

Sustainable design of chemical reagents for sensitive detection of pesticides using a machine learning workflow

SCS2025, Västerås, 16.06.2025
Analytical chemistry session

Henrik Hupatz

Postdoctoral researcher

Stockholm University
Department of Chemistry

henrik.hupatz@su.se

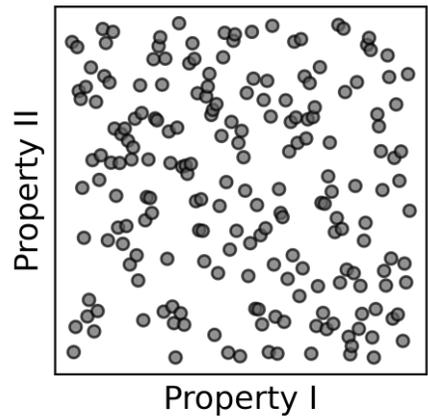
Kruve Lab



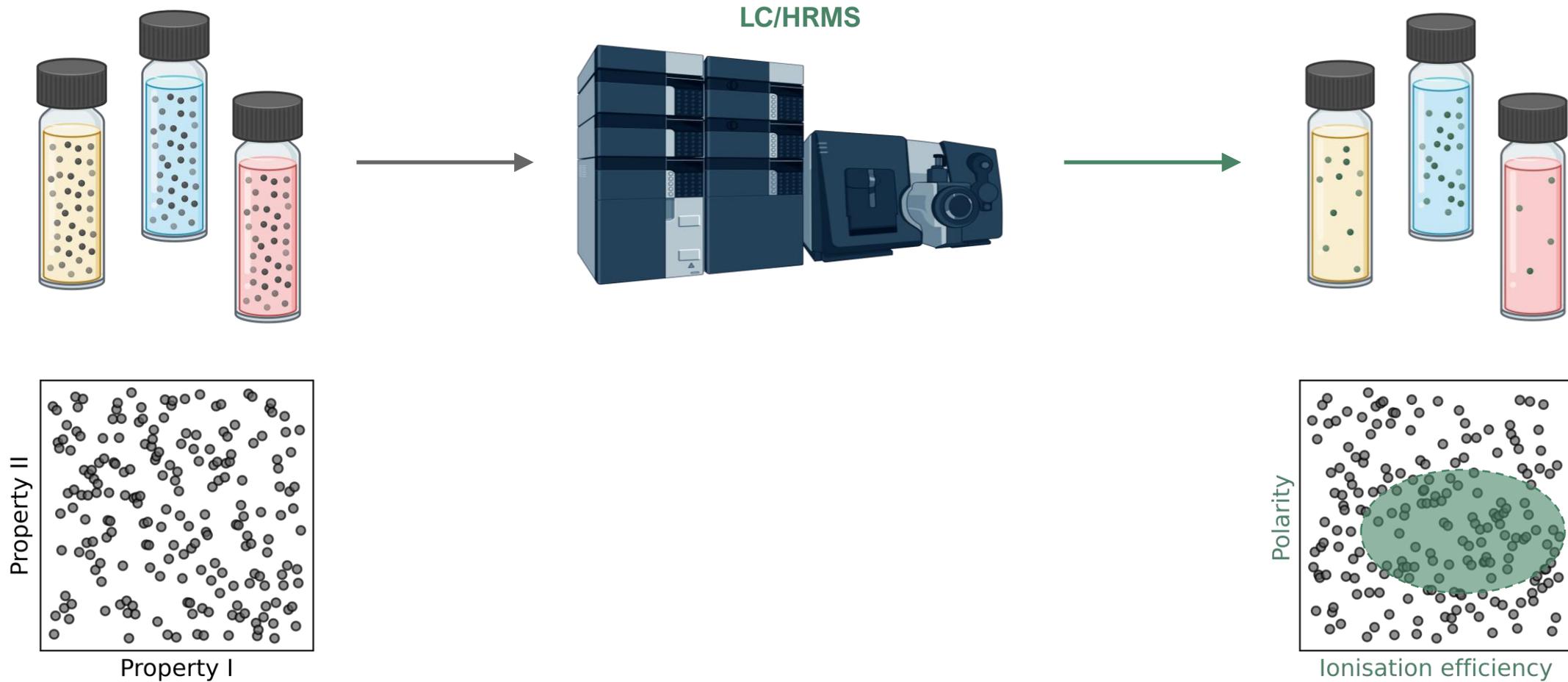
Stockholm
University

SUCCeSS

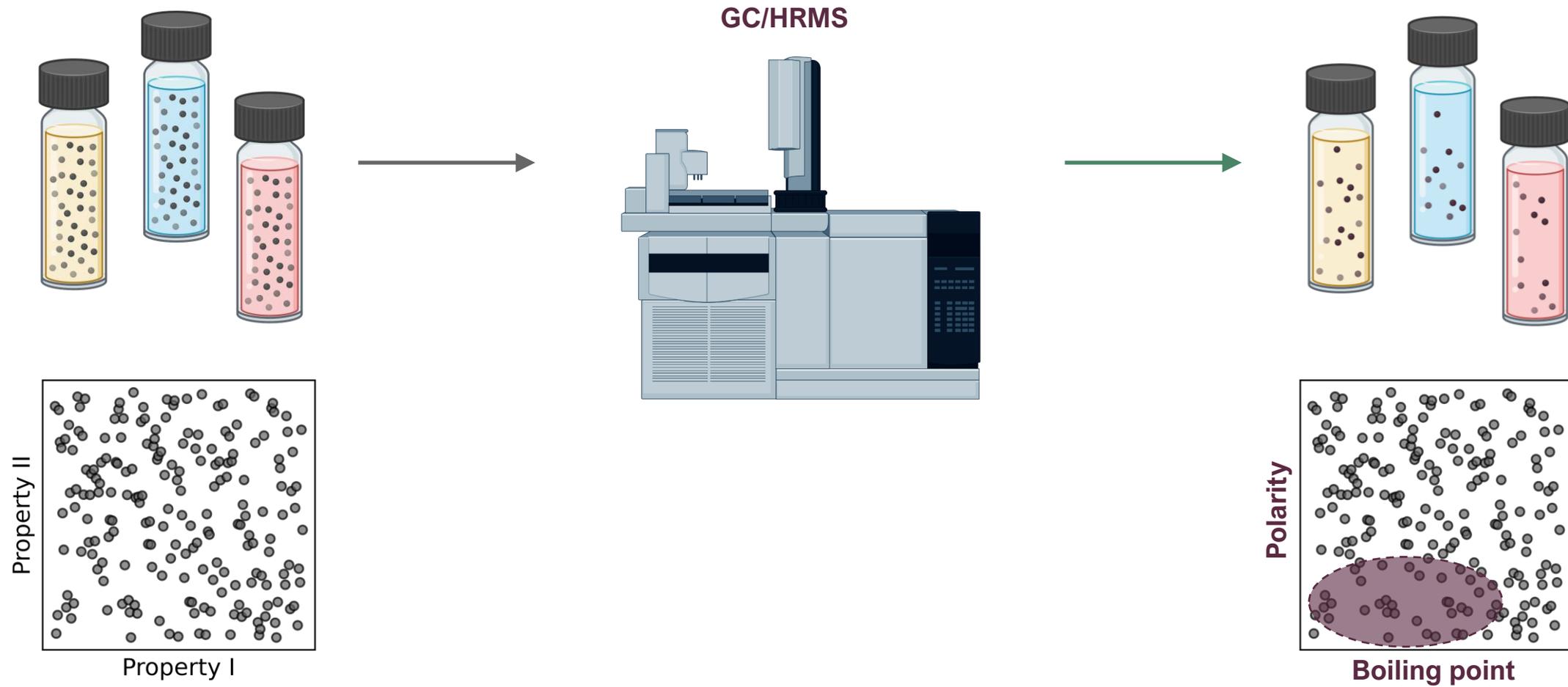
INTRODUCTION



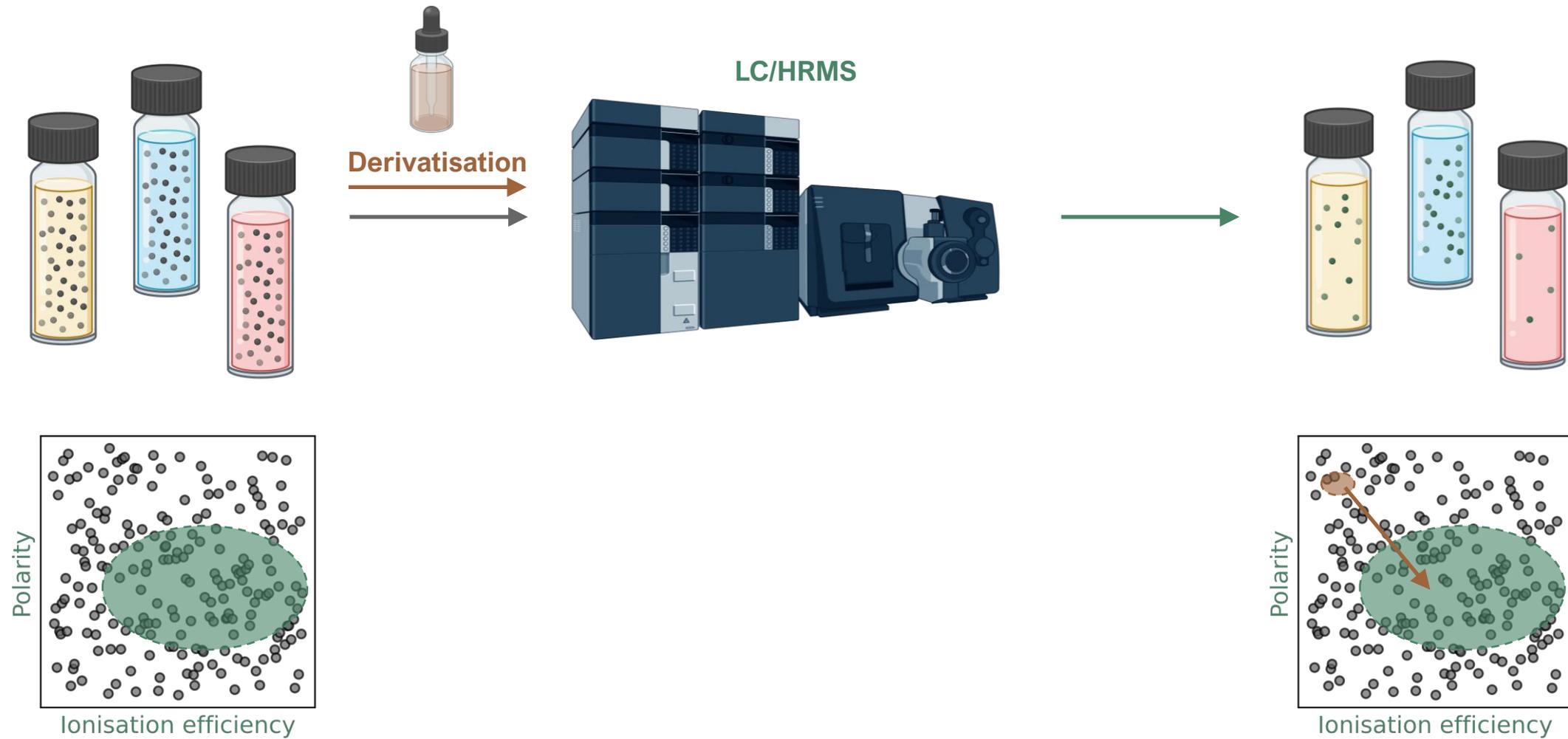
INTRODUCTION



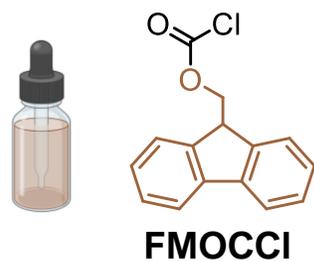
INTRODUCTION



INTRODUCTION



DERIVATISATION REAGENT DESIGN



DERIVATISATION REAGENT DESIGN



Experimental
experience
Literature review
Trial and error



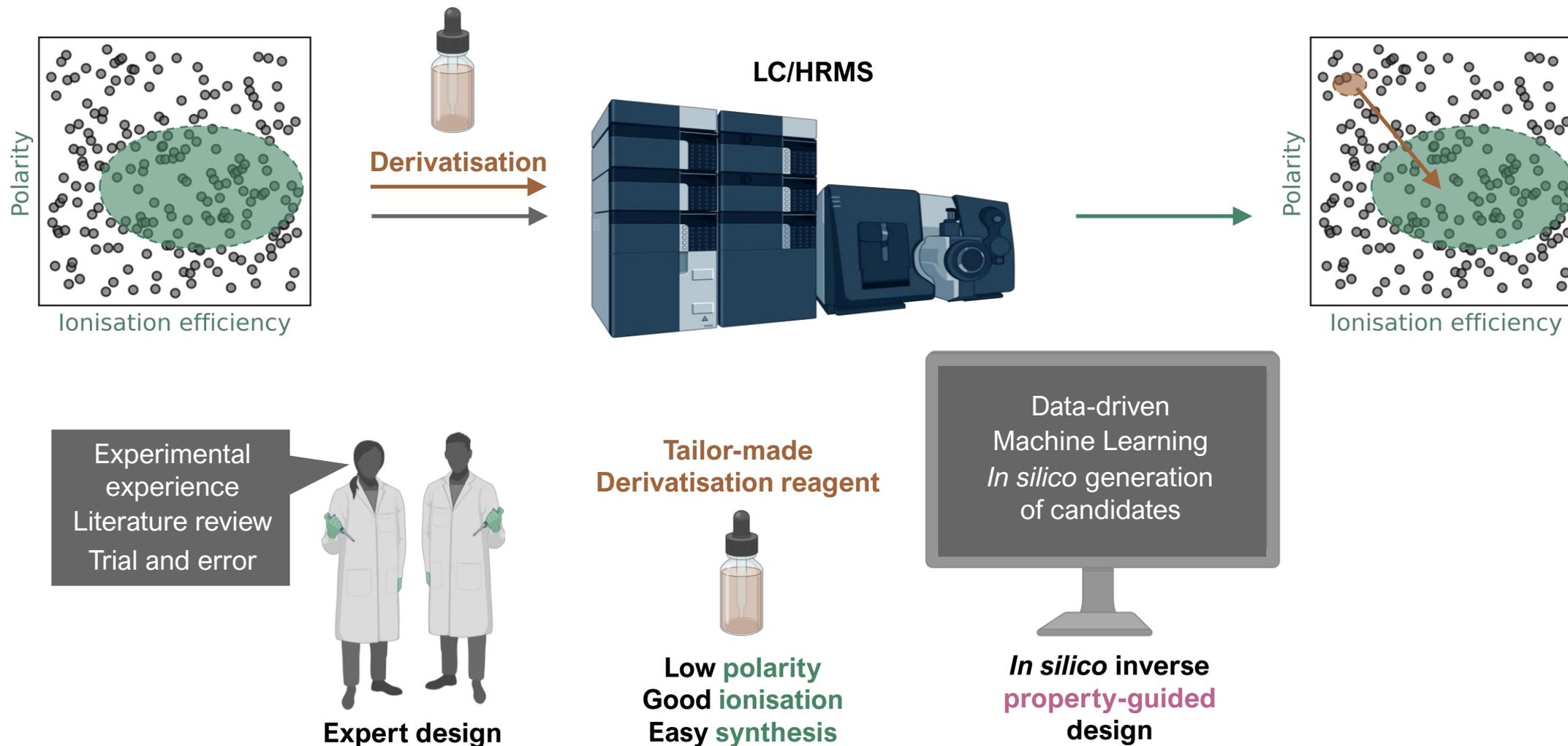
Expert design

Tailor-made
