

toxicity predictions of  
unidentified chemicals in water  
by HRMS & ML:  
how and why?

anneli kruve  
anneli.kruve@su.se  
kruvelab.com

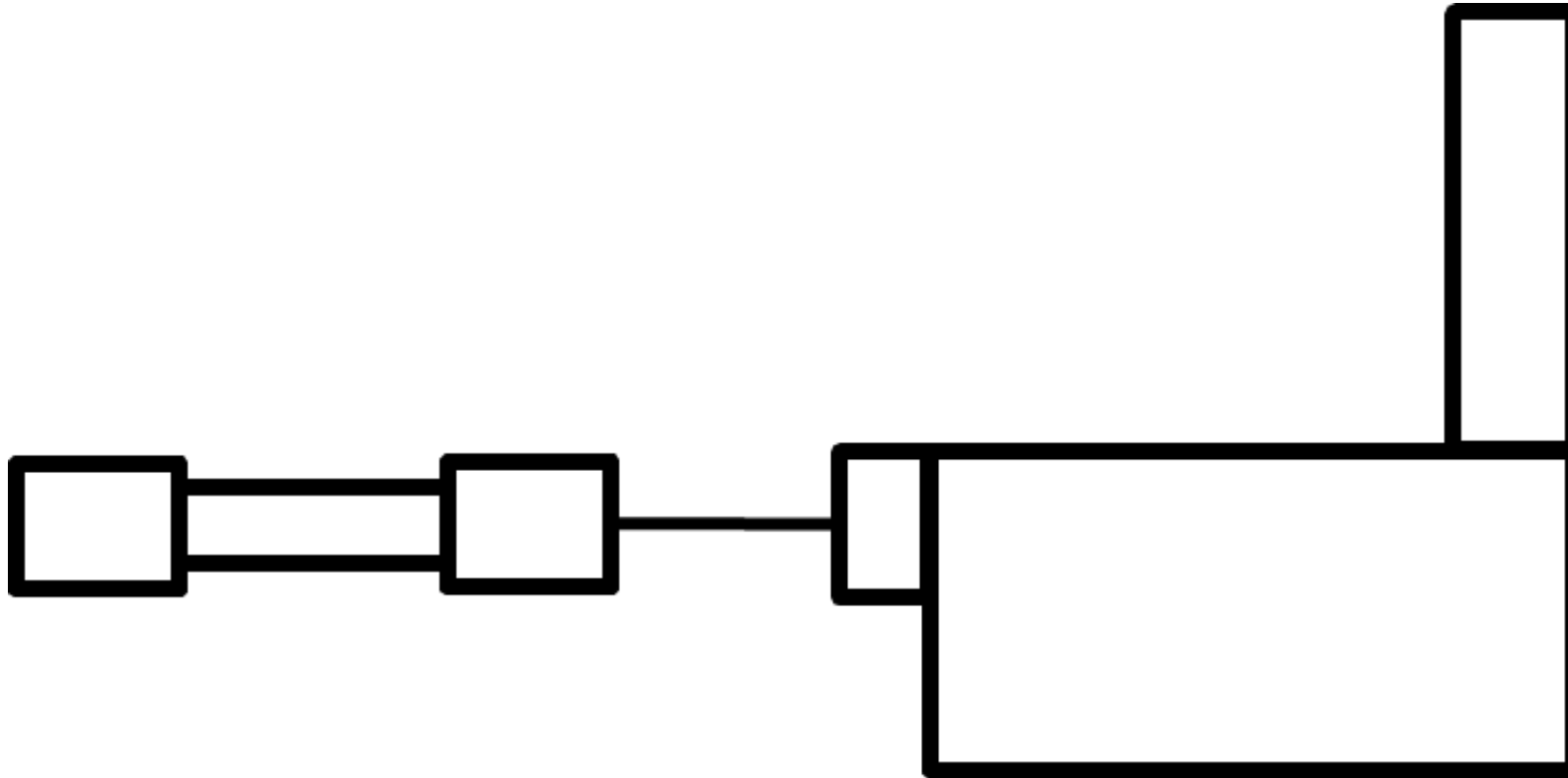


# water analysis

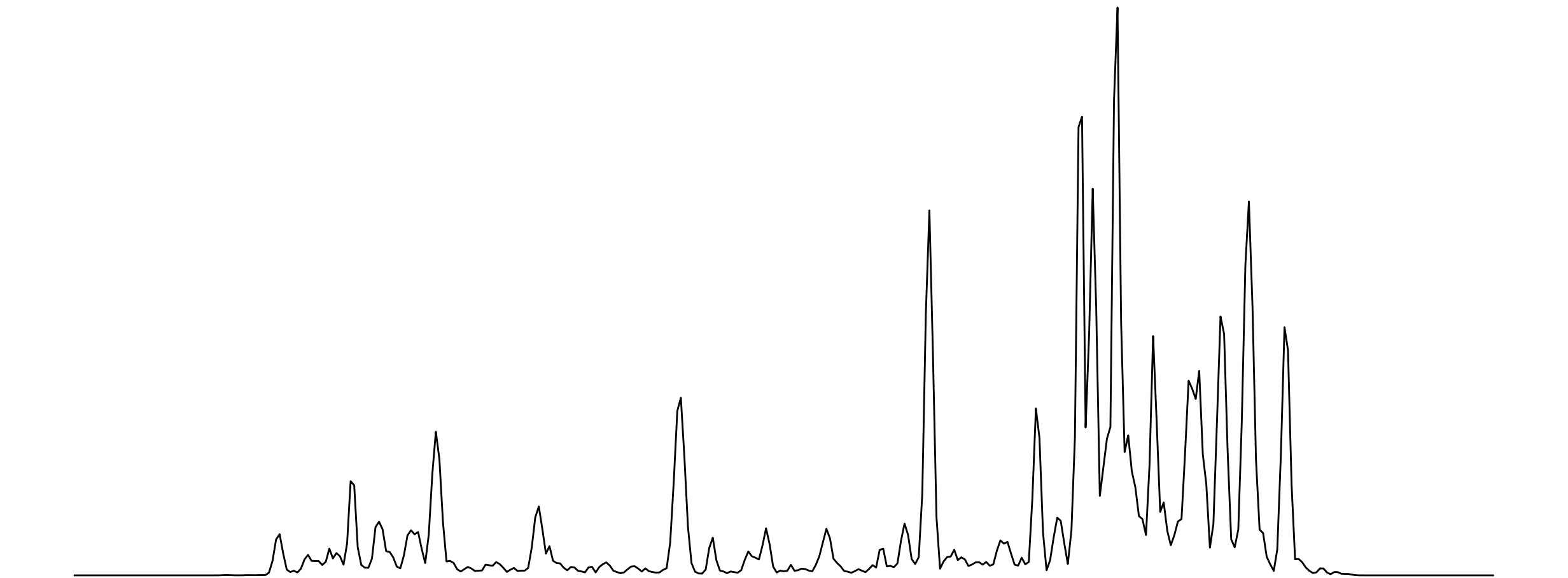




# nontarget screening with LC/HRMS



# nontarget screening with LC/HRMS



0

5

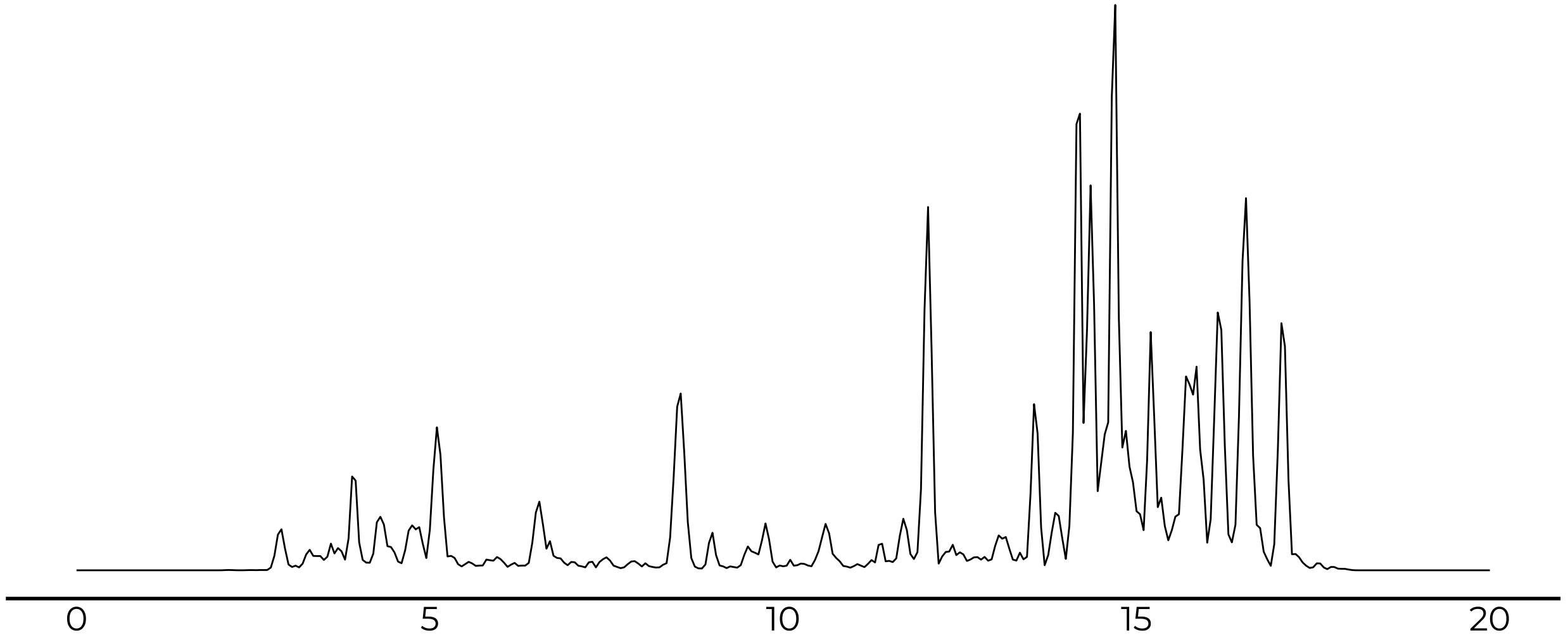
10

15

20

time

# what next?



# prioritization



toxicity

# prioritization



toxicity



concentration

