

# Artificial Intelligence Methods for Evaluating Mitochondrial Dysfunction

Exploring Various Chemical Notations Suitable for Neural  
Language Processing Models

ICANN

The 33rd International Conference on Artificial Neural Networks.

A conference of the European Neural Network Society

# Cardiovascular diseases



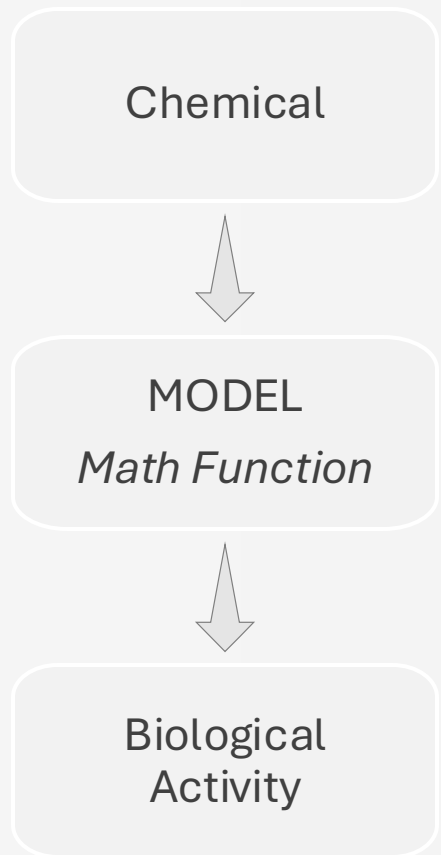
A **Cardiovascular disease** is a multifactorial disease that involves a combination of **genetic, environmental,** and **lifestyle** factors.

Each cardiovascular disease has a **different set of risk factors** and mechanisms of development.

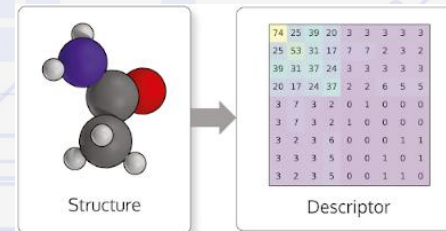
As a result, **accurately predicting** the onset or progression of cardiovascular diseases in humans can be **challenging**.

# In-silico methods: Quantitative Structure-Activity Relationship (QSAR)

With QSAR it is possible to predict the biological activity of a compound based on its chemical structure and other related properties



Chemicals are **encoded** in a suitable way for modeling.



Machine learning (ML) models to **predict interaction** with **biological targets**



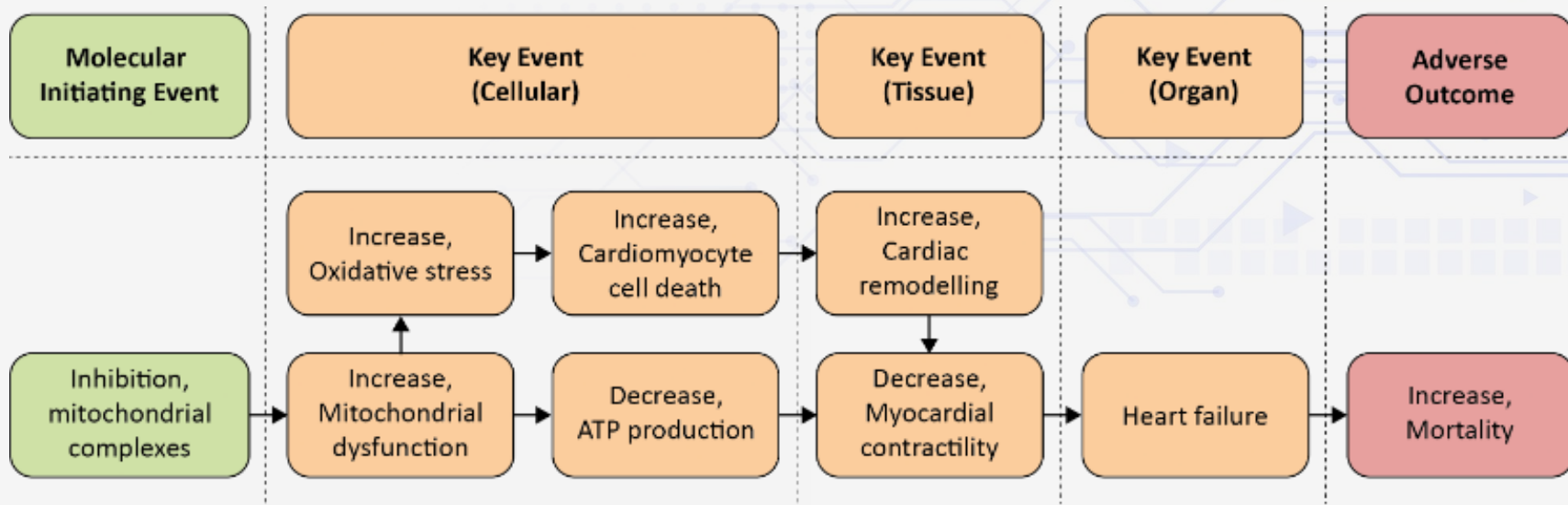
Biological targets identified based on the **AOPs** for cardiotoxicity



