

 HANDS-ON TUTORIAL

# RDKit and DeepChem

Comprehensive Guide to Cheminformatics and Machine Learning for Drug Discovery

## 1. Molecule Manipulation

---

### 1 Overview

Molecule manipulation is the foundation of computational chemistry. RDKit provides powerful tools for loading, parsing, and modifying molecular structures from various formats including SMILES, SDF, and MOL files. This enables researchers to programmatically work with chemical structures, perform substructure searches, and modify molecules for drug design.

# Key Operations



## Loading Molecules

Read molecular structures from SMILES strings, SDF files, or chemical databases



## Structure Analysis

Identify functional groups, ring systems, and chemical properties



## Modification

Add/remove atoms, modify bonds, and perform chemical transformations

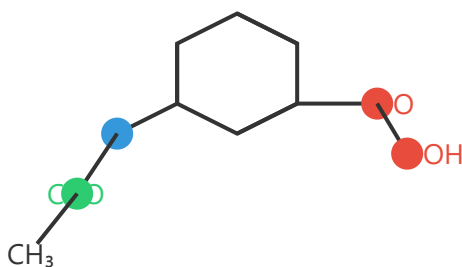


## Validation

Check molecular validity, sanitize structures, and ensure chemical correctness

## 3 Visual Example

### Aspirin (Acetylsalicylic Acid)



SMILES: CC(=O)Oc1ccccc1C(=O)O

## 4 Code Example

```
# Import RDKit library from rdkit import Chem from rdkit.Chem import Descriptors, AllChem # Load a molecule from SMILES smiles = 'CC(=O)Oc1ccccc1C(=O)O' # Aspirin mol = Chem.MolFromSmiles(smiles) # Validate and sanitize the molecule if mol is not None: Chem.SanitizeMol(mol) # Add explicit hydrogens mol_h = Chem.AddHs(mol) # Generate 3D coordinates AllChem.EmbedMolecule(mol_h) AllChem.MMFFOptimizeMolecule(mol_h) # Get molecular formula formula = Chem.rdMolDescriptors.CalcMolFormula(mol) print(f"Molecular Formula: {formula}") # C9H8O4 # Count atoms and bonds num_atoms = mol.GetNumAtoms() num_bonds = mol.GetNumBonds() print(f"Atoms: {num_atoms}, Bonds: {num_bonds}")
```

## 2. Descriptor Calculation

### 1 Overview

Molecular descriptors are numerical values that characterize chemical structures and their properties. These descriptors encode information about molecular size, shape, polarity, electronic properties, and more. They serve as input features for machine learning models in drug discovery, enabling quantitative structure-activity relationship (QSAR) studies.

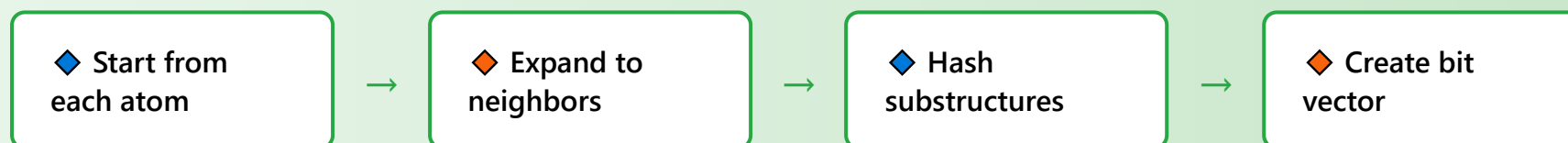
### 2 Types of Descriptors

Descriptor Type	Description	Examples
Physicochemical	Basic molecular properties	MW, LogP, TPSA, HBD/HBA

Descriptor Type	Description	Examples
<b>Topological</b>	Graph-based properties	Connectivity indices, Wiener index
<b>Fingerprints</b>	Structural bit/count vectors	ECFP, MACCS keys, Morgan FP
<b>3D Descriptors</b>	Geometry-dependent properties	PMI, Asphericity, Shape indices

### 3 Visual Example: Molecular Fingerprints

#### Circular Fingerprint (ECFP/Morgan) Generation



Morgan fingerprints encode molecular structure as fixed-length binary vectors (e.g., 2048 bits). Each bit represents the presence or absence of specific substructural features, enabling rapid similarity searches.

