MELLODDY-TUNER:
Data Standardization Framework for
Federated Machine Learning

**RDKit UGM 2020**
**Lukas Friedrich (Merck KGaA, Darmstadt)**

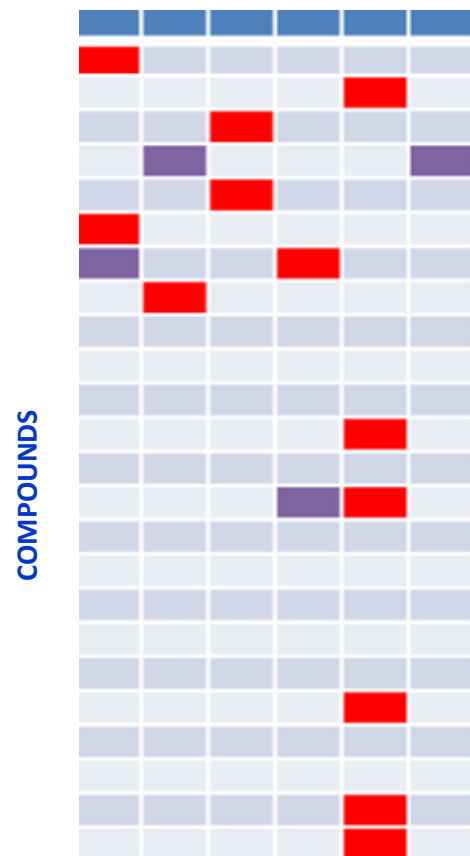**MACHINE LEARNING LEDGER ORCHESTRATION FOR DRUG**
**DISCOVERY**

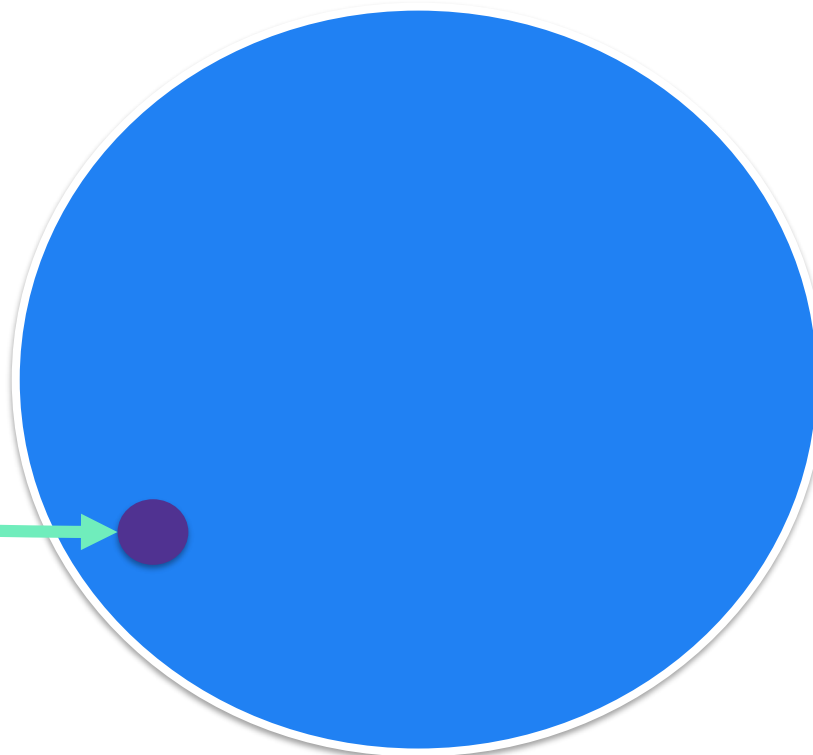# PREDICTIVE MODELING IN DRUG DISCOVERY

**ASSAYS**

**COMPOUNDS**

**Chemical space**

**Machine Learning (ML) Model**

» Model performance and applicability depend on amount of high quality data

# MACHINE LEARNING LEDGER
# ORCHESTRATION FOR DRUG DISCOVERY



PHARMA PARTNERS

MELLODDY

PUBLIC PARTNERS

# THE MELLODDY OBJECTIVES

**On average, bringing one drug to market costs €1.9 billion and 13 years[1].**

The virtualization of parts of drug discovery by machine learning (ML) is a promising approach to improve efficiencies.

MELLODDY aims to show predictive benefits of modelling across tasks, data types and partners at the largest achievable scale.

[1] DiMasi JA et al., 2016. Innovation in the pharmaceutical industry: new estimates of R&D costs. Journal of Health Economics 47, 20-33.

