# Multi-Instance Explainable Learning

## for decoding stereo-dependent biological effects

**Vasilii Fastovskii**

PhD Candidate,
Strasbourg University,
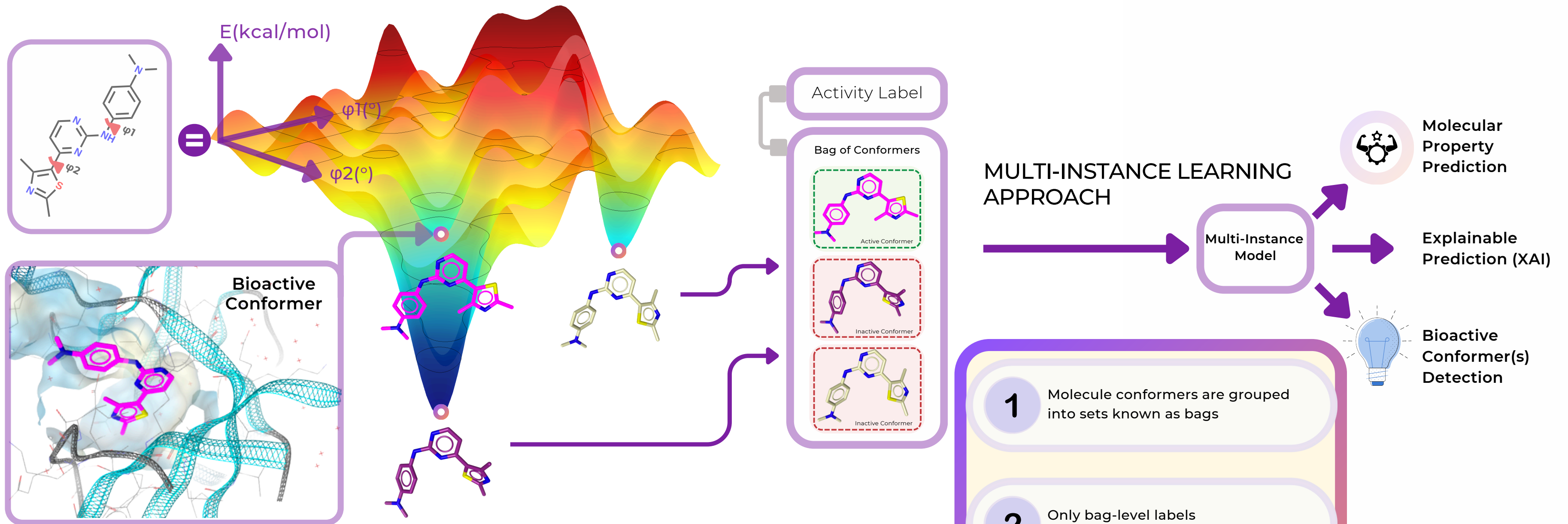Laboratory of Chemoinformatics

**Sanofi**

Dr Marc Bianciotto

Dr Christoph Grebner

**Strasbourg University,
Laboratory of Chemoinformatics**

Dr Gilles Marcou

Dr Prof Alexandre Varnek

E(kcal/mol)

φ1(°)

φ2(°)

Bioactive Conformer

Activity Label

Bag of Conformers

Active Conformer

Inactive Conformer

Inactive Conformer

MULTI-INSTANCE LEARNING APPROACH

Multi-Instance Model

Molecular Property Prediction

Explainable Prediction (XAI)

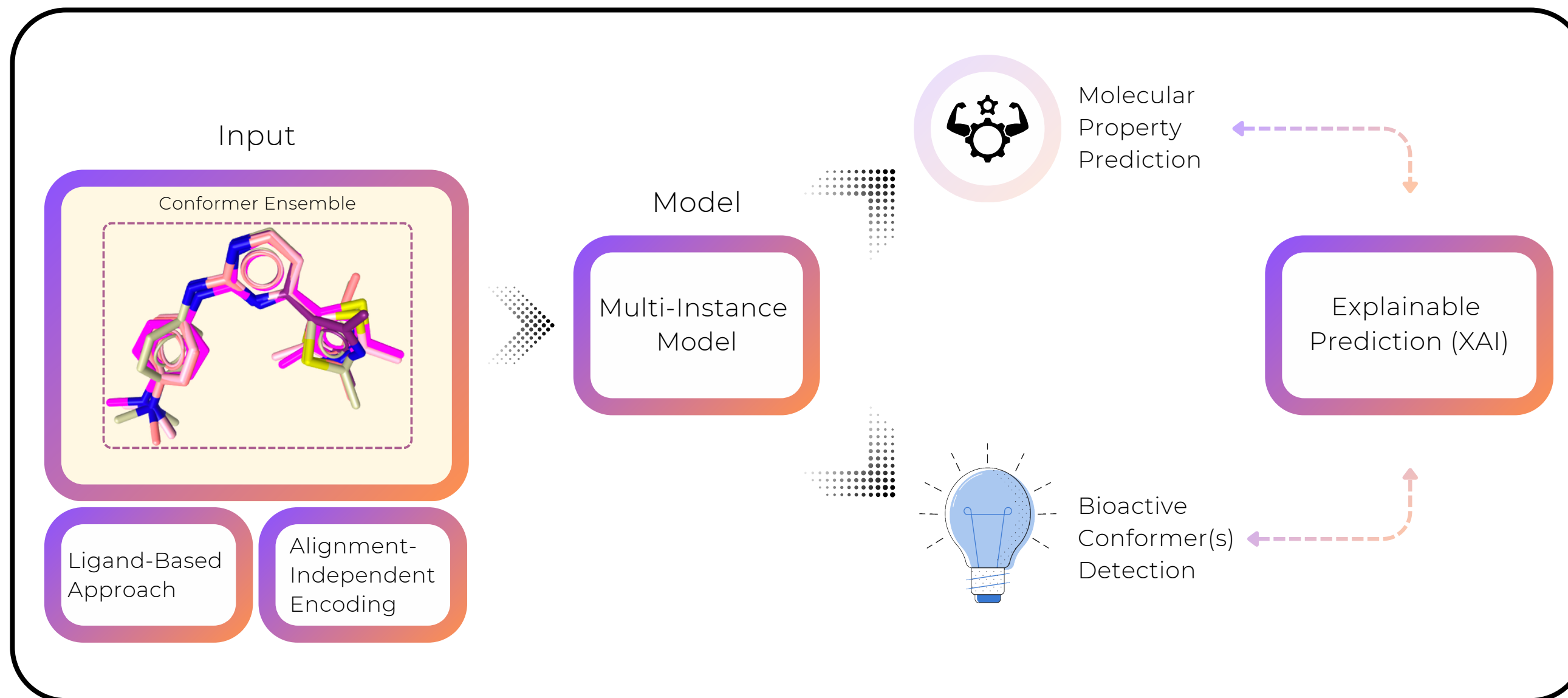Bioactive Conformer(s) Detection

**1** Molecule conformers are grouped into sets known as bags

**2** Only bag-level labels are available during training

**3** Standard MIL Task - Train a classifier/regressor that labels new bags

**4** Key Instance Detection (KID) Task - Identify the instances (conformers) most responsible for the bag's label
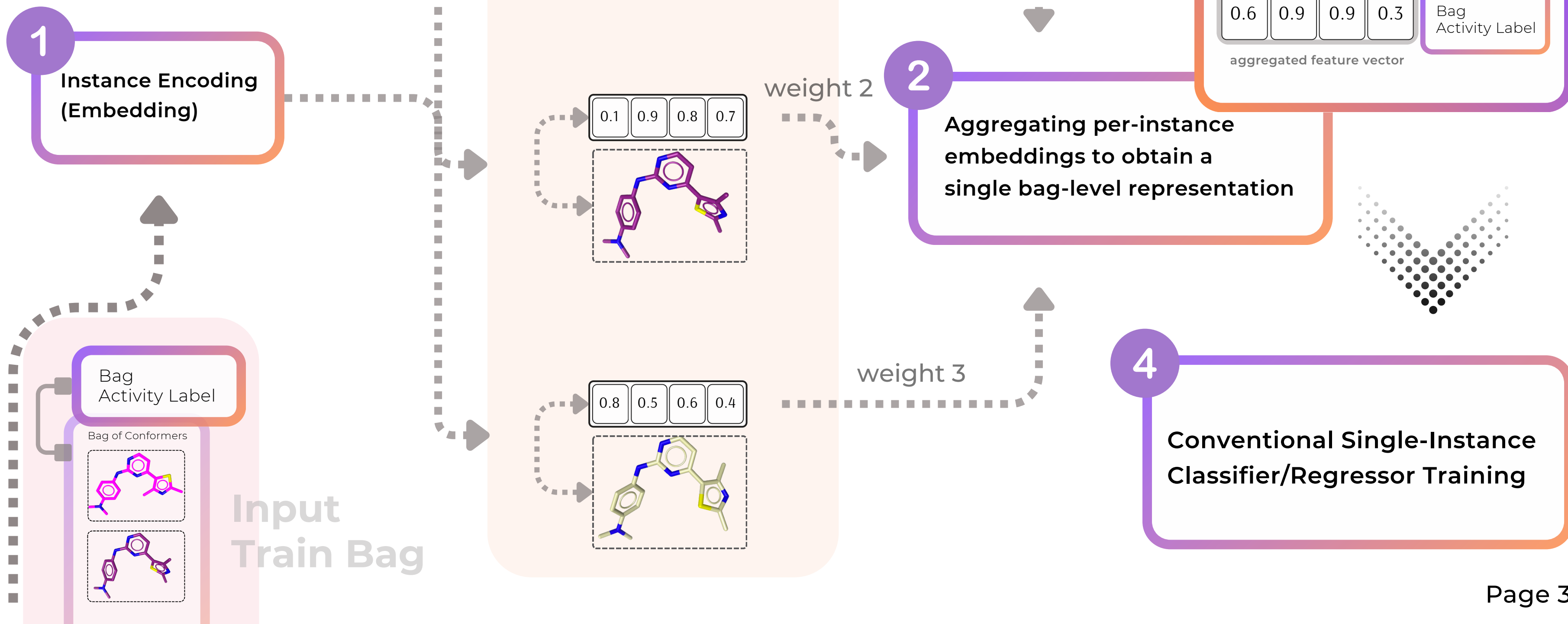
We aim to identify the exact molecular form(s) responsible for the observed/predicted properties to provide a data-driven three-dimensional shape hypothesis supporting a ligand-based model