### 4 Example Calculations

```
# Calculate various molecular descriptors
from rdkit import Chem
from rdkit.Chem import Descriptors, AllChem

smiles = 'CC(=O)Oc1ccccc1C(=O)O'
mol = Chem.MolFromSmiles(smiles)

# Physicochemical descriptors
mw = Descriptors.MolWt(mol) # Molecular weight
logp = Descriptors.MolLogP(mol) # Lipophilicity
tpsa = Descriptors.TPSA(mol) # Topological polar surface area
hbd = Descriptors.NumHDonors(mol) # H-bond donors
hba = Descriptors.NumHAcceptors(mol) # H-bond acceptors
rotatable = Descriptors.NumRotatableBonds(mol) #
```

```
Rotatable bonds print(f"MW: {mw:.2f}, LogP: {logp:.2f}, TPSA: {tpsa:.2f}") print(f"HBD: {hbd}, HBA: {hba},  
RotBonds: {rotatable}") # Generate Morgan fingerprint (2048 bits, radius 2) fp =  
AllChem.GetMorganFingerprintAsBitVect(mol, radius=2, nBits=2048) print(f"Fingerprint: {fp.ToBitString()  
[:50]}...") # Calculate molecular similarity using Tanimoto coefficient smiles2 =  
'CC(C)Cc1ccc(cc1)C(C)C(=O)O' # Ibuprofen mol2 = Chem.MolFromSmiles(smiles2) fp2 =  
AllChem.GetMorganFingerprintAsBitVect(mol2, radius=2, nBits=2048) similarity =  
DataStructs.TanimotoSimilarity(fp, fp2) print(f"Tanimoto Similarity: {similarity:.3f}")
```

💡 **Lipinski's Rule of Five:** Drug-like molecules typically have  $MW < 500$ ,  $LogP < 5$ ,  $HBD \leq 5$ ,  $HBA \leq 10$ . These descriptors help filter compound libraries for oral bioavailability.

## 3. Model Training with DeepChem

---

### 1 Overview

DeepChem is a powerful framework for building machine learning models in drug discovery. It provides pre-built architectures including graph convolutional networks (GCNs), message passing neural networks (MPNNs), and traditional ML models. DeepChem handles featurization, model training, and evaluation with minimal code, making it ideal for cheminformatics applications.

### 2 Available Model Architectures

### **Graph Convolutional Networks**

Learn from molecular graph structure directly without handcrafted features

### **Message Passing Networks**

Aggregate information from neighboring atoms through iterative message passing

### **Random Forest**

Ensemble of decision trees, robust and interpretable for QSAR

### **Multi-Task Networks**

Simultaneously predict multiple molecular properties with shared representations

## 3 Training Pipeline

 Load Data



 Featurize Molecules



 Split Dataset



 Train Model



 Evaluate Performance

## 4 Code Example

```
# Complete model training example with DeepChem
import deepchem as dc
from deepchem.models import GraphConvModel
import numpy as np

# 1. Load and featurize dataset
tasks = ['activity']
featurizer = dc.feat.ConvMolFeaturizer()
# Graph-based features
loader = dc.data.CSVLoader( tasks=tasks,
                             feature_field='smiles',
                             featurizer=featurizer )
dataset = loader.create_dataset('molecules.csv')

# 2. Split into train/validation/test sets
splitter = dc.splits.RandomSplitter()
train_dataset, valid_dataset, test_dataset = splitter.train_valid_test_split(
    dataset, frac_train=0.8, frac_valid=0.1, frac_test=0.1 )

# 3. Initialize and train model
model = GraphConvModel( n_tasks=len(tasks), mode='classification',
                        batch_size=64, learning_rate=0.001 )
# Train for 50 epochs
model.fit(train_dataset, nb_epoch=50)

# 4. Evaluate model performance
metric = dc.metrics.Metric(dc.metrics.roc_auc_score)
train_score = model.evaluate(train_dataset, [metric])
valid_score = model.evaluate(valid_dataset, [metric])
test_score = model.evaluate(test_dataset, [metric])

print(f"Training ROC-AUC: {train_score['roc_auc_score']:.3f}")
print(f"Validation ROC-AUC: {valid_score['roc_auc_score']:.3f}")
print(f"Test ROC-AUC: {test_score['roc_auc_score']:.3f}")

# 5. Make predictions on new molecules
predictions = model.predict(test_dataset)
```

## 4. Scaffold Splitting Strategy

### 1 Overview

Scaffold splitting is a crucial technique for creating realistic train/test splits in drug discovery. Unlike random splitting, scaffold splitting ensures that molecules with the same core structure (Bemis-Murcko scaffold) are grouped together. This prevents data leakage and provides a more rigorous test of a model's ability to generalize to novel chemical scaffolds, simulating real-world drug discovery scenarios.

