



# Augmenting Generative Models With Biomedical Knowledge Graphs Improves Targeted Drug Discovery

Aditya Malusare , Vineet Punyamoorthy , and Vaneet Aggarwal , *Senior Member, IEEE*

**Abstract**—Recent breakthroughs in generative modeling have demonstrated remarkable capabilities in molecular generation, yet the integration of comprehensive biomedical knowledge into these models has remained an untapped frontier. In this study, we introduce knowledge-driven embedding-augmented model (K-DREAM), a novel framework that leverages knowledge graphs to augment diffusion-based generative models for drug discovery. By embedding structured information from large-scale knowledge graphs, K-DREAM directs molecular generation toward candidates with higher biological relevance and therapeutic suitability. This integration ensures that the generated molecules are aligned with specific therapeutic targets, moving beyond traditional heuristic-driven approaches. In targeted drug design tasks, K-DREAM generates drug candidates with improved binding affinities and predicted efficacy, surpassing current state-of-the-art generative models. It also demonstrates flexibility by producing molecules designed for multiple targets, enabling applications to complex disease mechanisms. These results highlight the utility of knowledge-enhanced generative models in rational drug design and their relevance to practical therapeutic development.

**Impact Statement**—We introduce K-DREAM, a new approach to drug discovery that combines knowledge graphs with AI-driven drug design. Unlike conventional methods that focus mainly on chemical properties, our framework incorporates biological relationships to create more medically relevant drug candidates. When tested on various protein targets, our system produced potential drugs with better binding abilities while maintaining diversity and feasibility for manufacturing. Particularly promising is K-DREAM’s ability to design drugs that work on multiple targets simultaneously, a crucial advantage for treating complex diseases that involve multiple biological

pathways. By bringing biological context directly into the drug design process, our work could significantly reduce the time and resources needed to discover promising new treatments, potentially accelerating the journey from laboratory to patient.

**Index Terms**—Diffusion models, drug discovery, generative AI, knowledge representation.

## I. INTRODUCTION

GENERATIVE models have achieved significant success in fields such as text and image synthesis, establishing state-of-the-art performance in natural language processing, computer vision, and decision making [3], [4]. These successes have inspired the extension of generative models to other modalities, including graph-based representations of molecular structures, which hold considerable potential in drug discovery. Molecular generative models, in particular, have emerged as powerful tools for de novo drug design and optimization, leveraging graph-based representations to navigate the complex chemical space. Techniques such as variational autoencoders (VAEs) [5], generative adversarial networks (GANs), autoregressive models [6], and reinforcement learning [7], [8], [9] have been employed to explore molecular configurations and propose novel therapeutic candidates.

Despite their promise, these models are often constrained by their reliance on simplistic guidance mechanisms, typically involving heuristic chemical scores, such as synthetic accessibility (SA) [10] and quantitative estimate of drug-likeness (QED) [11], along with structural constraints such as validity and uniqueness. These conventional approaches largely overlook the extensive body of biomedical knowledge accumulated through years of research and clinical practice, which could significantly enhance the generative process.

Biomedical knowledge is systematically encoded into extensive knowledge graphs [1], [12], which span millions of relations and provide structured representations of interactions within biological systems. These graphs capture relationships across multiple biological scales, from individual molecular components, such as genes, proteins, and small molecules, to higher-order entities such as cells, tissues, and entire biological processes. The relationships in these graphs reflect diverse biological phenomena, including protein–protein interactions, gene regulation, metabolic pathways, and drug–target associations.

Received 25 March 2025; revised 23 July 2025 and 1 October 2025; accepted 10 October 2025. Date of publication 14 October 2025; date of current version 30 April 2026. This work was supported in part by the National Science Foundation under Grant FW-HTF-R-2129097. The work of Aditya Malusare was supported in part by Walther Cancer Foundation and in part by Purdue University Institute for Cancer Research under Grant P30CA023168. This article was recommended for publication by Associate Editor Dwarikanath Mahapatra upon evaluation of the reviewers’ comments. (Corresponding author: Aditya Malusare.)

Aditya Malusare and Vaneet Aggarwal are with Edwardson School of Industrial Engineering and the Institute of Cancer Research, Purdue University, West Lafayette, IN 47907 USA (e-mail: malusare@purdue.edu; vaneet@purdue.edu).

Vineet Punyamoorthy is with Elmore Family School of Electrical and Computer Engineering, Purdue University, West Lafayette, IN 47907 USA (e-mail: vpunyamo@purdue.edu).

This article has supplementary downloadable material available at <https://doi.org/10.1109/TAI.2025.3621169>, provided by the authors.

Digital Object Identifier 10.1109/TAI.2025.3621169

The structured format of biomedical knowledge graphs captures complex biological behaviors that arise from interactions between molecules. For example, cellular homeostasis (the ability of cells to maintain stable internal conditions) is driven by feedback loops within signaling pathways [13]. Similarly, phenotypic robustness, which refers to an organism’s ability to maintain a stable phenotype despite genetic changes or environmental fluctuations, often emerges from redundant genetic networks [14]. A major challenge in drug development is drug resistance, which can occur through adaptive mechanisms such as the overactivation of efflux pumps (proteins that expel drugs from cells) or mutations in drug target proteins [15]. Biomedical knowledge graphs also support drug repurposing by revealing hidden relationships between diseases and existing drugs [16], as demonstrated by the identification of baricitinib for COVID-19 treatment by linking its known mechanisms with pathways involved in viral infectivity [17]. By capturing these diverse and interconnected biological relationships, knowledge graphs provide a rich source of information for generative models to design drug candidates that are biologically meaningful. Existing work on integrating knowledge graphs primarily focuses on enhancing molecular property prediction tasks, such as predicting interactions [18], repurposing [19], and molecular property optimization [20]. However, the potential of knowledge graphs to guide the generative process itself remains largely unexplored.

In this article, we introduce a novel framework called knowledge-driven embedding-augmented model (K-DREAM), which bridges the gap between state-of-the-art generative models and the wealth of information contained within biomedical knowledge graphs (Fig. 1). Our approach aims to generate biologically relevant and therapeutically promising molecules by leveraging biomedical knowledge graphs to inform and guide molecular generation. Specifically, we integrate the embeddings derived from knowledge graphs, preserving the semantic integrity of the biological information, into a generative framework. By combining the power of diffusion-based generative modeling with knowledge graph embeddings, K-DREAM produces novel drug candidates that are not only chemically valid but also enriched with specific therapeutic properties, enhancing their potential in drug development. Prior work on applied generative models in this field uses conditional diffusion-based architectures to generate molecular graphs to maximize previously mentioned heuristic scores such as QED and SA, with techniques to adapt to optimizing for these nondifferentiable metrics. We demonstrate that the embeddings derived from biomedical knowledge graphs contain sufficient information to steer the generative process towards biologically relevant molecules, thereby improving the quality and relevance of the generated compounds without explicitly including quantitative scores or structural constraints.

To represent the biomedical knowledge in a form compatible with generative models, we use knowledge graph embedding (KGE) techniques, which transform entities and relationships in the graph into a continuous vector space while preserving semantic relationships. In this work, we employ the TransE model [2] due to its interpretability and efficiency. TransE represents relations as translations between entities in a

real-valued embedding space, making it suitable for conceptualizing molecular generation as a path through a unified biological embedding space. This additive nature of TransE allows us to reason intuitively about entity relationships and integrate these insights into the generative modeling process.

