

## TRIBUNAL EXAMINADOR

Este proyecto fue aprobado por el Tribunal Examinador de la carrera: Maestría Profesional en Ingeniería del Software con Énfasis en Inteligencia Artificial, requisito para optar por el título de grado de Maestría Profesional en Ingeniería del Software con Énfasis en Inteligencia Artificial, para los estudiantes: Hellen Pamela Aguilar Noguera y José Leonardo Araya Parajeles

---

*M.Sc. Rodrigo Herrera Garro*  
**Tutor**

---

*M.Sc. Ignacio Trejos Zelaya*  
**Lector 1**

---

*M.Sc. David Gerardo Alfaro Víquez*  
**Lector 2**

San José, Costa Rica, 23 de Junio de 2026

**Universidad CENFOTEC**

Maestría en Software con énfasis en Inteligencia Artificial

*Documento Final de Proyecto de Investigación Aplicada 2*

**Red Neuronal para Ranking de Similitud Molecular Multi-Métrica  
de Benzimidazoles**

Hellen Pamela Aguilar Noguera

Jose Leonardo Araya Parajeles

Universidad CENFOTEC

*En colaboración con la Universidad de Salamanca (USAL)*

**Tutor: Rodrigo Herrera Garro**

Marzo, 2026

Introducción	3
Materiales y Métodos	4
Conjunto de Datos	4
Preparación de Datos	5
Tokenización de SMILES	5
Arquitectura del Modelo	6
Procedimiento de Entrenamiento	6
Métricas de Evaluación	7
Implementación de la Aplicación Web	7
Resultados	8
Convergencia del Entrenamiento	8
Métricas de Evaluación	8
Validación Cualitativa por Inferencia	9
Discusión	10
Conclusiones	11
Referencias	12

### **Abstract**

This document presents the design, training, and evaluation of a neural network for simultaneously predicting four molecular similarity metrics between chemical compounds: 2D Tanimoto (ECFP4), Cosine similarity, 3D Shape similarity (OpenEye ROCS), and Combo Tanimoto. The model was trained on 48,462 molecular pairs derived from 20,095 ChEMBL compounds with real 3D similarity data generated using OpenEye OMEGA/ROCS, using 12 experimentally tested benzimidazoles as reference compounds. The convolutional neural network (CNN) architecture, with 399,300 trainable parameters, achieved an average  $R^2$  of 0.9353 and a Spearman correlation of 0.9724 on the validation set. The system was deployed as an interactive web application (Streamlit) that allows users to enter a SMILES string and obtain a similarity ranking against the reference compounds. This work continues PIA-01 by materializing the proposed functional prototype and demonstrating the feasibility of integrating artificial intelligence with cheminformatics to accelerate compound prioritization through molecular similarity.

**Palabras clave:** similitud molecular, red neuronal, benzimidazoles, Tanimoto, OpenEye ROCS, quimioinformática, CNN, SMILES, ranking molecular

### **Introducción**

La similitud molecular es un concepto central en la quimio-informática y el descubrimiento de fármacos. El principio fundamental establece que compuestos estructuralmente similares tienden a exhibir propiedades biológicas similares (Bajusz et al., 2015), lo cual sustenta el uso de métricas de similitud para priorizar compuestos candidatos en etapas tempranas de investigación. Sin embargo, la evaluación de similitud depende tradicionalmente de una sola métrica bidimensional, como el coeficiente de Tanimoto sobre fingerprints ECFP4, lo cual limita la captura de relaciones estructurales complejas.

En el Proyecto de Investigación Aplicada 1 (Aguilar Noguera y Araya Parajeles, 2025) se estableció un marco metodológico basado en similitud estructural mediante el coeficiente de Tanimoto y modelos supervisados como Random Forest, con aplicación a benzimidazoles ensayados en laboratorio. Se identificó como principal limitación la ausencia de resultados empíricos completos y la dependencia de una única métrica de similitud bidimensional.