Molecules can adopt multiple forms (conformers, tautomers, protonation states)
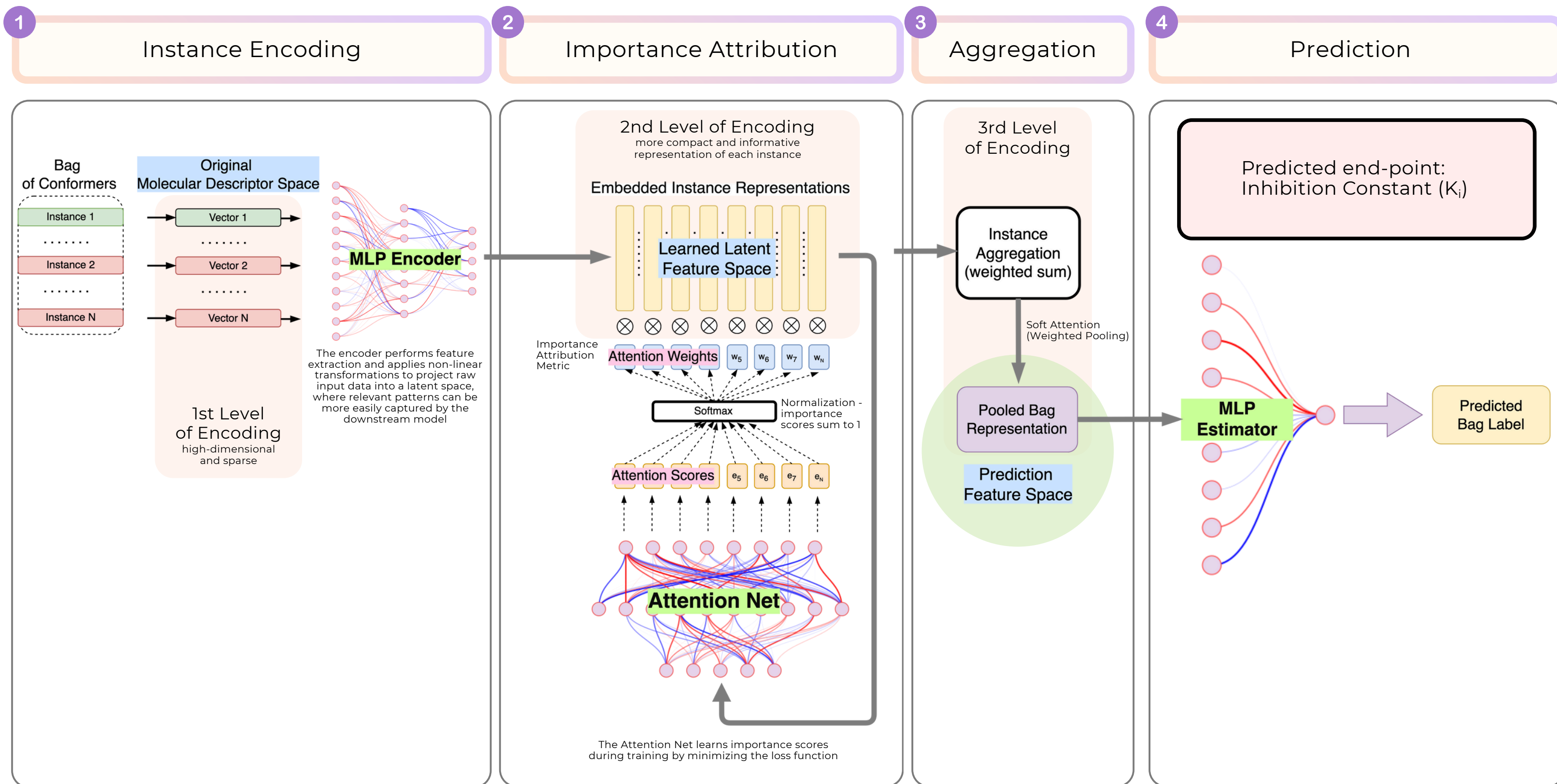
## 2. PROJECT GOALS

Input

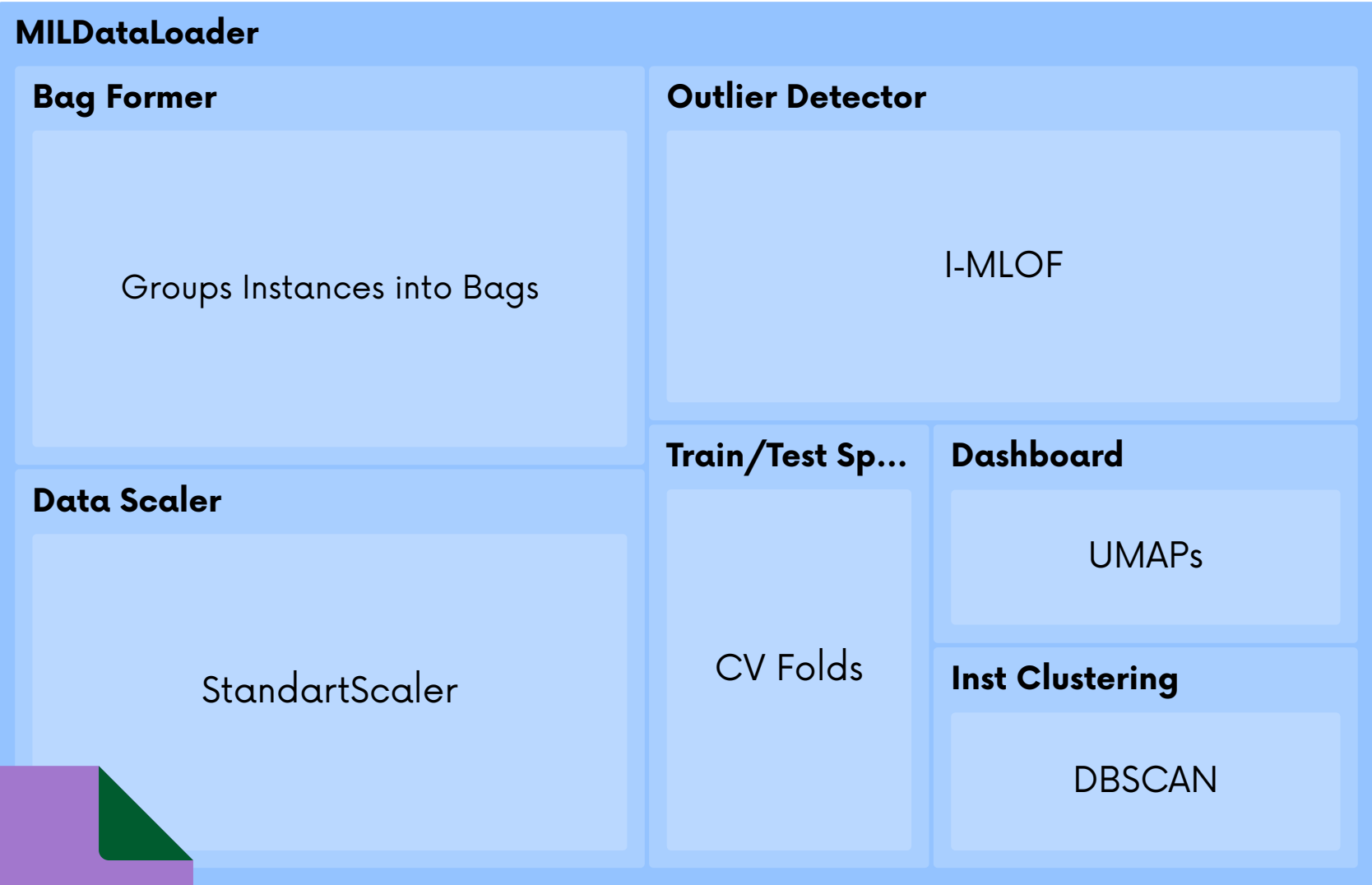Conformer Ensemble



Ligand-Based Approach

Alignment-Independent Encoding

Model

Multi-Instance Model

Molecular Property Prediction

Explainable Prediction (XAI)

