

# Machine learning tools can pinpoint high-risk water pollutants

Helen Sepman

[Helen.Sepman@aces.su.se](mailto:Helen.Sepman@aces.su.se)

*Kruve Lab*



Stockholm  
University



# Introduction





# Introduction

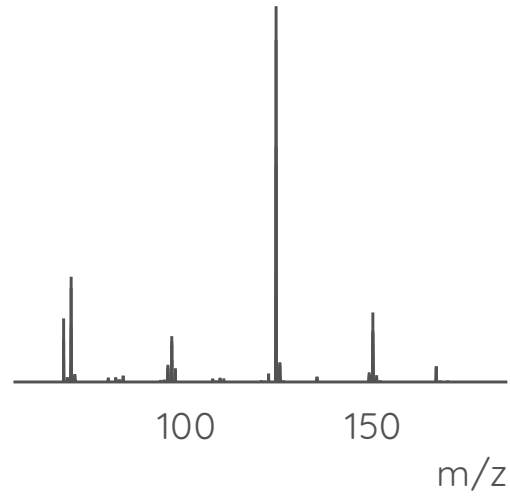
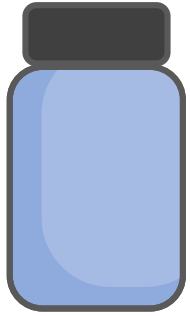


# Research questions

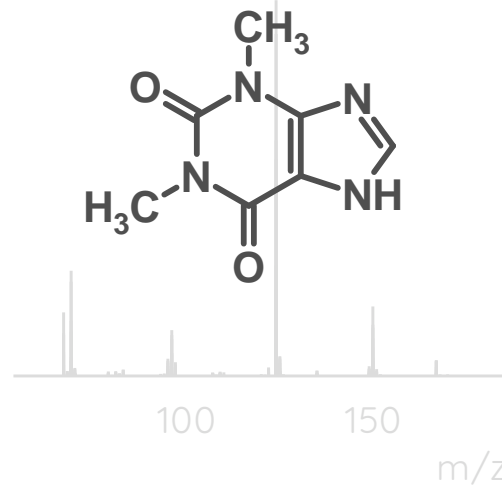
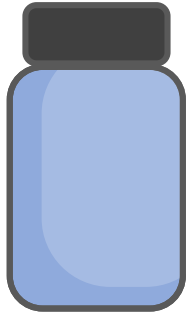
How to pinpoint high risk chemicals?

How to acquire high quality spectra for high risk chemicals?

# Assessing the risk of a detected chemical



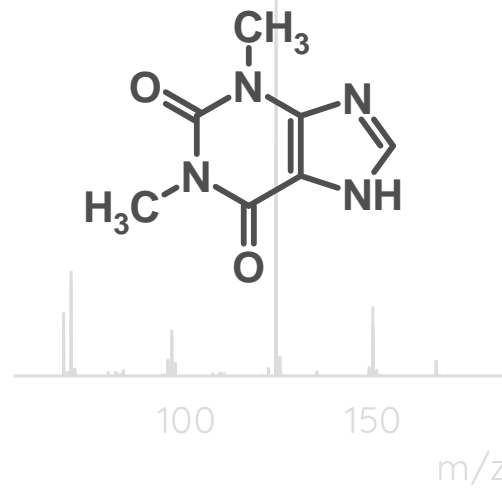
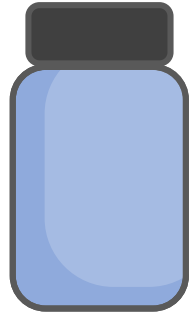
# Assessing the risk of a detected chemical



How much of this chemical is present in the environment?

How much of this chemical would start to cause adverse outcomes?

# Assessing the risk of a detected chemical



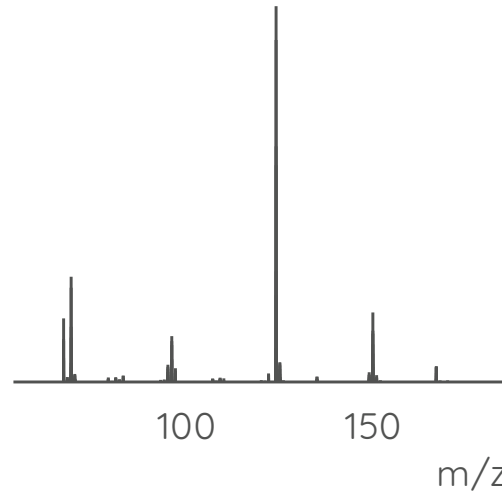
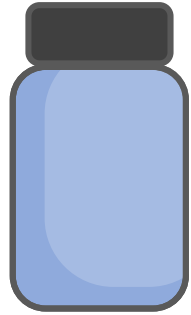
$C_{\text{environmental}} \gg C_{\text{toxic}}$

$C_{\text{environmental}} \ll C_{\text{toxic}}$

How much of this chemical is present in the environment?

How much of this chemical would start to cause adverse outcomes?

# Assessing the risk of a detected chemical



$C_{\text{environmental}} \gg C_{\text{toxic}}$

$C_{\text{environmental}} \ll C_{\text{toxic}}$

How much of this chemical is present in the environment?

How much of this chemical would start to cause adverse outcomes?

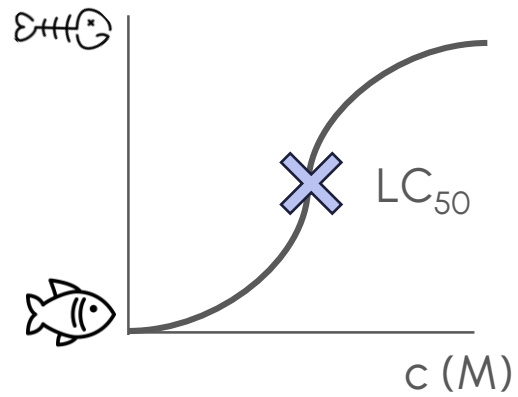


# Machine learning tools for risk estimation

# Machine learning tools for risk estimation



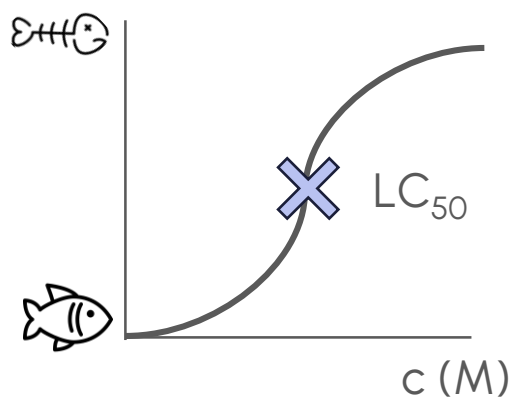
Lethal concentration, 50%  
( $LC_{50}$ )



# Machine learning tools for risk estimation

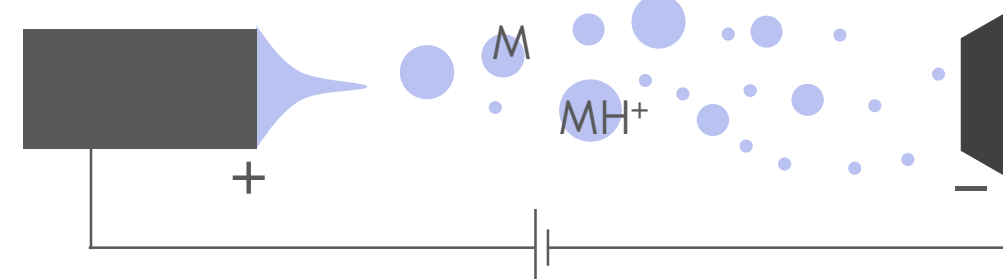


Lethal concentration, 50%  
( $LC_{50}$ )



Ionization efficiency  
(IE)

LC-ESI-MS

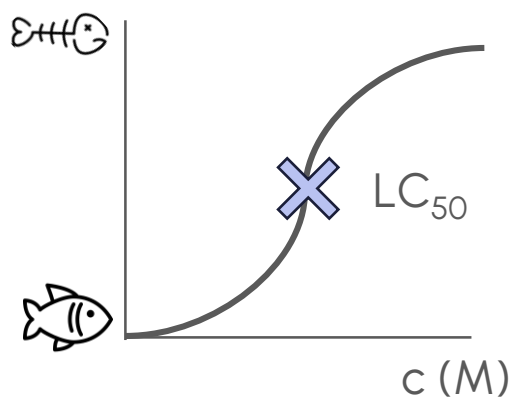




# Machine learning tools for risk estimation



Lethal concentration, 50%  
( $LC_{50}$ )

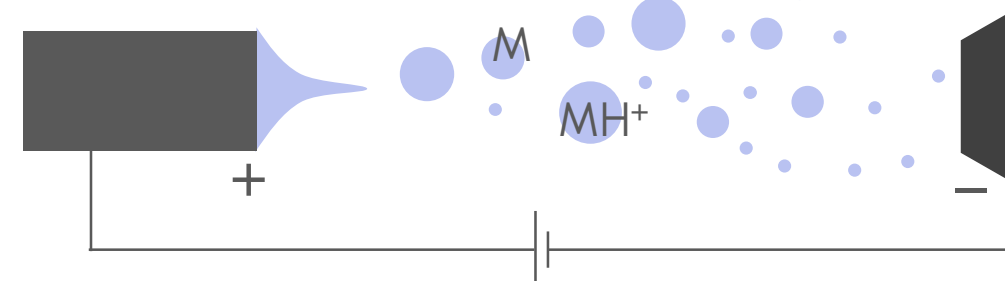


$$c_{\text{toxic}} \sim LC_{50}$$



Ionization efficiency  
(IE)

