# *Predicting hits with ML and limited data:*
# *3 new tricks*

**Jan H. Jensen**

Department of Chemistry,
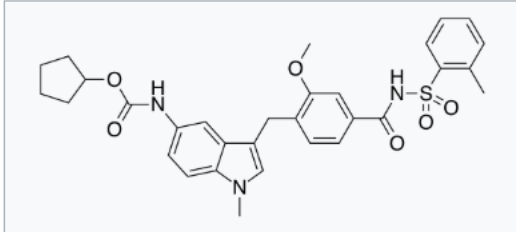University of Copenhagen

@janhjensen

# Target: some membrane protein (no X-ray structure)
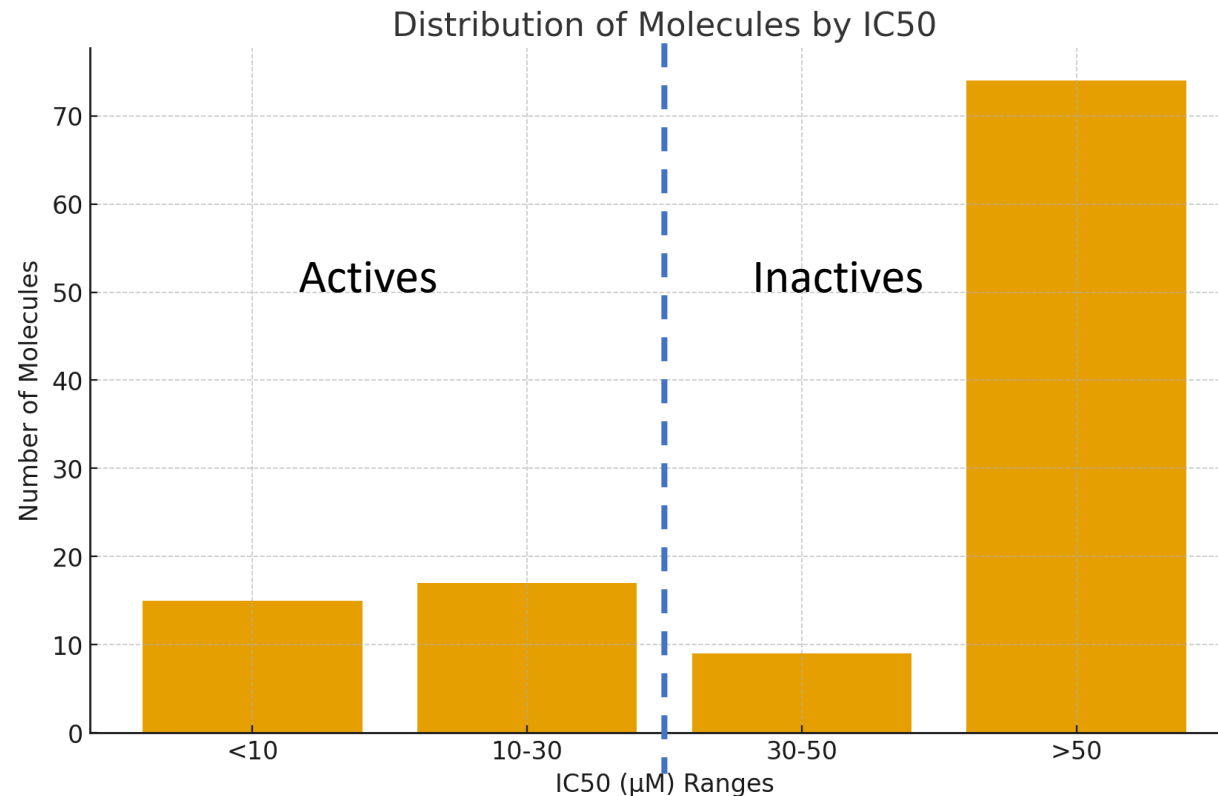
**Zafirlukast**
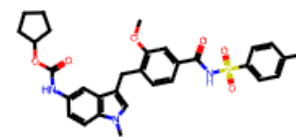


$IC_{50}$ 2 µM

~5 years
Some Pharma goodwill
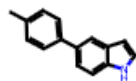BS and MSc students

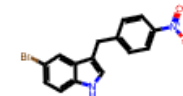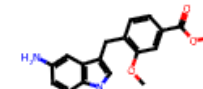1st-year PhD student*