K-DREAM outperforms state-of-the-art molecular generative models on targeted drug design tasks, producing compounds with improved biological relevance and therapeutic potential. In docking studies against specific protein targets, K-DREAM-generated molecules consistently achieved higher docking scores compared to other models, indicating greater binding likelihood to intended targets. The framework’s adaptability enables diverse generative tasks, including multitarget drug design for compounds with tailored polypharmacological profiles. We demonstrate this by generating molecules designed to interact with multiple protein targets simultaneously, exploring complex drug design scenarios potentially addressing diseases with multiple underlying mechanisms.

## II. METHODS

K-DREAM systematically bridges molecular generation with biomedical knowledge through four key components. We begin with foundational representations of both molecular structures and biomedical knowledge, establishing the mathematical framework for our approach. The molecular generation pipeline then builds from an unconditional base model to a knowledge-guided framework, incorporating information from biomedical relationships to steer the generative process. We develop specialized neural architectures to enable precise navigation of the resulting chemical-knowledge space, allowing targeted molecular design. Finally, we conclude with techniques for exploration and visualization of this unified space, followed by details of our experimental setup for reproducibility.

### A. Molecular Representation

Molecular structures are represented in this work using a planar graph  $\mathbf{G} = (\mathbf{X}, \mathbf{E})$  where  $\mathbf{X} \in \mathbb{R}^{N \times M}$  is a feature matrix for  $N$  nodes (heavy atoms) described by  $M$ -dimensional vectors encoding atom information, and  $\mathbf{E} \in \mathbb{R}^{N \times N}$  is the adjacency matrix indicating the presence of single, double or triple bonds between the nodes.

### B. Knowledge Graph Embeddings

We generate the knowledge graph embeddings (KGEs) for the PrimeKG dataset using TransE, which represents entities and relationships as embeddings with a linear translation relationship between them. Since PrimeKG includes a reversed triple  $(\mathbf{o}, \mathbf{r}, \mathbf{s})$  for each original triple  $(\mathbf{s}, \mathbf{r}, \mathbf{o})$ , we remove the former to preserve the directed nature of the graph for training TransE.

We utilize the implementations of the PrimeKG dataset and TransE model from PyKEEN [21]. The model is trained for 100 epochs with a learning rate of  $10^{-3}$ . Training a KGE model often involves generating negative triples to prevent under-fitting. One method for doing that is by assuming that every possible triple not present in the knowledge graph is incorrect (closed

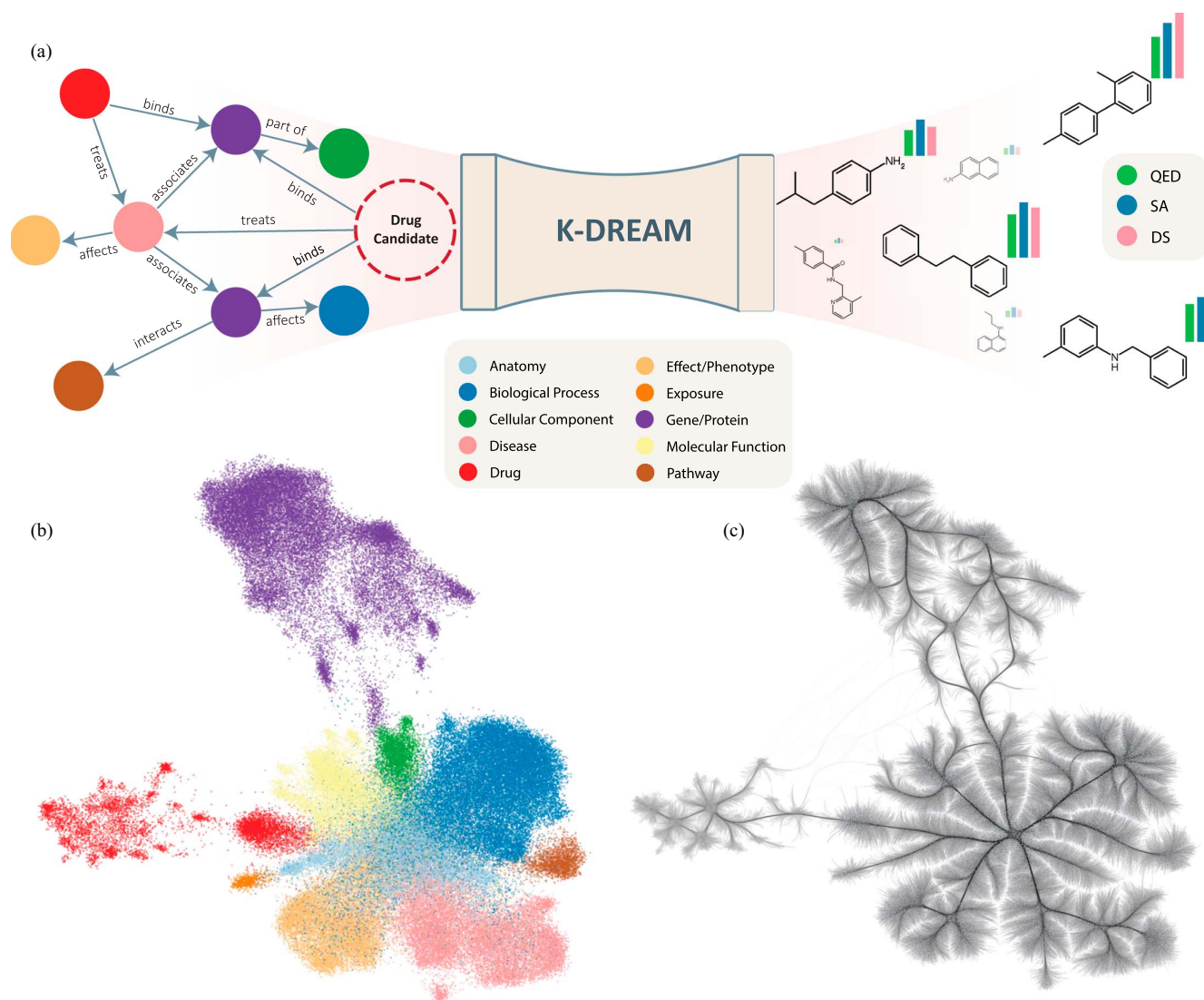


Fig. 1. *Knowledge-driven embedding-augmented model (K-DREAM)*. (a) Overview of the K-DREAM generative model for molecular structures, incorporating guidance information from embeddings derived from the knowledge graph. The diffusion process is guided to produce molecules that are both chemically valid and biologically relevant to a given target embedding. Molecules are evaluated using metrics such as docking score, quantitative estimate of drug-likeness (QED), and synthetic accessibility (SA). (b) Embeddings of the PrimeKG [1] knowledge graph created using the TransE [2] model, projected using the uniform manifold approximation projection (UMAP) algorithm. (c) Visualization of the relationships between UMAP-projected PrimeKG embeddings using edge bundling (a technique to reduce visual clutter in network visualizations by grouping edges that follow similar paths). Bundled edges represent relationships between entities, with thicker bundles indicating stronger or more numerous connections between related clusters.

world assumption). However, as knowledge graphs are inherently incomplete, this often leads to over-fitting. Therefore, we use the stochastic local closed world assumption (sLCWA) for training, which generates a random subset of all possible negative triples  $(s, r, o') \notin \text{PrimeKG}$  for each  $(s, r, o) \in \text{PrimeKG}$ . This method strikes a balance between under-fitting and the closed world approximation, while the stochasticity helps keep computational complexity lower than full LCWA.