## 2 Why Scaffold Splitting Matters

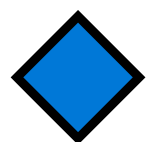
### ✗ Problem with Random Splitting

Random splitting can place similar molecules in both training and test sets, leading to overestimated model performance. The model may memorize structural patterns rather than learning generalizable chemical relationships.

### ✓ Advantage of Scaffold Splitting

Scaffold splitting ensures that test set molecules have different core structures from training molecules, providing a realistic evaluation of model generalization to novel chemical space.

## 3 Visual Example



**Scaffold A**

Benzene-based  
compounds



**Scaffold B**

Pyridine-based  
compounds



**Scaffold C**

Indole-based  
compounds

**Training Set**

**Validation Set**

**Test Set**

## 4 Implementation

```
# Implement scaffold splitting with DeepChem
import deepchem as dc from rdkit import Chem from
rdkit.Chem.Scaffolds import MurckoScaffold # Load dataset
tasks = ['activity'] featurizer =
```



```
dc.featurizer.CircularFingerprint(size=2048) loader = dc.data.CSVLoader( tasks=tasks, feature_field='smiles',
featurizer=featurizer ) dataset = loader.create_dataset('compounds.csv') # Apply scaffold splitting
scaffoldsplitter = dc.splits.ScaffoldSplitter() train, valid, test = scaffoldsplitter.train_valid_test_split(
dataset, frac_train=0.8, frac_valid=0.1, frac_test=0.1 ) print(f"Training set: {len(train)} molecules")
print(f"Validation set: {len(valid)} molecules") print(f"Test set: {len(test)} molecules") # Example: Extract
Bemis-Murcko scaffold from a molecule smiles = 'CCc1ccc(cc1)C(C)C(=O)O' mol = Chem.MolFromSmiles(smiles)
scaffold = MurckoScaffold.GetScaffoldForMol(mol) scaffold_smiles = Chem.MolToSmiles(scaffold)
print(f"Original SMILES: {smiles}") print(f"Scaffold SMILES: {scaffold_smiles}") # Compare with random
splitting randomsplitter = dc.splits.RandomSplitter() train_rand, valid_rand, test_rand =
randomsplitter.train_valid_test_split( dataset, frac_train=0.8, frac_valid=0.1, frac_test=0.1 )
```

**⚠ Important:** Models trained with scaffold splitting typically show lower test performance than those with random splitting. This is expected and reflects the true generalization capability to novel chemical structures.

## 5. Performance Evaluation

### 1 Overview

Proper model evaluation is critical in drug discovery to ensure reliable predictions. Different metrics are appropriate for different tasks: classification (active/inactive compounds), regression (binding affinity, IC50 values), or multi-task predictions. Selecting the right metrics and understanding their implications helps researchers make informed decisions about model deployment.

# Key Metrics for Different Tasks

Task Type	Recommended Metrics	Interpretation
Classification	ROC-AUC, Precision, Recall, F1-Score	Ability to distinguish active vs inactive compounds
Regression	RMSE, MAE, R <sup>2</sup> , Pearson Correlation	Accuracy of continuous property prediction
Ranking	Enrichment Factor, BEDROC	Early recognition of active compounds
Imbalanced Data	Balanced Accuracy, MCC, PR-AUC	Performance when classes are unequal