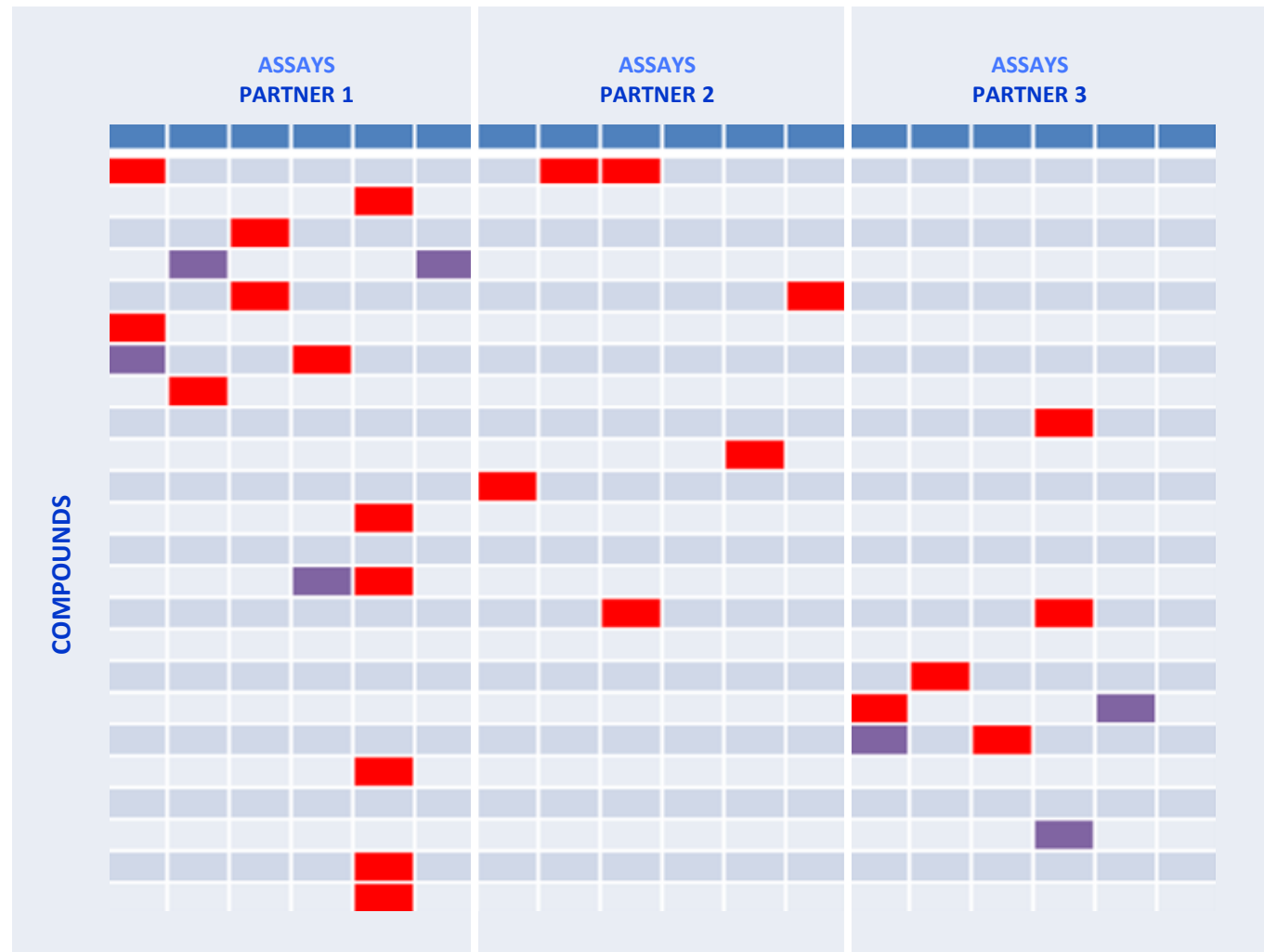**In three yearly runs, the increasingly sophisticated platform will learn from:**

- **> 10 million annotated small molecules**

- **> 1 billion assay biological activity labels**
- **Multiple high-complexity phenotypes at high throughput**