# Adverse Outcome Pathway (AOP): Mitochondrial Dysfunction

AOP framework highlights possible endpoints related to cardiotoxicity effects for modeling perspective:

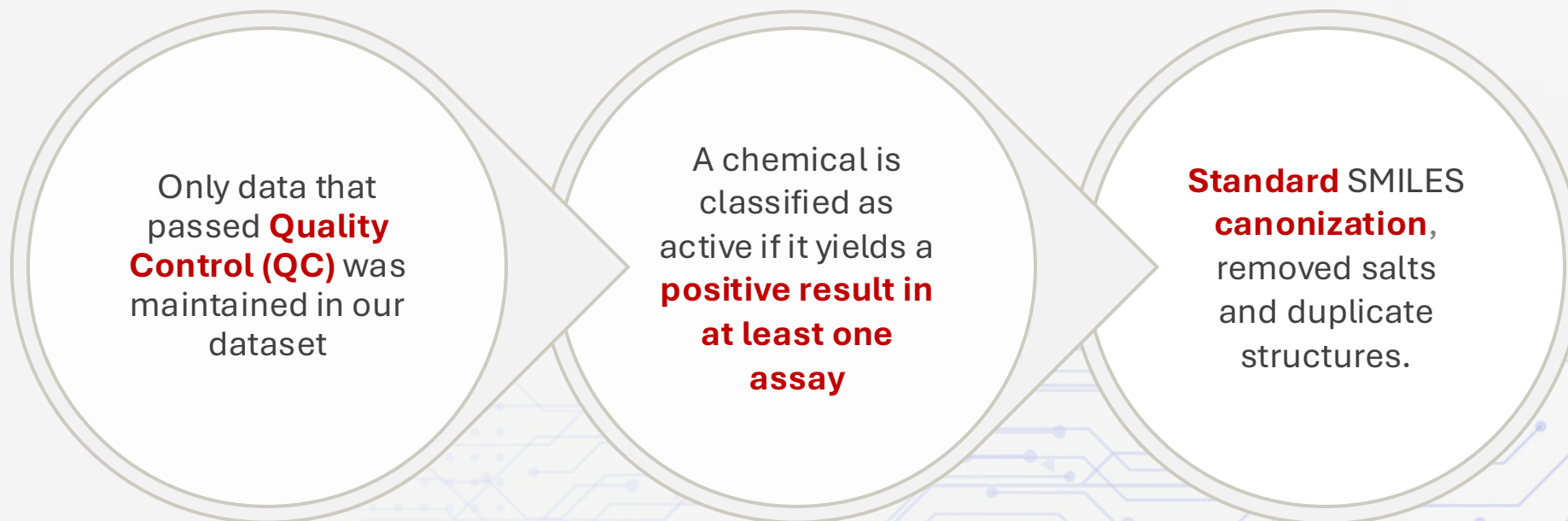
1. Structural effects
2. Contractile effects
3. Electrophysiological effects