derivatisation reagent

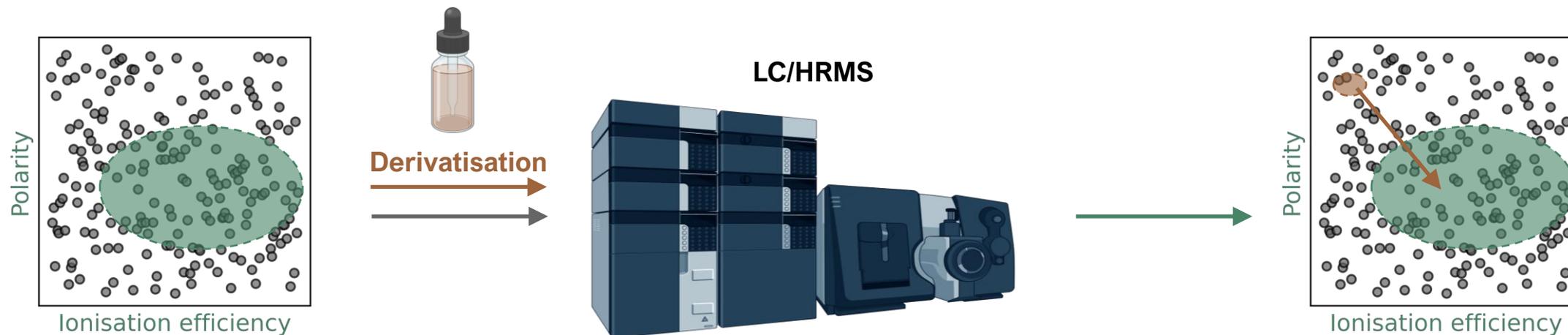


Low polarity
Good ionisation
Easy synthesis

DERIVATISATION REAGENT DESIGN



DERIVATISATION REAGENT DESIGN



Experimental experience
Literature review
Trial and error



Expert design

Tailor-made
derivatisation reagent



Low polarity
Good ionisation
Easy synthesis

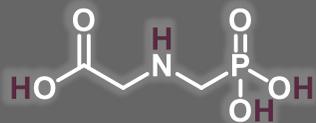
Data-driven
Machine Learning
In silico generation
of candidates

In silico inverse
property-guided
design

CAN A MACHINE
LEARNING WORKFLOW
DESIGN SUITABLE
DERIVATISATION
REAGENTS FOR
LC/HRMS?

DERIVATISATION REAGENT DESIGN

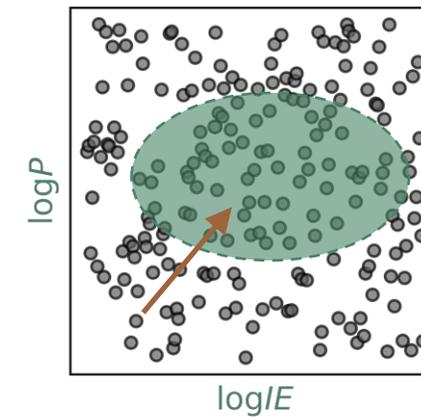
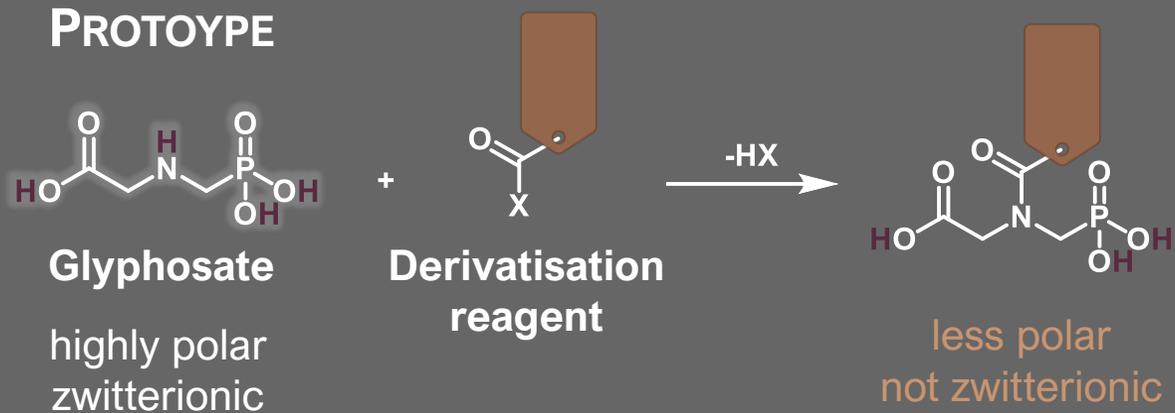
PROTOTYPE



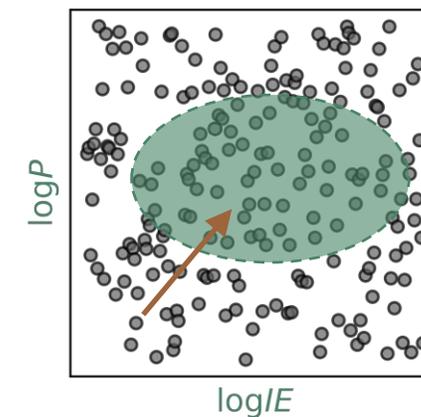
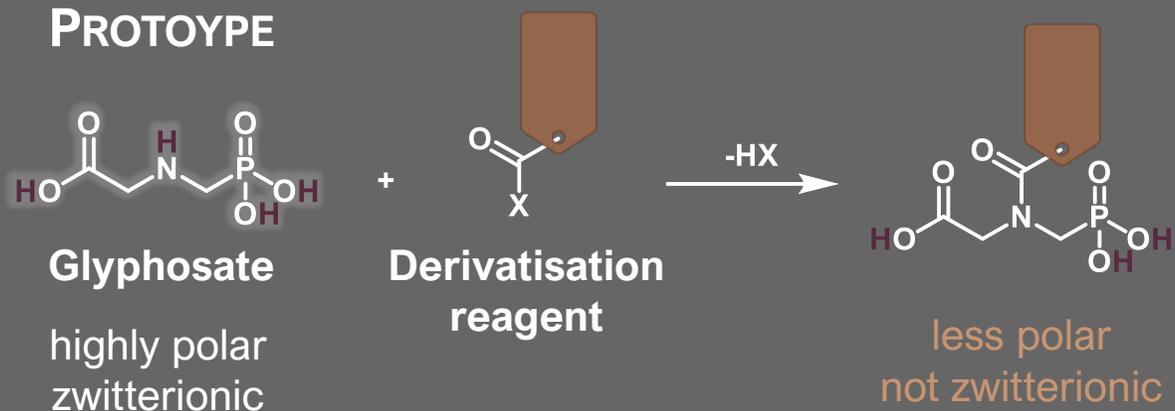
Glyphosate

