**frontiers**

# NeSyDPP4-QSAR: A Neuro-Symbolic AI Approach for Potent DPP-4-Inhibitor Discovery in Diabetes Treatment

1    **Delower Hossain[1], Ehsan Saghapour[2], Jake Y. Chen[1,2*]**

2    [1] Department of Computer Science, The University of Alabama at Birmingham, AL 35294, USA

3    [2] Department of Biomedical Informatics and Data Science, School of Medicine, The University of

4    Alabama at Birmingham, AL 35205, USA

5    **\* Correspondence:**

6    Corresponding Author

7    **jakechen@uab.edu**

8    **Keywords: DPP-4, Neuro-symbolic AI, Deep Learning**

9    **Abstract**

10    Diabetes Mellitus (DM) is a global epidemic and among the top ten leading causes of mortality (WHO,
11    2019), projected to rank seventh by 2030. The US National Diabetes Statistics Report (2021) states
12    that 38.4 million Americans have diabetes. Dipeptidyl Peptidase-4 (DPP-4) is an FDA-approved target
13    for type 2 diabetes mellitus (T2DM) treatment. However, current DPP-4 inhibitors are associated with
14    adverse effects, including gastrointestinal issues, severe joint pain (FDA safety warning),
15    nasopharyngitis, hypersensitivity, and nausea. Identifying novel inhibitors is crucial. Direct in vivo
16    DPP-4 inhibition assessment is costly and impractical, making in silico IC50 prediction a viable
17    alternative. Quantitative Structure-Activity Relationship (QSAR) modeling is a widely used
18    computational approach for chemical substance assessment.

19    We employ LTN, a neuro-symbolic approach, alongside DNN and transformers as baselines. DPP-4-
20    related data is sourced from PubChem, ChEMBL, BindingDB, and GTP, comprising 6,563 bioactivity
21    records (SMILES-based compounds with IC50 values) after deduplication and thresholding. A diverse
22    set of features including descriptors (CDK Extended-PaDEL), fingerprints (Morgan), chemical
23    language model embeddings (ChemBERTa2), LLaMa 3.2, and physicochemical properties is used to
24    train the NeSyDPP4-QSAR model.

25    The NeSyDPP4-QSAR model yielded the highest accuracy, incorporating CDKextended and Morgan
26    fingerprints, with an accuracy of 0.9725, an F1-score of 0.9723, an ROC AUC of 0.9719, and an MCC
27    of 0.9446. The performance was benchmarked against two standard baseline models: a deep neural
28    network and a transformer. To ensure fair comparisons, DNN models used the equivalent attributes
29    with the same dimension and network configuration as NeSyDPP4-QSAR. Our findings showed that
30    integrating the Neuro-symbolic strategy (neural network-based learning and symbolic reasoning) holds
31    immense potential for discovering drugs that can inhibit diabetes mellitus and classifying biological
32    activities that inhibit it.

33    **1    Introduction**

34    Diabetes Mellitus (DM) is a chronic metabolic disorder characterized by elevated blood glucose levels,
35    posing a significant global health burden. According to the World Health Organization (WHO) 2019

report, diabetes ranks among the top ten leading causes of mortality, with an estimated 1.6 million deaths worldwide [1-2]. In the United States, diabetes is a major public health challenge, affecting approximately 38 million people (11.3% of the population) and leading to $327 billion in medical expenses and lost wages annually [3]. Beyond economic costs, diabetes is associated with severe complications, including blindness, kidney failure, stroke, heart disease, and neuropathy.

DM is broadly classified into Type 1 Diabetes Mellitus (T1DM) and Type 2 Diabetes Mellitus (T2DM), with T2DM accounting for over 90% of all cases. One crucial therapeutic target for T2DM management is the Dipeptidyl Peptidase-4 (DPP-4) enzyme, which regulates glucose metabolism. DPP-4 inhibitors, a class of FDA-approved medications, help control blood sugar levels by inhibiting this enzyme. However, current DPP-4 inhibitors have been linked to adverse effects such as gastrointestinal issues, severe joint pain, nasopharyngitis, hypersensitivity, and nausea [4]. As a result, discovering safer and more effective DPP-4 inhibitors remains a critical research challenge.

Artificial Intelligence (AI) has revolutionized diabetes management and drug discovery over the past two decades. Early AI models focused on glucose level prediction, insulin dosage recommendations, and patient monitoring. In recent years, AI has expanded into de novo drug design, utilizing vast molecular datasets to identify new inhibitors and analyze complex relationships between genes, proteins, and disease mechanisms. In the field of DPP-4 inhibitor prediction, Quantitative Structure-Activity Relationship (QSAR) models have been widely employed using machine learning techniques such as Random Forest, Support Vector Machines, XGBoost, Gradient Boosting Machines, and Deep Neural Networks [5-10]. While these models have demonstrated high predictive performance, they suffer from limitations, including poor interpretability, data inefficiency, and a lack of reasoning capabilities. The black-box nature of deep learning models further complicates their use in critical healthcare applications, where transparency and explainability are essential.

To address these challenges, Neuro-Symbolic AI (NeSy) has emerged as a promising paradigm that combines neural networks with symbolic reasoning for more interpretable and data-efficient learning. Unlike traditional AI approaches, NeSy AI enables models to integrate domain knowledge and perform logical reasoning, making them particularly suited for bioactivity prediction in drug discovery. Several NeSy models have already demonstrated success in biomedical applications [11-13], such as: Protein Function Prediction (MultipredGO [14]), Gene Sequence Analysis (KBANN [15]), Diabetic Retinopathy Diagnosis (ExplainDR [16]), (Gene Sequence) KBANN [17], hERG-LTN [18], (Ontology) RRN [19], NSRL [20], Neuro-Fuzzy [21], FSKBANN [22], DeepMiRGO [23], NS-VQA [24], DFOL-VQA [25], LNN [26], NofM [27], PP-DKL [28], FSD [29], CORGI [30], NeurASP [31], XNMs [32], Semantic loss [33], NS-CL [34], LTN [35],. This study investigates the role of a hybrid Neuro-symbolic model integrating Logic Tensor Networks (LTN) for DPP-4 bioactivity prediction. Our objective is to identify potential DPP-4 inhibitors for T2DM treatment while improving prediction accuracy.