- **Multiple high-complexity phenotypes at high throughput**

**Privacy preservation of data and federated models is paramount.**
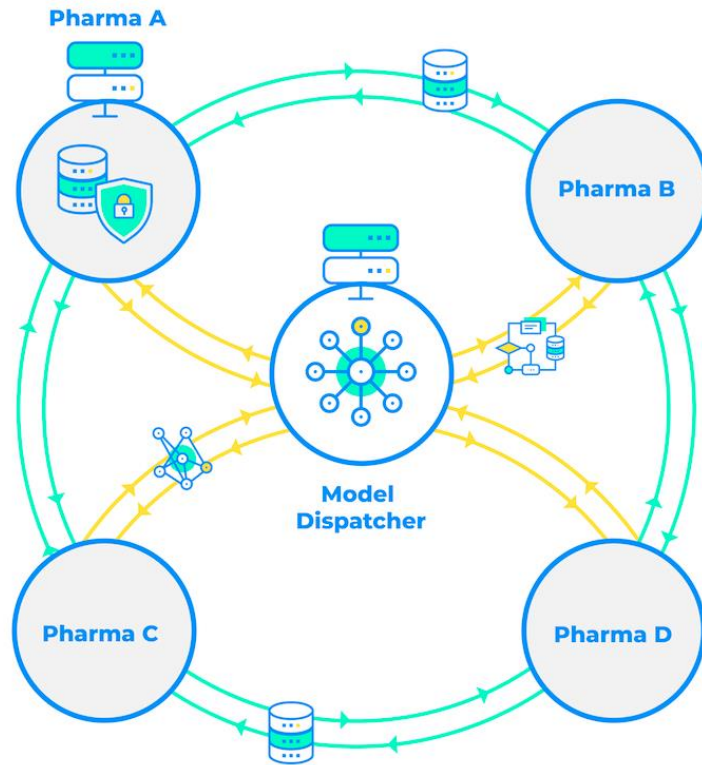
# MULTI-TASK LEARNING ACROSS PHARMA PARTNERS

# COMBINED PRIVACY-PRESERVING FEDERATED MACHINE LEARNING PLATFORM

Sensitive data and assay-specific models remain locked on each pharma's server

Lower level model components are securely exchanged and trained over the network

Complex but transparent pre-agreed access arrangements are strictly enforced



Pharma A

Pharma B

Model Dispatcher

Pharma C

Pharma D

Non sensitive Metadata for ML orchestration

Model updates

Model dispatcher

Data

Algorithm

IT infrastructures

IKTOS

KU LEUVEN

Kubermatic

MŰEGYETEM 1782

nVIDIA.

OWKIN

Substra Foundation

# MELLODDY-TUNER:
# DATA STANDARDIZATION FOR FEDERATED LEARNING

**MELLODDY-TUNER**
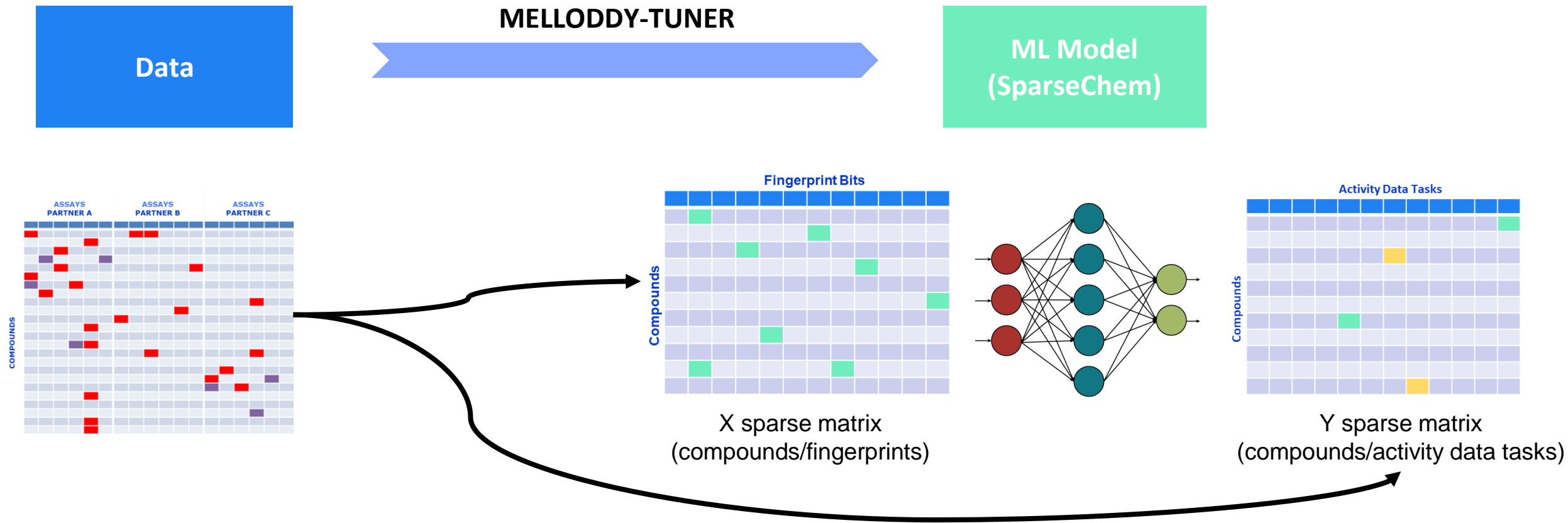
**Data**

**ML Model (SparseChem)**

**Federated Learning**

Standardize data locally at multiple partners

Provide suitable input files for machine learning algorithm

Guarantee uniform processing within consortium while preserving privacy of partner's data