El presente trabajo (PIA-02) aborda estas limitaciones mediante el desarrollo de una red neuronal capaz de predecir simultáneamente cuatro métricas de similitud molecular: Tanimoto 2D, Coseno, Shape 3D y Combo Tanimoto. A diferencia del enfoque anterior, este sistema incorpora información tridimensional generada por OpenEye ROCS (OpenEye Scientific Software, 2023), proporciona un ranking multi-métrica y se despliega como aplicación web interactiva.

Los benzimidazoles constituyen una familia química ampliamente estudiada en química medicinal (Oh et al., 2014; Mayence et al., 2011; Doğan et al., 2021). El sistema desarrollado utiliza 12 benzimidazoles ensayados en laboratorio como conjunto de referencia, permitiendo evaluar la similitud de cualquier compuesto candidato respecto a estas moléculas con propiedades verificadas experimentalmente.

Es importante destacar que este sistema evalúa exclusivamente similitud molecular estructural, no actividad biológica. Su propósito es servir como herramienta de priorización que reduzca el espacio de búsqueda química, aportando un ranking cuantitativo y reproducible basado en múltiples métricas complementarias.

## **Materiales y Métodos**

### **Conjunto de Datos**

Se utilizó un dataset consolidado de 20,095 compuestos obtenidos de la base de datos ChEMBL (Gaulton et al., 2017), cada uno con su representación SMILES (Weininger, 1988), siete descriptores moleculares (peso molecular, ALogP, donadores y aceptores de enlaces de hidrógeno, área de superficie polar topológica, enlaces rotables y anillos aromáticos) y métricas de similitud precalculadas.

Las métricas tridimensionales fueron generadas mediante el toolkit de OpenEye, para el cual se obtuvo una licencia académica gratuita otorgada por Cadence Molecular Sciences (OpenEye Scientific Software) en julio de 2025, a través del Prof. Rodrigo Herrera Garro, Director del Laboratorio de Inteligencia Artificial (AILAB) de la Universidad CENFOTEC. La licencia fue concedida en el marco de la investigación conjunta con el Departamento de Ciencias Farmacéuticas de la Universidad de Salamanca (USAL), con el compromiso de que los resultados sean compartidos al dominio público sin fines comerciales. Se utilizó OMEGA para la generación de conformeros 3D y ROCS (Rapid Overlay of Chemical Structures) para el cálculo de superposiciones tridimensionales, obteniendo Shape Tanimoto, Color Tanimoto y Combo Score. El 98.9% de los compuestos (19,880 de 20,095) generaron conformaciones válidas con OMEGA. Los 215 compuestos restantes fueron excluidos del entrenamiento. Es importante señalar que la licencia de OpenEye fue necesaria únicamente para la generación de los datos de entrenamiento;

el modelo entrenado opera de forma independiente sin requerir dicho software, lo cual permite su uso por investigadores sin acceso a licencias comerciales.

Como referencia se emplearon 12 benzimidazoles evaluados experimentalmente en laboratorio, proporcionados en colaboración con la Universidad de Salamanca.

### **Preparación de Datos**

Se aplicó un muestreo estratificado para balancear la representación por rangos de similitud Tanimoto 2D: alta ( $\geq 0.3$ ,  $n = 1,666$ ), media ( $0.15-0.3$ ,  $n = 1,666$ ) y baja ( $< 0.15$ ,  $n = 701$ ), resultando en 4,033 compuestos seleccionados. Para cada compuesto se generaron fingerprints Morgan (ECFP4, radio = 2, 2048 bits) mediante RDKit (Landrum, 2021). Un fingerprint molecular es una representación digital que codifica rasgos estructurales de una molécula y permite compararla computacionalmente con otras. Con estos fingerprints se calcularon las similitudes Tanimoto 2D y Coseno entre cada benzimidazol de referencia y cada compuesto de ChEMBL. El Combo Score de OpenEye (rango original 0-2) se normalizó al rango 0-1. Se produjeron 48,462 pares de entrenamiento (12 referencias  $\times$  4,033 compuestos más 66 pares inter-referencia).