*Key Contributions to This Study*

The significant key contribution of this study is: 1) we built a scalable, robust AI predictive model with immense accuracy improvement for T2DM inhibitors potency prediction. 2) A novel representation integrating data and rules (Knowledge) for DPP-4 inhibitor bio-activity classification 3) Acquired and utilized more diverse compound datasets with chemical embedding, descriptor, fingerprints, physiochemical properties that previous studies have not utilized. The proposed NeSyDPP4 can be used to discover novel DPP-4 active drugs by scanning large molecular datasets like ZINC, and identification of novel candidate compounds, accelerates de novo drug design. Additionally, it facilitates QSAR model downstream applications such as virtual screening, contraindications,

81  bioactivity indications, and other key elements of DPP-4 inhibitors therapy in the clinical setting
82  including docking, affinity prediction, ADMET analysis, and molecular dynamics (MD) studies for
83  DPP-4 clinical settings.

84  The remainder of this paper is structured as follows: Section II describes the methodology, Section III
85  presents the experimental results, and Section IV Discussion, and finally concludes with key findings
86  and future research directions. `

## 1.1    Data acquisition

88  The study constructed a new DPP-IV cohorts utilized four publicly available chemical compound
89  databases: ChEMBL [36] & BindingDB [37], PubChem, and GTP, The ChEMBL Database contains
90  more than 2 million compounds. We retrieved canonical SMILES related to DPP-4 inhibitor with the
91  target organism Homo Sapiens using ID: CHEMBL284 and standard type IC50. The data was extracted
92  using the ChEMBL Python API (chembl_webresource_client). The BindingDB manually uses DPP-4
93  string keywords (dipeptidyl peptidase-4) from their official site. In addition, PubChem in CSV format
94  with following(link), and GTP via the corresponding (link).

## 1.2    Data preparation and feature extraction

96  The initial bioactivity data remains various irrelevant attributes. We collected subsets focused on the
97  IC50 biological activity standard value, ChEMBL inhibitors ID, and Canonical SMILES. However,
98  numerical IC50 measurements in nM were given in ChEMBL, BindingDB, and GTP, but those in μM
99  were given in PubChem, and were harmonized all units into nanomolar (nM).  Subsequently, we
100 calculated pIC50 values from the IC50 values, applying a normalization step through log10 conversion
101 (equ. 1). Active and inactive label determined based on pIC50 by following previous DPP-iV chemical
102 research article [38]. Afterwards, a diverse array of features was extracted from SMILES representations,
103 encompassing Morgan fingerprints (512, 1024, and 2048 bits), CDKextended descriptors utilizing PaDELPy
104 [39], chemical embeddings generated via ChemBERTa2 and LLaMA3.2, as well as a comprehensive set of
105 physicochemical properties using RDkit [40].

106  Finally, ML trainable data comprised a total of 6563 upon dropping duplicates and NaN values.

107
$$\text{pIC}_{50} = -\log_{10}(\text{IC}_{50} \times 10^{-9}) \qquad (1)$$

## 1.3    LTN classification model

109 LTNs [35] were architected using two key components: a logic component and a neural network. The
110 visual architecture of the classification model can be found in Appendix A. The logical mechanism
111 contains a set of axioms or rules (explained in detail in the Knowledge-based setting). It's important to
112 note that LTN logical reasoning reveals through rules/axioms.  In our context Table 1 represents the
113 axioms and relevant knowledge base component. However, other network configuration parameters
114 can be found in Table 1.

115

| **Table 1:** LTN Knowledge-based Setting for DPP-IV Classification | |
|---|---|
| **Contents** | **Classification** |
| Define Axioms | • $\forall x_A, p(x_A, l_A)$: all the examples of class $A(Active = 0)$ should have a label $l_A$<br>• $\forall x_B, p(x_B, l_B)$: all the examples of class $B$ (Inactive = 1) should have a label $l_B$ |

| Axioms (rules, knowledge base) | $\mathcal{K} = \forall x_A p(x_A, l_A), \forall x_B p(x_B, l_B)$ |
|---|---|
| SatAgg is given by | $\text{SatAgg}_{\phi \in \mathcal{K}} \mathcal{G}_{\theta, x \leftarrow D}(\phi)$ |
| Learning & Loss | $\boldsymbol{L} = \left(1 - \text{SatAgg}_{\phi \in \mathcal{K}} \mathcal{G}_{\boldsymbol{\theta}, x \leftarrow \boldsymbol{B}}(\phi)\right)$ |
| **Note**: This table was developed inspired by the official LTN | |

116    Here,

117    The pMeanError aggregator

$$pME(u_1, \ldots, u_n) = 1 - \left(\frac{1}{n}\sum_{i=1}^{n}(1-u_i)^p\right)^{\frac{1}{p}} p \geq 1 \quad (2)$$

118  • SatAgg: This stands for "Satisfaction Aggregator"
119  • $\phi \in K$: This part indicates that $\phi$ (phi) belongs to the set K. $\phi$ is often used to represent a predicate.
120  • $\mathcal{G}(\theta)$: This is denoted by grounding ($\mathcal{G}$) with parameters $\theta$. $\theta$ represents a set of parameters or weights in a
121     model.
122  • $x \leftarrow D$: $D$ the data set of all examples (domain).
123  • The input to the functions SatAgg and $\mathcal{G}(\theta)$

124    In addition to experimenting with LTN, we conducted the simulation with DNN, and transformer with keras
125    integrated for the fair comparison of with LTN performance. Table 2 depicts the network configuration
126    parameters.