Bioactive Conformer(s) Detection

# EMBEDDED-SPACE (ES) MIL METHODS

**Train**

learned instance feature vector

| 0.2 | 0.3 | 0.7 | 0.1 |

weight 1

**Single bag-level representation (for each Bag in the Train Set)**

**3**

| 0.6 | 0.9 | 0.9 | 0.3 | Bag Activity Label |

aggregated feature vector

**1** Instance Encoding (Embedding)

| 0.1 | 0.9 | 0.8 | 0.7 |

weight 2

**2** Aggregating per-instance embeddings to obtain a single bag-level representation

| 0.8 | 0.5 | 0.6 | 0.4 |

weight 3

**4** Conventional Single-Instance Classifier/Regressor Training

**Input Train Bag**

Bag Activity Label

Bag of Conformers

# 3. PIPELINE

## 1 Instance Encoding

### Bag of Conformers

Instance 1
......
Instance 2
......
Instance N

### Original Molecular Descriptor Space

Vector 1
......
......
Vector 2
......
Vector N

**MLP Encoder**

The encoder performs feature extraction and applies non-linear transformations to project raw input data into a latent space, where relevant patterns can be more easily captured by the downstream model

### 1st Level of Encoding
high-dimensional and sparse

## 2 Importance Attribution

### 2nd Level of Encoding
more compact and informative representation of each instance

### Embedded Instance Representations

**Learned Latent Feature Space**

$\otimes$ $\otimes$ $\otimes$ $\otimes$ $\otimes$ $\otimes$ $\otimes$ $\otimes$

Importance Attribution Metric

**Attention Weights** | $w_5$ | $w_6$ | $w_7$ | $w_N$

Softmax

Normalization - importance scores sum to 1

**Attention Scores** | $e_5$ | $e_6$ | $e_7$ | $e_N$

**Attention Net**

The Attention Net learns importance scores during training by minimizing the loss function

## 3 Aggregation

### 3rd Level of Encoding

Instance Aggregation (weighted sum)

Soft Attention (Weighted Pooling)

Pooled Bag Representation

Prediction Feature Space

## 4 Prediction

Predicted end-point: Inhibition Constant ($K_i$)

**MLP Estimator**

Predicted Bag Label

Zankov D. et al. *J Chem Inf Model.* **2021**, 61, 10, 4913-4923.

Page 4

## 1. Data Loading and Preprocessing

**MILDataLoader**

**Bag Former**

Groups Instances into Bags

**Outlier Detector**

I-MLOF

**Data Scaler**

StandartScaler

**Train/Test Sp...**

CV Folds

**Dashboard**

UMAPs

**Inst Clustering**

DBSCAN

**CSV**

Descriptors

I-MLOF: Instance-Neighborhood based Multi-Instance Local Outlier Factor
UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction
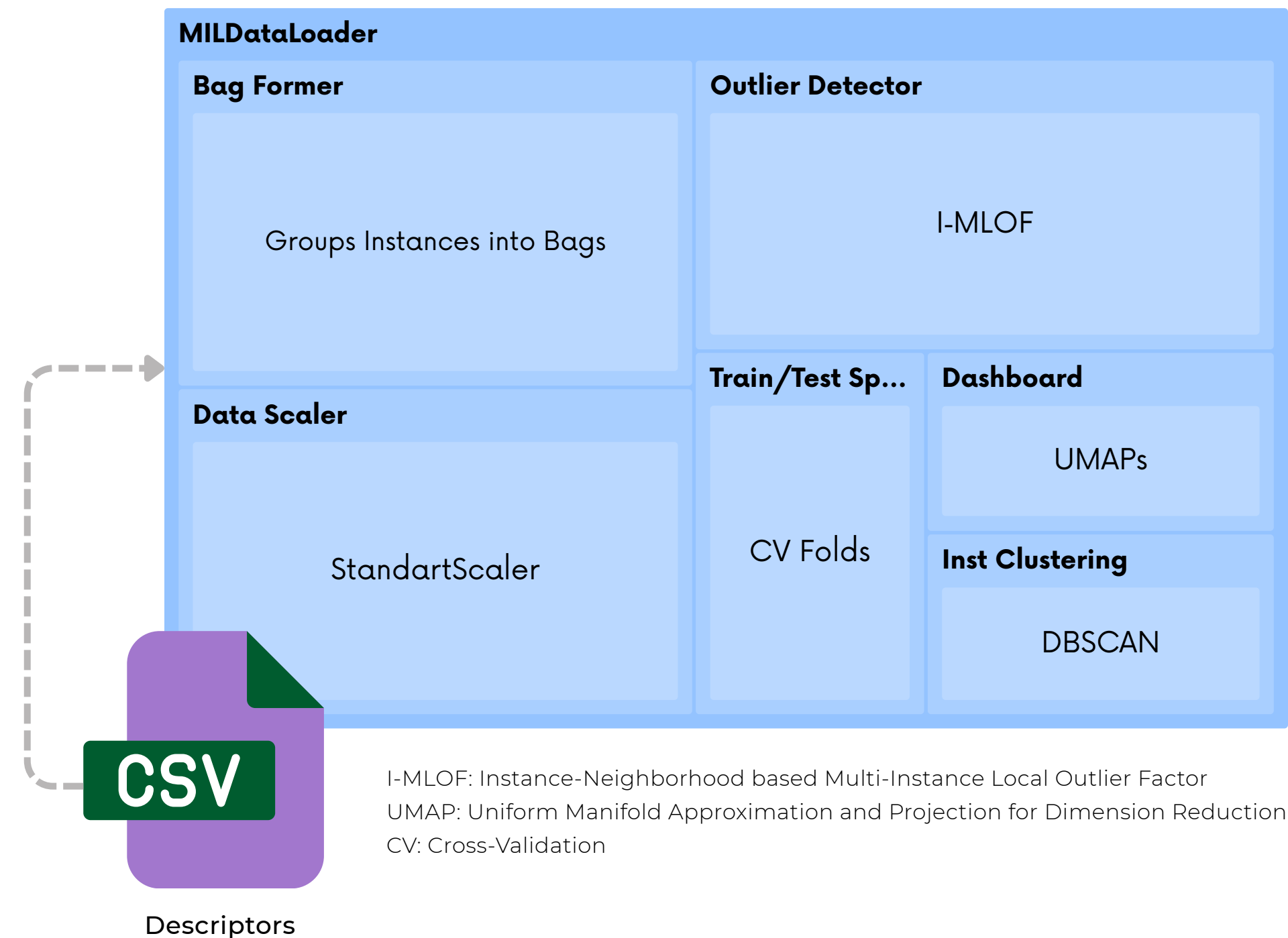CV: Cross-Validation

## Instance-Neighborhood Based
## Multi-Instance Local Outlier Factor (I-MLOF)

$$I\!-\!MLOF_{MinPts}(B) = \frac{\sum\limits_{C \in N_{MinPts}(B)} \frac{lrd_{MinPts}(C)}{lrd_{MinPts}(B)}}{|N_{MinPts}(B)|}$$

Wu, O., Li, B., Hu, W., Gao, J., & Zhu, M. (2010). Identifying Multi-instance Outliers. 430–441.
https://doi.org/10.1137/1.9781611972801.38

# MIL-Based QSAR Pipeline

## 1. Data Loading and Preprocessing

**MILDataLoader**

**Bag Former**

Groups Instances into Bags

**Data Scaler**

StandartScaler

**Outlier Detector**

I-MLOF

**Train/Test Sp...**

CV Folds

**Dashboard**

UMAPs

**Inst Clustering**

DBSCAN

**CSV**

Descriptors

I-MLOF: Instance-Neighborhood based Multi-Instance Local Outlier Factor
UMAP: Uniform Manifold Approximation and Projection for Dimension Reduction
CV: Cross-Validation

## Uniform Manifold Approximation and Projection (UMAP)



UMAP Projection

USRCAT, 3D FP
n_neighbors=20,
metric='cosine',
densmap=True

Label • PosInst_PosBag • NegBag • NegInst_PosBag • Outlier

An interactive Plotly Dash dashboard can be launched locally on your machine. The interface provides key dataset statistics and a visualization of the original descriptor space

# MIL–Based QSAR Pipeline

## 2. Baseline Evaluation

**MILBaselineBuilder**

**Classification [Linear Model - LR]**

| meanPool | maxPool |
|----------|---------|

instClf

**Classification [Nonlinear Model - RF]**

| meanPool | maxPool |
|----------|---------|

instClf



**MIL BL Performance [20-Fold CV (95% CI)]  USRCAT, 3D FP**

Metrics:
- accuracy
- precision
- recall
- f1_score
- roc_auc

The MILBaselineBuilder Class provides several baseline MIL strategies, such as mean pooling, max pooling, or instance-level classifier with max aggregation [instClf], with Logistic Regression [LR] or Random Forest [RF]

Target: Nitric-oxide synthase, CHEMBL3048
Regression Task



MORSE FP demonstrated the strongest performance in the MI setting, but performed poorly compared to other FPs in the SI setting.

