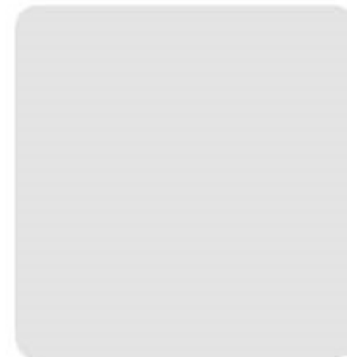# p$K_a$ predictions on top of the RDKit
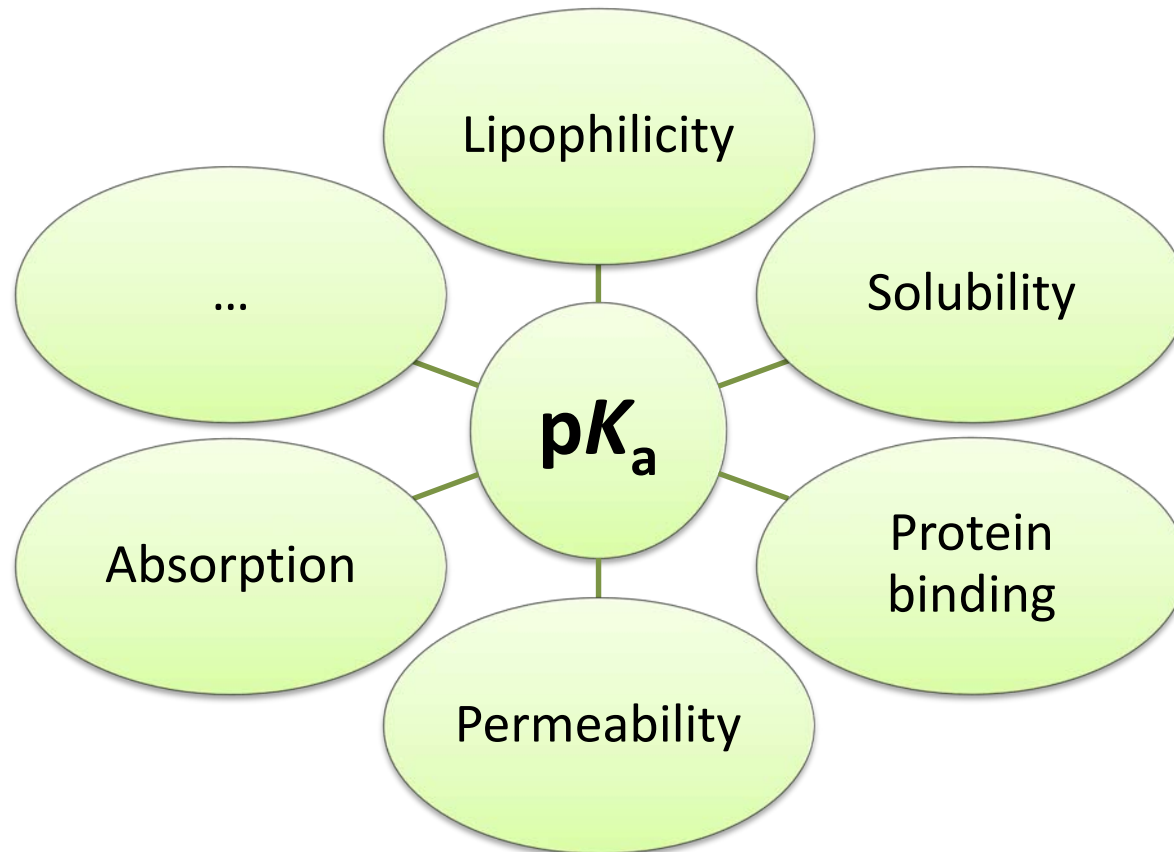
Marcel Baltruschat @CzodrowskiLab
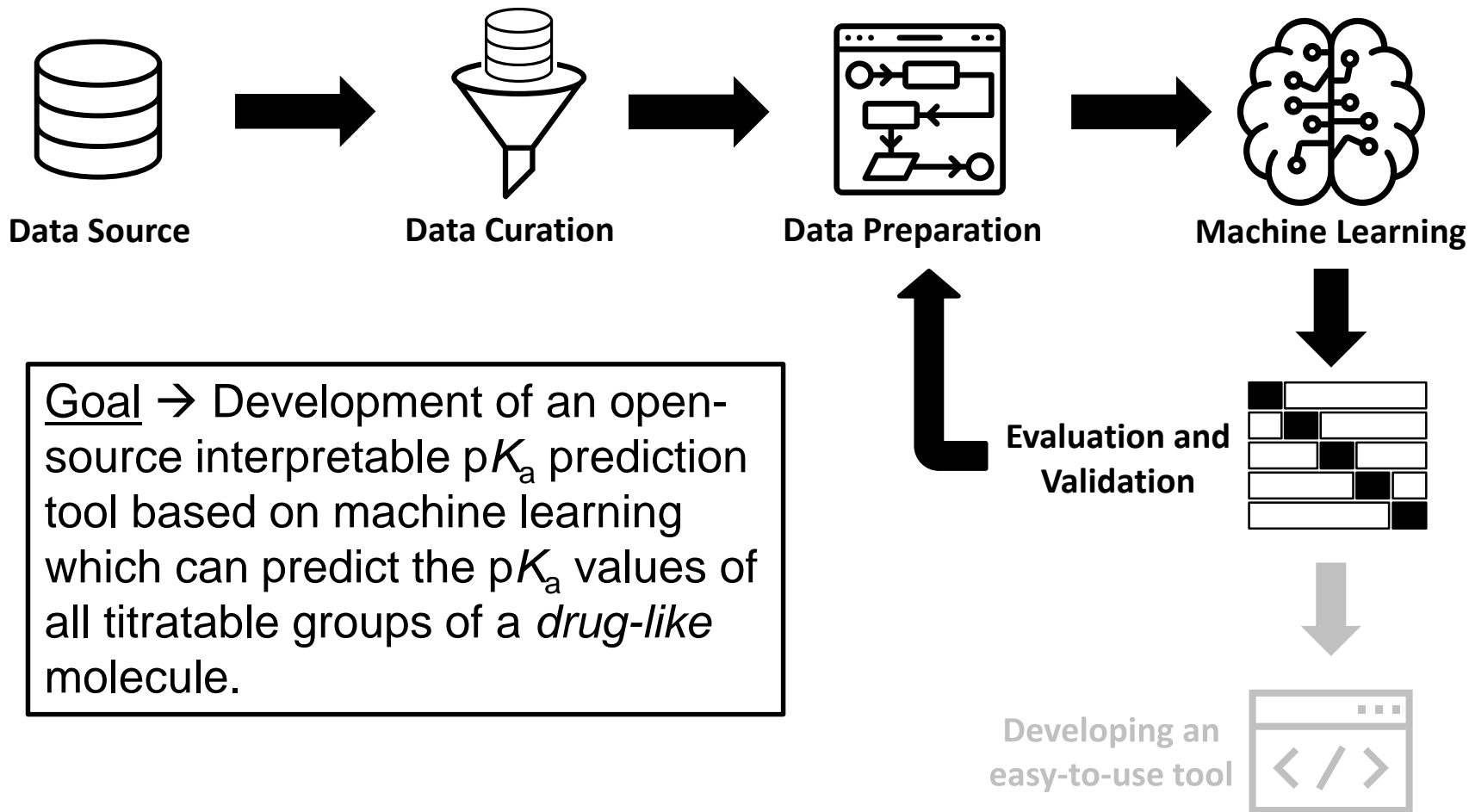
# Why p$K_a$?