### **Tokenización de SMILES**

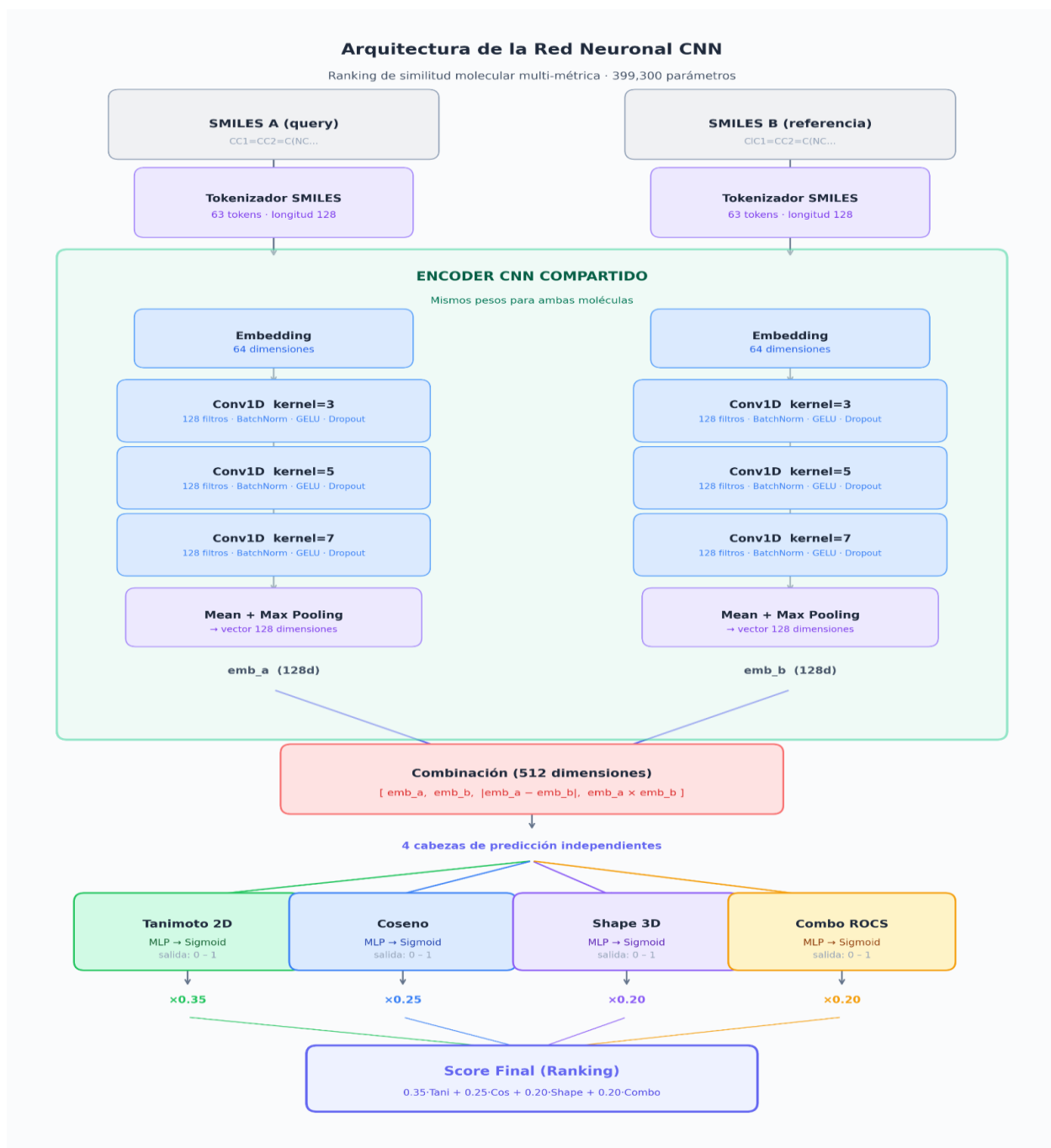
Para que las moléculas pudieran ser procesadas por la red neuronal, cada cadena SMILES fue transformada en una secuencia de índices numéricos mediante un tokenizador propio implementado para este proyecto. No se utilizó un tokenizador externo o preentrenado; la tokenización se definió mediante expresiones regulares basadas en la sintaxis de SMILES. Este enfoque fue escogido porque es simple, reproducible y suficiente para representar las estructuras químicas del conjunto de datos utilizado. El tokenizador reconoce unidades químicamente significativas como átomos de varios caracteres (Br, Cl), marcadores de estereoquímica (@, @@), átomos entre corchetes como [nH], símbolos de enlace, paréntesis y dígitos usados para representar cierres de anillos. Esto evita separar incorrectamente símbolos químicos relevantes y permite conservar mejor la información estructural de la molécula antes de ingresarla al modelo. El vocabulario se construyó a partir de los SMILES del conjunto de datos de ChEMBL y de los 12 benzimidazoles de referencia, obteniendo 63 tokens en total. Además, se incluyeron tokens especiales: PAD para completar secuencias a una longitud fija, UNK para símbolos no vistos, BOS para marcar el inicio de la secuencia y EOS para marcar el final. La longitud máxima se fijó en 128 tokens, cubriendo holgadamente el percentil 99 de las longitudes observadas, correspondiente a 45 tokens.