# prioritization



toxicity



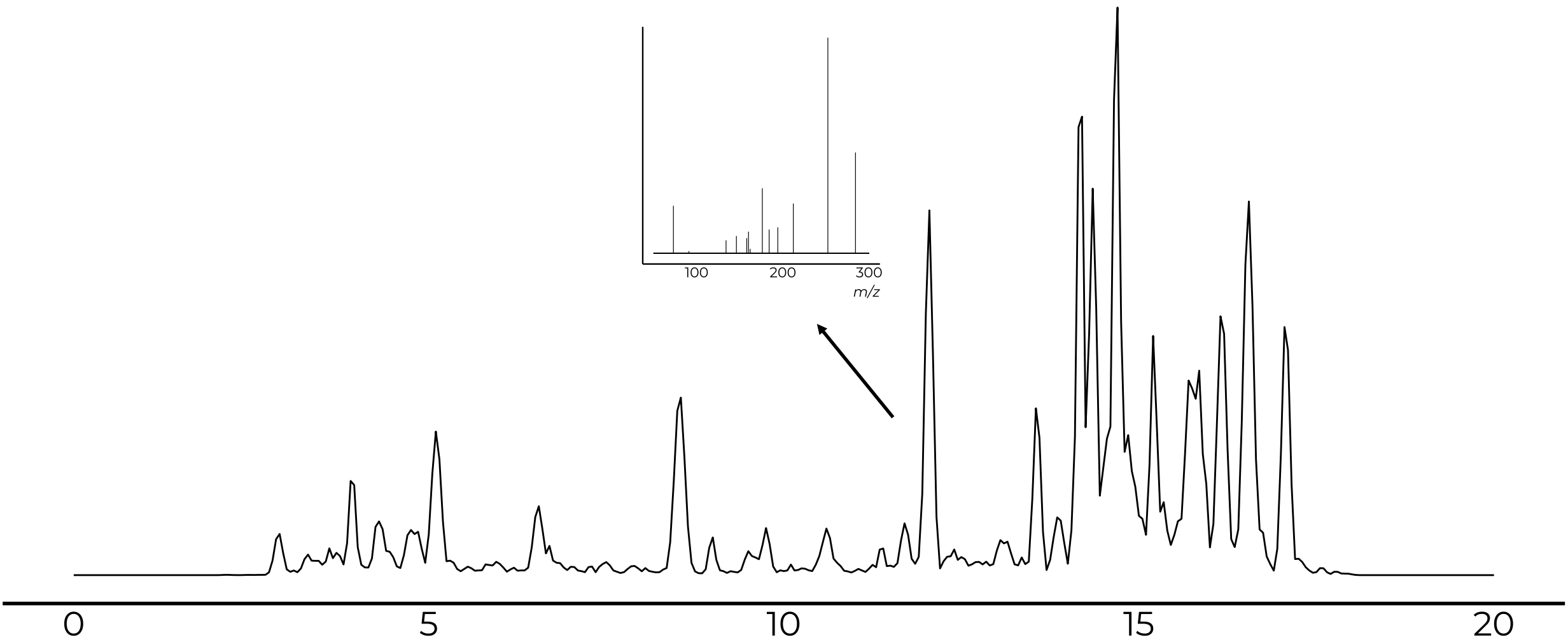
concentration



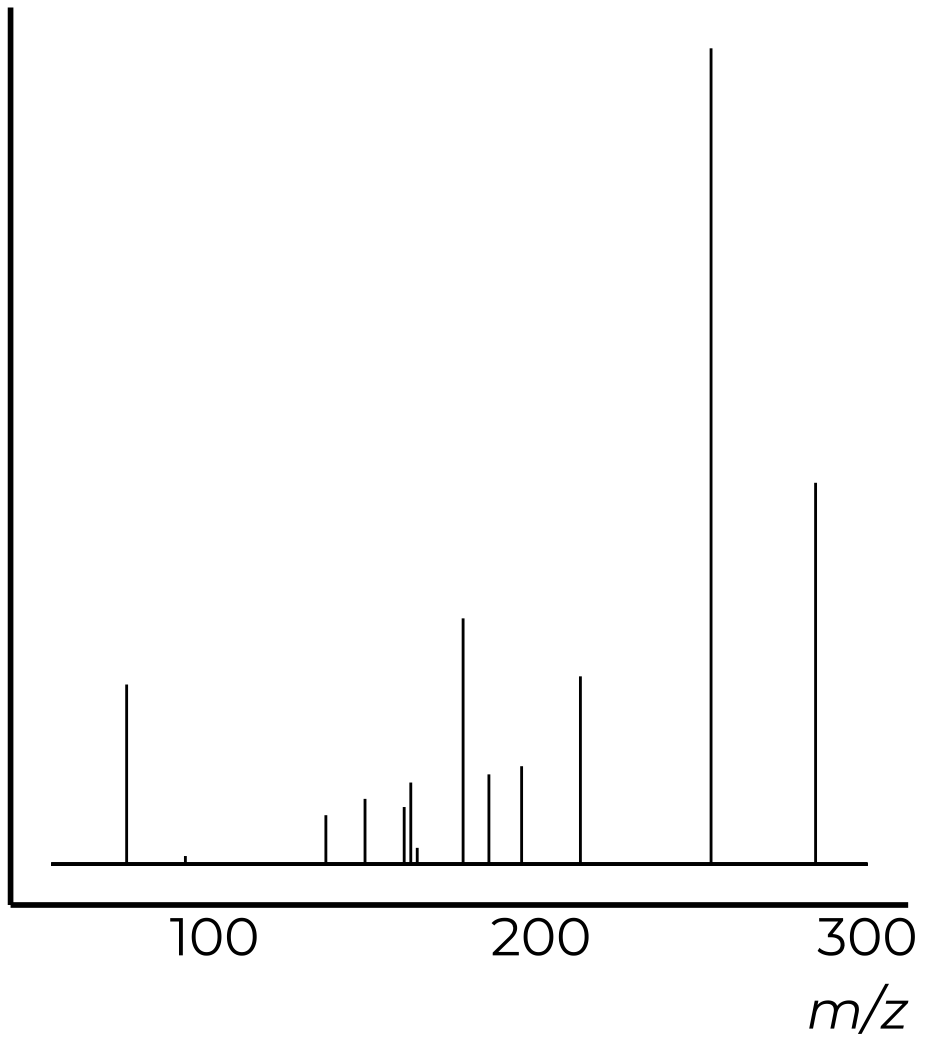
risk



# nontarget screening with LC/HRMS

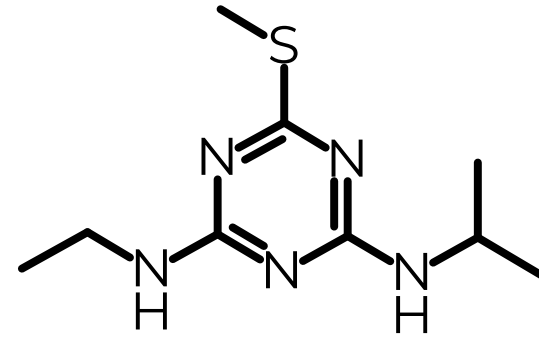
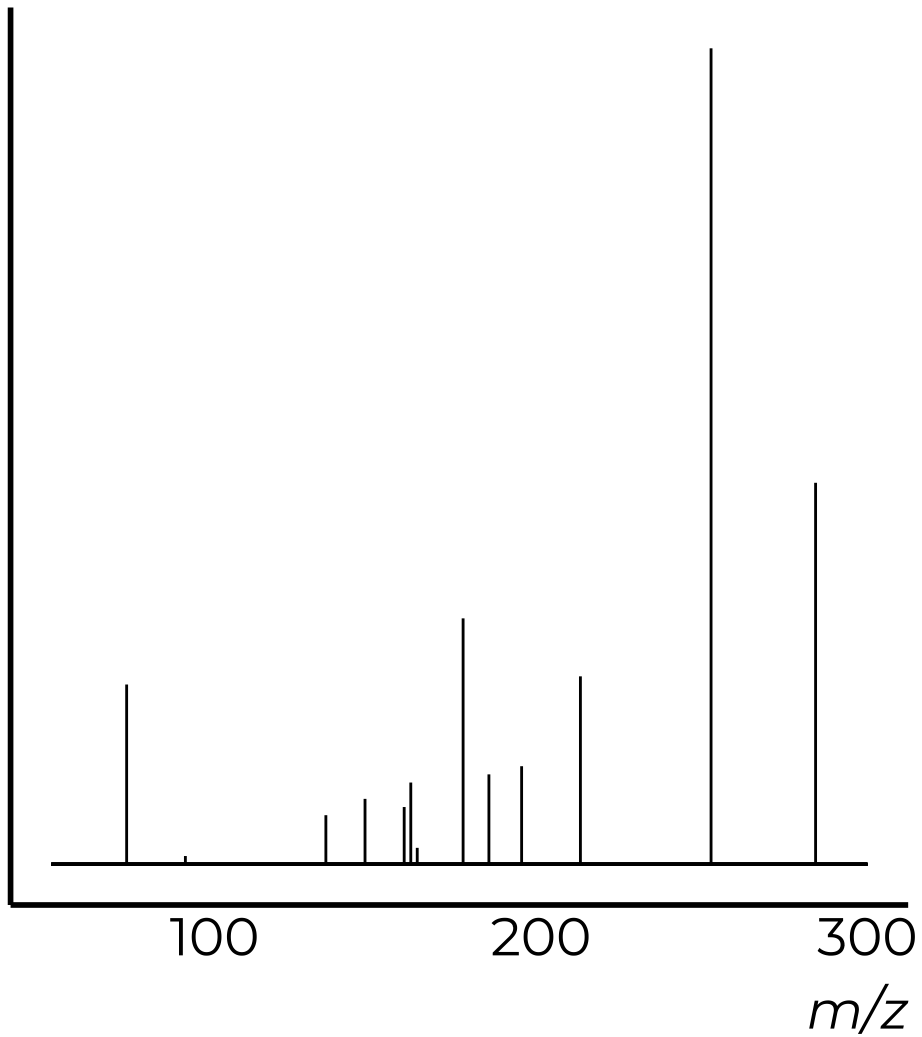


# toxicity assessment

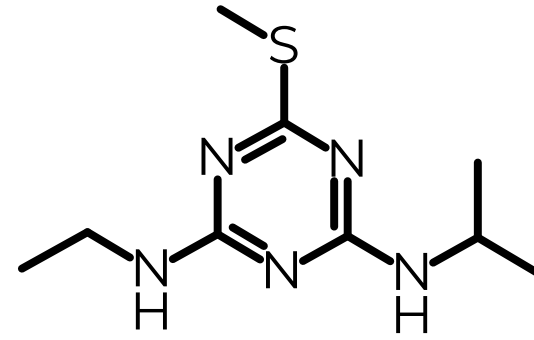
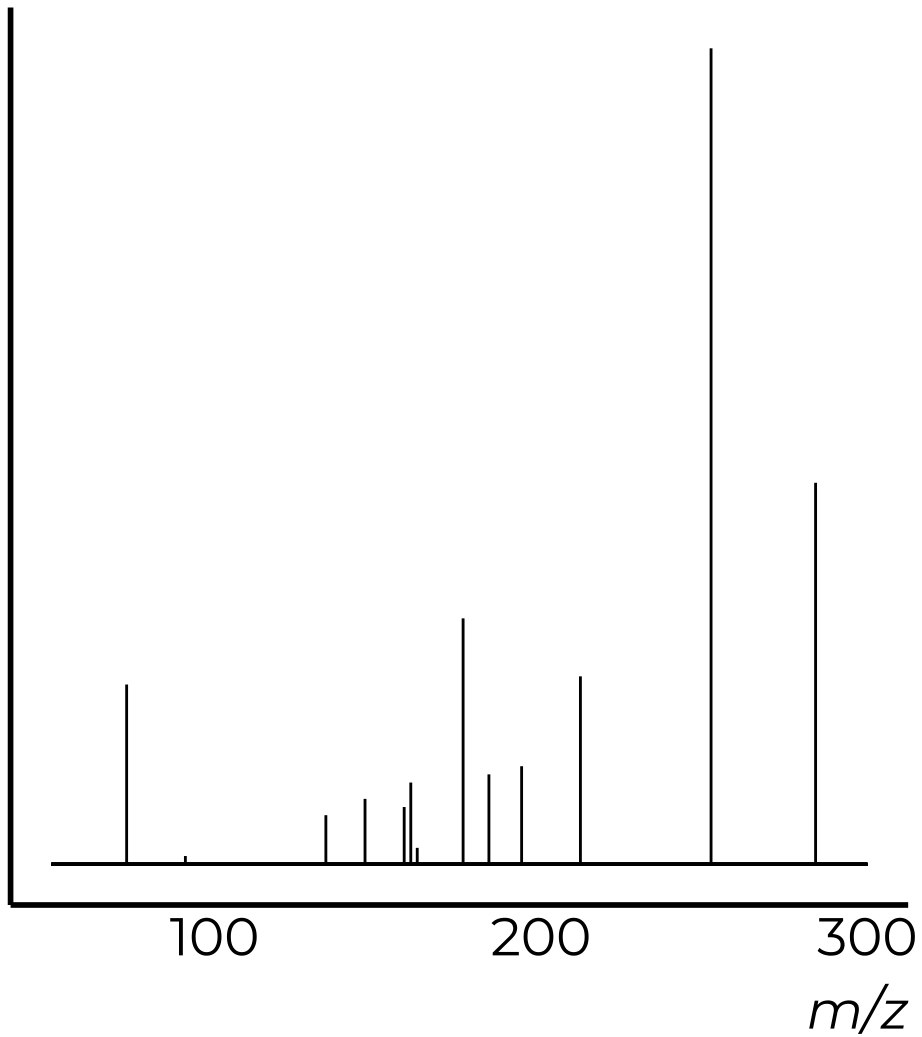




# toxicity assessment



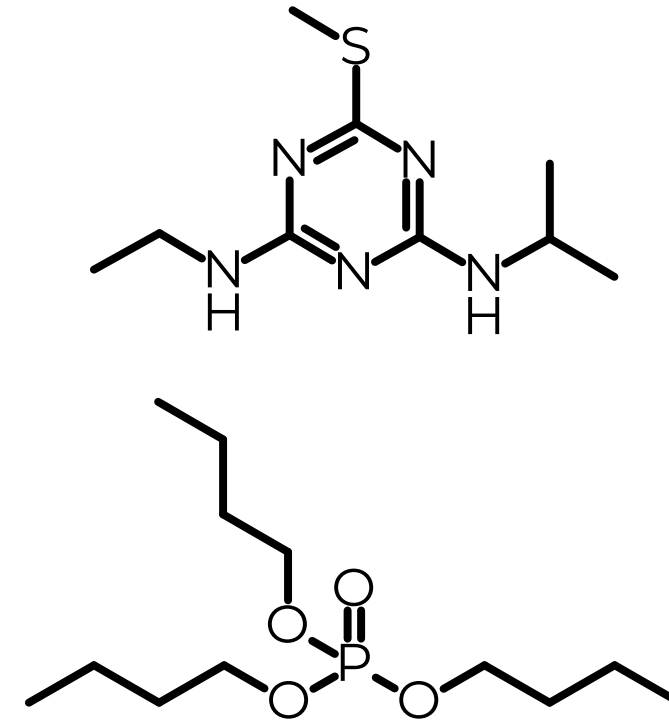
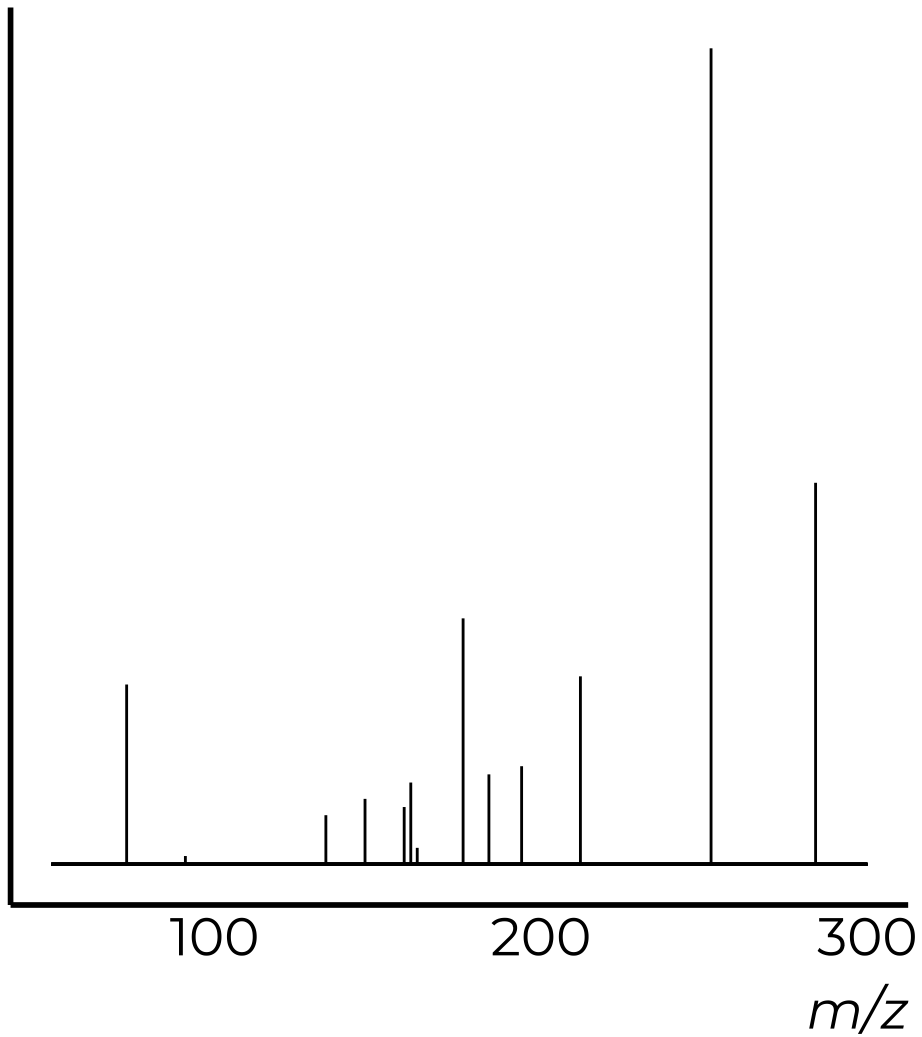
# toxicity assessment