- Latest published free p$K_a$ predictors can't reach the quality of commercial tools

- Free p$K_a$ predictors lack of features e.g. locating titratable groups

# Research Topic



**Data Source** → **Data Curation** → **Data Preparation** → **Machine Learning** → **Evaluation and Validation** → **Developing an easy-to-use tool**

> Goal → Development of an open-source interpretable p$K_a$ prediction tool based on machine learning which can predict the p$K_a$ values of all titratable groups of a *drug-like* molecule.
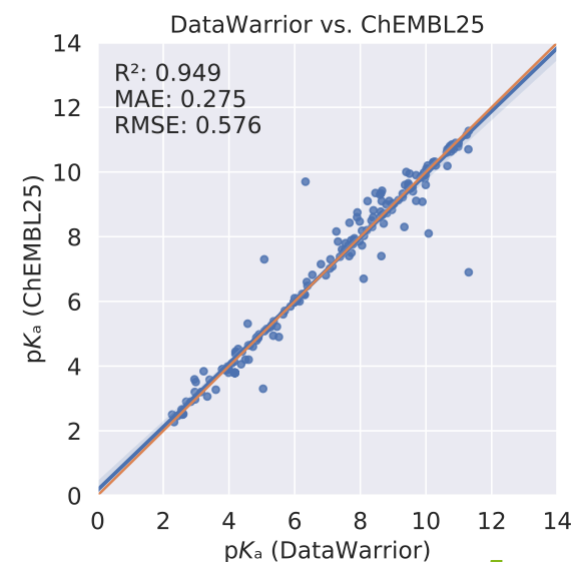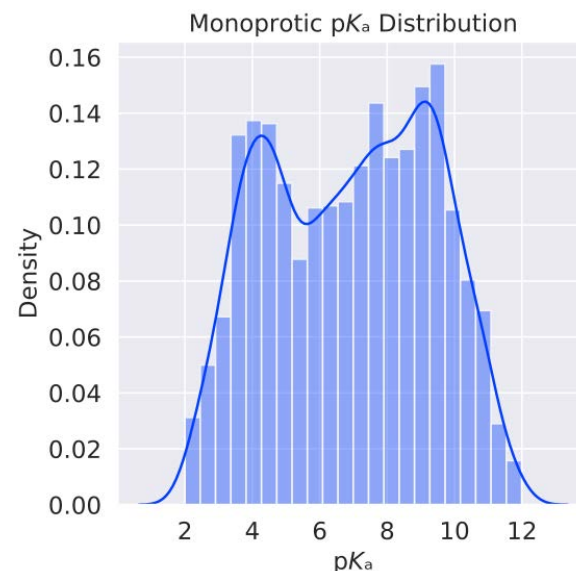
# Let's start with monoprotic molecules

Baltruschat M and Czodrowski P. Machine learning meets pKa [version 2; peer review: 2 approved].
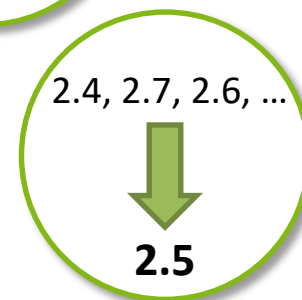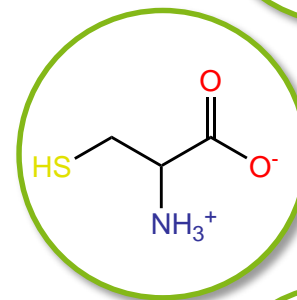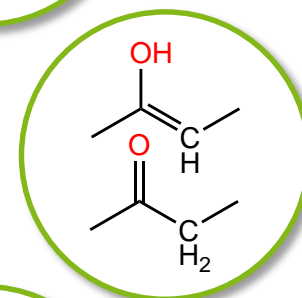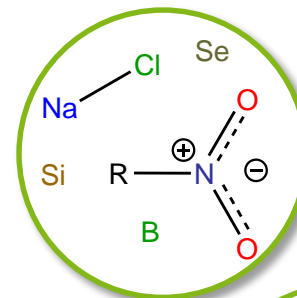*F1000Research* 2020, **9**(Chem Inf Sci):113 (https://doi.org/10.12688/f1000research.22090.2)

# Monoprotic Dataset



Monoprotic p$K_a$ Distribution

- Using the curated *ChEMBL25* and *DataWarrior* datasets

- **5994** curated unique monoprotic structures

- No source is specified for the values from *DataWarrior*

- *ChEMBL25* data points are completely taken from literature

- Good correlation of the intersection



DataWarrior vs. ChEMBL25

R²: 0.949
MAE: 0.275
RMSE: 0.576

# Data Curation

- Removal of salts, nitro groups, B, Se, Si

- Lipinski's rule of five (one violation allowed)

- $pK_a$ between 2 and 12

- Tautomer standardization

- Protonation at pH 7.4

- Combination of data points from duplicated structures while removing outliers

2.4, 2.7, 2.6, ...

**2.5**

# Machine Learning

## Algorithms

- Random Forest (1000 trees)
- SVR (gamma="auto"/"scale")
- Neural Network (MLP, 3 different architectures)
- XGB

## Training data

- 196/200 RDKit descriptors
- FeatureMorgan FP, radius 3, 4096 bits (FCFP6-like)
- Both combined
- Scaling for each of the three above

**42 model configurations**

**Evaluation through 5-fold cross validation and two external test sets**
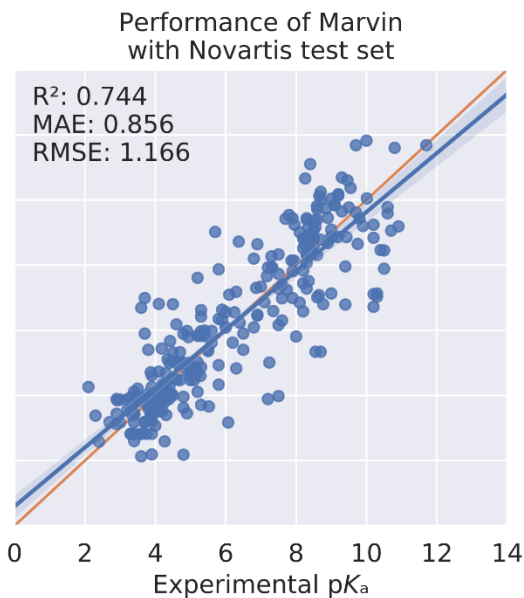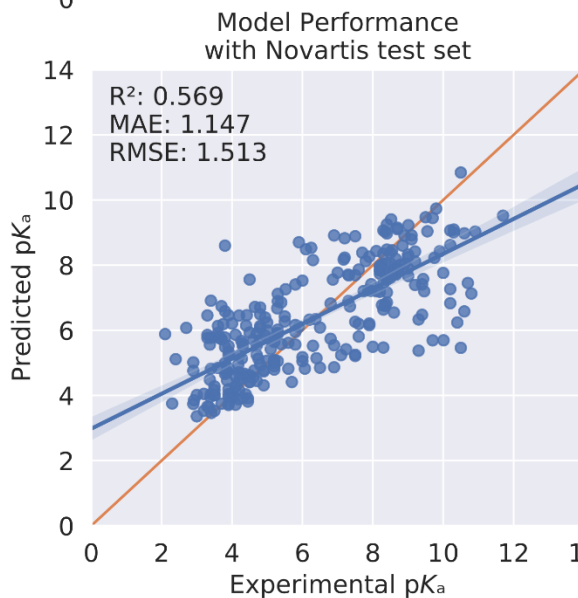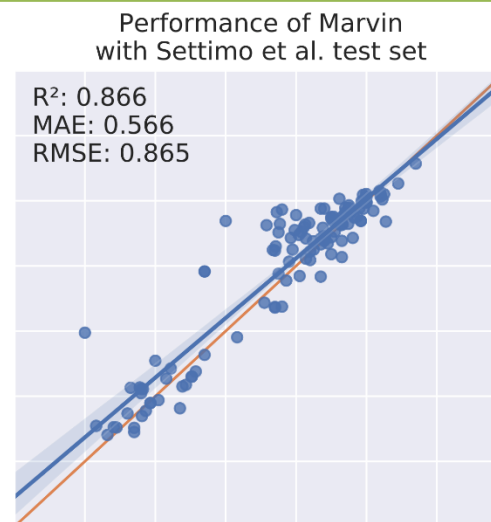
Settimo et. al.
123 mols