### C. Unconditional Generative Model

Score-based models define the forward diffusion  $q$  of a graph  $\mathbf{G}_t = (\mathbf{X}_t, \mathbf{E}_t)$  with a stochastic differential equation (SDE)

$$d\mathbf{G}_t = \mathbf{f}_t(\mathbf{G}_t) + g_t d\mathbf{w} \quad (1)$$

where  $\mathbf{w}$  is the standard Wiener process and  $\mathbf{f}_t$  and  $g_t$  are the coefficients of linear drift and scalar diffusion, respectively.

The generative process is initiated with a graph  $\mathbf{G}_0$  and iteratively diffused to generate  $\mathbf{G}_T$ .

The stochastic forward process described in (1) can be used for generation by solving its reverse-time version

$$d\mathbf{G}_t = [\mathbf{f}_t(\mathbf{G}_t) - g_t^2 \nabla_{\mathbf{G}_t} \log p_t(\mathbf{G}_t)] d\bar{t} + g_t d\bar{\mathbf{w}} \quad (2)$$

where  $\bar{t}$  and  $\bar{\mathbf{w}}$  represent a reverse time-step and stochastic process. A score network  $s_\theta$  is used to approximate  $\nabla_{\mathbf{G}_t} \log p_t(\mathbf{G}_t)$ , where  $p_t(\mathbf{G}_t)$  is the probability density of the diffused graph at time  $t$ , and simulate the reverse process in time to generate  $\mathbf{G}_{t-1}$ .

The parameters  $\theta$  are obtained using an adaptation of the loss function by Song et al. [22]

$$\theta^* = \arg \min_{\theta} \mathbb{E}_t \left\{ \lambda(t) \mathbb{E}_{\mathbf{G}_0} \mathbb{E}_{\mathbf{G}_t | \mathbf{G}_0} \left[ \|s_\theta(\mathbf{G}_t, t) - \nabla_{\mathbf{G}_t} \log p_t(\mathbf{G}_t | \mathbf{G}_0)\|_2^2 \right] \right\} \quad (3)$$

Our unconditioned model is based on the GDSS architecture [23], trained to sample compounds from the ZINC250k dataset [24].

#### D. Conditional Generative Model

Guiding the stochastic process can be achieved by adding the conditioning information  $\mathbf{y}$  at each diffusion step, converting  $p_t(\mathbf{G}_t)$  from (2) to  $p_t(\mathbf{G}_t|\mathbf{y})$ . The score network is then used to approximate the modified gradient  $\nabla_{\mathbf{G}_t} \log p_t(\mathbf{G}_t|\mathbf{y})$ . The conditional distribution is rearranged to give

$$\nabla_{\mathbf{G}_t} \log p_t(\mathbf{G}_t|\mathbf{y}) = \nabla_{\mathbf{G}_t} \log p_t(\mathbf{G}_t) + \nabla_{\mathbf{G}_t} \log p_t(\mathbf{y}|\mathbf{G}_t). \quad (4)$$

The term  $\nabla_{\mathbf{G}_t} \log p_t(\mathbf{y}|\mathbf{G}_t)$  in (4) steers the model towards optimizing for the condition, while the first term  $\nabla_{\mathbf{G}_t} \log p_t(\mathbf{G}_t)$  introduces variation into the trajectory and helps explore newer regions.

In order to estimate the conditional term above, we consider two noising processes  $q$  and  $\hat{q}$  for the unconditioned and conditioned cases, respectively.

Using a result by Dhariwal et al. [25], the conditional term can be written as follows:

$$\hat{q}(\mathbf{G}_{t-1}|\mathbf{G}_t, \mathbf{y}) \propto q(\mathbf{G}_{t-1}|\mathbf{G}_t)\hat{q}(\mathbf{y}|\mathbf{G}_{t-1}). \quad (5)$$

Using a first-order approximation for the reverse process, we obtain

$$\log \hat{q}(\mathbf{y}|\mathbf{G}_{t-1}) \approx \log \hat{q}(\mathbf{y}|\mathbf{G}_t) + \langle \nabla_{\mathbf{G}_t} \log \hat{q}(\mathbf{y}|\mathbf{G}_t), \mathbf{G}_{t-1} - \mathbf{G}_t \rangle. \quad (6)$$

We now train a regressor network  $P_\phi$  to estimate  $\mathbf{y}$  from a noised version of the graph  $\mathbf{G}_t$  and use the gradients of this network to guide the generative process. We make the assumption that the conditioned process is effectively a distribution centered at the estimated knowledge-based embedding  $P_\phi(\mathbf{G}_t)$ , i.e.,

$$\hat{q}(\mathbf{y}|\mathbf{G}_t) \sim \mathcal{N}(P_\phi(\mathbf{G}_t), \sigma_y^2 \mathbf{I}). \quad (7)$$

This leads to the overall approximation in (6) to be

$$\log \hat{q}(\mathbf{y}|\mathbf{G}_{t-1}) \approx \left\langle \nabla_{\mathbf{G}_t} \|\mathbf{y} - P_\phi(\mathbf{G}_t)\|_2^2, \mathbf{G}_{t-1} \right\rangle + \text{terms independent of } \mathbf{G}_{t-1}. \quad (8)$$

The implementation of the score-based network is based on MOOD [26]. The details of the regressor network  $P_\phi$  are described below.

#### E. Context Regressor Network

To guide the conditional generation process, we create a trainable network to estimate knowledge-based embeddings  $\mathbf{y}$  from a noised version of an input molecular graph  $\mathbf{G}_T$ .  $P_\phi(\mathbf{G}_T) \approx \mathbf{y}$  is used to implement a modified version of the classifier guidance algorithm by Sohl-Dickstein et al. [27]. While previous works by Lee et al. [26] and Vignac et al. [28] use a similar algorithm to guide conditional generative processes, ours is a different approach that utilizes a stack

of graph attention layers to estimate knowledge-based embeddings, effectively creating a map between chemical space and KGE space

$P_\phi(\mathbf{G}) = P_\phi(\mathbf{X}, \mathbf{E})$  is constructed by first passing the feature  $\mathbf{X}$  and adjacency matrices  $\mathbf{E}$  through an aggregation operation:  $H^1 = \sigma(\mathbf{E}\mathbf{X}\mathbf{W}_\phi^0)$

We then use a stack of self-attention layers

$$h_i^{l+1} = \sum_j \alpha_{ij} \cdot \mathbf{W}^l h_j^l \quad (9)$$

where

$$\alpha_{ij} = \text{SoftMax}(\sigma(\mathbf{W}_a^l \cdot [\mathbf{W}^l h_i^l \oplus \mathbf{W}^l h_j^l])). \quad (10)$$

Here,  $\mathbf{W}^l, \mathbf{W}_a^l$  are learnable parameters at the  $l$ th layer. The stack of attention layers produces a final output of dimension  $|\mathbf{y}|$ .

#### F. Interpolation in Chemical Space

Upon obtaining the space of embeddings for the knowledge graph entities, we devise an interpolation process between the embeddings of two entities to generate new embeddings, akin to generating new molecules by interpolating between two known molecules. This process is well explored in image-based inputs [29], and also in the molecular latent space generated by graph-based autoencoders [30]. Our work extends this to the knowledge graph space, where we interpolate between the embeddings of two entities to generate new embeddings and use them to in the generative process of molecules.

Given a set of target genes  $\mathbf{t}_i$ , we can devise a process to generate a drug that maximizes binding to all the targets. We can formulate this as a regression problem, where we aim to predict the drug embedding  $\mathbf{y}$  that minimizes the distance between the sum of the drug and target gene embeddings

$$\mathbf{y} = \arg \min_{\mathbf{d} \in \mathcal{D}} \sum_i \|\mathbf{d} + \mathbf{r}_i - \mathbf{t}_i\|_2^2. \quad (11)$$

Here,  $r$  can potentially be the embedding for the drug/protein relation, and  $\mathbf{t}_i$  is the embedding for the target gene  $i$ , and  $\mathcal{D}$  is the set of all possible linear combinations of known drug embeddings.

