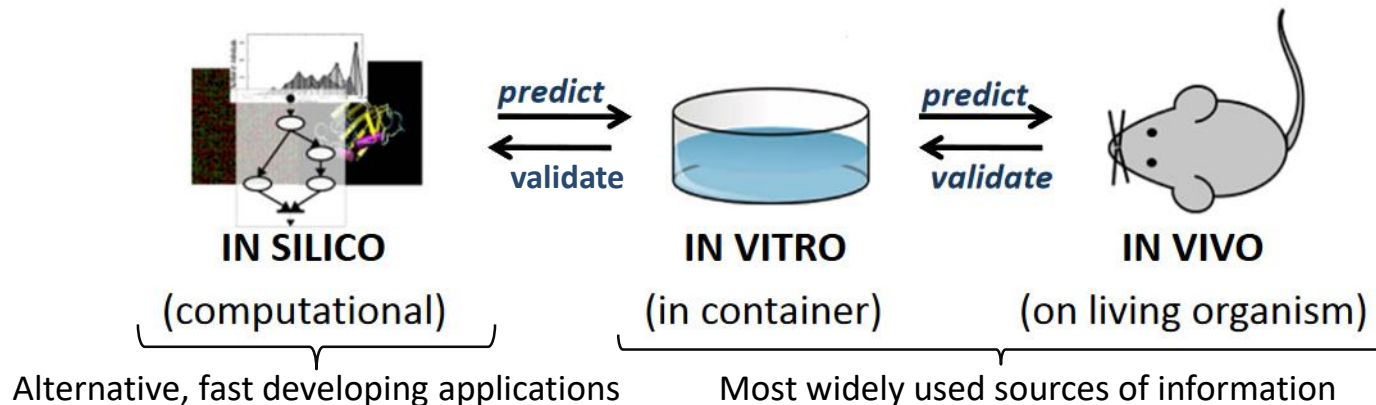


- All chemical substances need to be tested in terms of their toxicological, environmental or pharmaceutical properties before their use
  - Identify harmful effects on humans, animals, plants and the environment
  - Regulatory agencies require *in-vivo* testing for several toxic endpoints
  - ~ 2.9 million laboratory animals (Germany, 2010) and rising
- Establishment of alternative methods, reduction of animal testing



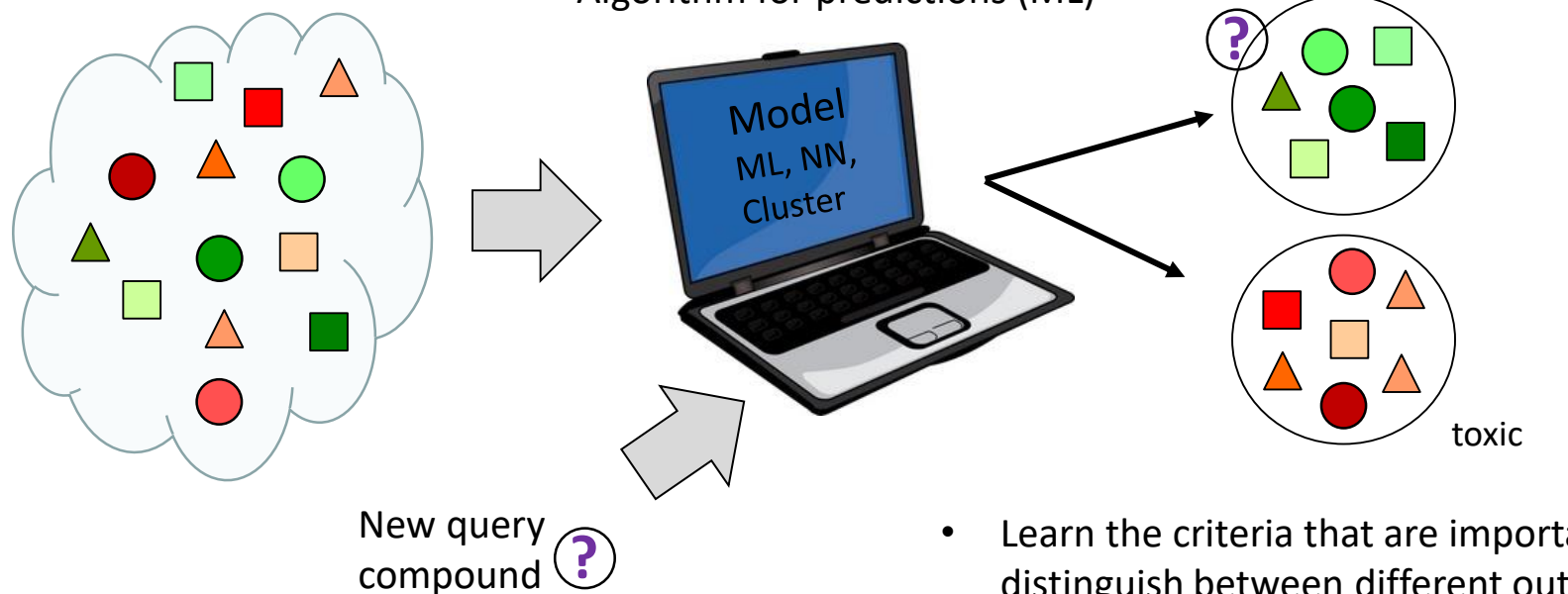
- Predict likely toxic compounds (QSAR, ADMETox, Machine learning)
- Reduce animal testing: a-priori exclusion from animal testing