# MELLODDY-TUNER: TECHNICAL OBJECTIVE



X sparse matrix (compounds/fingerprints)

Y sparse matrix (compounds/activity data tasks)

**Standardize structures & activity data to create sparse matrices compatible with SparseChem**

# MELLODDY-TUNER:
# DATA STANDARDIZATION FOR FEDERATED LEARNING

**Standardize smiles**

- Standardize structures:
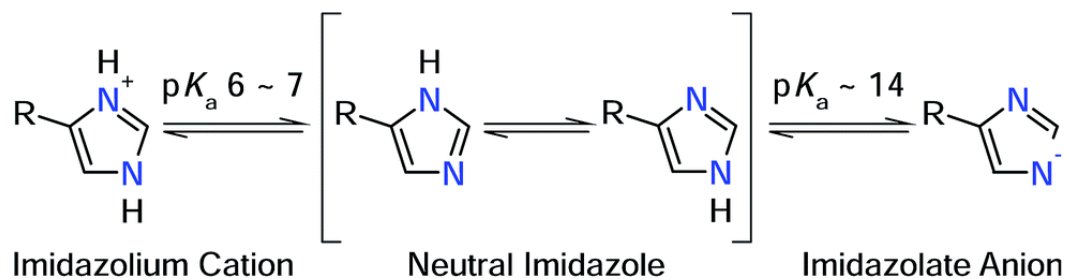
  charge_parent
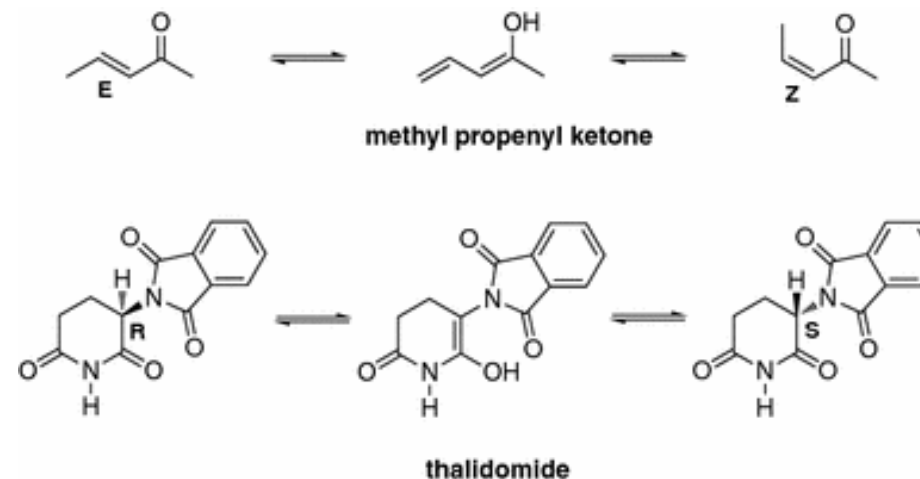
  isotope_parent

  stereo_parent

  **tautomer_parent**

# STRUCTURE STANDARDIZATION

**One structure may not be enough**



Horch et al., *RSC Adv.*, **4**, 54091-54095 (2014)

**Tautomerization can change stereochemistry**



Sitzmann, M. et al., *J Comput Aided Mol Des* **24,** 521–551 (2010)

**Objective: Make most consistent choice for standardization among several partners**
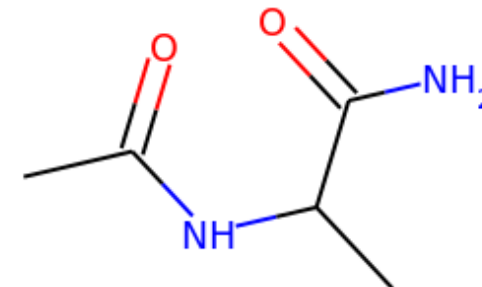
# STRUCTURE STANDARDIZATION
## TAUTOMERIZATION

*tautomer_parent*

**Updated "TautomerTransform"**

*tautomer_parent*

**No tautomerization of esters and amides**

TautomerTransform('1,3 (thio)keto/enol f', '[CX4!H0]-[C;!$([C]([CH1])(=[O,S,Se,Te;X1])-[N,O])]=[O,S,Se,Te;X1]')