### III. RESULTS

K-DREAM builds upon recent developments in diffusion models by incorporating a regression-based guidance mechanism informed by biomedical knowledge graphs. The model employs a high-dimensional regression approach that maps molecular structures to a latent space encoded with biomedical information derived from the graph. We extend the typical formulation of classifier-based diffusion models that often rely on simpler scalar metrics [26] or discrete classifications for guidance. The regression network creates a mapping between chemical structures and knowledge-derived embedding space, providing a means to guide the diffusion process based on biological context. By operating in a continuous, high-dimensional latent space, K-DREAM allows for nuanced control over the generated molecules' properties.

We trained the knowledge graph embedding (KGE) model, TransE, on the PrimeKG [1] dataset to generate knowledge graph embeddings. The dataset was preprocessed by removing duplicate entries and reversed triples to preserve the directed nature of the graph. Using the stochastic local closed world assumption (sLCWA) to generate negative triples, we balanced between underfitting and computational efficiency. The resulting embeddings were visualized using uniform manifold approximation and projection (UMAP) [31], revealing clear separation between different entity classes in the projected space. This separation demonstrates the model's ability to capture meaningful semantic relationships between entities, with clusters forming around related biological concepts such as genes, proteins, diseases, and drug compounds. We use edge bundling [32] to visualize relations between the embeddings, enabling us to see the connections between entities more clearly, highlighting how certain entities are interlinked through various biological processes. They signify important biological interactions, such as drug-target associations, gene regulation pathways, and protein-protein interactions.

To bridge the gap between chemical structures and knowledge graph embeddings, we developed a neural network with graph attention layers called the context regressor network (CRN). Constructed using a stack of graph attention layers, this network maps molecular graphs to their corresponding knowledge-based embeddings. Crucially, the training process involved augmenting the input data by creating noised versions of each chemical structure. These noised variants were trained to map to the same embedding as the original molecule. This approach ensures that structurally similar molecules in chemical space—those within a “neighborhood” of a given chemical structure—are mapped to the same point in the embedding space. This design choice preserves local chemical relationships while allowing for improved guidance in the generative process.

The underlying unconditioned generative model is a score-based model that samples the molecular distribution of the ZINC250k dataset, a collection of 250 k commercially available chemical compounds [24]. The objective of the guidance mechanism is to minimize the Euclidean distance between the embedding of the molecule at the current diffusion step and the target embedding. This is modeled as a multidimensional Gaussian function centered at the target embedding. The contribution of the guidance term to the score is controlled by a hyperparameter  $\lambda_X$ , which determines the weight of the guidance term relative to the underlying score-based model. In our figures, we display the negated docking scores for clarity, as lower (more negative) scores indicate stronger binding affinity.

#### A. K-DREAM Achieves State-of-the-Art Docking Scores in Targeted Drug Design Tasks

The evaluation of targeted drug design is conducted on five protein targets selected from the DUD-E database [34]: PARP1, JAK2, FA7, 5HT1B, and BRAF. These proteins, chosen for their high AUROC in docking score prediction using QuickVina 2 [35], represent critical targets in various therapeutic areas. PARP1, a DNA repair enzyme, is implicated in cancer therapy

resistance. JAK2, a tyrosine kinase, plays a crucial role in hematopoiesis and is a target for myeloproliferative disorders. FA7, or coagulation factor VII, is essential in the blood clotting cascade and a potential target for anticoagulant therapies. 5HT1B, a serotonin receptor subtype, is involved in mood regulation and migraine pathophysiology. BRAF, a serine/threonine kinase, is frequently mutated in melanoma and other cancers, making it a prime target for oncology drug development.

Docking scores are computed using QuickVina 2 through molecular docking simulations to provide a quantitative measure of the binding affinity between a ligand and a protein target. These scores predict the likelihood of a compound's interaction with its intended target. Lower (more negative) docking scores indicate stronger predicted binding, suggesting potentially higher efficacy of the drug candidate.

We compare K-DREAM against several categories of molecular generative models. Variational autoencoder-based methods include HierVAE [5], which generates molecular graphs using structural motifs as building blocks through hierarchical encoding from atoms to connected motifs. Diffusion-based approaches include GDSS [23], which models joint node-edge distributions via stochastic differential equations, and MOOD [26], a score-based diffusion model that incorporates out-of-distribution control and uses property prediction network gradients to guide generation toward desired chemical properties. Reinforcement learning methods include REINVENT [8], MORLD [7], FREED [9], and FREED-QS [9], which optimize molecular properties through reward-based learning in chemical space. We also compare against KARL [33], another knowledge-augmented approach, and evaluate K-DREAM with zero guidance to demonstrate the contribution of knowledge graph embeddings.

K-DREAM generated 3000 molecules for each protein target. These molecules were subsequently ranked based on their docking scores, with the top 5% selected for analysis. The mean and standard deviation of docking scores for this subset were calculated for each protein target (Fig. 2). Following previous work, we applied filters based on quantitative estimate of drug-likeness (QED) [11], synthetic accessibility (SA) [10] score, and Tanimoto similarity [36] to ensure the drug-likeness and synthetic feasibility of the generated molecules. The QED threshold of 0.5 ensures favorable physicochemical properties, while the SA score cutoff of 0.44 promotes synthetic feasibility. A maximum Tanimoto similarity of 0.4 to known actives maintains novelty while preserving desirable structural features. These criteria collectively ensure that K-DREAM generates not only potent but also practically viable novel drug candidates. The mean validity and uniqueness of the generated molecules is greater than 99.8%.

K-DREAM achieved mean docking scores of  $-12.13$  kcal/mol for PARP1,  $-10.39$  kcal/mol for JAK2,  $-9.08$  kcal/mol for FA7,  $-11.48$  kcal/mol for 5HT1B, and  $-11.41$  kcal/mol for BRAF. We compare against baselines that use reinforcement learning [7], [8], [9], variational autoencoders [5], and score-based models [26], [37]. K-DREAM outperforms these models across all protein targets, demonstrating its superior ability to generate molecules