Assumption: Structurally similar molecules have similar properties  
and, thus, similar biological activity

Library of known toxic and non-toxic compounds against a specific toxic endpoint

Requirements:

- Machine representation of the molecules
- Quantitative measure of similarity between the molecular descriptors
- Algorithm for predictions (ML)



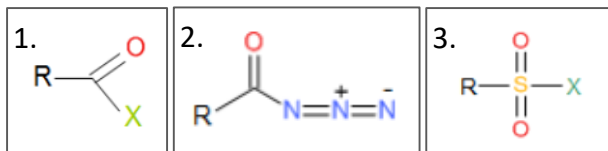
## 1. Data set

- 10,000 compounds x 12 toxic endpoints (Tox21 challenge)
- Fingerprint representation

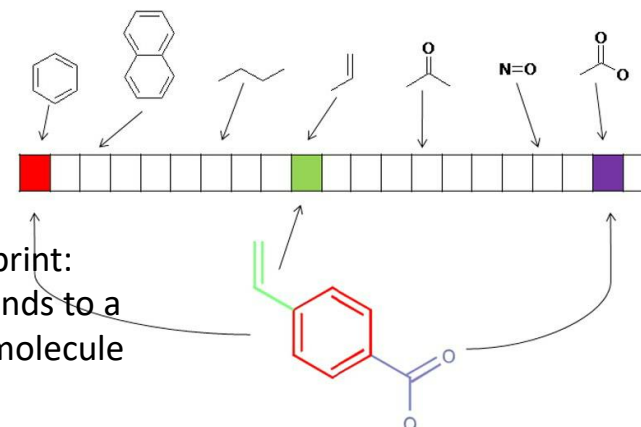
## 2. Methods

- Machine learning (ML), i.e., support vector machine, random forest, neuronal nets, deep learning
- Split test and training data
- Model evaluation

## 3. Toxicity prediction on new compounds (and extraction of novel toxicophores)

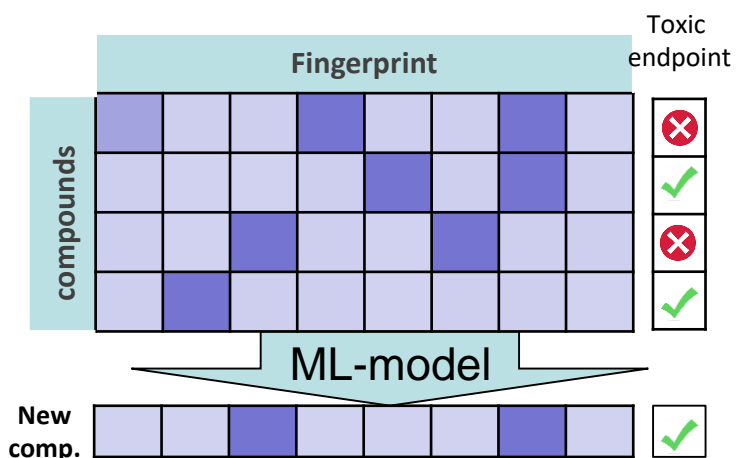


Skin sensitization: Acid halides, acid azides, sulphonyl halides



Molecular fingerprint:  
Each bit corresponds to a fragment of the molecule

Data matrix





### Prepare data set

- Collect compounds (Tox21 data set)
- Convert compounds to fingerprint (MACCS or other)
- Collect classification per compound (toxic/non-toxic)

Data matrix

		Fingerprint							Toxic endpoint
compounds									✗
									✓
									✗
									✓

### Machine learning, e.g. Random forest

#### Split data in train and test set

```
In [24]: from sklearn.cross_validation import train_test_split
X_train, X_test, y_train, y_test = train_test_split(fps, y, test_size=0.20)
```

#### Train the model

See <http://scikit-learn.org/stable/modules/generated/sklearn.ensemble.RandomForestClassifier.html> for an explanation of the parameter

```
In [22]: from sklearn.ensemble import RandomForestClassifier
from sklearn import metrics

forest = RandomForestClassifier(n_jobs=-1, n_estimators=100, max_features=100)
forest.fit(X_train, y_train) # Build a forest of trees from the training set
```

#### Test performance of the model

```
In [23]: y_pred = forest.predict(X_test) # Predict class for X
accuracy = metrics.accuracy_score(y_test, y_pred)
roc_auc = metrics.roc_auc_score(y_test, y_pred)
print "Accuracy: %.2f; AUC: %.2f" % (accuracy, roc_auc)
```

Accuracy: 0.90; AUC: 0.92

## ■ Query page

- Input structure of interest
- Smiles/mol format

## ■ Result page

- Predicted toxicity
- Most similar compounds
- Identified toxicophores

Smiles string      Go!

Or draw compound

Marvin JS

H  
C  
N  
O  
S  
F

Query:

Predicted Toxicity

Similar compounds from data set

Identified toxicophores and endpoints

- Insides into structural bioinformatics/cheminformatics
- Programming
  - Mainly Python language
  - Using libraries such as RDKit, sklearn, pandas, ...
  - Some Web design
- Soft skills
  - Project management
  - Presentation
  - Scientific writing
- Work place
  - Campus Charité Mitte

Contact:

[andrea.volkamer@charite.de](mailto:andrea.volkamer@charite.de)