### **Arquitectura del Modelo**

Se diseñó una red neuronal (Bromley et al., 1993) con encoder CNN (red neuronal convolucional) compartido y cuatro cabezas de similitud independientes. La red utiliza el mismo encoder con pesos compartidos para procesar ambas moléculas del par, garantizando que la función de similitud sea simétrica:  $\text{sim}(A, B) = \text{sim}(B, A)$ . La Figura 1 presenta la arquitectura completa del modelo.

### **Figura 1**

## Arquitectura de la red neuronal CNN para ranking de similitud molecular multi-métrica



*Nota. La figura muestra cómo el modelo compara dos moléculas. Primero, cada molécula se representa mediante su cadena SMILES y esta se transforma en una secuencia numérica mediante el tokenizador. Luego, ambas secuencias pasan por el mismo encoder CNN, es decir,*

*por la misma red con los mismos pesos. Este encoder genera para cada molécula un vector numérico llamado embedding, que resume los patrones estructurales aprendidos por la red, como combinaciones de átomos, enlaces y subestructuras. Estos embeddings no provienen de una base externa, sino que son generados por el propio modelo durante el entrenamiento. Finalmente, los embeddings de ambas moléculas se comparan y se utilizan para predecir cuatro métricas de similitud: Tanimoto 2D, Coseno, Shape 3D y Combo ROCS.*

El encoder consta de una capa de embedding (64 dimensiones) seguida de tres bloques convolucionales 1D con kernels de tamaño 3, 5 y 7 respectivamente (128 filtros cada uno), normalización por lotes (BatchNorm), activación GELU y dropout (0.15). Las convoluciones 1D recorren la secuencia de tokens detectando patrones locales de diferentes tamaños, análogo a la detección de subestructuras químicas. La salida se agrega mediante mean pooling y max pooling enmascarados, concatenados y proyectados a un espacio de 128 dimensiones.

Para cada par de moléculas, los embeddings (emb\_a, emb\_b) se combinan en una representación de 512 dimensiones: [emb\_a, emb\_b, |emb\_a – emb\_b|, emb\_a × emb\_b]. Esta representación concatenada captura tanto la información individual de cada molécula como su diferencia absoluta e interacción elemento a elemento. Dicha representación alimenta cuatro cabezas de similitud independientes (perceptrón multicapa con activación Sigmoid en la capa de salida), cada una especializada en predecir una métrica: Tanimoto 2D, Coseno, Shape 3D y Combo ROCS. El modelo total tiene 399,300 parámetros entrenables.

