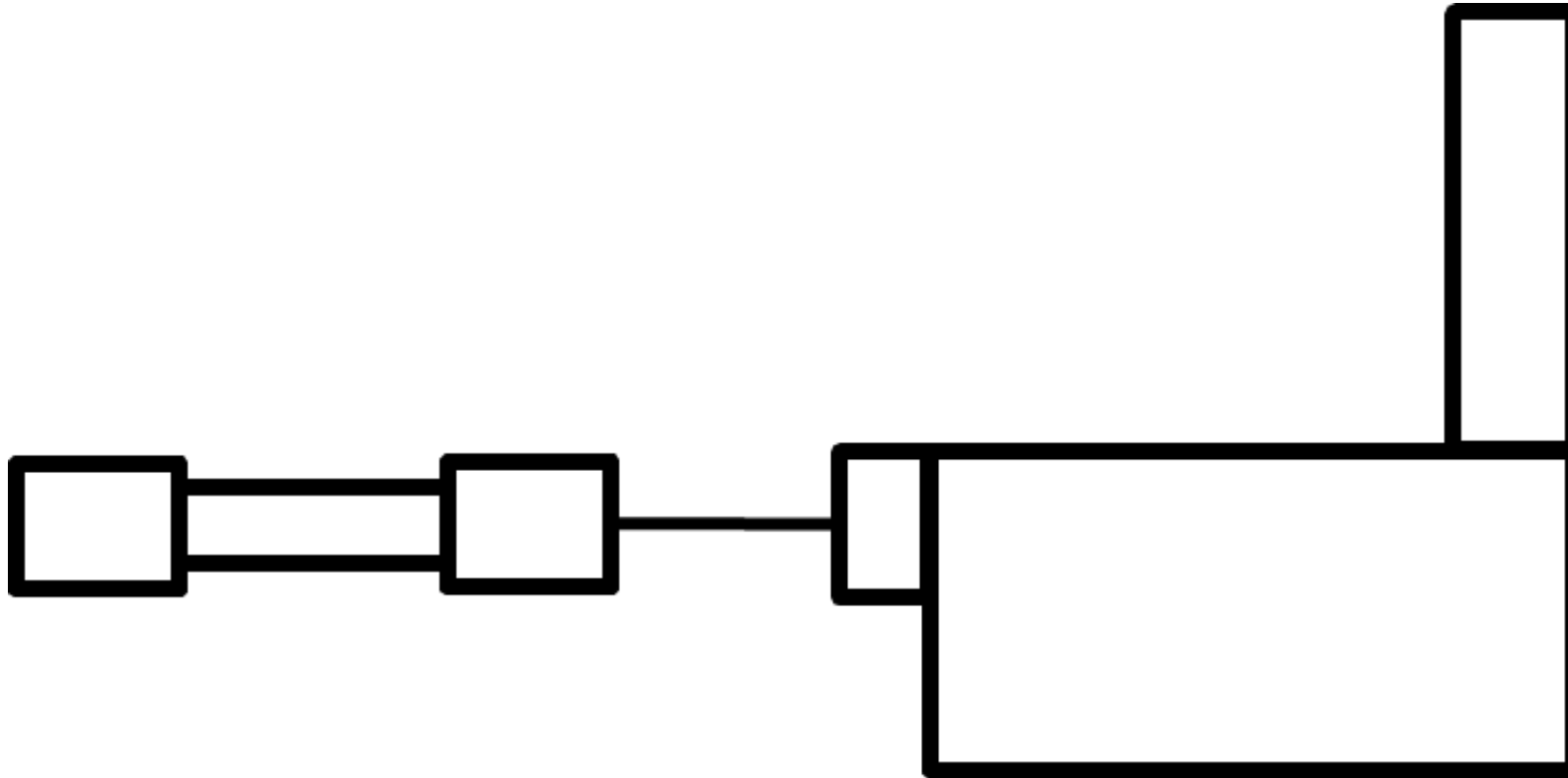# how can machine learning help us to evaluate the risk possessed by emerging contaminants?