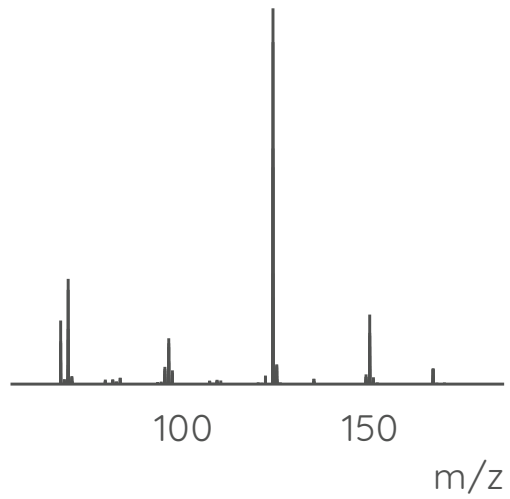
LC-ESI-MS



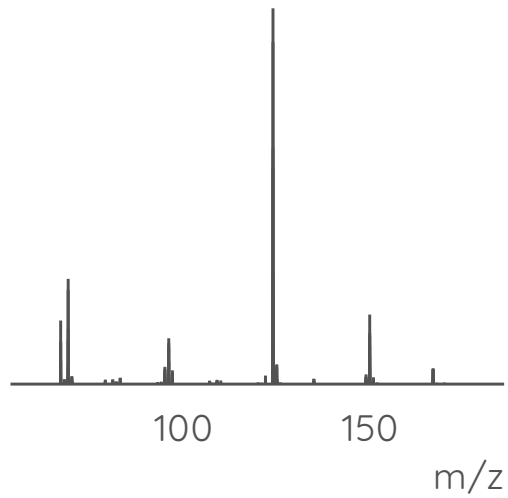
$$c_{\text{environmental}} \sim \frac{\text{Intensity}}{\text{IE}}$$