Novartis
280 mols

# Best Results

| No. | Model | Train data | MAE (CV) | RMSE (CV) | R² (CV) |
|---|---|---|---|---|---|
| #1 | | DESC+MF3 (scaled) | 0.682 | 1.032 | 0.820 |
| #2 | RF (1000 trees) | DESC+MF3 | 0.683 | 1.032 | 0.820 |
| #3 | | MF3 | 0.708 | 1.094 | 0.797 |
| #4 | | MF3 (scaled) | 0.708 | 1.094 | 0.797 |

# Best Results



Model Performance with Settimo et al. test set

R²: 0.889
MAE: 0.532
RMSE: 0.785

Performance of Marvin with Settimo et al. test set

R²: 0.866
MAE: 0.566
RMSE: 0.865

Model Performance with Novartis test set

R²: 0.569
MAE: 1.147
RMSE: 1.513
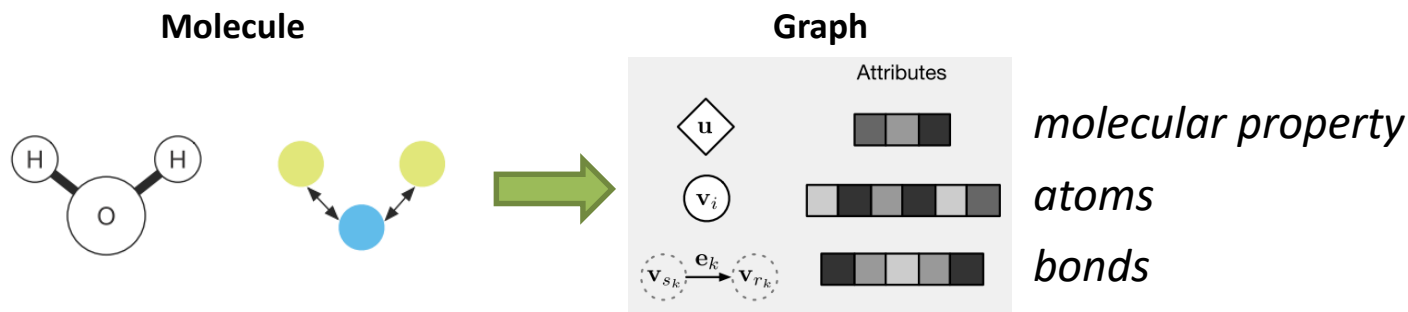
Performance of Marvin with Novartis test set

R²: 0.744
MAE: 0.856
RMSE: 1.166

# Graph Convolutional Networks (GCN) and QML

With David Bushiri and Prof. Dr. Enrico Tapavicza

- *PyTorch Geometric* module for GCNs



**Molecule**

**Graph**

Attributes

$u$ — *molecular property*

$v_i$ — *atoms*

$v_{s_k} \xrightarrow{e_k} v_{r_k}$ — *bonds*
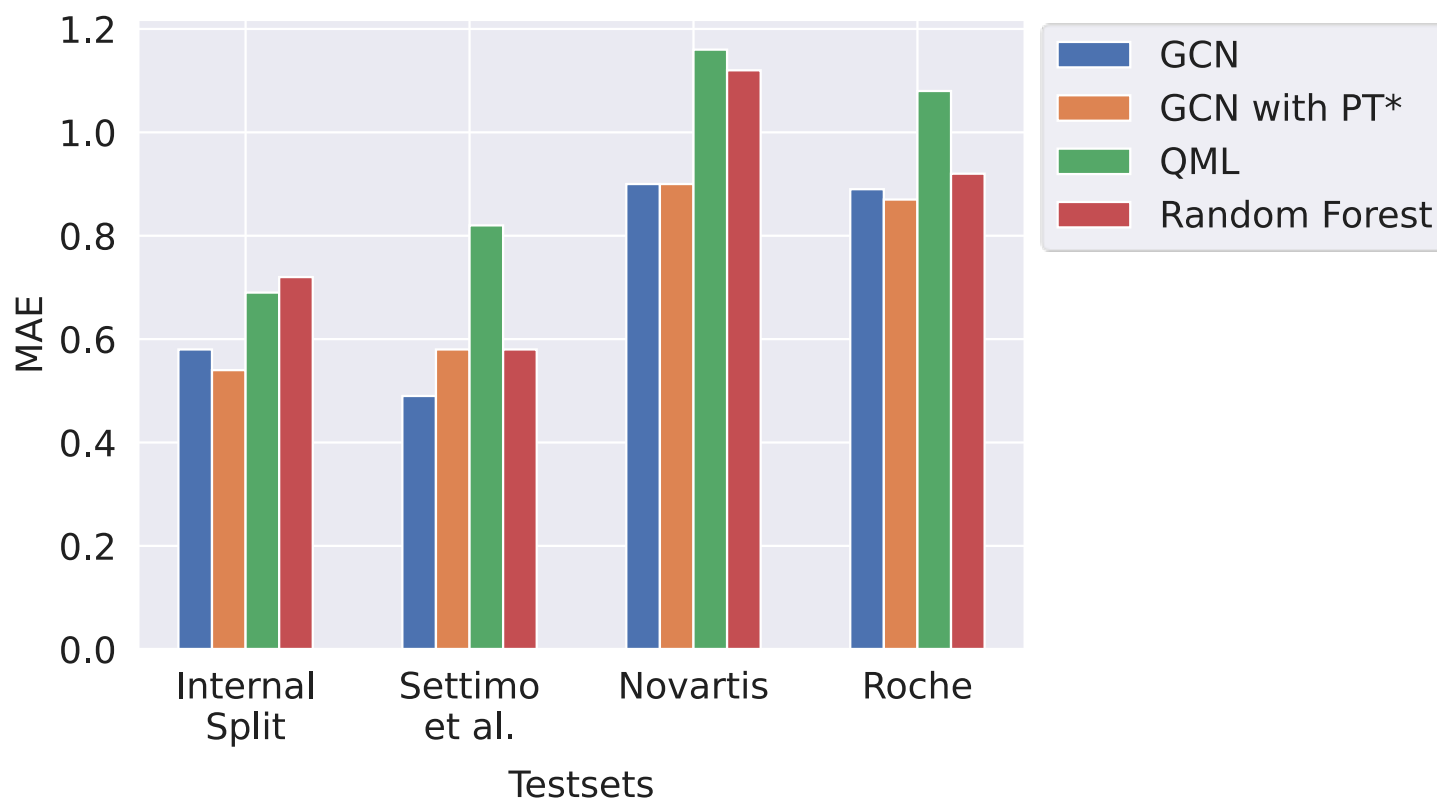
- *QML*: Kernel-ridge regression based method