anneli kruve

anneli.kruve@su.se

kruvelab.com
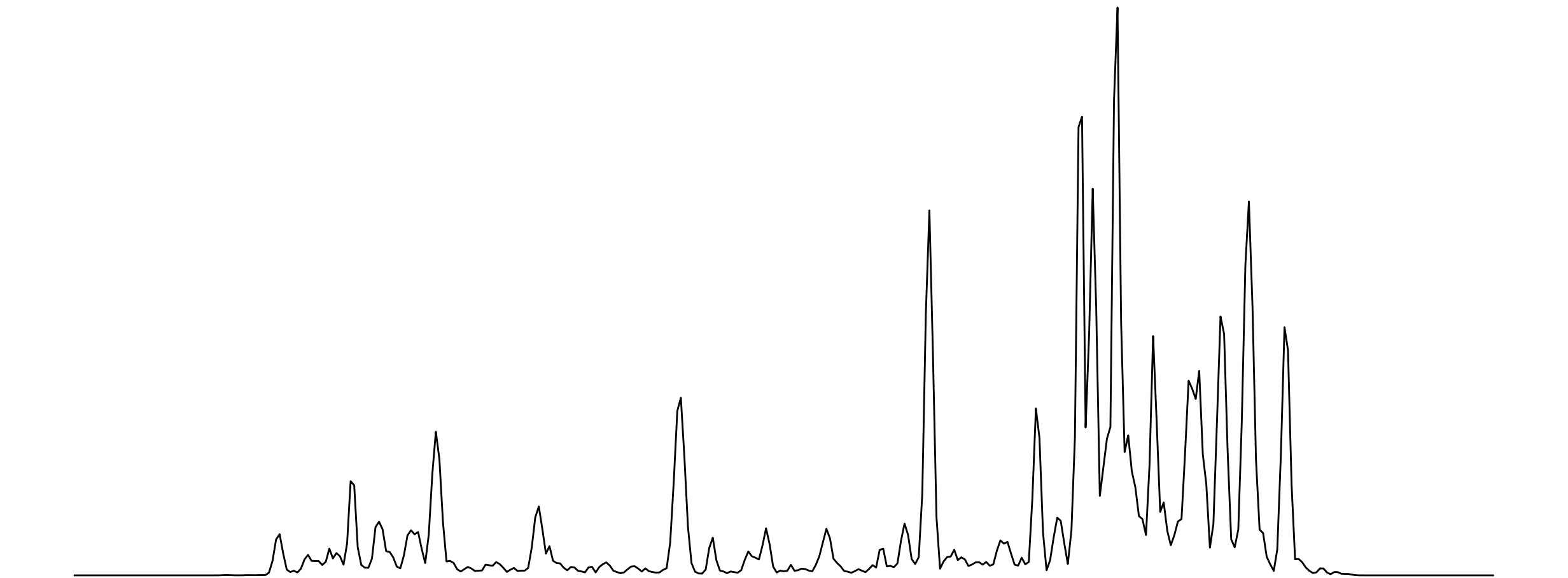
# water analysis

anneli kruve

anneli.kruve@su.se

# nontarget screening with LC/HRMS

# nontarget screening with LC/HRMS



time

anneli kruve

anneli.kruve@su.se

# what next?



time

anneli kruve

anneli.kruve@su.se

# prioritization

toxicity

# prioritization

 toxicity

 concentration

anneli kruve anneli.kruve@su.se

# prioritization

toxicity

concentration
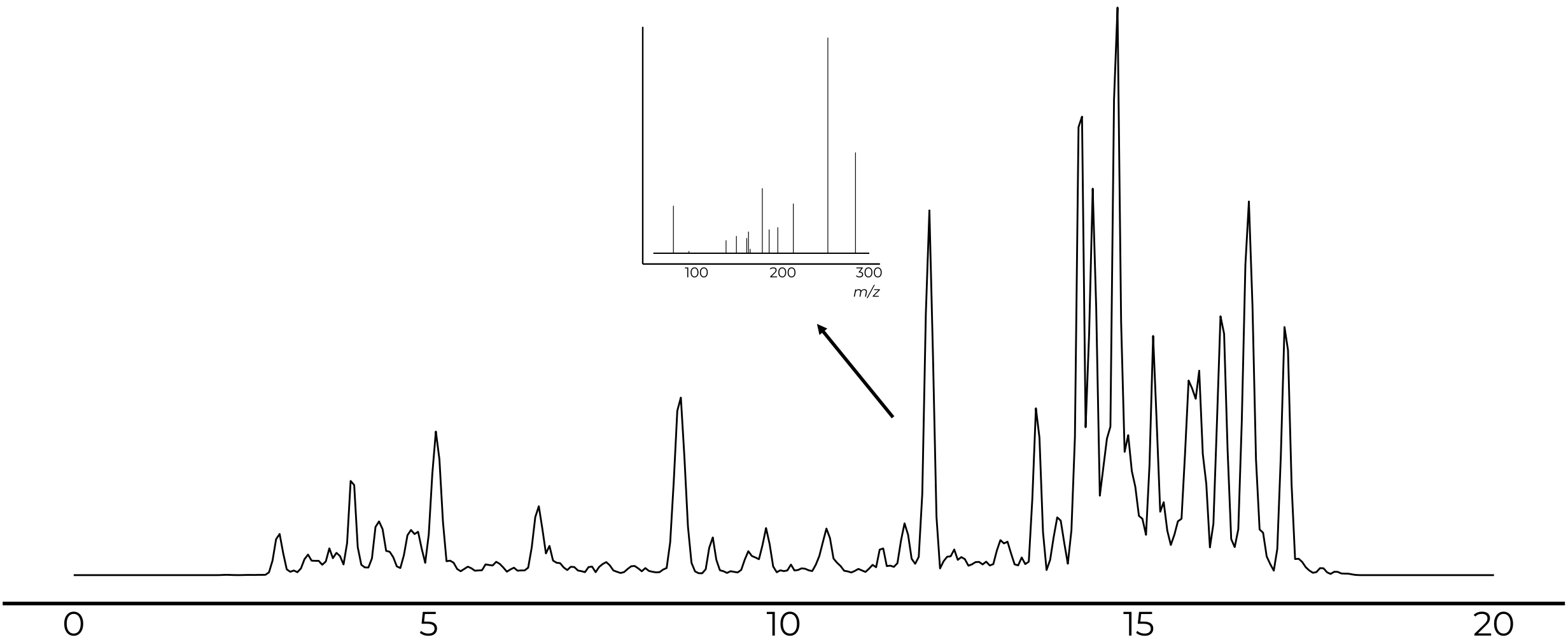
risk

# prioritization

toxicity

concentration

risk

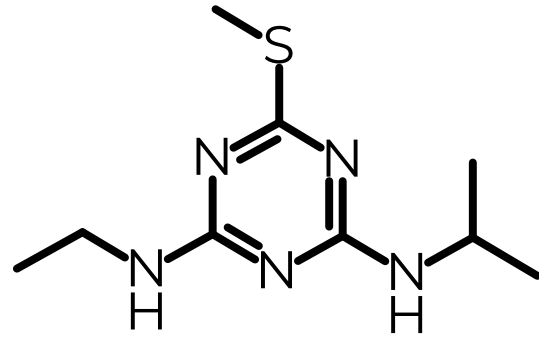$$\text{PriorityScore} = \frac{c_{\text{predicted}}}{\text{AC}_{50}^{\text{5th percentile}}}$$

anneli kruve

anneli.kruve@su.se

# nontarget screening with LC/HRMS



100   200   300
*m/z*

time

anneli kruve

anneli.kruve@su.se

# toxicity assessment



anneli kruve                                                    anneli.kruve@su.se