## ASSAYS Mitochondrial Dysfunction

APR\_HepG2\_MitoMass\_24h\_dn  
 APR\_HepG2\_MitoMass\_24h\_up  
 APR\_HepG2\_MitoMass\_72h\_dn  
 APR\_HepG2\_MitoMass\_72h\_up  
 APR\_HepG2\_MitoMembPot\_24h\_dn  
 APR\_HepG2\_MitoMembPot\_24h\_up  
 APR\_HepG2\_MitoMembPot\_72h\_dn  
 APR\_HepG2\_MitoMembPot\_72h\_up  
 ATG\_XTT\_Cytotoxicity\_up  
 TOX21\_MMP\_ratio\_down  
 TOX21\_MMP\_ratio\_up  
 TOX21\_MMP\_rhodamine

# Results of Data Collection and Data Curation



	Total Chemicals	Training set	Validation set	Holdout set
All	5004	4052	451	501
Active	1147	929	103	115
Inactive	3857	3123	348	386

# Baseline Machine Learning: state of the art for Mitochondrial Dysfunction

Machine Learning methods	External Validation						Descriptors	Data Numerosity External Validation	References
	Balanced Accuracy	Precision	Sensitivity	Specificity	MCC	F1-Score			
Gradient Boosting	0.708	0.573	0.467	0.948	0.454	0.515	Atom Pair FP	893	DOI: 10.1002/minf.202000005
Random Forrest	0.743	0.279	0.793	0.692	0.338	0.413	RDKit mol.desc.	893	DOI: 10.1002/minf.202000005
Neural Network	-	0.45	0.68	0.88	0.48	0.54	CDDD	761	DOI: <a href="https://doi.org/10.1021/acs.chemrestox.3c00086">10.1021/acs.chemrestox.3c00086</a>
Extreme GB	0.742	0.650	0.600	0.883	0.485	0.602	CDDD	1001	DOI: <a href="https://doi.org/10.3390/toxics12010087">https://doi.org/10.3390/toxics12010087</a>

Results taken from the original papers; not recalculated.



# ENCODING CHEMICAL INFORMATION

There are many **different ways to encode chemical information**, meaning that chemical properties can be represented in **various forms suitable for machine learning**. But which method is the best?

## Machine Learning Models

SVC

Logistic Regression

Decision Tree

Random Forest

KNN

GaussianNB

## Artificial Intelligence: Deep Learning

GCNN

Deep Neural Network

Multitask Models

Multimodal Models

**Molecular Descriptor**

MW	AMW	Mv	Me	Mp	Mi	GD	nTA	nBM
206.3	6.252	0.581	0.988	0.629	1.126	0.143	5	7

**Morgan Fingerprint**

**Graph**

**Word Embedding**

**CDDD**

Chiarified text: in C C C C C C [ N + ] [ - 0 ] [ 0 - ]

Embedding:

```
[[[-0.02392174  0.02877975  0.04714444 ... -0.03017147 -0.00659642
  0.02174393]
 [-0.02392174  0.02877975  0.04714444 ... -0.03017147 -0.00659642
  0.02174393]
 [-0.02392174  0.02877975  0.04714444 ... -0.03017147 -0.00659642
  0.02174393]
 ...
 [-0.02493726 -0.02160301 -0.02045706 ... 0.04390224 -0.0471576
  0.0295709 ]
 [-0.02493726 -0.02160301 -0.02045706 ... 0.04390224 -0.0471576
  0.0295709 ]
 [-0.02493726 -0.02160301 -0.02045706 ... 0.04390224 -0.0471576
  0.0295709 ]]]]
```

Character embedding shape: (1, 159, 64)

[Open Access](#) [Editor's Choice](#) [Article](#)

**Artificial Intelligence and Machine Learning Methods to Evaluate Cardiotoxicity following the Adverse Outcome Pathway Frameworks**

by Edoardo Luca Viganò<sup>1\*</sup>, Davide Ballabio<sup>2</sup> and Alessandra Roncaglioni<sup>1</sup>

<sup>1</sup> Laboratory of Environmental Toxicology and Chemistry, Department of Environmental Health Sciences, Istituto di Ricerche Farmacologiche Mario Negri IRCCS, 20156 Milan, Italy

<sup>2</sup> Milano Chemometrics and QSAR Research Group, Department of Earth and Environmental Sciences, University of Milano-Bicocca, 20126 Milan, Italy

\* Author to whom correspondence should be addressed.