➢ **Using 5196 DFT-optimized structures**

# Best GCN and QML Results



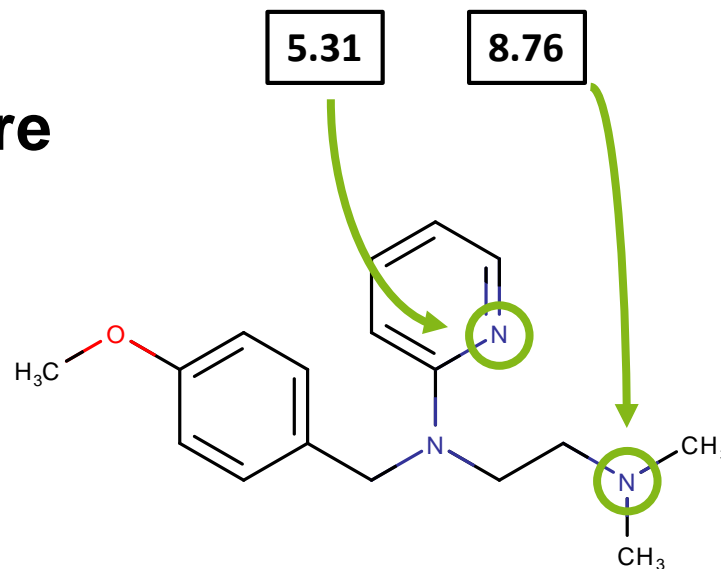*Pretrained with 900 000 protomers from the ZINC dataset with polar desolvation energy used
as target

# Let's go multiprotic!

# p$K_a$ Predictions of Multiprotic Molecules:
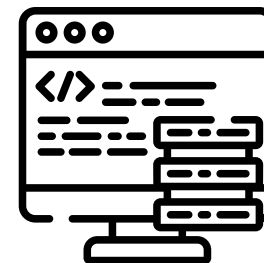## → The first two problems to solve

- **Identify and locate titratable groups without licensed software**
  - **Must be done for training and every prediction**

- Assign the p$K_a$ values from the datasets to the related titratable groups
  - Must be done only for training set

5.31    8.76

# The Idea

1. Identify the titratable groups with available tools

2. Generate a hardcoded list of SMARTS pattern from the tool results that covers all major groups

- Used Tools:

**Marvin Suite**

**Dimorphite-DL**

| # | SMARTS |
|---|--------|
| 1 | C#C |
| 2 | C(=O)O |
| 3 | *C(=O)[OH] |
| 4 | C(=O)[F,Cl,Br,I] |
| 5 | [#8X1] |
| 6 | [X3]=[!O] |
| 7 | [c] |

# Datasets

| Source | p$K_a$ Values | Unique Molecules |
| --- | --- | --- |
| ChEMBL26 | 8503 | 6617 |
| DataWarrior | 7911 | 7463 |
| Hunt et al. | 2488 | 2277 |
| Settimo et al. | 612 | 511 |
| Literature Compilation | 1765 | 1353 |
| SAMPL6 | 31 | 24 |
| Novartis | 1025 | 646 |
| Roche | 1762 | 1738 |
| OpenEye | 55322 | 23875 |
| **Total (curated)** | **49349** | **17538** |