115 indoles

**\*Niels Guldager**
**Christian M. Pedersen**
**Anders Aa. Rehfeld**

Distribution of Molecules by IC50

Actives          Inactives

Number of Molecules

<10        10-30        30-50        >50

IC50 (µM) Ranges

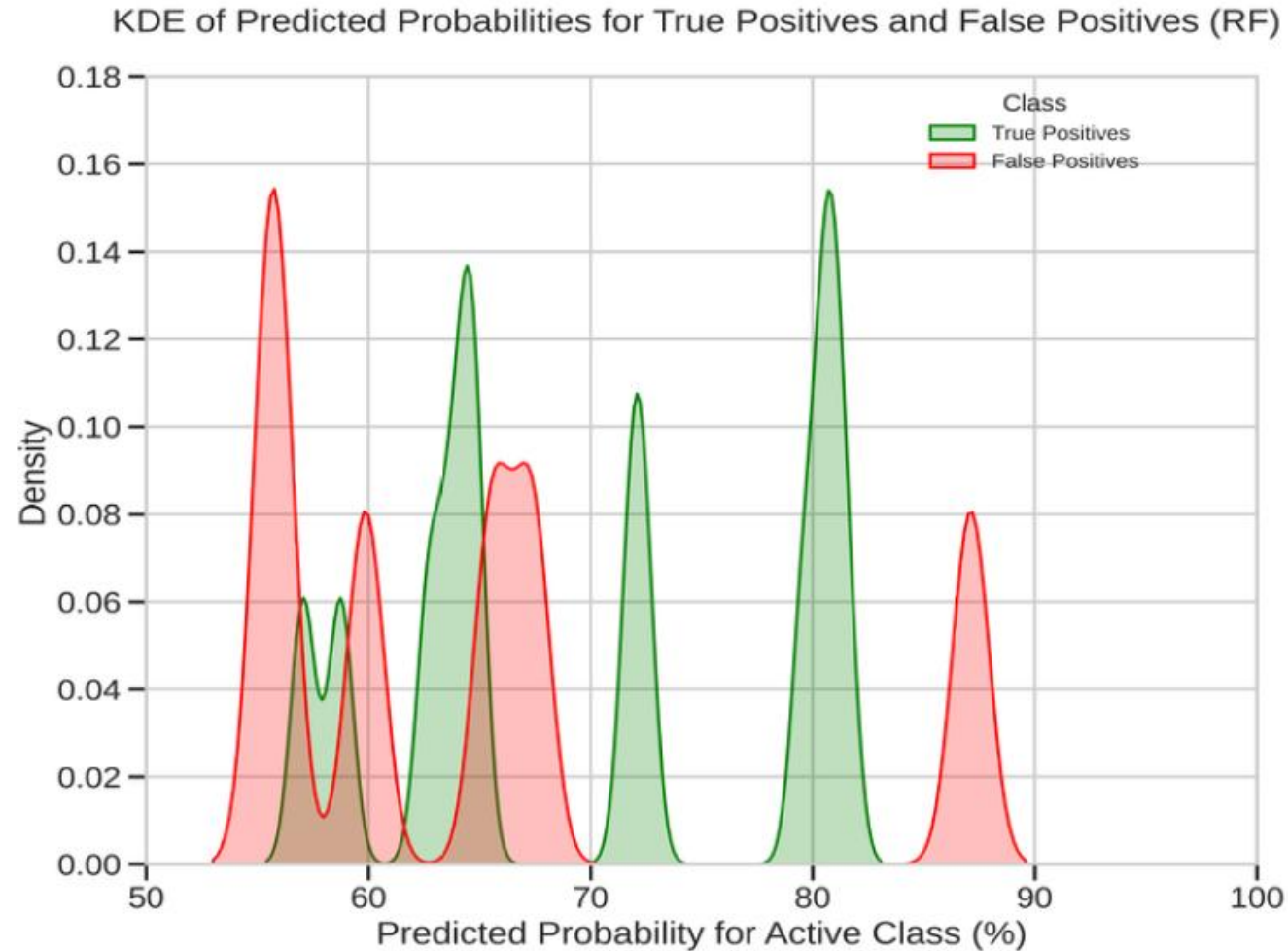# 66% actives are zafirkulast derivatives

zafirlukast

# Very challenging data set: Businness as usual doesn't work

ECFP4-based RF* classifier

**Henry Teahan**
**Maria H. Rasmussen**
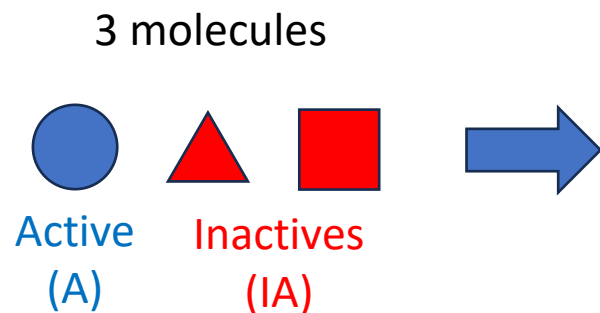
Leave-one-out



KDE of Predicted Probabilities for True Positives and False Positives (RF)

*bagged trees

# Trick # 1: PDL

Pairwise Difference Learning
for Classification

Mohamed Karim Belaid[1,2(✉)] ⓘ, Maximilian Rabus[2] ⓘ,
and Eyke Hüllermeier[3] ⓘ

## Training

9 pairs

3 molecules

Active
(A)

Inactives
(IA)

1 (same category)

0 (different category)

0

0

1

1 (same category)

0

1

1

**Trick # 1: PDL**

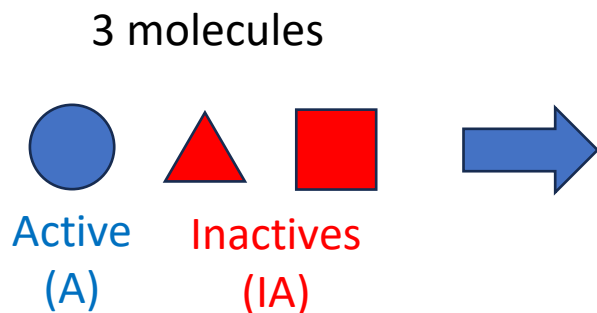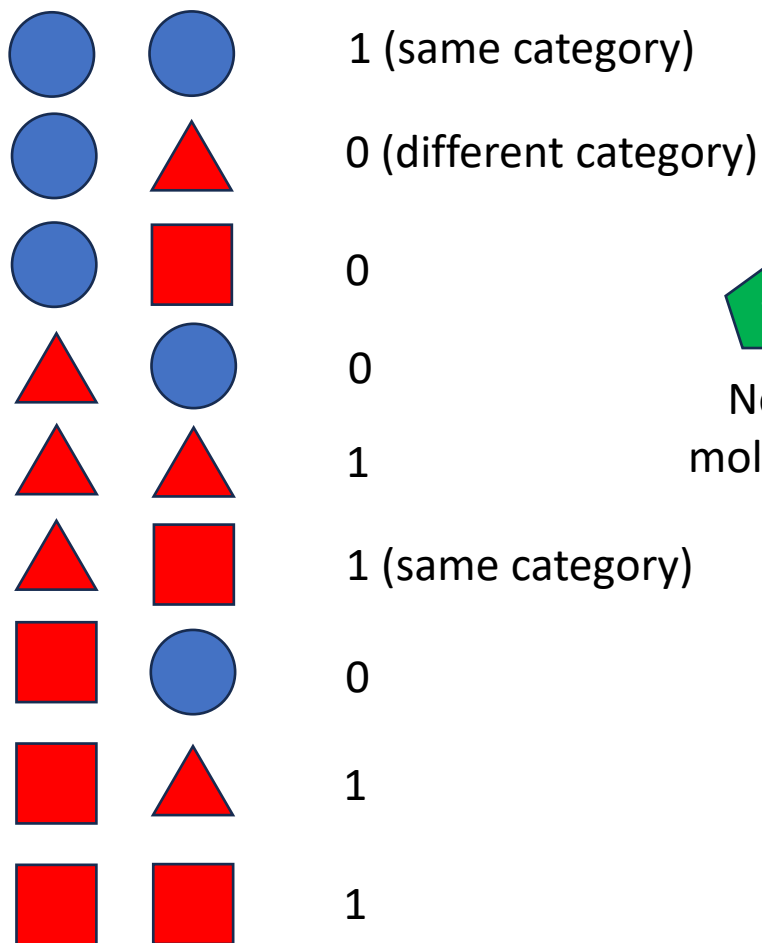Pairwise Difference Learning
for Classification