LC<sub>50</sub> = 9.3 mg/L

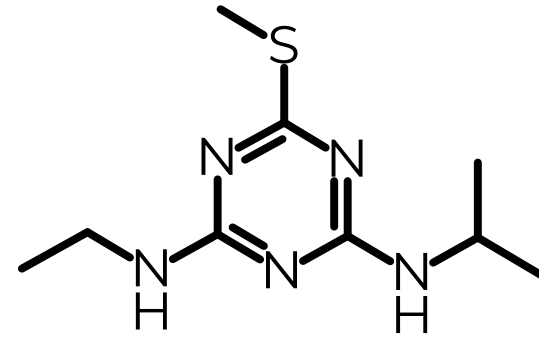
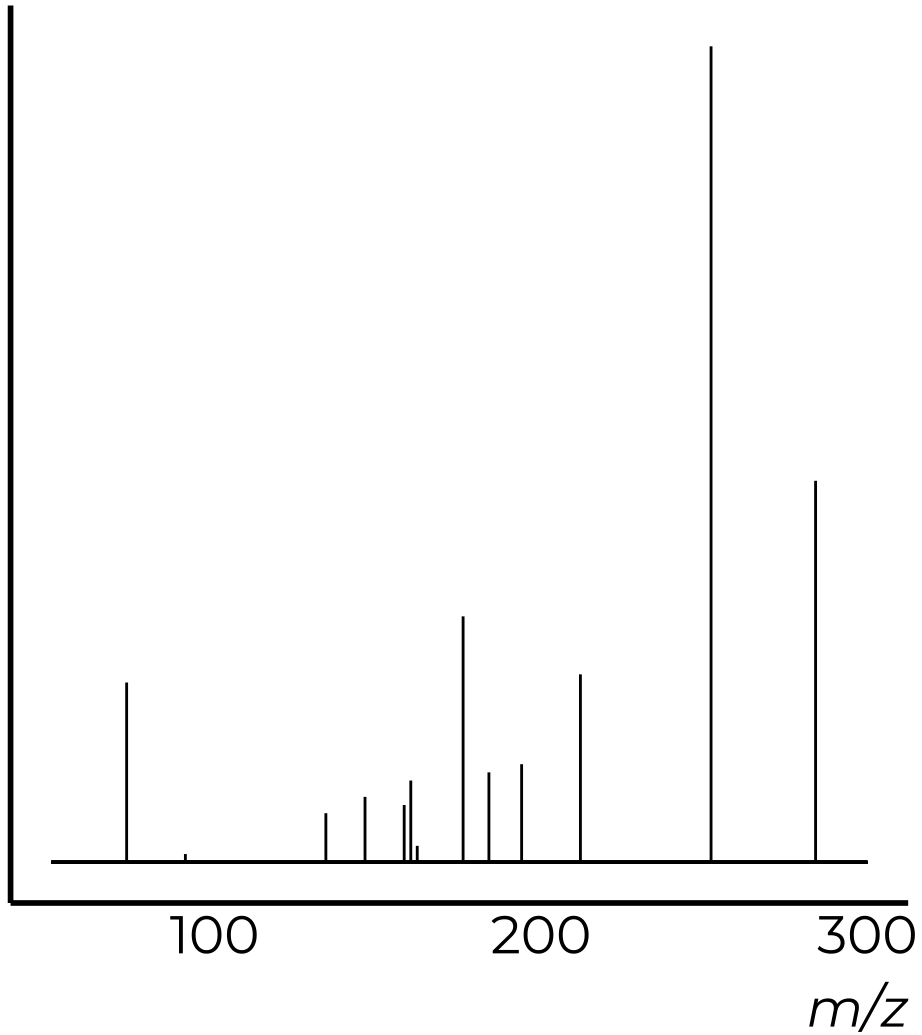


# toxicity assessment

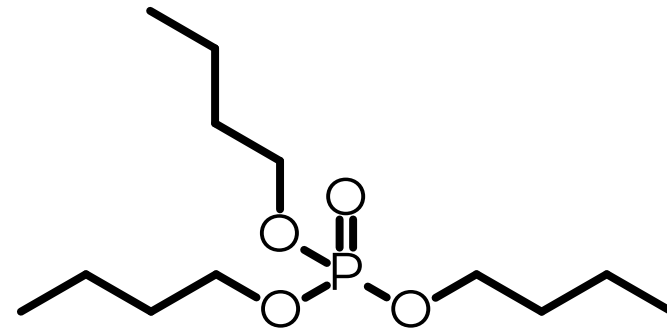


LC<sub>50</sub> = 9.3 mg/L

# toxicity assessment



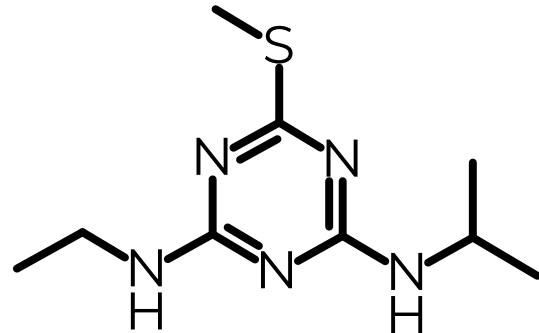
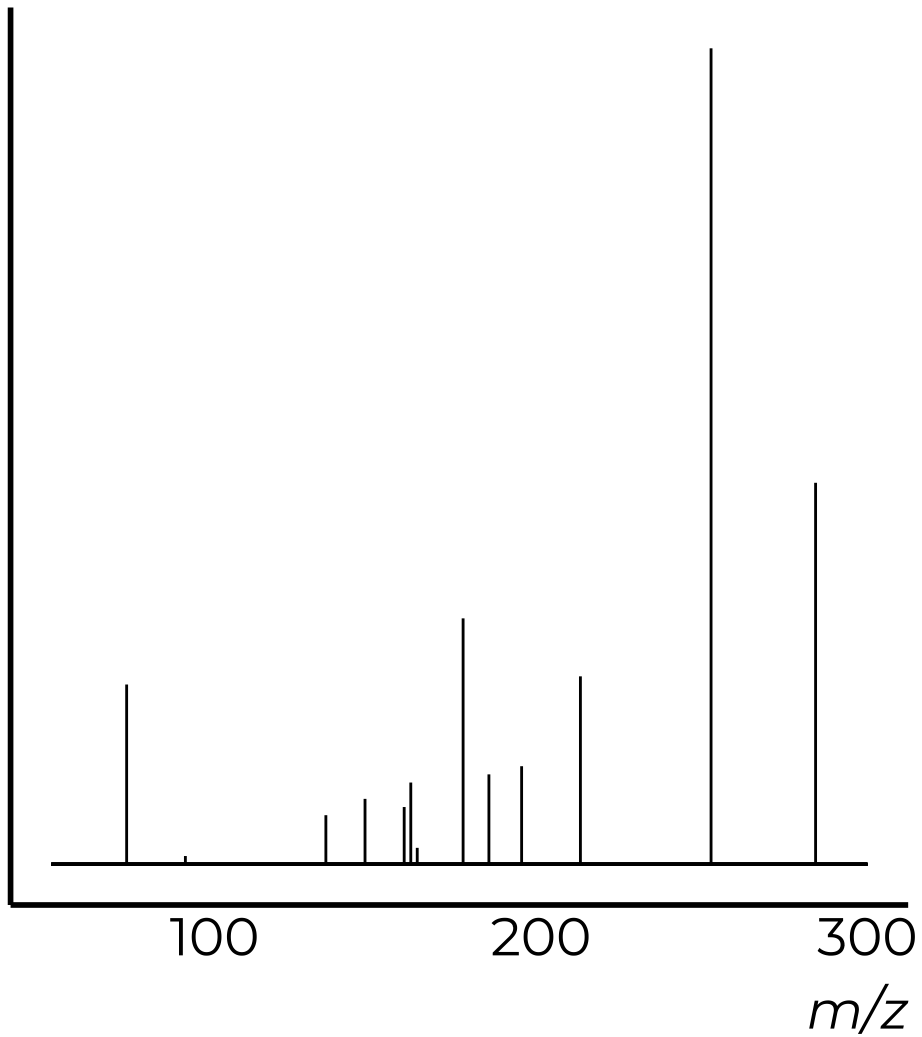
$LC_{50} = 9.3 \text{ mg/L}$



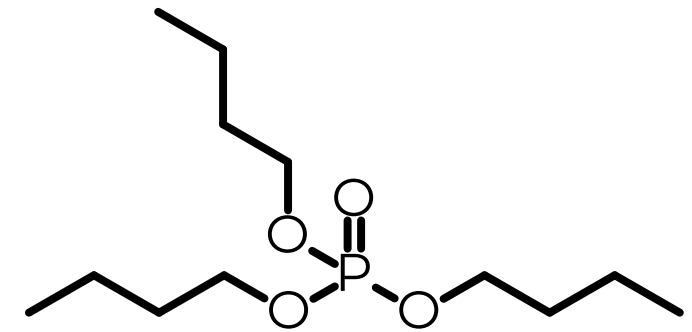
$LC_{50} = ? \text{ mg/L}$



# toxicity assessment



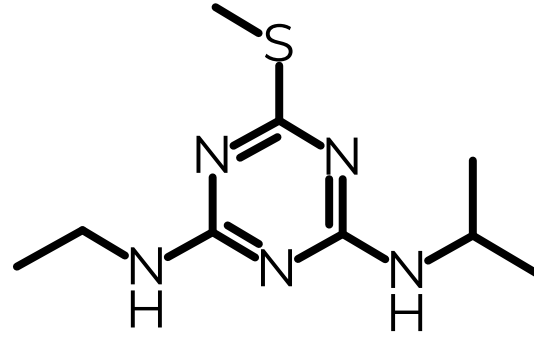
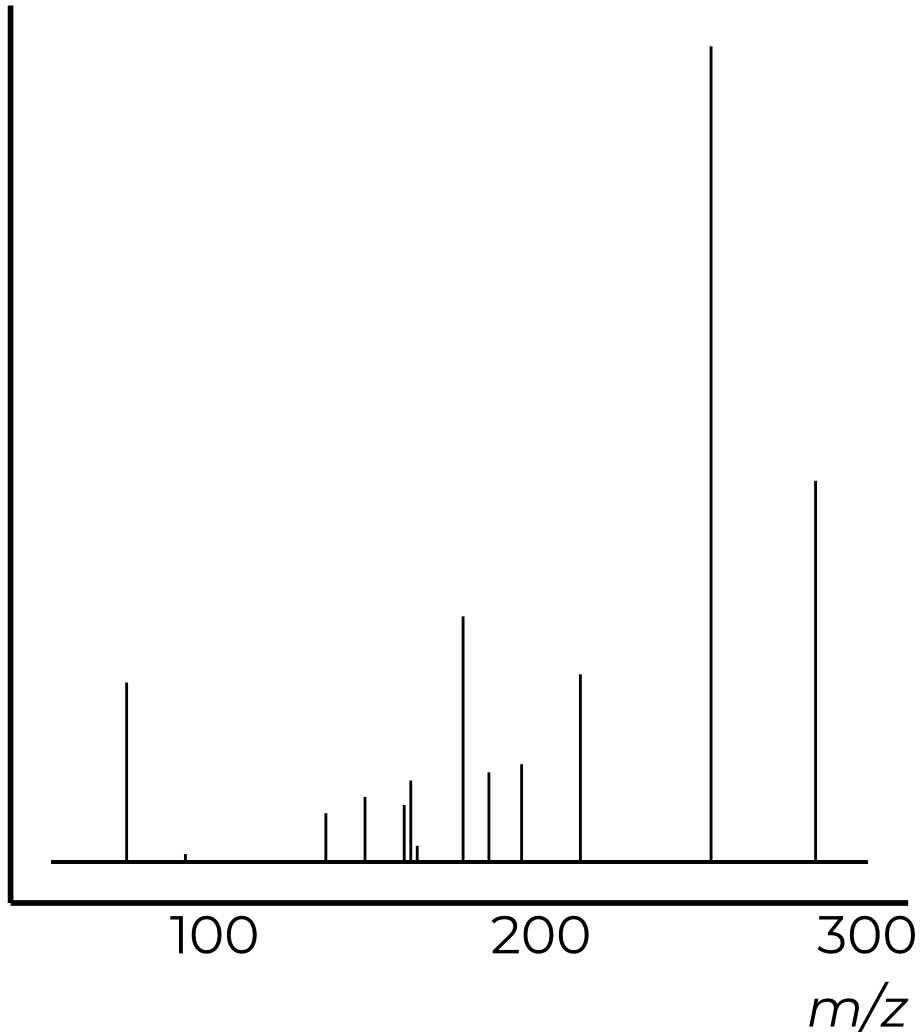
LC<sub>50</sub> = 9.3 mg/L



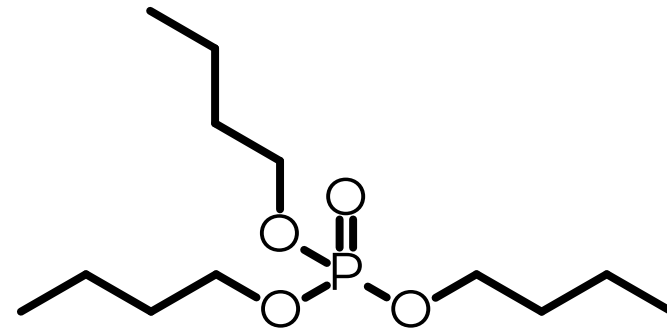
LC<sub>50</sub> = ? mg/L

?

# toxicity assessment



$LC_{50} = 9.3 \text{ mg/L}$

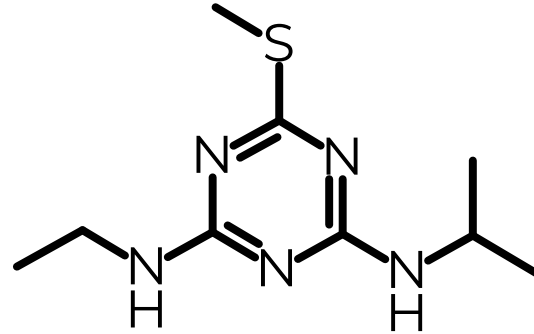


$LC_{50} = ? \text{ mg/L}$

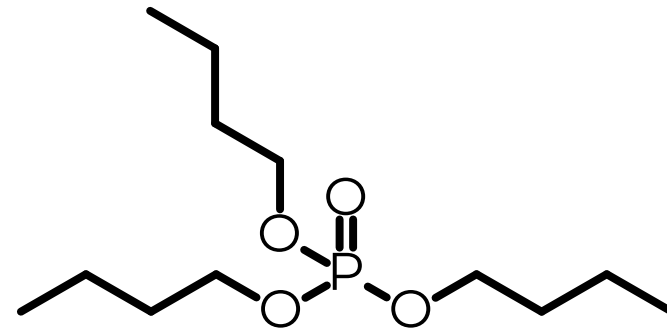
?

$LC_{50} = ? \text{ mg/L}$

# toxicity assessment



LC<sub>50</sub> = 9.3 mg/L



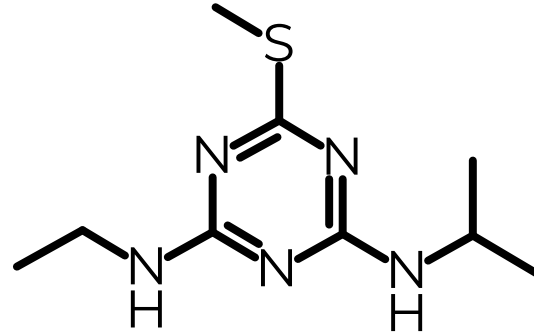
LC<sub>50</sub> = ? mg/L

?

