax_ps_ampv</i>	Pseudo-data, energy conductance (in cm/d)	continuous
<i>tax_ps_ampkap</i>	Pseudo-data, allocation fraction to soma	continuous
<i>tax_ps_amppm</i>	Pseudo-data, volume-specific somatic maintenance cost (in J/d · cm³)	continuous
Phylogenetic distance stored in a separate file	Phylogenetic distance	continuous

Table 2 Hyperparameters for the four models. \*For the Gaussian process regression hyperparameter, no default values is given in the implementation.

Model	Hyperparameter	Default	Gridsearch values
LASSO	alpha	1.0	[np.round(i, 5) for i in np.logspace(-5, 0, num=26)]
RF	n_estimators	100	50, 100, 150, 300
	max_depth	None	50, 100, 200
	min_samples_split	2	2, 5, 10
	max_samples	None	0.25, 0.5, 1.0
	max_features	1.0	'sqrt', 1
XGBoost	n_estimators	100	50, 100
	eta	0.3	0.1, 0.2, 0.3
	gamma	0	0, 1, 10
	max_depth		3, 6, 9, 12
	min_child_weight	1	1, 3, 5
	subsample	1	0.5, 1.
GP	n_inducing	*	100, 250, 500, 1000

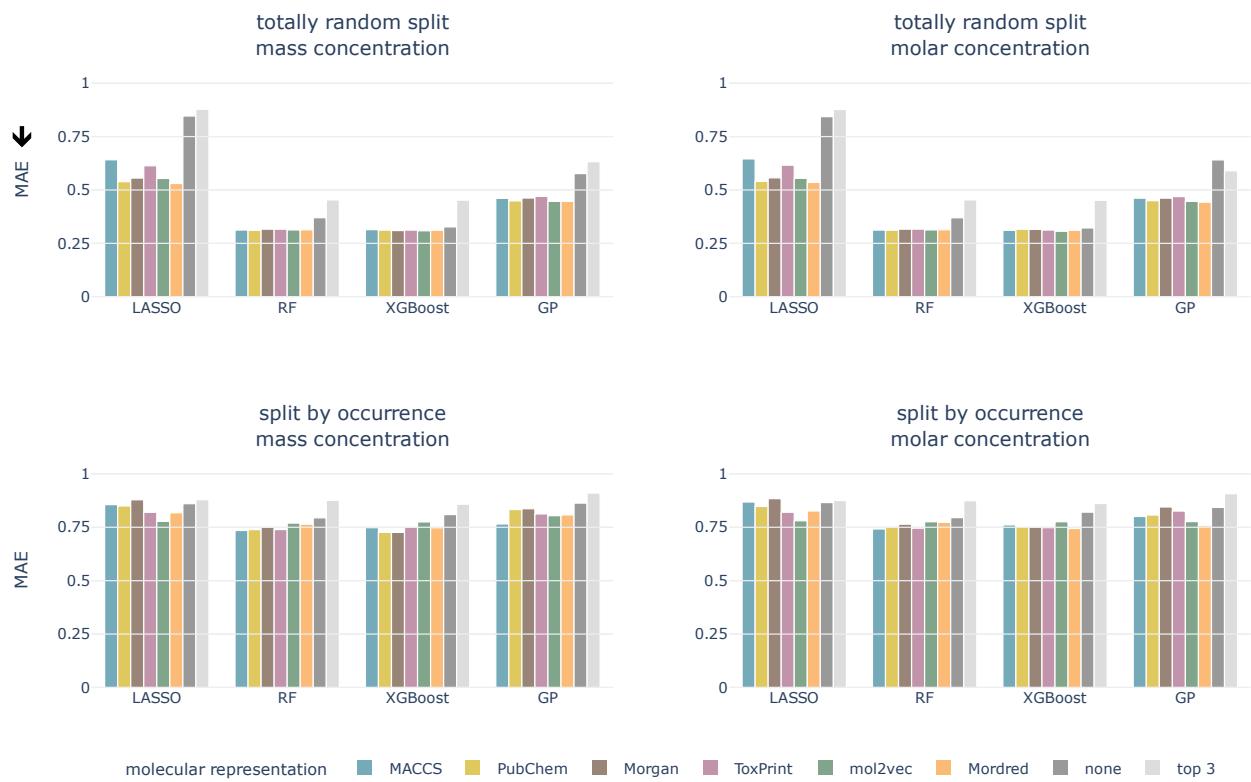


Fig. 1 Cross-validated MAE for both data splittings, concentration types, all models, and molecular representations. Arrow indicates the lower the better.