Toxics 2024, 12(1), 87; <https://doi.org/10.3390/toxics12010087>

# NLP Methods For Molecular Property Prediction

NLP methods offer a **data-driven** and computationally efficient approach to toxicological screening.

NLP models learn the chemical **grammar** and **semantics** of chemical notation.



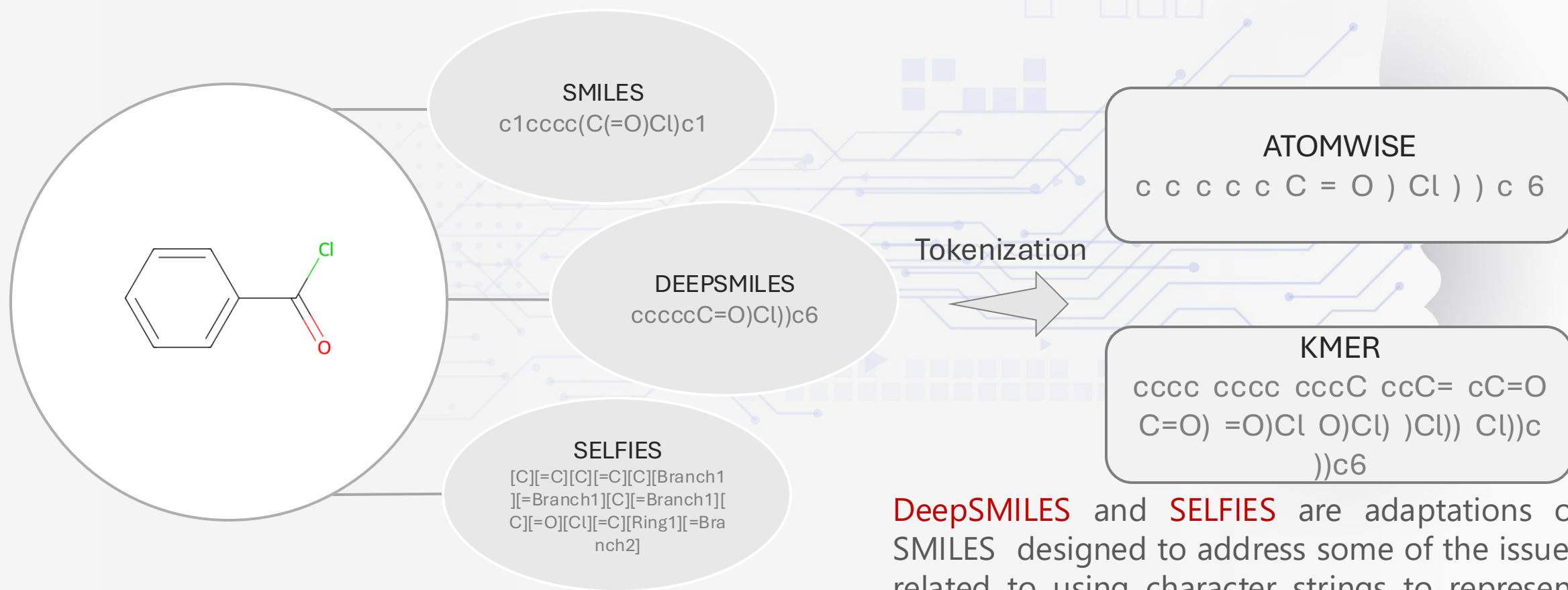
**New model** developed in-house with a custom architecture

Large-Scale Self-Supervised Pretraining for Molecular Property Prediction models:  
**ChemBERTa, RoBERTa.**