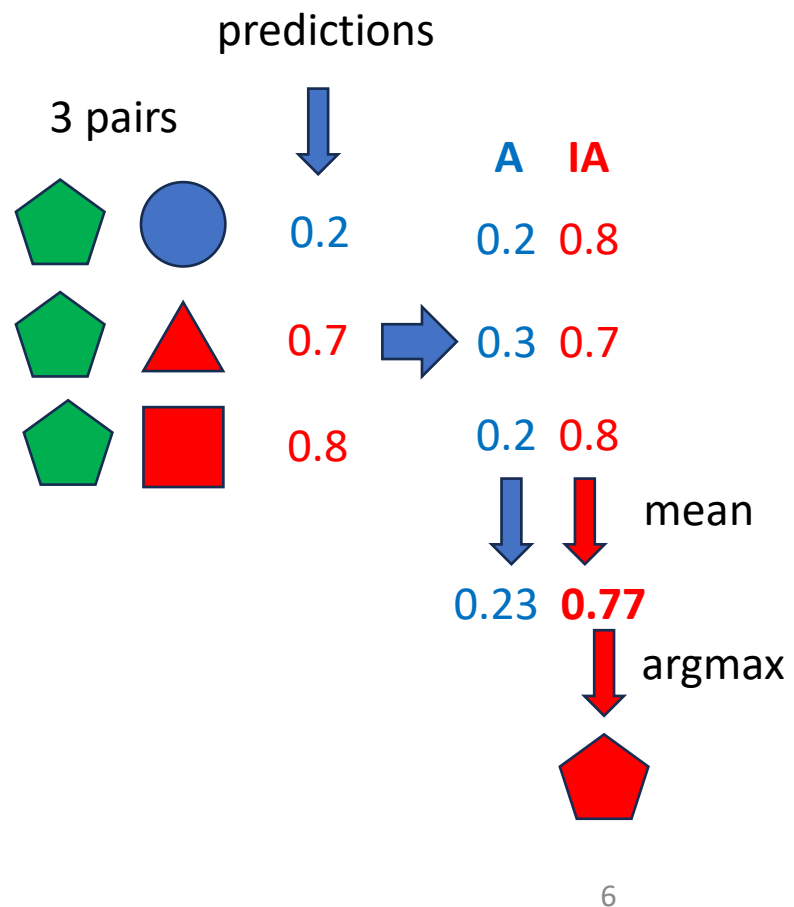Mohamed Karim Belaid[1,2(✉)] ⓘⒹ, Maximilian Rabus[2] ⓘⒹ,
and Eyke Hüllermeier[3] ⓘⒹ

# PDL works much better



KDE of Predicted Probabilities for True Positives and False Positives (PDL-RF)
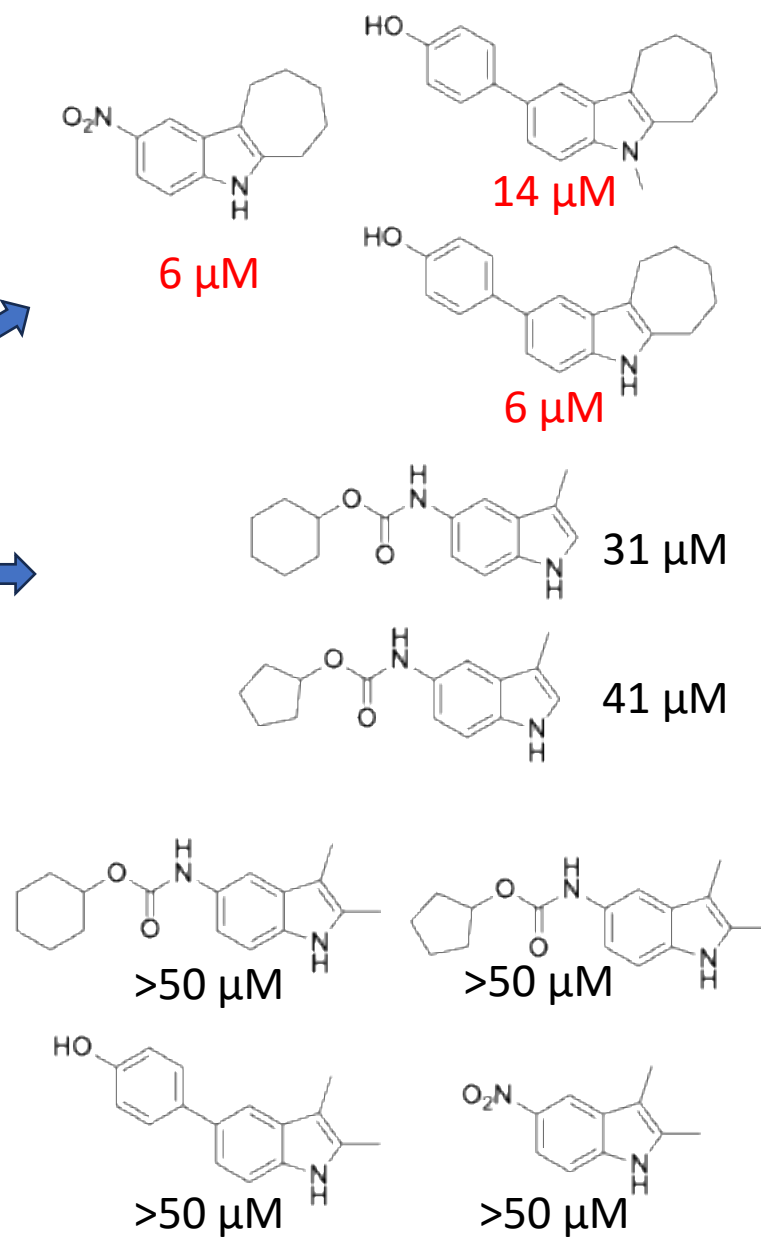
PDL fingerprint   **FP1** + **FP2** + **(FP1-FP2)**