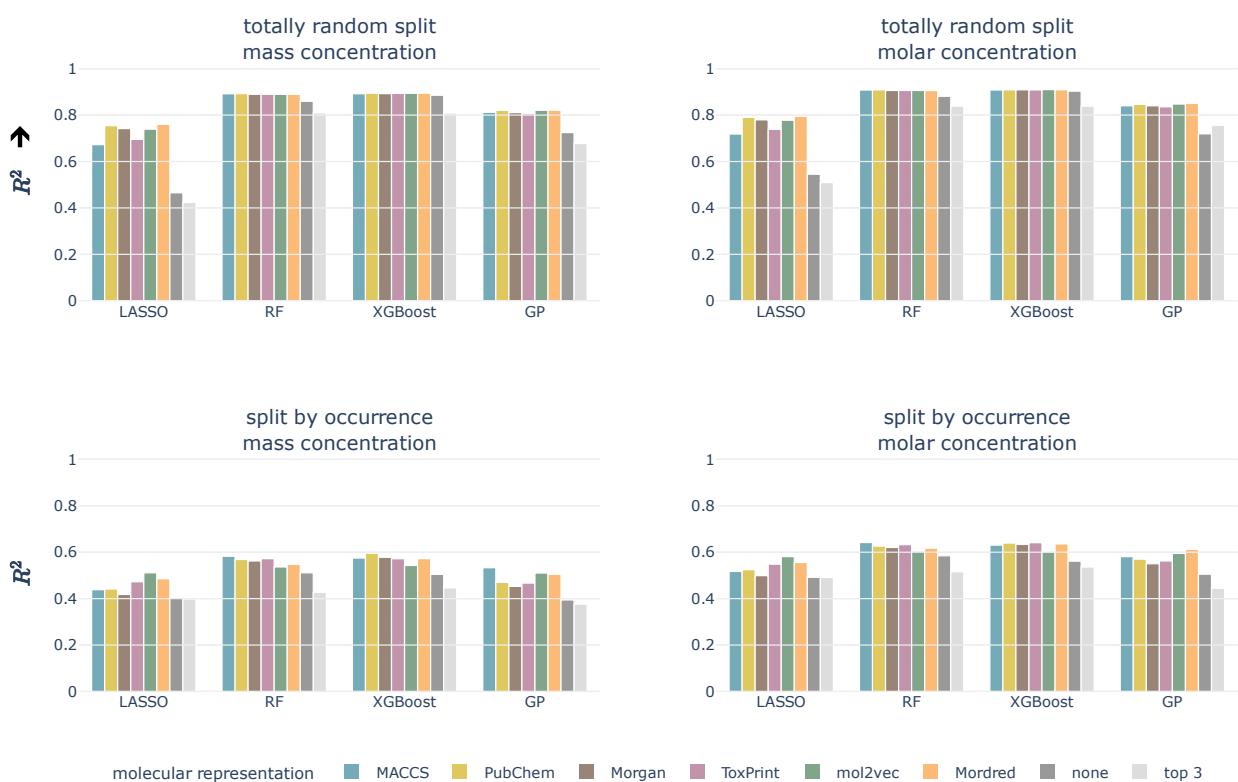


Fig. 2 Cross-validated  $R^2$  for both data splittings, concentration types, all models, and molecular representations. Arrow indicates the higher the better.

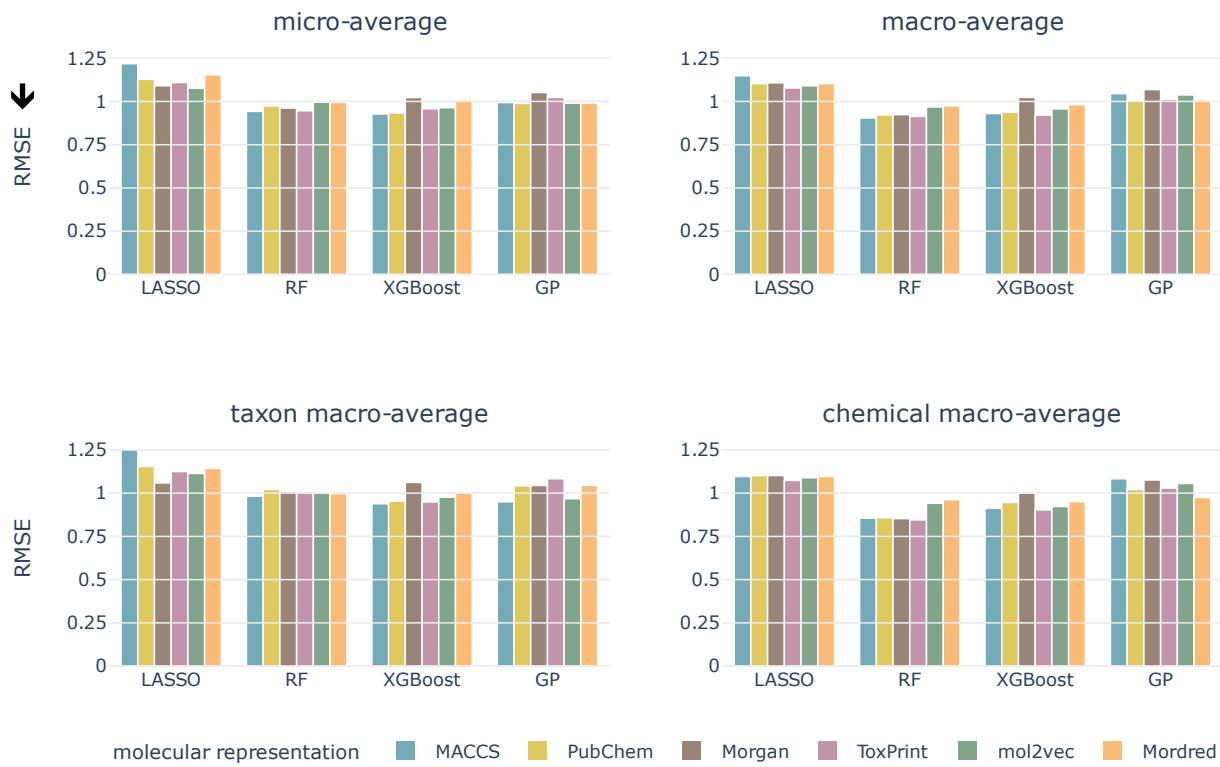


Fig. 3 Test micro- and macro-averaged RMSE for molar LC50, split by occurrence of chemical compounds, all models, and molecular representations. Arrow indicates the lower the better.



