# Prioritizing candidate structures in non-targeted LC/ESI/HRMS analysis by combining machine learning predictions

**Wei-Chieh (Harry) Wang**

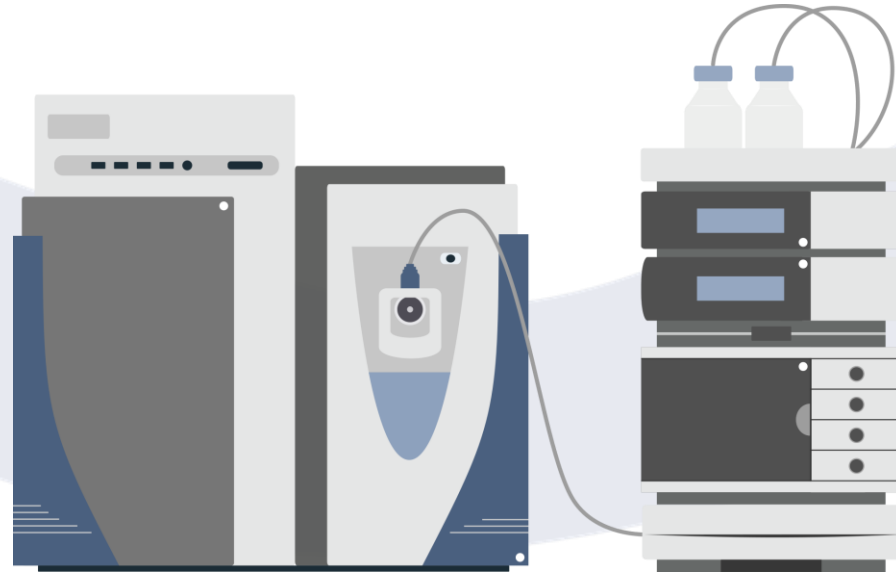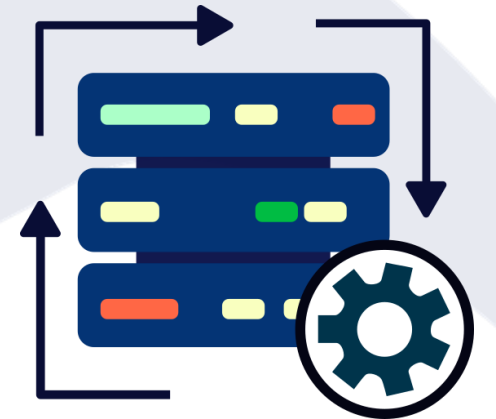**wei-chieh.wang@su.se**

**Stockholm University**

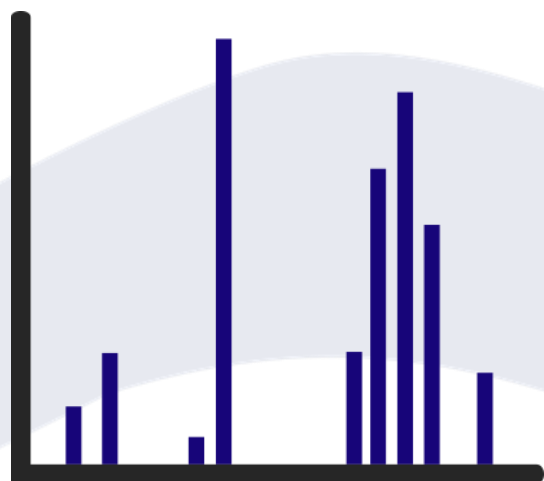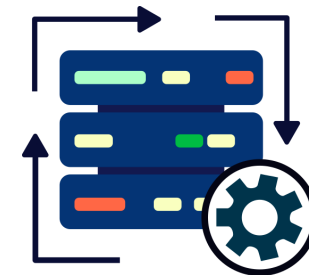# Non-targeted screening (NTS)