**Table 2: LTN, DNN, Transformer Models Parameters Summary Classification**

| Parameters | LTN | DNN |
|---|---|---|
| **Activation** | ReLU | ReLU |
| **Units** | (768,384,192,2) | (768,384,192,2) |
| **No of Dense layers** | 4 | 4 |
| **Seed** | 42 | 42 |
| **Batch Size** | 128 | 128 |
| **Training Epochs** | 100 | 100 |
| **Learning Rate** | 0.00001 | 0.00001 |
| **Loss Function** | LTN pMeanError | Sparse Categorical Crossentropy |
| **Optimizer** | Adam | Adam |

Note: The input depends on selected features (e.g., fingerprint, embedding, etc).
Transformer configuration parameter can be found on this project Github.

## 1.4   Model Training and Validation Phase

128    LTN, DNN, and transformer models were trained and tested using TensorFlow 2.15.1 Python 3.10.16 on UAB
129    server, NVIDIA A100 80GB PCIe, other dependency packages can be found on this project GitHub
130    environment.yml. In the training phase, we did partition the data as 80:10:20 ratios over 100 epochs during.
131    while following metrics, such as Accuracy, F-score (F), ROC AUC Score, and Mathew Correlation
132    Coefficient (MCC), were used to assess the trained model's performance, and the misclassified classes
133    can appear in the Fig. 2.

4

**Equation 1** Accuracy

$$\text{Accuracy} = \frac{TP + TN}{TP + TN + FP + FN} \tag{3}$$

**Equation 2** F1 Score

$$F_1 = \frac{2 \; x \; Precision \; x \; Recall}{Precision + Recall} \tag{4}$$

**Equation 3** ROC AUC Score

$$\text{ROC AUC} = \int_0^1 TPR \; d(FPR) \tag{5}$$

**Equation 4** MCC

$$\text{MCC} = \frac{TP \times TN - FP \times FN}{\sqrt{(TP + FP) \times (TP + FN) \times (TN + FP) \times (TN + FN)}} \tag{6}$$

134

## 2    Result

136  Here, we describe the performance of the developed NeSyDPP4 Model, for revealing DPP4 potential
137  inhibitors leveraging LTN architecture (rules Integration into the neural network). DNN, and
138  transformer since raw data is string format. We computed diverse features with the respective
139  smiles/drugs such as morgan fingerprint, CDKExtended descriptor, Chemical foundation language
140  model embedding using ChemBERTa2, LLaMA3.2 embedding, Physiochemical properties using
141  RDkit. There are three tables in this section. Table 3 shows all the features separated and combination
142  input results, Table 4 exposes the fair comparison with baseline DNN, and transformer architecture
143  performance.

144  In Table 3 depicts the different input performance of LTN. The best-performing feature set is
145  combining CDKExtended + ECFP, which yielded the highest Accuracy (97.25%), F1-score (97.23%),
146  AUC-ROC (97.19%), and MCC (94.46, while physicochemical features alone yield the lowest
147  performance. ChemBERTa2 and Llama3.2 performed comparably but were lower than fingerprint-
148  based methods of Accuracy (73.49%), F1-score (73.16%), AUC-ROC (73.09%), and MCC (46.38%),
149  Overall, suggesting that physicochemical properties alone are insufficient for effective bioactivity
150  classification.

151  Furthermore, Table 4 demonstrates the comparison performance and efficiency of three models, LTN
152  and DNN, and transformer for predicting the properties of molecules associated with T2DM DPP-IV
153  inhibitors.

154  Upon obtaining the accurate features of chemical compounds from Table 1 experiment, we proceeded
155  experiment with DNN utilizing same architecture and similar input as shows Table. The LTN model

156 with CDKExtended + ECFP features outperforms the others, achieving 97.25% accuracy and 94.46%
157 MCC, demonstrating the effectiveness of neuro-symbolic reasoning. The DNN model, using the same
158 features, performs slightly lower (96.95% accuracy, 93.85% MCC), which indicates that LTN's logical
159 constraints enhance predictions. In contrast, the Transformer model with SMILES embeddings shows
160 the lowest performance (78.21% accuracy, 56.41% MCC), suggesting that fingerprint-based features
161 are more effective than SMILES-based embeddings for bioactivity classification. However, the Fig
162 illustrated the highest misclassification occurred by transformers since performance is lowest
163 compared to three model simulation.

**Table 4:** LTN DPP4 Bio-activity Classification Result Summary

| Model | Features | Input | Acc | F1 | AUC ROC | MCC |
|---|---|---|---|---|---|---|
| | CDKExetended + ECFP | 1024+(512+1024+2048) | **0.9725** | **0.9723** | **0.9719** | **0.9446** |
| | ECFP | 1024 | 0.9687 | 0.9684 | 0.9680 | 0.9370 |
| | ECFP | 2048 | 0.9657 | 0.9654 | 0.9650 | 0.9308 |
| | ECFP | 512 | 0.9649 | 0.9646 | 0.9643 | 0.9293 |
| LTN | Combined All | 7430 | 0.9634 | 0.9631 | 0.9632 | 0.9262 |
| | CDKExetended | 1024 | 0.9504 | 0.9499 | 0.9492 | 0.9001 |
| | ChemBERTa2 | 768 | 0.8956 | 0.8944 | 0.8935 | 0.7892 |
| | Llama3.2 | 2048 | 0.8933 | 0.8926 | 0.8933 | 0.7854 |
| | Physiochemical | 6 | 0.7349 | 0.7316 | 0.7309 | 0.4638 |