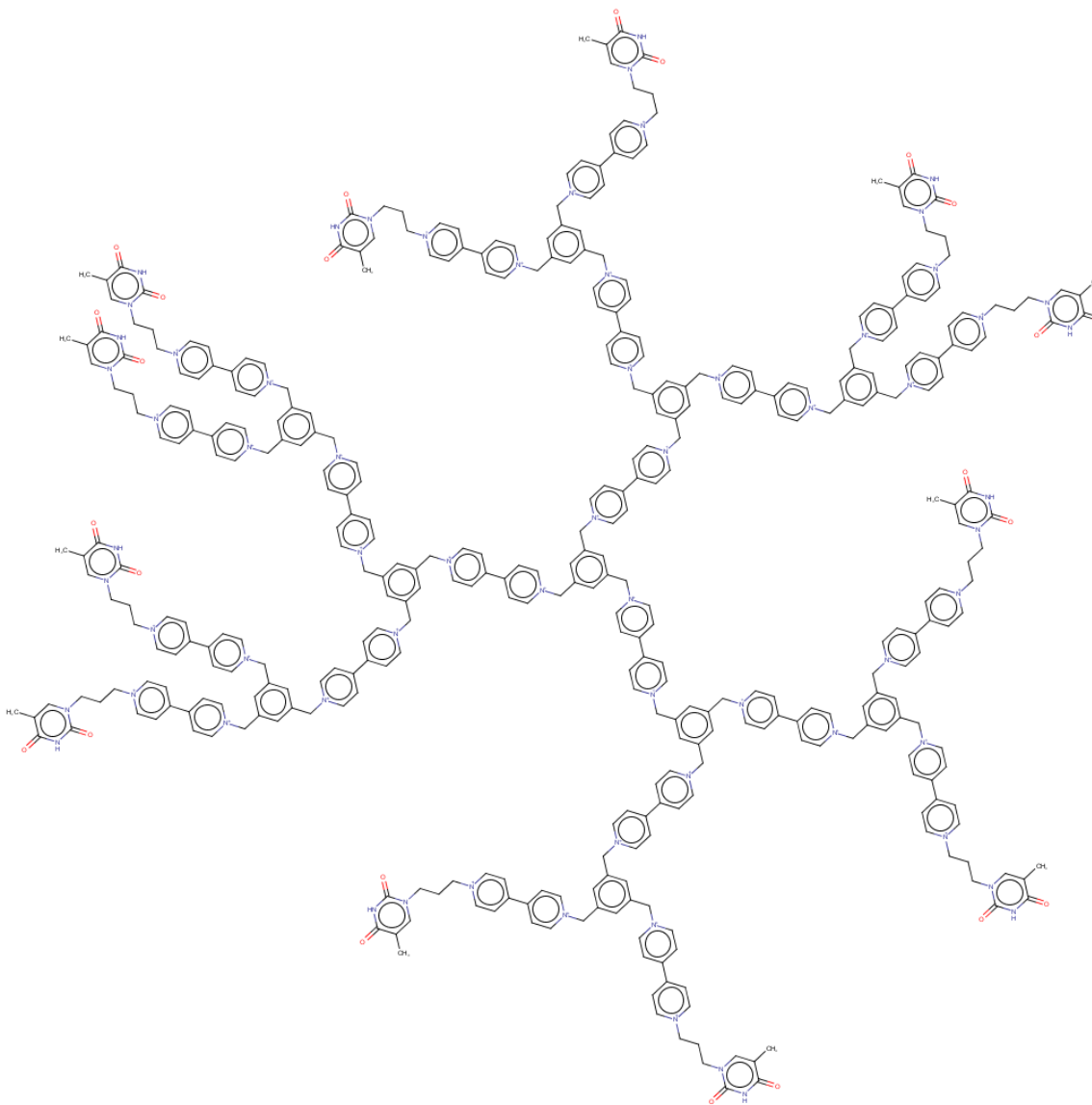Moving Hydrogen from first atom to last atom