highly polar
zwitterionic

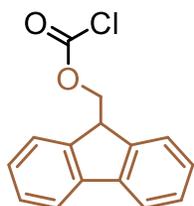
DERIVATISATION REAGENT DESIGN



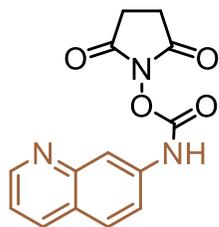
DERIVATISATION REAGENT DESIGN



COMMERCIAL REAGENTS

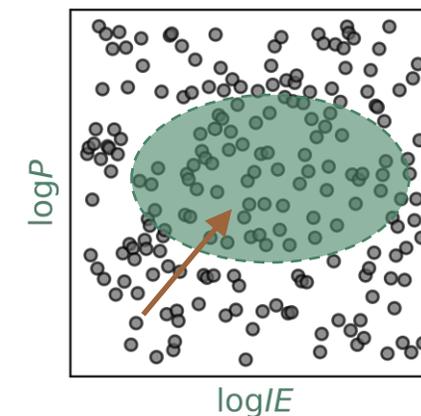
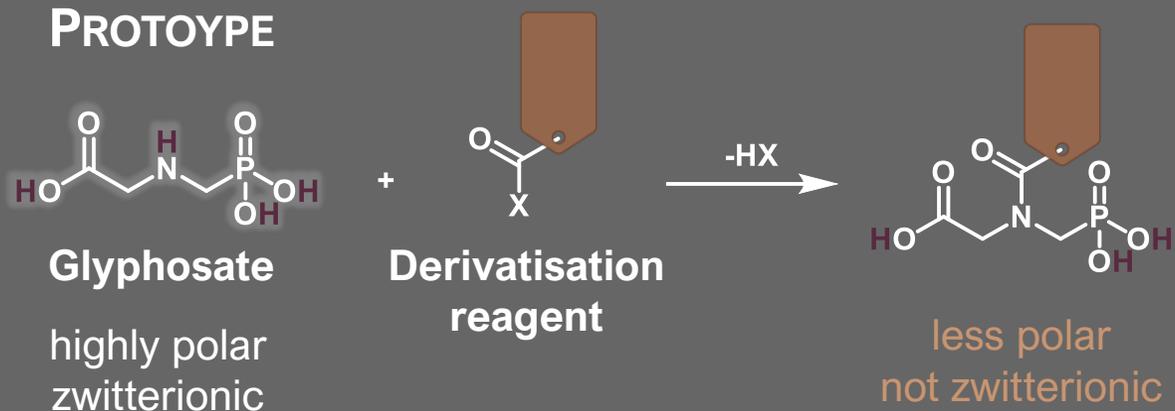


FMOCCI

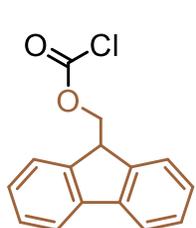


ACQ

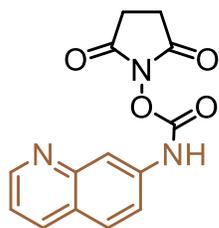
DERIVATISATION REAGENT DESIGN



COMMERCIAL REAGENTS



FMOCCI



ACQ

EXPERIMENTAL FOUNDATION



Reagent synthesis
Easy and cheap



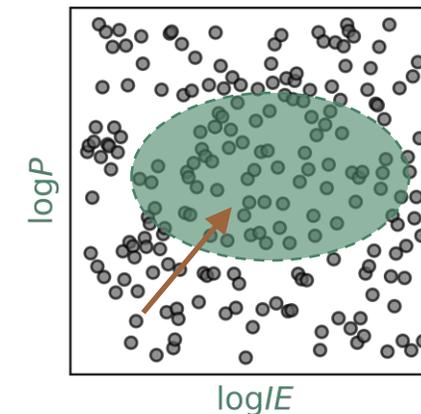
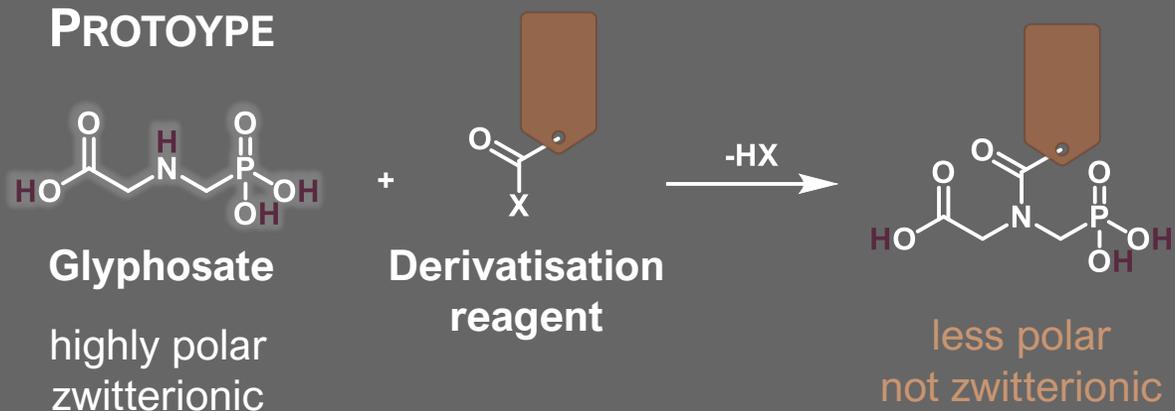
Sample preparation
Suitable reactivity



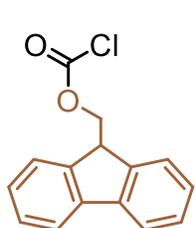
Separation
Suitable retention time

Detection
High ionization efficiency

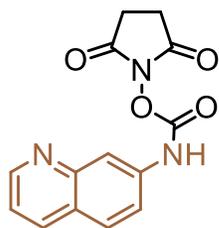
DERIVATISATION REAGENT DESIGN



COMMERCIAL REAGENTS



FMOCCI



ACQ

EXPERIMENTAL FOUNDATION



Reagent synthesis
Easy and cheap

SAScore



Sample preparation
Suitable reactivity

-COOH group



Separation
Suitable retention time

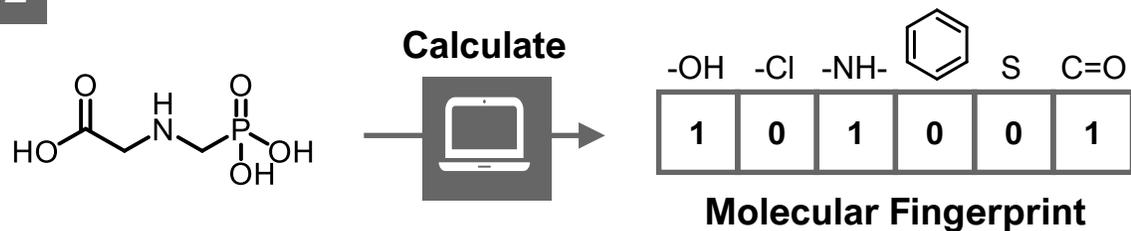
logP

Detection
High ionization efficiency

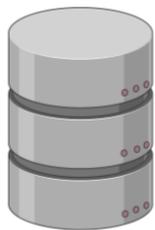
log/E

IONIZATION EFFICIENCY PREDICTION

log*E*



log*E*
dataset

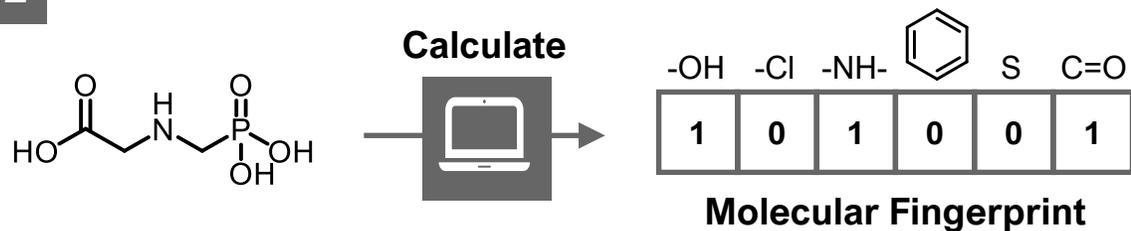


6120
datapoints

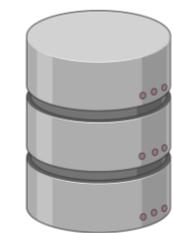
different
experimental
conditions

IONIZATION EFFICIENCY PREDICTION

log*E*

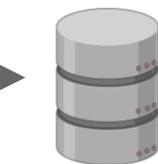


log*E*
dataset

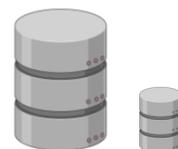
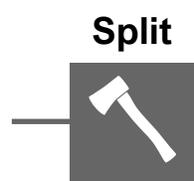