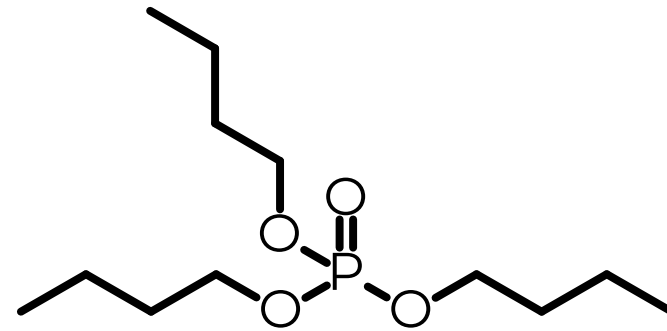
LC<sub>50</sub> = ? mg/L

# toxicity assessment

<1%



LC<sub>50</sub> = 9.3 mg/L



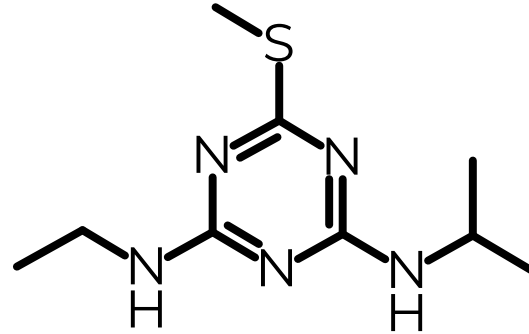
LC<sub>50</sub> = ? mg/L

?

LC<sub>50</sub> = ? mg/L

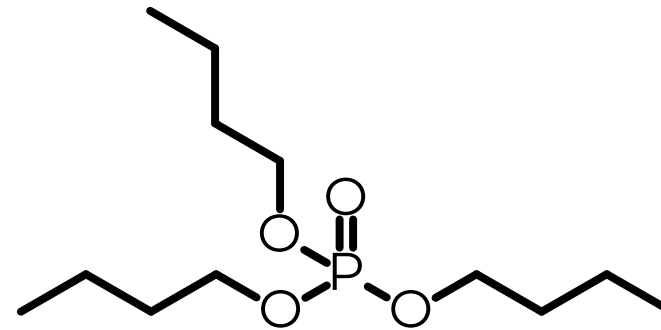
# toxicity assessment

<1%



LC<sub>50</sub> = 9.3 mg/L

<2%



LC<sub>50</sub> = ? mg/L

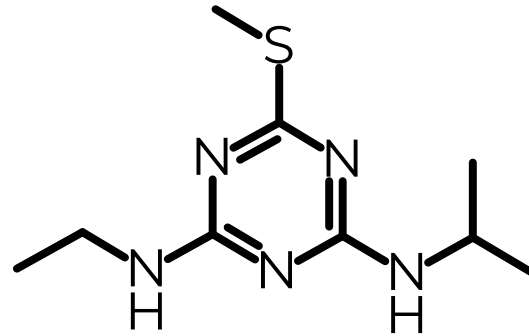
?

LC<sub>50</sub> = ? mg/L



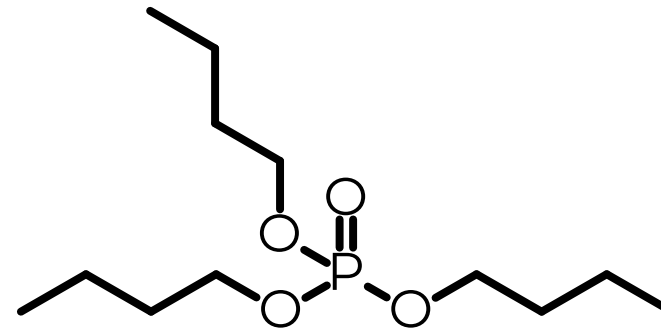
# toxicity assessment

<1%



LC<sub>50</sub> = 9.3 mg/L

<2%



LC<sub>50</sub> = ? mg/L

~98%

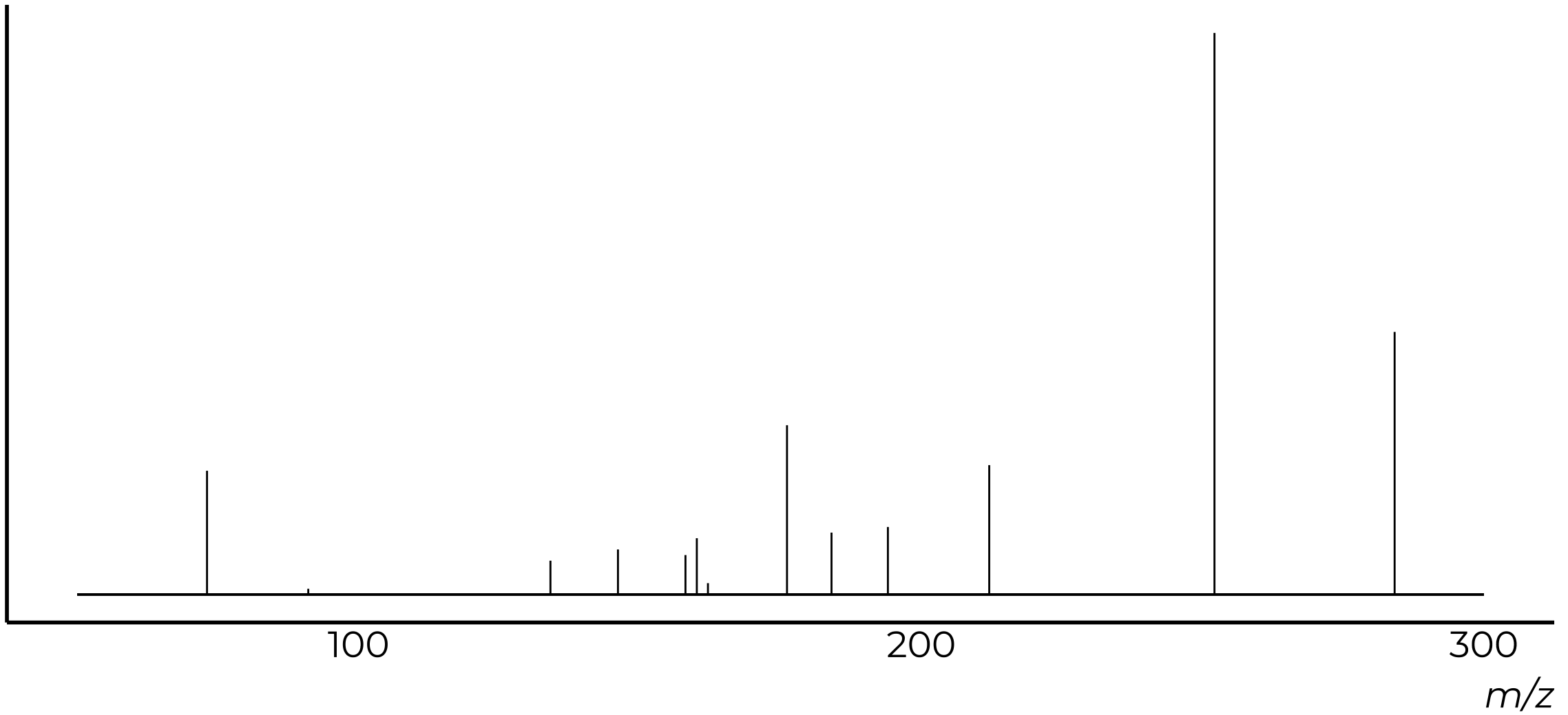
?

LC<sub>50</sub> = ? mg/L

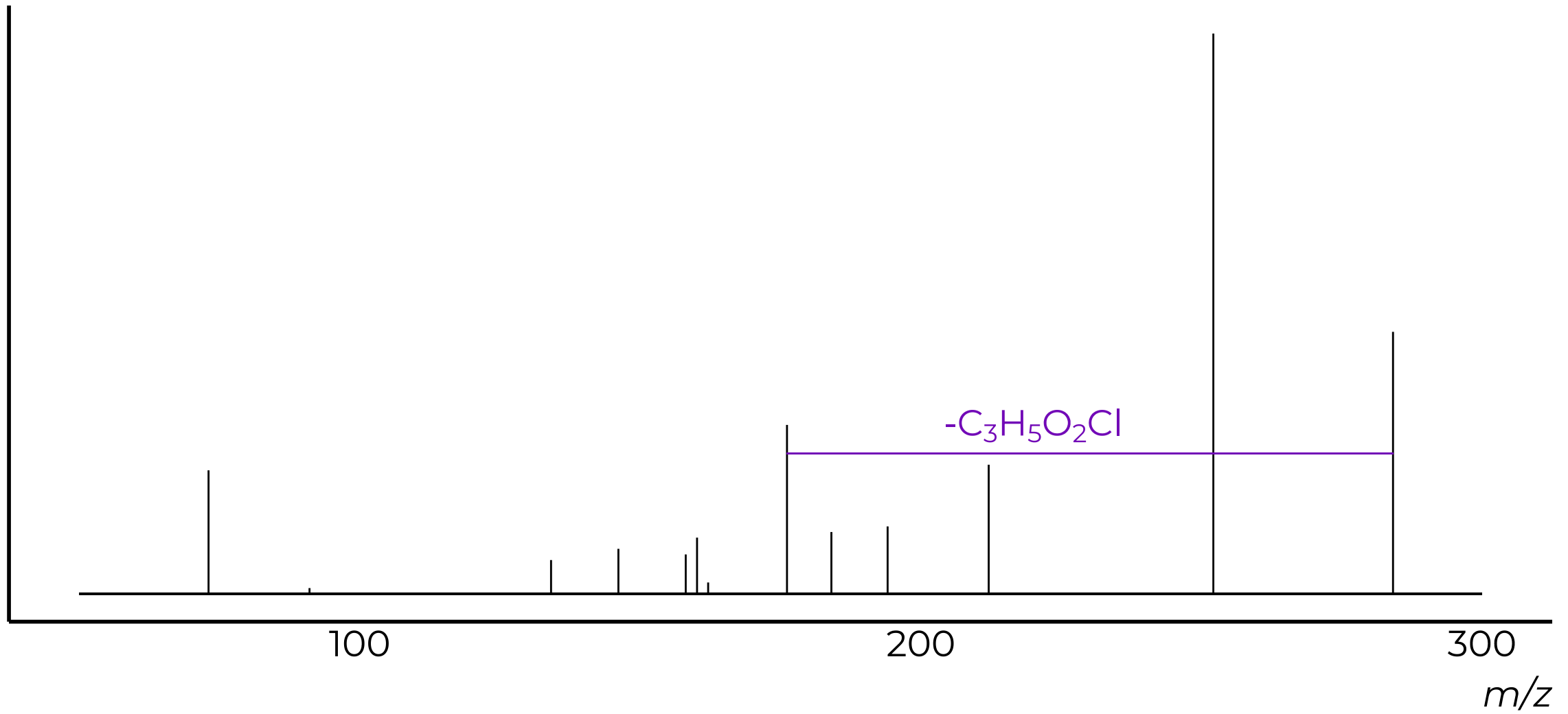
information available

in MS<sup>2</sup> spectra

# MS<sup>2</sup> spectra



# MS<sup>2</sup> spectra



modelling toxicity



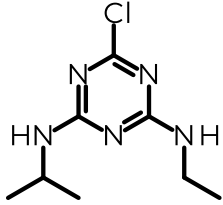
selected endpoint

# selected endpoint



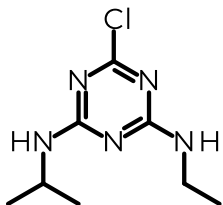
fathead minnow, bluegill and rainbow trout

