# Scikit-Mol – RDKit integration in Scikit-Learn: Updates and New Features

# Brief Summary
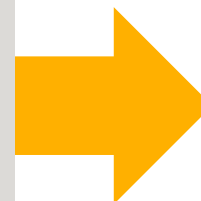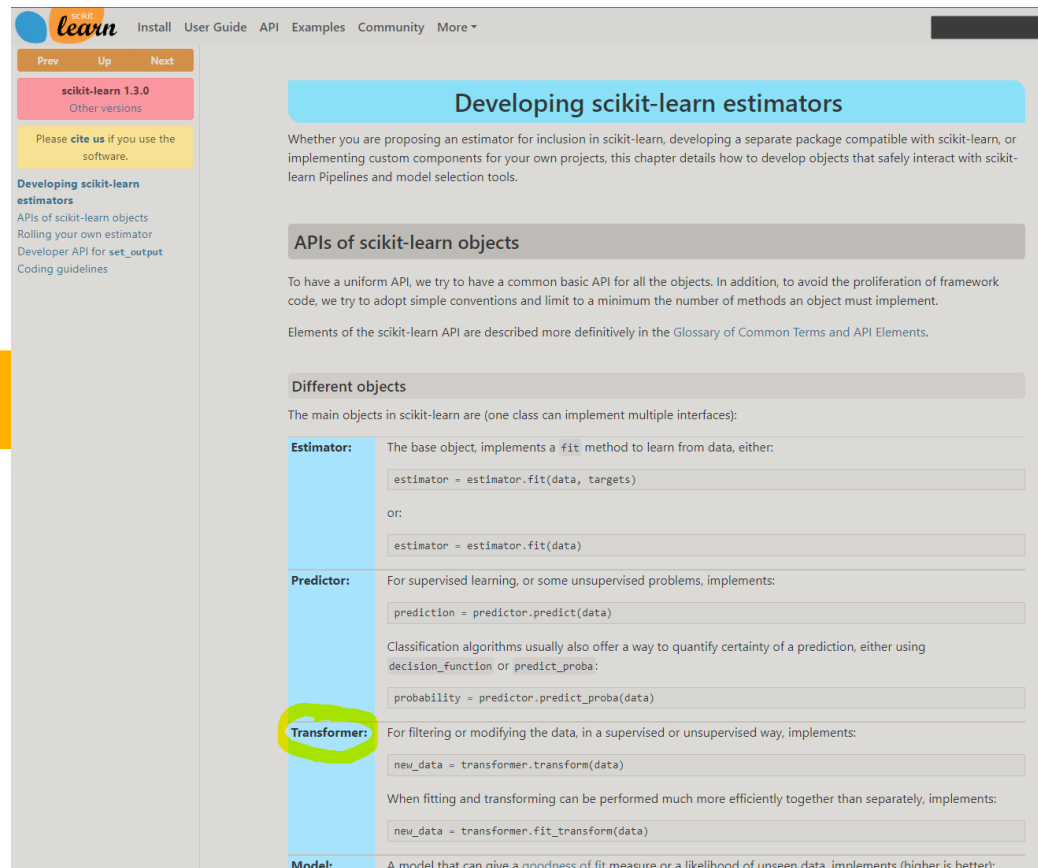
Project started at **2022 RDKit UGM Hackathon**