### **Procedimiento de Entrenamiento**

El modelo se entrenó con una función de pérdida Huber (SmoothL1Loss) calculada por separado para cada una de las cuatro métricas de similitud y luego combinada en una suma ponderada. Esta función fue seleccionada porque el problema corresponde a una regresión de

valores continuos de similitud en escala 0–1, y Huber ofrece un balance entre precisión y robustez: penaliza adecuadamente errores pequeños, pero reduce el efecto de errores grandes o valores atípicos. Esto es útil en datos de similitud molecular, donde pueden existir variaciones asociadas a la generación de conformaciones 3D y a las diferencias entre métricas. En particular, la pérdida total se definió como  $1.5 \cdot \text{Loss\_tanimoto\_2d} + 1.2 \cdot \text{Loss\_coseno} + 0.8 \cdot \text{Loss\_shape} + 1.0 \cdot \text{Loss\_combo}$ , de modo que los errores en las métricas consideradas más confiables o prioritarias influyeran más en la optimización. Se utilizó el optimizador AdamW (Loshchilov y Hutter, 2019) con tasa de aprendizaje 0.001 y un coeficiente de regularización de  $1 \times 10^{-4}$ , implementado mediante el parámetro `weight_decay`, con el fin de penalizar pesos excesivamente grandes y favorecer una mejor generalización. Además, se aplicó `gradient clipping` (norma máxima 1.0) y `early stopping` con paciencia de 15 épocas. La partición de datos fue 90% para entrenamiento (43,616 pares) y 10% para validación (4,846 pares), con semilla fija para reproducibilidad.

### **Métricas de Evaluación**

Para evaluar el desempeño del modelo se utilizaron cinco métricas complementarias: error cuadrático medio (MSE), error absoluto medio (MAE), coeficiente de determinación ( $R^2$ ), correlación de Pearson ( $r$ ) y correlación de Spearman ( $\rho$ ). Estas métricas permiten medir tanto la precisión numérica de las predicciones como la capacidad del modelo para preservar el orden de similitud entre compuestos. El MSE calcula el promedio de los errores elevados al cuadrado y penaliza con mayor fuerza las predicciones con errores grandes. El MAE mide el error promedio en la misma escala de las métricas de similitud, por lo que permite interpretar directamente la magnitud del error. El  $R^2$  indica qué proporción de la variabilidad de los valores reales es explicada por el modelo, mientras que Pearson mide la relación lineal entre valores reales y predichos. Finalmente, Spearman evalúa si el modelo conserva el orden relativo de los compuestos. Esta

métrica es especialmente importante porque el objetivo del sistema es generar un ranking de similitud; por tanto, no basta con aproximar los valores numéricos, sino también ubicar en mejores posiciones a los compuestos más similares.

### **Implementación de la Aplicación Web**

Se desarrolló una interfaz web con Streamlit que permite ingresar uno o más compuestos en notación SMILES, visualizar la estructura molecular renderizada por RDKit, obtener el ranking de similitud contra los 12 benzimidazoles de referencia con las cuatro métricas simultáneamente, y exportar los resultados en formato CSV. El score final de ranking se calcula como combinación lineal ponderada:  $\text{Score} = 0.35 \cdot \text{Tanimoto\_2D} + 0.25 \cdot \text{Coseno} + 0.20 \cdot \text{Shape} + 0.20 \cdot \text{Combo}$ .

## **Resultados**

### **Convergencia del Entrenamiento**

El modelo alcanzó la mejor pérdida de validación (0.001256) en la época 19 de un máximo de 40. A partir de ese punto, aunque la pérdida de entrenamiento continuó disminuyendo, la pérdida de validación no logró valores menores que el obtenido en la época 19. Esto indica un ligero sobreajuste, ya que el modelo seguía ajustándose mejor a los datos de entrenamiento, pero

no mejoraba su capacidad de generalización sobre datos no usados para actualizar sus pesos. Para mitigarlo, se utilizaron varias estrategias de regularización: dropout de 0.15 en la arquitectura CNN, gradient clipping con norma máxima 1.0,  $\text{weight\_decay} = 1 \times 10^{-4}$  en AdamW y early stopping con paciencia de 15 épocas. Además, se conservó como modelo final el correspondiente a la mejor pérdida de validación (época 19), en lugar del último modelo entrenado.