# toxicity assessment

# toxicity assessment



$LC_{50} = 9.3$ mg/L
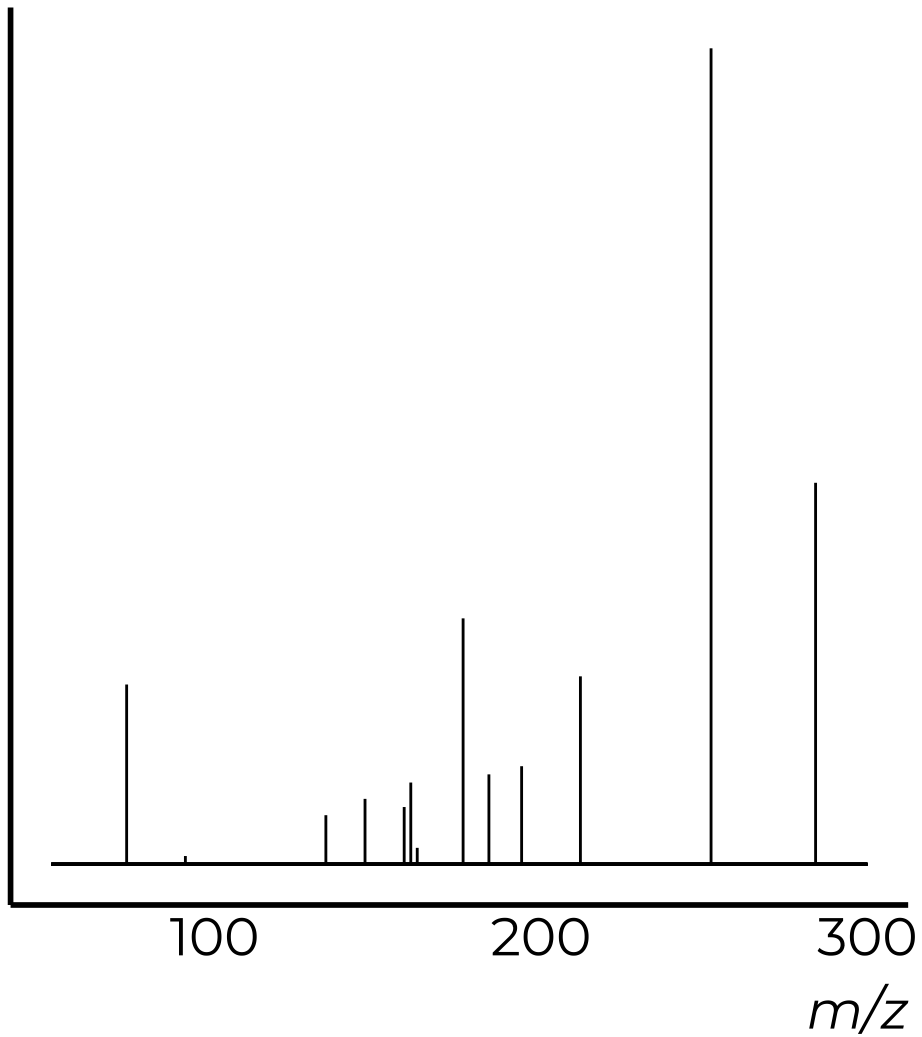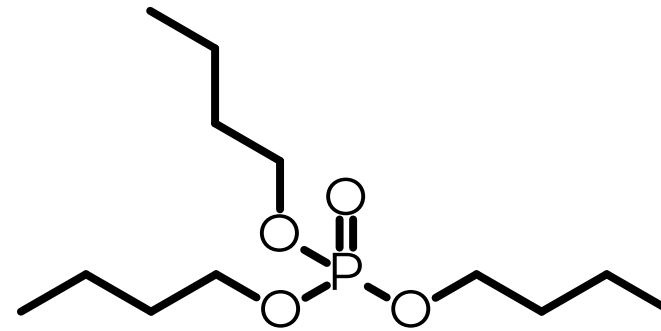
anneli kruve

anneli.kruve@su.se

# toxicity assessment



$LC_{50} = 9.3$ mg/L

# toxicity assessment



$LC_{50}$ = 9.3 mg/L

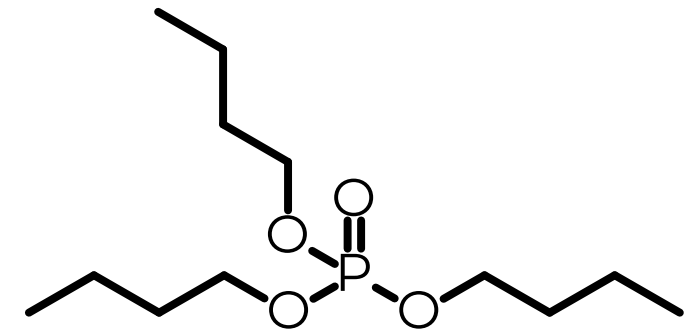$LC_{50}$ = ? mg/L

anneli kruve                    anneli.kruve@su.se

# toxicity assessment



$LC_{50} = 9.3$ mg/L

$LC_{50} = ?$ mg/L

?

anneli kruve                    anneli.kruve@su.se

# toxicity assessment



LC$_{50}$ = 9.3 mg/L

LC$_{50}$ = ? mg/L

LC$_{50}$ = ? mg/L

anneli kruve

anneli.kruve@su.se

# toxicity assessment

LC$_{50}$ = 9.3 mg/L

LC$_{50}$ = ? mg/L

**?**

LC$_{50}$ = ? mg/L

anneli kruve

anneli.kruve@su.se

# toxicity assessment

<1%



LC$_{50}$ = 9.3 mg/L

LC$_{50}$ = ? mg/L

?