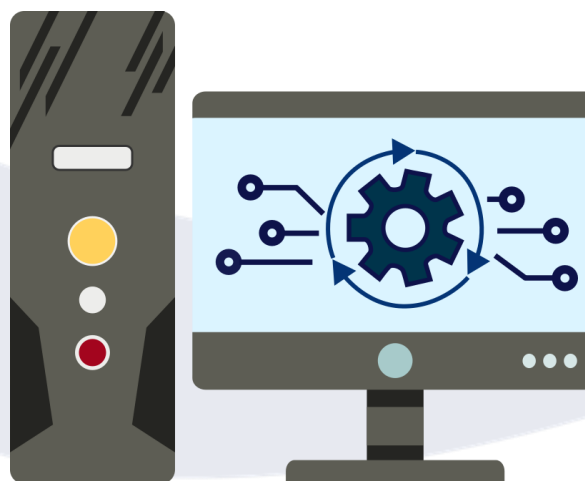**Sample preparation**

**LC/ESI/IM/HRMS
measurement**

**Data processing**

# Data processing

**MS spectra**

**Annotation**

**Candidate lists**

# Candidate validation
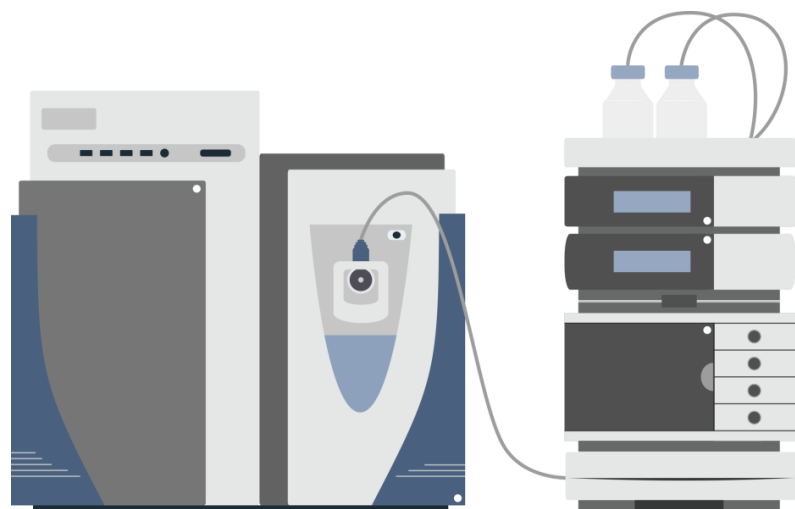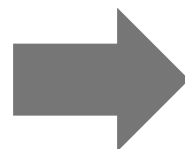
**Candidate validation**

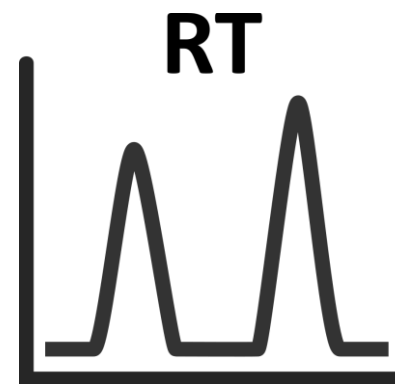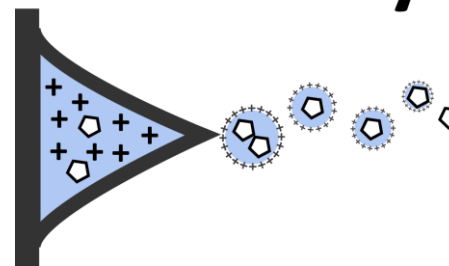**Lack of chemical standards**