**Tabla 1**

*Progresión de la pérdida durante el entrenamiento del modelo CNN*

Época	Train Loss	Val Loss	Tanimoto	Coseno	Shape	Combo
1	0.01017	0.00494	0.00063	0.00109	0.00195	0.00112
5	0.00220	0.00224	0.00056	0.00062	0.00046	0.00030
10	0.00123	0.00161	0.00042	0.00046	0.00035	0.00015
19*	—	0.00126	—	—	—	—
25	0.00078	0.00136	0.00032	0.00045	0.00029	0.00012
34	Early stop	—	—	—	—	—

*Nota.* \*Mejor modelo seleccionado. Elaboración propia.

## Métricas de Evaluación

El conjunto de evaluación correspondió al subconjunto de validación separado durante el entrenamiento. A partir de los 48,462 pares generados, se aplicó una partición aleatoria reproducible con semilla fija ( $\text{random\_state} = 42$ ), destinando el 90% de los pares al entrenamiento (43,616) y el 10% restante a validación/evaluación (4,846). Este subconjunto no fue utilizado para actualizar los pesos del modelo, sino para monitorear la generalización durante el entrenamiento, seleccionar el mejor modelo y calcular las métricas finales de desempeño. Cada par del conjunto de validación contiene dos moléculas en formato SMILES y sus cuatro valores reales de similitud: Tanimoto 2D, Coseno, Shape 3D y Combo ROCS. Por tanto, los resultados reportados corresponden a una evaluación interna sobre el conjunto de validación, no a una prueba externa con nuevos compuestos experimentales. Los resultados se presentan en la Tabla 2. El modelo

alcanzó un  $R^2$  promedio de 0.9353, explicando el 93.5% de la varianza en las predicciones de similitud. La correlación de Spearman promedio de 0.9724 indica que el ranking predicho preserva casi perfectamente el orden real.

**Tabla 2**

*Métricas de evaluación por tipo de similitud sobre el conjunto de validación*

Métrica	MSE	MAE	$R^2$	Pearson	Spearman
Tanimoto 2D	0.000556	0.0187	0.8972	0.9634	0.9605
Coseno	0.000772	0.0214	0.9234	0.9650	0.9609
Shape 3D	0.000583	0.0190	0.9481	0.9841	0.9840
Combo ROCS	0.000282	0.0128	0.9724	0.9879	0.9841
<i>Promedio</i>	<i>0.000548</i>	<i>0.0180</i>	<i>0.9353</i>	<i>0.9751</i>	<i>0.9724</i>

*Nota.* Elaboración propia. MSE = error cuadrático medio; MAE = error absoluto medio.

La métrica mejor predicha fue Combo ROCS ( $R^2 = 0.9724$ ,  $\rho = 0.9841$ ), seguida de Shape 3D ( $R^2 = 0.9481$ ). Tanimoto 2D presentó el  $R^2$  más bajo (0.8972), pero con un error absoluto medio de apenas 0.019, equivalente a menos del 2% de error promedio en la escala 0–1.

### **Validación Cualitativa por Inferencia**

Se realizaron pruebas de inferencia con cuatro compuestos representativos para verificar si el ranking producido por el modelo era coherente desde el punto de vista químico (Tabla 3). En este contexto, se esperaba que el compuesto de referencia (BZ-6) se reconociera a sí mismo como su coincidencia principal, que compuestos estructuralmente relacionados (benzimidazoles genéricos) obtuvieran scores intermedios-altos frente a referencias de la misma familia, y que un compuesto no relacionado (Aspirina) presentara un score menor. Por tanto, estos resultados no constituyen una validación estadística adicional, sino una comprobación cualitativa de que el sistema discrimina razonablemente entre compuestos relacionados y no relacionados.

**Tabla 3**

*Score final de similitud para compuestos de prueba (top-1 de cada query)*

Query	Top-1	Score	Tanimoto 2D	Shape 3D
BZ-6 (referencia)	BZ-6	0.6247	0.7074	0.3672
2-fenilbenzimidazol	141	0.5099	0.4067	0.6445
Benzimidazol base	21	0.4531	0.3360	0.5812
Aspirina (control neg.)	141	0.3601	0.2768	0.4608

*Nota.* El score final es la combinación ponderada de las cuatro métricas. Elaboración propia.