LC$_{50}$ = ? mg/L

anneli kruve                                anneli.kruve@su.se

# toxicity assessment

<1%

$LC_{50} = 9.3$ mg/L

<2%

$LC_{50} = ?$ mg/L

?

$LC_{50} = ?$ mg/L

anneli kruve

anneli.kruve@su.se

# toxicity assessment

<1%   LC$_{50}$ = 9.3 mg/L

<2%   LC$_{50}$ = ? mg/L

~98%   LC$_{50}$ = ? mg/L

?

anneli kruve                                    anneli.kruve@su.se

# predicting toxicity

for detected chemicals

# workflow

MS$^2$ spectra

structure as SMILES

molecular descriptors

predict toxicity

anneli kruve      anneli.kruve@su.se

# workflow

MS$^2$ spectra

molecular descriptors

predict toxicity

anneli kruve                    anneli.kruve@su.se

# information available

in MS$^2$ spectra

# MS² spectra



*m/z*

anneli kruve anneli.kruve@su.se

# MS² spectra



-C₃H₅O₂Cl

$m/z$

# data for machine learning models

anneli kruve                                    anneli.kruve@su.se

# data for machine learning models

CompTox

all toxicity
values

# data for machine learning models

CompTox

all toxicity
values

one species

# data for machine learning models

CompTox

all toxicity
values

LC$_{50}$

one species

# data for machine learning models



CompTox

all toxicity
values

same
conditions

$LC_{50}$

one species
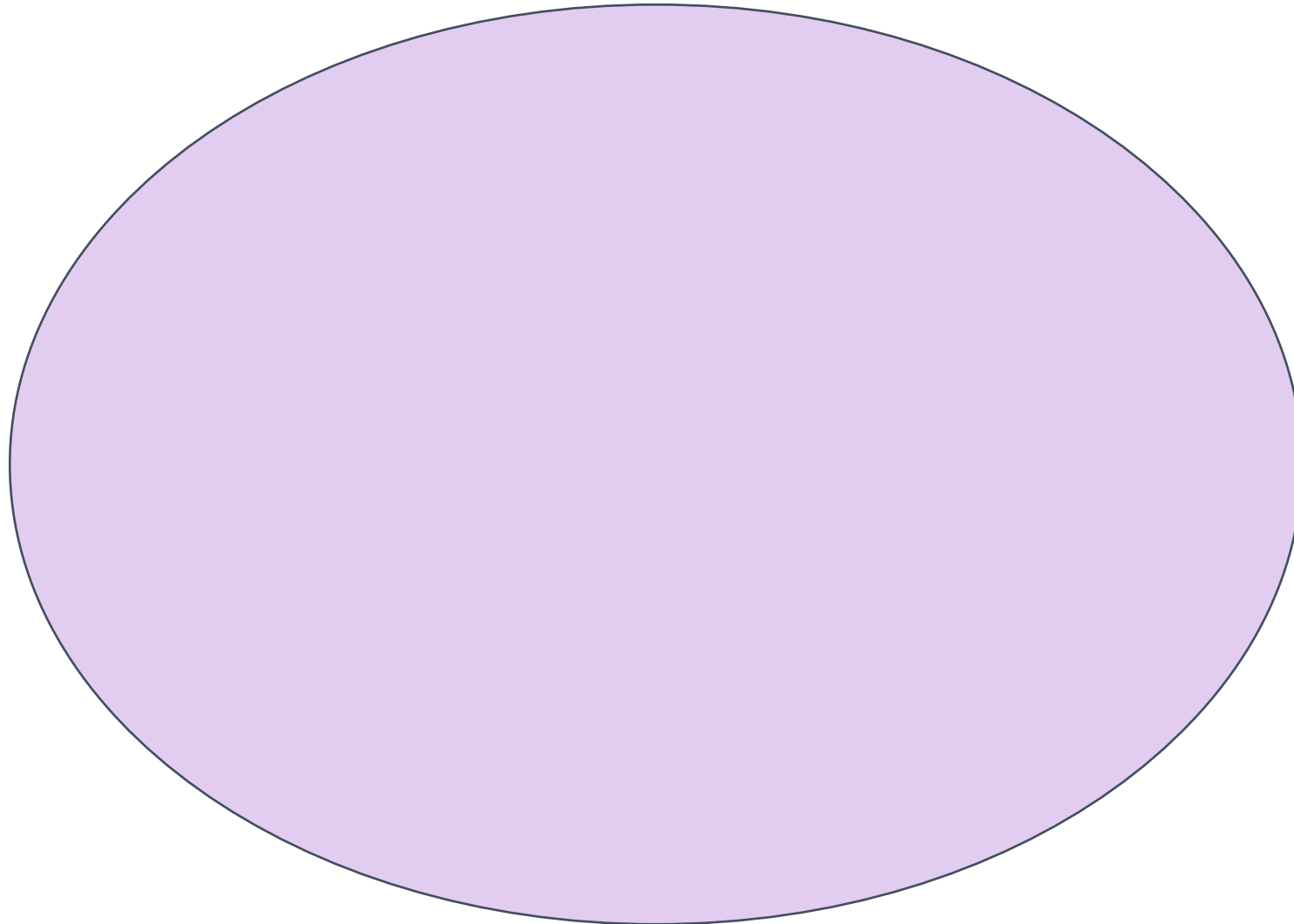
anneli kruve

anneli.kruve@su.se

# data for machine learning models

MassBank

all
spectra

anneli kruve                    anneli.kruve@su.se

# data for machine learning models



MassBank

LC/HRMS

all
spectra

anneli kruve                                            anneli.kruve@su.se

# data for machine learning models



CompTox

MassBank

all toxicity values

same conditions

$LC_{50}$

one species

LC/HRMS

all spectra

anneli kruve

anneli.kruve@su.se

# data for machine learning models



CompTox

MassBank

same
conditions

all toxicity
values

$LC_{50}$

LC/HRMS

all
spectra

one species

anneli kruve

anneli.kruve@su.se

# predicting toxicity

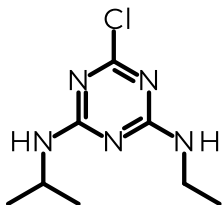from the structure

# workflow

structure as SMILES

molecular fingerprints

machine learning for predicting toxicity

anneli kruve                                    anneli.kruve@su.se

# selected endpoint

anneli kruve                anneli.kruve@su.se

# selected endpoint

fathead minnow, bluegill, and rainbow trout

anneli kruve                                                                    anneli.kruve@su.se