6120
datapoints

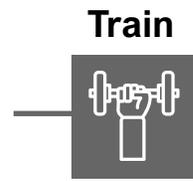
different
experimental
conditions



419
chemicals

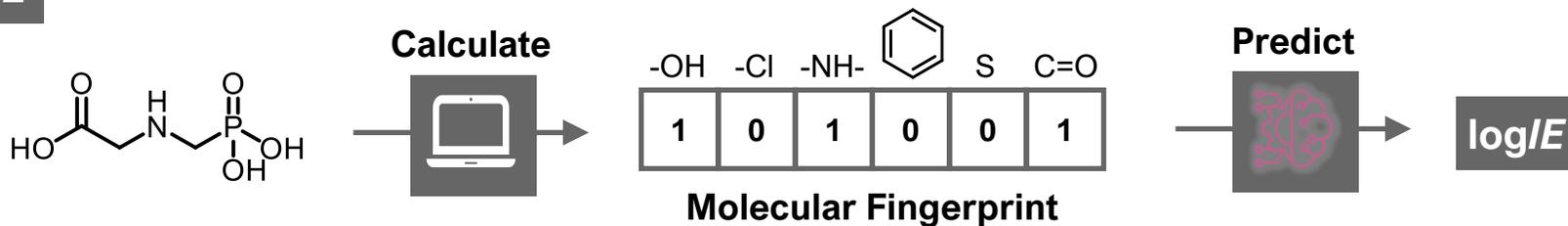


335:82
train/test set

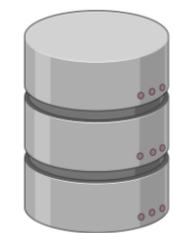


IONIZATION EFFICIENCY PREDICTION

log*E*



log*E*
dataset



6120
datapoints

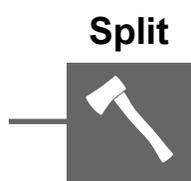
*different
experimental
conditions*



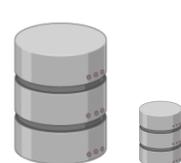
Clean



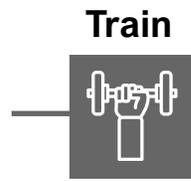
419
chemicals



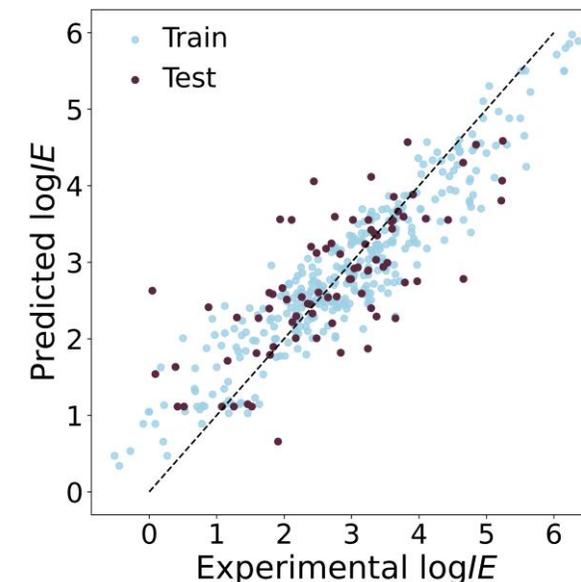
Split



335:82
train/test set



Train



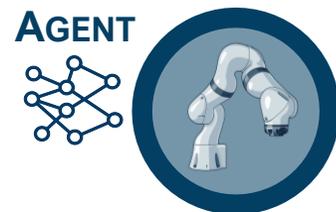
Selected

Algorithm: **XGBoost**

Representation: **MACCS fingerprints**

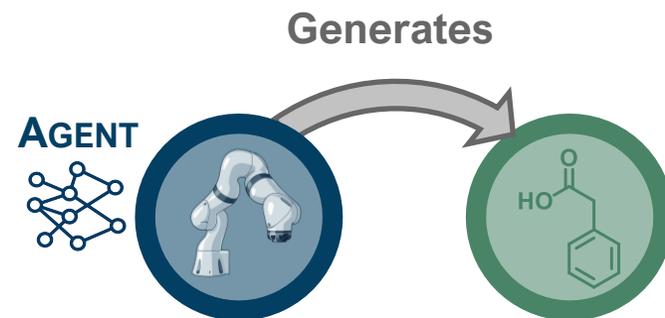
RMSE on test set: **0.78**

REINFORCEMENT LEARNING



Data-driven
Machine Learning
In silico generation
of candidates

REINFORCEMENT LEARNING

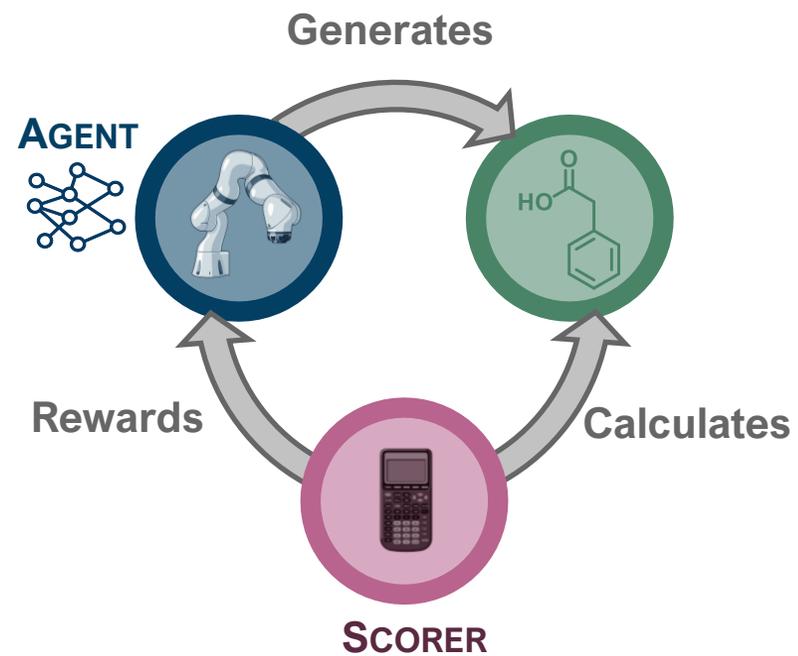


Data-driven
Machine Learning
In silico generation
of candidates

REINFORCEMENT LEARNING

SCORER

- $7 = \log P$

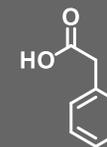
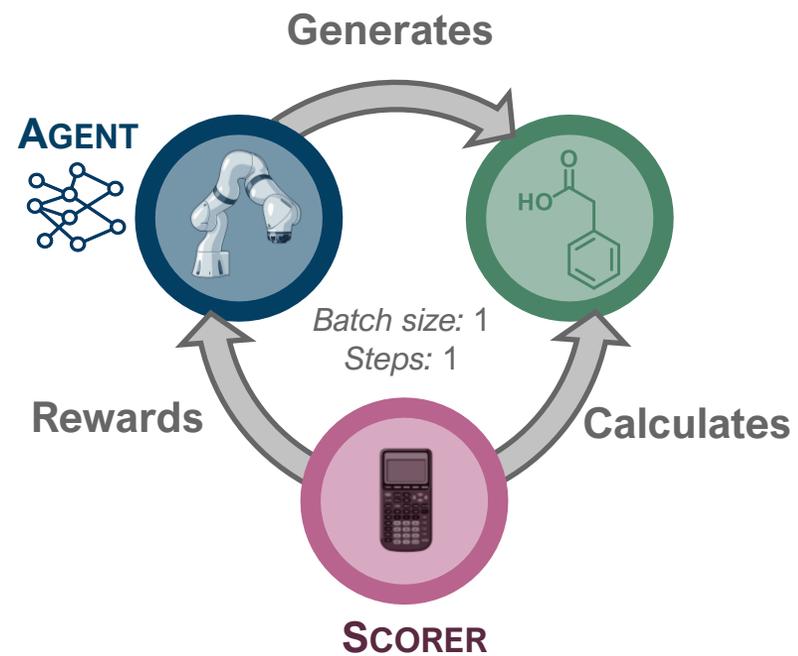


Data-driven
Machine Learning
In silico generation
of candidates

REINFORCEMENT LEARNING

SCORER

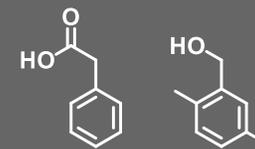
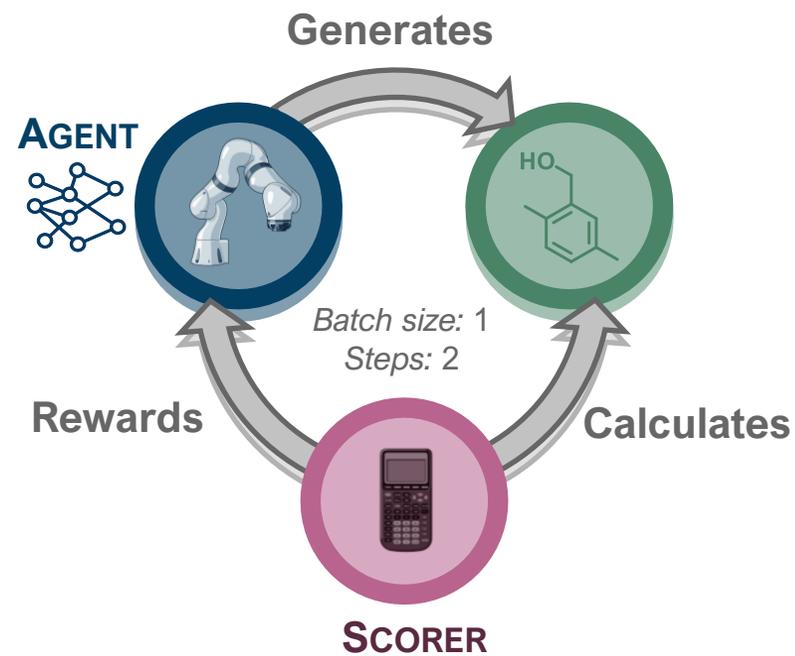
- $7 = \log P$



REINFORCEMENT LEARNING

SCORER

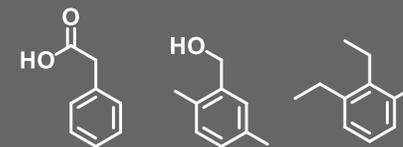
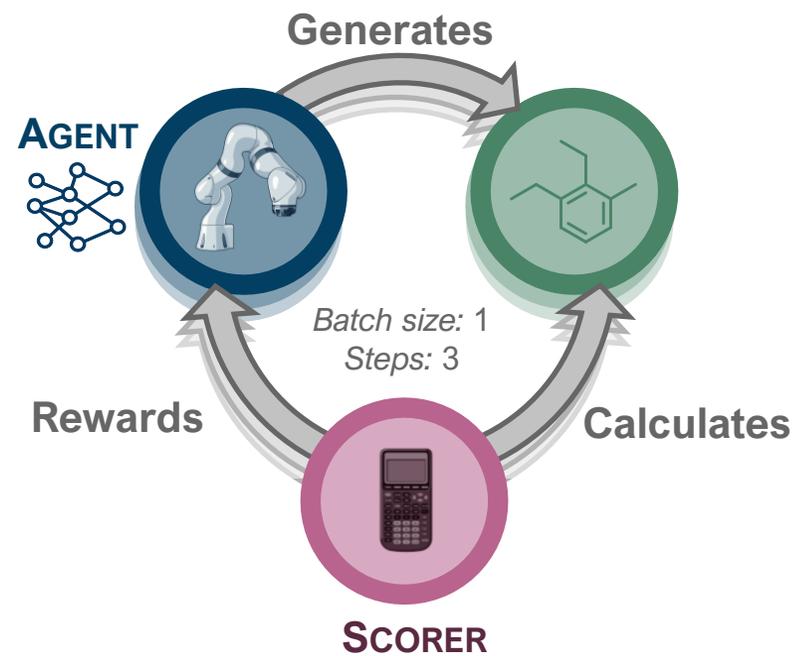
- $7 = \log P$