164

**Table 4:** LTN DPP4 Bio-activity Classification Result Summary

| Model | Features | Input Dimension | Acc | F1 | AUC ROC | MCC |
|---|---|---|---|---|---|---|
| **LTN** | CDKExetended + ECFP | 1024+(512+1024+2048) | **0.9725** | **0.9723** | **0.9719** | **0.9446** |
| **DNN** | CDKExetended + ECFP | 1024+(512+1024+2048) | 0.9695 | 0.9692 | 0.9691 | 0.9385 |
| **Transformer** | SMILES/Emb | 212 | 0.7821 | 0.7306 | 0.8549 | 0.5641 |

165

6

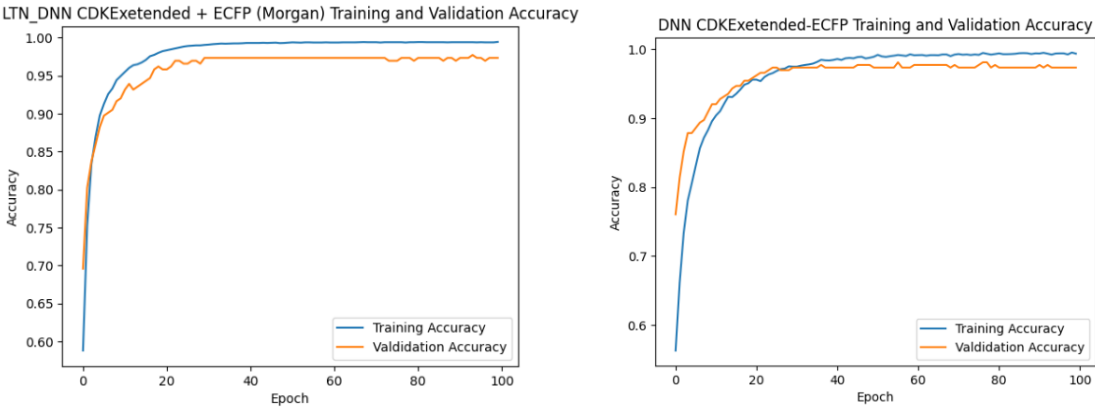| Model | Author | Metrics | result | Ref |
|---|---|---|---|---|
| NeSyDPP4 | Hossain *et al* | Accuracy | 0.9725 | |
| Random Forest | Oky Hermansyah *et al* | Accuracy | 0.9221 | 42 |
| DNN | Haris Hamzah *et al* | Accuracy | 0.9060 | 43 |
| QSAR-DNN | Alhadi Bustamam *et al* | Accuracy | 0.9040 | 9 |
| NB | Jie Cai *et al* | Accuracy | 0.8720 | 44 |
| Conv1D-LSTM | Adawiyah Ulfa *et al* | Accuracy | 0.8618 | 38 |
| XGBoost | Oky Hermansyah *et al* | Accuracy | 0.8164 | 45 |

166



167

Fig.1, Epoch and Accuracy curve during the training and validation phase of LTN and DNN model
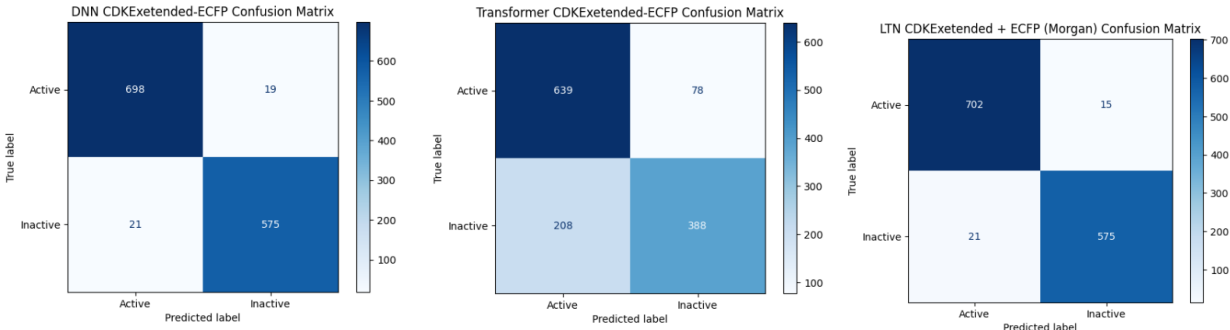
169



170

Fig.2, Classification matrix of DNN, Transformer, and LTN using CDKextended+Morgan all bit's features, it depicts that Transformer misclassified highest number of the samples.

## 3    Discussion

This article aimed to employ neuro-symbolic modeling (LTN), an integration of data and a logic-driven approach, for predicting diabetes mellitus DPP-4 inhibition. The study's findings provide valuable insights into the applicability and robustness of the LTN model in predicting inhibitor bioactivity behavior. As an illustration, the utilization of this advanced machine learning technique (LTN)

7

178 surpassed the state-of-the-art performance compared to other models with classification tasks, the LTN
179 model demonstrates superior accuracy of 0.9725 and an MCC score of 0.9446 for the DPP-4 inhibitors,
180 while other studies shows the QSAR-DNN model by Bustamam et al. [1] achieved an accuracy of
181 0.9040, Ulfa et al. [2] reported an accuracy of 0.8618 using Conv1D-LSTM. Random Forest by
182 Hermansyah et al. [3] yielded an accuracy of 0.9221. DNN by hamzah et al. [4] obtained an accuracy
183 of 0.9060. NB by Cai et al. [5] gained an accuracy of 0.8720. XGBoost by Hermansyah et al. [6]
184 achieved an accuracy of 0.8164.

185 The implications drawn from this research are profound. The utilization of neuro-symbolic modeling
186 (LTN), blending data-driven and knowledge-driven methodologies has shown remarkable potential in
187 predicting diabetes mellitus through DPP-4 inhibitors activity classification. Thus, this research tiles
188 the way for advanced machine learning applications in diabetes prediction and marks a significant step
189 forward in understanding inhibitor behavior and its implications for DM. These findings advocate for
190 the transformative potential of LTN in diabetes prediction and emphasize the value of further
191 exploration and implementation of neuro-symbolic strategies in healthcare research and applications.