# STRUCTURE STANDARDIZATION

**Standardization can take time:**

1. Enumerate all possible tautomers

2. Score all enumerated tautomers

3. Return canonical tautomer

**»** **Limit molecule size
(max. number heavy atoms)
Limit number of enumerated tautomers**



**Molecule from ChEMBL25
Standardization time: ~30 min**

# MELLODDY-TUNER:
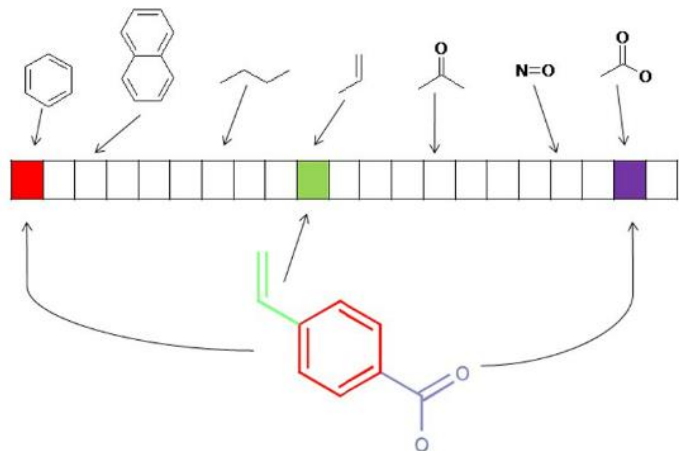## DATA STANDARDIZATION FOR FEDERATED LEARNING

**Standardize smiles**

**Calculate descriptors**

- Calculate fingerprint

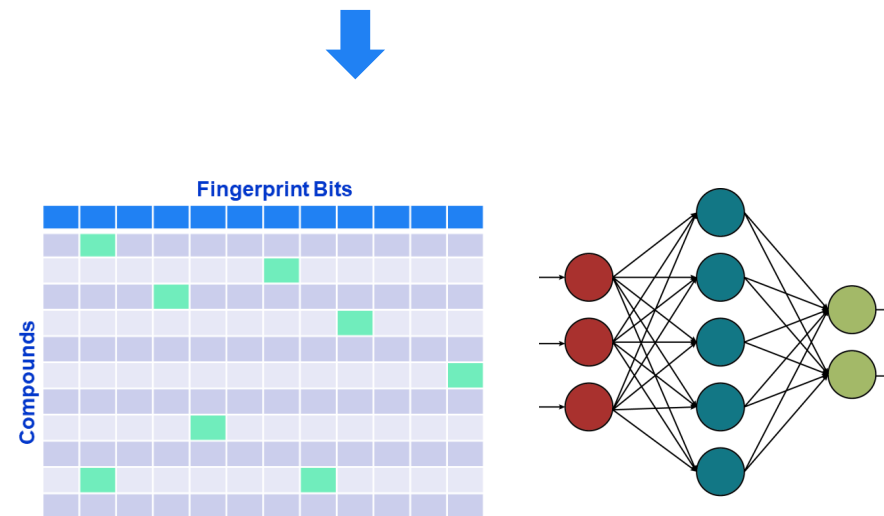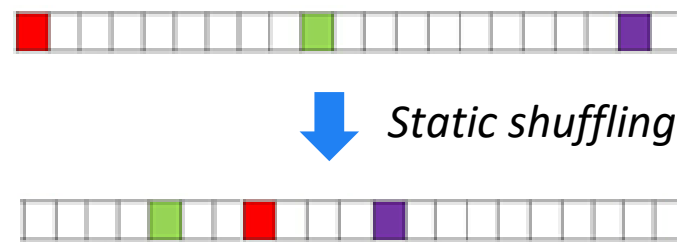- **Cluster fingerprints into folds using locality-sensitive hashing (LSH)**

# FINGERPRINT

**1. Representation of molecules as bit vectors**

**(Morgan Fingerprint with certain length):**



➔ **Fingerprint can be „reverse-engineered"**
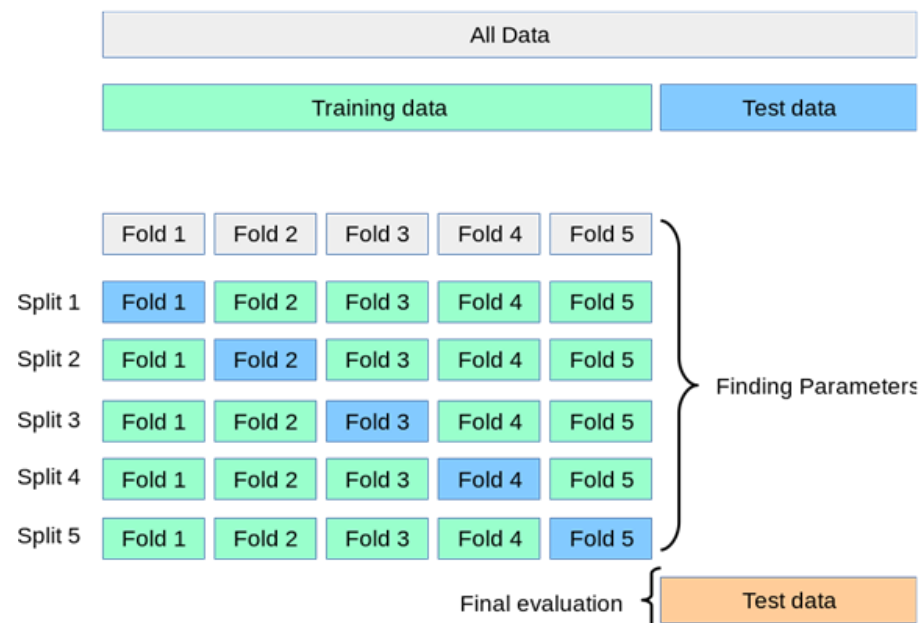(Tuan Le' et al., ChemRxiv (2020))