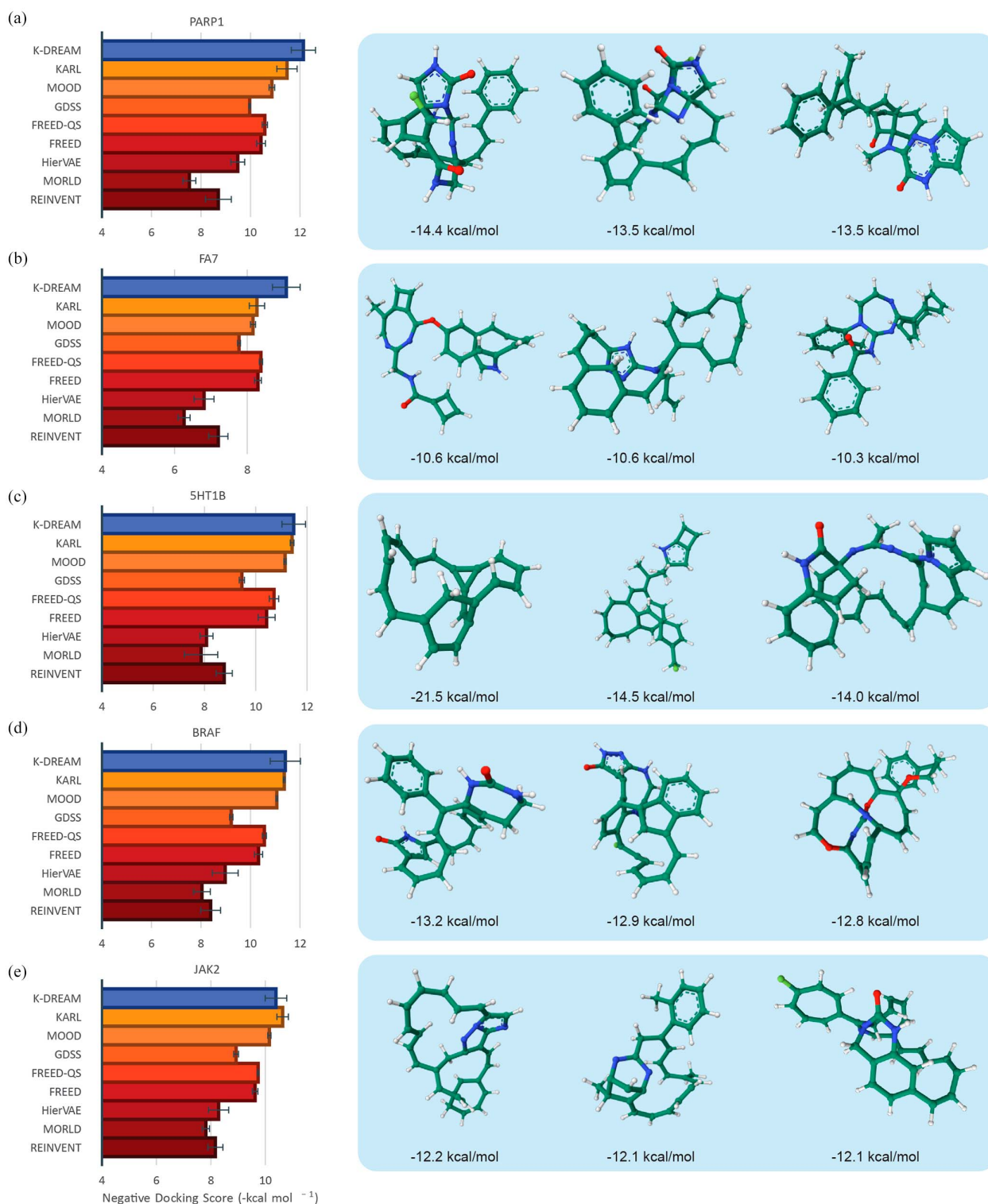


Fig. 2. *Docking scores.* The mean docking scores of the top 5% molecules generated by K-DREAM are compared against baseline models for five protein targets (left). For each protein, we also show the top 3 molecules generated by K-DREAM, along with their docking scores (right).

with high predicted binding affinity. We note the improved performance of K-DREAM relative to the KARL model, which also uses contextual information from knowledge graphs, with an explicit reward function for heuristic scores. The design of the context regressor network, more robust and

interpretable KGE model, and the resulting improved predicted target embeddings contribute to the gains in performance. Additional details regarding the baseline models as well as the generation hyperparameters can be found in the Supplementary Material.

TABLE I  
TOP 5% DOCKING SCORE (KCAL/MOL)

Method	Target Protein				
	parp1	fa7	5ht1b	braf	jak2
REINVENT [8]	-8.702 ( $\pm$ 0.523)	-7.205 ( $\pm$ 0.264)	-8.770 ( $\pm$ 0.316)	-8.392 ( $\pm$ 0.400)	-8.165 ( $\pm$ 0.277)
MORLD [7]	-7.532 ( $\pm$ 0.260)	-6.263 ( $\pm$ 0.165)	-7.869 ( $\pm$ 0.650)	-8.040 ( $\pm$ 0.337)	-7.816 ( $\pm$ 0.133)
HierVAE [5]	-9.487 ( $\pm$ 0.278)	-6.812 ( $\pm$ 0.274)	-8.081 ( $\pm$ 0.252)	-8.978 ( $\pm$ 0.525)	-8.285 ( $\pm$ 0.370)
FREED [9]	-10.427 ( $\pm$ 0.177)	-8.297 ( $\pm$ 0.094)	-10.425 ( $\pm$ 0.331)	-10.325 ( $\pm$ 0.164)	-9.624 ( $\pm$ 0.102)
FREED-QS [9]	-10.579 ( $\pm$ 0.104)	-8.378 ( $\pm$ 0.044)	-10.714 ( $\pm$ 0.183)	-10.561 ( $\pm$ 0.080)	-9.735 ( $\pm$ 0.022)
GDSS [23]	-9.967 ( $\pm$ 0.028)	-7.775 ( $\pm$ 0.039)	-9.459 ( $\pm$ 0.101)	-9.224 ( $\pm$ 0.068)	-8.926 ( $\pm$ 0.089)
MOOD [26]	-10.865 ( $\pm$ 0.113)	-8.160 ( $\pm$ 0.071)	-11.145 ( $\pm$ 0.042)	-11.063 ( $\pm$ 0.034)	-10.147 ( $\pm$ 0.060)
KARL [33]	-11.475 ( $\pm$ 0.410)	-8.270 ( $\pm$ 0.211)	-11.149 ( $\pm$ 0.072)	-11.364 ( $\pm$ 0.035)	-10.636 ( $\pm$ 0.212)
K-DREAM (zero guidance)	-10.918 ( $\pm$ 0.592)	-8.390 ( $\pm$ 0.352)	-9.679 ( $\pm$ 0.345)	-10.402 ( $\pm$ 0.585)	-10.057 ( $\pm$ 0.483)
K-DREAM	<b>-12.137</b> ( $\pm$ 0.490)	<b>-9.078</b> ( $\pm$ 0.381)	<b>-11.484</b> ( $\pm$ 0.459)	<b>-11.407</b> ( $\pm$ 0.612)	<b>-10.391</b> ( $\pm$ 0.395)

Note: Results are the means and the standard deviations of 5 runs. Previous results sourced from [26]. Bold entries indicate the best value, where a lower value is better.

### B. K-DREAM Interpolates in Chemical Space to Generate Multitarget Compounds

Multitarget drug design presents a promising avenue for developing more effective therapeutics, particularly for complex diseases with interconnected pathways. Previous work [38] in this the field has cast this problem into the form of multiobjective latent space optimization, which uses combinations of *in-silico* predictions to produce better candidates. Our work improves upon this foundational idea by leveraging the rich contextual information of the biomedical knowledge graph. K-DREAM’s ability to interpolate between target embeddings enables the generation of molecules with tailored multitarget profiles. From our set of five target proteins from the previous experiment, we select JAK2 and PARP1 for a multitarget design task. The coactivation of JAK2 and PARP1 in hepatitis virus-associated liver cancer underscores the need for multitarget therapeutic approaches [39]. JAK2 is upregulated in virus-associated hepatocellular carcinoma (vHCC) and contributes to tumor progression. Concurrently, PARP1, a DNA repair enzyme, promotes therapy resistance in various cancers, including HCC [40]. Multitarget drugs addressing both JAK2 and PARP1 could potentially offer a more streamlined approach to combat the complex molecular landscape of HCC while mitigating the challenges associated with combination therapies.