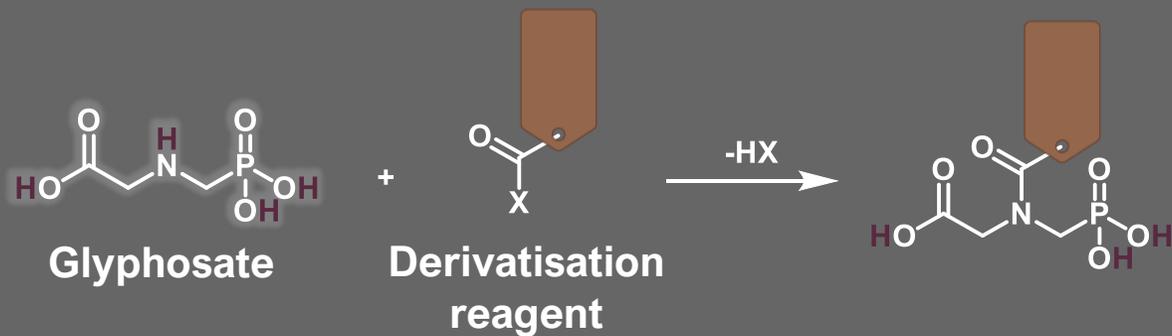
REINFORCEMENT LEARNING

SCORER

- $7 = \log P$



GENERATIVE MODELLING



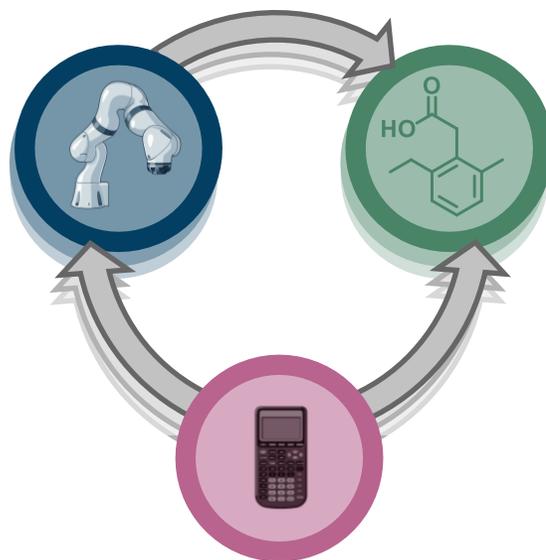
SCORER

$\log P$

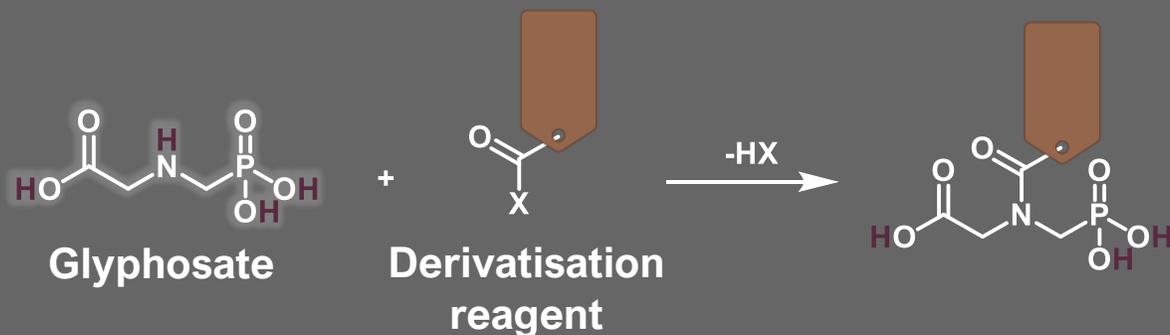
SAScore

$\log E$

-COOH group



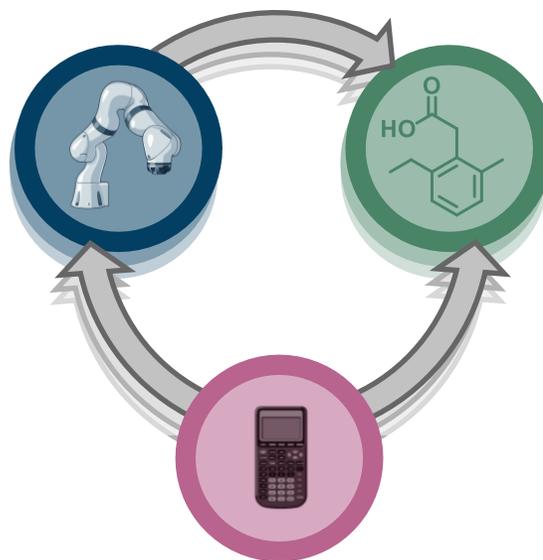
GENERATIVE MODELLING



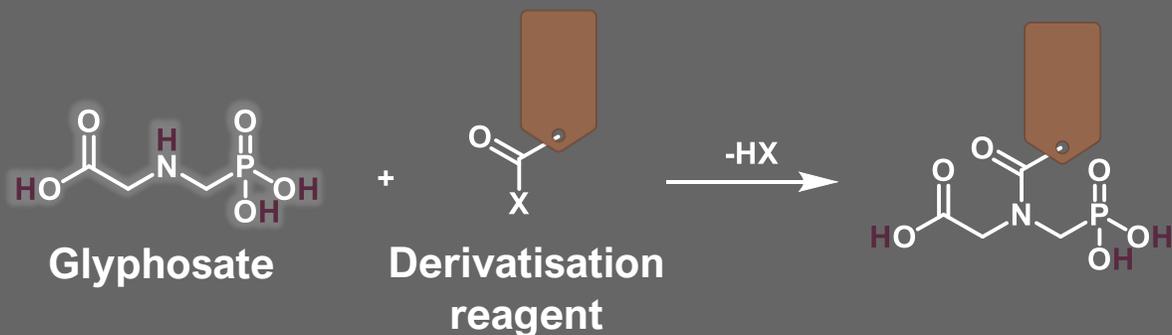
		Glyph	GlyphFMOC	GlyphACQ
predicted	log<i>E</i>	1.4	3.6	3.7
calculated	log<i>P</i>	-1.2	2.5	1.3

SCORER

- 1 – 3 **log*P***
- < 2.5 **SAScore**
- > 3.8 **log*E***
- 1 **-COOH group**
- no prim./sec. amine

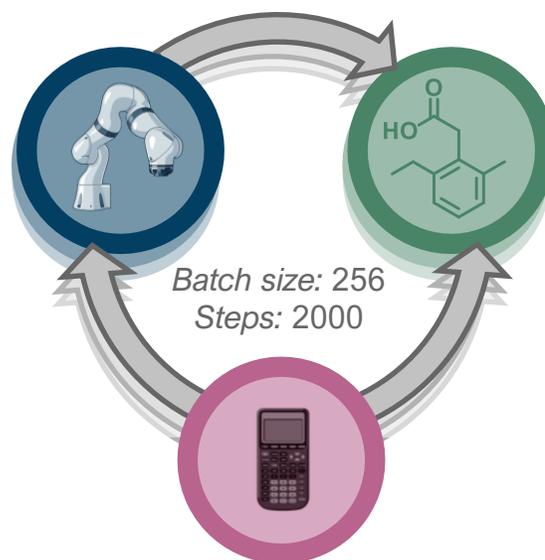


GENERATIVE MODELLING

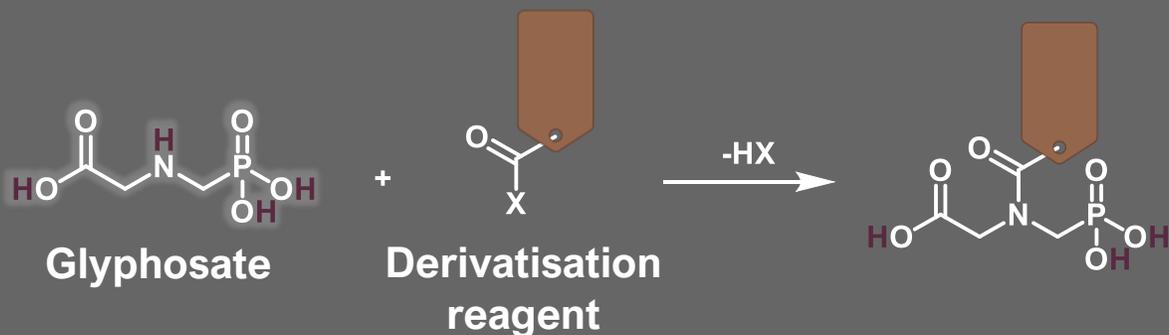


SCORER

- 1 – 3 **logP**
- < 2.5 **SAScore**
- > 3.8 **log/E**
- 1 **-COOH group**
- no prim./sec. amine



GENERATIVE MODELLING

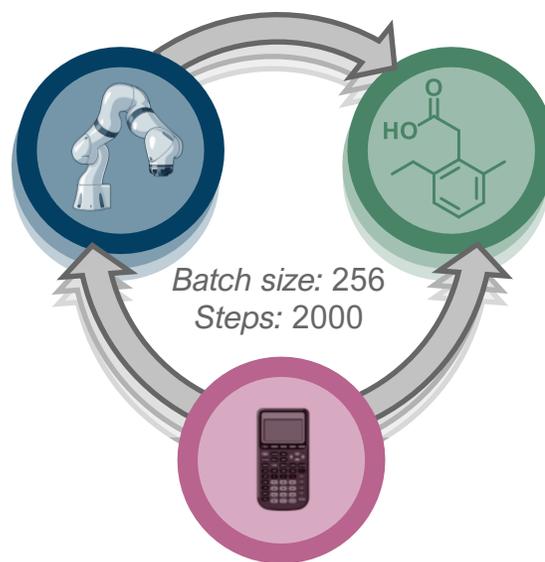


512,000 structures generated:

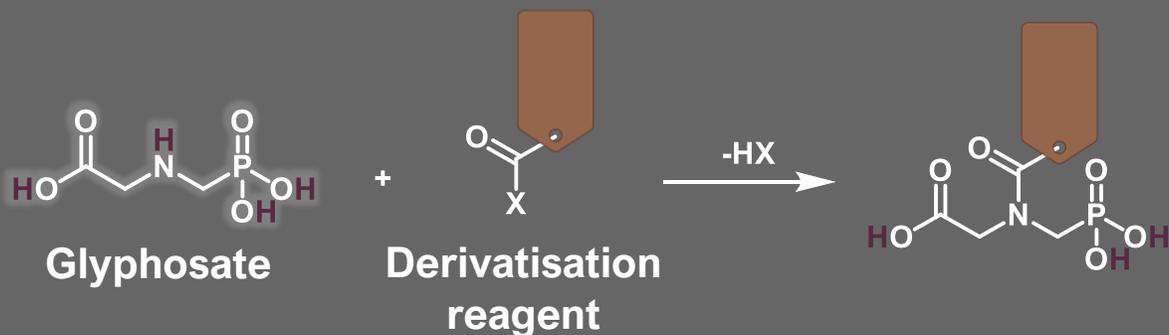
High-scoring candidates	SAS of top 1%	min SAS	log/E of top 1%	max log/E
36,312	1.49	1.26	4.09	4.24

SCORER

- 1 – 3 **logP**
- < 2.5 **SAScore**
- > 3.8 **log/E**
- 1 **-COOH group**
- no prim./sec. amine

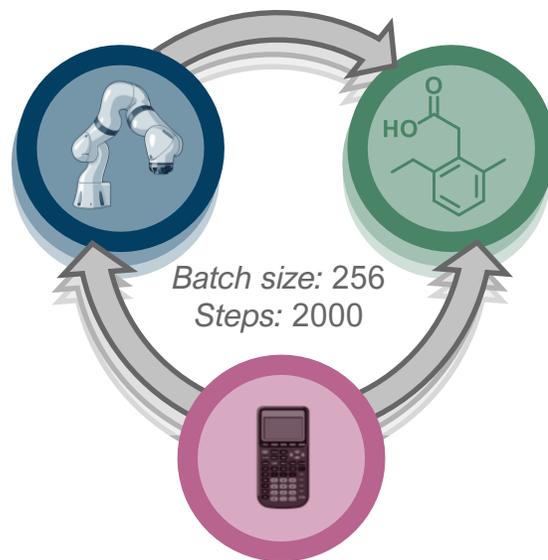


GENERATIVE MODELLING



SCORER

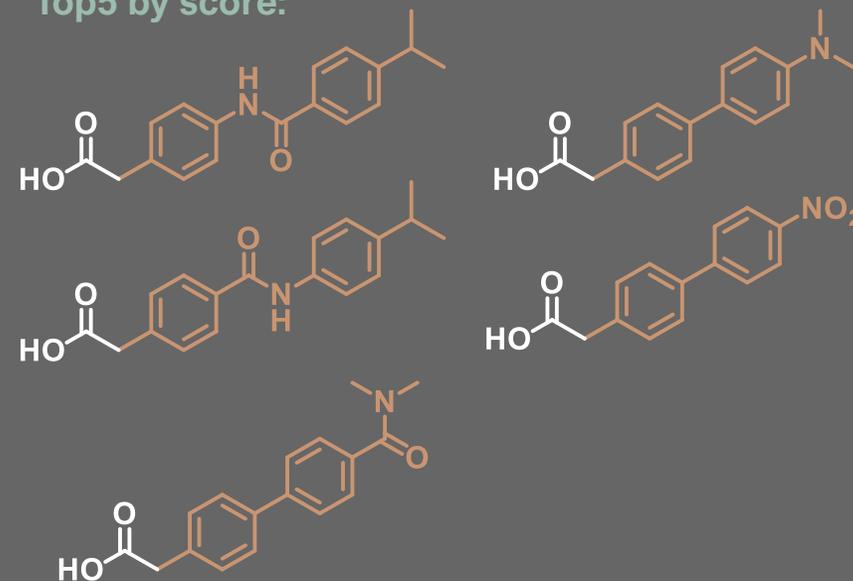
- 1 – 3 **logP**
- < 2.5 **SAScore**
- > 3.8 **log/E**
- 1 **-COOH group**
- no prim./sec. amine



