# A DEEP LEARNING MODEL FOR MOLECULAR FINGERPRINTING

## MATTIA CORDIOLI

DEPARTMENT OF ELECTRICAL, COMPUTER AND BIOMEDICAL ENGINEERING

MASTER'S DEGREE IN BIOENGINEERING

Advisor:

**PROF. RICCARDO BELLAZZI**

Co-advisor:

**PROF. BLAŽ ZUPAN**

biolab

Univerza *v Ljubljani*
Fakulteta *za računalništvo*
*in informatiko*

BIO-MEDICAL INFORMATICS
"*Mario Stefanelli*"

UNIVERSITÀ DI PAVIA

# Outline

- Introduction to Chemoinformatics
  - Molecule representation techniques
  - Fingerprints

- A novel approach: Deep Learning fingerprints
  - Convolutional Neural Networks

- Results on different datasets

- Software development

- Conclusions and future developments

- Standard representations:
  - **not sig...**
  - **2D stru...**
  - **molecu...**

- Molecular ...
  - **connec...**
  - encode ...

- Linear nota...
  - **more c...**
  - useful molecul...
  - **SMILES...** Entry Specification

**Aspirin - $C_9H_8O_4$**
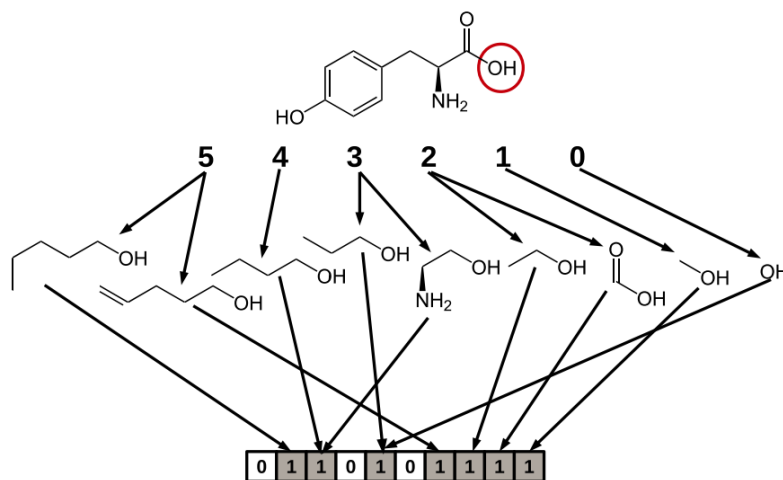


CC(=O)Oc1ccccc1C(=O)O

# Fingerprints

- SMILES are **not enough** in Chemoinformatics applications:
  - **Similarity** and **substructure search**
  - **Virtual screening**
  - **QSAR** and **machine learning** models

- **Fingerprints**: binary vectors of fixed length

**Substructure keys-based**

**Topological / Hashed**

# Deep Learning for Fingerprinting

- Standard fingerprints **limits:**
  - Necessity of a **fragments dictionary** for substructure keys-based FPs
  - Topological FPs are usually **really long** (1024 – 2048 bits)
  - Binary, **not real-valued**
  - **Not trainable** for target-specific tasks

- Novel approach: molecular **embedding** through **Deep Learning**
  - **Deep Neural Networks** learn and abstract powerful representations of input data

- Literature approaches:
  - CNN for **molecular graphs convolution**
  - CNN applied directly to **SMILES strings**

# Aims of the Thesis

- Development of a deep learning model for **molecular fingerprinting**:

  - **Simple CNN** architecture

  - **SMILES strings** as input

- Creation of a new **Chemoinformatics Add-on** for Orange:

  - To provide a tool for easily analyse and work with chemistry data

  - Implementation of the model in a tool for **molecular embedding**

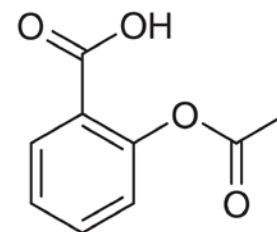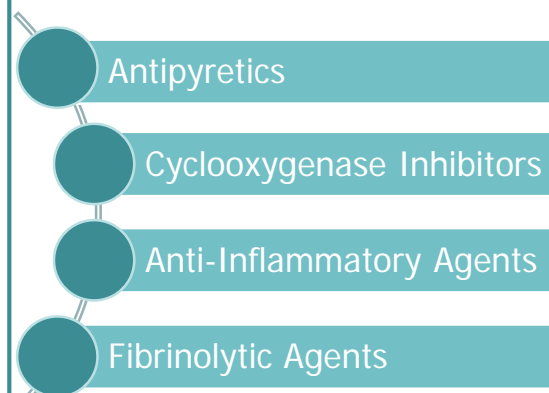  - Implementation of other useful tools, e.g. **molecules visualization**

# Data Retrieval

**PubChem**

- Open chemistry database at NIH
  - **~92 million** compounds with information about structure, chemical/physical properties, identifiers, pharmacology, toxicity, patents, …

- **MeSH** Ontology terms for **pharmacological actions**

- **PubChemAPI** Python library:
  - Programmatic access to PubChem to retrieve data
  - Linking to MeSH Ontology DB to retrieve associated pharmacological actions
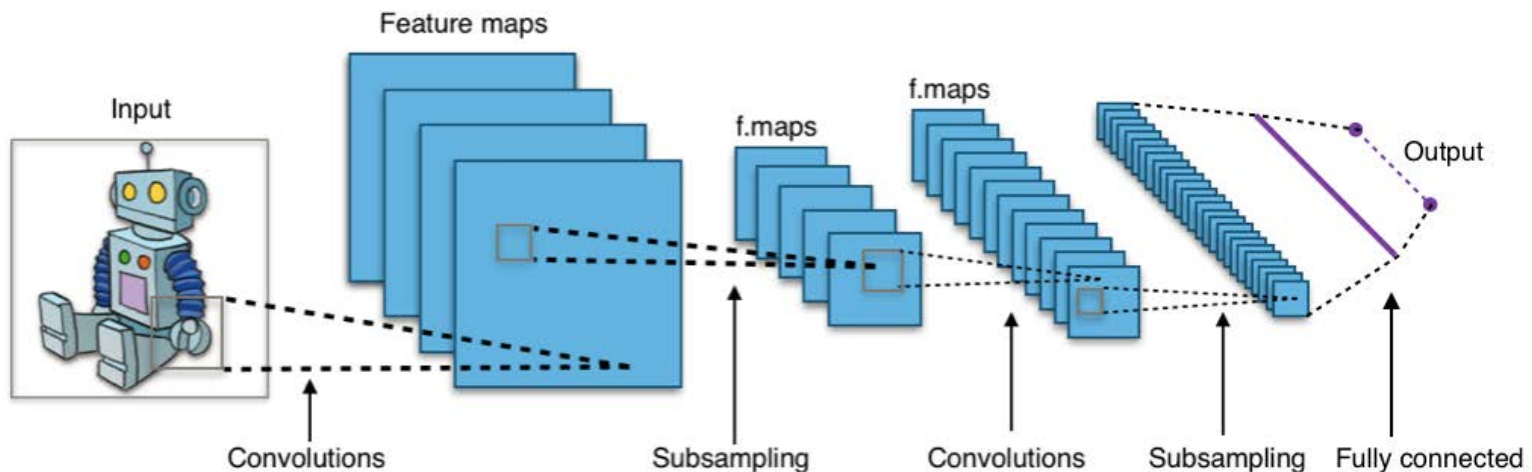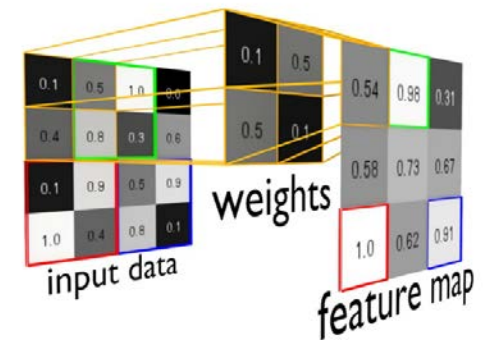
**Aspirin Pharma Actions:**

- Antipyretics
- Cyclooxygenase Inhibitors
- Anti-Inflammatory Agents
- Fibrinolytic Agents