# workflow



structure as SMILES

# workflow

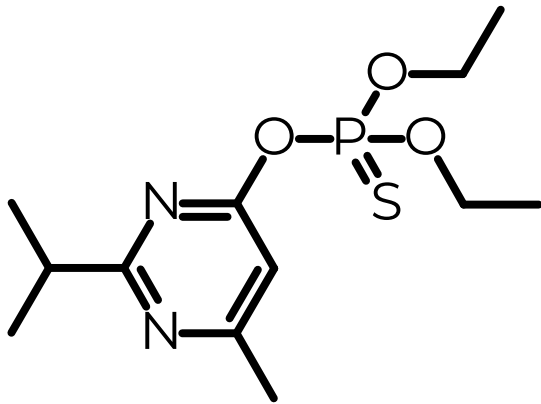


structure as SMILES



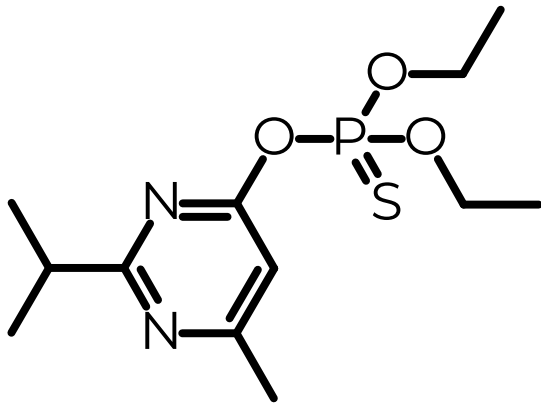
molecular descriptors

# structural fingerprints



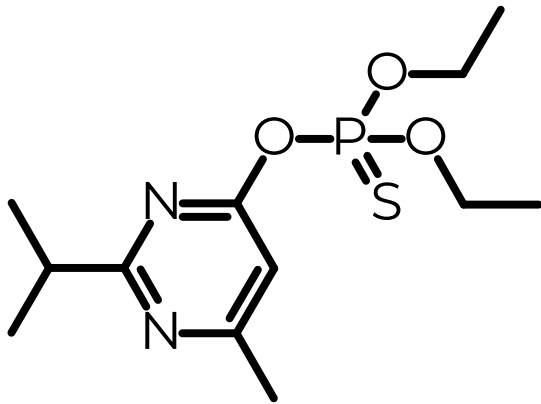


# structural fingerprints



R: rcdk  
→

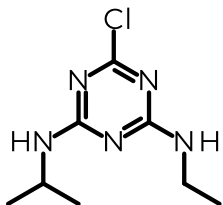
# structural fingerprints



R: rcdk  
→

0	
1	
1	
0	
1	

# workflow



structure as SMILES



molecular descriptors



machine learning for predicting  $LC_{50}$

# 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

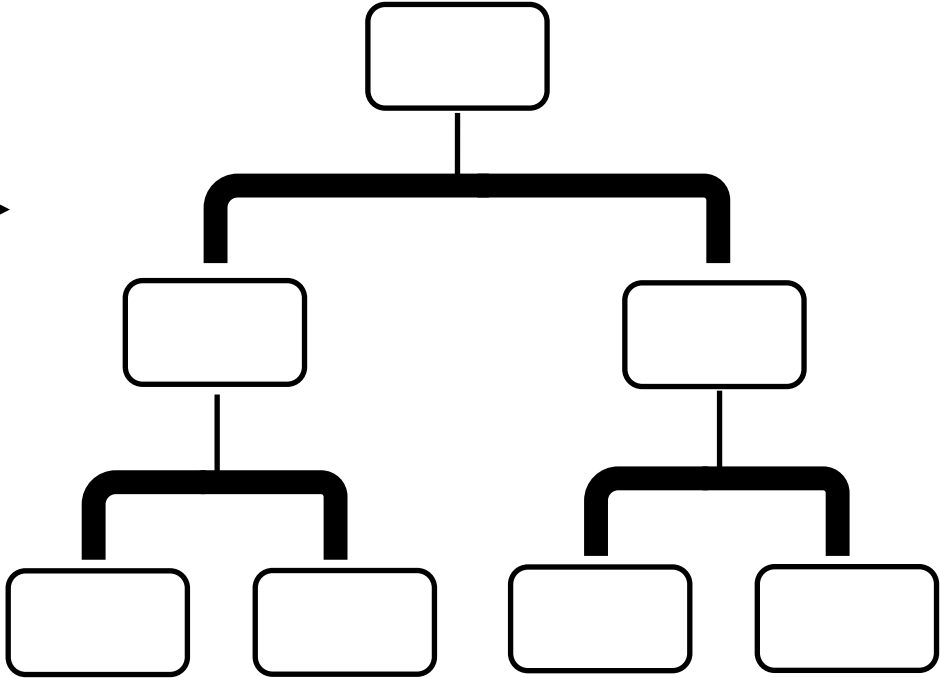
# model training

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

test set  
130  
chemicals

gradient  
boosting



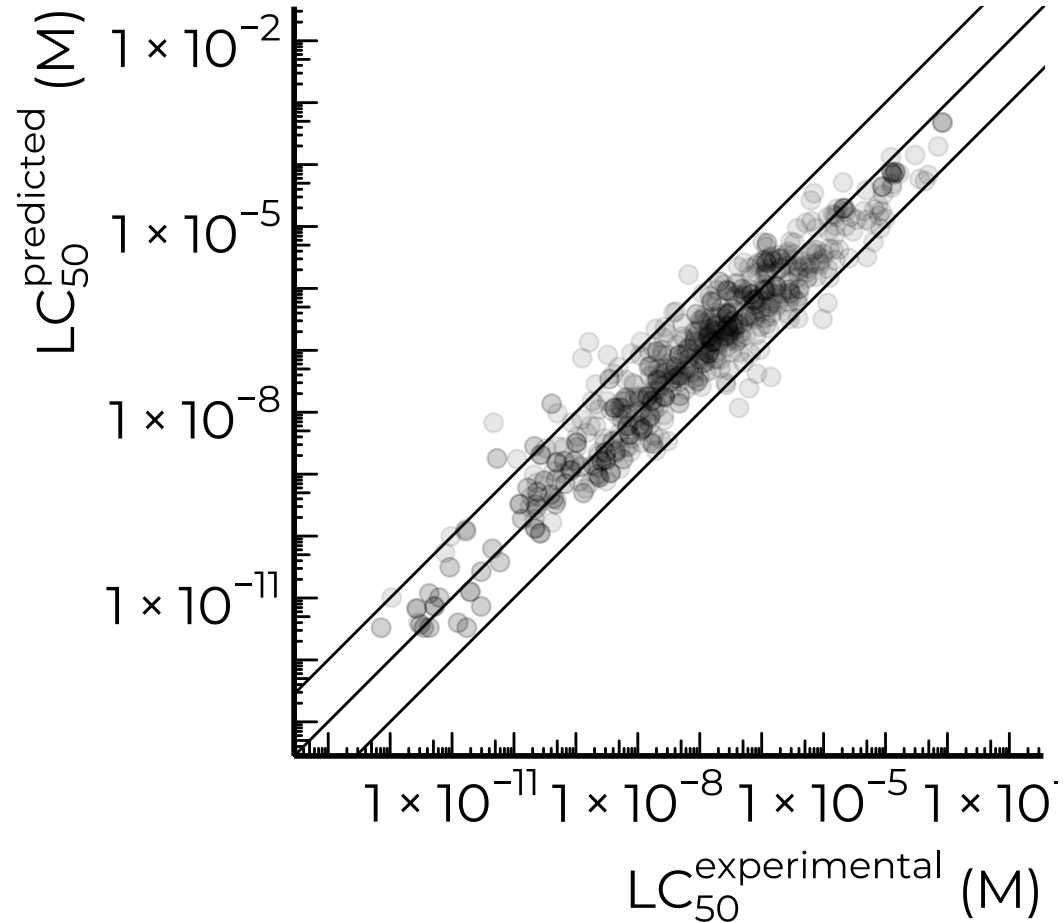
# performance

of LC<sub>50</sub> predictions with molecular fingerprints

# LC<sub>50</sub> predictions

Peets et al. ES&T 2022

fish LC<sub>50</sub>



training set

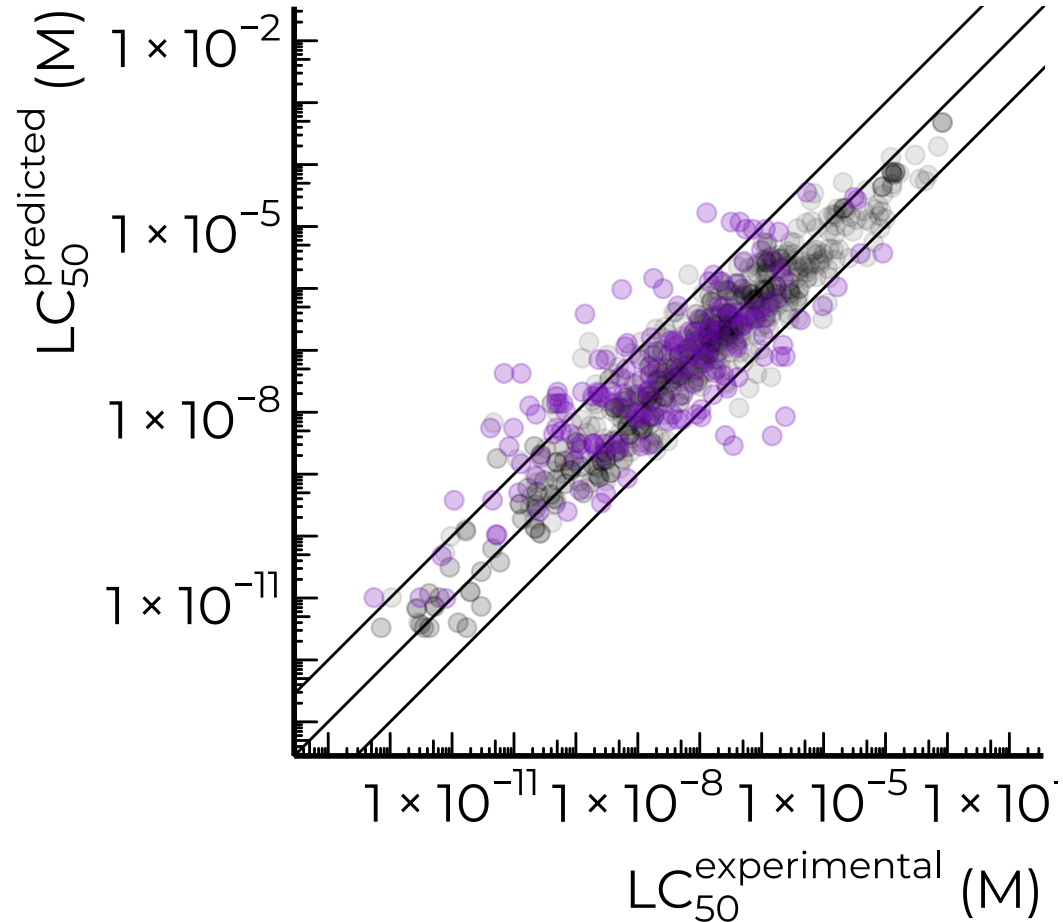
RMSE 0.52 log(M)



# LC<sub>50</sub> predictions

Peets et al. ES&T 2022

fish LC<sub>50</sub>



training set

RMSE 0.52 log(M)

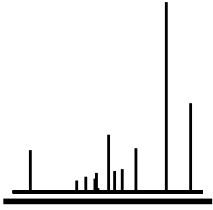
test set

RMSE 0.78 log(M)

unidentified chemicals

from MS<sup>2</sup> spectra

# workflow



MS<sup>2</sup> spectra

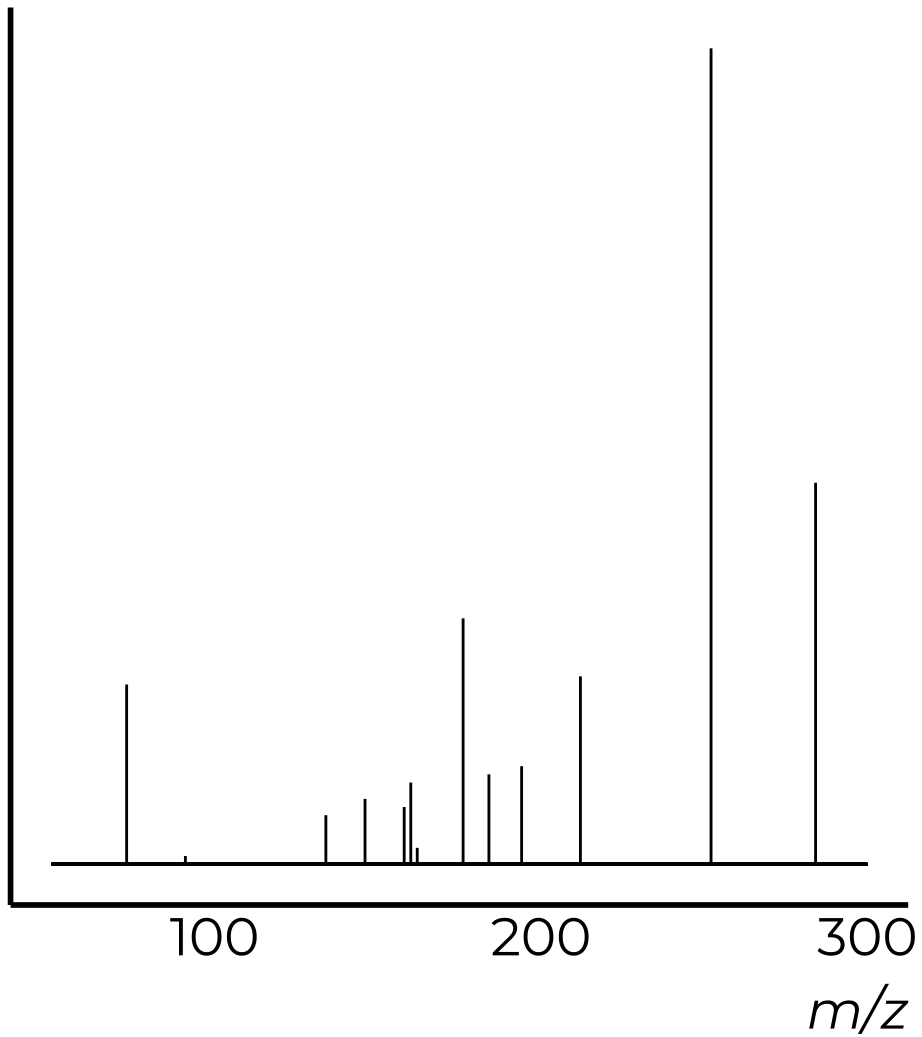


molecular fingerprints with SIRIUS+CSI:FingerID



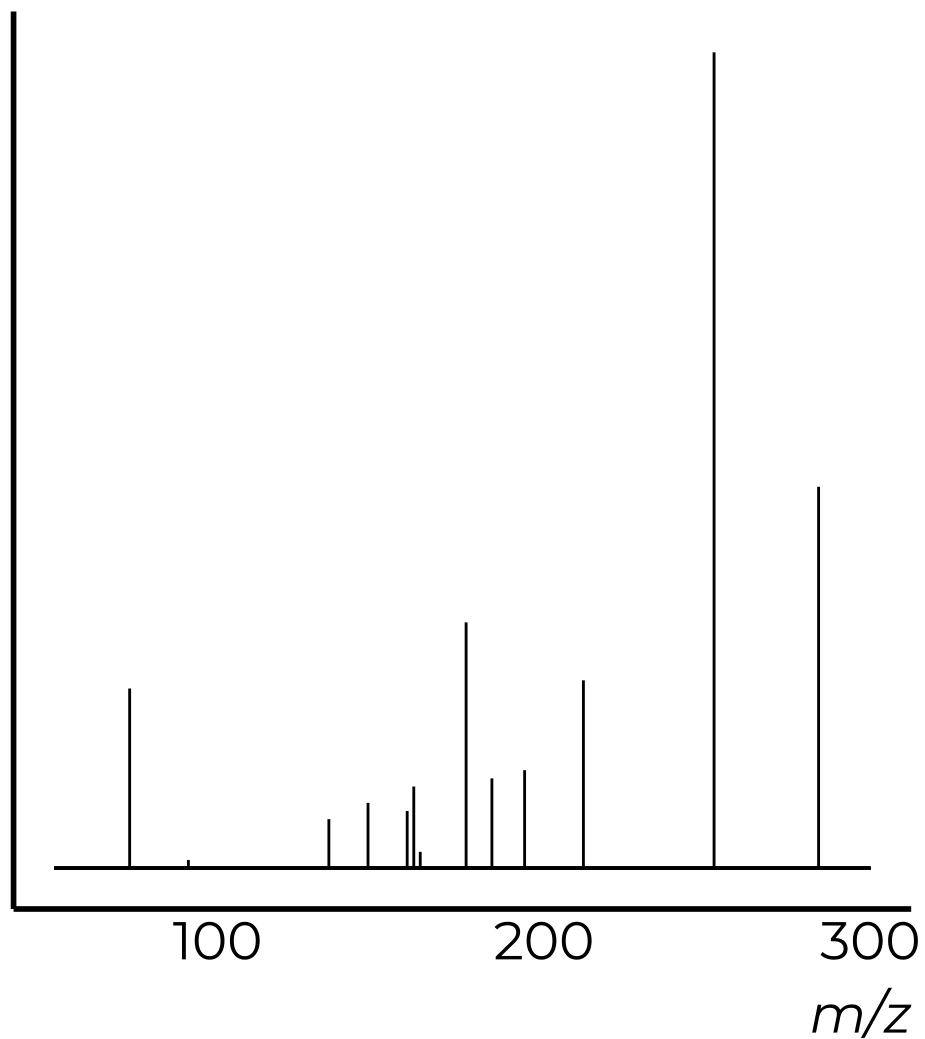
predict LC<sub>50</sub>

# predict for unknown chemicals

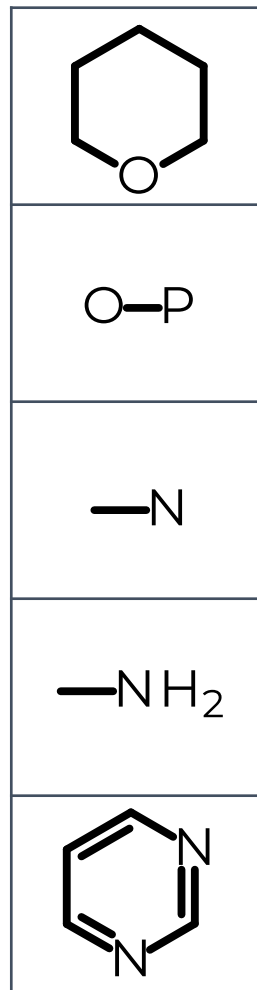


?