*20-Fold Cross-Validation
**3D Single-Instance (SI) Model was built on the minimum energy conformers (ETKDGv3+MMFF94)
***3D Multi-Instance (MI) Model was built on molecular conformer ensembles (50 conformers/mol, RMSD threshold 1Å, ETKDGv3+MMFF94)

# MIL-Based QSAR Pipeline

## 3. MIL Model Construction

### 3.1. MIL Embedder

**MIL Embedder**

| **MLPEmbedder** | **Pre-trained Embedder** |

**MILModelBuilder**

```python
class MILModelBuilder:
    def __init__(self,
                 embedder_type="mlp",
                 aggregator_type="max",
                 predictor_type="linear",
                 input_dim=128):
        self.embedder_type = embedder_type
        self.aggregator_type = aggregator_type
        self.predictor_type = predictor_type
        self.input_dim = input_dim

    def build(self):
        # 1) Build Embedder
        if self.embedder_type not in EMBEDDER_DICT:
            raise ValueError(f"Unknown embedder_type: {self.embedder_type}")
        embedder_cls = EMBEDDER_DICT[self.embedder_type]
        embedder = embedder_cls(input_dim=self.input_dim)

        # 2) Build Aggregator
        if self.aggregator_type not in AGGREGATOR_DICT:
            raise ValueError(f"Unknown aggregator_type: {self.aggregator_type}")
        aggregator_cls = AGGREGATOR_DICT[self.aggregator_type]
        aggregator = aggregator_cls(input_dim=embedder.get_output_dim())

        # 3) Build Predictor
        if self.predictor_type not in PREDICTOR_DICT:
            raise ValueError(f"Unknown predictor_type: {self.predictor_type}")
        predictor_cls = PREDICTOR_DICT[self.predictor_type]
        predictor = predictor_cls(input_dim=aggregator.get_output_dim())

        # 4) Compose final model
        model = MILModel(embedder, aggregator, predictor)
        return model
```
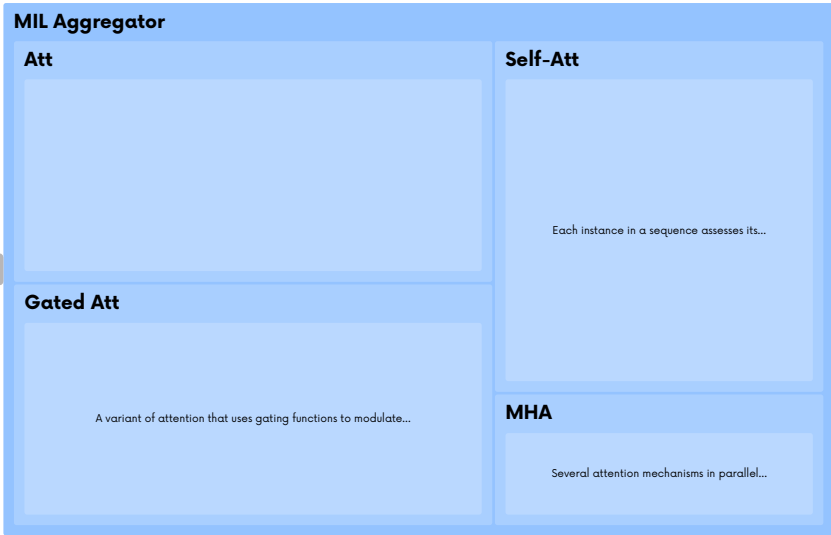
### 3.2. MIL Aggregator

**MIL Aggregator**

| **Att** | **Self-Att** |
|---|---|
| | Each instance in a sequence assesses its... |
| **Gated Att** | **MHA** |
| A variant of attention that uses gating functions to modulate... | Several attention mechanisms in parallel... |