# Data Preprocessing

- **15 474 compounds** retrieved, annotated with **489 terms**
    - **CID, SMILES, name, formula, MeSH terms + tree numbers**

- **Preprocessing:**
    - **Duplicate rows** (same SMILES and terms, different names)
    - Terms appearing **<20 times**
    - Terms with tree number **not starting with 'D27.505'**

- **Final dataset:**
    - **9 174 records**
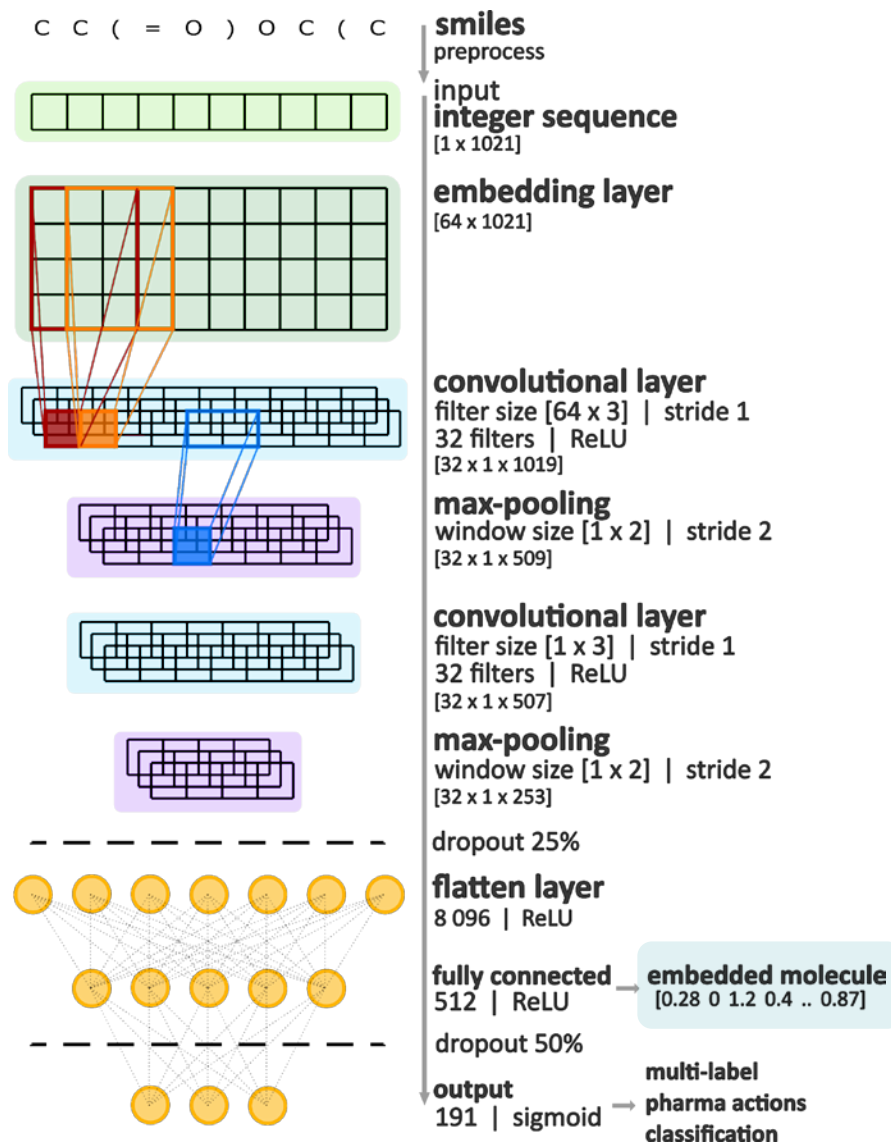    - **191 terms**

# Embedding Model: CNNs

- **Convolutional Neural Networks:**

    - Input data convoluted with **kernels** to obtain **feature maps**
    - **MAX Pooling** layers to reduce dimensionality
    - **Fully connected** layers on top
    - **Dropout** to prevent overfitting

# Embedding Model: Architecture



- **Keras** to design and train the model

- **Supervised** learning
  - Pharmacological actions **multi-label classification**