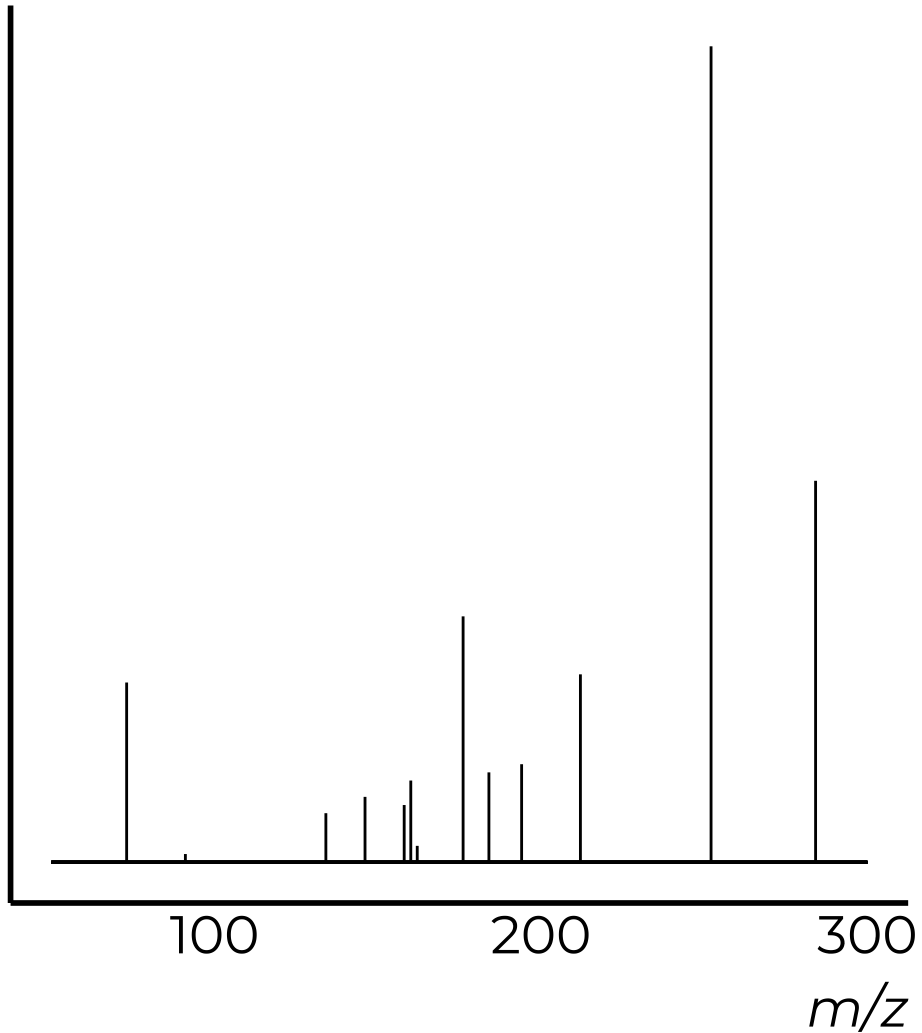
# predict for unknown chemicals



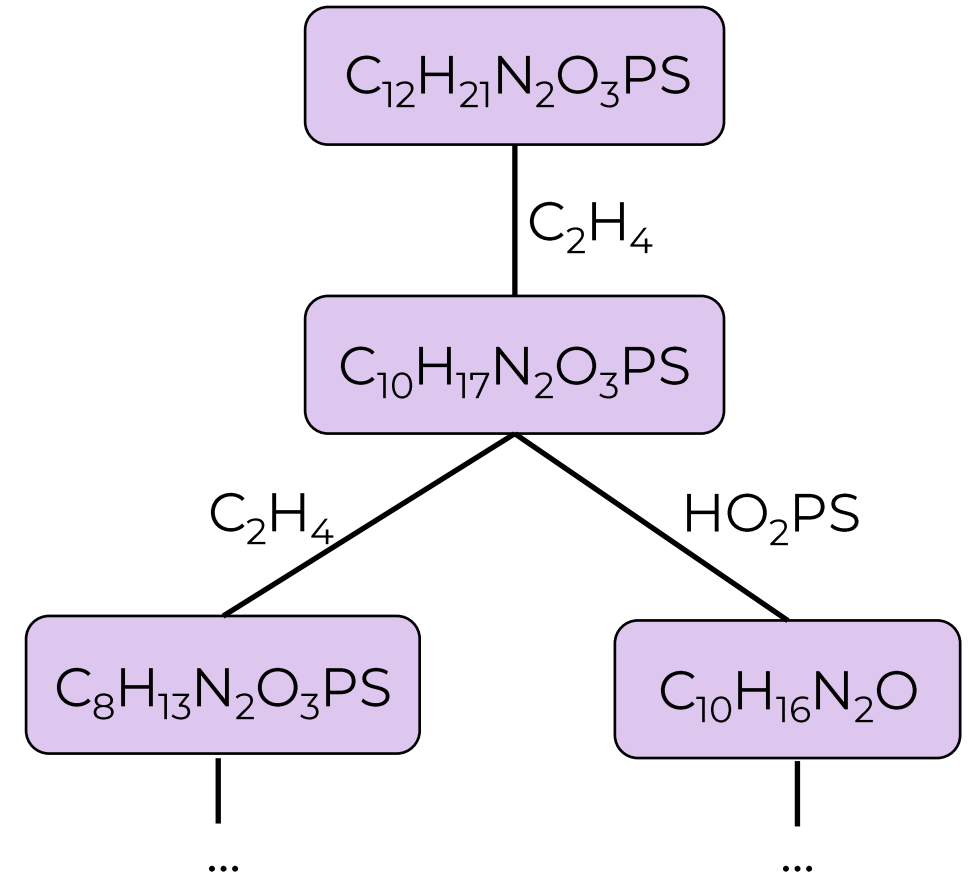
SIRIUS+  
CSI:FingerID  
→



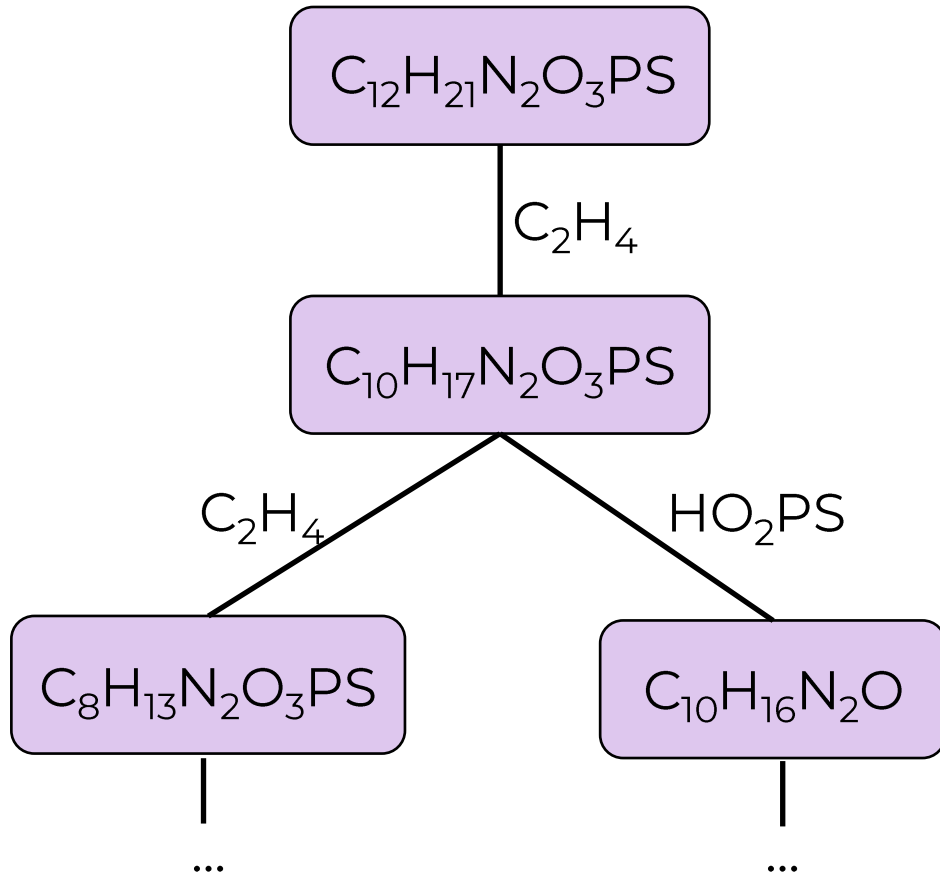
# predict for unknown chemicals



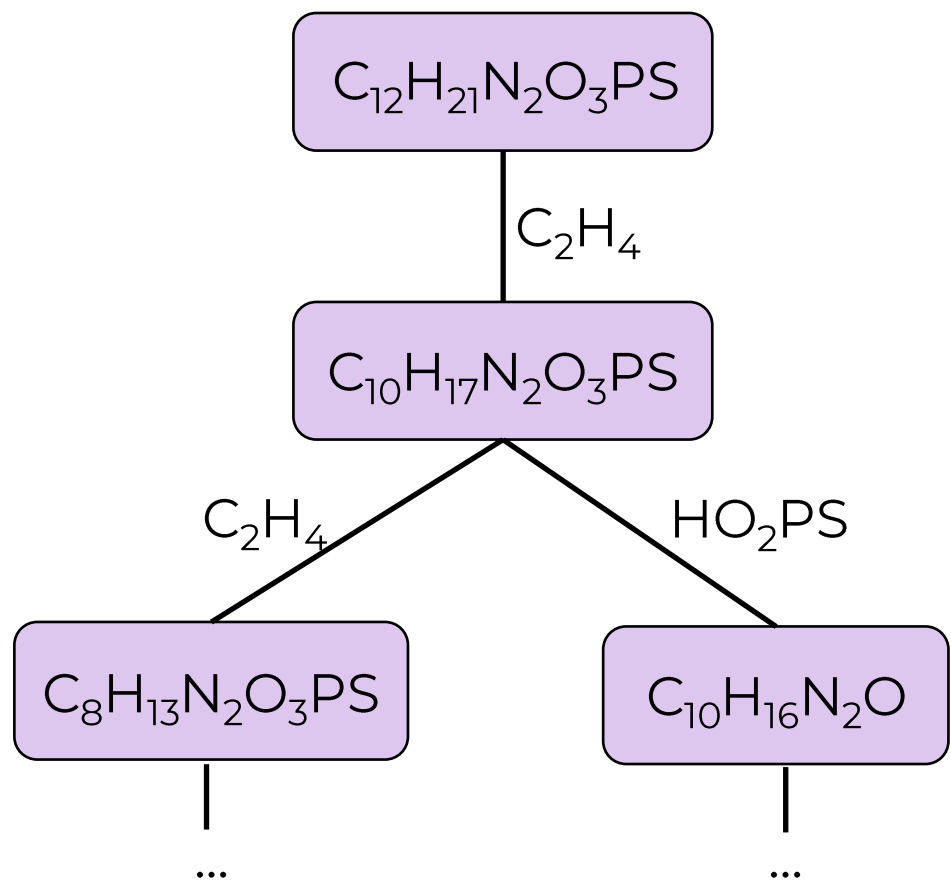
SIRIUS+  
CSI:FingerID  
→



# predict for unknown chemicals



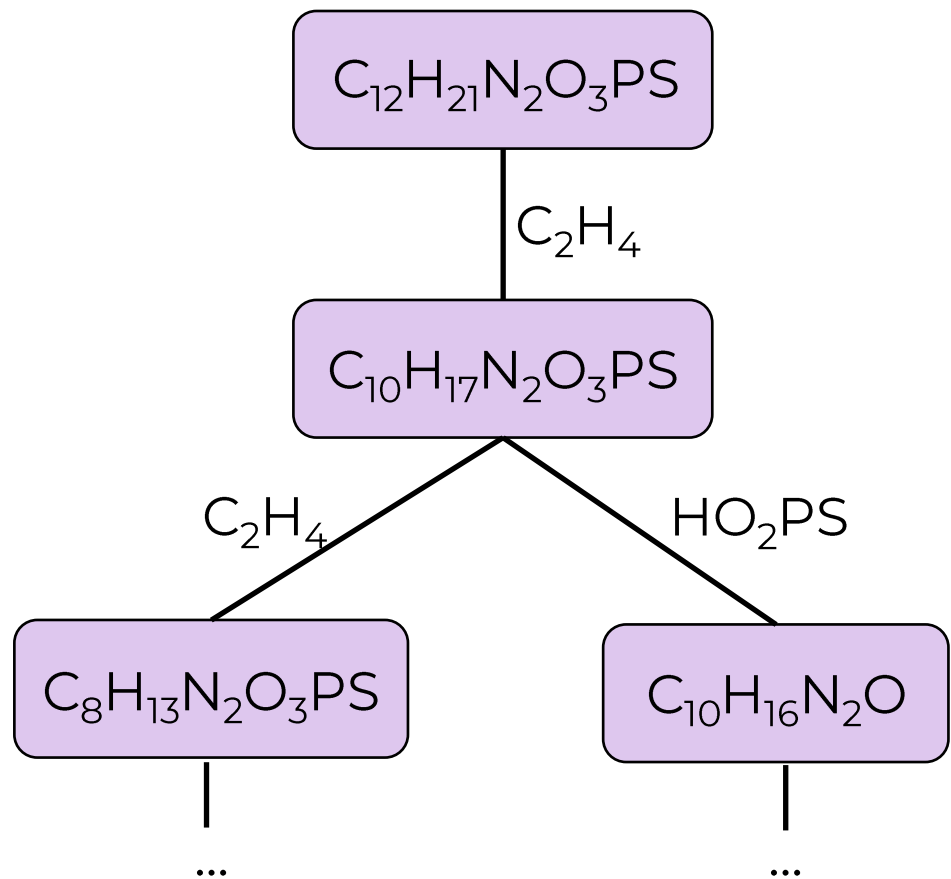
# predict for unknown chemicals



0.001	
0.999	$O-P$
0.999	$-N$
0.198	$-NH_2$
0.988	

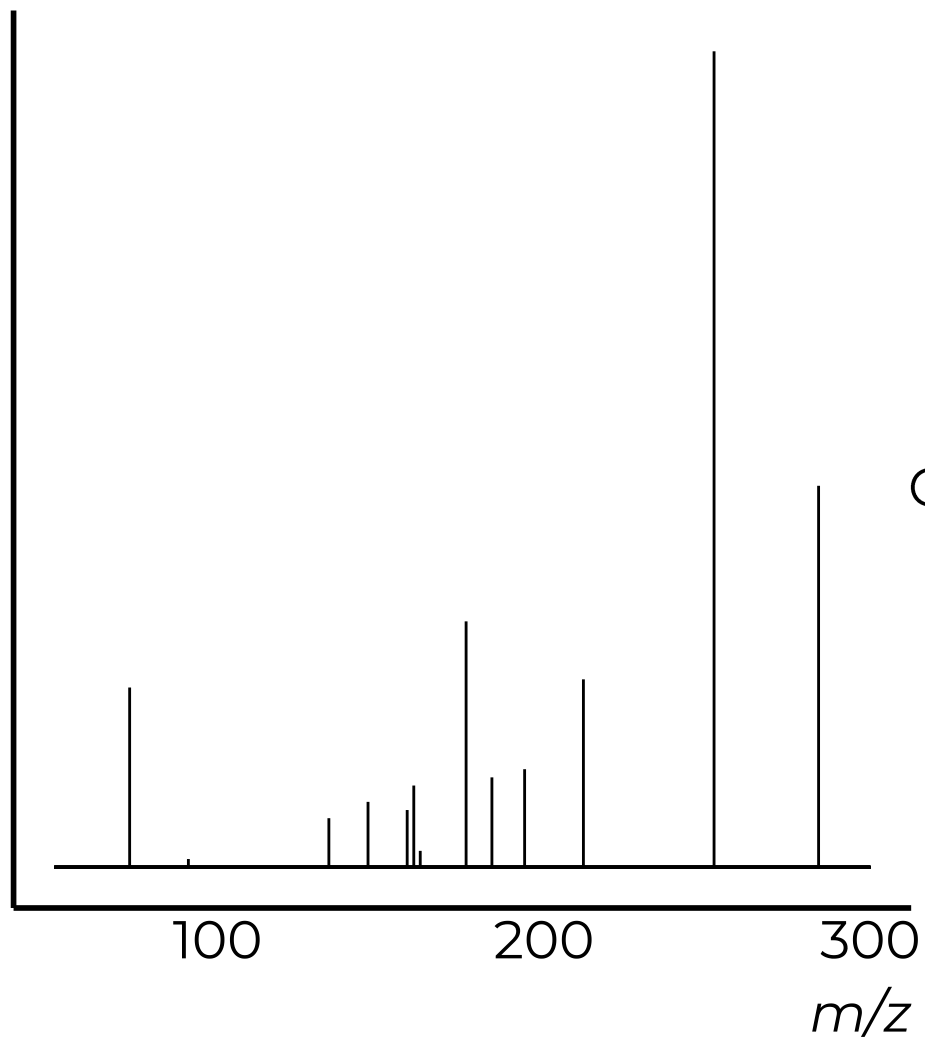


# predict for unknown chemicals

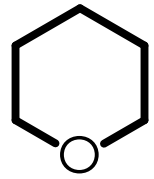



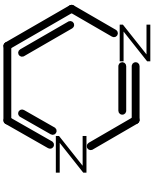


0	
1	
1	
0	
1	

# predict for unknown chemicals

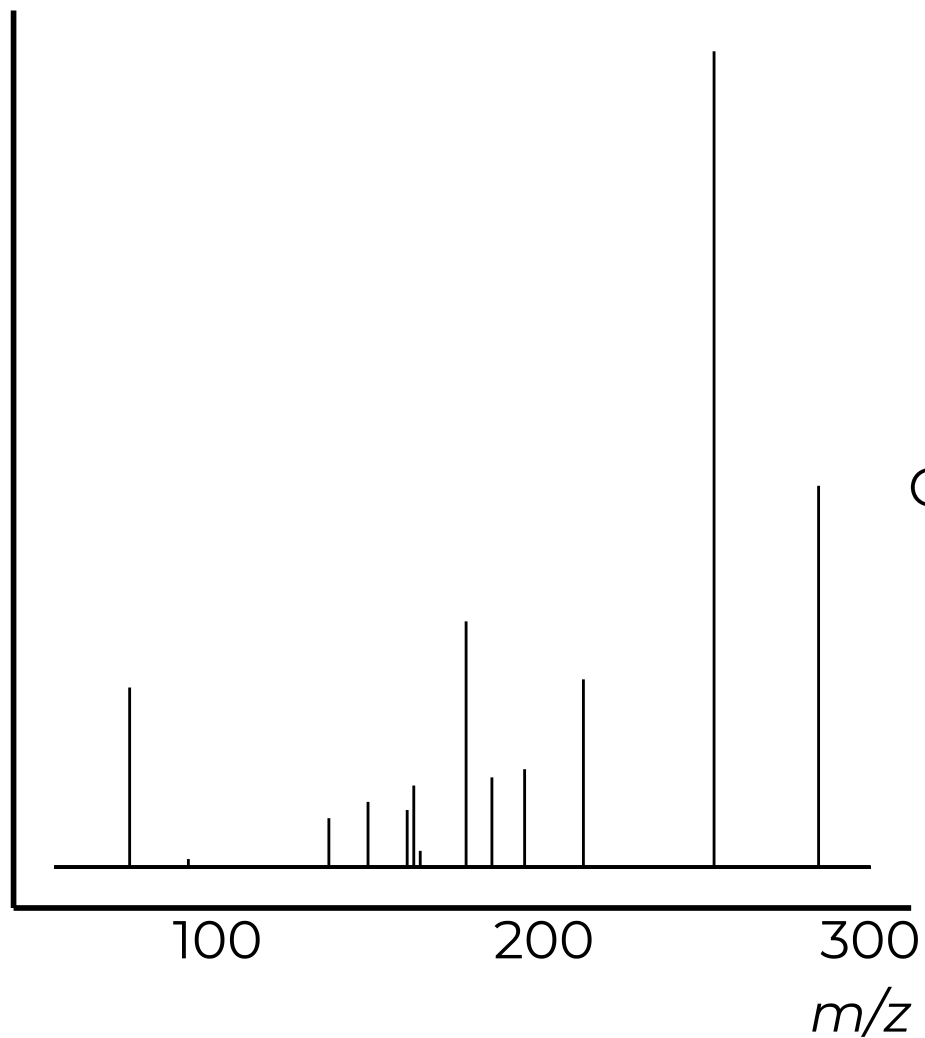