# ChemAxon Marvin

Marvin provides atom ids for titratable groups during p$K_a$ calculation

➡ Extract environment around all atom ids

➡ Group by environment and count

# Dimorphite-DL

- Calculates all possible microstates of a molecule in a specified pH range

# Overview

- For both Marvin and Dimorphite-DL

  – Investigate environment
    distributions with radius 0 to 6

  – Looking at the saturation curves,
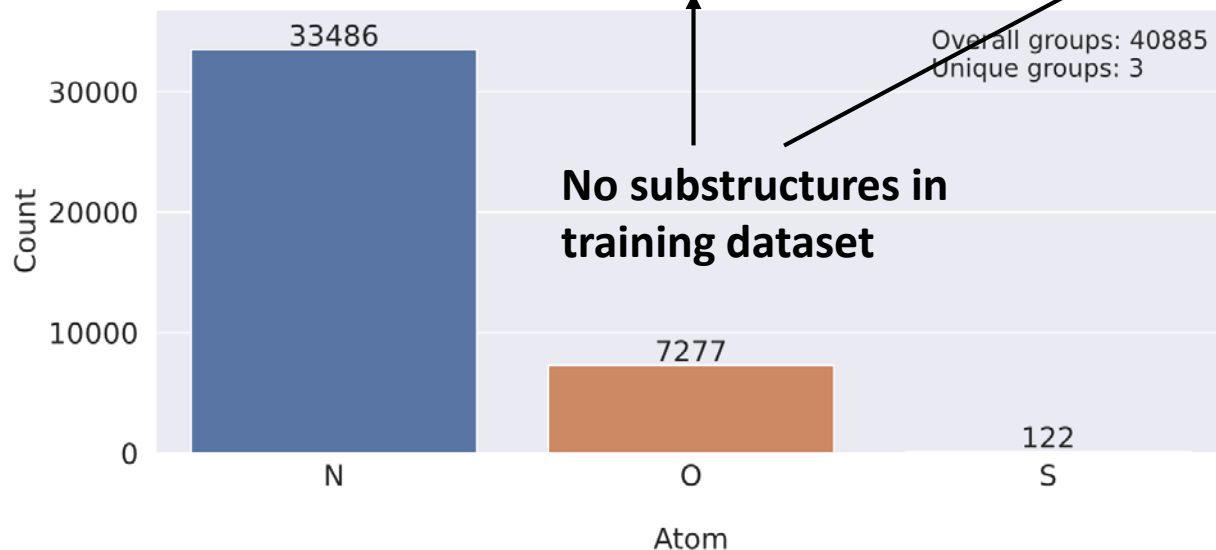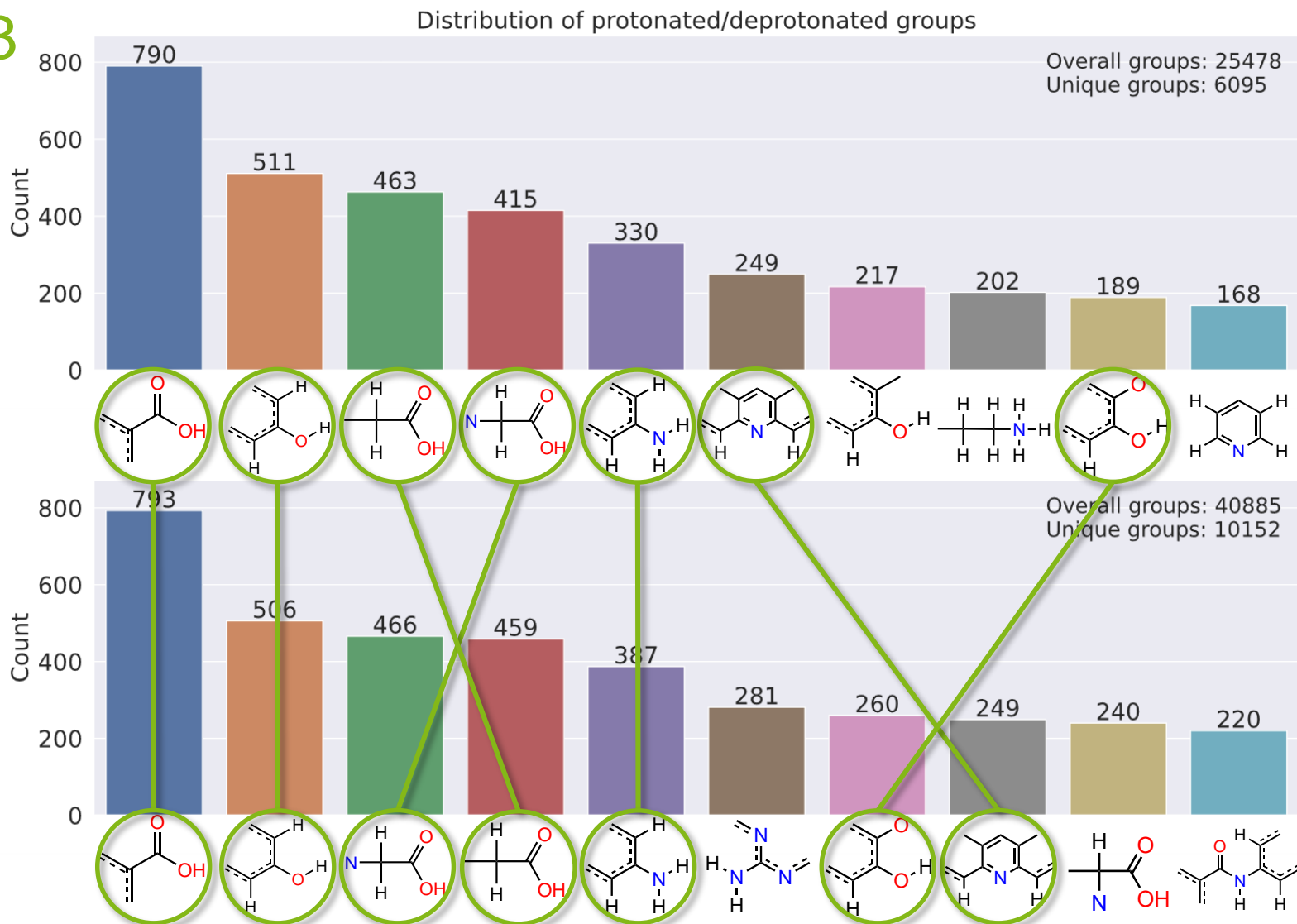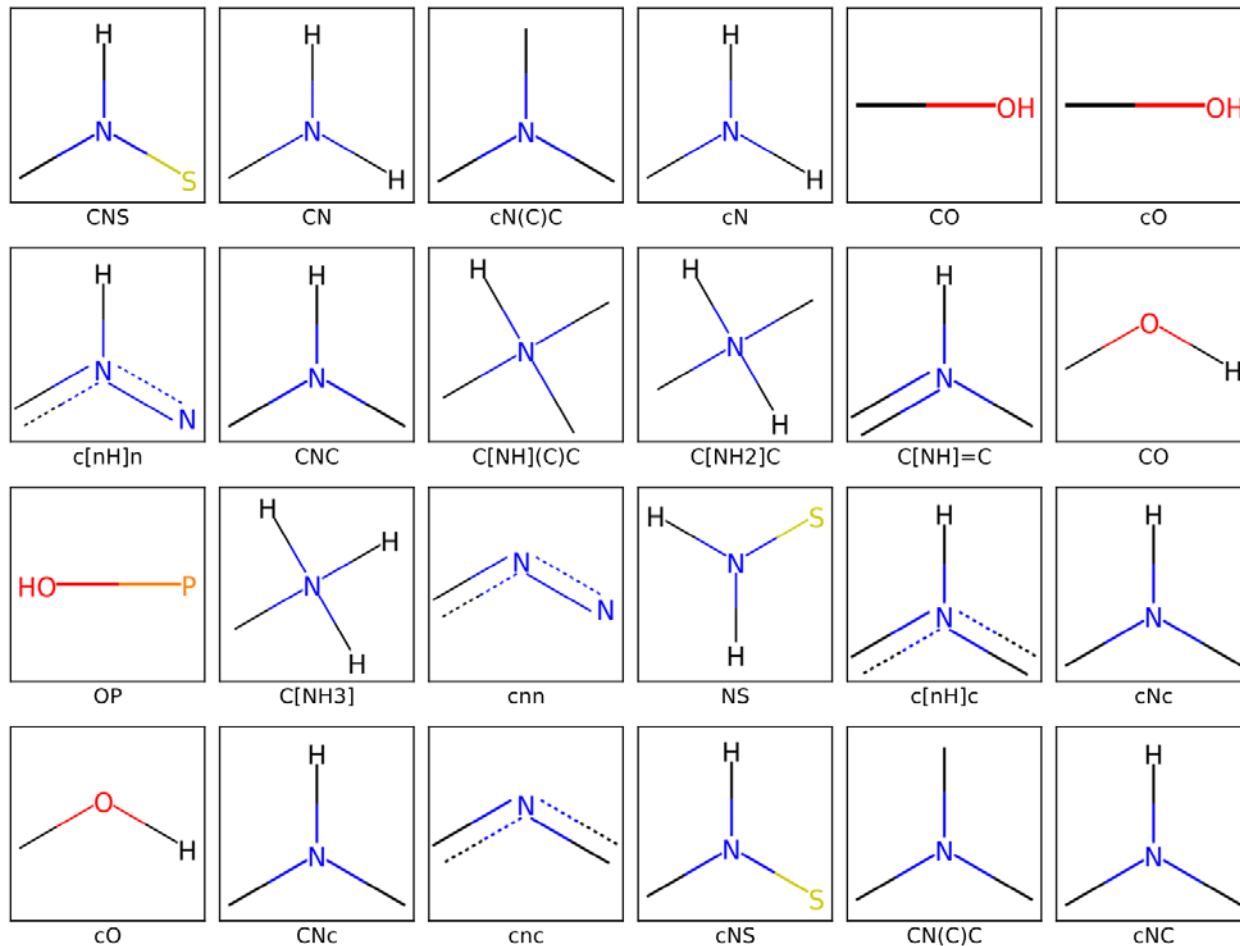    how many environments / groups do
    we *really* need?