$$\text{PriorityScore} = \frac{\frac{\text{Intensity}}{\text{IE}}}{LC_{50}}$$

# Machine learning tools for risk estimation



# Machine learning tools for risk estimation

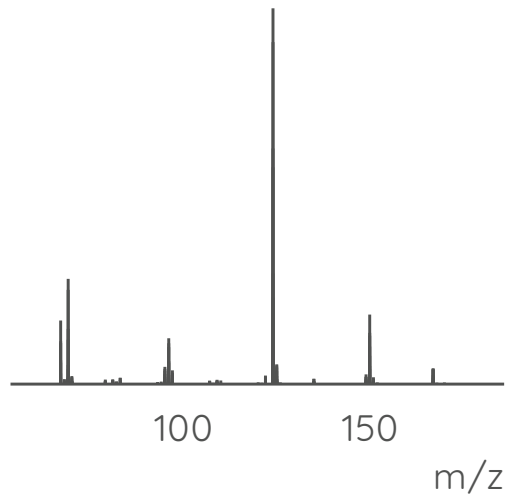


probability

C=O	0.3
CCN	0.9
Cl	0.0
F	0.1
C=C	0.2
CCO	0.8

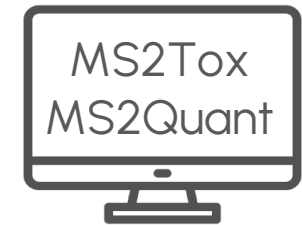


# Machine learning tools for risk estimation

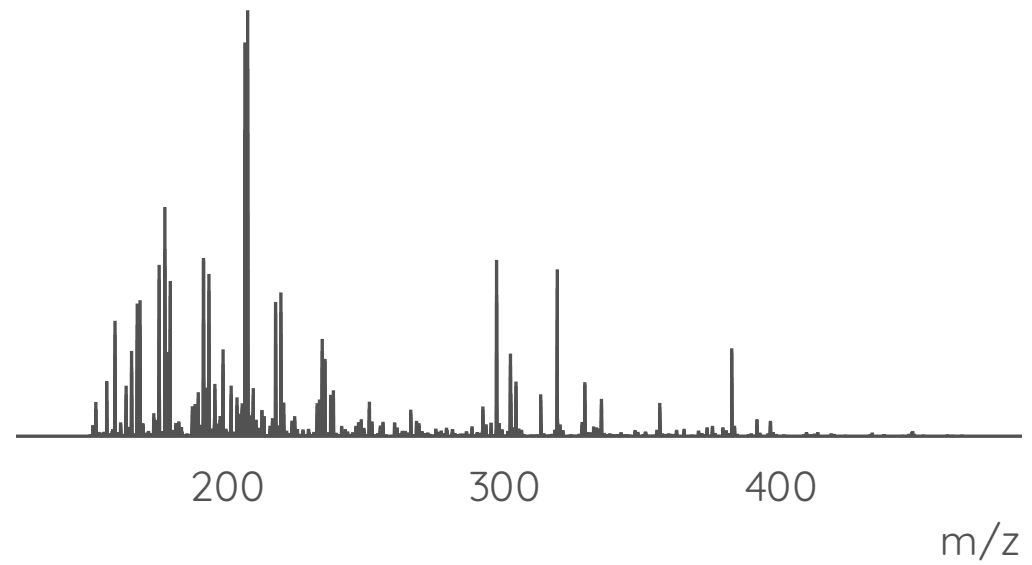


probability

C=O	0.3
CCN	0.9
Cl	0.0
F	0.1
C=C	0.2
CCO	0.8



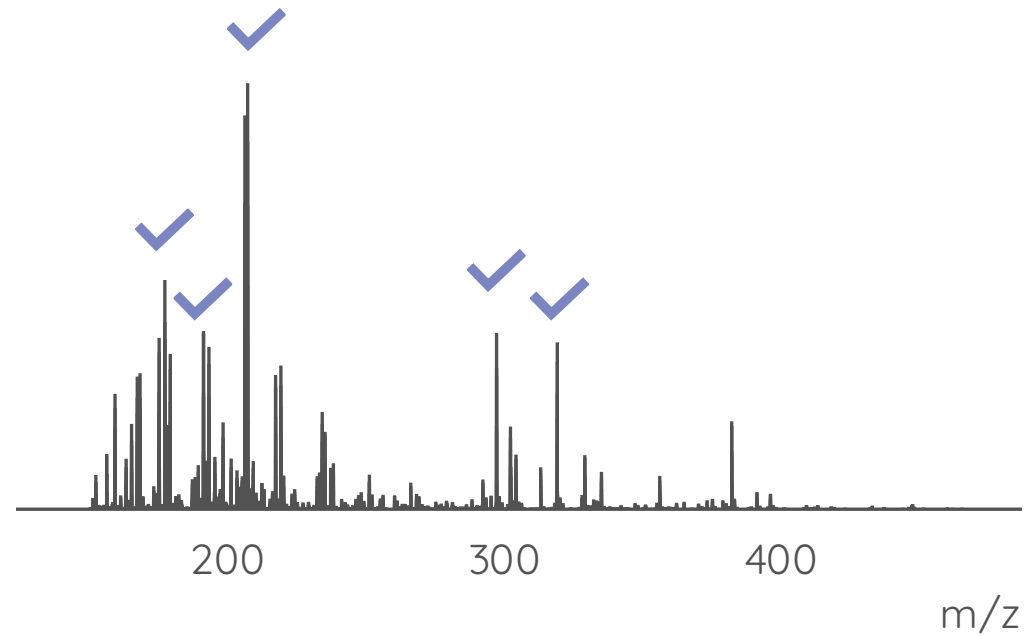
# Data acquisition approaches



# Data acquisition approaches

Top5

Top 5 highest  
Intensity peaks





# Data acquisition approaches

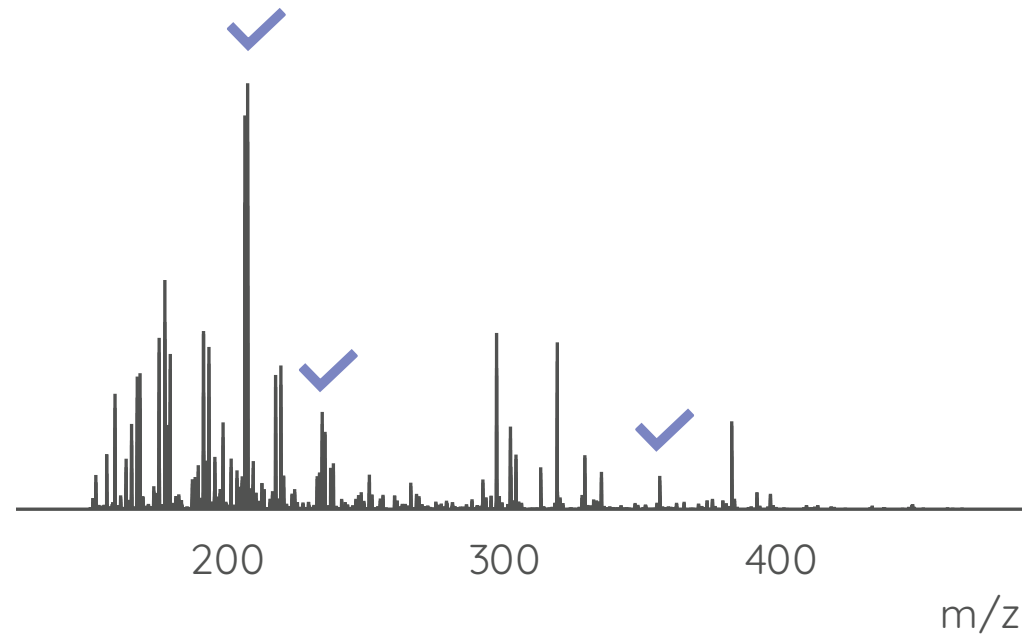
Top5

SWS

Surface water suspect  
list as inclusion list



523 compound  
masses



# Data acquisition approaches

Top5

SWS

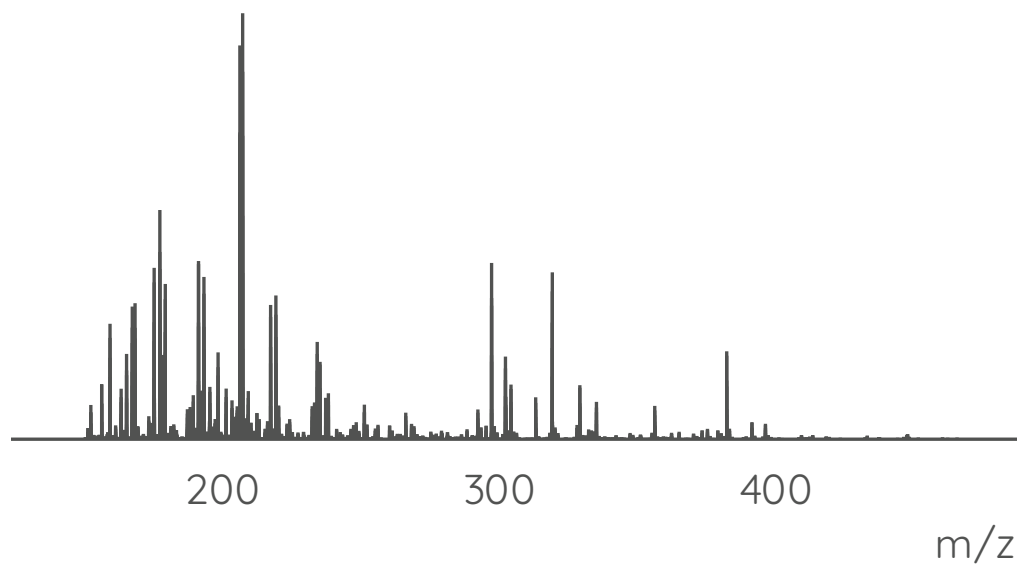
PCL1000

1000 high risk chemicals  
from PubChemLite

PubChemLite

SMILES	<i>m/z</i>
...	...
...	...
...	...
...	...
...	...

~450 K chemicals



# Data acquisition approaches

Top5

SWS

PCL1000

1000 high risk chemicals  
from PubChemLite

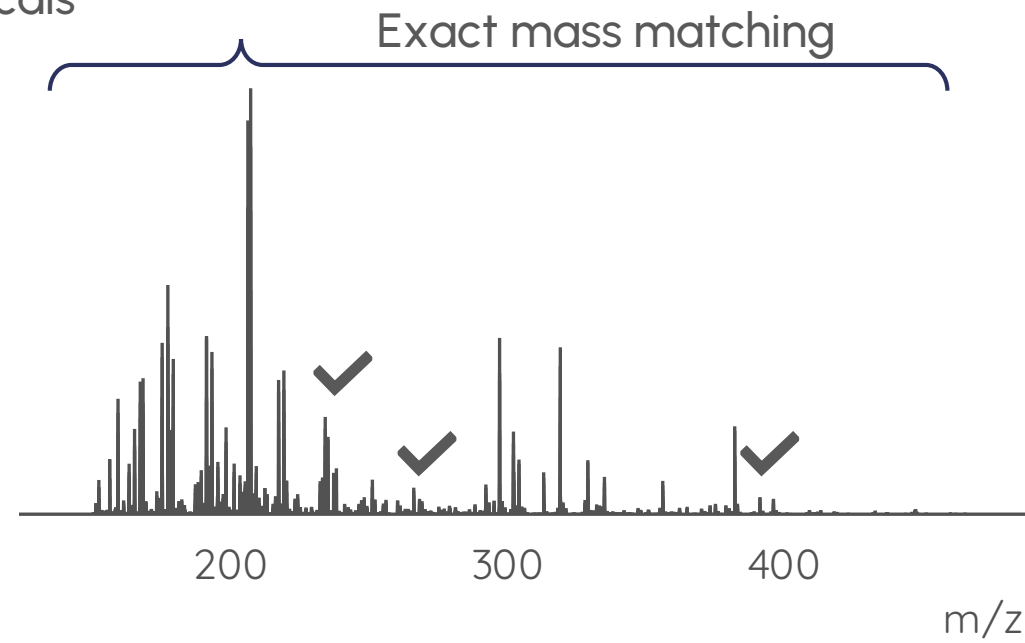
PubChemLite

SMILES	m/z
...	...
...	...
...	...
...	...
...	...

~450 K chemicals

Match

✓  
✗  
✓  
✓  
✗





# Data acquisition approaches

Top5

SWS

PCL1000

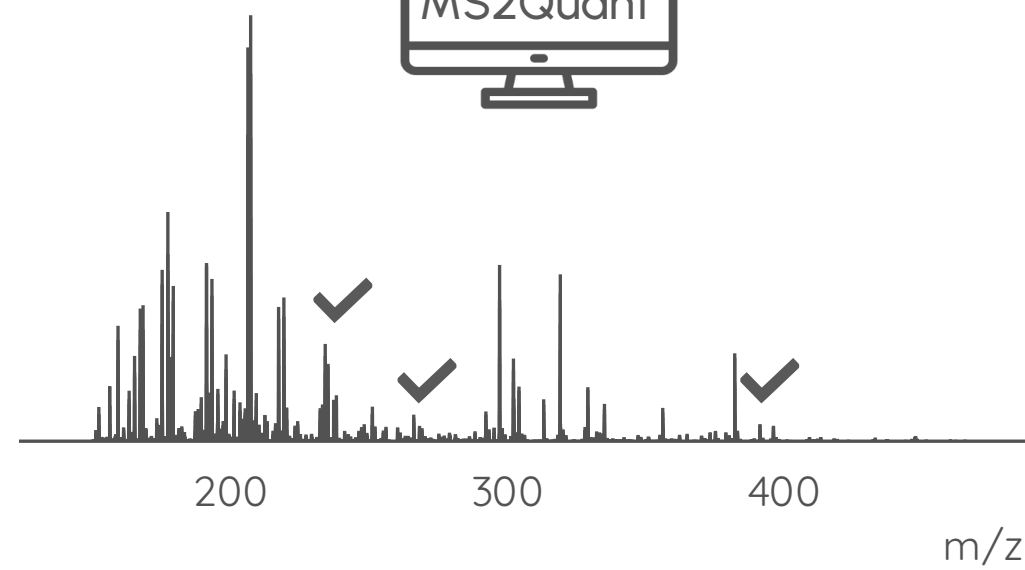
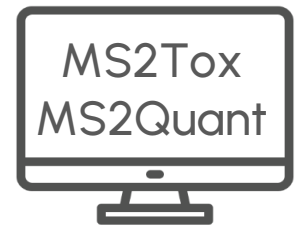
1000 high risk chemicals  
from PubChemLite

PubChemLite

SMILES	m/z
...	...
...	...
...	...
...	...
...	...

~450 K chemicals

Match	IE	LC <sub>50</sub>
✓	...	...
✗	...	...
✓	...	...
✓	...	...
✗	...	...



# Data acquisition approaches

Top5

SWS

PCL1000

1000 high risk chemicals  
from PubChemLite

PubChemLite

SMILES	m/z
...	...
...	...
...	...
...	...
...	...

~450 K chemicals

Match

✓  
✗  
✓  
✓  
✗

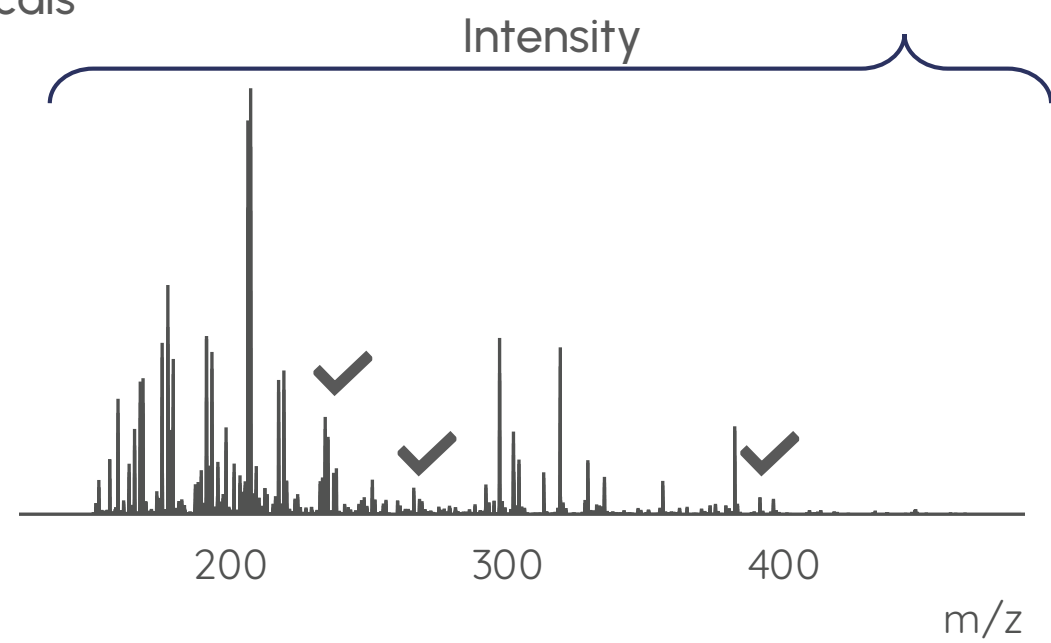
IE

...  
...  
...  
...