# selected endpoint

fathead minnow, bluegill, and rainbow trout

water flea

# selected endpoint

fathead minnow, bluegill, and rainbow trout

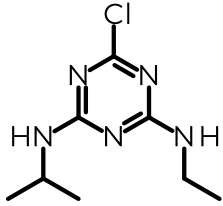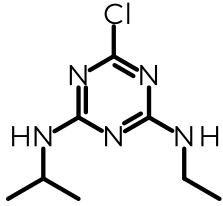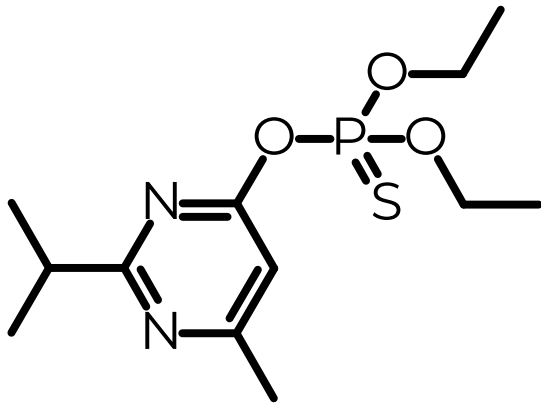water flea

algae

# workflow

structure as SMILES

anneli.kruve@su.se

# workflow

structure as SMILES

molecular fingerprints

anneli.kruve@su.se

# structural fingerprints

# structural fingerprints

R: rcdk

# structural fingerprints

R: rcdk →

| 0 | (tetrahydropyran ring with O) |
|---|---|
| 1 | O—P |
| 1 | —N |
| 0 | —NH$_2$ |
| 1 | (pyrimidine ring) |

anneli.kruve@su.se

# workflow

structure as SMILES

molecular fingerprints

machine learning for predicting $LC_{50}$

anneli kruve

anneli.kruve@su.se

# model training

| mass (Da) | fp1 | ... | fp243 |
|-----------|-----|-----|-------|
| 317.32000 | 0 | ... | 0 |
| 208.26100 | 1 | ... | 0 |
| 240.21499 | 1 | ... | 0 |
| 300.57998 | 0 | ... | 0 |
| 201.22500 | 0 | ... | 0 |

# model training

| mass (Da) | fp1 | ... | fp243 |
|-----------|-----|-----|-------|
| 317.32000 | 0 | ... | 0 |
| 208.26100 | 1 | ... | 0 |
| 240.21499 | 1 | ... | 0 |
| 300.57998 | 0 | ... | 0 |
| 201.22500 | 0 | ... | 0 |

training set
517
chemicals

test set
130
chemicals

# model training

| mass (Da) | fp1 | ... | fp243 |
|-----------|-----|-----|-------|
| 317.32000 | 0 | ... | 0 |
| 208.26100 | 1 | ... | 0 |
| 240.21499 | 1 | ... | 0 |
| 300.57998 | 0 | ... | 0 |
| 201.22500 | 0 | ... | 0 |

training set
517
chemicals

gradient
boosting

test set
130
chemicals

anneli kruve

anneli.kruve@su.se

# performance

of $LC_{50}$ predictions with molecular fingerprints

# LC$_{50}$ predictions

## Peets et al. ES&T 2022

fish LC$_{50}$



training set

RMSE 0.52 log(M)

anneli kruve

anneli.kruve@su.se

# LC$_{50}$ predictions

Peets et al. ES&T 2022

fish LC$_{50}$



training set

RMSE 0.52 log(M)

test set

RMSE 0.78 log(M)

anneli kruve

anneli.kruve@su.se

# unidentified chemicals

from MS$^2$ spectra

# workflow

$MS^2$ spectra

molecular fingerprints with SIRIUS+CSI:FingerID

predict $LC_{50}$ with pretrained gradient boosting

anneli kruve

anneli.kruve@su.se

# predict for unknown chemicals



$\longrightarrow$ ?

anneli kruve                                    anneli.kruve@su.se

# predict for unknown chemicals

# predict for unknown chemicals



SIRIUS+
CSI:FingerID

$C_{12}H_{21}N_2O_3PS$

$C_2H_4$

$C_{10}H_{17}N_2O_3PS$

$C_2H_4$     $HO_2PS$

$C_8H_{13}N_2O_3PS$     $C_{10}H_{16}N_2O$

...     ...