SIRIUS+  
CSI:FingerID  
→

0	
1	
1	
0	
1	

gradient  
→  
boosting

# predict for unknown chemicals



SIRIUS+  
CSI:FingerID  
→

0	
1	
1	
0	
1	

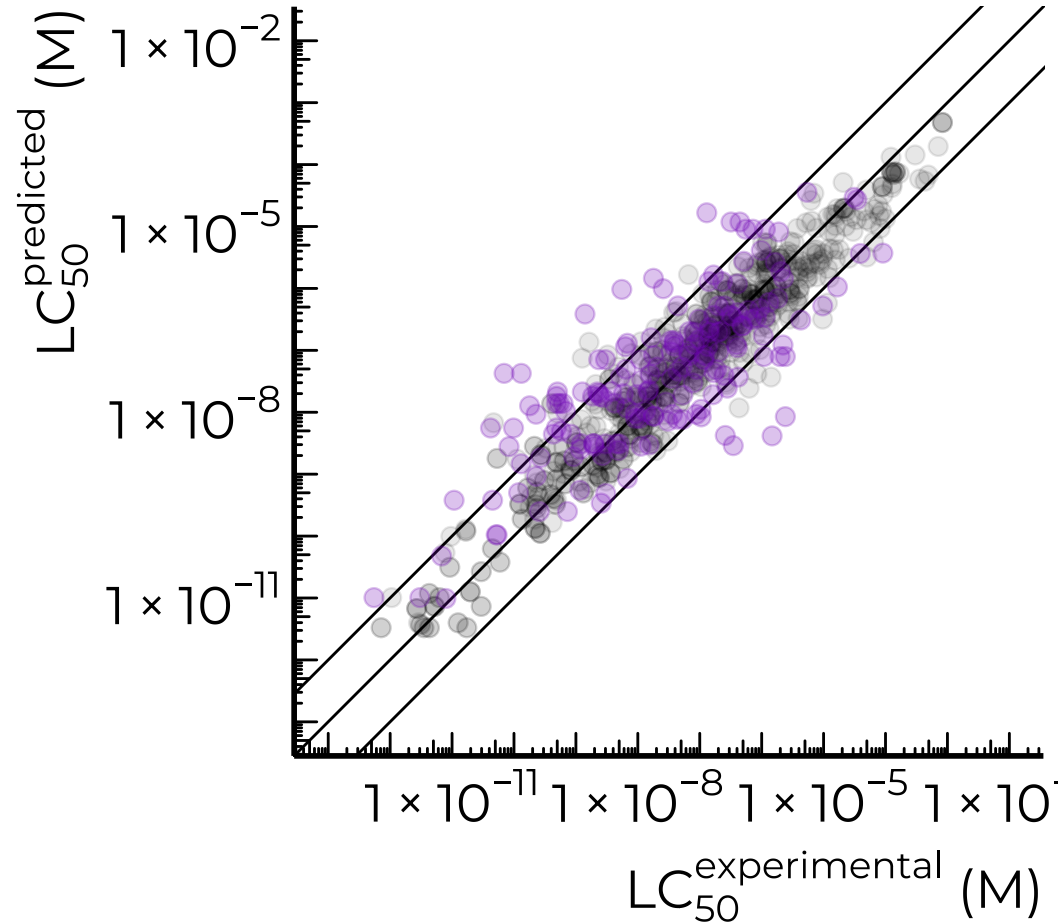
gradient  
boosting →  $LC_{50} = -2.2 \log(\text{mM})$

# LC<sub>50</sub> predictions

# LC<sub>50</sub> predictions

Peets et al. ES&T 2022

fish LC<sub>50</sub>



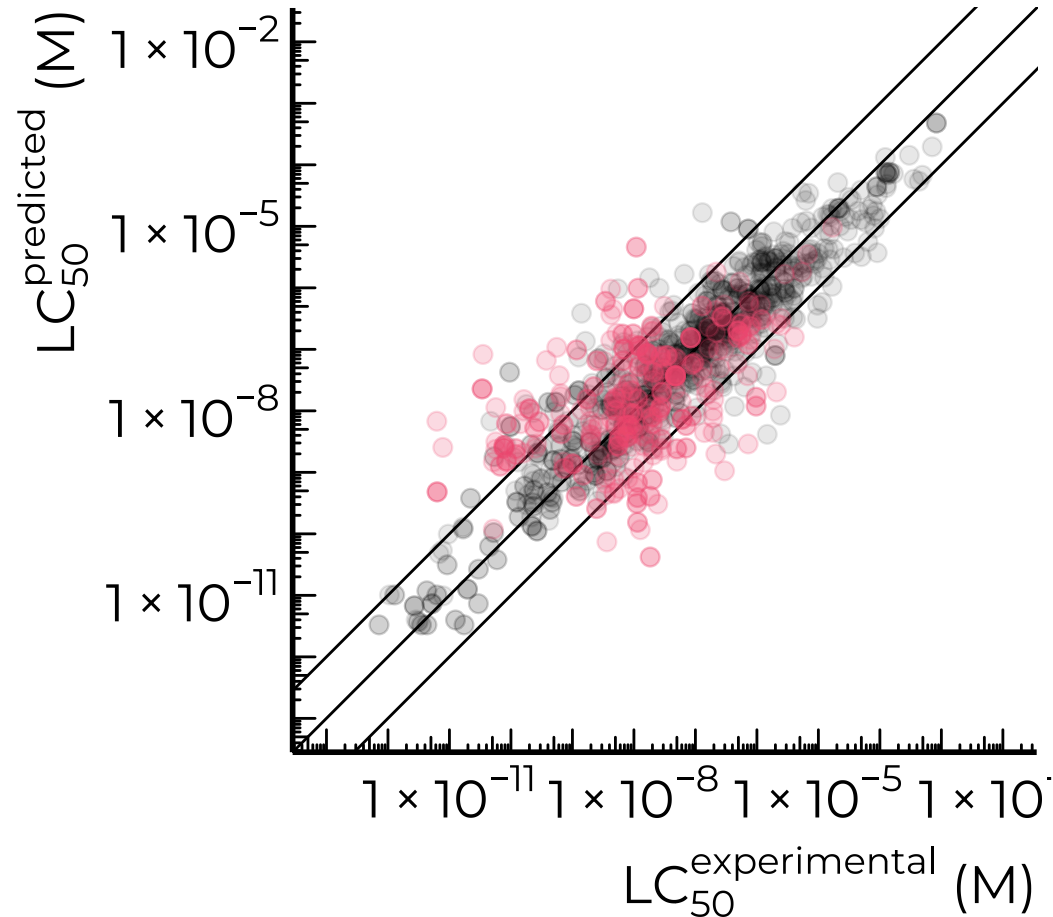
test set on structures

RMSE 0.78 log(M)

# LC<sub>50</sub> predictions

Peets et al. ES&T 2022

fish LC<sub>50</sub>



test set on structures

RMSE 0.78 log(M)

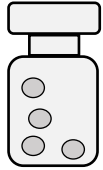
validation on MassBank

RMSE<sub>model</sub> 0.88 log(M)

SD<sub>experimental</sub> 0.44 log(mM)