LC<sub>50</sub>

...  
...  
...  
...

$$\text{PriorityScore} = \frac{\frac{\text{Intensity}}{\text{IE}}}{\text{LC}_{50}}$$



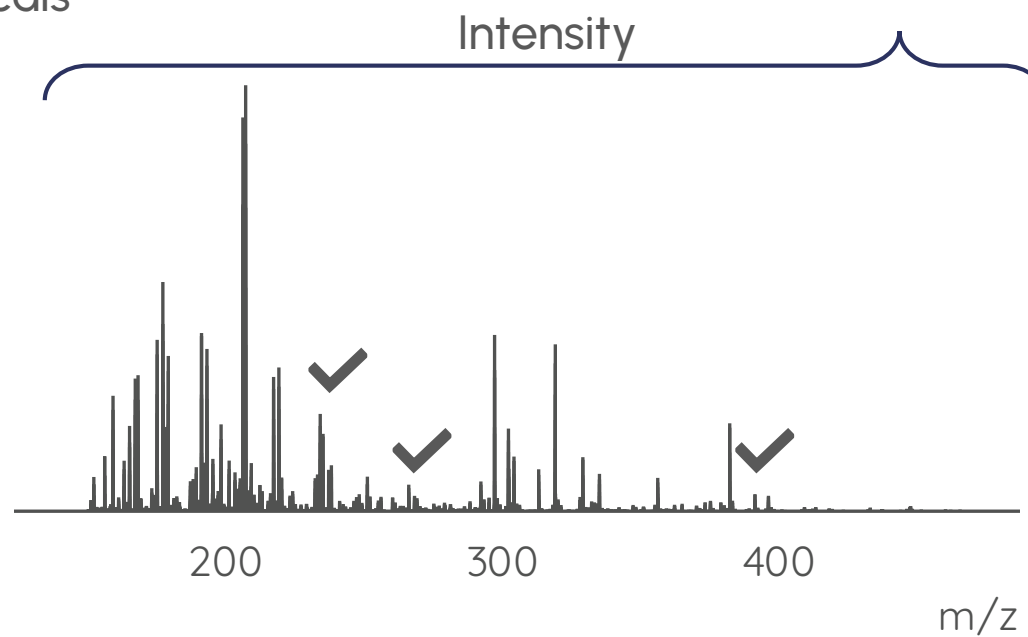
# Data acquisition approaches

Top5

SWS

PCL1000

1000 high risk chemicals  
from PubChemLite



# Data acquisition approaches

Top5

SWS

PCL1000

1000 high risk chemicals  
from PubChemLite

PubChemLite

SMILES	m/z
...	...
...	...
...	...
...	...
...	...

~450 K chemicals

Match



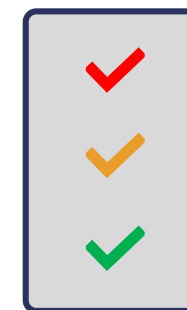
IE

...

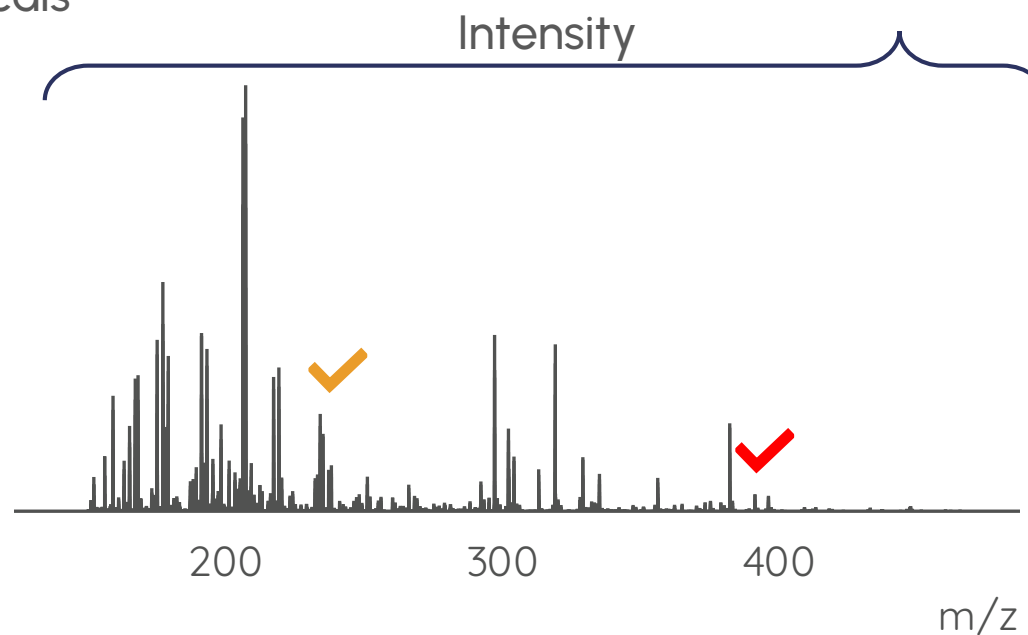
LC<sub>50</sub>

...