https://scikit-learn.org/stable/developers/develop.html

https://github.com/EBjerrum/scikit-mol

# Quick RDKit Integration into Scikit-Learn Pipelines

```python
from sklearn.linear_model import Ridge
from sklearn.pipeline import make_pipeline

from scikit_mol.conversions import SmilesToMolTransformer
from scikit_mol.fingerprints import MorganFingerprintTransformer

model = make_pipeline(SmilesToMolTransformer(),
                      MorganFingerprintTransformer(),
                      Ridge(alpha=10))
print(model)

model.fit(smiles_list_train, y_train)

print(f"Train score is :{model.score(smiles_list_train, y_train):0.2F}")
print(f"Test score is  :{model.score(smiles_list_test, y_test):0.2F}")
```
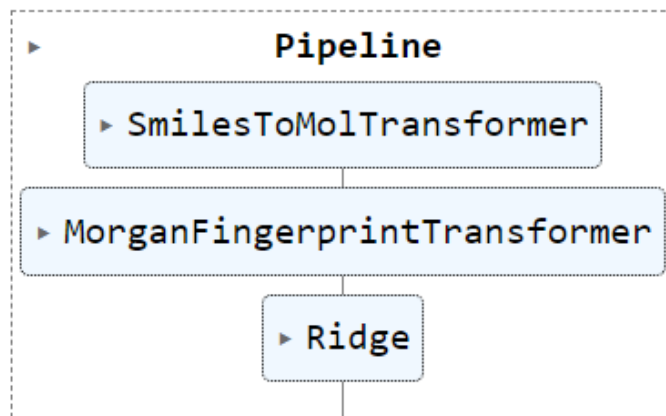
```
Pipeline(steps=[('smilestomoltransformer', SmilesToMolTransformer()),
                ('morganfingerprinttransformer',
                 MorganFingerprintTransformer()),
                ('ridge', Ridge(alpha=10))])
Train score is :0.74
Test score is  :0.63
```
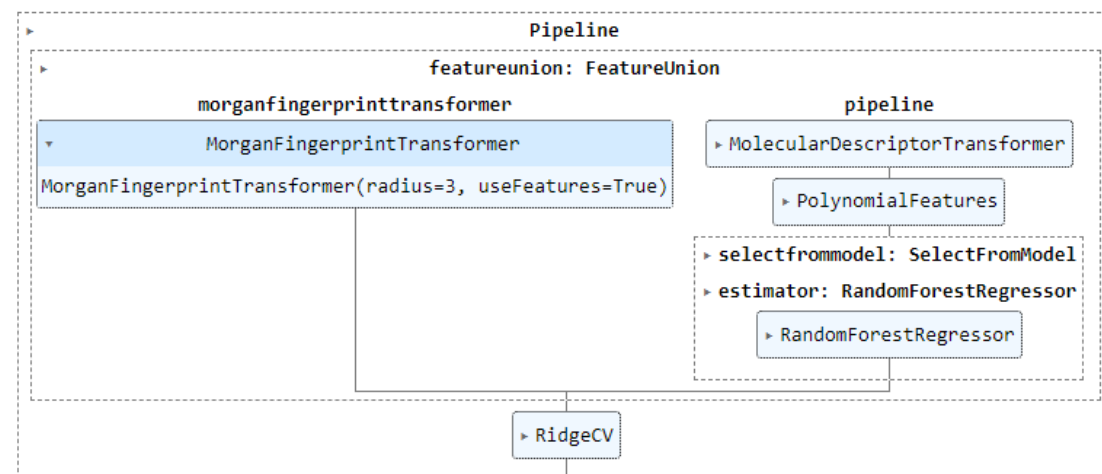
- Featureization from RDKit molecules or SMILES strings
- Integrate into Scikit-Learn models and pipelines
- Fit Scikit-Learn models directly on molecules or SMILES
- Easily get hyperparameter tuning of fingerprinting as compatible with e.g. SKOpt
- Self-documenting and self-contained models
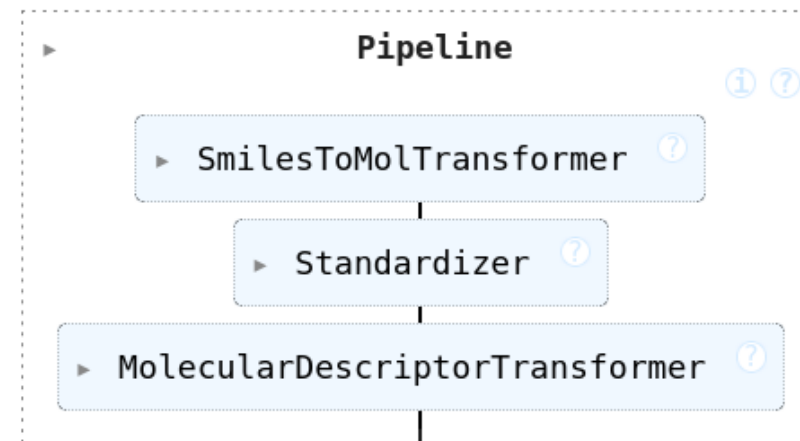- Use standard scikit-learn classes to build more complex models

# Maintenance and Project Improvements

- Online documentation[https://scikit-mol.readthedocs.io/en/latest/](https://scikit-mol.readthedocs.io/en/latest/)
- Switch to RDKit generators
- Conda package:[https://anaconda.org/conda-forge/scikit-mol](https://anaconda.org/conda-forge/scikit-mol)
- uv build framework

- Improved parallelism using joblib
- Improved handling of parallelism for ad-hoc custom feature transformers

# Pandas Input and Output

- `.set_output(transform="pandas")`
- Very handy when doing feature importance analysis!
- Column names are checked for consistency when predicting
- `model.predict(features[model.feature_names_in_])`



|  | MaxAbsEStateIndex | MaxEStateIndex | MinAbsEStateIndex | MinEStateIndex | qed | SPS | MolWt |
|---|---|---|---|---|---|---|---|
| 0 | 13.448610 | 13.448610 | 0.056985 | -0.432587 | 0.353101 | 14.289474 | 522.591980 |
| 1 | 12.863074 | 12.863074 | 0.026212 | -0.050849 | 0.682187 | 16.033333 | 425.558014 |
| 2 | 13.424788 | 13.424788 | 0.266700 | -0.413763 | 0.443905 | 15.852942 | 465.588013 |
| 3 | 12.725823 | 12.725823 | 0.052996 | -0.052996 | 0.577709 | 17.812500 | 478.467987 |
| 4 | 6.356910 | 6.356910 | 0.898244 | 0.898244 | 0.658108 | 13.052631 | 246.313004 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 154 | 6.217065 | 6.217065 | 0.175664 | 0.175664 | 0.916154 | 35.700001 | 312.239990 |
| 155 | 0.458245 | 0.458245 | 0.420212 | 0.420212 | 0.278112 | 21.714285 | 465.644080 |