192 *Limitation*

193 Acknowledging the limitations of our study, we state that while LTN has demonstrated significant
194 promise, it may be uncapable to incorporate external biological additional knowledge with neural
195 networks.

196 **4    Conclusion**

197 Diabetes Mellitus is a vital global health concern, and discovering effective chemical substances is
198 crucial to tackling this epidemic. This study intend to develop QSAR system for the therapeutic
199 potential of DPP-4 inhibitors employing a novel approach called the LTN (Neuro-symbolic AI) that
200 integrates domain-specific knowledge into neural networks. The study is a pioneer in applying Neuro-
201 symbolic strategy in the DM domain and provides new insights showing groundbreaking performance
202 for revealing DPP-4 potential inhibitors. The root cause of achieving such performance could be
203 upholding learning and reasoning principles and training neural networks with rules. Furthermore, we
204 experimented with DNN, an NLP Transformer model, whereas LTN also attained prominent Accuracy.

205 In conclusion, the findings of this study prove that LTN is among the state-of-the-art models for
206 uncovering potential DPP-4 inhibitors. We aim to deploy the model within a real-time prediction
207 application to identify the right therapeutic agent that could promptly benefit ML practitioners,
208 academics, and industry researchers. However, an ideal next step could involve integrating additional
209 potential Neuro-symbolic strategies, such as Semantic Loss, DeepProblog on GLP-1, IDO, and PTP1B
210 DM inhibitors extracting a variety of new descriptors, and fingerprints with different datasets
211 (PubChem, Protein Data Bank) focusing Regression Task.

212 **5    Conflict of Interest**

213 The authors declare that they have no conflicts of interests in this work.

214 **6    Author Contributions**

215 The author, Delower Hossain, designed, implemented, and wrote the manuscripts, and Ehsan
216 Saghapour worked together to edit and review. Dr. Jake Chen has been actively guided as project
217 administrator.

218

## 9    Data Availability Statement

The dataset that utilized in this study can be found here link
And experimented code repo can be found here

## 10    References

1. World Health Organization: WHO. (2024, August 7). *The top 10 causes of death*. https://www.who.int/news-room/fact-sheets/detail/the-top-10-causes-of-death#:~:text=Of%20the%2056.9%20million%20deaths%20worldwide%20in%202016%2C,of%20death%20globally%20in%20the%20last%2015%20years.

2. World Health Organization. (2020). *World health statistics 2020: monitoring health for the SDGs, sustainable development goals*. https://iris.who.int/bitstream/handle/10665/332070/9789240005105-eng.pdf

3. *Methods for the National Diabetes Statistics Report*. (2024, May 15). Diabetes. https://www.cdc.gov/diabetes/php/data-research/methods.html?CDC_AAref_Val=https://www.cdc.gov/diabetes/data/statistics-report/index.html

4. Huang, J., Jia, Y., Sun, S., & Meng, L. (2020). Adverse event profiles of dipeptidyl peptidase-4 inhibitors: data mining of the public version of the FDA adverse event reporting system. *BMC Pharmacology and Toxicology*, *21*(1). https://doi.org/10.1186/s40360-020-00447-w

5. Yanuar, A., Hermansyah, O., & Bustamam, A. (2019). QSAR Modeling for Prediction of pIC50 DPP-4 Inhibitors with Machine Learning Method. *Conference*. https://conference.ui.ac.id/afps/AFPS-ICAPPS2019/paper/view/25310

6. Hermansyah, O., Bustamam, A., & Yanuar, A. (2021). Virtual screening of dipeptidyl peptidase-4 inhibitors using quantitative structure–activity relationship-based artificial intelligence and molecular docking of hit compounds. *Computational Biology and Chemistry*, *95*, 107597. https://doi.org/10.1016/j.compbiolchem.2021.107597

7. Ojo, O. A., Ojo, A. B., Okolie, C., Abdurrahman, J., Barnabas, M., Evbuomwan, I. O., Atunwa, O. P., Atunwa, B., Iyobhebhe, M., Elebiyo, T. C., Nwonuma, C. O., Adegboyega, A. E., Qusti, S., Alshammari, E. M., Hetta, H. F., & Batiha, G. E. S. (2021). Elucidating the interactions of compounds identified from Aframomum melegueta seeds as promising candidates for the management of diabetes mellitus: A computational approach. *Informatics in Medicine Unlocked*, *26*, 100720. https://doi.org/10.1016/j.imu.2021.100720

250    8.  Septiawan, N. R. R., Prakoso, N. B. H., & Kurniawan, N. I. (2022). DPP IV Inhibitors Activities Prediction as An Anti-
251         Diabetic Agent using Particle Swarm Optimization-Support Vector Machine Method. *Jurnal RESTI (Rekayasa Sistem Dan*
252         *Teknologi Informasi)*, *6*(6), 974–980. https://doi.org/10.29207/resti.v6i6.4470

253    9.  Bustamam, A., Hamzah, H., Husna, N. A., Syarofina, S., Dwimantara, N., Yanuar, A., & Sarwinda, D. (2021). Artificial
254         intelligence paradigm for ligand-based virtual screening on the drug discovery of type 2 diabetes mellitus. *Journal of Big*
255         *Data*, *8*(1). https://doi.org/10.1186/s40537-021-00465-3

256    10. Ajiboye, B. O., Iwaloye, O., Owolabi, O. V., Ejeje, J. N., Okerewa, A., Johnson, O. O., Udebor, A. E., & Oyinloye, B. E.
257         (2021). Screening of potential antidiabetic phytochemicals from Gongronema latifolium leaf against therapeutic targets of
258         type 2 diabetes mellitus: multi-targets drug design. *SN Applied Sciences*, *4*(1). https://doi.org/10.1007/s42452-021-04880-2

259    11.  Hossain D, Chen JY. A Study on Neuro-Symbolic Artificial Intelligence: Healthcare Perspectives. arXiv preprint
260         arXiv:2503.18213. 2025 Mar 23.