Fig. 4 Including phylogenetic distances: cross-validated RMSE for Gaussian process models, both data splittings, concentration types, and all molecular representations. Arrow indicates the lower the better.

Table 3 Best hyperparameters for LASSO.

concentration	data split	mol. repr.	alpha
molar	totally random	MACCS	0.00001
molar	totally random	PubChem	0.00004
molar	totally random	Morgan	0.00001
molar	totally random	ToxPrint	0.00002
molar	totally random	mol2vec	0.00025
molar	totally random	Mordred	0.00016
molar	occurrence	MACCS	0.01585
molar	occurrence	PubChem	0.02512
molar	occurrence	Morgan	0.01000
molar	occurrence	ToxPrint	0.00631
molar	occurrence	mol2vec	0.01585
molar	occurrence	Mordred	0.10000
mass	totally random	MACCS	0.00001
mass	totally random	PubChem	0.00004
mass	totally random	Morgan	0.00010
mass	totally random	ToxPrint	0.00002
mass	totally random	mol2vec	0.00010
mass	totally random	Mordred	0.00010
mass	occurrence	MACCS	0.02512
mass	occurrence	PubChem	0.02512
mass	occurrence	Morgan	0.01000
mass	occurrence	ToxPrint	0.00631
mass	occurrence	mol2vec	0.01000
mass	occurrence	Mordred	0.10000

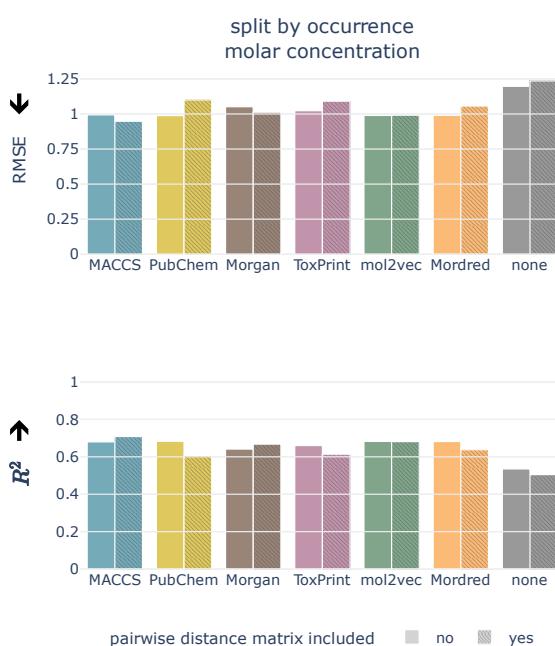


Fig. 5 Including phylogenetic distances: test RMSE and  $R^2$  for Gaussian process models, molar LC50, split by occurrence of chemical compounds, and all molecular representations. Arrows indicate the lower/higher the better.