# Radius 0

**ChemAxon Marvin**

**Dimorphite-DL**



Distribution of protonated/deprotonated atoms

Overall groups: 25478
Unique groups: 5

Overall groups: 40885
Unique groups: 3

**No substructures in training dataset**

Radius 3

Distribution of protonated/deprotonated groups

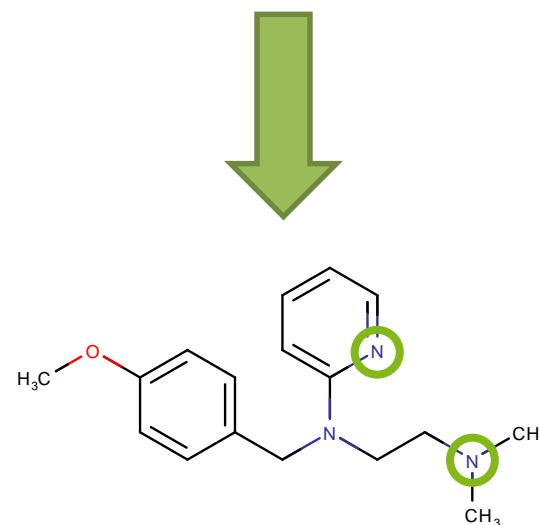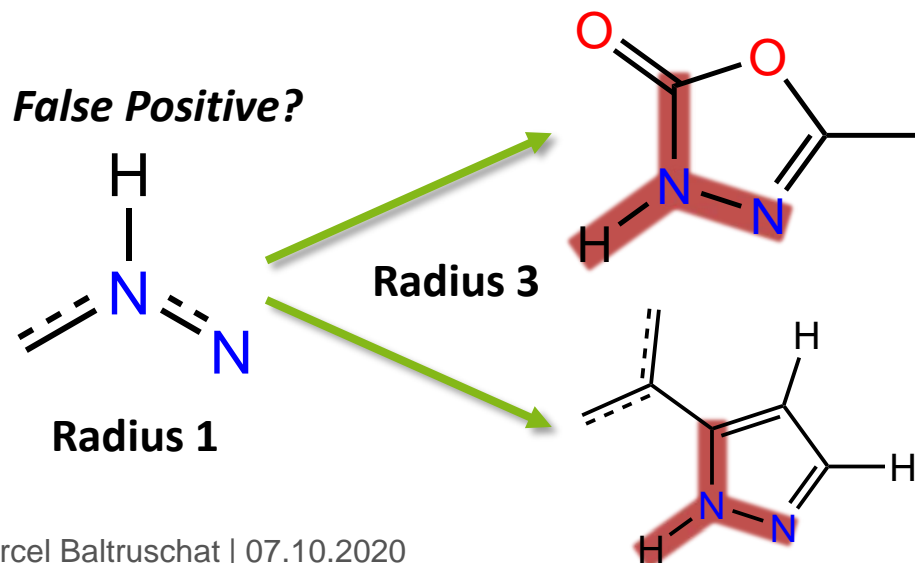# Saturation Curve
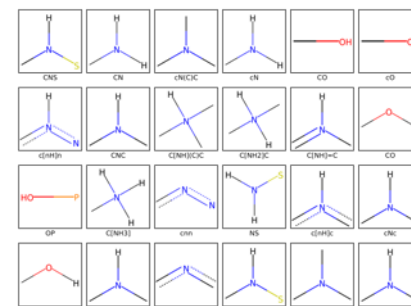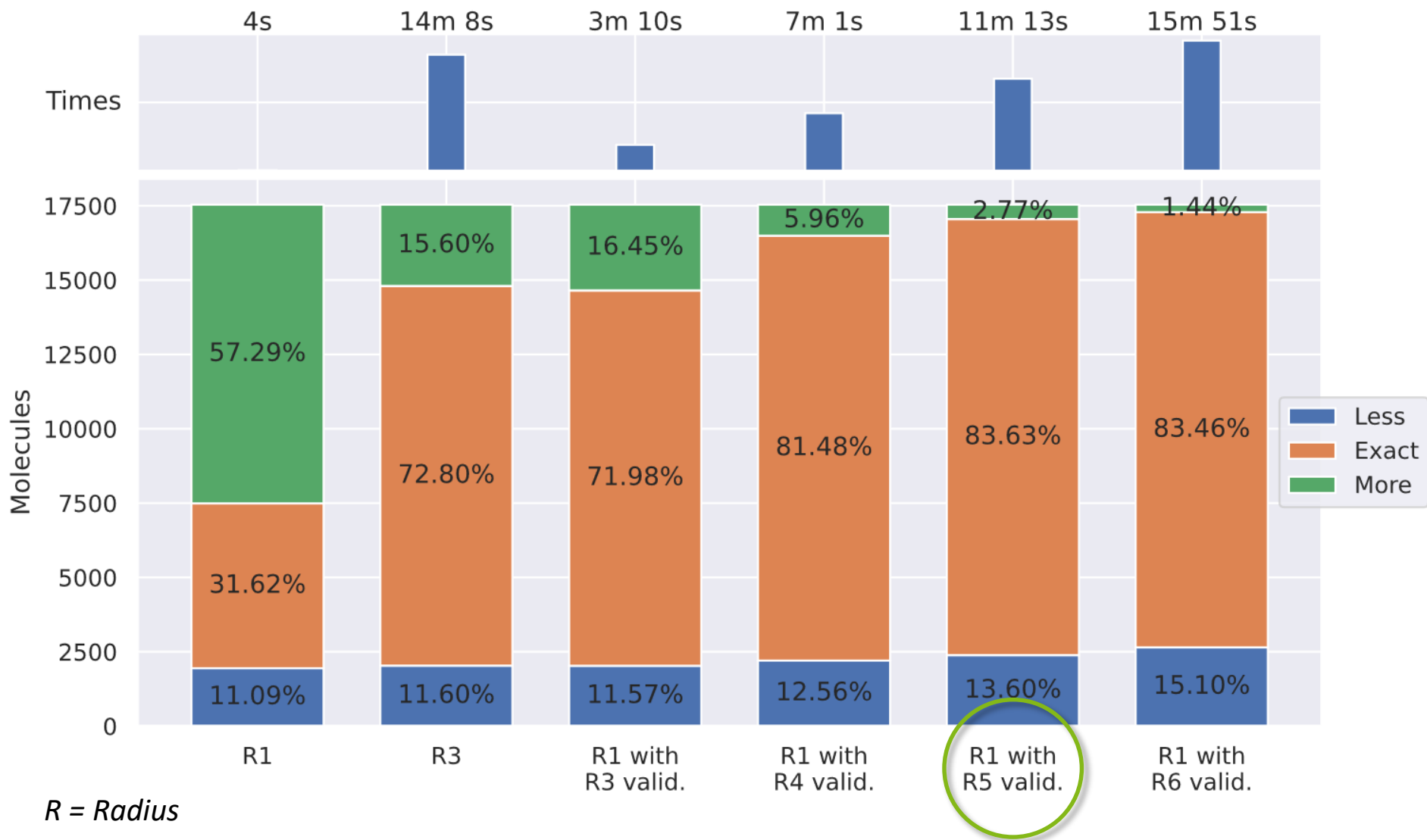
**ChemAxon Marvin**

**Dimorphite-DL**

# Result → 24 Titratable Fragments

# Validation

- Find the locations with the extracted titratable fragments
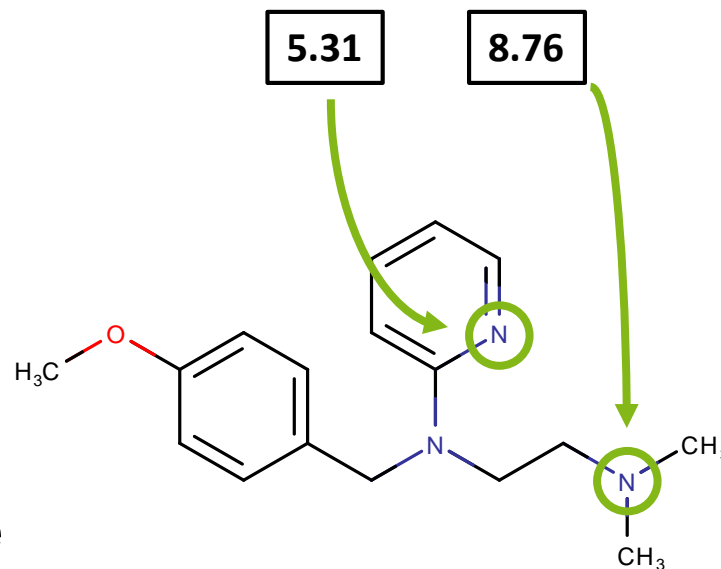
- Validate for all radii