PriorityScore



1000  
chemicals





# Data acquisition approaches

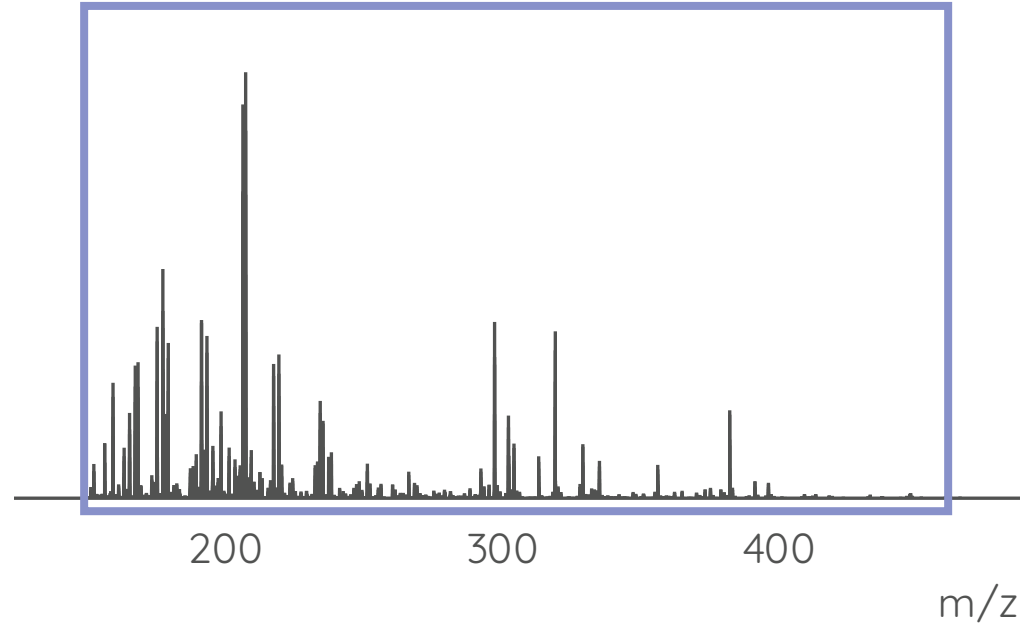
Top5

SWS

PCL1000

DIA

Data independent  
acquisition



# Comparison of methods - standards



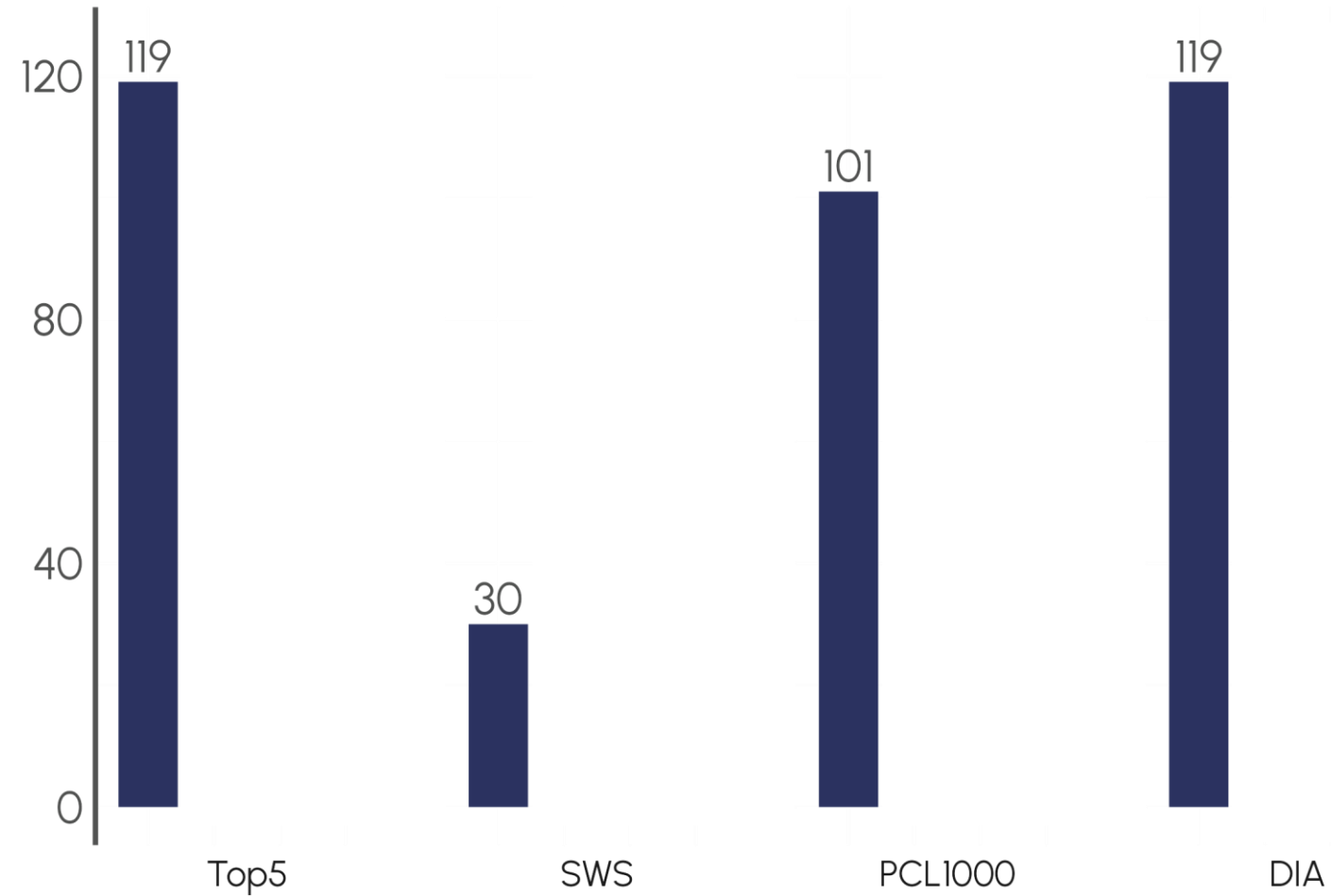
Standard mix  
119 chemicals

# Comparison of methods - standards

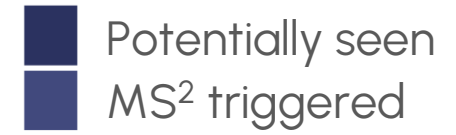
■ Potentially seen



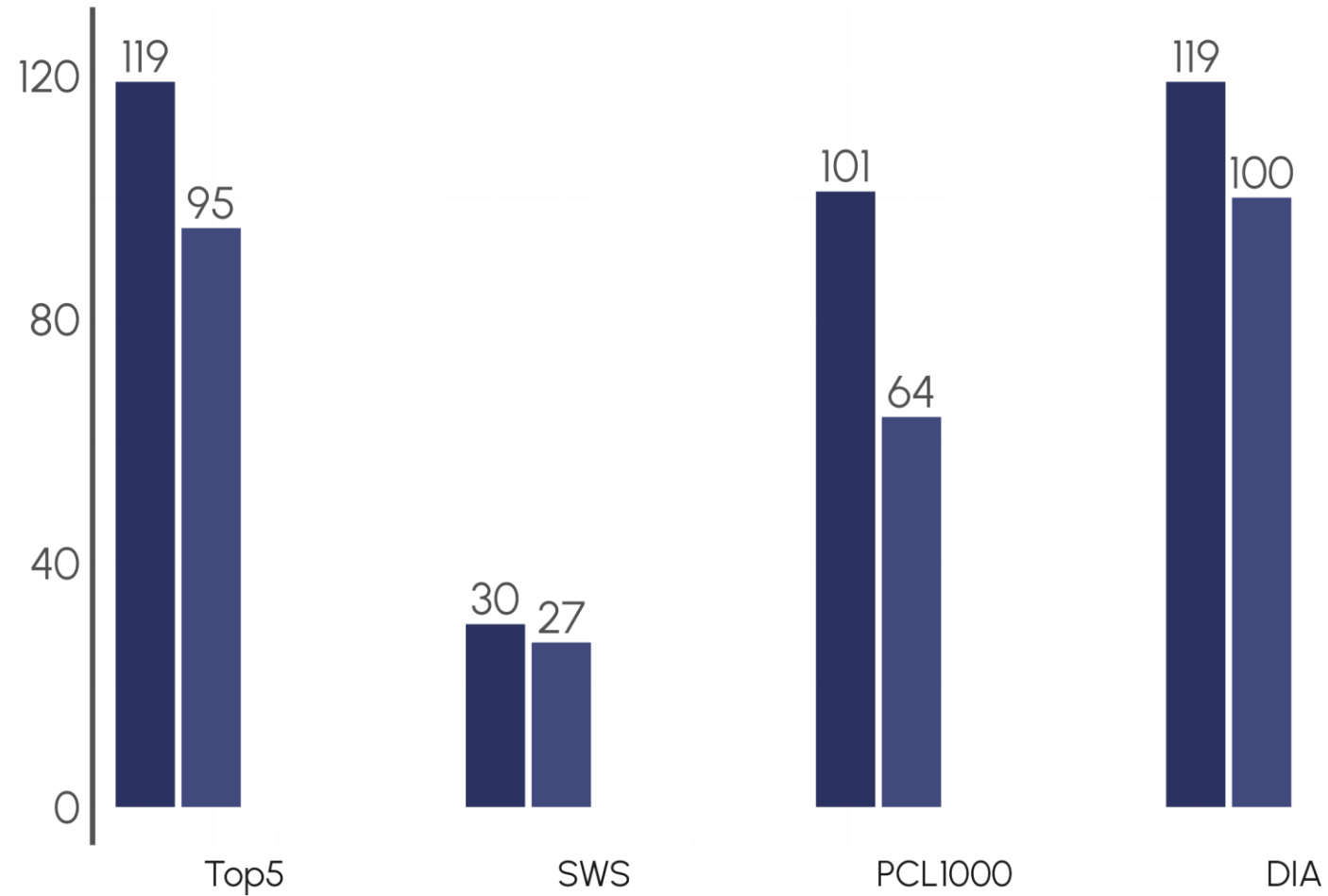
Standard mix  
119 chemicals



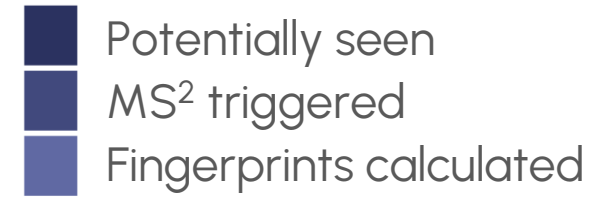