**Screening Enamines Building Blocks set (1288 indoles)**

**Predicted**

**Tested**

14 μM

6 μM

6 μM

31 μM

41 μM

>50 μM

>50 μM

>50 μM

>50 μM

>50 μM

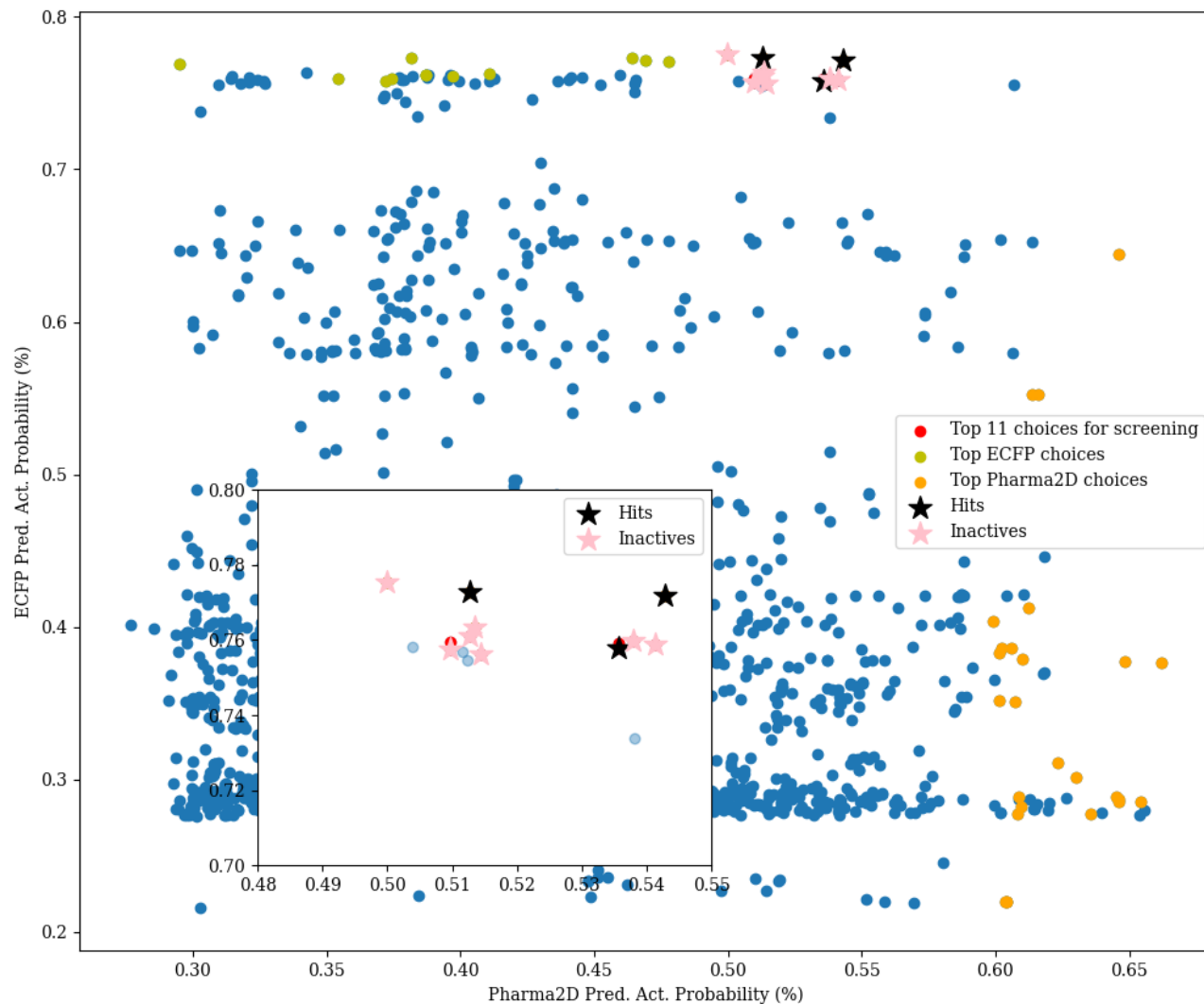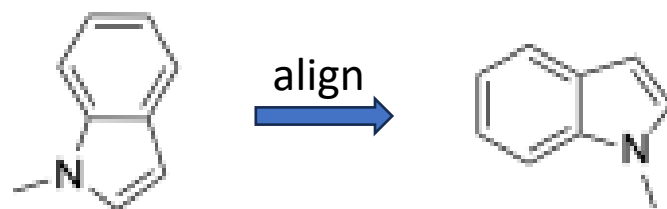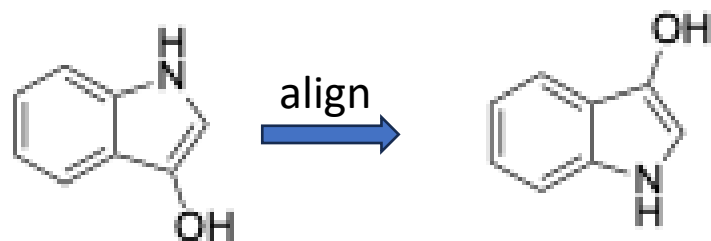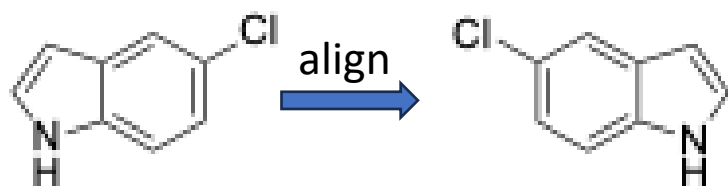8