- Penultimate layer activations
  - **512** bits real-valued **fingerprint**

- Training
  - **on the entire dataset**
  - **GeForce GTX TITAN X GPU**

- Model saved to be used as **embedder**

# Validation: CNN Performance

- Assessing CNN performance:
  - **70/30 train/test split**

- Metric:
  - Area Under ROC Curve (**AUC**) **for each term** separately

| Minimum AUC | Maximum AUC | Mean AUC |
|:---:|:---:|:---:|
| 0.62 | 0.99 | 0.87 |

# Comparison: MeSH Terms Prediction

- Comparison with **ECFP**:
  - Circular/topological fingerprint
  - **Standard in QSAR**
  - **512-bits** version

- **Pharmacological Actions** prediction:
  - **Logistic Regression**
  - **One-Vs-All** approach
  - 10-Fold Cross Validation
  - AUC

| Fingerprint | Mean AUC |
|:-----------:|:--------:|
| CNNFP | **0.99** |
| ECFP | 0.92 |

- **Non-Pharma** terms prediction:
  - 1 091 compounds discarded in the preprocessing phase

| Fingerprint | Mean AUC |
|:-----------:|:--------:|
| CNNFP | 0.83 |
| ECFP | **0.92** |

# Comparison: Other QSAR Datasets

- Comparison on datasets obtained from MoleculeNet:
  - **Logistic Regression** and **Random Forest**
  - CNNFP, ECFP, CNNFP+ECFP
  - 10-Fold CV AUC

- **ClinTox**:
  - **1 491** compounds
  - Clinical Trial toxicity
  - FDA approval status

| Task | Classifier | ECFP | CNNFP | CNNFP+ECFP |
|------|------------|------|-------|------------|
| CT Toxicity | LR | 0.72 | **0.93** | **0.95** |
| | RF | 0.74 | **0.94** | **0.96** |
| FDA Approval | LR | 0.74 | **0.92** | **0.95** |
| | RF | 0.74 | **0.94** | **0.97** |

# Comparison: Other QSAR Datasets - 2

- **BACE**:
  - **1 522** β-secretase-1 inhibitors
  - Binding results

- **BBBP:**
  - **2 000** compounds
  - Blood-brain barrier permeability

| Dataset | Classifier | ECFP | CNNFP | CNNFP+ECFP |
|---------|-----------|------|-------|------------|
| BACE | LR | **0.85** | 0.72 | 0.81 |
| | RF | **0.87** | 0.79 | 0.86 |
| BBBP | LR | **0.83** | 0.79 | **0.85** |
| | RF | 0.88 | **0.89** | **0.91** |

# Comparison: t-SNE Visualization

- **Data**:
  - Compounds related to the 5 most frequent MeSH Terms
  - ClinTox
  - BACE
  - BBBP

- **t-SNE**:
  - Non-linear dimensionality reduction
  - **PCA** (**100 components**) applied before

| | Explained Variance | |
|---|---|---|
| | **ECFP** | **CNNFP** |
| MeSH Terms | 52% | **91%** |
| ClinTox | 65% | **90%** |
| BACE | 84% | **94%** |
| BBBP | 67% | **87%** |

# Comparison: t-SNE Visualization - 2

# Comparison: t-SNE Visualization - 3



ClinTox   BACE   BBBP

CNNFP · CNNFP · CNNFP

ECFP · ECFP · ECFP

# Orange

- **Machine Learning and Data Mining suite** developed by **Bioinformatics Lab** at **University of Ljubljana**

# Chemoinformatics Add-on

- **Molecule Embedding Widget**
  - Embedding on Orange server, using Keras and TITAN X GPU

# Chemoinformatics Add-on - 2

- **Molecule Viewer Widget**

# Conclusions and Future Developments

- **Conclusions:**
  - Novel deep learning model for molecular fingerprinting
  - Short real-valued fingerprint (512 bits) with high representative power
  - Simple architecture, using simple input representation
  - Good capability of generalitazion
  - Trainable for target specific applications

- **Future Developments:**
  - Optimize a tool for pharmacological actions prediction
    - Drug repurposing
  - Chemical interpretation of the learned features
  - Extension of Chemoinformatics Add-on functionalities

*Grazie per l'attenzione*