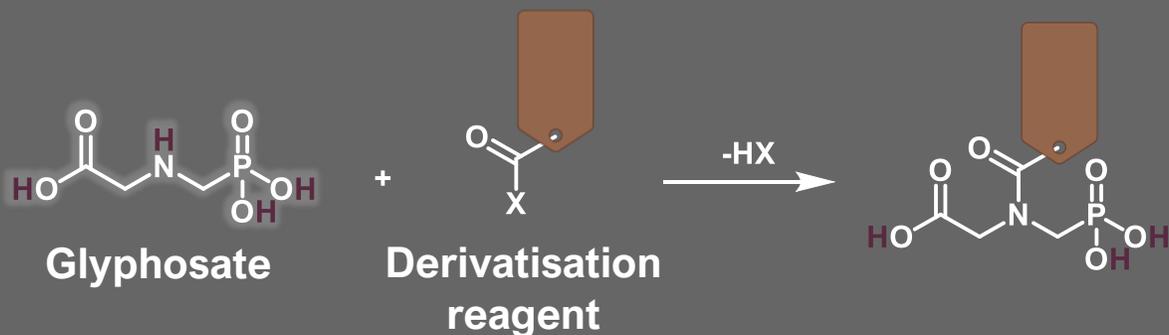
512,000 structures generated:

High-scoring candidates	SAS of top 1%	min SAS	log/E of top 1%	max log/E
36,312	1.49	1.26	4.09	4.24

Top5 by score:

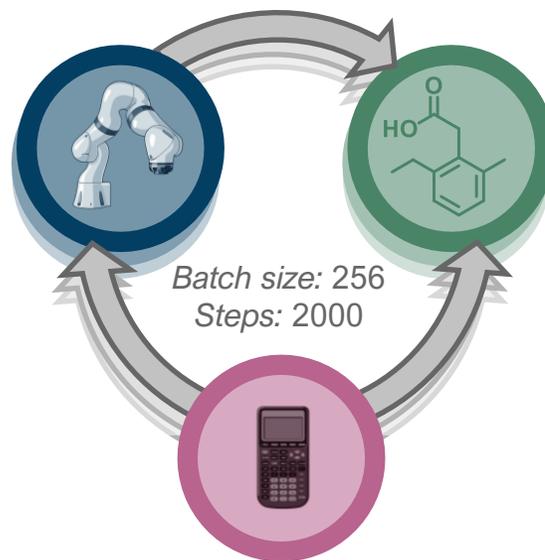


GENERATIVE MODELLING



SCORER

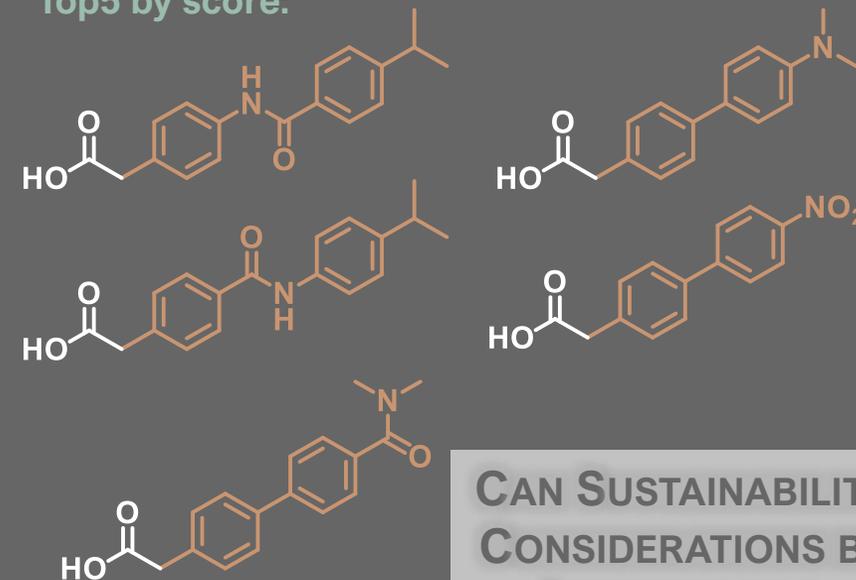
- 1 – 3 **logP**
- < 2.5 **SAScore**
- > 3.8 **log/E**
- 1 **-COOH group**
- no prim./sec. amine



512,000 structures generated:

High-scoring candidates	SAS of top 1%	min SAS	log/E of top 1%	max log/E
36,312	1.49	1.26	4.09	4.24

Top5 by score:



CAN SUSTAINABILITY CONSIDERATIONS BE INVOLVED IN A GENERATIVE MODELLING WORKFLOW?

SUSTAINABILITY CONSIDERATIONS

PRINCIPLES OF GREEN CHEMISTRY



Use of
renewable
feedstocks



Designing
safer
chemicals

SUSTAINABILITY CONSIDERATIONS

PRINCIPLES OF GREEN CHEMISTRY



Use of
renewable
feedstocks



Designing
safer
chemicals

EXPERIMENTAL FOUNDATION



Lignin
*Included in
synthesis*

Vanillin



Reagent synthesis
Easy and cheap

SAScore



Sample preparation
Low Hazard Suitable reactivity

Hazard score

-COOH group



Separation
*Suitable
retention time*

Detection
*High ionization
efficiency*

logP

log/E

SUSTAINABILITY CONSIDERATIONS

PRINCIPLES OF GREEN CHEMISTRY



Use of
renewable
feedstocks



Designing
safer
chemicals

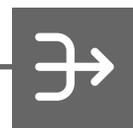


Predict



35 Hazard endpoints

Combined



Hazard score

EXPERIMENTAL FOUNDATION



Lignin
*Included in
synthesis*

Vanillin



Reagent synthesis
Easy and cheap

SAScore



Sample preparation
Low Hazard Suitable reactivity

Hazard score

-COOH group



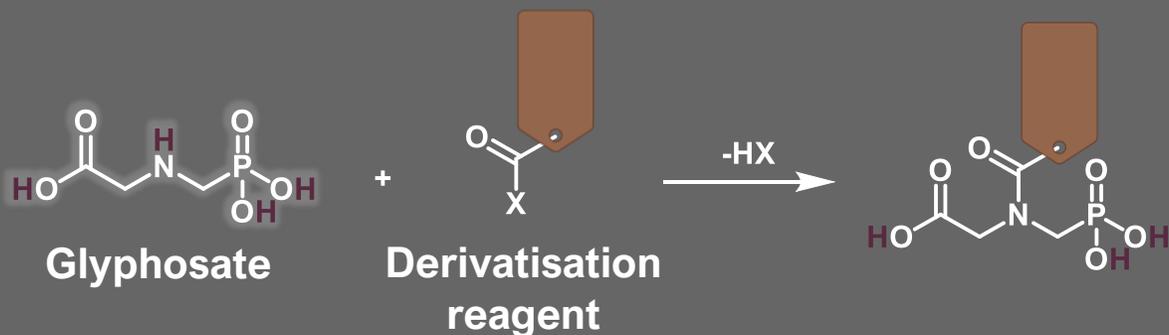
Separation
*Suitable
retention time*

Detection
*High ionization
efficiency*

logP

log/E

GENERATIVE MODELLING

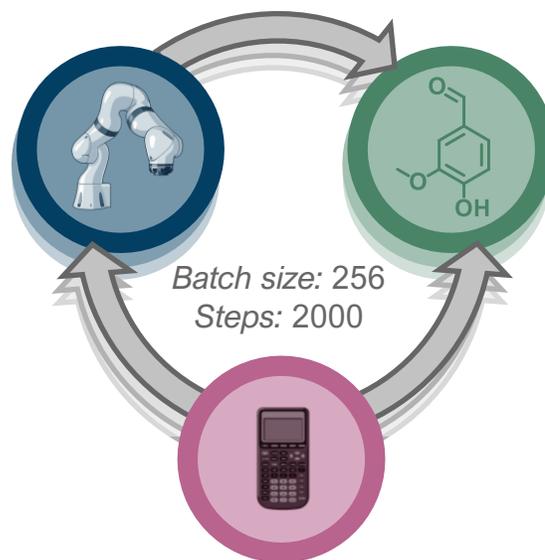


512,000 structures generated:

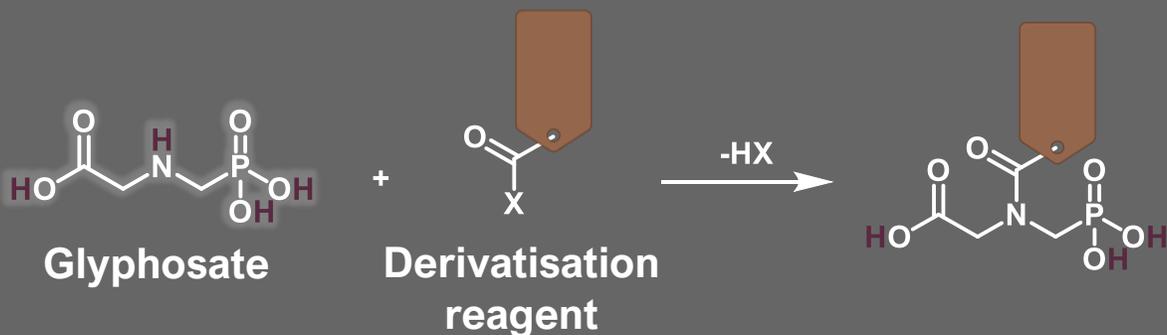
Vanillin	High-scoring candidates	SAS of top 1%	min SAS	log/E of top 1%	max log/E
No	36,312	1.49	1.26	4.09	4.24
Yes	29,120	1.71	1.51	4.20	4.29

SCORER

- 1 – 3 **logP**
- < 2.5 **SAScore**
- > 3.8 **log/E**
- 1 **-COOH group**
- no prim./sec. amine
- 1 **Vanillin**



GENERATIVE MODELLING

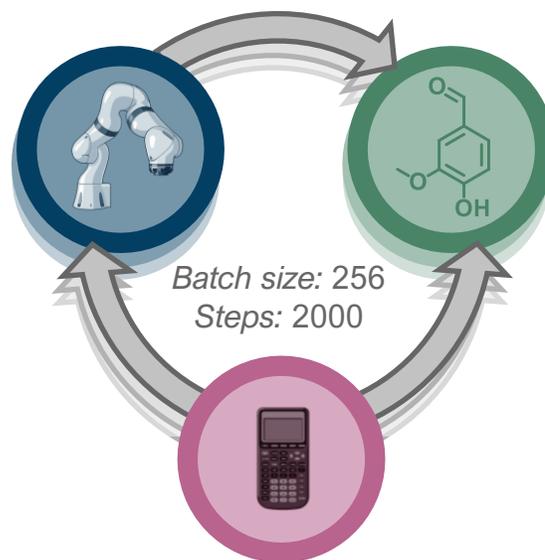


512,000 structures generated:

Vanillin	High-scoring candidates	SAS of top 1%	min SAS	log/E of top 1%	max log/E
Yes	29,120	1.49	1.26	4.09	4.24