*m/z*

100     200     300

anneli kruve                                        anneli.kruve@su.se

# predict for unknown chemicals



SIRIUS+
CSI:FingerID

$C_{12}H_{21}N_2O_3PS$

$C_2H_4$

$C_{10}H_{17}N_2O_3PS$

$C_2H_4$         $HO_2PS$

$C_8H_{13}N_2O_3PS$      $C_{10}H_{16}N_2O$

...        ...

| 0.001 | |
|---|---|
| 0.999 | O—P |
| 0.999 | —N |
| 0.198 | —NH$_2$ |
| 0.988 | |

anneli kruve                                    anneli.kruve@su.se

# predict for unknown chemicals



$C_{12}H_{21}N_2O_3PS$

$C_2H_4$

$C_{10}H_{17}N_2O_3PS$

SIRIUS+
CSI:FingerID

$C_2H_4$          $HO_2PS$

$C_8H_{13}N_2O_3PS$          $C_{10}H_{16}N_2O$

...          ...

| | |
|---|---|
| 0 | |
| 1 | O—P |
| 1 | —N |
| 0 | —NH$_2$ |
| 1 | |

# predict for unknown chemicals



SIRIUS+
CSI:FingerID

| | |
|---|---|
| 0 | ⬡O |
| 1 | O—P |
| 1 | —N |
| 0 | —NH$_2$ |
| 1 | (pyrimidine ring) |

gradient
boosting

LC$_{50}$ =
-2.2 log(mM)

anneli kruve                                                    anneli.kruve@su.se

# LC$_{50}$ predictions

# LC$_{50}$ predictions

Peets et al. ES&T 2022

fish LC$_{50}$



test set on structures

RMSE 0.78 log(M)

anneli kruve                    anneli.kruve@su.se

# LC$_{50}$ predictions

Peets et al. ES&T 2022

fish LC$_{50}$



test set on structures

RMSE 0.78 log(M)

validation on MassBank

RMSE$_{model}$ 0.88 log(M)

SD$_{experimental}$ 0.44 log(mM)
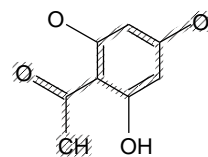
# model interpretation

anneli kruve                    anneli.kruve@su.se

# model interpretation

# model interpretation



shap

0    250    500    750

*m/z*

anneli kruve                                        anneli.kruve@su.se
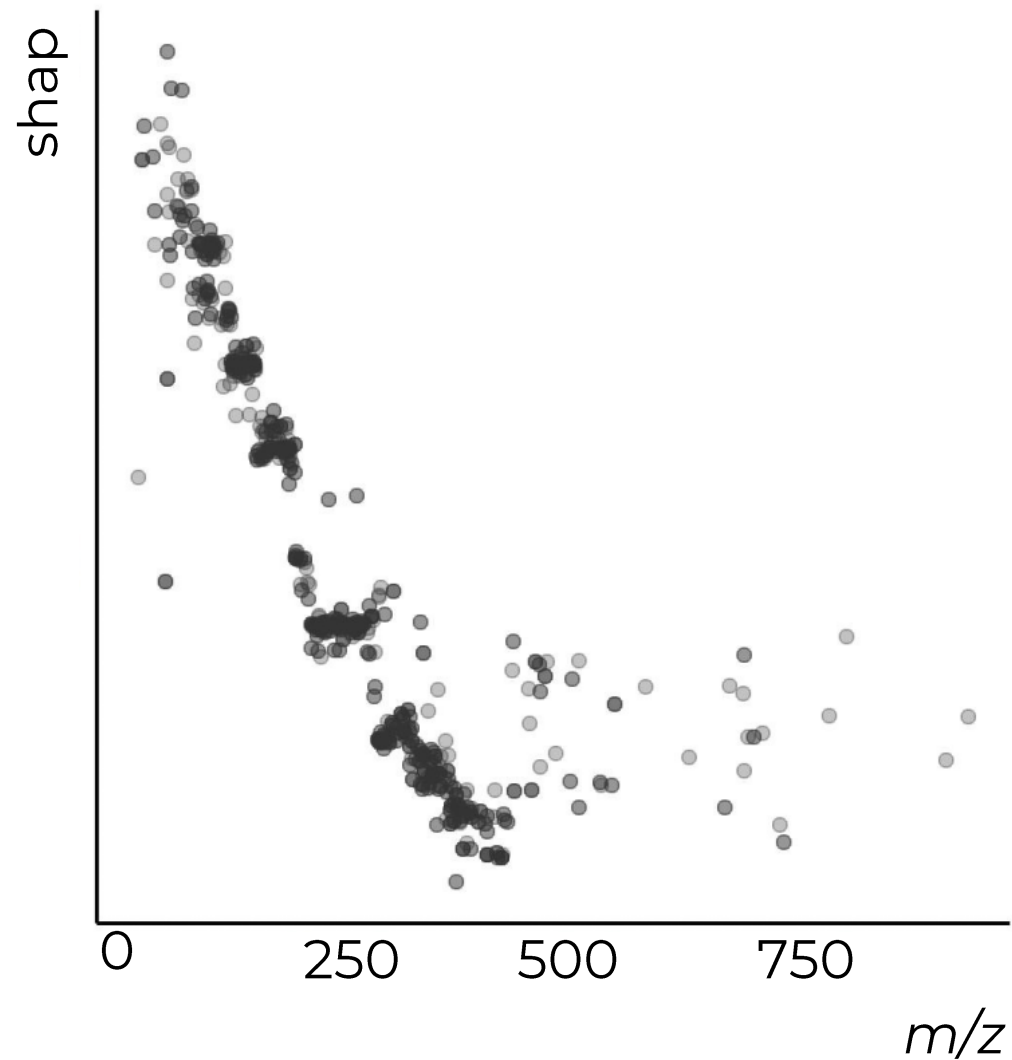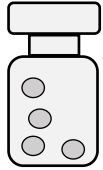
# model interpretation