**2. Static shuffling of bit positions using secret key**



*Static shuffling*

**Molecular fingerprints are part of X matrix for SparseChem**

# FINGERPRINT
# TRAINING, VALIDATION AND TEST SETS

**Evaluation of ML model performance requires data split into training, validation & test sets**
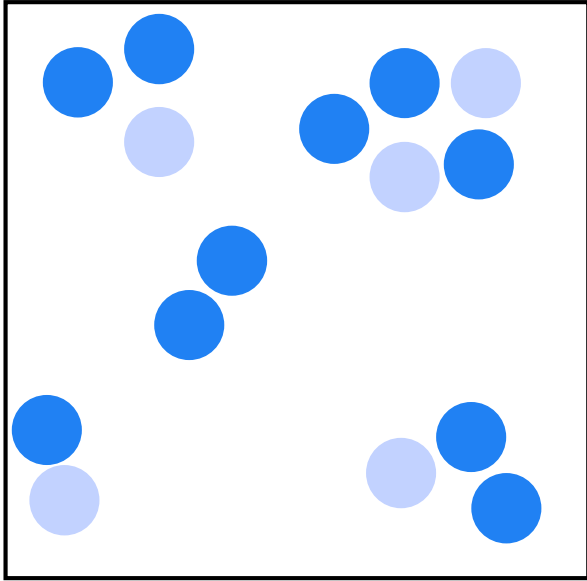


Scikit-learn.org/stable/modules/cross_validation.html

**How can we consistently assign compounds to folds across multiple partners?**
**How to guarantee that identical compounds from different partners land in the same fold?**

# TRAIN/TEST SPLIT:
## RANDOM VS CLUSTER BASED SPLIT



● Training data     ● Test data     ○ Training data w/o assay data     ○ Test data w/o assay data

**Random Split:**
Overly optimistic
performance assessment

**Cluster based split:**
More realistic
performance assessment

**Cluster based split:**
Uneven distribution of
assay data among clusters

» **perfect clustering not required, but privacy-preserving is necessary**

# MELLODDY-TUNER:
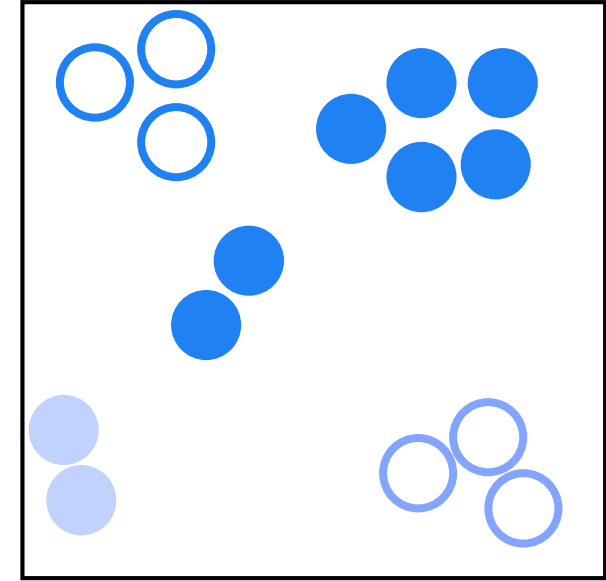# DATA STANDARDIZATION FOR FEDERATED LEARNING

**Standardize smiles** → **Calculate descriptors** → **Format activity data**

- Remove data from failed compounds

- Aggregate replicates

- Filter out tasks not fulfilling minimum number of actives/inactives

# ACTIVITY DATA
## BIT VECTOR REPRESENTATION



**Activity fingerprints are part of Y matrix for SparseChem**

| Assay | A | B | C | D |
|---|---|---|---|---|
| **%Ctrl @ 10µM** | 70 | 8 | 82 | 100 |
| Activity class | 0 | 1 | 0 | 0 |

# ACTIVITY DATA
## DATA AGGREGATION AND FILTERING

| Remove data of „failed" compounds | Aggregate data of replicates | Filte out tasks not fulfilling criteria for SparseChem |

„Failed" compounds like

m1 = Chem.MolFromSmiles('CO(C)C')