# Chemical Notations

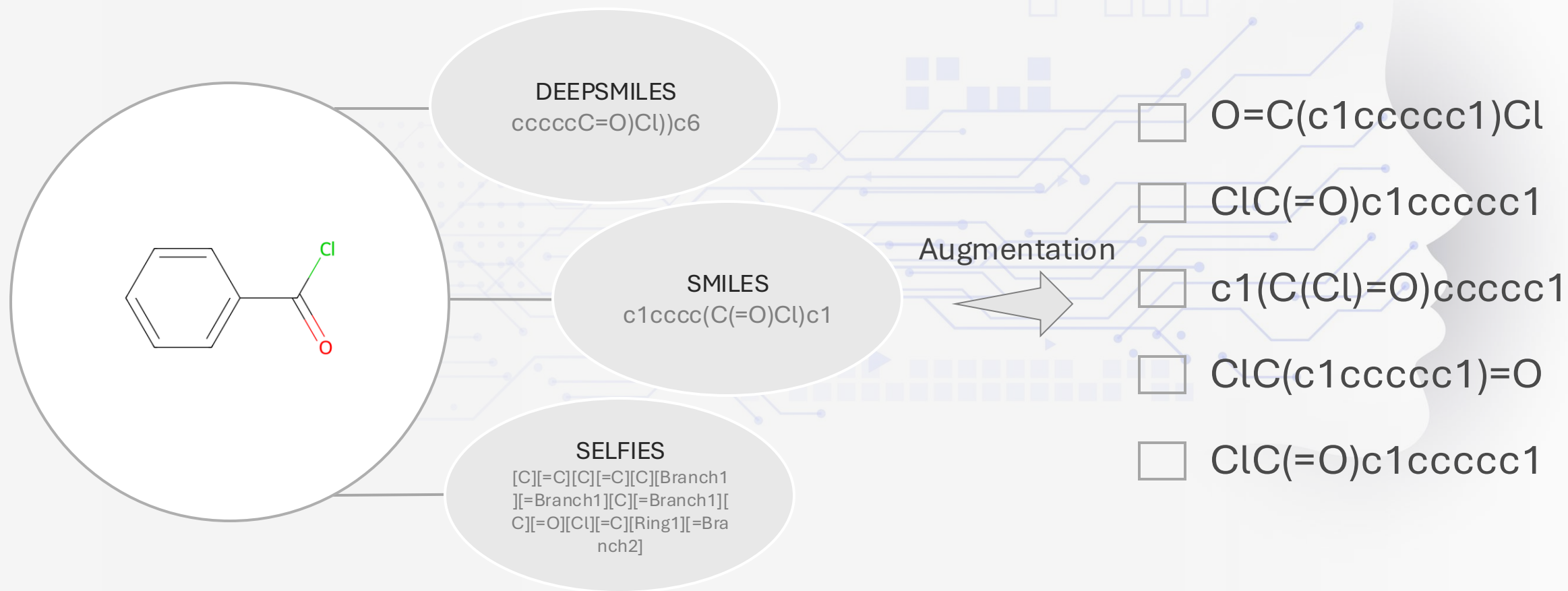
The **Simplified Molecular Input Line Entry System** (SMILES) is a widely used chemical notation system that represents the structure of molecules as a linear string of characters.