- Test a hierarchical structure

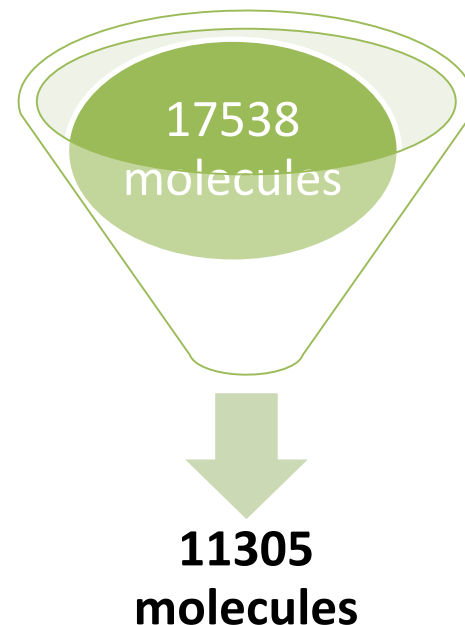# p$K_a$ Predictions of Multiprotic Molecules:
## → The first two problems to solve

- Identify and locate titratable groups without licensed software
  - Must be done for training and every prediction

- **Assign the p$K_a$ values from the datasets to the related titratable groups**
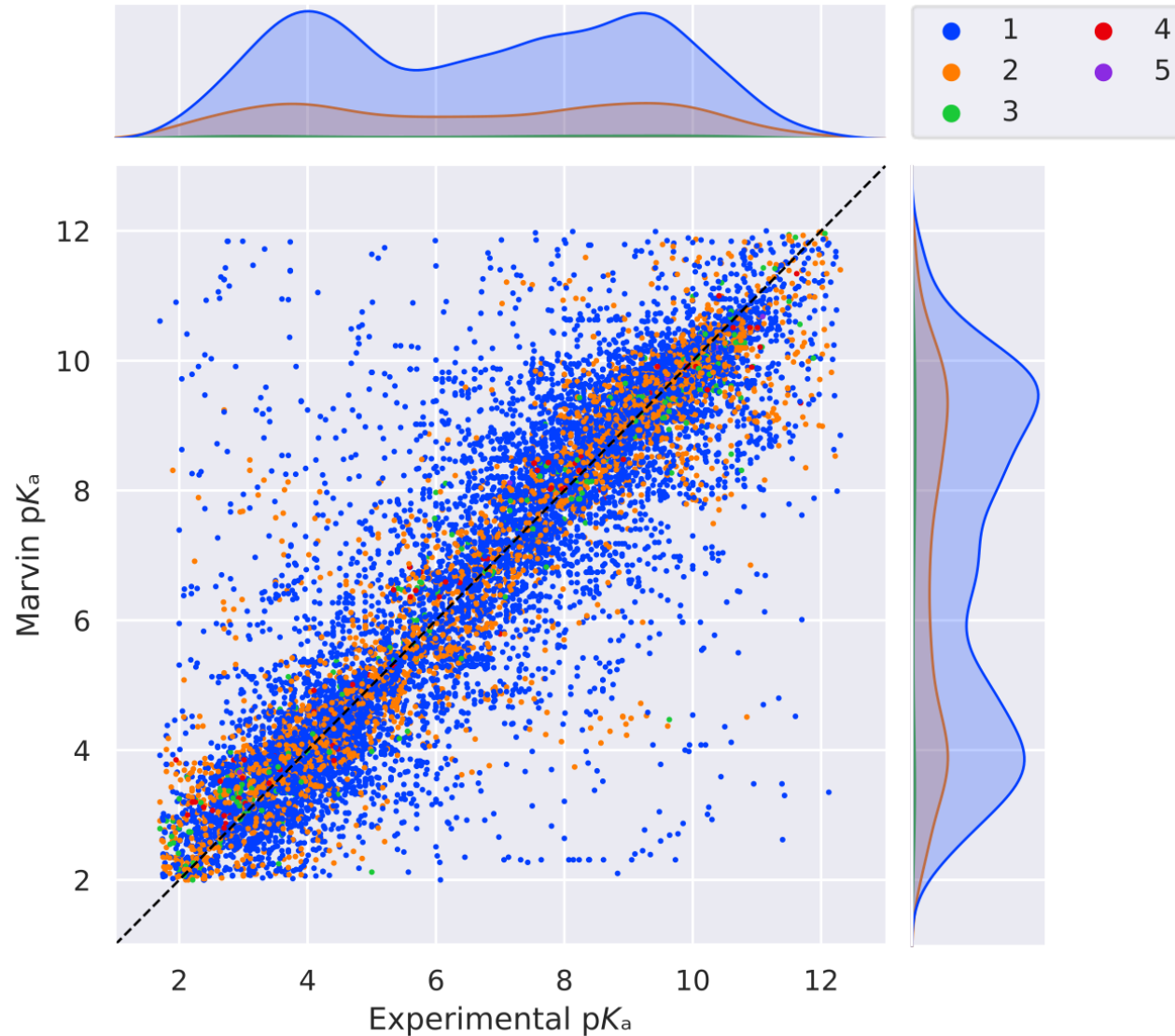  - **Must be done only for training set**
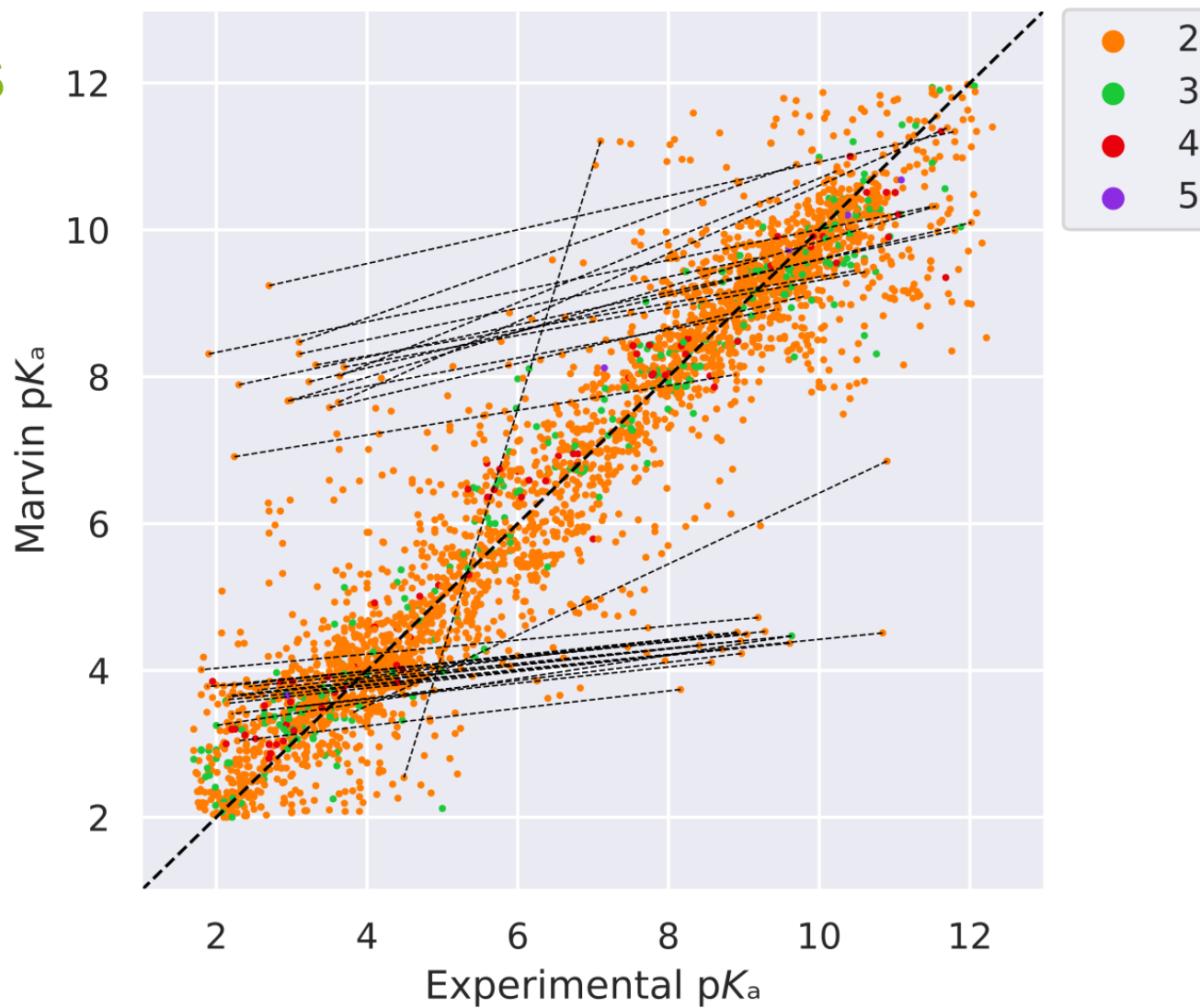
5.31