## 3 Example Performance Dashboard

### Model Performance Summary

ROC-AUC Score

0.87

Accuracy

82%

Precision

0.79

Recall

0.85

F1-Score

0.82

MCC

0.64

## 4 Evaluation Code

```
# Comprehensive model evaluation with multiple metrics
import deepchem as dc from sklearn.metrics import (
    roc_auc_score, accuracy_score, precision_score, recall_score, f1_score, matthews_corrcoef, confusion_matrix,
    classification_report )
import numpy as np # Assume model and test_dataset are already defined
# 1. Get predictions
y_pred_proba = model.predict(test_dataset)
y_pred = (y_pred_proba > 0.5).astype(int)
y_true = test_dataset.y
# 2. Calculate classification metrics
roc_auc = roc_auc_score(y_true, y_pred_proba)
accuracy = accuracy_score(y_true, y_pred)
precision = precision_score(y_true, y_pred)
recall = recall_score(y_true,
```

```

y_pred) f1 = f1_score(y_true, y_pred) mcc = matthews_corrcoef(y_true, y_pred) print("=== Classification
Performance ===") print(f"ROC-AUC Score: {roc_auc:.3f}") print(f"Accuracy: {accuracy:.3f}")
print(f"Precision: {precision:.3f}") print(f"Recall: {recall:.3f}") print(f"F1-Score: {f1:.3f}") print(f"MCC:
{mcc:.3f}") # 3. Confusion matrix cm = confusion_matrix(y_true, y_pred) print("\n=== Confusion Matrix ===")
print(f"True Negatives: {cm[0,0]}") print(f"False Positives: {cm[0,1]}") print(f"False Negatives: {cm[1,0]}")
print(f"True Positives: {cm[1,1]}") # 4. Detailed classification report print("\n=== Detailed Report ===")
print(classification_report(y_true, y_pred, target_names=['Inactive', 'Active'])) # 5. For regression tasks
from sklearn.metrics import mean_squared_error, mean_absolute_error, r2_score # If predicting continuous
values (e.g., IC50, binding affinity) # rmse = np.sqrt(mean_squared_error(y_true, y_pred)) # mae =
mean_absolute_error(y_true, y_pred) # r2 = r2_score(y_true, y_pred) # 6. Cross-validation for robust
evaluation from sklearn.model_selection import cross_val_score # cv_scores = cross_val_score(model, X, y,
cv=5, scoring='roc_auc') # print(f"CV ROC-AUC: {cv_scores.mean():.3f} (+/- {cv_scores.std():.3f})")

```

## 5 Interpretation Guidelines



### ROC-AUC > 0.8

Excellent discrimination between active and inactive compounds



### High Precision

Few false positives - important for reducing experimental cost



### High Recall

Few false negatives - important for not missing active compounds



### MCC > 0.5

Good balance, especially valuable for imbalanced datasets



**Best Practice:** Always evaluate models on multiple metrics and compare performance across different splitting strategies (random vs scaffold). Report confidence intervals using cross-validation or bootstrap resampling.



# Summary and Best Practices

---

## 1 Complete Workflow Integration

### 1 Molecule Manipulation

Load & validate structures



### 2 Descriptor Calculation

Compute features & fingerprints



### 3 Model Training

Build ML/DL models



### 4 Scaffold Splitting

Create realistic test sets



## 5 Performance Evaluation

Validate & interpret results

## 2 Key Takeaways

### ✓ Do's

- ✓ Sanitize molecules before analysis
- ✓ Use scaffold splitting for evaluation
- ✓ Calculate multiple descriptors
- ✓ Report multiple evaluation metrics
- ✓ Validate with cross-validation

### ✗ Don'ts

- ✗ Rely only on random splitting
- ✗ Use single metric for evaluation
- ✗ Ignore data imbalance issues
- ✗ Skip molecule validation steps
- ✗ Overfit to training data

## 3 Additional Resources



### Recommended Learning Materials

- **RDKit Documentation:** [rdkit.org/docs](https://rdkit.org/docs) - Comprehensive API reference
- **DeepChem Tutorials:** [deepchem.io/tutorials](https://deepchem.io/tutorials) - Hands-on examples
- **MoleculeNet:** Large-scale benchmark datasets for molecular ML
- **Papers:** "Molecular graph convolutions" (Duvenaud et al., 2015)
- **Community:** RDKit mailing list and DeepChem forum



## Ready to Start Your Drug Discovery Journey!

Combine RDKit's molecular manipulation with DeepChem's machine learning capabilities to accelerate your research