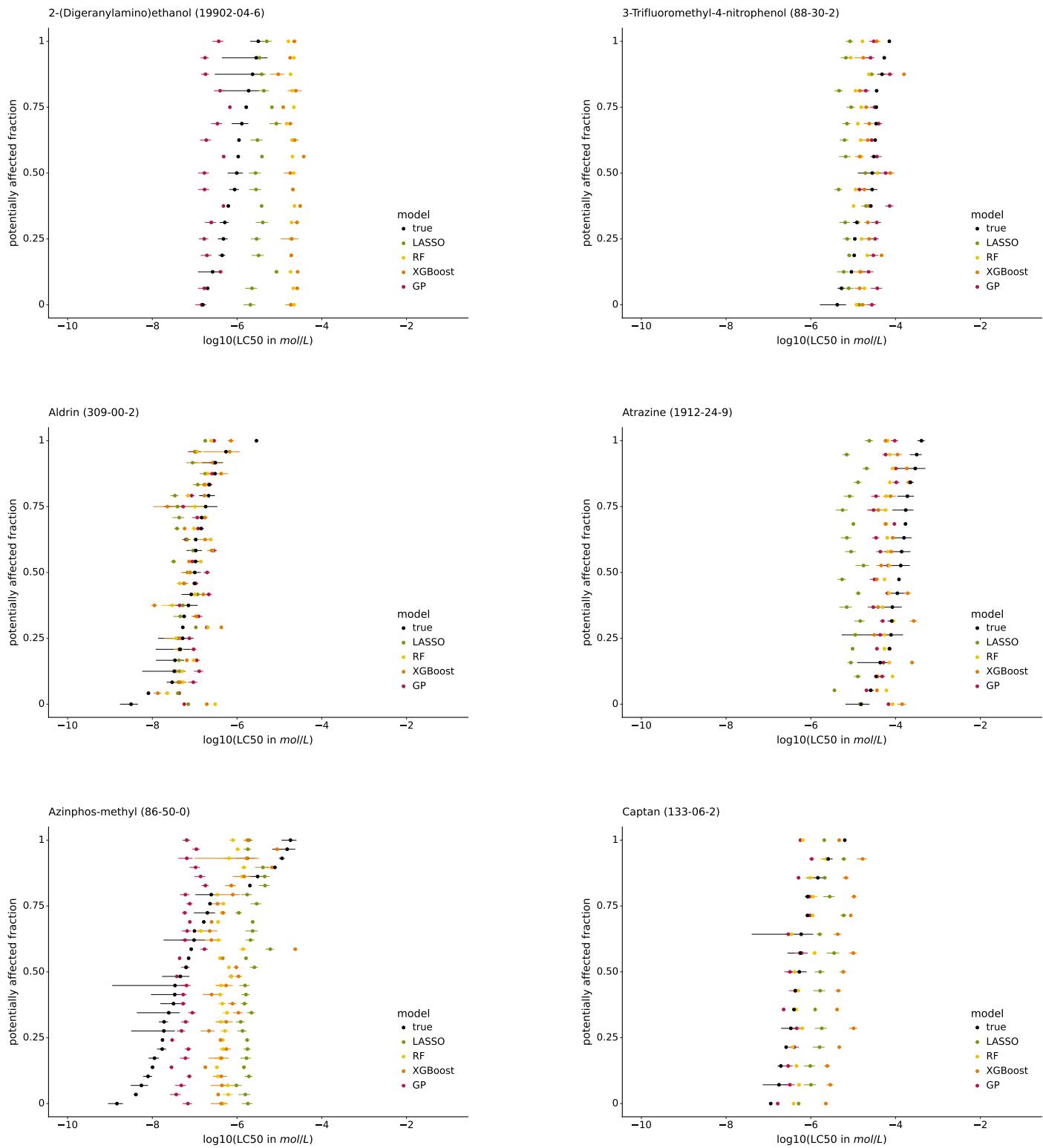


Fig. 6 Species sensitivity distributions (SSDs) of 2-(Digeranylamo)ethanol, 3-Trifluoromethyl-4-nitrophenol, Aldrin, Atrazine, Azinphos-methyl, and Captan.

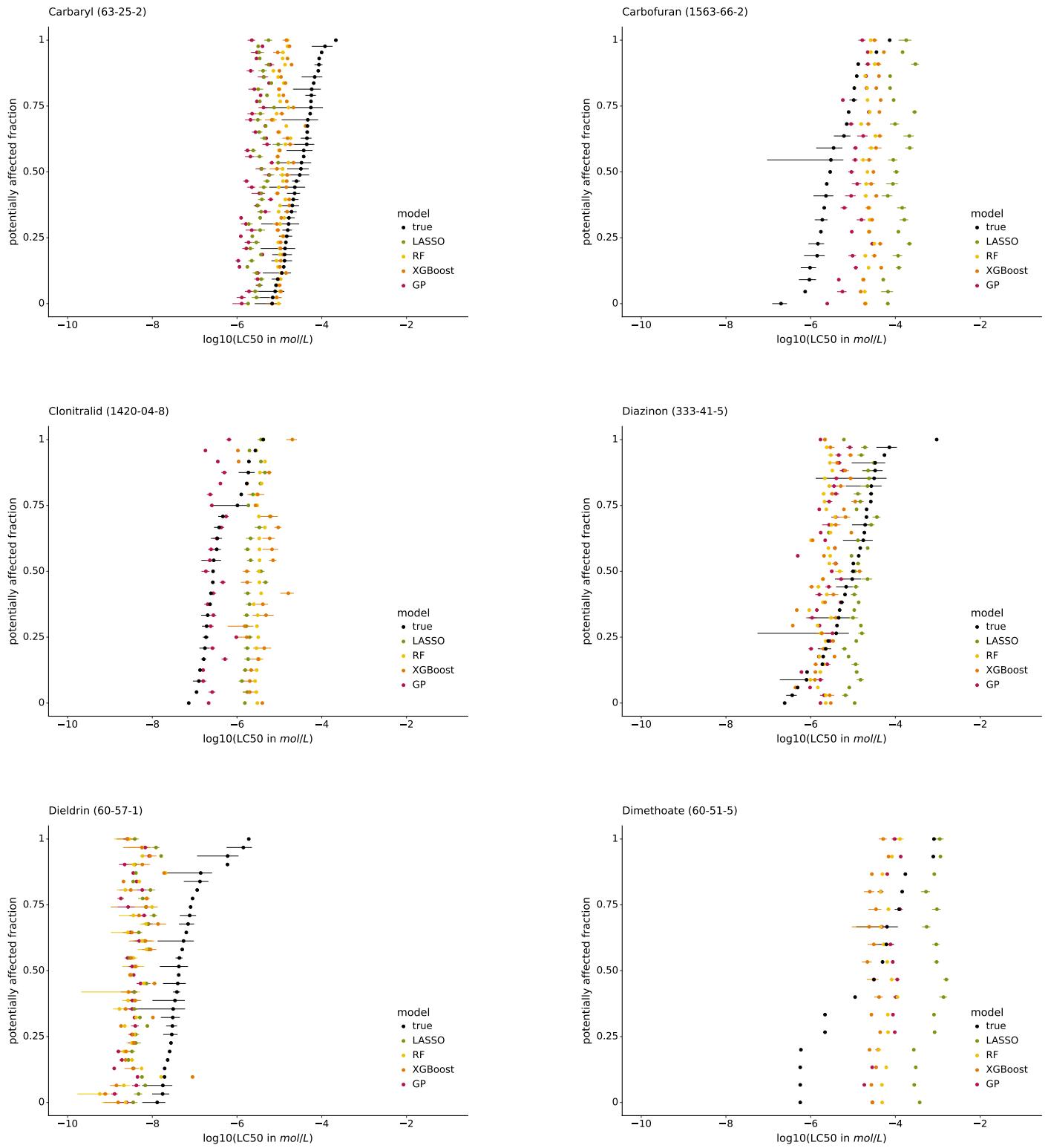


Fig. 7 Species sensitivity distributions (SSDs) of Carbaryl, Carbofuran, Clonitralid, Diazinon, Dieldrin, and Dimethoate.