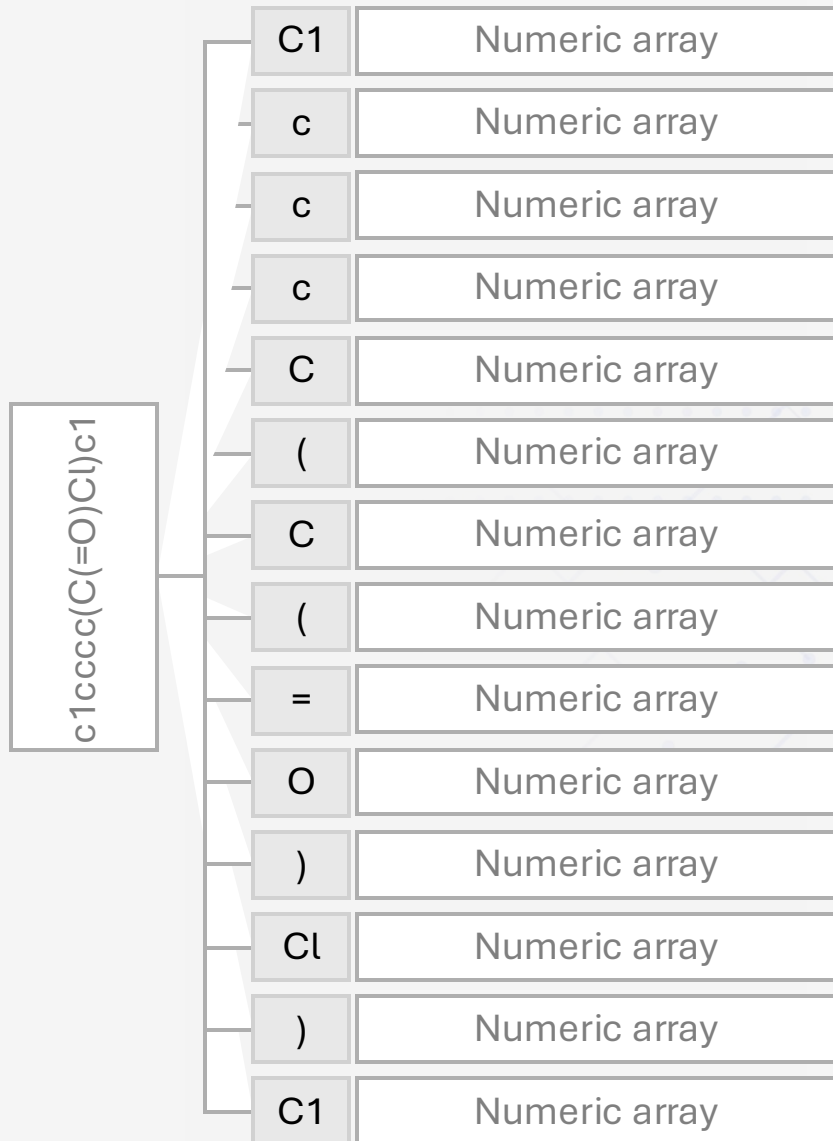
**DeepSMILES** and **SELFIES** are adaptations of SMILES designed to address some of the issues related to using character strings to represent chemicals in machine learning.

# SMILES Augmentation

Multiple SMILES represent the same molecule. This feature about chemical notation is explored as a technique for data **augmentation** of a molecular QSAR dataset.



# Character Embedding and Model Architecture

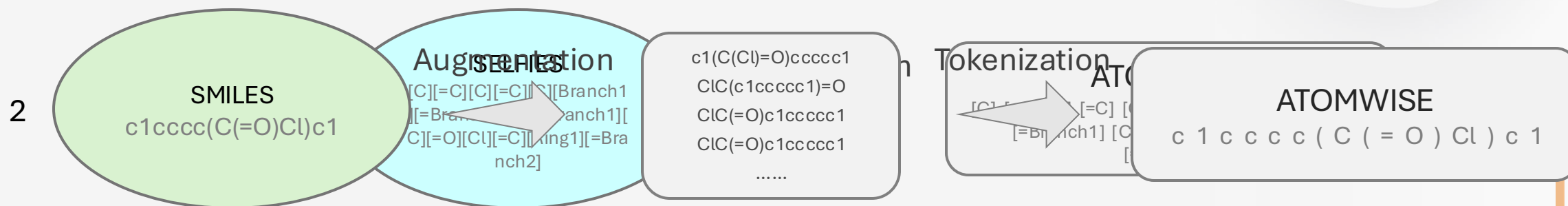


Layer (type)	Output Shape	Param #
Input Layer	[(None, 1)]	0
Text Vectorization	(None, 54)	0
Embedding	(None, 54, 128)	2091776
Conv1D	(None, 54, 128)	328192
Conv1D	(None, 54, 64)	655616
Bidirectional	(None, 54, 64)	394240
Bidirectional	(None, 54, 32)	123648
GlobalMaxPooling1D	(None, 128)	0
Dense	(None, 128)	16512
Batch Normalization	(None, 128)	512
Dropout	(None, 128)	0
Dense	(None, 128)	1056
Batch Normalization	(None, 128)	128
Dropout	(None, 128)	0
Dense	(None, 64)	528
Dense	(None, 1)	17