# model interpretation

# toxic chemicals

in wastewater

# case study on wastewater

wastewater samples

anneli.kruve@su.se

# case study on wastewater

wastewater samples

LC/HRMS analysis

anneli kruve                    anneli.kruve@su.se
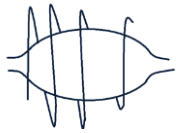
# case study on wastewater

anneli kruve

anneli.kruve@su.se

# case study on wastewater

wastewater samples

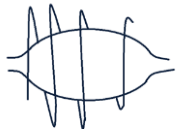LC/HRMS analysis

molecular fingerprints with SIRIUS+CSI:FingerID

anneli kruve
anneli.kruve@su.se

# case study on wastewater

wastewater samples

LC/HRMS analysis

molecular fingerprints with SIRIUS+CSI:FingerID

predict $LC_{50}$ with pretrained gradient boosting

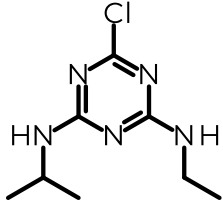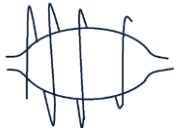anneli kruve                                    anneli.kruve@su.se

quality control

# quality control

216 analytical standard

# quality control

216 analytical standard

DIA and DDA MS$^2$ data

anneli kruve                                    anneli.kruve@su.se

# quality control

216 analytical standard

DIA and DDA MS$^2$ data

comparison with experimental LC$_{50}$

anneli kruve                                        anneli.kruve@su.se

# DDA



lab water · groundwater · surface water · wastewater

$LC_{50}^{predicted}$ (M) vs $LC_{50}^{experimental}$ (M)

RMSE = 0.95 log-mM · 0.74 log-mM · 0.86 log-mM · 0.47 log-mM

anneli kruve

anneli.kruve@su.se

# DIA



lab water    groundwater    surface water    wastewater

$LC_{50}^{predicted}$ (M)

$LC_{50}^{experimental}$ (M)

RMSE = 0.85 log-mM    1.09 log-mM    1.18 log-mM    1.03 log-mM
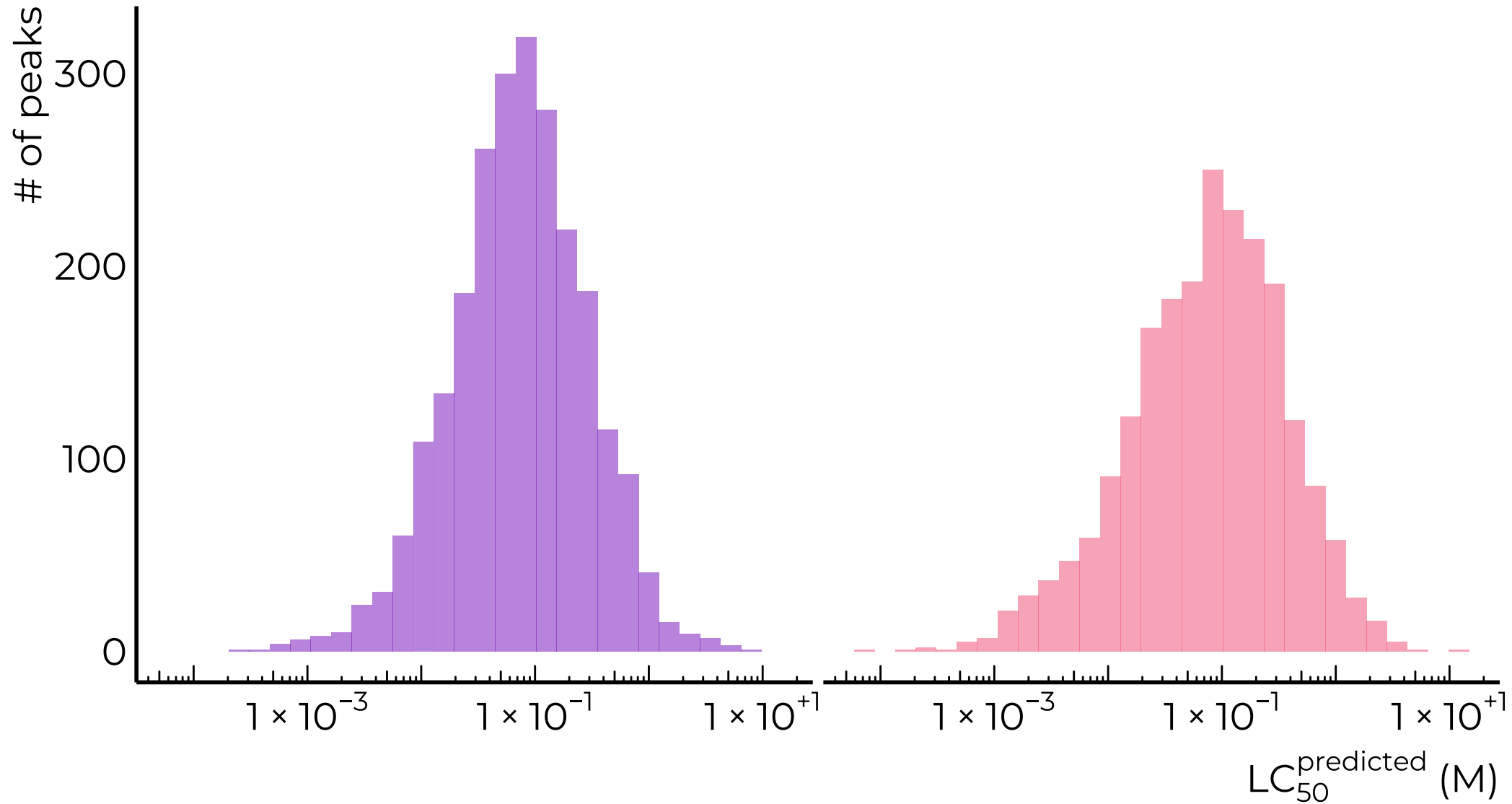
anneli kruve                                    anneli.kruve@su.se