pinpointing toxic chemicals

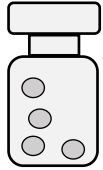
# case study on wastewater



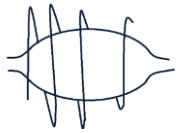
wastewater samples



# case study on wastewater

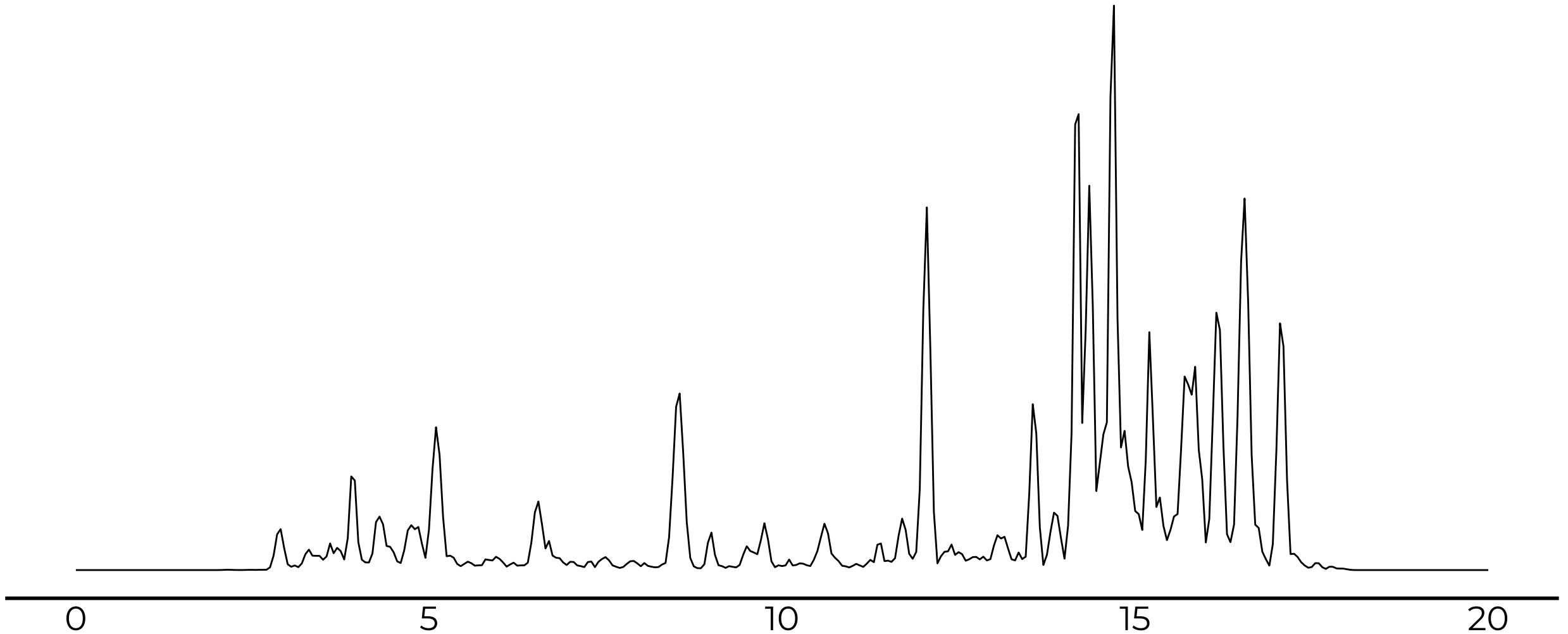


wastewater samples

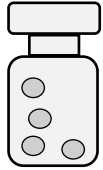


LC/HRMS analysis

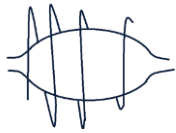
# case study on wastewater



# case study on wastewater



wastewater samples

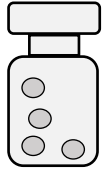


LC/HRMS analysis

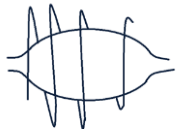


MS-DIAL peak picking

# case study on wastewater



wastewater samples



LC/HRMS analysis

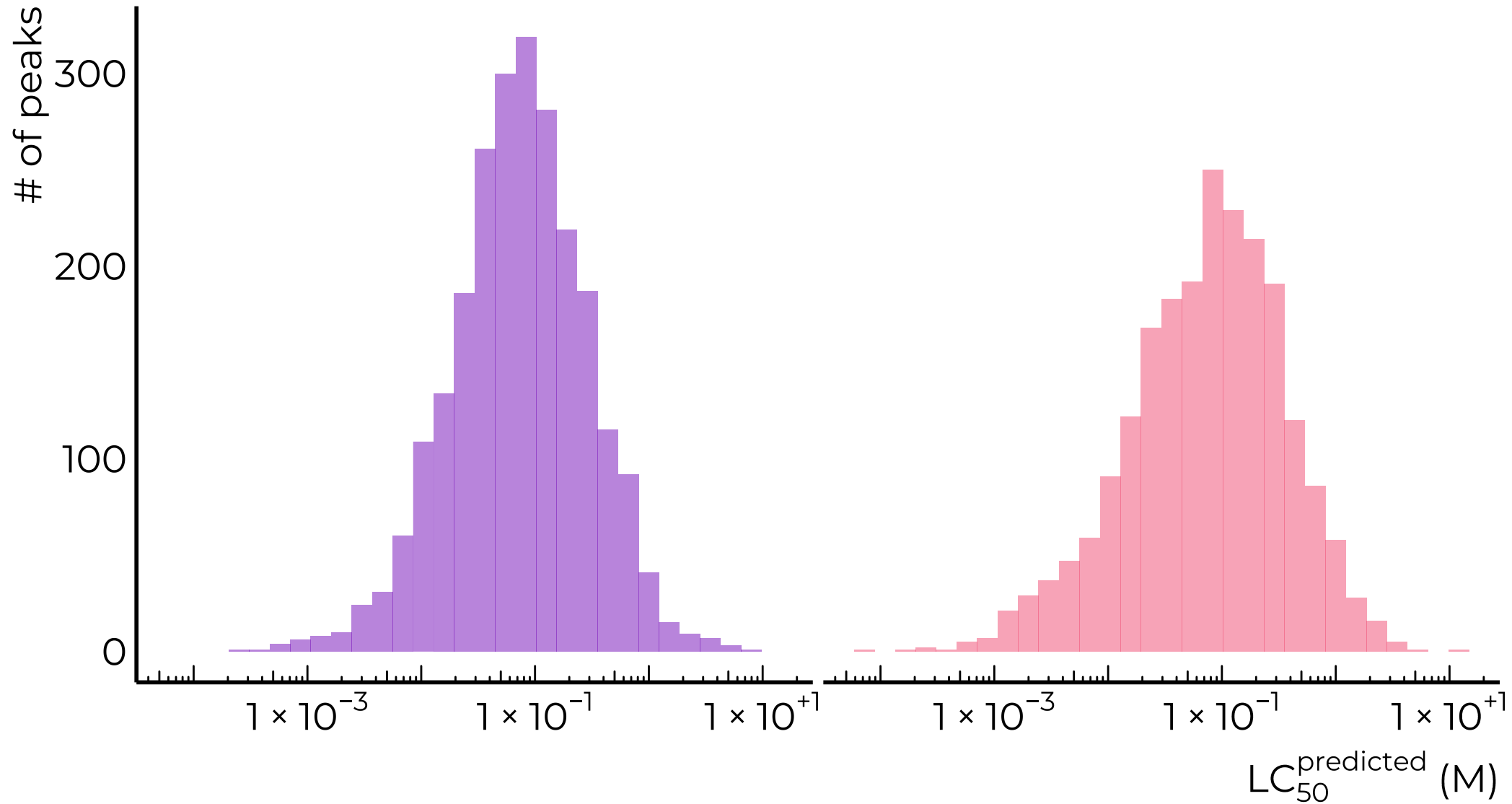


MS-DIAL peak picking

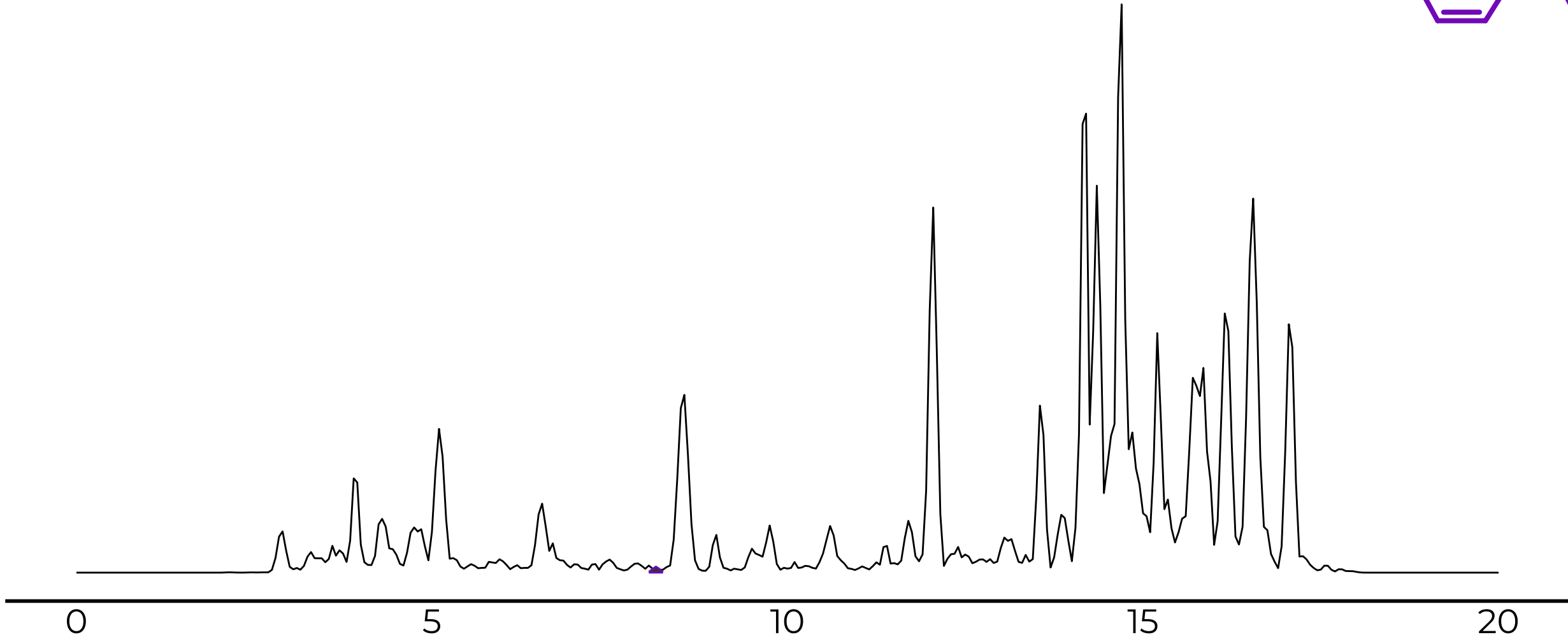
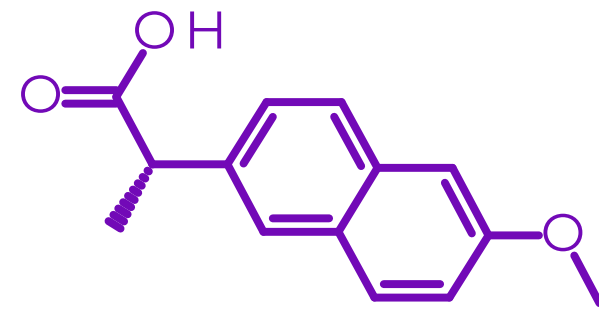


SIRIUS+CSI:FingerID

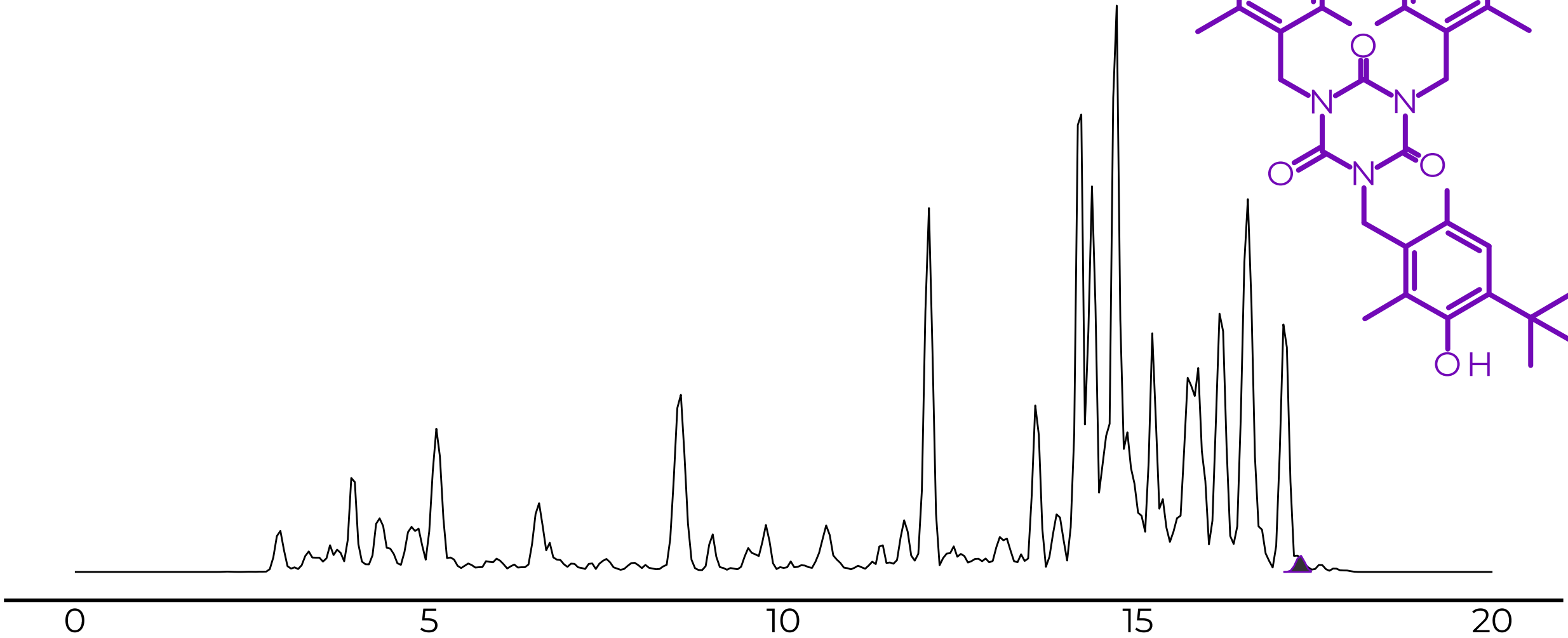
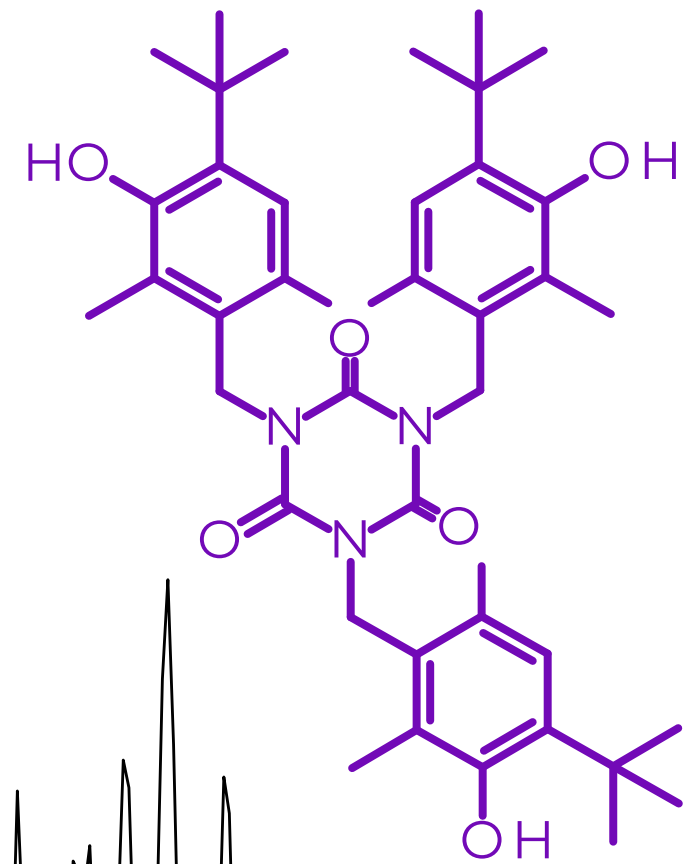
# LC<sub>50</sub> distribution



# naproxen



# cyanox CY 1790

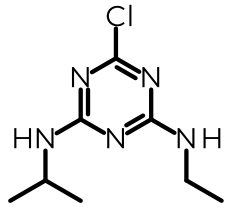


summary

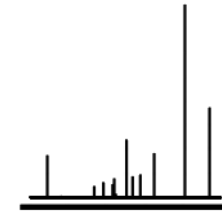


# prioritization in NTS

toxicity



structure



MS<sup>2</sup> spectrum



[kruvelab.com](http://kruvelab.com)

[anneli.kruve@su.se](mailto:anneli.kruve@su.se)