# 11 new predictions from Enamine Hit Locator Set (7649 indoles)
## (10 were tested -> 0 hits)

# Trick # 2    The "2D atomic pharmacophore model"

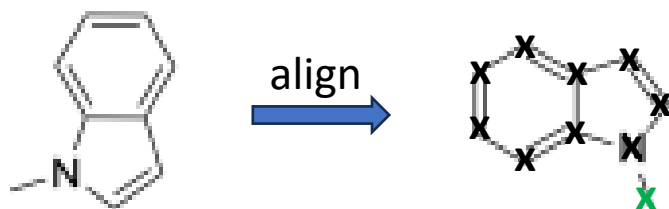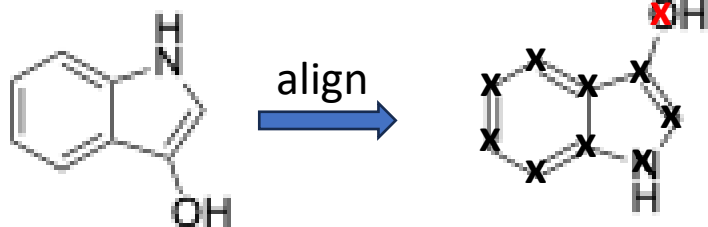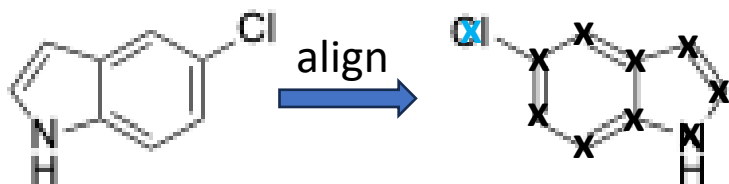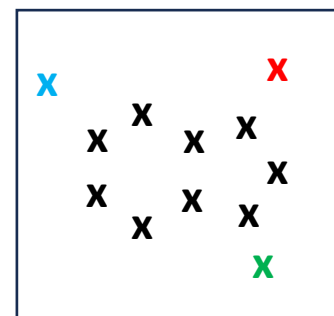**Training set**

# Trick # 2

# The "2D atomic pharmacophore model"

**Training set**



Combine
& cluster
2D coords

(x, … , x, x, x, x)

Descriptor vector

x's are…
Presence (binary)
HBD (binary)
HBA (binary)
Atomic logP values (float)

# Trick # 2      The "2D atomic pharmacophore model"

**Training set**



align

(x, ... , x, x, 0, 0)

align

(x, ... , x, 0, x, 0)

align

(x, ... , x, 0, 0, x)

Combine
& cluster
2D coords

(x, ... , x, x, x, x)