pinpointing toxic chemicals

# case study on wastewater



anneli kruve

anneli.kruve@su.se
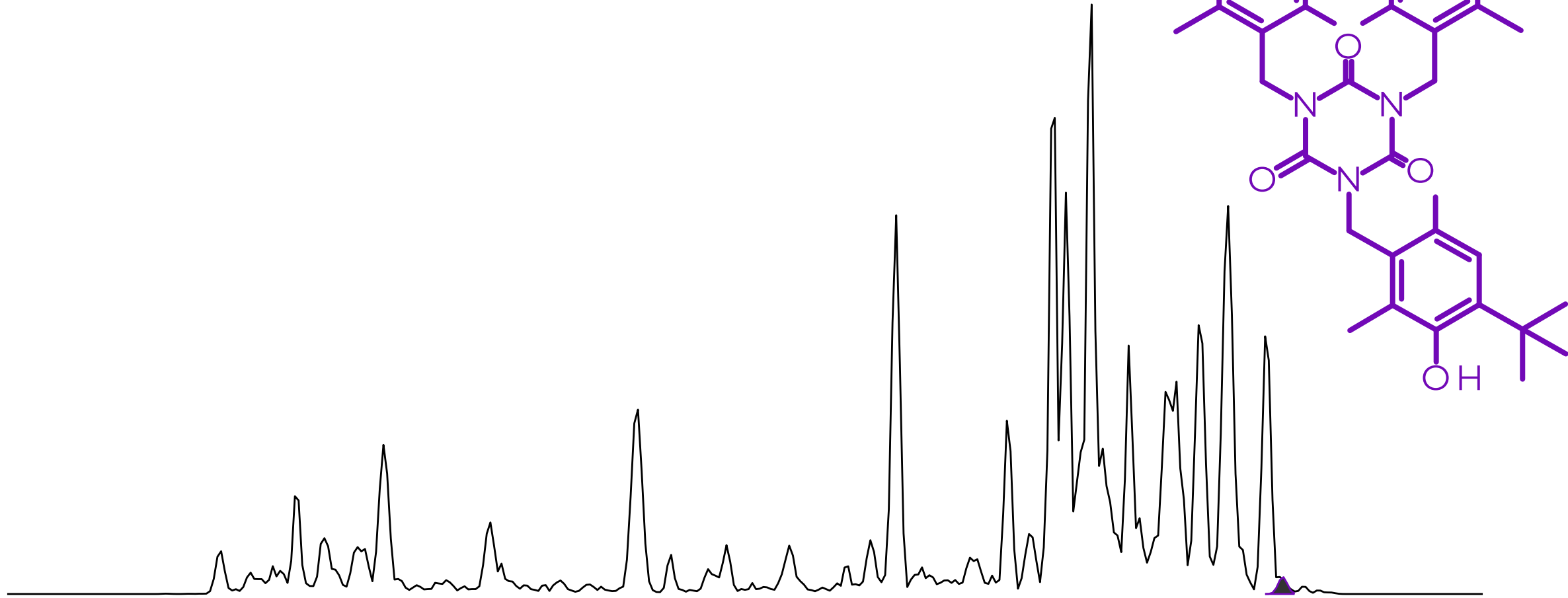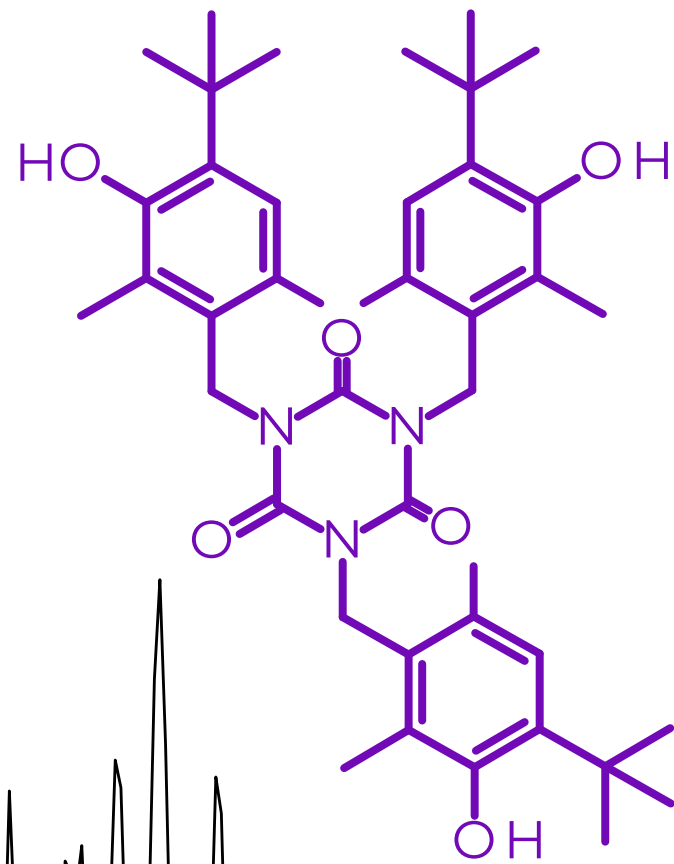
# LC$_{50}$ distribution



anneli kruve

anneli.kruve@su.se

# naproxen

anneli kruve

anneli.kruve@su.se

time

# cyanox CY 1790



anneli kruve

anneli.kruve@su.se

summary

# prioritization in NTS

toxicity

concentration

risk

anneli kruve                                    anneli.kruve@su.se

# prioritization in NTS
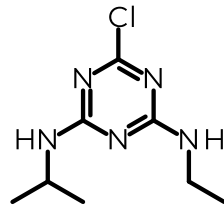
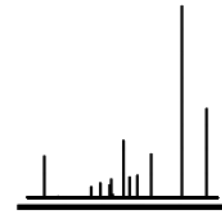toxicity                    concentration                    risk



structure                    MS$^2$ spectrum

anneli kruve                                    anneli.kruve@su.se

kruvelab.com

anneli.kruve@su.se