### C. Ablation Studies

Table I compares the performance of K-DREAM without the guidance information from the knowledge graph for all five protein targets. We observe a significant decrease in docking scores across all targets, demonstrating the effectiveness of the embeddings in the guidance process

For quantitative evaluation of multitarget interpolation, we generated three compound sets consisting of 3000 molecules each: PARP1-targeted ( $P$ ), JAK2-targeted ( $J$ ), and interpolated-target ( $I$ ) (Fig. 3). PARP1-targeted compounds scored ( $\mu_{P,P} = -11.84$ ,  $\sigma = 0.54$ ) when evaluated against PARP1 and ( $\mu_{P,J} = -7.53$ ,  $\sigma = 1.84$ ) against JAK2. JAK2-targeted compounds achieved ( $\mu_{J,P} = -10.48$ ,  $\sigma = 0.95$ ) and ( $\mu_{J,J} = -10.21$ ,  $\sigma = 0.57$ ), respectively. Interpolated compounds attained ( $\mu_{I,P} = -11.06$ ,  $\sigma = 0.94$ ) and

( $\mu_{I,J} = -9.89$ ,  $\sigma = 0.79$ ). Post-hoc Hotelling’s  $T^2$  tests confirmed significant distinctions between all pairs ( $P$  versus  $I$ :  $T^2 = 354.22$ ,  $F = 176.67$ ;  $J$  versus  $I$ :  $T^2 = 65.52$ ,  $F = 32.68$ ;  $P$  versus  $J$ :  $T^2 = 935.85$ ,  $F = 466.95$ ; all  $p < 0.001$ ). Mahalanobis distances between centroids ( $D_{P,I} = 2.49$ ,  $D_{J,I} = 0.67$ ,  $D_{P,J} = 3.01$ ) demonstrated the interpolated set’s intermediate position. Interpolated compounds’ PARP1 and JAK2 scores showed a slight negative correlation ( $r = -0.16$ ,  $p < 0.05$ ). The proportion of compounds exceeding predefined thresholds ( $\text{PARP1} \leq -11.5$ ,  $\text{JAK2} \leq -10.0$ ) was highest for interpolated compounds (15.4%) compared to PARP1-targeted (9.6%) and JAK2-targeted (12.5%) sets, validating K-DREAM’s ability to generate compounds with balanced dual-inhibition profiles. (All docking scores are stated in kcal/mol.)

The Supplementary Material contains additional ablation studies with reference to the embedding dimension.

## IV. DISCUSSION

K-DREAM achieved state-of-the-art performance in targeted drug design, surpassing baseline models in generating molecules with superior docking scores. Unlike baseline models like GDSS, MOOD and others that directly incorporate docking score information, K-DREAM leverages knowledge graph representations to guide molecule generation. This approach demonstrates that biomedical knowledge graphs contain sufficient information to steer the generative process towards biologically relevant molecules. The model’s ability to generate high-scoring molecules without explicit docking information suggests it captures complex relationships between molecular structures and biological targets, potentially leading to more diverse and innovative drug candidates.

While docking scores provide standardized computational metrics, they have significant limitations as primary evaluation measures, relying on static protein structures and simplified scoring functions that cannot capture dynamic conformational changes or cellular complexity. Docking scores also fail to predict pharmacokinetic properties, selectivity, or functional activity essential for therapeutic efficacy. Future validation of

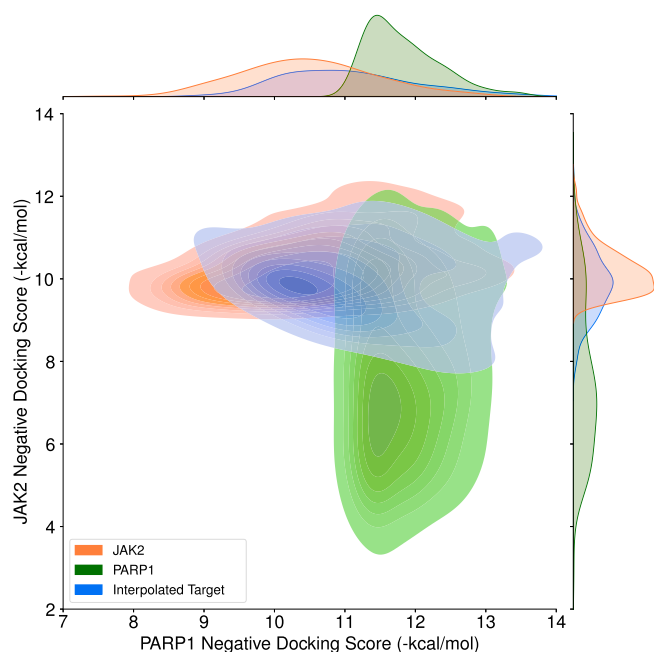


Fig. 3. *Multitarget drug design.* The distribution of the top 10% of molecules (ranked by docking score) generated by K-DREAM for PARP1-targeted (green), JAK2-targeted (orange), and an interpolated target (blue), evaluated by their docking scores with PARP1 (x-axis) and JAK2 (y-axis). We can see that the interpolated molecules achieve a balance between the two target proteins with a high mean score on both axes.

K-DREAM requires experimental studies including biochemical binding assays, cell-based functional activity tests, and *in vivo* pharmacokinetic studies to confirm that knowledge graph guidance translates to genuine improvements in biological activity and therapeutic potential beyond computational predictions.

In creating multitarget compounds, the existence of the IFNGR-JAK-STAT-PARP1 pathway between our chosen targets PARP1 and JAK2 in the knowledge graph exemplifies the interconnected nature of biological systems, which traditional drug discovery methods often struggle to address comprehensively. K-DREAM's integration of knowledge graphs with generative models represents a significant advancement in capturing and leveraging these complex relationships. By encoding pathway information and protein interactions in the knowledge graph, K-DREAM can generate compounds that simultaneously modulate multiple targets within a relevant biological context. Statistical analysis validates this capability: interpolated compounds achieve balanced binding to both targets ( $\mu_{I,P} = -11.06$ ,  $\mu_{I,J} = -9.89$ ) with only minimal trade-offs between PARP1 and JAK2 binding ( $r = -0.16$ ,  $p < 0.05$ ). Most importantly, the interpolated compounds show superior dual-targeting with 15.4% exceeding both binding thresholds compared to 9.6% and 12.5% for single-target optimization. This approach enables the exploration of chemical space guided by system-level understanding, potentially yielding molecules with optimized efficacy and reduced off-target effects. The ability to interpolate between target embeddings further allows for fine-tuning of multitarget profiles, offering a powerful tool

for rational drug design in complex diseases where multiple pathways contribute to pathogenesis.

A key dependence of the model is in the weight of the guidance mechanism relative to the unconditioned generative process, controlled by the hyperparameter  $\lambda_X$ , which quantifies the contribution of the multidimensional Gaussian distribution centered at the target embedding. Varying  $\lambda_X$  has the effect of changing the variance of the Gaussian, with higher  $\lambda_X$  values corresponding to a sharper peak around the target, leading to improved docking scores. We observe an optimal value for  $\lambda_X$ , since an increasingly sharper peak might lead to suboptimal guidance in the initial stages of the diffusion process (Fig. 4).

We note that the increased standard deviation in K-DREAM's scores, despite higher mean scores, could be attributed to the curse of dimensionality in high-dimensional Gaussian distributions [41]. This becomes a factor due to our guidance mechanism, which uses a Gaussian distribution to guide the generative process towards target embeddings. In high-dimensional spaces, the probability mass of a multivariate Gaussian concentrates in a thin shell away from the mean. This concentration effect leads to a larger variety of molecules, as the generative process is guided towards points on this shell rather than converging to a single central point. The volume of this shell grows exponentially with dimensionality, providing a vast space for diverse, high-scoring candidates. Consequently, while K-DREAM consistently produces molecules with improved docking scores, the diversity within this high-scoring set contributes to larger standard deviations. This behavior aligns with the model's objective of exploring a broader, biologically relevant chemical space, rather than converging to a narrow set of solutions.