El modelo reconoció correctamente al compuesto BZ-6 como el más similar a sí mismo (score 0.6247), con scores descendentes para compuestos estructuralmente relacionados. La Aspirina, utilizada como control negativo, obtuvo el score más bajo (0.3601), demostrando una separación clara entre compuestos benzimidazólicos y no relacionados.

### Discusión

Los resultados demuestran que una red neuronal CNN relativamente compacta (399,300 parámetros) puede aprender a predecir múltiples métricas de similitud molecular simultáneamente con alta precisión. El  $R^2$  promedio de 0.9353 y la correlación de Spearman de 0.9724 indican que el modelo captura tanto la magnitud absoluta como el orden relativo de las similitudes, aspecto crítico para la generación de rankings confiables.

Es notable que las métricas tridimensionales (Shape y Combo) fueron mejor predichas que las bidimensionales (Tanimoto 2D). Esto puede explicarse porque la similitud de forma es una propiedad más global de la molécula y presenta distribuciones más suaves, mientras que los fingerprints binarios pueden producir cambios discontinuos ante pequeñas modificaciones estructurales, fenómeno conocido como activity cliffs (Stumpfe y Bajorath, 2012).

La arquitectura con encoder compartido garantiza la simetría de la función de similitud:  $\text{sim}(A, B) = \text{sim}(B, A)$ , propiedad matemáticamente deseable que no se garantiza en todos los

enfoques de aprendizaje profundo para similitud molecular. Las cuatro cabezas independientes permiten que cada métrica capture aspectos diferentes de la relación estructura-similitud sin interferencia mutua, funcionando como un sistema de aprendizaje multi-tarea (multi-task learning).

Respecto al PIA-01 (Aguilar Noguera y Araya Parajeles, 2025), este trabajo materializa tres de las recomendaciones planteadas: (a) la incorporación de información tridimensional mediante datos reales de OpenEye ROCS, (b) el despliegue como herramienta web interactiva, y (c) la obtención de resultados empíricos cuantitativos con métricas formales de evaluación. El enfoque evolucionó de Random Forest con Tanimoto unidimensional a una red neuronal multi-métrica que integra información 2D y 3D.

Una limitación del presente trabajo es que la similitud molecular no implica actividad biológica directa. El sistema está diseñado explícitamente para ranking de similitud estructural. Adicionalmente, el entrenamiento se realizó en CPU debido a restricciones de infraestructura, lo cual limitó la exploración de arquitecturas más profundas como Transformers. No obstante, los resultados obtenidos con la arquitectura CNN son competitivos y demuestran que modelos compactos pueden ser altamente efectivos para esta tarea.

## Conclusiones

Se diseñó, entrenó y evaluó exitosamente una red neuronal CNN para la predicción multi-métrica de similitud molecular, alcanzando un  $R^2$  promedio de 0.9353 y Spearman de 0.9724 sobre datos de validación. El sistema integra información bidimensional (fingerprints ECFP4) y tridimensional (OpenEye ROCS) en un marco unificado con cuatro cabezas de predicción independientes.

El prototipo funcional desplegado como aplicación web permite a investigadores ingresar compuestos candidatos y obtener rankings de similitud contra benzimidazoles de referencia con propiedades verificadas experimentalmente, reduciendo significativamente el tiempo de priorización frente a la búsqueda manual.

Este trabajo demuestra la viabilidad de combinar inteligencia artificial con quimioinformática para acelerar la evaluación de similitud molecular. Como trabajo futuro se recomienda la incorporación de predicciones farmacocinéticas, la validación experimental de los compuestos mejor rankeados, y la extensión del sistema a otras familias químicas.