# Comparison of methods - standards



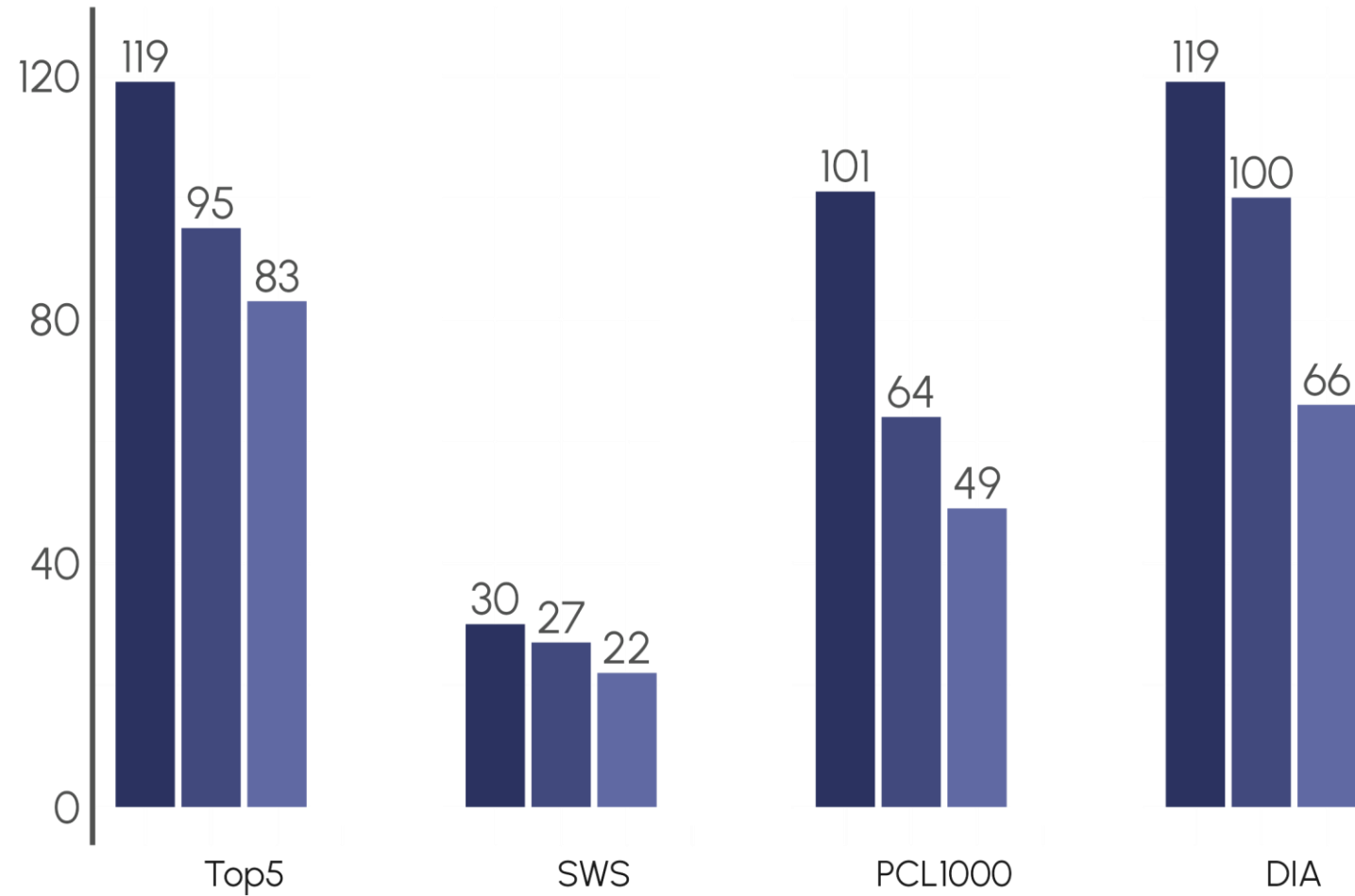
Standard mix  
119 chemicals



# Comparison of methods - standards



Standard mix  
119 chemicals

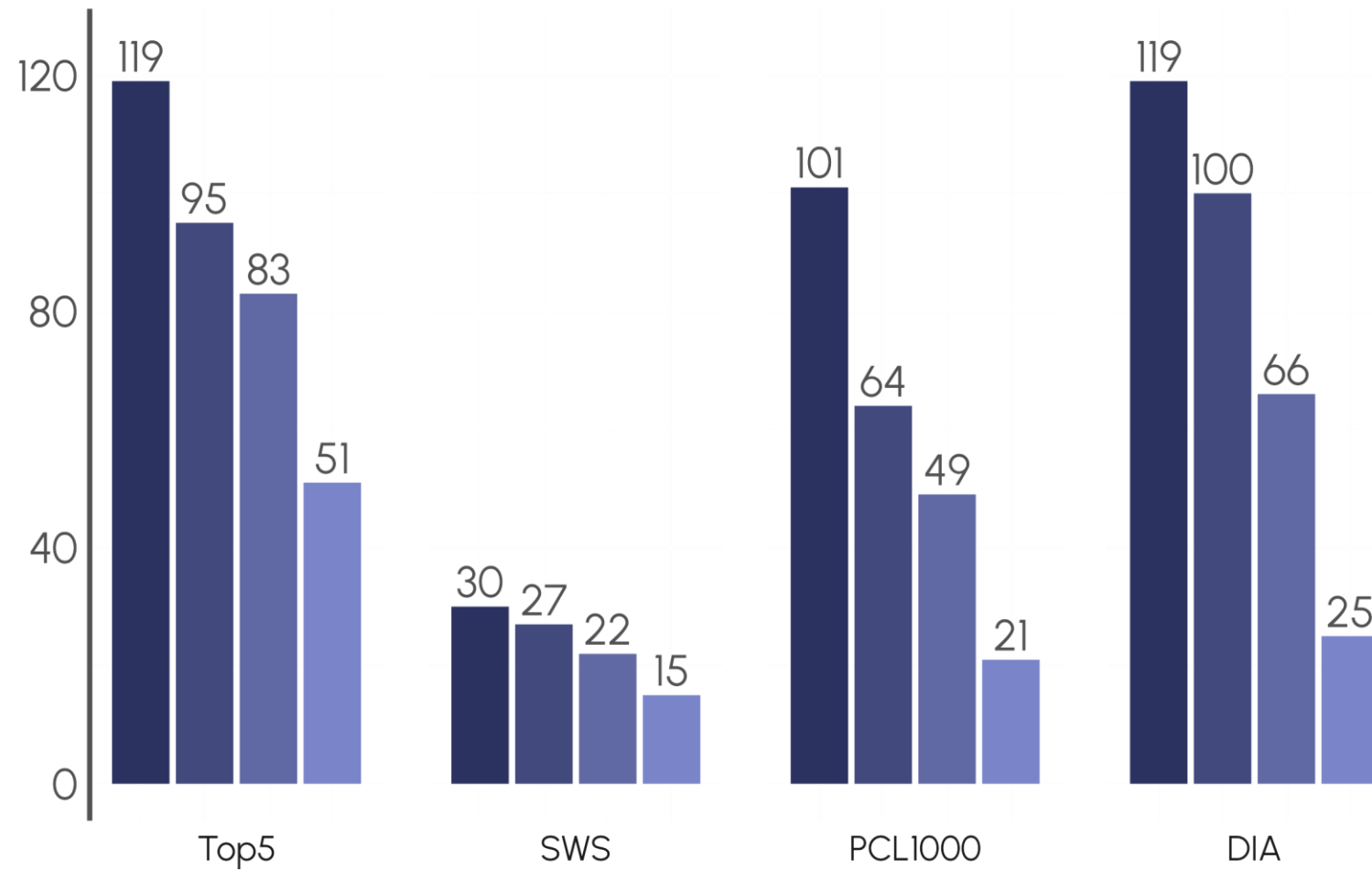
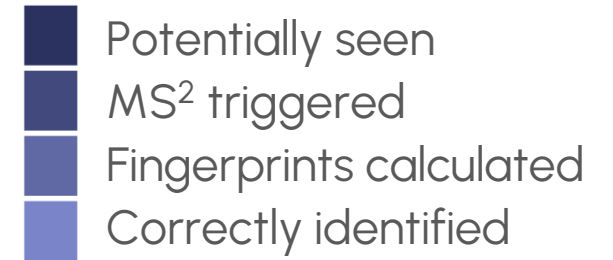




# Comparison of methods - standards



Standard mix  
119 chemicals

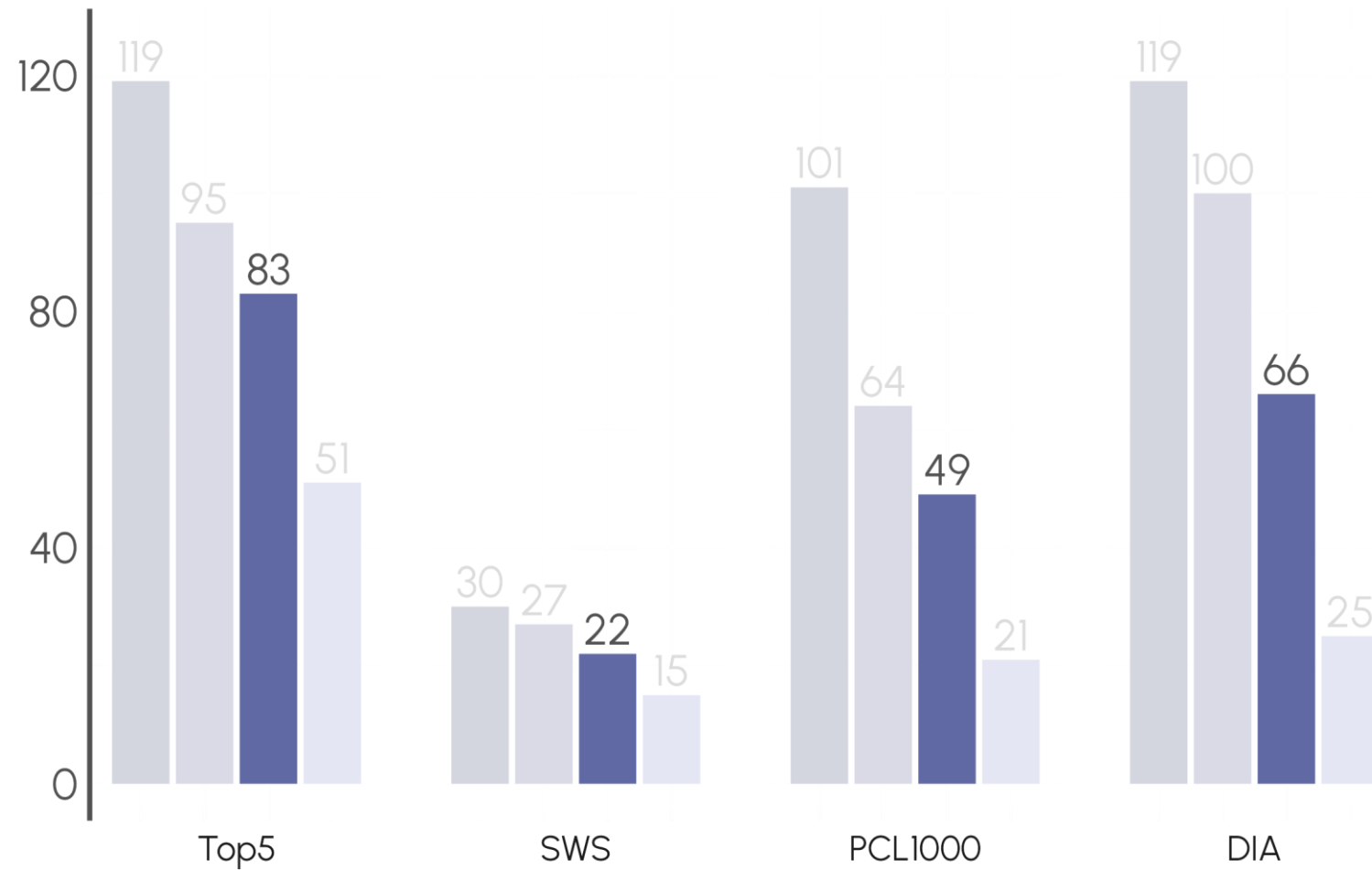


# Comparison of methods - standards



Standard mix  
119 chemicals

- Potentially seen
- MS<sup>2</sup> triggered
- Fingerprints calculated
- Correctly identified