261    12.  Yu, D., Yang, B., Liu, D., Wang, H., & Pan, S. (2023). A survey on neural-symbolic learning systems. *Neural Networks*,
262         *166*, 105–126. https://doi.org/10.1016/j.neunet.2023.06.028

263    13. Wang, W., Yang, Y., & Wu, F. (2024). Towards Data-and Knowledge-Driven AI: A survey on Neuro-Symbolic Computing.
264         *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 1–22. https://doi.org/10.1109/tpami.2024.3483273

265    14. Giri, S. J., Dutta, P., Halani, P., & Saha, S. (2020). MulTiPREDGO: Deep Multi-Modal Protein function prediction by
266         amalgamating protein structure, sequence, and interaction information. *IEEE Journal of Biomedical and Health Informatics*,
267         *25*(5), 1832–1838. https://doi.org/10.1109/jbhi.2020.3022806

268    15. Towell, G. G., & Shavlik, J. W. (1994). Knowledge-based artificial neural networks. *Artificial Intelligence*, *70*(1–2), 119–
269         165. https://doi.org/10.1016/0004-3702(94)90105-8

270    16. Jang, S.-I., Girard, M. J. A., & Thiery, A. H. (2021). Explainable Diabetic Retinopathy Classification Based on Neural-
271         Symbolic Learning. *CEUR-WS*, 104–114. http://ceur-ws.org/Vol-2986/paper8.pdf

272    17. Maclin, R., & Shavlik, J. W. (1994). Refining algorithms with knowledge-based neural networks: improving the Chou-
273         Fasman algorithm for protein folding. *Conference on Learning Theory*, 249–286. http://dl.acm.org/citation.cfm?id=188535

274    18.  Hossain D, Chen JY, Abir FA. hERG-LTN: A New Paradigm in hERG Cardiotoxicity Assessment Using Neuro-Symbolic
275         and Generative AI Embedding (MegaMolBART, Llama3. 2, Gemini, DeepSeek) Approach. bioRxiv. 2025:2025-02.

276    19. Yang, F., Yang, Z., & Cohen, W. W. (2017). Differentiable Learning of Logical Rules for Knowledge Base Reasoning.
277         *NeurIPS*. https://proceedings.neurips.cc/paper_files/paper/2017/hash/0e55666a4ad822e0e34299df3591d979-Abstract.html

278    20. Hohenecker, P., & Lukasiewicz, T. (2020). Ontology Reasoning with Deep Neural Networks. *Journal of Artificial*
279         *Intelligence Research*, *68*. https://doi.org/10.1613/jair.1.11661

280    21. Yang, Z., Ishay, A., & Lee, J. (2020). NeurASP: Embracing Neural Networks into Answer Set Programming. *Proceedings of*
281         *the Twenty-Ninth International Joint Conference on Artificial Intelligence (IJCAI-20)*, 1755.
282         https://www.ijcai.org/proceedings/2020/0243.pdf

283  22. Kora, P., Meenakshi, K., Swaraja, K., Rajani, A., & Islam, M. K. (2019). Detection of Cardiac arrhythmia using fuzzy logic.

284      *Informatics in Medicine Unlocked*, *17*, 100257. https://doi.org/10.1016/j.imu.2019.100257

285  23. Wang, J., Zhang, J., Cai, Y., & Deng, L. (2019). DEEPMIR2GO: Inferring functions of human MicroRNAs using a deep

286      Multi-Label Classification model. *International Journal of Molecular Sciences*, *20*(23), 6046.

287      https://doi.org/10.3390/ijms20236046

288  24. Yi, K., Wu, J., Gan, C., Torralba, A., Kohli, P., & Tenenbaum, J. B. (2018). Neural-Symbolic VQA: Disentangling

289      Reasoning from Vision and Language Understanding. *arXiv (Cornell University)*, *31*, 1031–1042.

290      http://arxiv.org/pdf/1810.02338.pdf

291  25. Amizadeh, S., Palangi, H., Polozov, O., Huang, Y., & Koishida, K. (2020). Neuro-Symbolic Visual Reasoning:

292      Disentangling "Visual" from "Reasoning." *International Conference on Machine Learning*, *1*, 279–290.

293      http://proceedings.mlr.press/v119/amizadeh20a.html

294  26. Riegel, R., Gray, A. G., Luus, F. P. S., Khan, N., Makondo, N., Akhalwaya, I. Y., Qian, H., Fagin, R., Barahona, F., Sharma,

295      U., Ikbal, S., Karanam, H., Neelam, S., Likhyani, A., & Srivastava, S. K. (2020). Logical neural networks. *arXiv (Cornell*

296      *University)*. https://doi.org/10.48550/arxiv.2006.13155

297  27. Towell, G., & Shavlik, J. W. (1991). Interpretation of Artificial Neural Networks: Mapping Knowledge-Based Neural

298      Networks into Rules. *Neural Information Processing Systems*, *4*, 977–984. http://papers.nips.cc/paper/546-interpretation-of-

299      artificial-neural-networks-mapping-knowledge-based-neural-networks-into-rules.pdf

300  28. Lavin, A. (2022). Neuro-Symbolic neurodegenerative disease modeling as probabilistic programmed deep kernels. In *Studies*

301      *in computational intelligence* (pp. 49–64). https://doi.org/10.1007/978-3-030-93080-6_5

302  29. Dobosz, K., & Duch, W. (2008). Fuzzy symbolic dynamics for neurodynamical systems. In *Lecture notes in computer*

303      *science* (pp. 471–478). https://doi.org/10.1007/978-3-540-87559-8_49

304  30. Arabshahi, F., Lee, J., Gawarecki, M., Mazaitis, K., Azaria, A., & Mitchell, T. (2021). Conversational Neuro-Symbolic

305      commonsense reasoning. *Proceedings of the AAAI Conference on Artificial Intelligence*, *35*(6), 4902–4911.