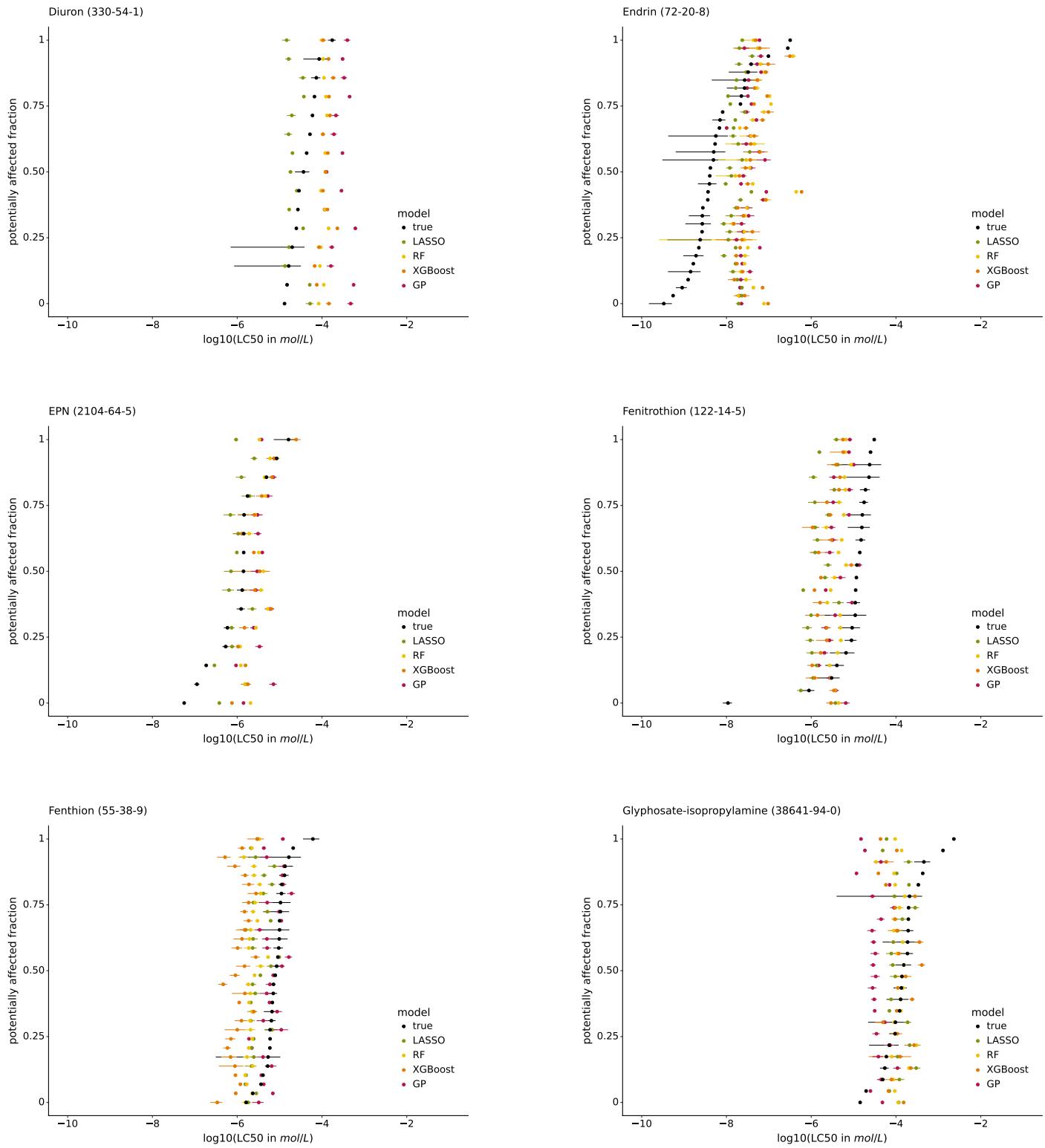


Fig. 8 Species sensitivity distributions (SSDs) of Diuron, Endrin, EPN (Ethyl p-nitrophenyl), Fenitrothion, Fenthion, and Glyphosate-isopropylamine.