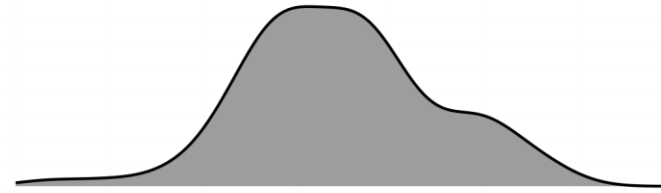
# Comparison of methods - standards



Standard mix  
119 chemicals

$$\text{PriorityScore} = \frac{\frac{\text{Intensity}}{\text{IE}}}{\text{LC}_{50}}$$

Structure



Top5



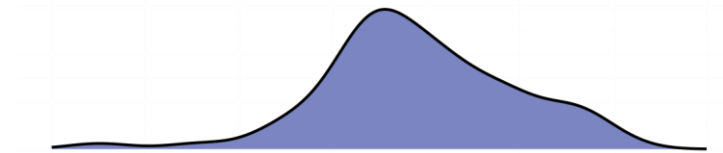
SWS



PCL1000



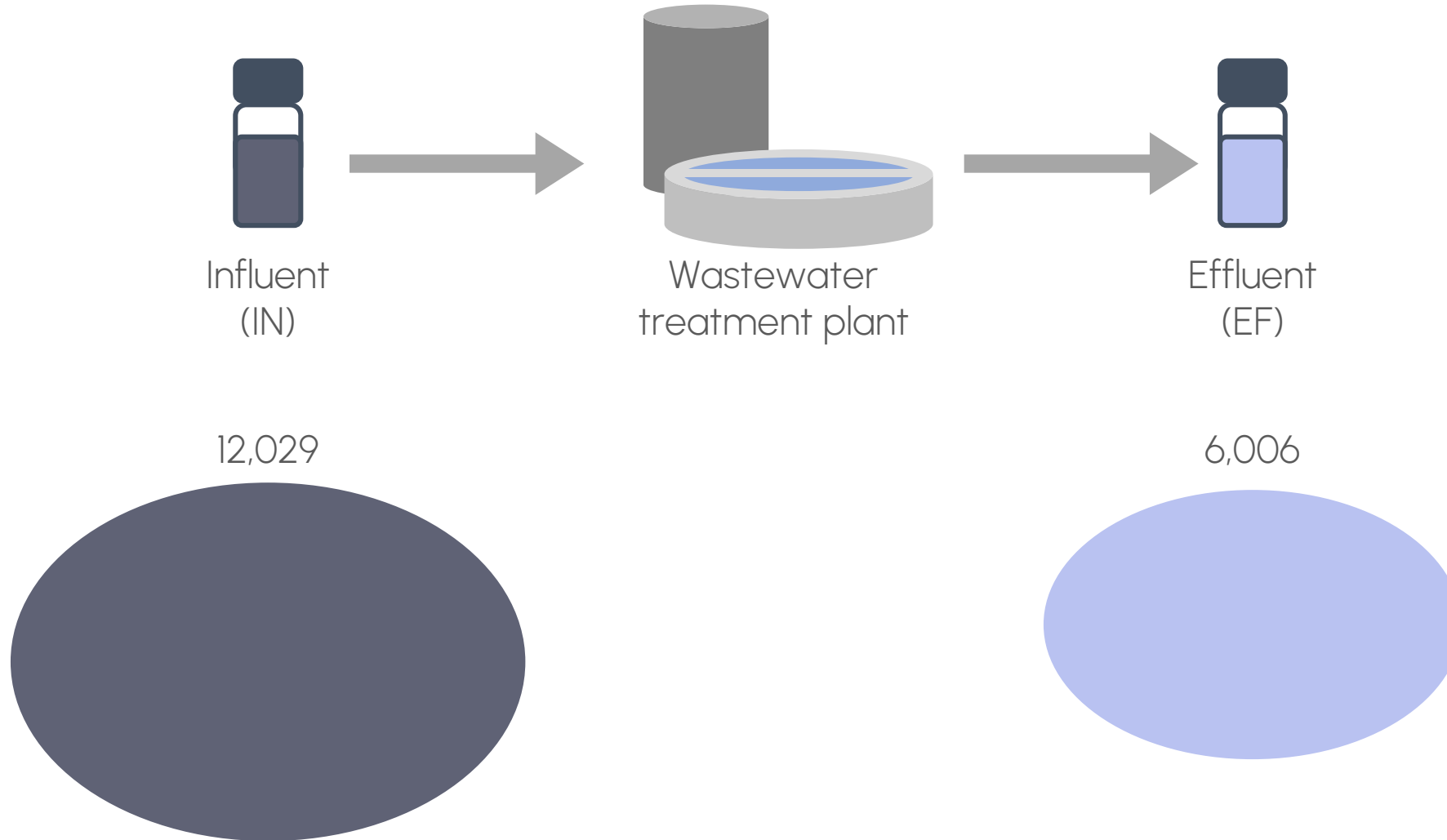
DIA



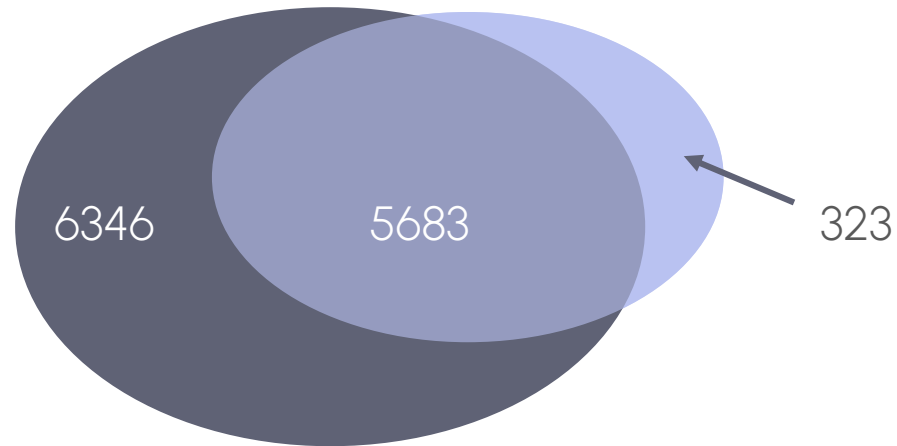
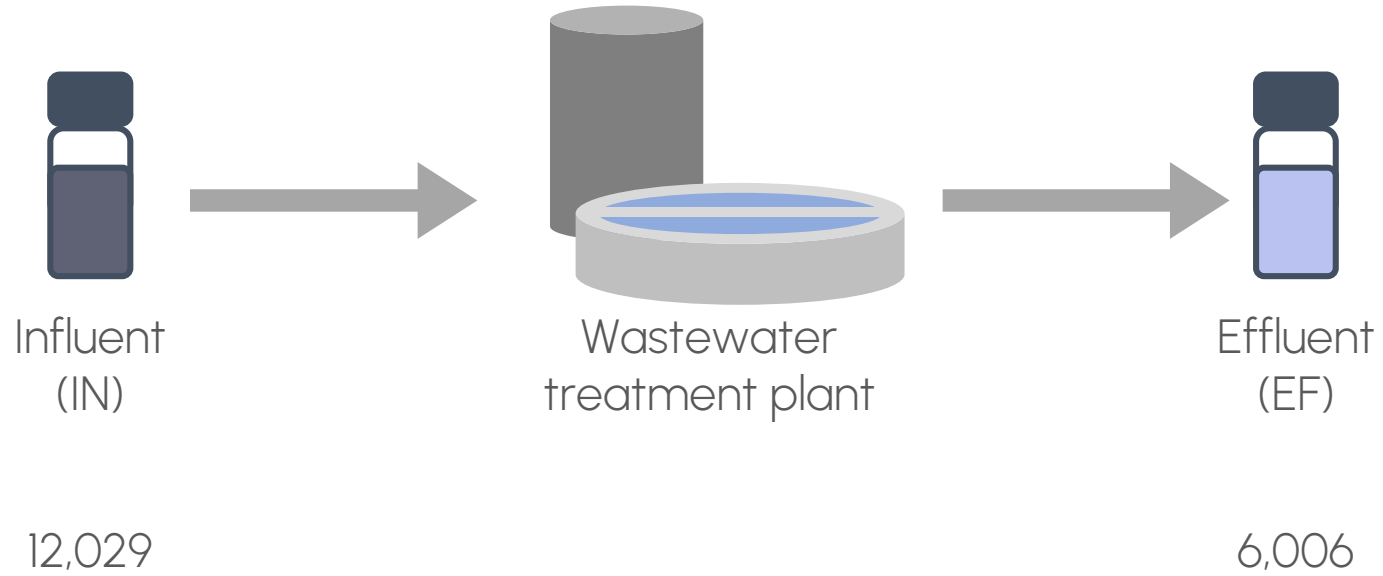
4 6 8 10