306      https://doi.org/10.1609/aaai.v35i6.16623

307  31.  Shi, J., Zhang, H., & Li, J. (2019). Explainable and explicit visual reasoning over scene graphs. *2022 IEEE/CVF Conference*

308      *on Computer Vision and Pattern Recognition (CVPR)*, 8368–8376. https://doi.org/10.1109/cvpr.2019.00857

309  32. Teru, K., Denis, E., & Hamilton, W. (2020). Inductive relation prediction by subgraph reasoning. *International Conference*

310      *on Machine Learning*, *1*, 9448–9457. http://proceedings.mlr.press/v119/teru20a/teru20a.pdf

311  33. Xu, J., Zhang, Z., Friedman, T., Liang, Y., & Broeck, G. (2018, July 3). *A Semantic Loss Function for Deep Learning with*

312      *Symbolic Knowledge*. PMLR. https://proceedings.mlr.press/v80/xu18h.html

313  34. Mao, J., Gan, C., Kohli, P., Tenenbaum, J. B., & Wu, J. (2019). The Neuro-Symbolic Concept Learner: Interpreting scenes,

314      words, and sentences from natural supervision. *International Conference on Learning Representations*.

315      https://openreview.net/pdf?id=rJgMlhRctm

11

316    35.  Badreddine, S., Garcez, A. D., Serafini, L., & Spranger, M. (2021). Logic Tensor networks. *Artificial Intelligence*, *303*,
317         103649. https://doi.org/10.1016/j.artint.2021.103649

318    36.  Gaulton, A., Bellis, L. J., Bento, A. P., Chambers, J., Davies, M., Hersey, A., Light, Y., McGlinchey, S., Michalovich, D.,
319         Al-Lazikani, B., & Overington, J. P. (2011). ChEMBL: a large-scale bioactivity database for drug discovery. *Nucleic Acids*
320         *Research*, *40*(D1), D1100–D1107. https://doi.org/10.1093/nar/gkr777

321    37.  Gilson, M. K., Liu, T., Baitaluk, M., Nicola, G., Hwang, L., & Chong, J. (2015). BindingDB in 2015: A public database for
322         medicinal chemistry, computational chemistry and systems pharmacology. *Nucleic Acids Research*, *44*(D1), D1045–D1053.
323         https://doi.org/10.1093/nar/gkv1072

324    38.  Ulfa, A., Bustamam, A., Yanuar, A., Amalia, R., & Anki, P. (2021). Model QSAR Classification Using Conv1D-LSTM of
325         Dipeptidyl Peptidase-4 Inhibitors. *IEEExplore*, *pp.160-163*, 1–6. https://doi.org/10.1109/aims52415.2021.9466083

326    39.  *PaDEL*. (n.d.). GitHub. https://github.com/dataprofessor/padel.git
327         Ecrl. (n.d.). *GitHub - ecrl/padelpy: A Python wrapper for PaDEL-Descriptor software*. GitHub.
328         https://github.com/ecrl/padelpy/tree/master

329    40.  *Installation — The RDKit 2024.09.6 documentation*. (n.d.). https://www.rdkit.org/docs/Install.html

330    41.  Hermansyah, O., Bustamam, A., & Yanuar, A. (2021). Virtual screening of dipeptidyl peptidase-4 inhibitors using
331         quantitative structure–activity relationship-based artificial intelligence and molecular docking of hit compounds.
332         *Computational Biology and Chemistry*, *95*, 107597. https://doi.org/10.1016/j.compbiolchem.2021.107597

333    42.  Hermansyah, O., Bustamam, A., & Yanuar, A. (2020). Virtual screening of DPP-4 inhibitors using QSAR-Based artificial
334         intelligence and molecular docking of HIT compounds to DPP-8 and DPP-9 enzymes. *Research Square (Research Square)*.
335         https://doi.org/10.21203/rs.2.22282/v2

336    43.  Hamzah, H., Bustamam, A., Yanuar, A., & Sarwinda, D. (2020). Predicting The Molecular Structure Relationship and The
337         Biological Activity of DPP-4 Inhibitor Using Deep Neural Network with CatBoost Method as Feature Selection.
338         *IEEEXplore*, 101–108. https://doi.org/10.1109/icacsis51025.2020.9263204

339    44.  Cai, J., Li, C., Liu, Z., Du, J., Ye, J., Gu, Q., & Xu, J. (2017). Predicting DPP-IV inhibitors with machine learning
340         approaches. *Journal of Computer-Aided Molecular Design*, *31*(4), 393–402. https://doi.org/10.1007/s10822-017-0009-6

341    45.  Identification of DPP-4 inhibitor active compounds using machine learning classification. (2023). *ResearchGate*.

342    46.  FDA approved Dipeptidyl Peptidase IV (DPP IV) Inhibitors
343         https://www.ncbi.nlm.nih.gov/books/NBK542331/#:~:text=DPP%2D4%20inhibitors%2C%20known%20as,sax
344         agliptin%2C%20linagliptin%2C%20and%20alogliptin

345    47.  Wikipedia DPP-4 Inhibitors https://en.wikipedia.org/wiki/Dipeptidyl_peptidase-4_inhibitor

346

347     **Appendix A**: LTN Model Architecture for multiclass classification.