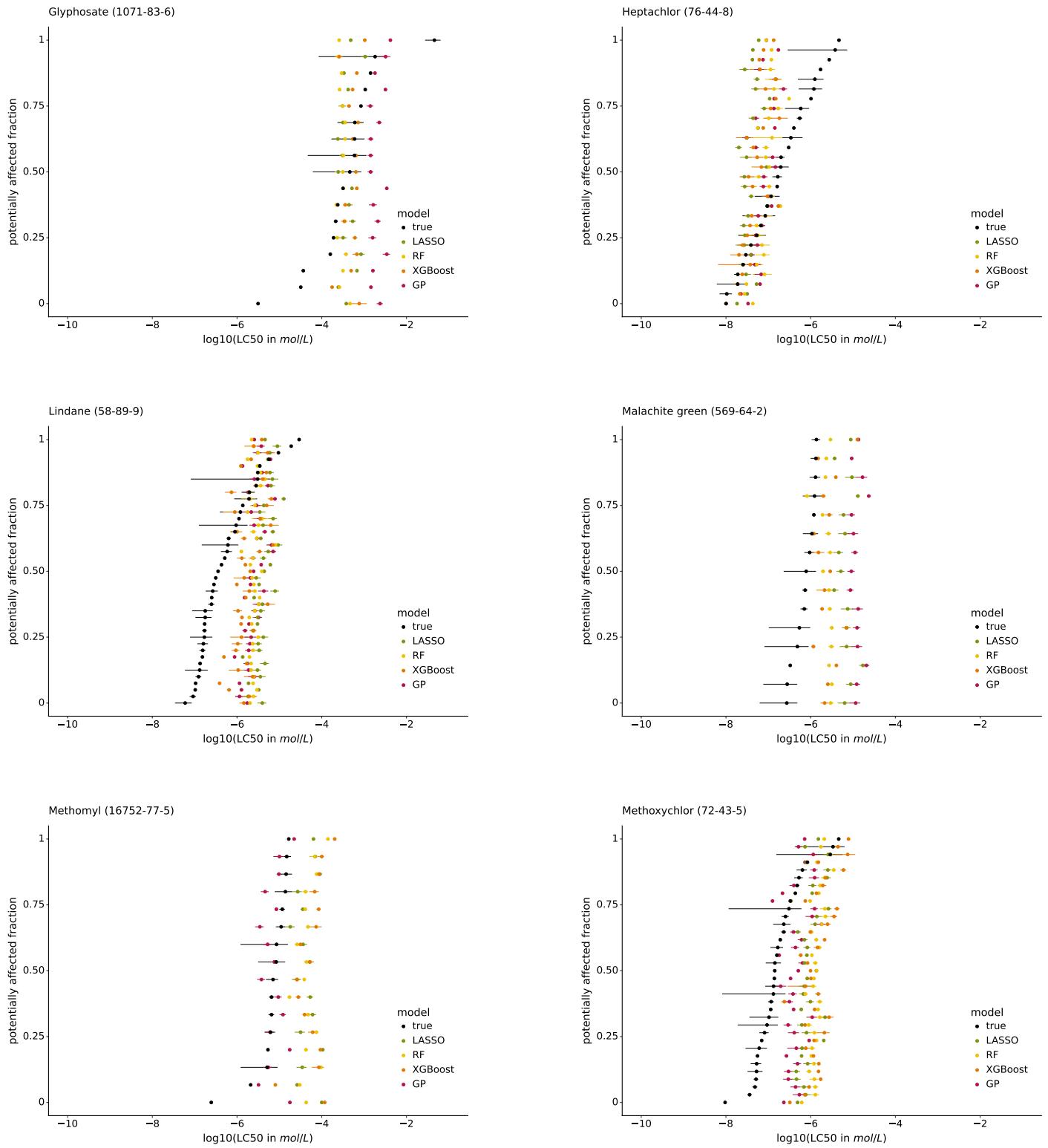


Fig. 9 Species sensitivity distributions (SSDs) of Glyphosate, Heptachlor, Lindane, Malachite green, Methomyl, and Methoxychlor.