## V. CONCLUSION

This study demonstrates that biomedical knowledge graphs contain valuable information that can effectively guide molecular generation toward biologically relevant compounds. By embedding this knowledge into a diffusion-based generative framework, K-DREAM produces drug candidates with improved binding potential across multiple therapeutic targets while maintaining chemical validity and diversity. The framework's performance in targeted drug design tasks consistently exceeded baseline models across all tested protein targets. Unlike previous approaches that directly optimize for docking scores, K-DREAM leverages the rich contextual information within knowledge graphs to steer the generative process toward promising regions of chemical space without explicitly incorporating target-specific scoring functions. Perhaps most significantly, K-DREAM's ability to interpolate between target embeddings enables the generation of molecules with balanced multitarget profiles. This capability addresses a critical need in modern drug discovery, where many diseases involve multiple interconnected pathways that cannot be effectively addressed by single-target approaches. The integration of knowledge graph embeddings with generative models creates a unified space where chemical structures and biological interactions can be jointly explored. This paradigm shift moves beyond traditional

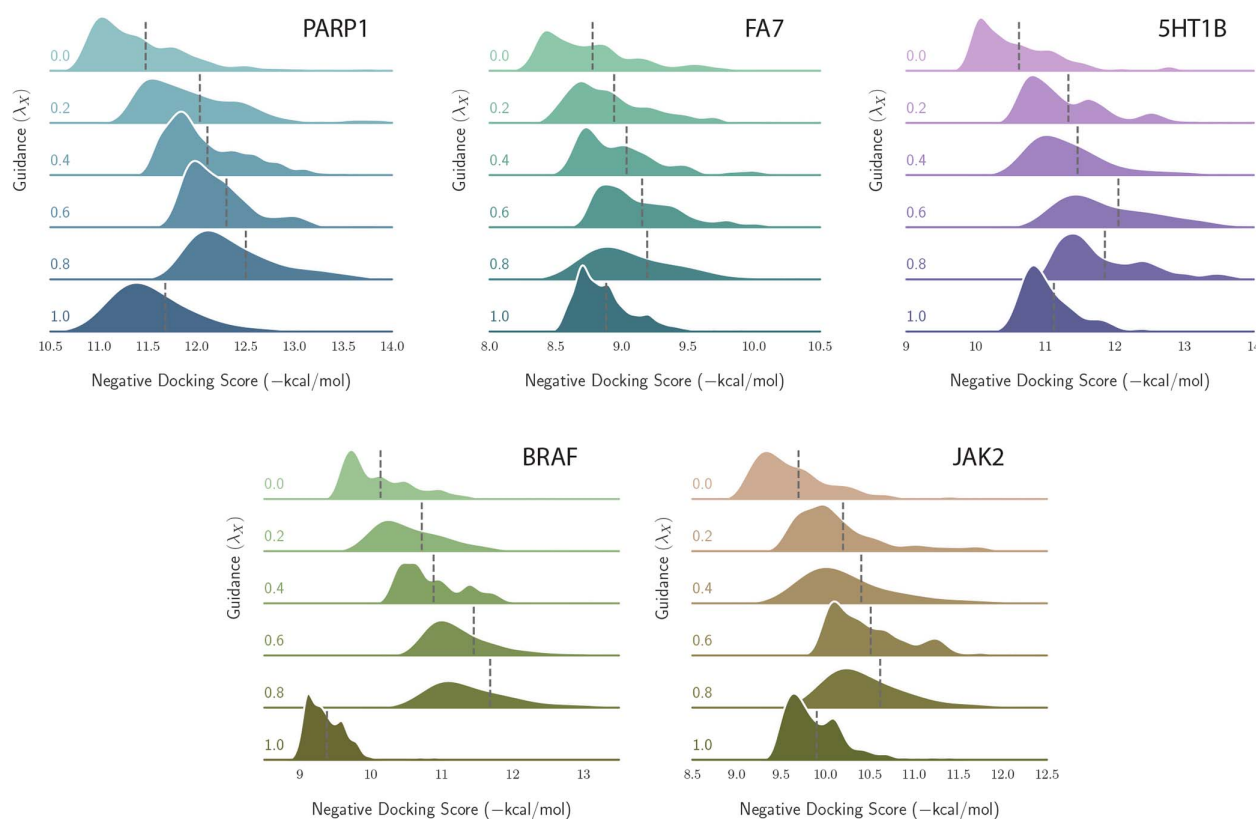


Fig. 4. Molecular docking score distribution with varying guidance levels. The extent of guidance is controlled with a hyperparameter  $\lambda_X$  that determines the weight of the guidance term in the loss function.

property-based optimization toward a more comprehensive, biologically informed approach to drug design.

Future research could further enhance this framework by incorporating additional biomedical knowledge sources, exploring alternative embedding techniques, and extending the methodology to address other drug design challenges such as optimizing pharmacokinetic properties alongside target binding. By continuing to bridge the gap between computational generation and biological understanding, knowledge-enhanced generative models may significantly accelerate the discovery of novel therapeutics for unmet medical needs.

#### ACKNOWLEDGMENT

The authors acknowledge the use of ChatGPT and Claude to refine human-written text in this manuscript.

#### REFERENCES

- [1] P. Chandak, K. Huang, and M. Zitnik, "Building a knowledge graph to enable precision medicine," *Sci. Data*, vol. 10, Feb. 2023.
- [2] A. Bordes, N. Usunier, A. Garcia-Duran, J. Weston, and O. Yakhnenko, "Translating embeddings for modeling multi-relational data," in *Advances in Neural Information Processing Systems*, C. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K. Weinberger, Eds., vol. 26. New York, NJ, USA: Curran Associates, Inc., 2013.
- [3] H. Cao et al., "A survey on generative diffusion models," *IEEE Trans. Knowl. Data Eng.*, vol. 36, no. 7, pp. 2814–2830, Jul. 2024.
- [4] J. Chen, B. Ganguly, Y. Xu, Y. Mei, T. Lan, and V. Aggarwal, "Deep generative models for offline policy learning: Tutorial, survey, and perspectives on future directions," *Trans. Mach. Learn. Res.*, vol. 24, no. 6, pp. 131–145, 2024.
- [5] W. Jin, R. Barzilay, and T. Jaakkola, "Hierarchical generation of molecular graphs using structural motifs," in *Proc. Int. Conf. Mach. Learn.*, PMLR, 2020, pp. 4839–4848.
- [6] L. Kong, J. Cui, H. Sun, Y. Zhuang, B. A. Prakash, and C. Zhang, "Autoregressive diffusion model for graph generation," in *Proc. Int. Conf. Machine Learn.*, PMLR, 2023, pp. 17391–17408.
- [7] W. Jeon and D. Kim, "Autonomous molecule generation using reinforcement learning and docking to develop potential novel inhibitors," *Sci. Rep.*, vol. 10, 2020, Art. no. 12.
- [8] M. Olivecrona, T. Blaschke, O. Engkvist, and H. Chen, "Molecular de-novo design through deep reinforcement learning," *J. Cheminf.*, vol. 9, no. 1, pp. 1–14, 2017.
- [9] S. Yang, D. Hwang, S. Lee, S. Ryu, and S. J. Hwang, "Hit and lead discovery with explorative RL and fragment-based molecule generation," *Adv. Neural Inf. Process. Syst.*, vol. 34, no. 4, pp. 7924–7936, 2021.
- [10] P. Ertl and A. Schuffenhauer, "Estimation of synthetic accessibility score of drug-like molecules based on molecular complexity and fragment contributions," *J. Cheminf.*, vol. 1, pp. 1–11, 2009.
- [11] G. R. Bickerton, G. V. Paolini, J. Besnard, S. Muresan, and A. L. Hopkins, "Quantifying the chemical beauty of drugs," *Nat. Chem.*, vol. 4, no. 2, pp. 90–98, 2012.
- [12] D. S. Himmelstein et al., "Systematic integration of biomedical knowledge prioritizes drugs for repurposing," *Elife*, vol. 6, 2017, Art. no. e26726.
- [13] S. L. Harris and A. J. Levine, "The p53 pathway: Positive and negative feedback loops," *Oncogene*, vol. 24, no. 17, pp. 2899–2908, 2005.
- [14] A. Wagner, "Robustness against mutations in genetic networks of yeast," *Nat. Genet.*, vol. 24, no. 4, pp. 355–361, 2000.
- [15] L. Fernández and R. E. Hancock, "Adaptive and mutational resistance: Role of porins and efflux pumps in drug resistance," *Clin. Microbiol. Rev.*, vol. 25, no. 4, pp. 661–681, 2012.
- [16] S. K. Mohamed, V. Nováček, and A. Nounu, "Discovering protein drug targets using knowledge graph embeddings," *Bioinformatics*, vol. 36, no. 2, pp. 603–610, 2020.
- [17] P. J. Richardson, B. W. S. Robinson, D. P. Smith, and J. Stebbing, "The AI-assisted identification and clinical efficacy of baricitinib in the treatment of covid-19," *Vaccines*, vol. 10, no. 6, 2022.