Descriptor vector
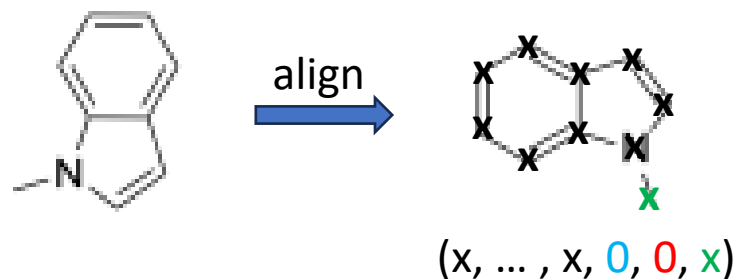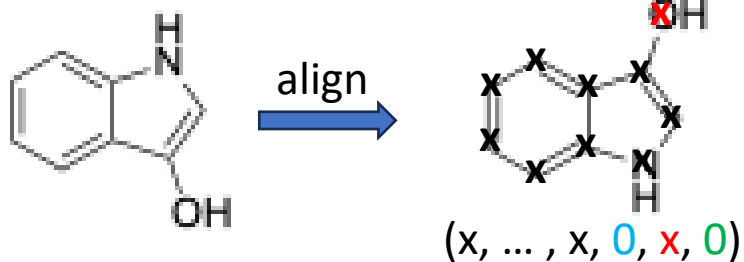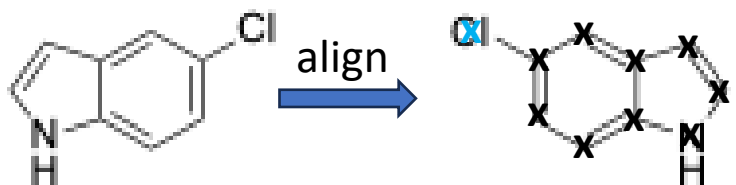
x's are...
Presence (binary)
HBD (binary)
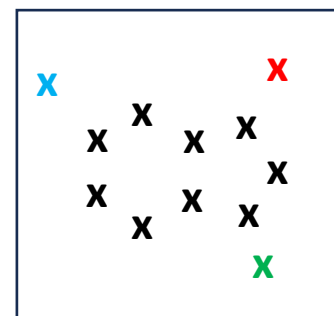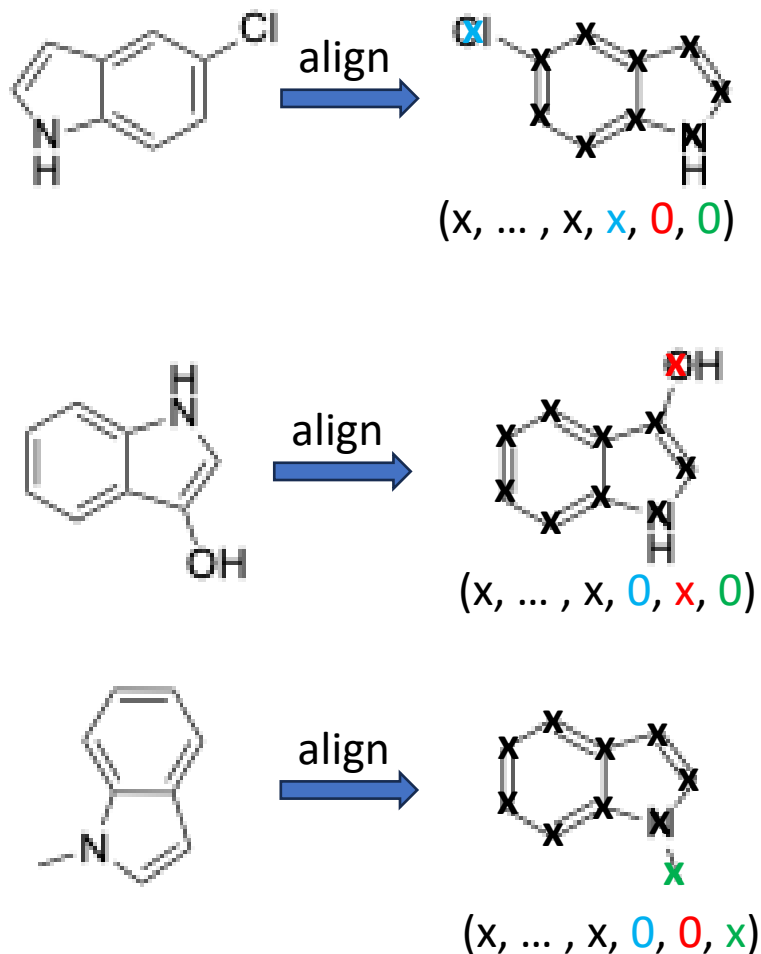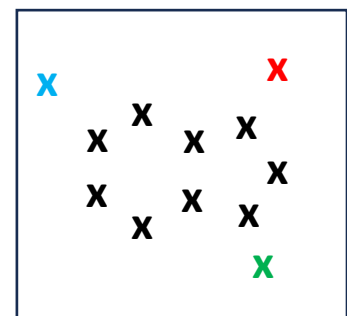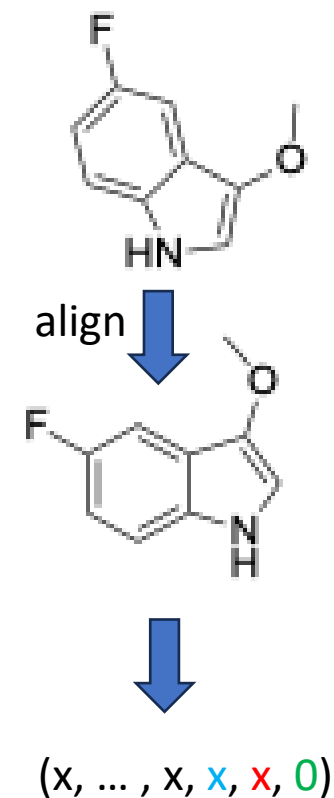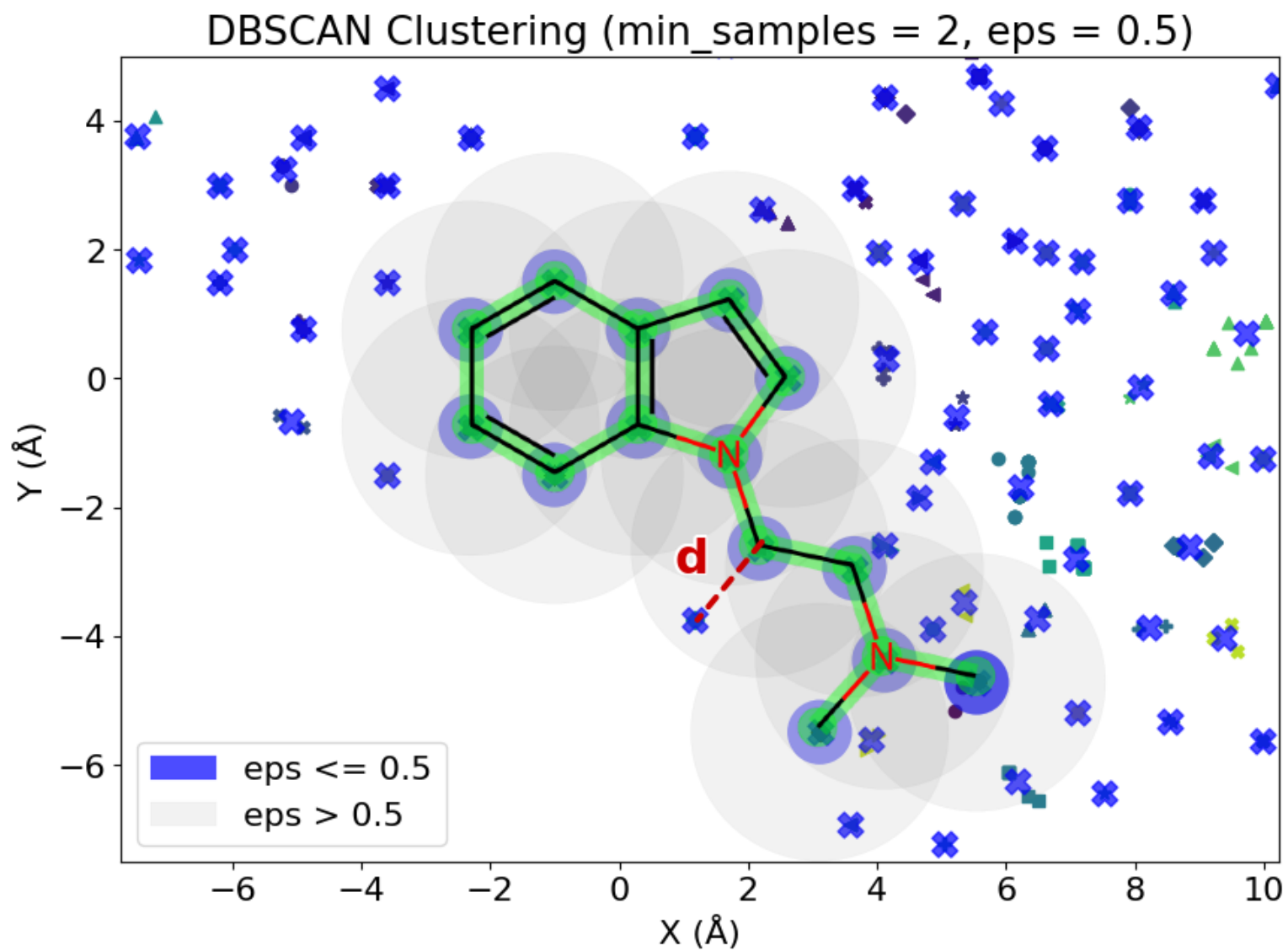HBA (binary)
Atomic logP values (float)

