

# Supplementary data - “Using Molecular Embeddings in QSAR Modeling: Does it Make a Difference?”

María Virginia Sabando,<sup>\*,†,‡</sup> Ignacio Ponzoni,<sup>†,‡</sup> Evangelos E. Milios,<sup>¶</sup> and Axel J. Soto<sup>†,‡</sup>

<sup>†</sup>*Institute for Computer Science and Engineering (UNS-CONICET), Bahía Blanca, Argentina*

<sup>‡</sup>*Department of Computer Science and Engineering, Universidad Nacional del Sur, Bahía Blanca, Argentina*

<sup>¶</sup>*Faculty of Computer Science, Dalhousie University, Halifax, NS, Canada*

E-mail: virginia.sabando@cs.uns.edu.ar

## Contents of the Repository

- *DeepSMILES formulas*: a folder containing one .csv file per labeled dataset (SR-ARE, SR-MMP, SR-ATAD5, HIV, PCBA-686978, ESOL, FreeSolv, Lipophilicity), each corresponding to the *DeepSMILES*<sup>1</sup> canonicalization of the compounds and the activity values.
- *RDKit SMILES formulas*: a folder containing one .csv file per labeled dataset (SR-ARE, SR-MMP, SR-ATAD5, HIV, PCBA-686978, ESOL, FreeSolv, Lipophilicity), each corresponding to the RDKit<sup>2</sup> canonicalization of the compounds and the activity values.

- *Trained embedding models*: a folder containing five .zip files. Each .zip file contains trained embedding models for each of the five embedding methods analyzed in our paper: *SMILESTVec*,<sup>3</sup> *Mol2Vec*,<sup>4</sup> *Seq2Seq*,<sup>5</sup> *SA-BiLSTM*<sup>6</sup> and *PaccMann*.<sup>7</sup> These trained models can be loaded and used to extract embeddings for new datasets besides the ones covered in our paper, except for the supervised models *SA-BiLSTM* and *PaccMann* which, because of having undergone a supervised training, are only meant to be used with the datasets used in this study.
- *Source code*: a folder containing all the source code needed in order reproduce our approach. We provide scripts for training the supervised methods *SA-BiLSTM* and *PaccMann*; scripts for training all classification and regression methods tested in our paper (Naïve Bayes, Support Vector Machine, Random Forest, Ridge regression, Gradient Boosting Regression and Feed-Forward Neural Networks); simple scripts showing how to extract embeddings from the pre-trained models provided in this repository and scripts for computing traditional molecular representations.

Please note that we do not provide the source code to train the unsupervised models *SMILESTVec*<sup>i</sup>, *Mol2Vec*<sup>ii</sup> and *Seq2Seq*<sup>iii</sup>, since we neither own nor have contributed to the development of such source code. Instead, please refer to the corresponding GitHub repositories provided by the original authors of such techniques.

- *Traditional\_molecular\_representations.zip*: a .zip file containing all traditional molecular representations computed for each of the eight labeled datasets under study (SR-ARE, SR-MMP, SR-ATAD5, HIV, PCBA-686978, ESOL, FreeSolv, Lipophilicity).

---

<sup>i</sup><https://github.com/hkmztrk/SMILESTVecProteinRepresentation/tree/master/source/word2vec>.

<sup>ii</sup><https://github.com/samoturk/mol2vec>.

<sup>iii</sup><https://github.com/XericZephyr/seq2seq-fingerprint>.

# Datasets

- *SR-ARE*, a bioassay for small molecule agonists of the antioxidant response element (ARE) signaling pathway<sup>iv</sup>. This assay of data is contained in the *Tox21 challenge* dataset, consisting of qualitative toxicity measurements on twelve biological targets<sup>v</sup>.
- *SR-MMP*, a stress response assay for small molecule disruptors of the mitochondrial membrane potential (MMP)<sup>vi</sup>. This dataset is also included in the *Tox21 challenge* dataset.
- *SR-ATAD5*, a set of small molecules that induce genotoxicity in human embryonic kidney cells expressing luciferase-tagged ATAD5<sup>vii</sup>. This dataset also belongs to the *Tox21 challenge* dataset.
- *HIV*, a dataset introduced by the Drug Therapeutics Program (DTP) AIDS Antiviral Screen containing information about molecular ability to inhibit HIV replication. Screening results were categorized as confirmed inactive (CI), confirmed active (CA) and confirmed moderately active (CM). We merged the compounds categorized under the latter two labels, which yielded two classes: inactive (CI) and active (CA and CM)<sup>viii</sup>.
- *PCBA-686978*, a PubChem bioassay containing information about molecular ability to inhibit the human tyrosyl-DNA phosphodiesterase 1 (TDP1)<sup>ix</sup>.
- *ESOL*, a dataset consisting of water solubility data for 1128 compounds, used to train models that estimate solubility directly from chemical structures.<sup>8</sup>