# Safe Inference Mode

- Problematic molecules or SMILES ~~can~~ *will* be sent for models by ~~users~~GenAI

- **safe_inference** mode ensures that the whole batch doesn't fail

- Scikit-Mol transformers handle **"False" objects** (like False, 0, InvalidMol) gracefully

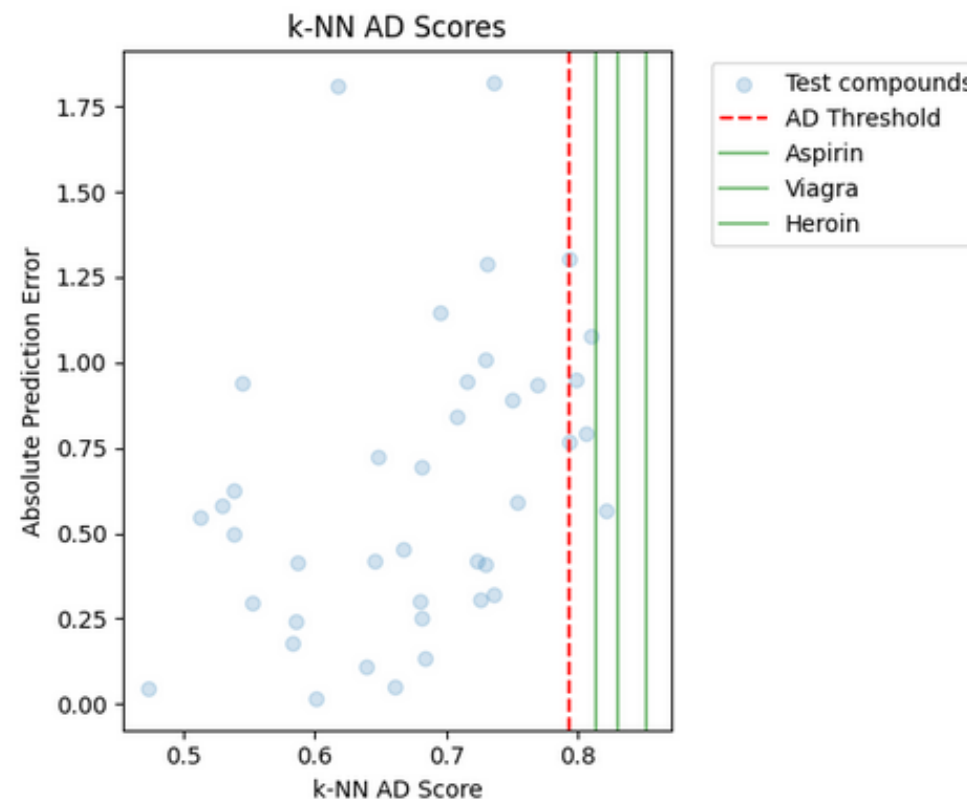- Scikit-learn estimators should be wrapped in **SafeInferenceWrapper**

```python
from scikit_mol.safeinference import set_safe_inference_mode
from sklearn.pipeline import Pipeline


pipe = Pipeline(
    [
        ("smi2mol", SmilesToMolTransformer()),
        ("mfp",
            MorganFingerprintTransformer(radius=2, fpSize=25)),
        ("safe_regressor",
        SafeInferenceWrapper(LogisticRegression())),
    ]
)

set_safe_inference_mode(pipe, True)
```

# Feature Based Applicability Domain Estimators

- **Applicability domain estimators** now available
- Can be fit to molecular set with a **percentile cutoff** (i.e. validation set)
- Slight deviation from sklearn api
  - `.transform()` → returns score
  - `.predict()` → returns 0 or 1
  - `.score_transform()` → returns 0 to 1 (soft boundary)
- bounding box, convex hull, hotelling, isolation forest, kernel density, kNN, leverage, local outlier, mahalanobis, topkat

# Conclusion and Acknowledgements

- Scikit-Mol is the original RDKit hackathon initiated project
- New features still being added as contributions from individuals, academia and industry
- Plenty of notebooks and examples at [https://scikit-mol.readthedocs.io/en/latest/](https://scikit-mol.readthedocs.io/en/latest/)
- **pip/conda install scikit-mol**

- **Thanks for your attention!**

E. J. Bjerrum et al., "Scikit-Mol brings cheminformatics to Scikit-Learn," ChemRxiv, Dec. 2023, doi: 10.26434/chemrxiv-2023-fzqwd.