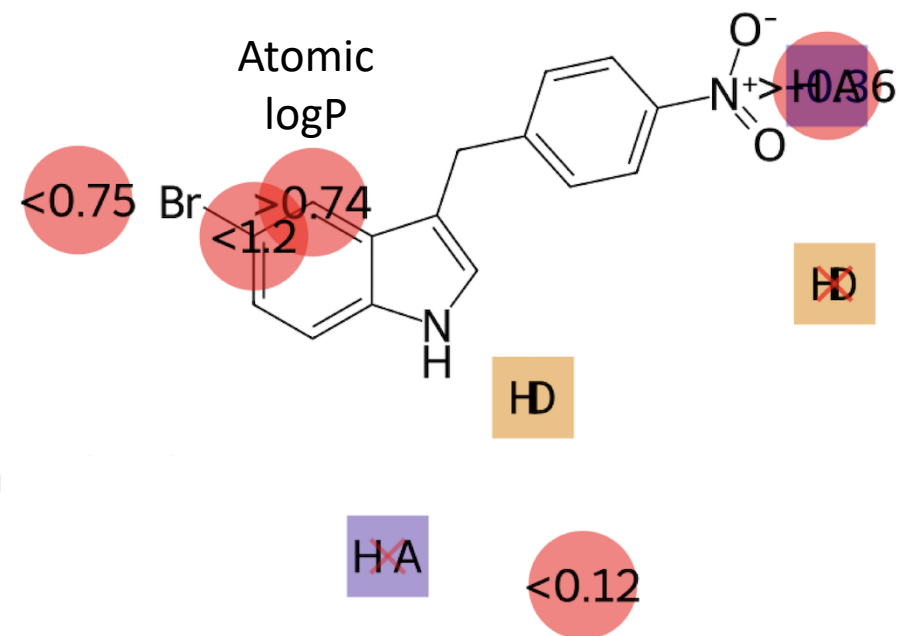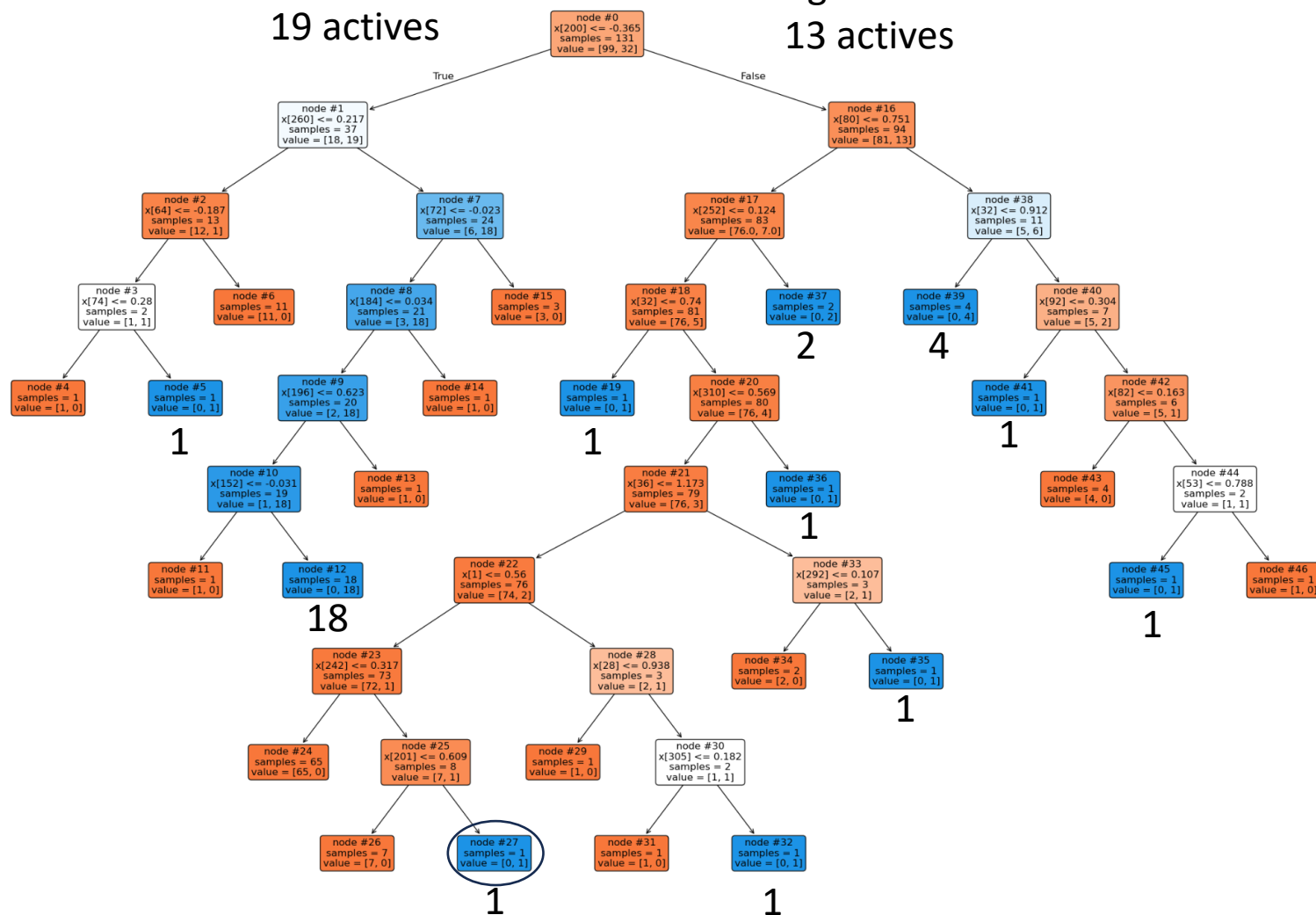# Trick # 2  The "2D atomic pharmacophore model"

DBSCAN Clustering (min_samples = 2, eps = 0.5)