8.76

# Assign Values to Groups

- Combine values that apparently belong to the same titratable group

  – Error range of 0.3 p$K_a$ units

- Find the experimental value that comes closest to the corresponding *Marvin* prediction
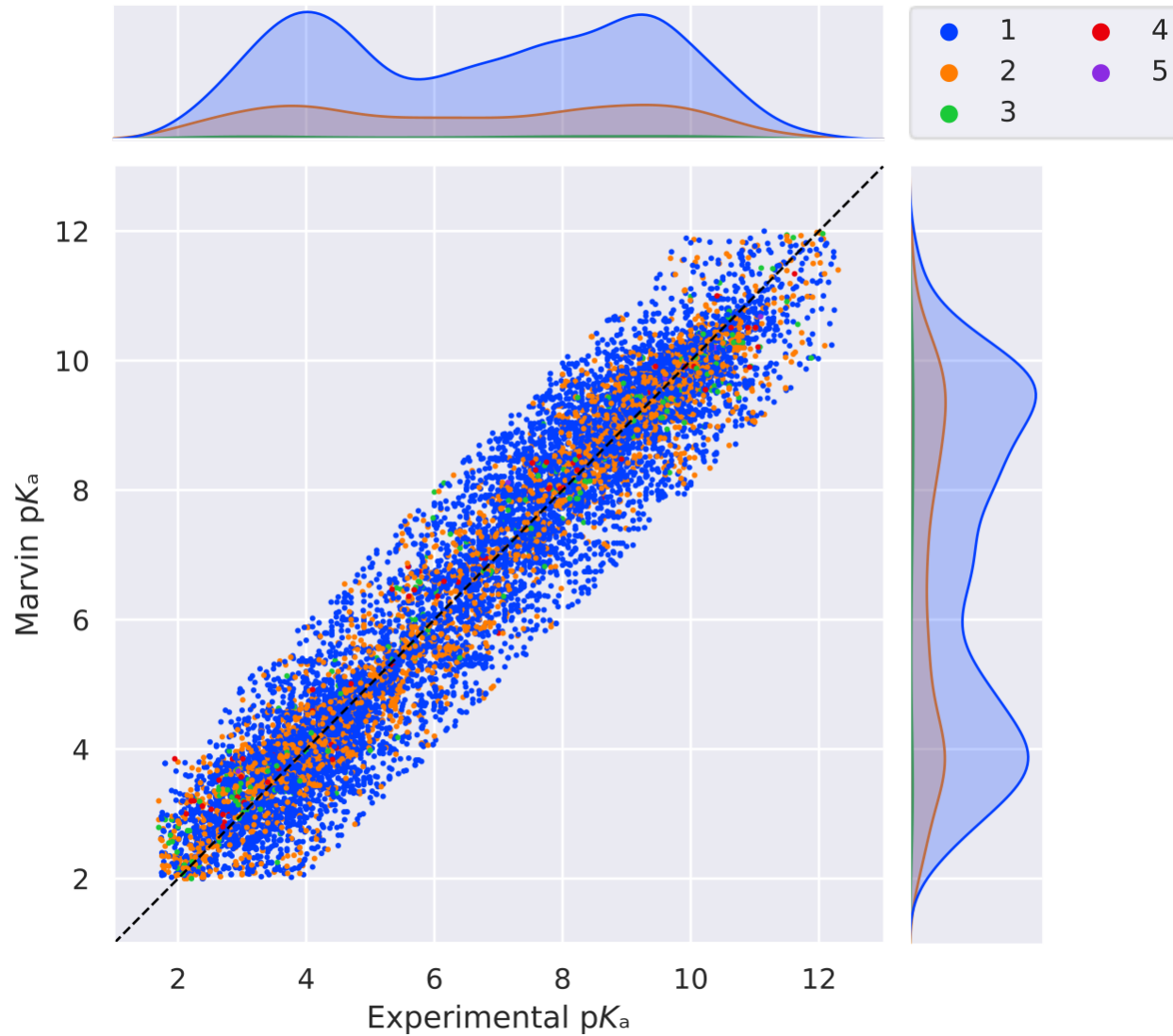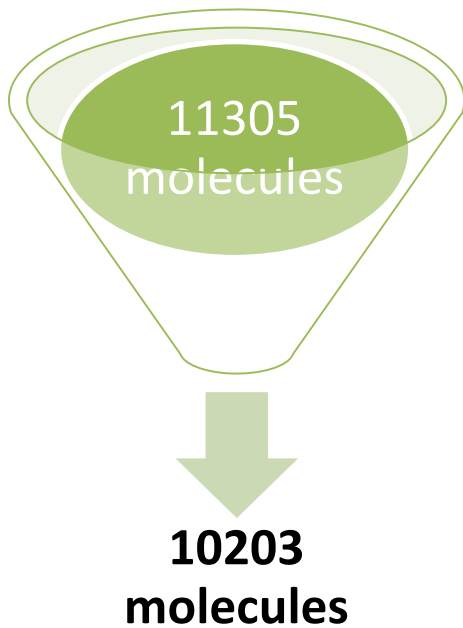
- Only consider "exact matches" for now

17538 molecules

**11305 molecules**

# Results

# Results

# Results

- Cut at max. error = 2



11305
molecules

**10203
molecules**

# Outlook

Further investigation of the results and testing with other prediction tools

Improvement of SMARTS pattern extraction

Reducing amount of rejects through value assignment

Replace OpenEye tautomers with RDKit integrated MolVS

Start with machine learning for multiprotic molecules

Develop and publish an easy-to-use toolkit

# Acknowledgements

**CzodrowskiLab**
Paul Czodrowski
David Bushiri

**California State University, USA**
Enrico Tapavicza

**Novartis Pharma AG**
Richard A. Lewis
Stephane Rodde

**Roche Pharma AG**
Christian Kramer

**Bayer AG**
Michael E. Beck