$\log_{10}(\text{PriorityScore})$

# Comparison of methods - wastewater



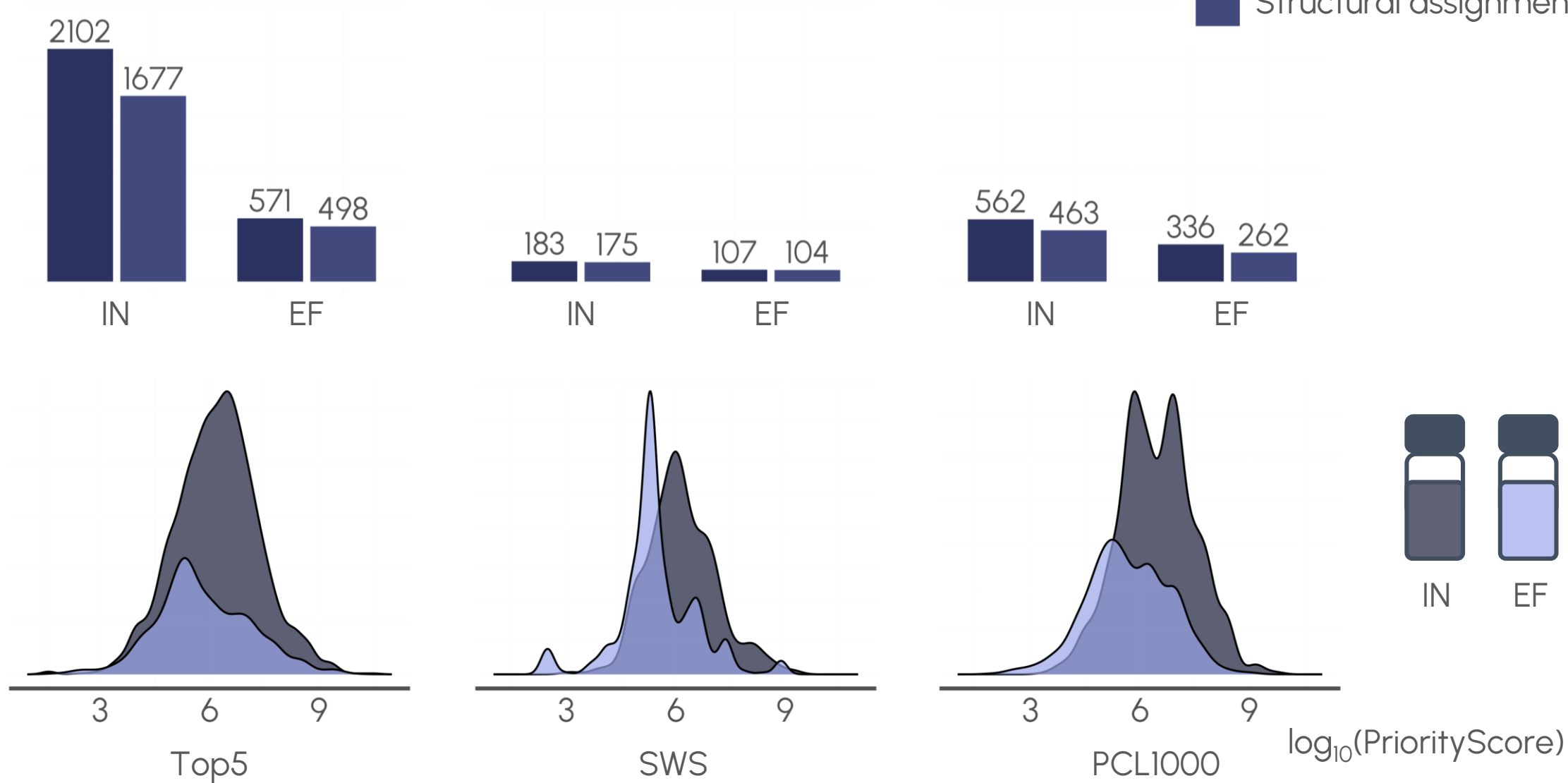
# Comparison of methods - wastewater





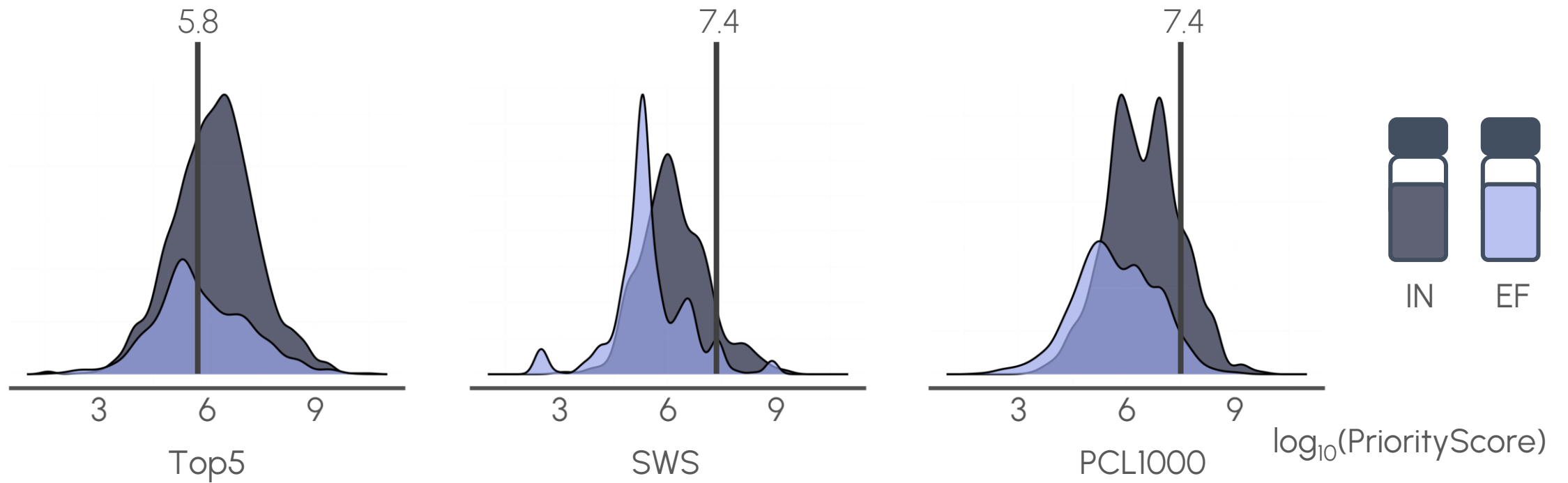
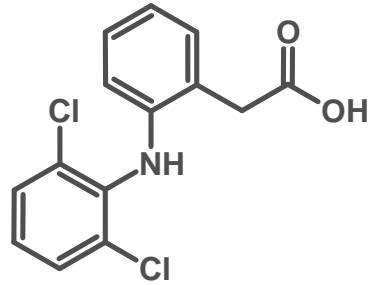
# Comparison of methods - wastewater

Fingerprints calculated  
Structural assignment



# Comparison of methods - wastewater

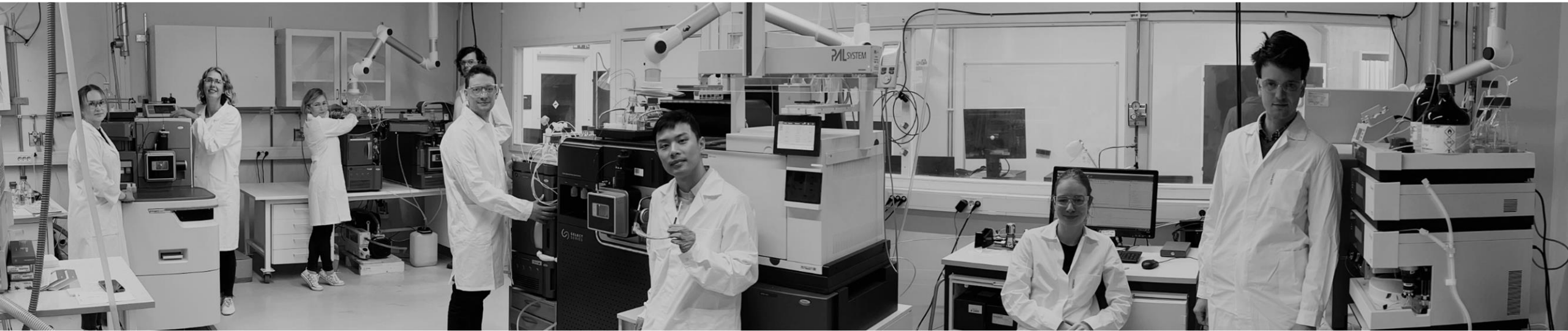
Diclofenac



# Preliminary conclusions and future prospects

- Chemicals triggered by different approaches cover similar range of priority scores
- Acquisition methods are complementary
- Confirmation of detected chemicals

# Acknowledgements



Lisa Jonsson, Louise Malm, Pilleriin Peets, Malte Posselt, Michael McLachlan, Matthew MacLeod, Magnus Breitholtz, Jonathan Martin, Anneli Kruve

Department of Materials and Environmental Chemistry and Department of Environmental Science in Stockholm University

FORMAS RapMixTox #2021-01511

Helen Sepman      Helen.Sepman@aces.su.se

*Kruvelab.com*