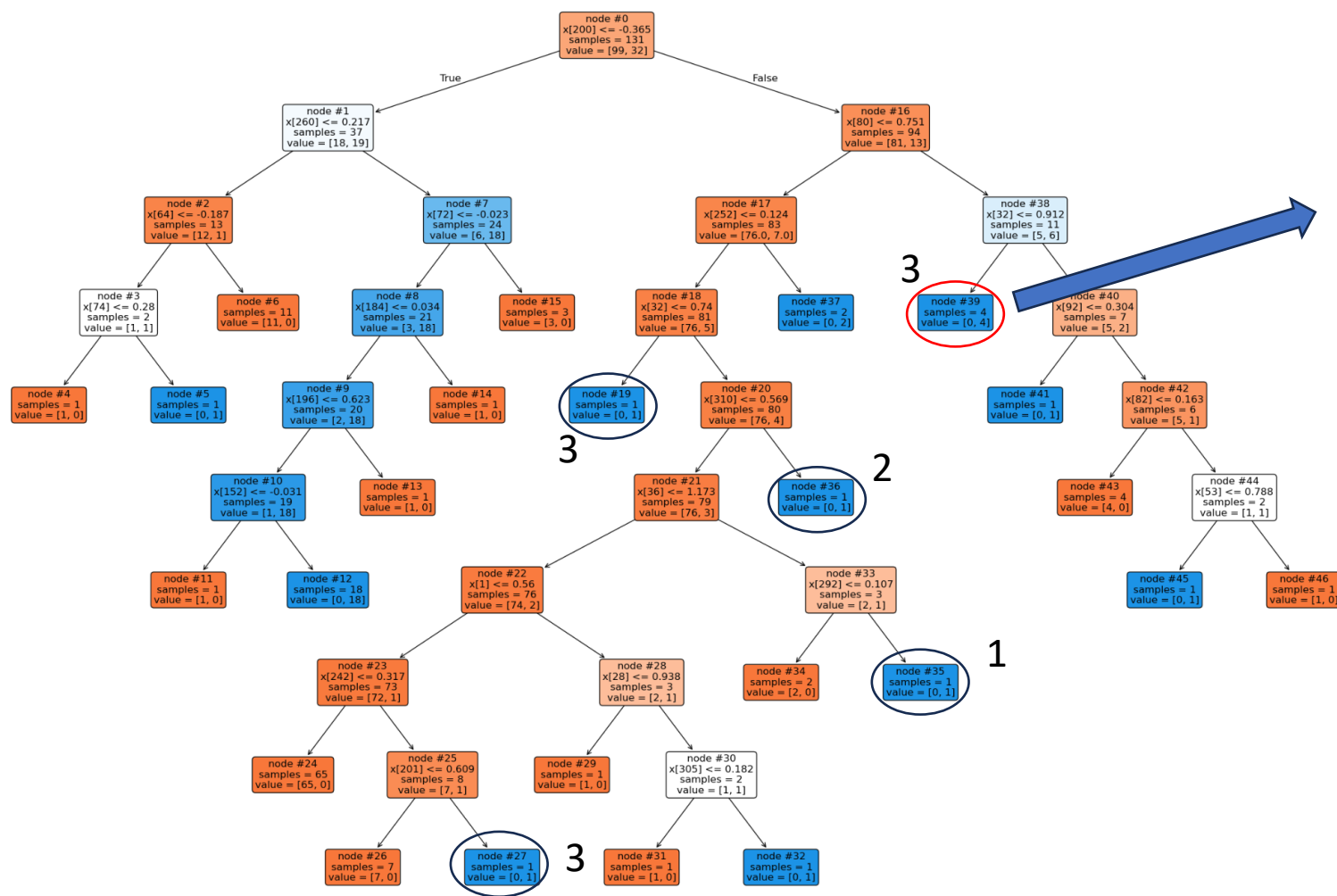eps > 0.5 → property scaled with distance

# Trick # 3

## Better descriptor → simpler model
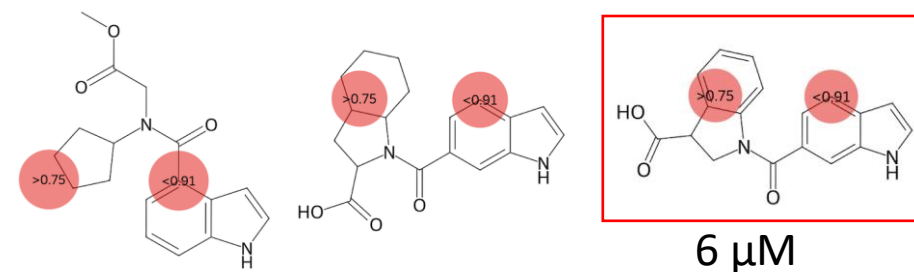### Decision tree offers interpretability



Left side:
19 actives
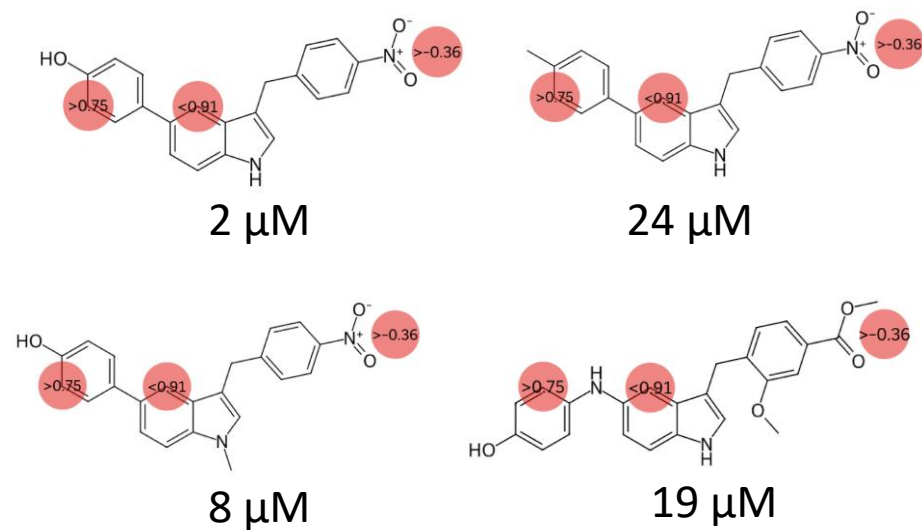
Right side:
13 actives

Atomic logP

# 12 new predictions from Enamine Hit Locator Set
## (11 were tested → 1 hit)

We picked these molecules …

… because of these molecules

# 3 tricks

Pair difference learning for classification*

2D atomic pharmacophore model

Decision tree for insight/analysis

*also works for regression