*Same fingerprint for different compounds (e.g. stereochemistry not considered)*



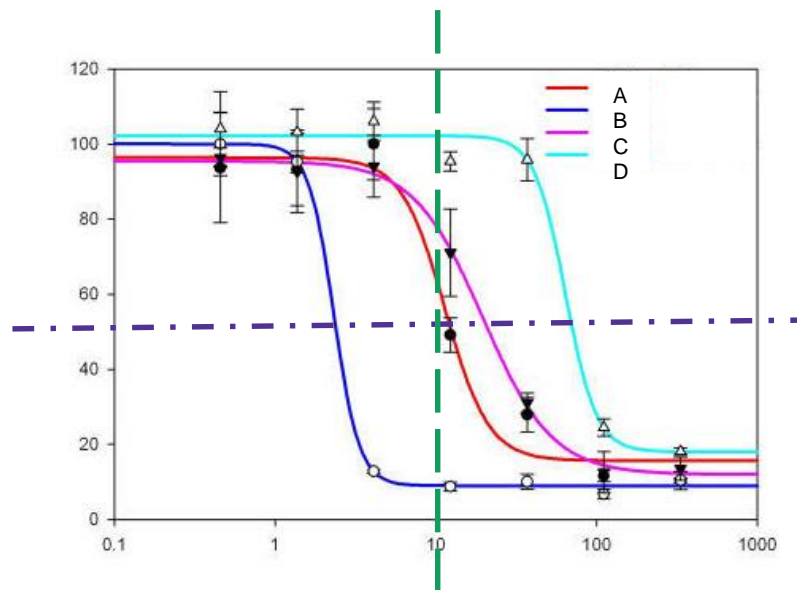| Assay | A | B | C | D |
|---|---|---|---|---|
| %Ctrl @ 10µM | 70 | 8 | 82 | 100 |
| Activity class | 0 | 1 | 0 | 0 |

| Assay | A | B | C | D |
|---|---|---|---|---|
| %Ctrl @ 10µM | 90 | 23 | 15 | 100 |
| Activity class | 0 | 0 | 1 | 0 |

**Majority voting approach:**

„Majority" class wins for tasks with multiple datapoints

„Minority" class wins for tasks with draws



Guarantee sufficient amount of both classes in all folds

**Analysis of ML performance is possible**

# MELLODDY-TUNER:
# DATA STANDARDIZATION FOR FEDERATED LEARNING

**Standardize smiles** → **Calculate descriptors** → **Format activity data** → **Convert to matrices**

- Create X sparse matrix (compound/fingerprint)

- Create Y sparse matrix (compound/activity data tasks)

# SPARSE MATRICES



**Dataframes**

**Sparse matrices**

# SPARSE MATRICES
# MELLODDY-TUNER & SPARSECHEM

**MELLODDY-TUNER**

**SparseChem**

run_name/
results/
results_tmp/
files_4_ML/

Fingerprint Bits

Compounds

Activity Data Tasks

Compounds

MELLODDY-TUNER provides dataframes of standardized data (results),
mapping tables & excluded data (results_tmp) and
SparseChem-compatible matrices (files_4_ML)

# MELLODDY-TUNER
## SUMMARY

**user-defined parameters**

**example_parameters.json**

| Standardize smiles | Calculate descriptors | Format activity data | | Convert to matrices |
|---|---|---|---|---|

**Standardize smiles**

- Standardize structures:

  charge_parent

  isotope_parent

  stereo_parent

  tautomer_parent

**Calculate descriptors**

- Calculate fingerprint
- Cluster fingerprints into folds using locality-sensitive hashing (LSH)

**Format activity data**

- Remove data from failed compounds
- Aggregate replicates
- Filter out tasks not fulfilling minimum number of actives/inactives

**Consistency check**

**Convert to matrices**

- Create compound/ fingerprint matrix
- Create compound/ activity data matrix

**standardize_smiles.py**     **calculate_descriptors.py**     **activity_data_formatting.py**     **csv_2_mtx.py**

**prepare_4_melloddy.py**

# MELLODDY-TUNER
# CONSISTENCY CHECK

## pre-defined parameters

| Reference set | Standardize smiles | Calculate descriptors | Hash result files |
|---|---|---|---|
| • 10 reference structures to be processed | • Standardize structures: charge_parent isotope_parent stereo_parent tautomer_parent | • Calculate fingerprint<br>• Cluster fingerprints into folds using locality-sensitive hashing (LSH) | • Hash generated output files and parameter file<br>• Compare hash with a distributed reference hash |

**standardize_smiles.py**     **calculate_descriptors.py**

**hash_reference_set.py**

# MELLODDY-TUNER: TECHNICAL DETAILS

## MELLODDY-TUNER@Github

- bin
- melloddy_tuner
- tests/structure_preparation_test
- unit_test
- .dockerignore
- .gitignore
- .gitlab-ci.yml
- Dockerfile
- Dockerfile_alternative
- LICENSE
- README.md
- environment_melloddy_tuner.yml
- environment_melloddy_tuner_generic...
- install_environment.sh
- setup.py

**Python 3.6 or higher**

**Conda environment**

**Docker image available**

## Machine Learning Code (local version):

### SparseChem@Github

*Press release Sept. 17th , 2020*

Hugo Ceulemans (Janssen Pharmaceuticals):
*"[…] Over the next year we'll turn our focus on studying the hypothesis that multi-partnered modeling will yield superior predictive models for drug discovery."*

*Press release Sept. 17th , 2020*

*MELLODDY*

# ACKNOWLEDGEMENTS