### 3.3. MIL Predictor

**MIL Predictor**

**MLPPredictor**

**Config file**

The MILModelBuilder class is a modular constructor for an MIL model. It dynamically builds three main components - an embedder, an aggregator, and a predictor - by selecting their implementations based on the provided config file
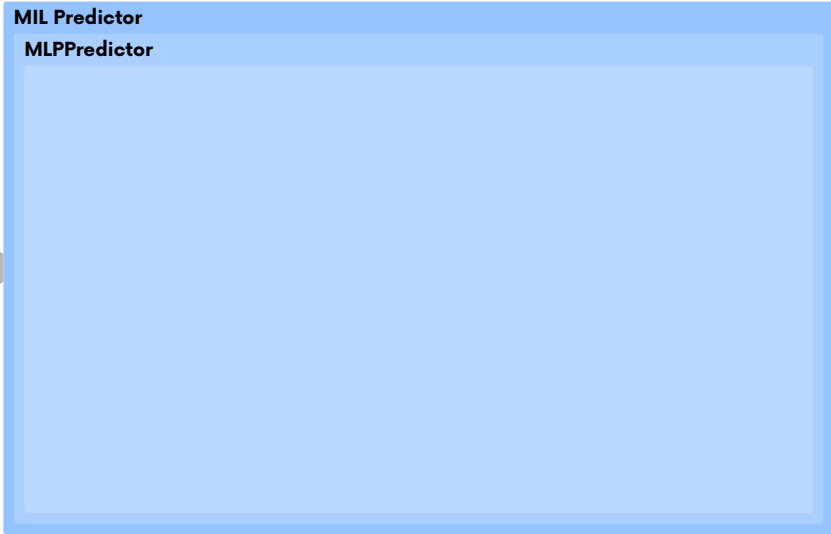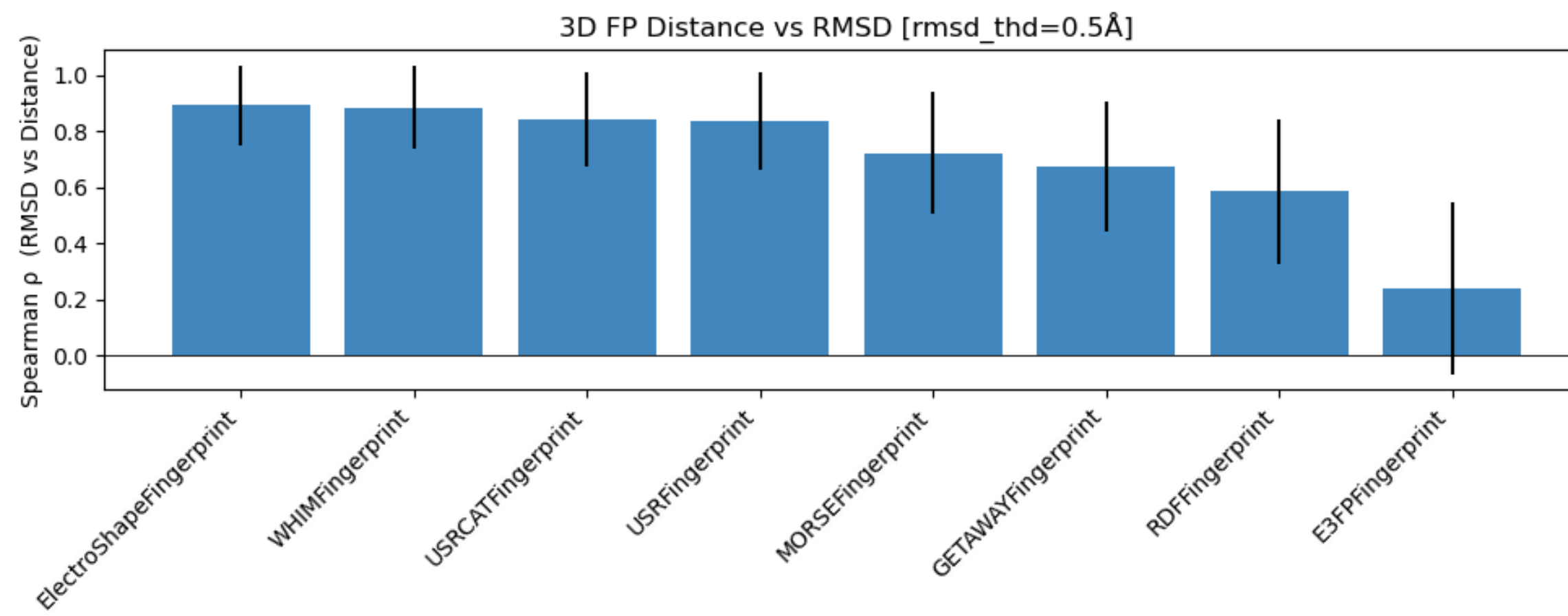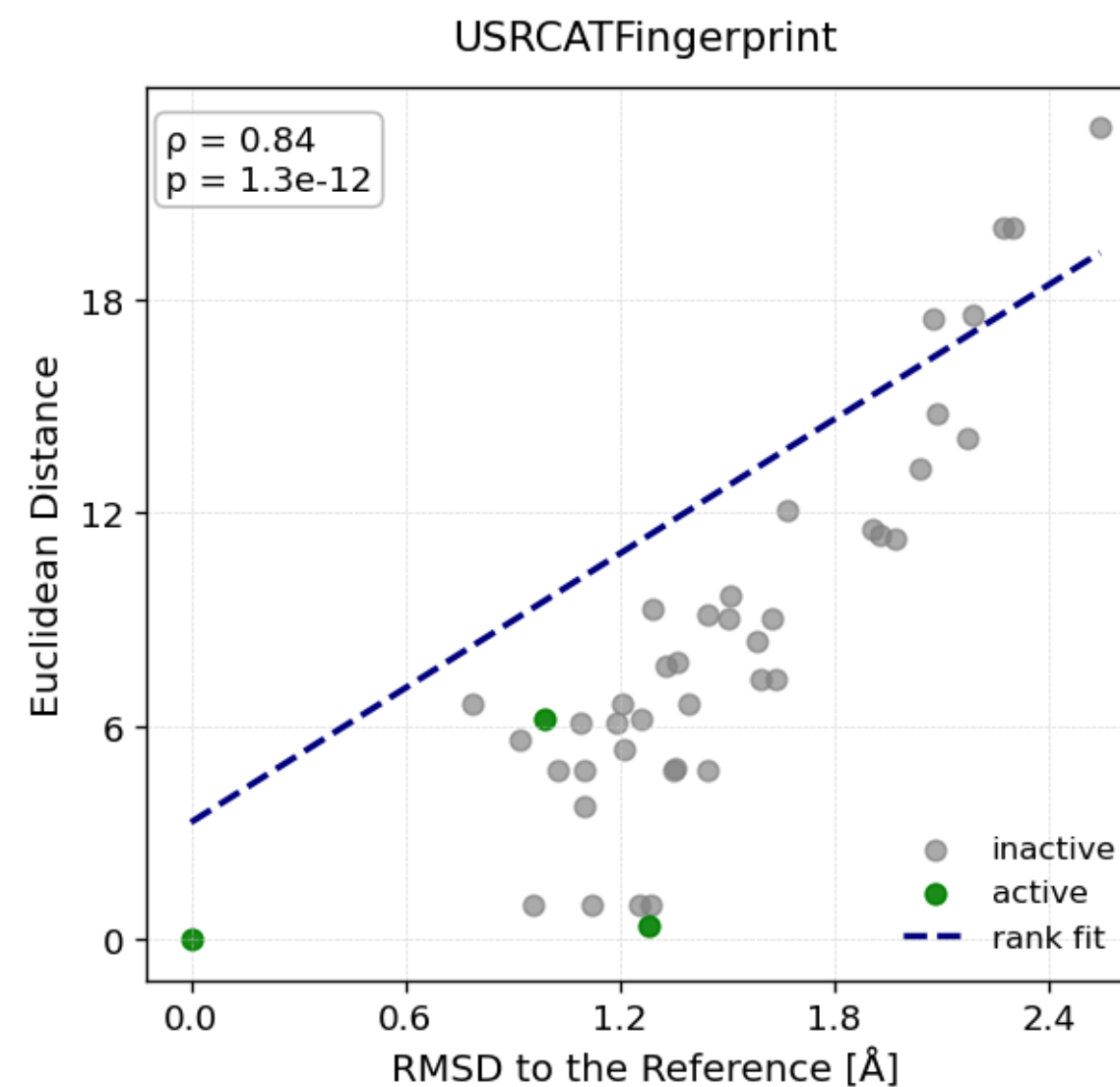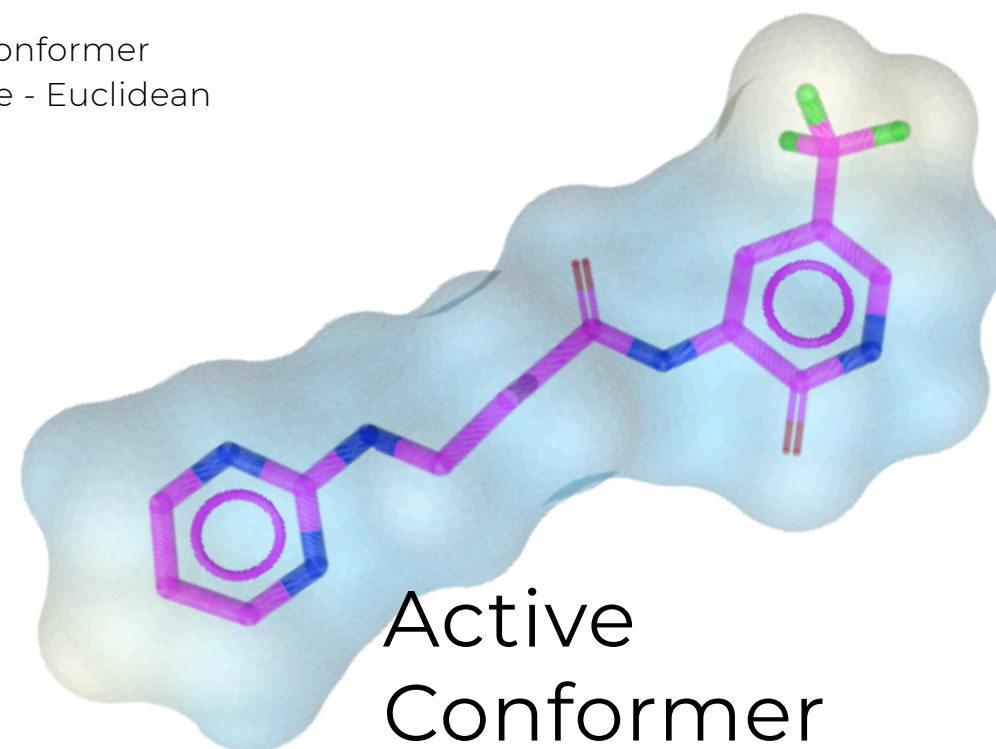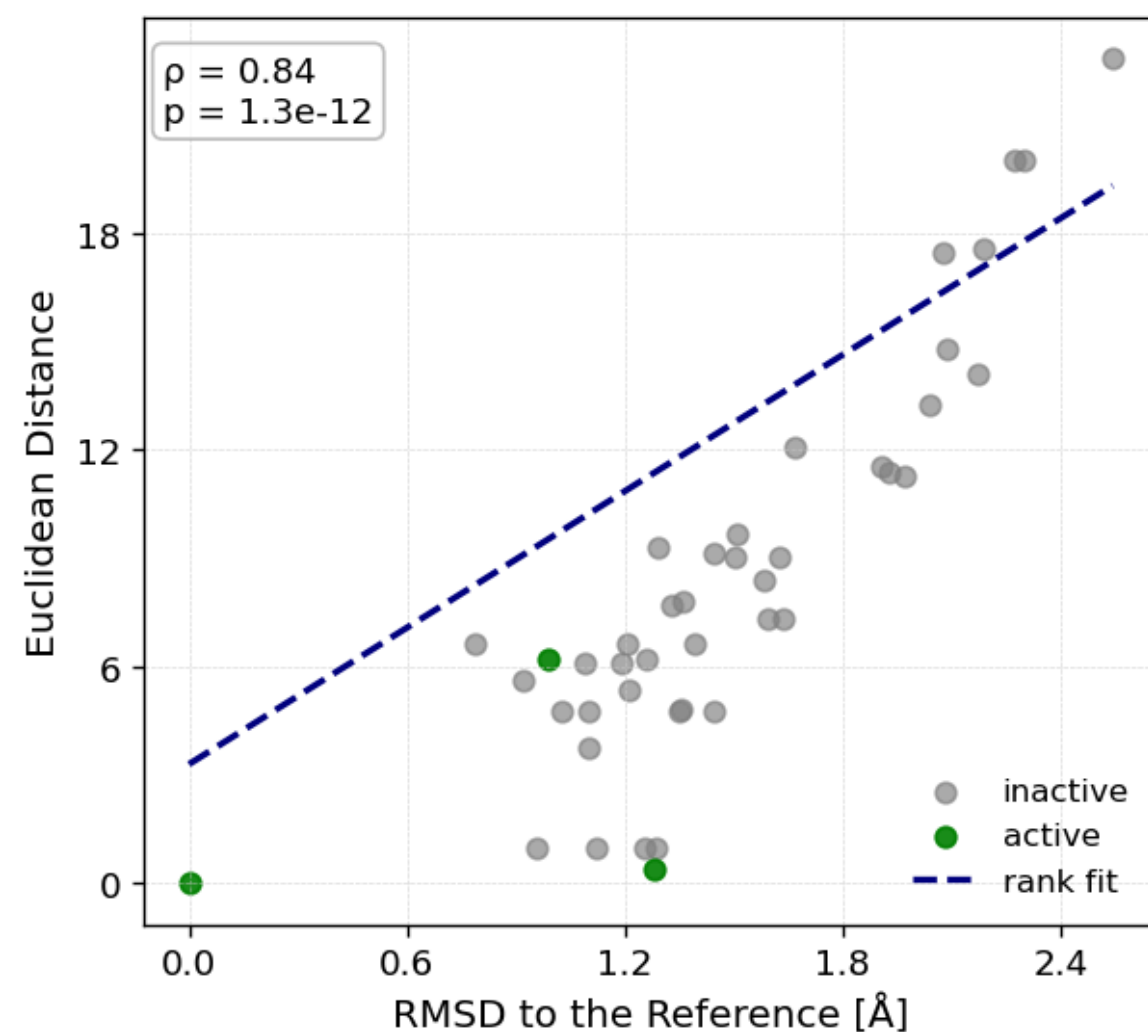
# MIL-Based QSAR Pipeline: Encoding Choice