**Costly**

**Time-consuming**

# Information from measurements

**CCS**

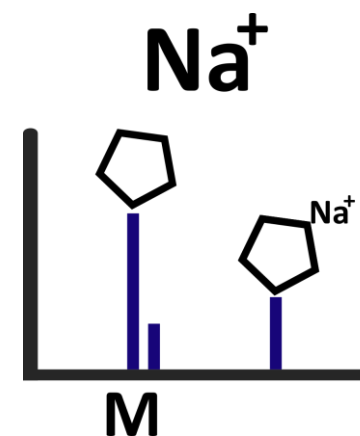**Ionizability**

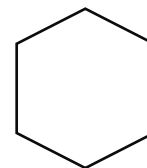**MS spectra**

**RT**

**Na⁺**

**M**

**LC/ESI/IM/HRMS measurement**
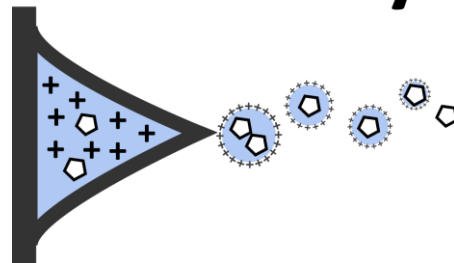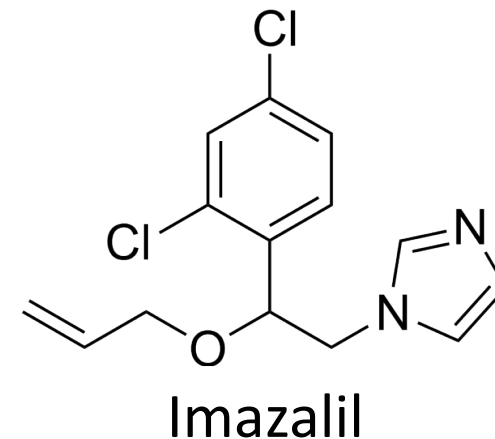
# Prioritization

**Ionizability**

Cyclohexane

Imazalil

**Properties of the candidates** ✖

✔ **Observed experimental properties**

Cyclohexane

**De-prioritized**

**Prioritized**

Imazalil

# Machine learning (ML) models

**Ionizability**

**Binary behavior**

**Continuous predictions**

**Threshold**

**RT**

**Continuous values**

**Continuous predictions**

**Error ranges (RMSE)**

# Uncertainty

- ## Model-based uncertainty
  - ### Model prediction errors (RMSE)


- ## Compound-based uncertainty
  - ### Conformal prediction system (CPS)

G. Shafer, V. Vovk. J. Machine Learning Research 9, 371-421 (2008).

# ML-supported prioritization
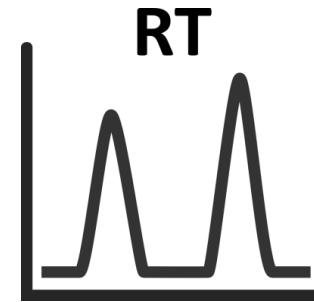


$$Accuracy = \frac{Number\ of\ true\ positive}{Numbre\ of\ the\ LC/HRMS\ features}$$

$$Efficiency = \frac{Number\ of\ the\ eliminated\ candidates}{Number\ of\ the\ total\ candidates}$$

A. Souihi, M.P. Mohai, E. Palm et al. J. Chroma 1666, 462867 (2022).

# Annotation performance from SIRIUS

Wei-Chieh (Harry) Wang

Dührkop, K., Fleischauer, M., Ludwig, M. et al. Nat Methods 16, 299–302 (2019).

# Results for prioritization

- RT prediction:
  - Efficiency: 84.02%
  - Accuracy: 13.64%

- Ionizability prediction
  - Recall-Efficiency curve

- Combining two models
  - Efficiency: 99.55%
  - Accuracy: 5.36 %

# Current challenges for combining predictions from various ML models

- Different application domains
  - The model was trained in different chemical spaces.



**Target set**

**Training set**

- No compound-based uncertainty available

# Conclusions & Future perspectives

- A strict combination of machine learning models led to an undesired removal of true positives.


- Incorporate additional machine learning prediction models.

- Retrain models using data from the same chemical space.

- Estimate compound-based uncertainty using a conformal prediction system.