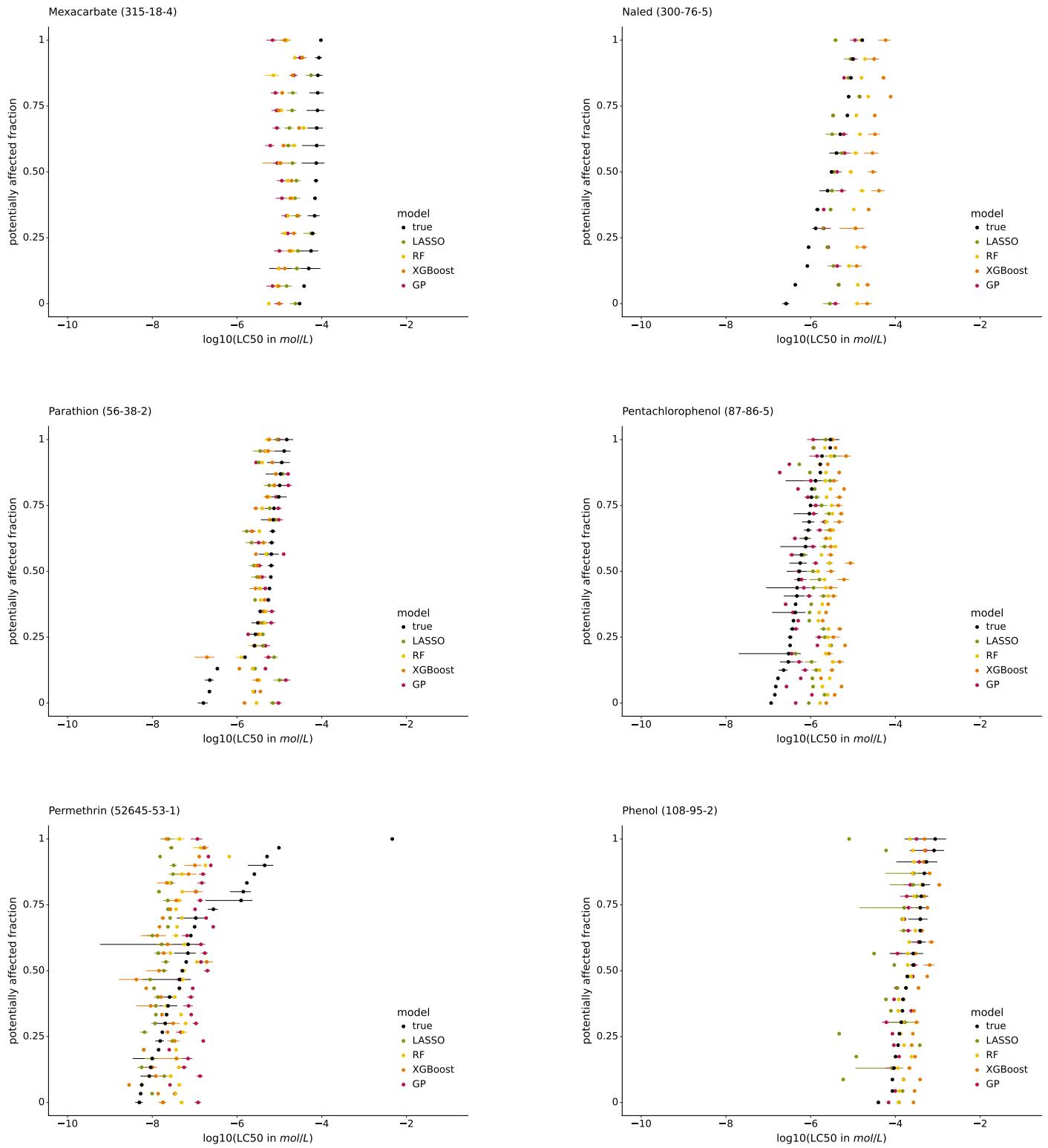


Fig. 10 Species sensitivity distributions (SSDs) of Mexacarbate, Naled, Parathion, Pentachlorophenol, Permethrin, and Phenol.

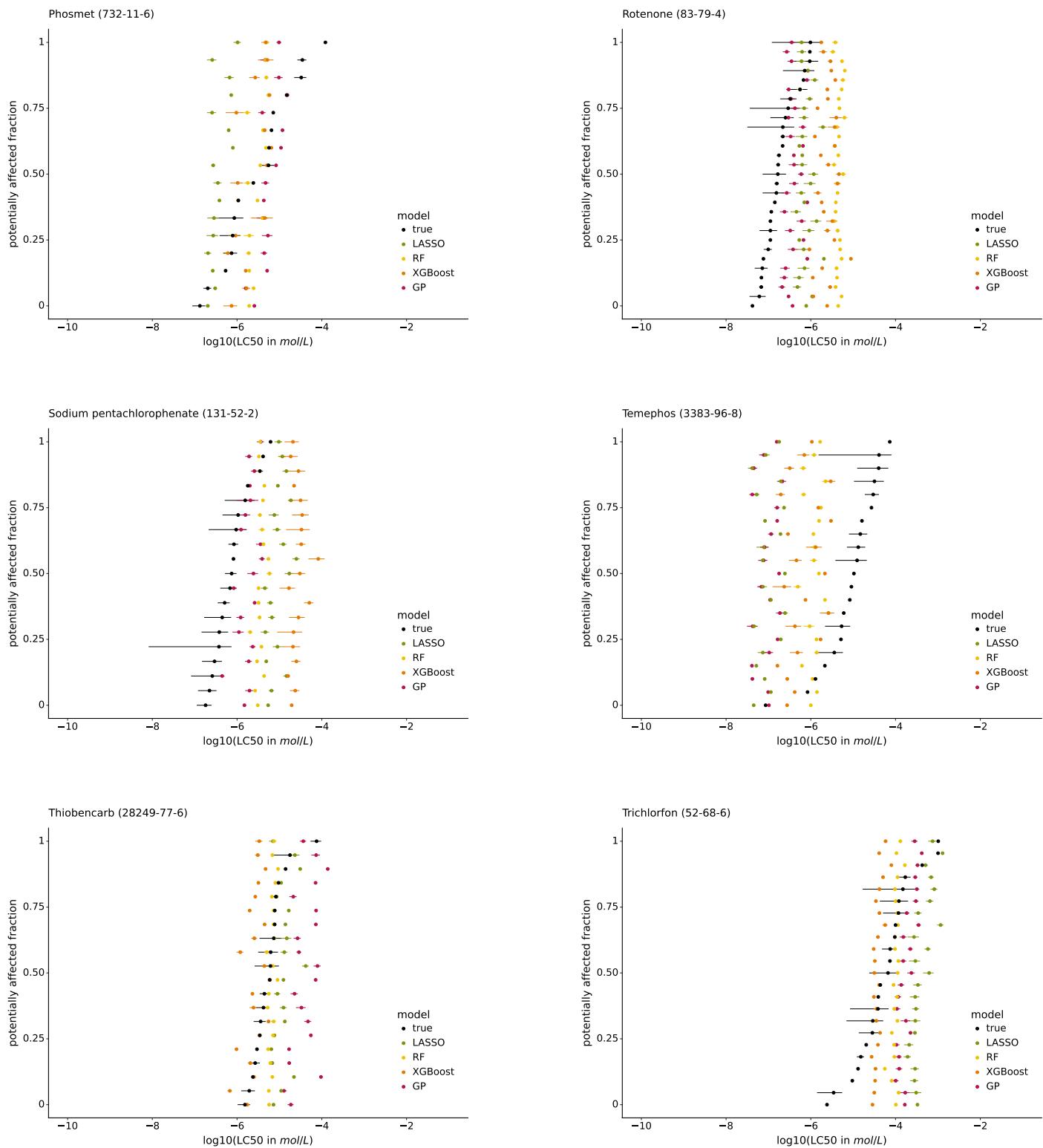


Fig. 11 Species sensitivity distributions (SSDs) of Phosmet, Rotenone, Sodium pentachlorophenate, Temephos, Thiobencarb, and Trichlorfon.