*Reference - min energy active conformer
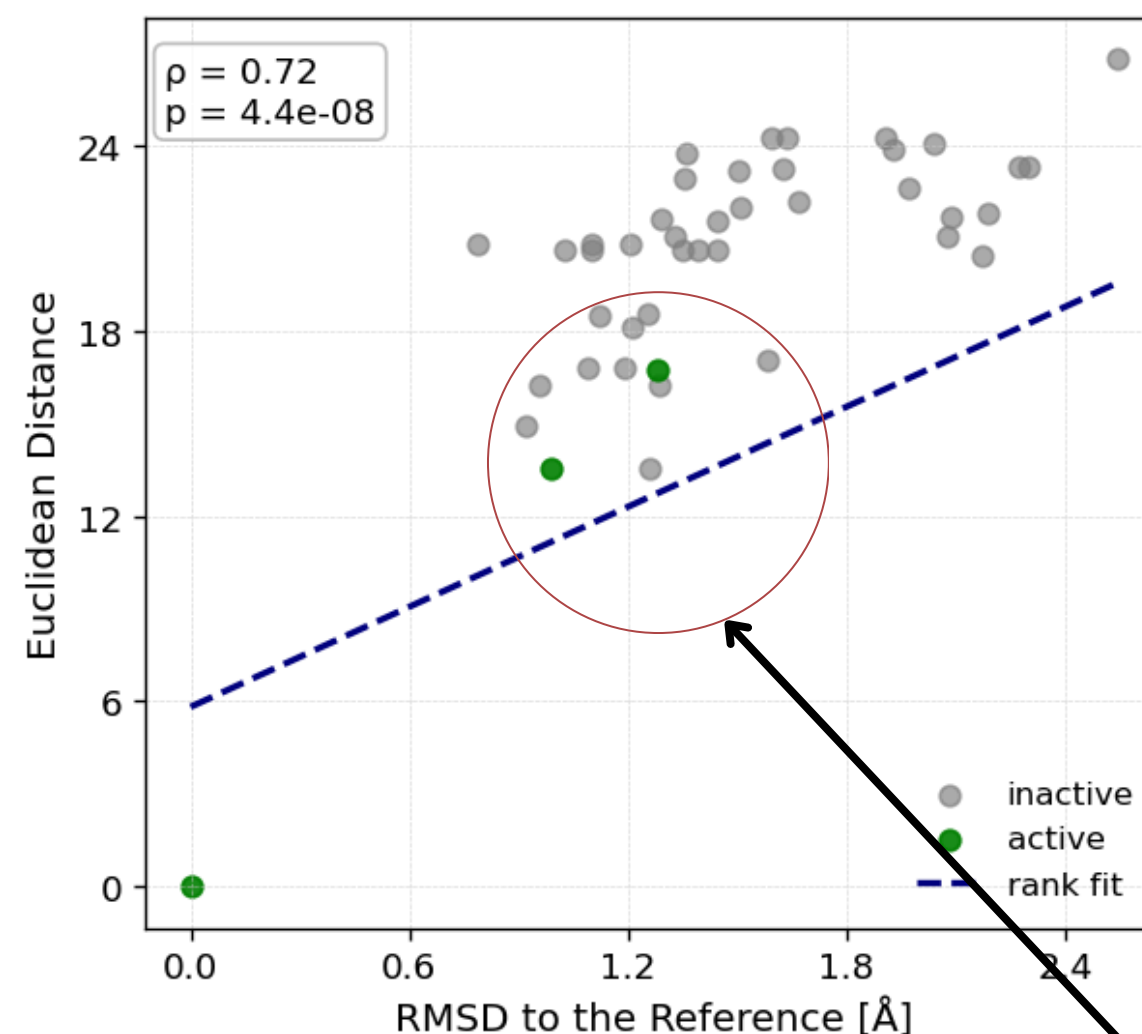**Distance metric in the feature space - Euclidean

Active Conformer

## USRCATFingerprint

$\rho = 0.84$
$p = 1.3e\text{-}12$

Euclidean Distance

RMSD to the Reference [Å]

- inactive
- active
- - rank fit

## 3D FP Distance vs RMSD [rmsd_thd=0.5Å]

Spearman $\rho$ (RMSD vs Distance)

ElectroShapeFingerprint, WHIMFingerprint, USRCATFingerprint, USRFingerprint, MORSEFingerprint, GETAWAYFingerprint, RDFFingerprint, E3FPFingerprint

What requirements do we have for the encoding we utilize?



Clustering of bioactive conformers

Reference - min energy active conformer
*Distance metric in the feature space - Euclidean
**RMSD threshold between conformers in the Bag = 0.5 Å
***All plots were made for the same molecule

(1KE8) [A] LS4299

LEU83A

GLU81A

ASP86A

HOH402A

LYS89A

PHE80A

ALA 144A

LEU 134A

ILE10A

ALA 31A    VAL18A

VAL64A

**2D Pharmacophore Model**

*7 Pharmacophore Features were considered
in **Virtual Screening (VS) Procedure

Target: **CDK2 (1KE8)**
Pharmacophore: *7 features (3HBA, 2 HBD, 2Ar)
**VS: 399 actives, 399 decoys

**Actives**: Min. 1 Conformer matches a Pharmacophore (SMI)
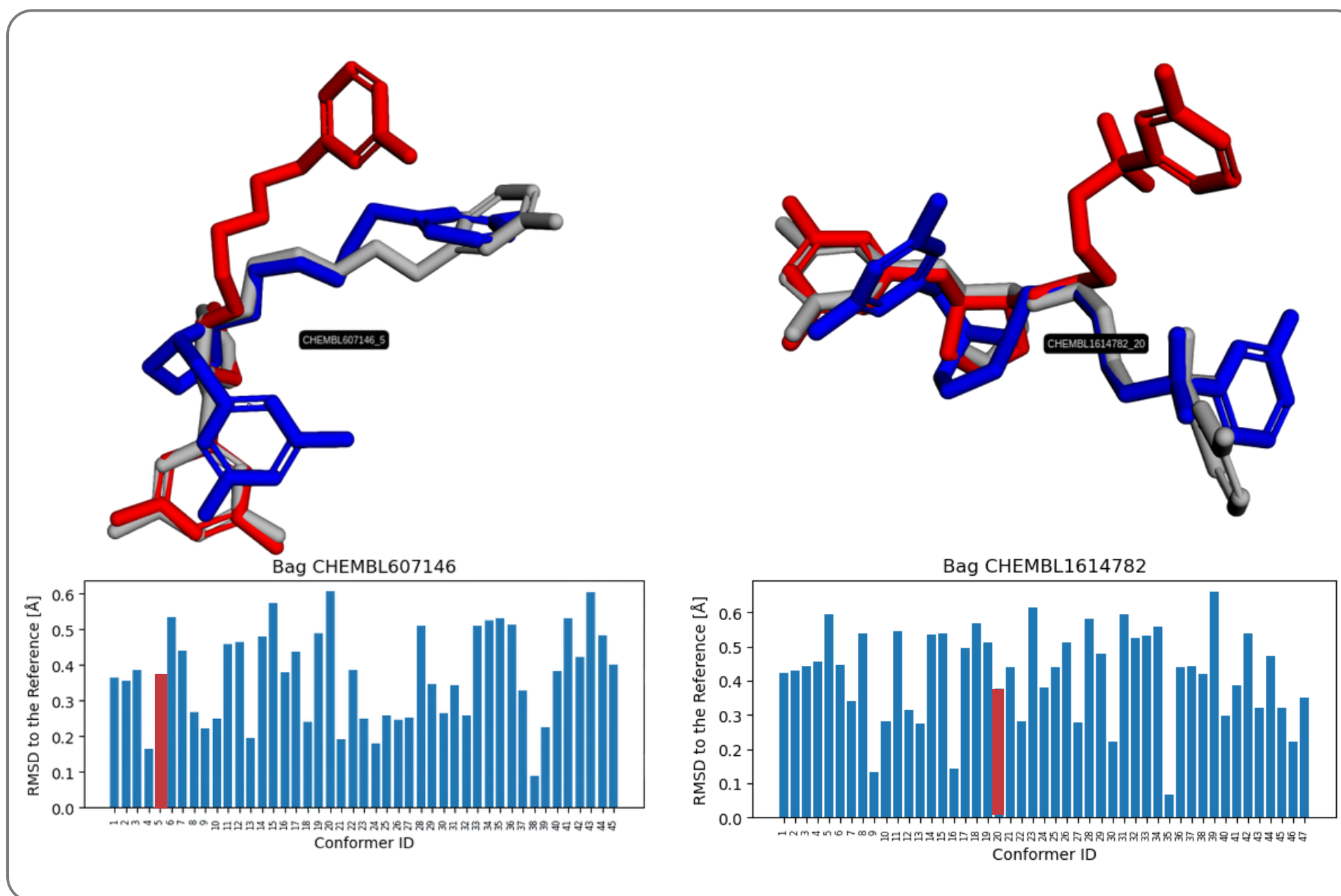**Decoys**: No conformers match the Pharmacophore