SCORER

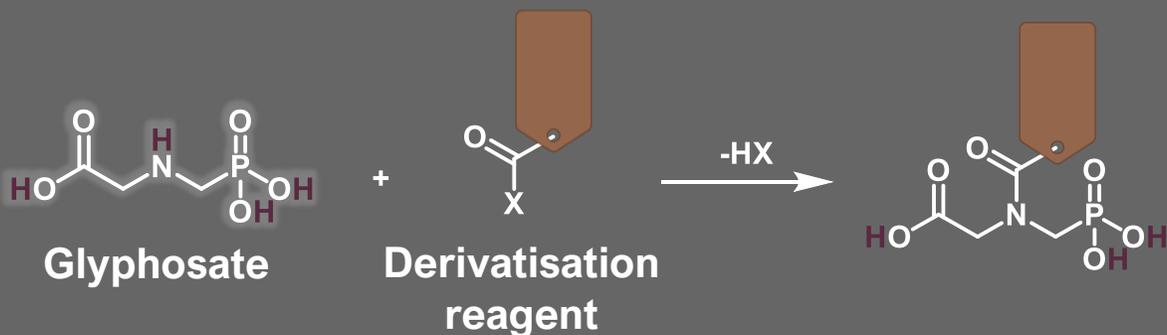
- 1 – 3 **logP**
- < 2.5 **SAScore**
- > 3.8 **log/E**
- 1 **-COOH group**
- no prim./sec. amine
- 1 **Vanillin**



FILTER

1. 10k highest score
2. 1k lowest **Hazard score**

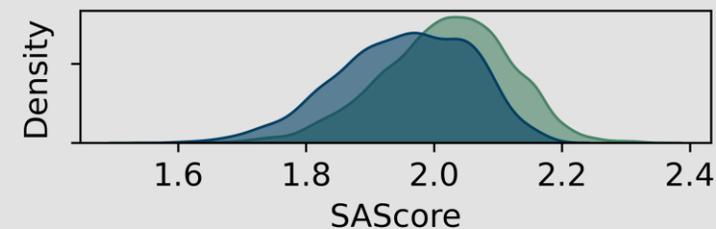
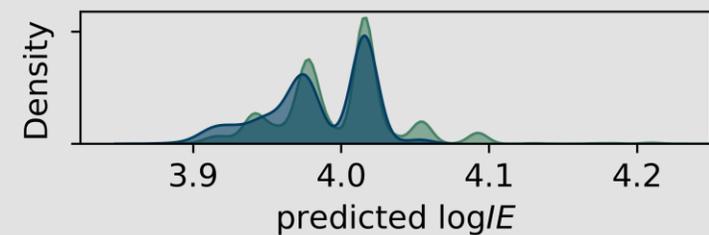
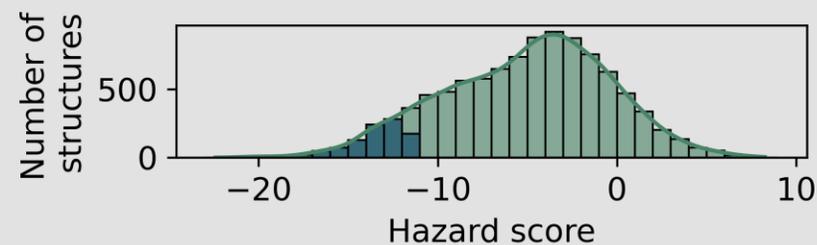
GENERATIVE MODELLING



512,000 structures generated:

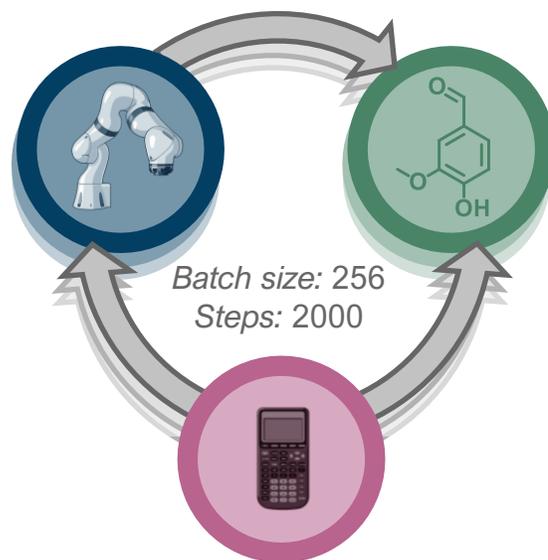
Vanillin	High-scoring candidates	SAS of top 1%	min SAS	log/E of top 1%	max log/E
Yes	29,120	1.49	1.26	4.09	4.24

Density plots **before** and **after** Hazard score filter:

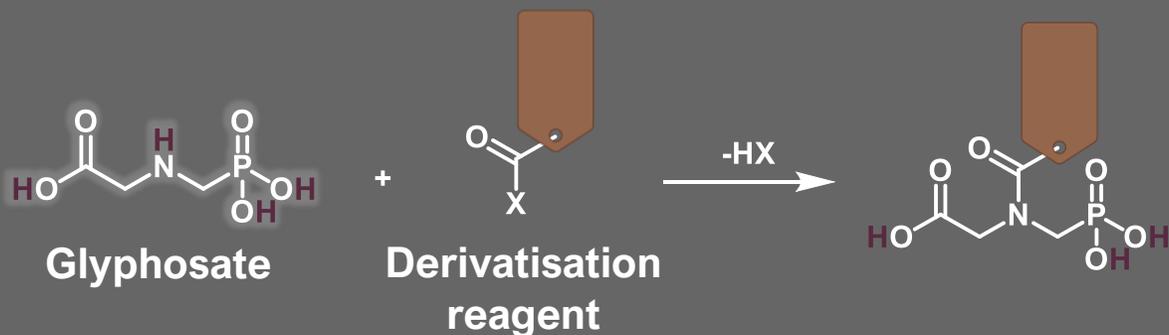


SCORER

- 1 – 3 **logP**
- < 2.5 **SAScore**
- > 3.8 **log/E**
- 1 **-COOH group**
- no prim./sec. amine
- 1 **Vanillin**

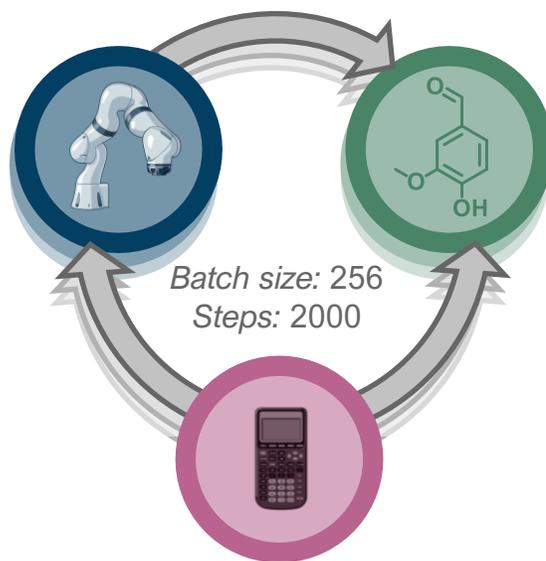


GENERATIVE MODELLING



SCORER

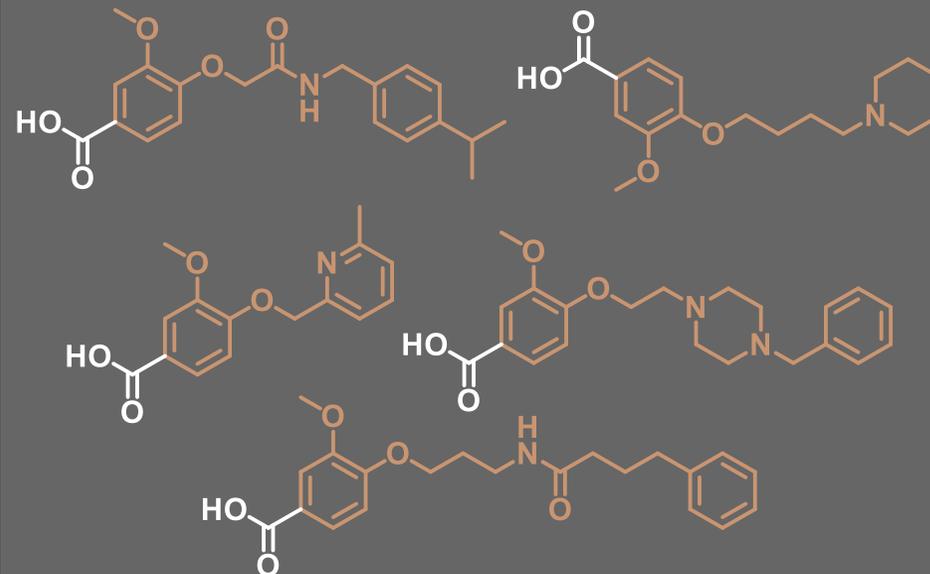
- 1 – 3 **logP**
- < 2.5 **SAScore**
- > 3.8 **log/E**
- 1 **-COOH group**
- no prim./sec. amine
- 1 **Vanillin**



512,000 structures generated:

Vanillin	High-scoring candidates	SAS of top 1%	min SAS	log/E of top 1%	max log/E
Yes	29,120	1.49	1.26	4.09	4.24

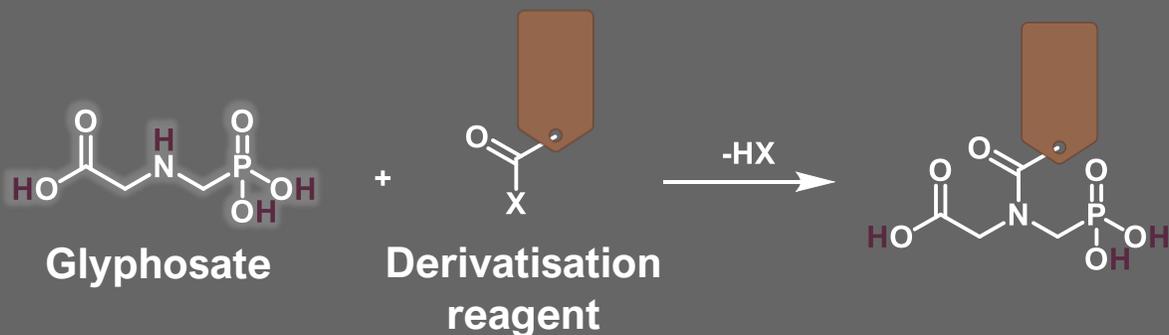
After filter:



FILTER

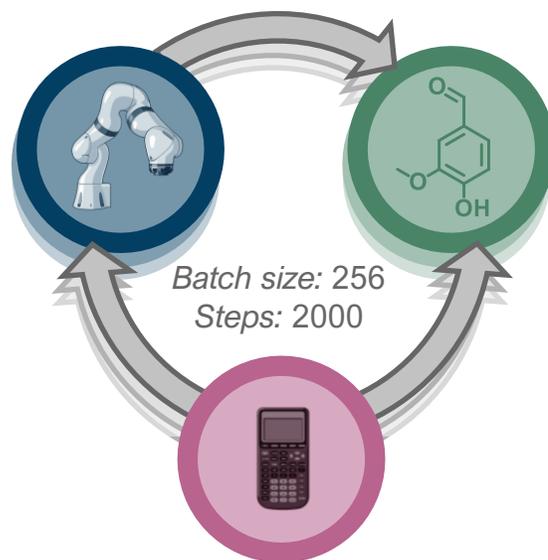
1. 10k highest score
2. 1k lowest **Hazard score**
3. 5 lowest SAScore and diverse