348

349

350
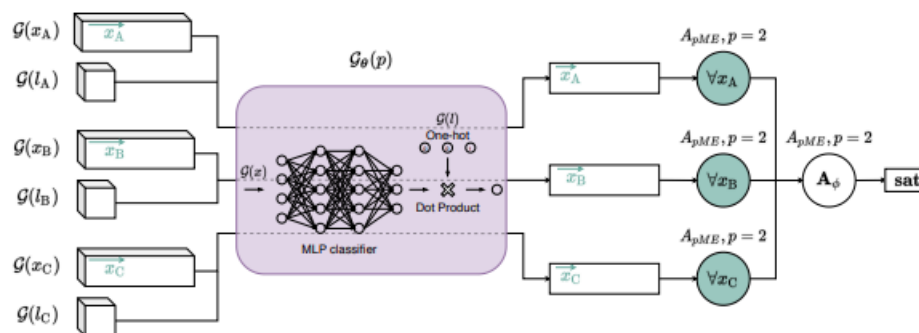
351

352

353

354



Fig. 7: LTN Classification Architecture [36]

355     **Appendix B**: A list of FDA, EU, EMA (European Medicines Agency), JAPAN, and KOREN BODY
356     approved DPP4 inhibitor's structure and respective 3D compound structures images as below.

357

| ChEMBL ID | Target | Approved Body | Smiles | Ref |
|---|---|---|---|---|
| CHEMBL376359 | Alogliptin | FDA | Cn1c(=O)cc(N2CCC[C@@H](N)C2)n(Cc2ccccc2C#N)c1=O | [46] |
| CHEMBL1929396 | Anagliptin | Japan | Cc1cc2ncc(C(=O)NCC(C)(C)NCC(=O)N3CCC[C@H]3C#N)cn2n1 | [46-47] |
| CHEMBL3707235 | Gemigliptin | Korea | N[C@@H](CC(=O)N1CCc2c(nc(C(F)(F)F)nc2C(F)(F)F)C1)CN1CC(F)(F)CCC1=O | [46-47] |
| CHEMBL237500 | Linagliptin | FDA | CC#CCn1c(N2CCC[C@@H](N)C2)nc2c1c(=O)n(Cc1nc(C)c3ccccc3n1)c(=O)n2C | [46] |
| CHEMBL385517 | Saxagliptin | FDA | N#C[C@@H]1C[C@@H]2C[C@@H]2N1C(=O)[C@@H](N)C12CC3CC(CC(O)(C3)C1)C2 | [46] |
| CHEMBL1422 | Sitagliptin | FDA | N[C@@H](CC(=O)N1CCn2c(nnc2C(F)(F)F)C1)Cc1cc(F)c(F)cc1F | [46] |
| CHEMBL2147777 | Teneligliptin | Japan | Cc1cc(N2CCN([C@@H]3CN[C@H](C(=O)N4CCSC4)C3)CC2)n(-c2ccccc2)n1 | **Error! Reference source not found.**[46-47] |
| CHEMBL142703 | Vildagliptin | EMA | N#C[C@@H]1CCCN1C(=O)CNC12CC3CC(CC(O)(C3)C1)C2 | [46-47] |

358

359

360

13

361

**Appendix C**: LTN / Knowledge-based Setting

The construction of all the axioms components conceived from the official LTN framework [35].

Classification:

- *Domains*
  - *items*, denoting the examples from the DPP-4 dataset
  - *labels*, representing the class labels (IC50 values)
- *Define Variables*
  - $x_{active}, x_{inactive},$ , for the positive examples of classes $A$ *and* $B$
  - $x$ for all examples
  - $D(x_A) = D(x_B) = D(x) = items$
- *Define Constants*
  - $L_{active}, L_{inactive}$ the labels of classes $A(0)$ *and* $B(1)$ Respectively.
  - $D(l_A) = D(l_B) = labels$ (active inactive pic50 based)
- *Define the P predicate.*

  - $\rho(x, l)$ Denoting the fact that item $x$ is classified as $l$;
  - $D_{in}(P) = items, labels$.


- *Connectives:*
  - *For All $\forall$, And $\wedge$, Not $\neg$, Or $\vee$, Implies $\Rightarrow$*
- *Axiom*
  - $\forall x_A, p(x_A, l_A)$: all the examples of class $A(active)$ should have a label $l_A$

  - $\forall x_B, p(x_B, x_B)$: all the examples of class $B$ (*Inactive*) should have a label $l_B$


Notice that rules about exclusiveness, such as $\forall \left( P(x, l_A) \Rightarrow \left( \neg P(x, l_B) \wedge, \neg P(x, l_C) \right) \right)$ They are omitted since such constraints are already imposed by the grounding of $P$, below, more specifically by the *softmax* function.

- Grounding:
  - $\mathcal{G}(items) = R^N$, items are described by $N$ features:
  - $\mathcal{G}(labels) = N^2$, We use an encoding to represent classes.
  - $\mathcal{G}(x_{active}) \in R^{m_1 \times N}$, that is, $\mathcal{G}(x_{active})$ is a sequence of $m_1$ examples of class active;

14

407    ○   $\mathcal{G}(x_{inactive}) \in R^{m_2 \times N}$, that is, $\mathcal{G}(x_{inactive})$ is a sequence of $m_2$ examples of
408         class inactive;
409

410    ○   $\mathcal{G}(x) \in R^{(m_1+m_2) \times N}$, that is, $\mathcal{G}(x)$ It is a sequence of all the examples.
411

412    ○   $\mathcal{G}(l_A) = 0, \mathcal{G}(l_B) = 1$;
413

414         $\mathcal{G}(P \mid \theta): x, l \mapsto l^\top \cdot softmax(MLP_\theta(x))$, where $MLP$ has two output neurons
415         corresponding to as many classes, notably in our cases, two classes as we explained
416         early, and $\cdot$ denotes the dot product as a way of selecting an output for $\mathcal{G}(P \mid \theta)$.
417         Multiplying the $MLP$ output by the probability. $l^\top$ Gives the probability corresponding
418         to the class denoted by $l$.

419

420

421