MIL Performance [CV ±95% CI]

Metric Value

Metrics
- accuracy
- precision
- recall
- f1
- roc_auc
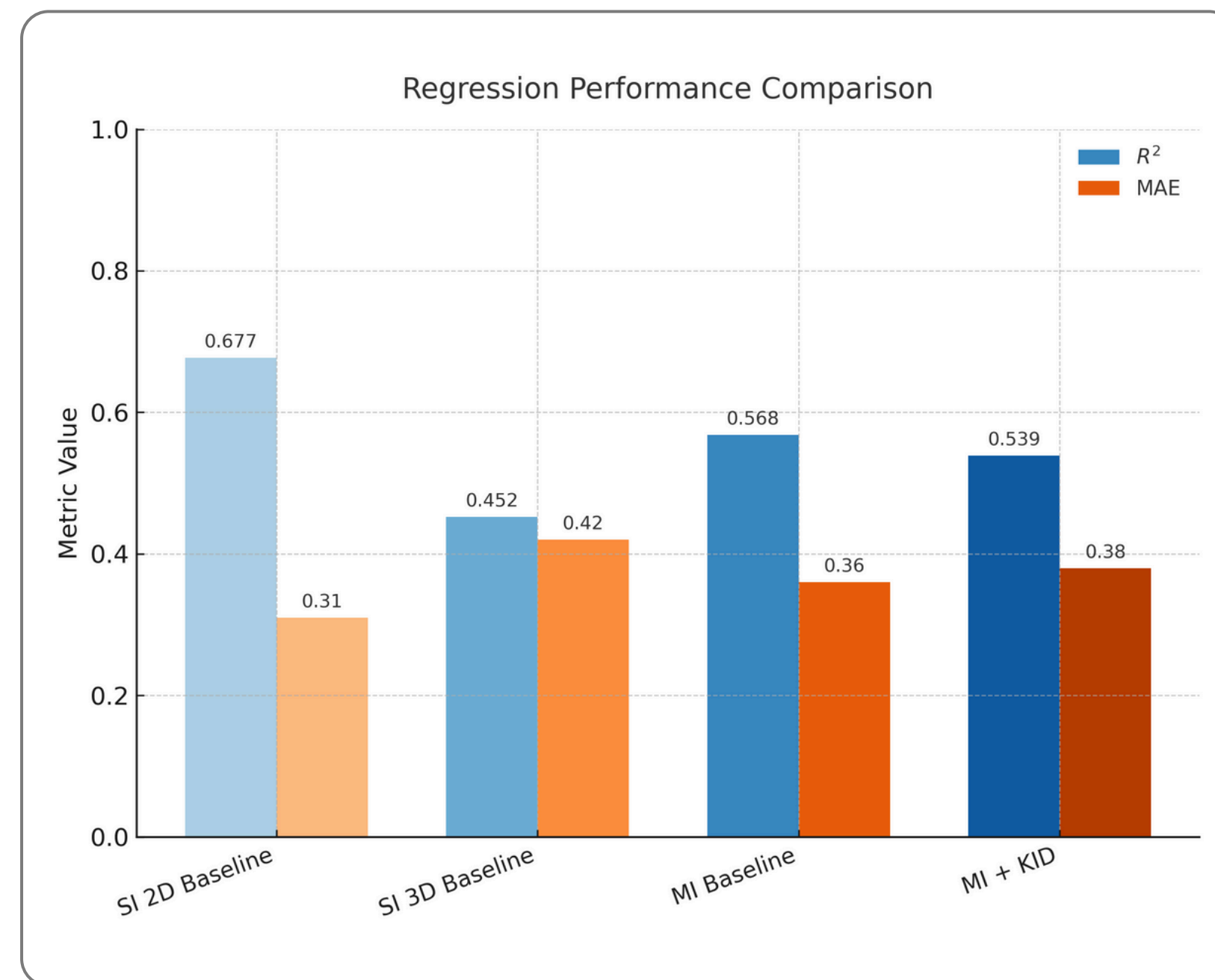- mAP
- MRR
- mean_auc
- mean_ndcg

**Multi-Head Attention MIL**

Regression task,
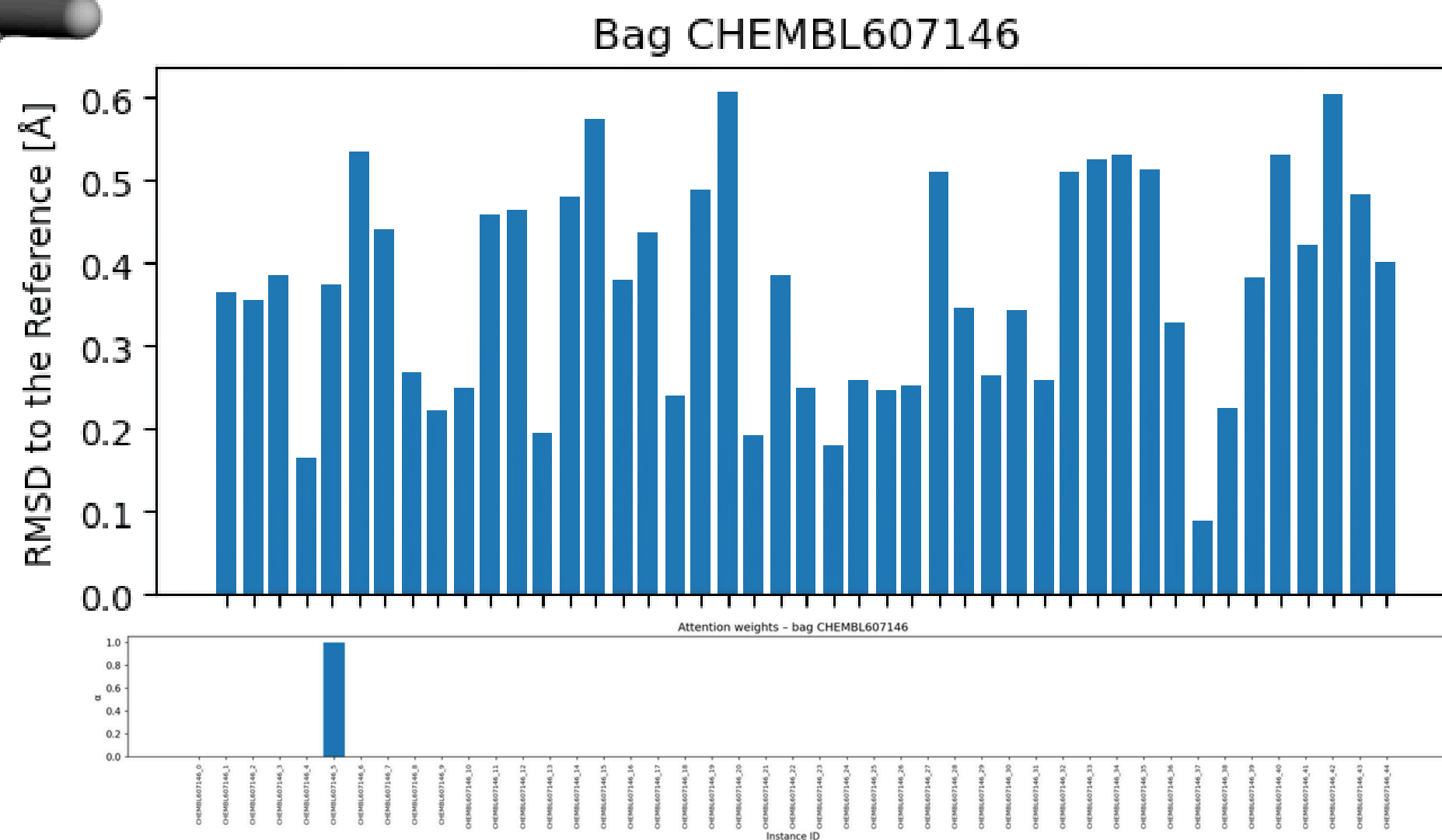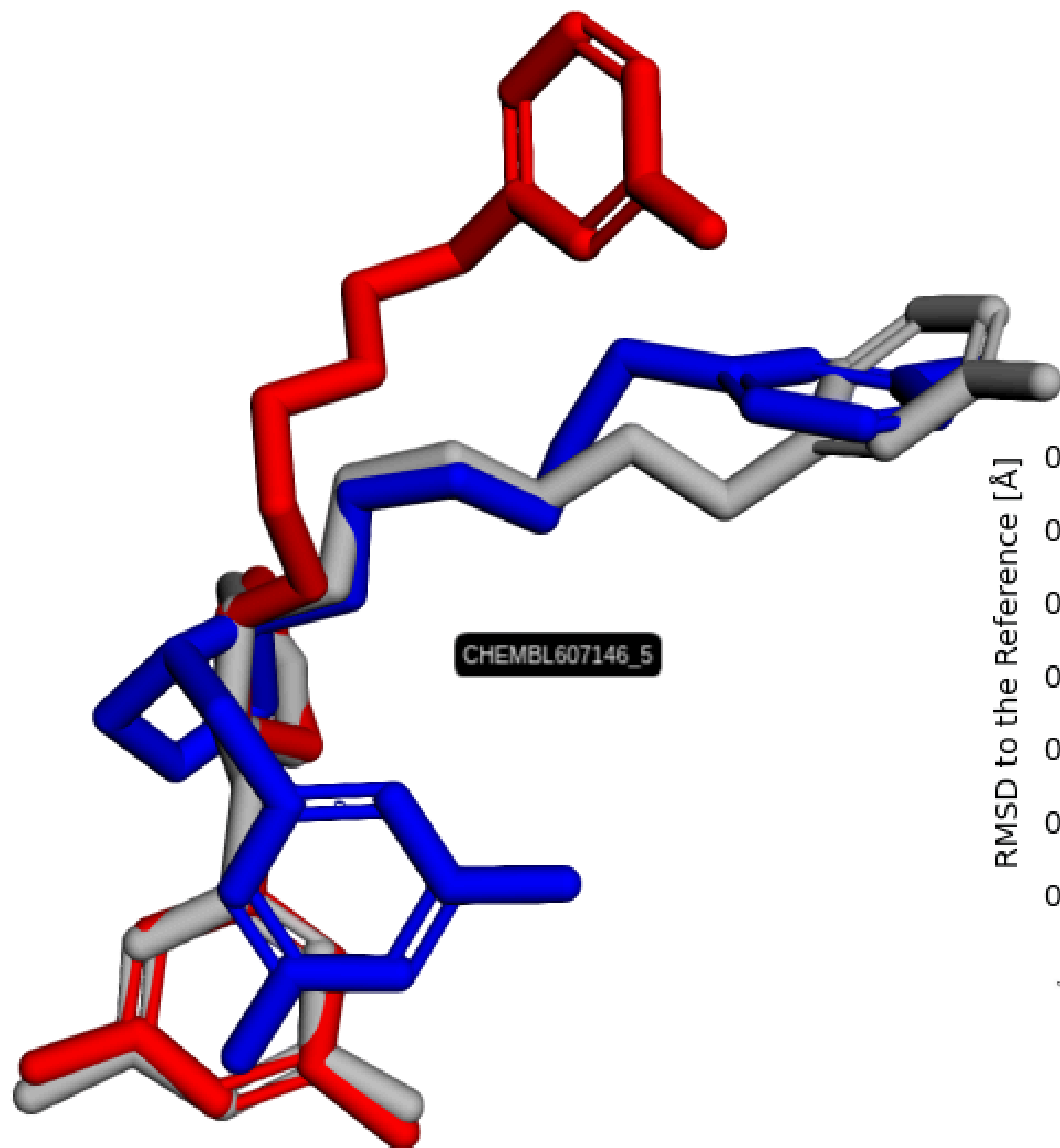target - nitric oxide synthase, CHEMBL3048