---

<sup>iv</sup><https://pubchem.ncbi.nlm.nih.gov/bioassay/743219>

<sup>v</sup><https://tripod.nih.gov/tox21/challenge/about.jsp>

<sup>vi</sup><https://pubchem.ncbi.nlm.nih.gov/bioassay/720637>

<sup>vii</sup><https://pubchem.ncbi.nlm.nih.gov/bioassay/720516>

<sup>viii</sup><https://wiki.nci.nih.gov/display/NCIDTPdata/AIDS+Antiviral+Screen+Data>

<sup>ix</sup><https://pubchem.ncbi.nlm.nih.gov/bioassay/686978>

- *FreeSolv*—the *Free Solvation Database*—, a dataset comprising experimental and calculated hydration free energy values for small molecules in water.<sup>9</sup>
- *Lipophilicity*, a dataset curated from ChEMBL database<sup>10</sup> including experimental results of octanol/water distribution coefficient (*LogD* 7.4) for 4200 compounds<sup>x</sup>.

## Contact us

Please refer to our paper for further information:

- **“Using Molecular Embeddings in QSAR Modeling: Does it Make a Difference?”**, authored by M.V. Sabando, I. Ponzoni, E.E. Milios and A .J. Soto. Preprint available at arXiv.org.

For further details or inquiries regarding the resources in this repository please contact the corresponding author.

## References

- (1) Dalke, A. DeepSMILES: An Adaptation of SMILES for Use in Machine-Learning of Chemical Structures. 2018.
- (2) Landrum, G. RDKit: open-source cheminformatics <http://www.rdkit.org>. 2016.
- (3) Öztürk, H.; Ozkirimli, E.; Özgür, A. A novel methodology on distributed representations of proteins using their interacting ligands. *Bioinformatics* **2018**, *34*, i295–i303.
- (4) Jaeger, S.; Fulle, S.; Turk, S. Mol2vec: unsupervised machine learning approach with chemical intuition. *Journal of Chemical Information and Modeling* **2018**, *58*, 27–35.

---

<sup>x</sup>[https://www.ebi.ac.uk/chembl/document\\_report\\_card/CHEMBL3301361/](https://www.ebi.ac.uk/chembl/document_report_card/CHEMBL3301361/)

- (5) Xu, Z.; Wang, S.; Zhu, F.; Huang, J. Seq2seq fingerprint: An unsupervised deep molecular embedding for drug discovery. Proceedings of the 8th ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics. 2017; pp 285–294.
- (6) Zheng, S.; Yan, X.; Yang, Y.; Xu, J. Identifying structure–property relationships through SMILES syntax analysis with self-attention mechanism. *Journal of Chemical Information and Modeling* **2019**, *59*, 914–923.
- (7) Oskooei, A.; Born, J.; Manica, M.; Subramanian, V.; Sáez-Rodríguez, J.; Rodríguez Martínez, M. PaccMann: Prediction of anticancer compound sensitivity with multi-modal attention-based neural networks. *arXiv e-prints* **2018**, arXiv:1811.06802.
- (8) Delaney, J. S. ESOL: estimating aqueous solubility directly from molecular structure. *Journal of chemical information and computer sciences* **2004**, *44*, 1000–1005.
- (9) Mobley, D. L.; Guthrie, J. P. FreeSolv: a database of experimental and calculated hydration free energies, with input files. *Journal of computer-aided molecular design* **2014**, *28*, 711–720.
- (10) Bento, A. P.; Gaulton, A.; Hersey, A.; Bellis, L. J.; Chambers, J.; Davies, M.; Krüger, F. A.; Light, Y.; Mak, L.; McGlinchey, S., et al. The ChEMBL bioactivity database: an update. *Nucleic acids research* **2014**, *42*, D1083–D1090.