Tot. param: 2303081 (8.79 MB)  
 Trainable params: 2302953 (8.79 MB)  
 Non-trainable params: 128 (512.00 Byte)

# Results

	Notation	Holdout Set					Validation	Tokenizer	AUG	
		BA	Prec	Sens	Spec	MCC	F1-Score			F1-Score
2	SMILES atomwise	0.766	0.660	0.625	0.906	0.542	0.642	0.744	atomwise	no
	AUG SMILES atomwise	0.810	0.552	0.812	0.808	0.550	0.657	0.803	atomwise	yes
	SMILES kmer	0.747	0.536	0.661	0.834	0.461	0.592	0.857	kmer	no
	AUG SMILES kmer	0.761	0.717	0.589	0.932	0.561	0.647	0.613	kmer	yes
1	DeepSMILES atomwise	0.785	0.557	0.741	0.829	0.519	0.636	0.731	atomwise	no
	AUG DeepSMILES atomwise	0.764	0.503	0.741	0.787	0.469	0.599	0.745	atomwise	yes
	DeepSMILES kmer	0.764	0.503	0.741	0.787	0.469	0.599	0.822	kmer	no
	AUG DeepSMILES kmer	0.742	0.503	0.679	0.805	0.439	0.578	0.757	kmer	yes
1	<b>SELFIES atomwise</b>	<b>0.775</b>	<b>0.707</b>	<b>0.625</b>	<b>0.924</b>	<b>0.575</b>	<b>0.664</b>	<b>0.688</b>	<b>atomwise</b>	<b>no</b>
	AUG SELFIES atomwise	0.749	0.605	0.616	0.883	0.495	0.611	0.887	atomwise	yes
	SELFIES kmer	0.715	0.438	0.688	0.742	0.375	0.535	0.812	kmer	no
	AUG SELFIES kmer	0.738	0.551	0.625	0.851	0.456	0.586	0.613	kmer	yes



# Conclusions

- Previous works suggest to us that **NLP** are very promising approaches to predict **mitochondrial dysfunction**.
- We tested different approaches to encode chemical structures as strings of characters (**SMILES, SELFIES, DEEPSMILES**) and explored various tokenization methods (Atomwise, KMER).
- These methods **outperform** the machine learning models' capability to predict mitochondrial toxicity.

<i>External Set</i>	<i>Best ML: XGB</i>	<i>Neural Network 1</i>	<i>Neural Network 2</i>	<i>Advanced AI Mechanistic (CDDD)</i>	<i>Advanced AI Multi-task NN</i>	<i>Advanced AI Multimodal</i>	<b>NLP</b>	<i>Ensemble (undersampled gradient boosting models)</i>
MCC	0.49	0.48	0.53	0.54	0.56	0.56	<b>0.58</b>	0.61
Performance Comparison with NLP	<b>+18.4%</b>	<b>+20.8%</b>	<b>+9.4%</b>	<b>+7.4%</b>	<b>+3.5%</b>	<b>+3.5%</b>	-	<b>-4.9%</b>
Reference	DOI: <a href="https://doi.org/10.3390/toxics12010087">https://doi.org/10.3390/toxics12010087</a>	DOI: <a href="https://doi.org/10.1021/acs.cchemrestox.3c00086">10.1021/acs.cchemrestox.3c00086</a>	DOI: <a href="https://doi.org/10.1002/minf.202000005">10.1002/minf.202000005</a>	DOI: <a href="https://doi.org/10.1021/acs.cchemrestox.3c00086">10.1021/acs.cchemrestox.3c00086</a>	DOI: <a href="https://doi.org/10.1021/acs.cchemrestox.3c00086">10.1021/acs.cchemrestox.3c00086</a>	DOI: <a href="https://doi.org/10.1021/acs.cchemrestox.3c00086">10.1021/acs.cchemrestox.3c00086</a>	10.3390/toxics12010087	<a href="https://doi.org/10.1016/j.comtox.2021.100189">https://doi.org/10.1016/j.comtox.2021.100189</a>

Results taken from the original papers; not recalculated.