### Referencias

- ◆ Aguilar Noguera, H. P., y Araya Parajeles, J. L. (2025). Modelos de inteligencia artificial para predecir la actividad biológica de compuestos químicos contra *Leishmania donovani* [Proyecto de Investigación Aplicada 1]. Universidad CENFOTEC.
- ◆ Bajusz, D., Rácz, A., y Héberger, K. (2015). Why is Tanimoto index an appropriate choice for fingerprint-based similarity calculations? *Journal of Cheminformatics*, 7(1), 20. <https://doi.org/10.1186/s13321-015-0069-3>
- ◆ Bromley, J., Guyon, I., LeCun, Y., Säckinger, E., y Shah, R. (1993). Signature verification using a time delay neural network. *Advances in Neural Information Processing Systems*, 6, 737–744.

- ◆ Brown, N. (1998). Chemoinformatics: A new name for an old problem. *Journal of Chemical Information and Computer Sciences*, 38(6), 973–975. <https://doi.org/10.1021/ci980302o>
- ◆ Doğan, B., Ülgen, K. Ö., y Yelekçi, K. (2021). Benzimidazole derivatives as potential antileishmanial agents: Design, synthesis and biological evaluation. *European Journal of Medicinal Chemistry*, 224, 113712. <https://doi.org/10.1016/j.ejmech.2021.113712>
- ◆ Gaulton, A., Hersey, A., Nowotka, M., Bento, A. P., Chambers, J., Mendez, D., ... y Leach, A. R. (2017). The ChEMBL database in 2017. *Nucleic Acids Research*, 45(D1), D945–D954. <https://doi.org/10.1093/nar/gkw1074>
- ◆ Gómez-Bombarelli, R., Wei, J. N., Duvenaud, D., Hernández-Lobato, J. M., Sánchez-Lengeling, B., Sheberla, D., ... y Aspuru-Guzik, A. (2018). Automatic chemical design using a data-driven continuous representation of molecules. *ACS Central Science*, 4(2), 268–276. <https://doi.org/10.1021/acscentsci.7b00572>
- ◆ Herwaldt, B. L. (1999). Leishmaniasis. *The Lancet*, 354(9185), 1191–1199. [https://doi.org/10.1016/S0140-6736\(98\)10178-4](https://doi.org/10.1016/S0140-6736(98)10178-4)
- ◆ Kingma, D. P., y Welling, M. (2019). An introduction to variational autoencoders. *Foundations and Trends in Machine Learning*, 12(4), 307–392. <https://doi.org/10.1561/22000000056>
- ◆ Landrum, G. (2021). RDKit: Open-source cheminformatics software. <https://www.rdkit.org/>
- ◆ Loshchilov, I., y Hutter, F. (2019). Decoupled weight decay regularization. *Proceedings of the International Conference on Learning Representations (ICLR)*. <https://arxiv.org/abs/1711.05101>

- ◆ Mayence, A., Huang, J., y Leung, P. P. (2011). Synthesis and antileishmanial activity of benzimidazole derivatives. *Bioorganic & Medicinal Chemistry Letters*, 21(18), 5447–5450. <https://doi.org/10.1016/j.bmcl.2011.07.052>
- ◆ Oh, H. J., Kim, J. H., y Lee, S. K. (2014). Antileishmanial activity of novel benzimidazole derivatives. *Parasitology International*, 63(5), 723–728. <https://doi.org/10.1016/j.parint.2014.05.003>
- ◆ OpenEye Scientific Software. (2023). ROCS—Rapid Overlay of Chemical Structures. Cadence Molecular Sciences. <https://www.eyesopen.com/rocs>
- ◆ Organización Mundial de la Salud. (2023). Enfermedades tropicales desatendidas: Leishmaniasis. <https://www.who.int/es/news-room/fact-sheets/detail/leishmaniasis>
- ◆ Stumpfe, D., y Bajorath, J. (2012). Exploring activity cliffs in medicinal chemistry. *Journal of Medicinal Chemistry*, 55(7), 2932–2942. <https://doi.org/10.1021/jm201706b>
- ◆ Weininger, D. (1988). SMILES, a chemical language and information system. *Journal of Chemical Information and Computer Sciences*, 28(1), 31–36. <https://doi.org/10.1021/ci00057a005>