GENERATIVE MODELLING



SCORER

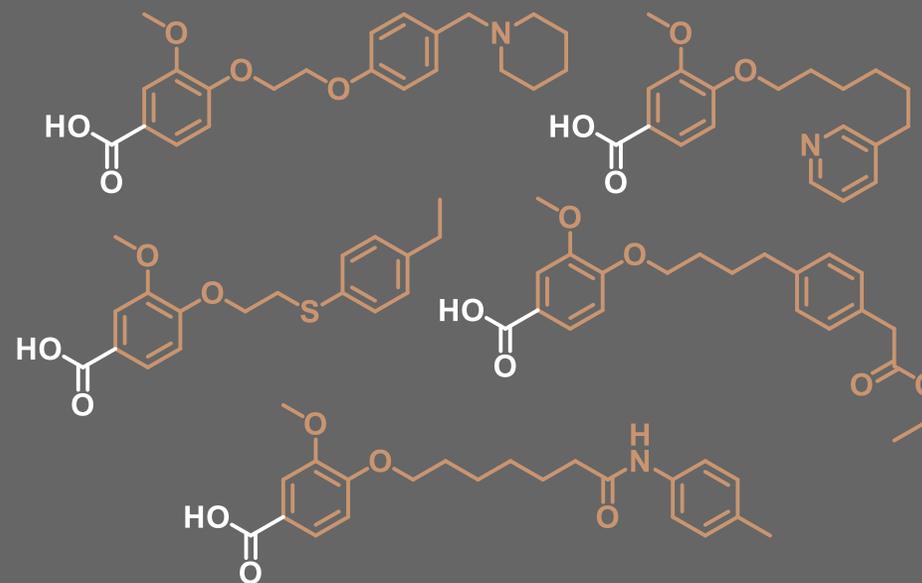
- 2.7 – 4.7 **logP**
- < 2.5 **SAScore**
- > 3.8 **log/E**
- 1 **-COOH group**
- no prim./sec. amine
- 1 **Vanillin**



512,000 structures generated:

Vanillin	High-scoring candidates	SAS of top 1%	min SAS	log/E of top 1%	max log/E
Yes	15,449	1.78	1.65	4.12	4.25

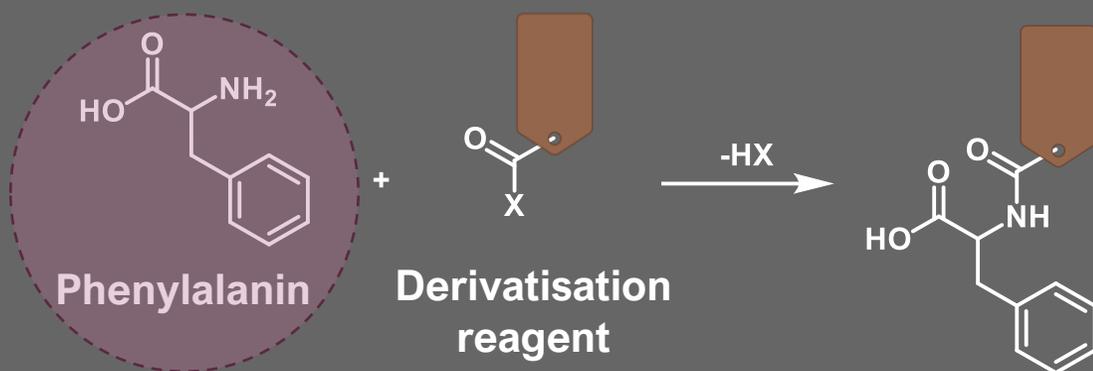
After filter:



FILTER

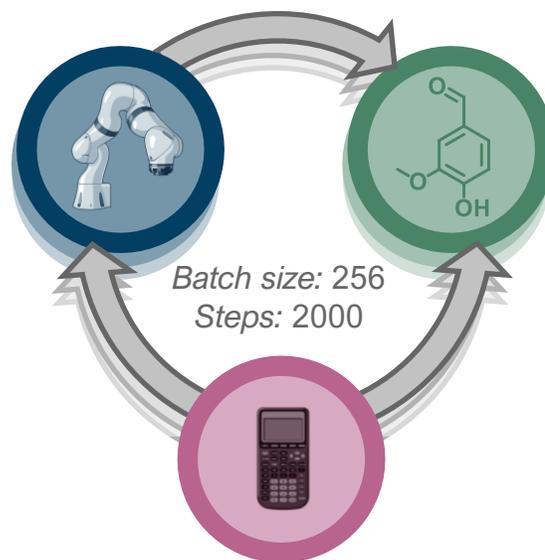
1. 10k highest score
2. 1k lowest **Hazard score**
3. 5 lowest SAScore and diverse

GENERATIVE MODELLING



SCORER

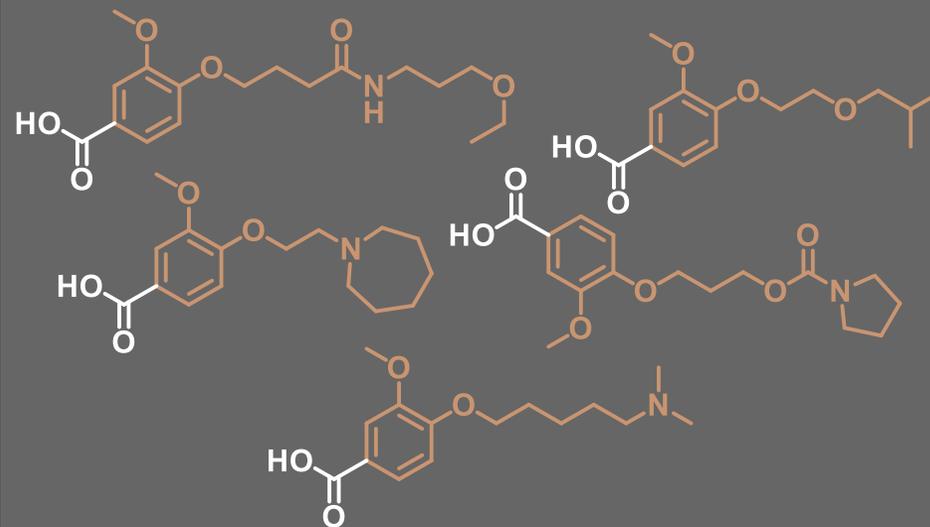
- 2.7 – 4.7 **logP**
- < 2.5 **SAScore**
- > 3.8 **log/E**
- 1 **-COOH group**
- no prim./sec. amine
- 1 **Vanillin**



512,000 structures generated:

Vanillin	High-scoring candidates	SAS of top 1%	min SAS	log/E of top 1%	max log/E
Yes	36,158	1.84	1.67	4.61	4.75

After filter:

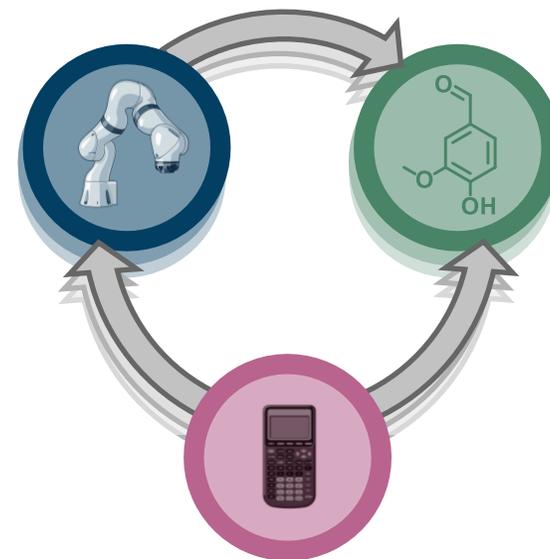


FILTER

1. 10k highest score
2. 1k lowest **Hazard score**
3. 5 lowest SAScore and diverse

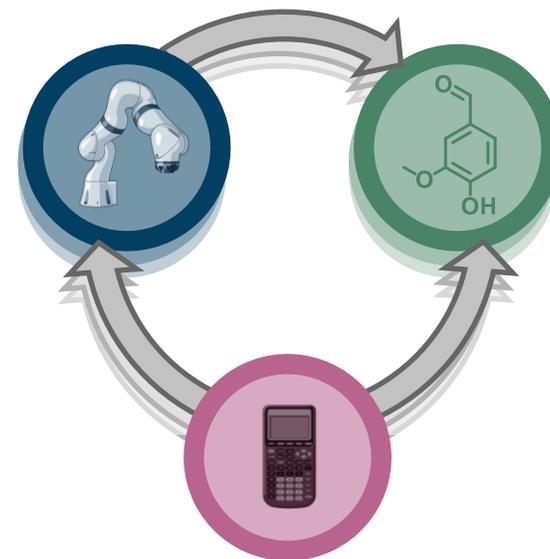
CONCLUSION

- **Machine learning workflow** generates promising candidates for **tailor-made derivatisation reagents** based on an **experimental foundation**.
- Workflow allows the incorporation of **sustainability considerations**.
- Filtering by **Hazard score** leads to a minor decrease in **analytical performance values**.
- **Adaptable** to different **analytical conditions** and analytes.

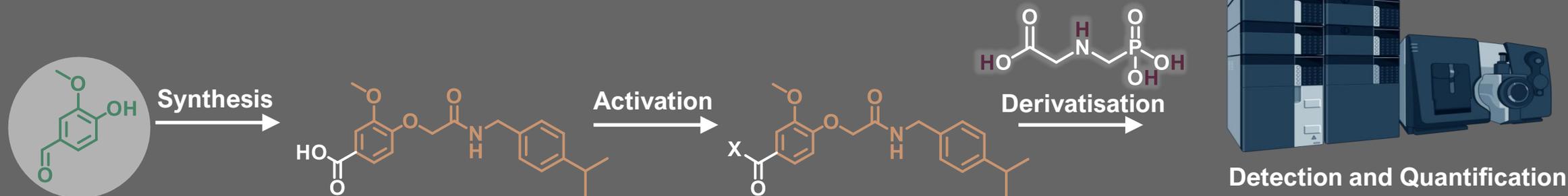


CONCLUSION

- **Machine learning workflow** generates promising candidates for **tailor-made derivatisation reagents** based on an **experimental foundation**.
- Workflow allows the incorporation of **sustainability considerations**.
- Filtering by **Hazard score** leads to a minor decrease in **analytical performance values**.
- **Adaptable** to different **analytical conditions** and analytes.



EXPERIMENTAL VERIFICATION ONGOING:



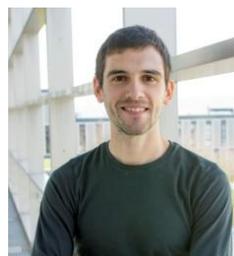
ACKNOWLEDGEMENTS



Anneli Kruve



Berit Olofsson Miguel Rivero-Crespo



Helen Sepman



Ida Rahu



Ziye Zhang



Stockholm
University

SUCCeSS



Kruve lab

JOIN TOMORROW:
SUCCeSS sessions
with interdisciplinary
topics!

Visit posters:



Louise Malm

NR. 65
*Implementing rules for
improved quantification
of transformation
products and
metabolites*



Iris Hättestrand

NR. 43
*Evaluation of
generative models for
structural elucidation of
chemicals from mass
spectrometry data*