- [18] X. Lin, Z. Quan, Z.-J. Wang, T. Ma, and X. Zeng, "KGNN: Knowledge graph neural network for drug-drug interaction prediction," *IJCAI*, vol. 380, pp. 2739–2745, 2020.
- [19] D. Bang, S. Lim, S. Lee, and S. Kim, "Biomedical knowledge graph learning for drug repurposing by extending guilt-by-association to multiple layers," *Nat. Commun.*, vol. 14, no. 1, p. 3570, Jan. 2023.
- [20] X. Zeng, X. Tu, Y. Liu, X. Fu, and Y. Su, "Toward better drug discovery with knowledge graph," *Curr. Opin. Struct. Biol.*, vol. 72, pp. 114–126, 2022.
- [21] M. Ali et al., "PyKEEN 1.0: A Python library for training and evaluating knowledge graph embeddings," *J. Mach. Learn. Res.*, vol. 22, no. 82, pp. 1–6, 2021.
- [22] Y. Song, L. Shen, L. Xing, and S. Ermon, "Solving inverse problems in medical imaging with score-based generative models," in *Proc. Int. Conf. Learn. Representations*, 2022, pp. 167–191.
- [23] J. Jo, S. Lee, and S. J. Hwang, "Score-based generative modeling of graphs via the system of stochastic differential equations," in *Proc. 39th Int. Conf. Mach. Learn.*, 2022, pp. 169–187.
- [24] J. J. Irwin, T. Sterling, M. M. Mysinger, E. S. Bolstad, and R. G. Coleman, "Zinc: A free tool to discover chemistry for biology," *J. Chem. Inf. Model.*, vol. 52, no. 7, pp. 1757–1768, 2012.
- [25] P. Dhariwal and A. Q. Nichol, "Diffusion models beat GANs on image synthesis," in *Proc. Adv. Neural Inf. Process. Syst.*, 2021, pp. 157–164.
- [26] S. Lee, J. Jo, and S. J. Hwang, "Exploring chemical space with score-based out-of-distribution generation," in *Proc. 40th Int. Conf. Mach. Learn.*, 2023, pp. 18872–18892.
- [27] J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, "Deep unsupervised learning using nonequilibrium thermodynamics," in *Proc. Int. Conf. Machine Learn.* PMLR, 2015, pp. 2256–2265.
- [28] C. Vignac, I. Krawczuk, A. Siraudin, B. Wang, V. Cevher, and P. Frossard, "Digress: Discrete denoising diffusion for graph generation," in *Proc. 11th Int. Conf. Learn. Representations*, 2023, pp. 224–239.
- [29] A. Oring, Z. Yakhini, and Y. Hel-Or, "Autoencoder image interpolation by shaping the latent space," in *Proc. 38th Int. Conf. Mach. Learn.*, 2021, pp. 8281–8290.
- [30] Y. Du, X. Liu, N. M. Shah, S. Liu, J. Zhang, and B. Zhou, "Chemspace: Interpretable and interactive chemical space exploration," *Trans. Mach. Learn. Res.*, vol. 27, pp. 3579–3581, 2023.
- [31] L. McInnes, J. Healy, N. Saul, and L. Großberger, "Umap: Uniform manifold approximation and projection," *J. Open Source Softw.*, vol. 3, no. 29, 2018, Art. no. 861.
- [32] H. Zhou, P. Xu, X. Yuan, and H. Qu, "Edge bundling in information visualization," *Tsinghua Sci. Technol.*, vol. 18, no. 2, pp. 145–156, 2013.
- [33] A. Malusare and V. Aggarwal, "Improving molecule generation and drug discovery with a knowledge-enhanced generative model," *IEEE Trans. Comput. Biol. Bioinf.*, vol. 22, no. 1, pp. 375–381, Jan. 2025.
- [34] M. M. Mysinger, M. Carchia, J. J. Irwin, and B. K. Shoichet, "Directory of useful decoys, enhanced (dud-e): Better ligands and decoys for better benchmarking," *J. Med. Chem.*, vol. 55, pp. 6582–6594, Jul. 2012.
- [35] A. Alhossary, S. D. Handoko, Y. Mu, and C.-K. Kwok, "Fast, accurate, and reliable molecular docking with quickvina 2," *Bioinformatics*, vol. 31, no. 13, pp. 2214–2216, 2015.
- [36] D. Bajusz, A. Rácz, and K. Héberger, "Why is Tanimoto index an appropriate choice for fingerprint-based similarity calculations?" *J. Cheminf.*, vol. 7, pp. 2789–2798, May 2015.
- [37] J. Jo, S. Lee, and S. J. Hwang, "Score-based generative modeling of graphs via the system of stochastic differential equations," in *Proc. Int. Conf. Mach. Learn.* PMLR, 2022, pp. 10362–10383.
- [38] R. Winter, F. Montanari, A. Steffen, H. Briem, F. Noé, and D.-A. Clevert, "Efficient multi-objective molecular optimization in a continuous latent space," *Chem. Sci.*, vol. 10, pp. 8016–8024, 2019.
- [39] Y.-G. Cherng et al., "Induced mitochondrial alteration and DNA damage via IFNGR-JAK2-STAT1-PARP1 pathway facilitates viral hepatitis associated hepatocellular carcinoma aggressiveness and stemness," *Cancers*, vol. 13, Jun. 2021, Art. no. 3755.
- [40] M. Yu et al., "PARG inhibition limits HCC progression and potentiates the efficacy of immune checkpoint therapy," *J. Hepatol.*, vol. 77, no. 1, pp. 140–151, 2022.
- [41] V. Berisha et al., "Digital medicine and the curse of dimensionality," *NPJ Digit. Med.*, vol. 4, no. 1, 2021, Art. no. 153.