# Next steps

## 1. Fine tuning Large-Scale Self-Supervised Pretraining for Molecular Property Prediction

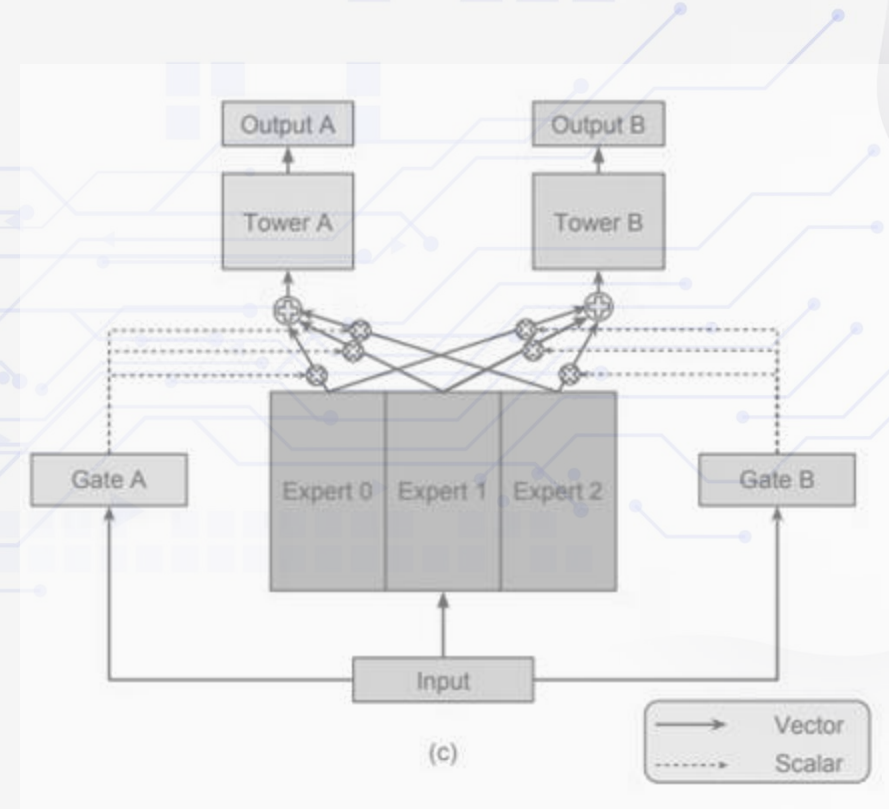


- **ChemBERTa** is Large-Scale Self-Supervised Pretraining with 12 attention heads and 6 layers, resulting in 72 distinct attention mechanisms
- **Fine-tuning** this type of model, which is trained on millions of SMILES, could help to better **generalize chemical assessments**

# Next steps

2. Explore Multitask methods to expanded the coverage for predicting possible toxic interaction between chemicals and biological target that bring to cardiotoxicity

	<i>N° compounds</i>	<i>Active %</i>	
<b>Mode Of Actions</b>	Cardiomyocyte Myocardial Injury	5320	29
	Change Action Potential	415	27
	Change in Inotropy	922	23
	Change In Vasoactivity	4969	20
	Endothelial Injury Coagulation	5634	43
	Valvular Injury Proliferation	268	34
<b>Apical Cardiotoxicity</b>	Clinical Data	848	73
<b>Key Events</b>	hERG inhibition	8462	51
	KE1 Increased Oxidative Stress	636	30
	KE2 Increased Mitochondrial Dysfunction	5004	23



# ACKNOWLEDGEMENTS AND DISCLAIMER

- This work has received funding from the European Union's Horizon 2020 research and innovation program under grant agreement No 101037090, project ALTERNATIVE.
- This presentation reflects the author's view, only, and the Commission is not responsible for any use that may be made of the information provided.



This project has received funding from the EU's Horizon 2020 research and innovation programme under grant agreement No 101037090





**Thanks for your attention!**