**Figure 1.** Key Instance Detection (KID). Top: Superposed conformers after RDKit GetO3A alignment: the lowest-energy conformer generated by RDKit ETKDGv3 (red), the experimentally observed conformer (grey), and the conformer identified as the key instance (blue). Bottom: Root-mean-square deviation (RMSD) values of all generated conformers to the experimental reference. The red bar highlights the conformer that received the highest attention weight.
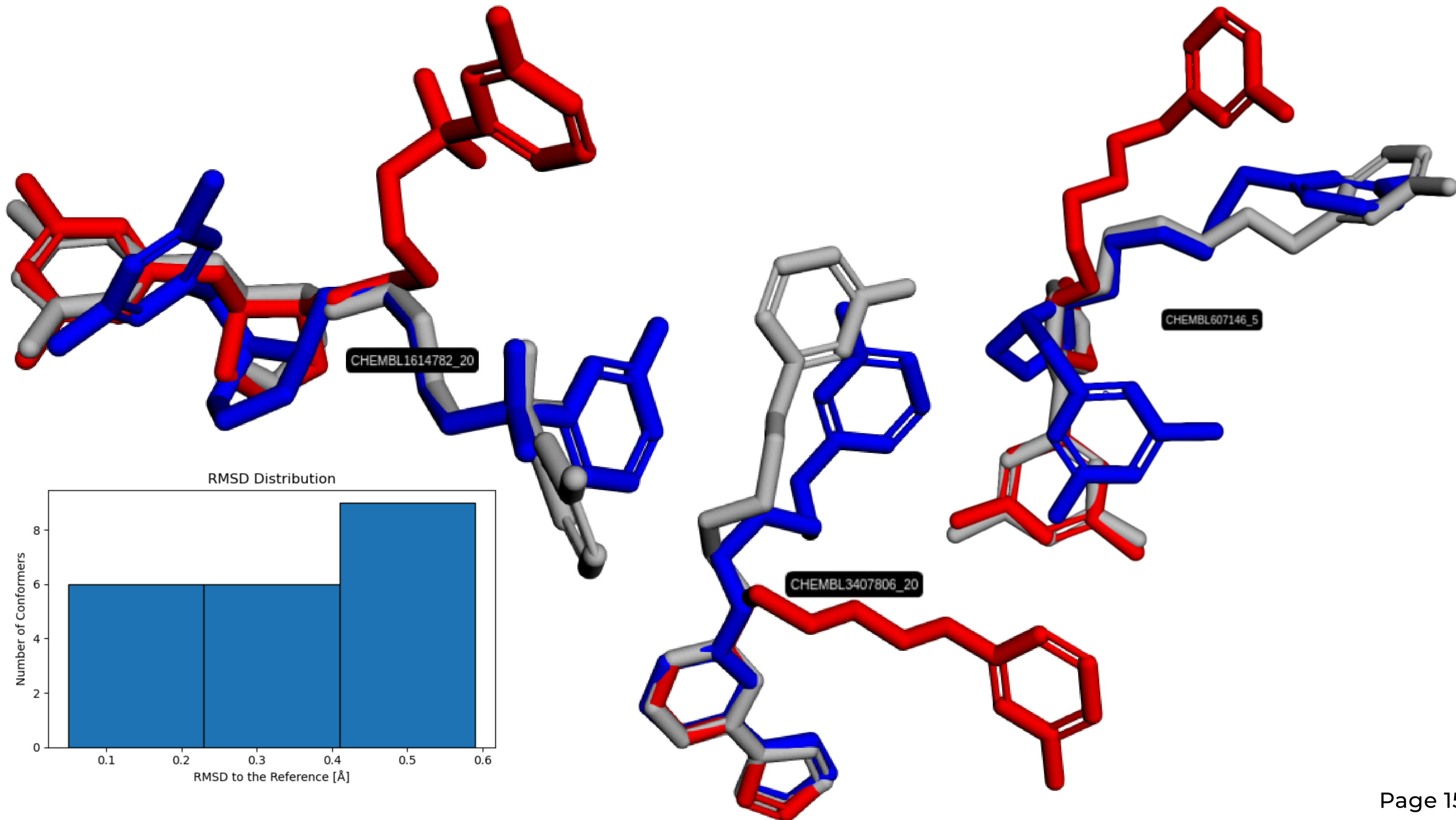


**Figure 2.** External validation accuracy. Regression task, target - nitric oxide synthase, CHEMBL3048: (1) SI-2D Baseline - SI RF regressor, 2D Morgan FP; (2) SI-3D Baseline - SI RF regressor trained on the minimum-energy conformers (one per molecule), 3D MoRSE FP; (3) MI Baseline - MI RF regressor; predictions obtained by mean-pooling 3D MoRSE descriptors across conformers whithin a bag; (4) MI + KID - developed method, described in the Pipeline section.

Abbreviations: RF, Random Forest; MAE, mean absolute error; MoRSE, Molecule Representation of Structures based on Electron diffraction

Bag CHEMBL607146

- Minimum-energy conformer (red)
- Experimentally observed conformer (grey)
- Predicted Key Instance (blue)

RMSD Distribution

## Conclusions & Outlook

- The model demonstrated the ability to prioritize conformers with lower RMSD to the reference bioactive structure, suggesting effective identification of relevant molecular shapes
- Validation of the model's attribution mechanism remains the primary goal of the project

[1] Dieterich, T. G., et al. Artificial Intelligence, 89(1–2), 31–71. (1997)

[2] Zankov, D. V., et al. Journal of Chemical Information and Modeling, 61(10), 4913–4923. (2021)

[3] Gomez, A., et al. arxiv.1706.03762. (2017)