Table 4 Best hyperparameters for random forest. The hyperparameter `min_samples_leaf` was fixed to the default value of 1.

concentration	data split	mol. repr.	n_estimators	max_depth	max_samples	min_samples_split	max_features
molar	totally random	MACCS	300	50	1.0	2	sqrт
molar	totally random	PubChem	150	200	1.0	2	sqrт
molar	totally random	Morgan	300	50	1.0	2	sqrт
molar	totally random	ToxPrint	300	100	1.0	2	sqrт
molar	totally random	mol2vec	300	200	1.0	2	sqrт
molar	totally random	Mordred	300	200	1.0	2	sqrт
molar	occurrence	MACCS	300	200	1.0	2	sqrт
molar	occurrence	PubChem	150	50	1.0	2	sqrт
molar	occurrence	Morgan	300	100	1.0	2	sqrт
molar	occurrence	ToxPrint	50	200	1.0	2	sqrт
molar	occurrence	mol2vec	300	50	1.0	2	sqrт
molar	occurrence	Mordred	50	50	0.25	2	sqrт
mass	totally random	MACCS	300	100	1.0	2	sqrт
mass	totally random	PubChem	300	50	1.0	2	sqrт
mass	totally random	Morgan	300	200	1.0	2	sqrт
mass	totally random	ToxPrint	300	50	1.0	2	sqrт
mass	totally random	mol2vec	300	200	1.0	2	sqrт
mass	totally random	Mordred	300	100	1.0	2	sqrт
mass	occurrence	MACCS	50	200	1.0	2	sqrт
mass	occurrence	PubChem	100	200	0.5	2	sqrт
mass	occurrence	Morgan	150	50	1.0	2	sqrт
mass	occurrence	ToxPrint	100	50	1.0	2	sqrт
mass	occurrence	mol2vec	50	200	1.0	2	sqrт
mass	occurrence	Mordred	150	200	1.0	2	sqrт

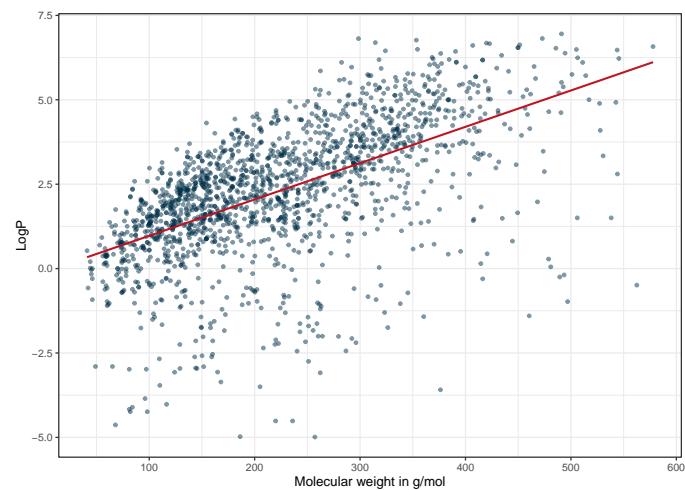


Fig. 12 Correlation of molecular weight and logP for the chemicals in the ADORE "t-F2F" challenge. The red line indicates a linear regression model fit to the data points.

Table 5 Best hyperparameter for XGBoost.

concentration	data split	mol. repr.	n_estimators	eta	gamma	max_depth	min_child_weight	subsample
molar	totally random	MACCS	100	0.100000	0	12	1	1.0
molar	totally random	PubChem	100	0.100000	0	12	3	1.0
molar	totally random	Morgan	100	0.200000	0	12	5	1.0
molar	totally random	ToxPrint	100	0.200000	0	12	3	1.0
molar	totally random	mol2vec	100	0.100000	0	12	3	1.0
molar	totally random	Mordred	100	0.100000	0	12	3	1.0
molar	occurrence	MACCS	100	0.100000	1	9	5	0.5
molar	occurrence	PubChem	100	0.300000	0	3	3	1.0
molar	occurrence	Morgan	100	0.300000	1	3	3	1.0
molar	occurrence	ToxPrint	100	0.200000	0	6	5	1.0
molar	occurrence	mol2vec	100	0.100000	1	6	3	1.0
molar	occurrence	Mordred	100	0.100000	0	6	1	0.5
mass	totally random	MACCS	100	0.100000	0	12	1	0.5
mass	totally random	PubChem	100	0.300000	0	9	1	1.0
mass	totally random	Morgan	100	0.200000	0	12	3	1.0
mass	totally random	ToxPrint	100	0.200000	0	12	5	1.0
mass	totally random	mol2vec	100	0.200000	0	9	1	1.0
mass	totally random	Mordred	100	0.300000	0	9	5	1.0
mass	occurrence	MACCS	100	0.200000	0	6	3	1.0
mass	occurrence	PubChem	100	0.300000	0	3	5	1.0
mass	occurrence	Morgan	100	0.300000	1	3	3	0.5
mass	occurrence	ToxPrint	100	0.200000	10	3	5	0.5
mass	occurrence	mol2vec	100	0.100000	1	6	3	0.5
mass	occurrence	Mordred	100	0.100000	1	6	5	1.0

Table 6 Best hyperparameters for GP regression.

concentration	data split	mol. repr.	n_inducing
molar	totally random	MACCS	1000
molar	totally random	PubChem	1000
molar	totally random	Morgan	1000
molar	totally random	ToxPrint	1000
molar	totally random	mol2vec	1000
molar	totally random	Mordred	1000
molar	occurrence	MACCS	250
molar	occurrence	PubChem	500
molar	occurrence	Morgan	100
molar	occurrence	ToxPrint	500
molar	occurrence	mol2vec	250
molar	occurrence	Mordred	1000
mass	totally random	MACCS	1000
mass	totally random	PubChem	1000
mass	totally random	Morgan	1000
mass	totally random	ToxPrint	1000
mass	totally random	mol2vec	1000
mass	totally random	Mordred	1000
mass	occurrence	MACCS	250
mass	occurrence	PubChem	250
mass	occurrence	Morgan	100
mass	occurrence	ToxPrint	500
mass	occurrence	mol2vec	100
mass	occurrence	Mordred	1000