# Acknowledgement

Kruve lab
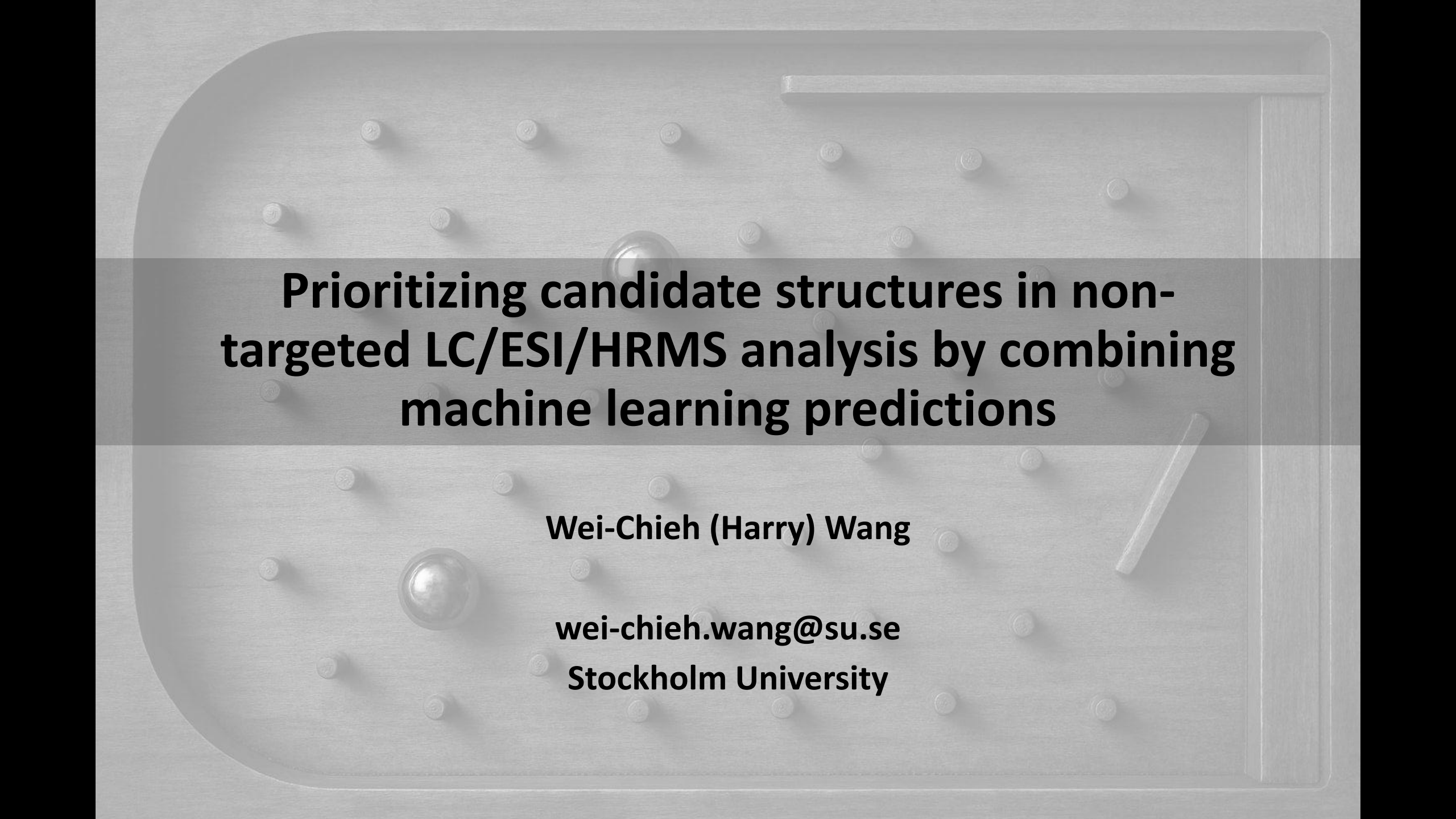

Stockholm University


European Research Council
Established by the European Commission


Swedish Research Council

# Prioritizing candidate structures in non-targeted LC/ESI/HRMS analysis by combining machine learning predictions

**Wei-Chieh (Harry) Wang**

**wei-chieh.wang@su.se**

**Stockholm University**