

Digital Discovery

Accepted Manuscript

This article can be cited before page numbers have been issued, to do this please use: J. Song, K. Yang, Y. Xiong, K. Tao, J. Ji, P. Cao and L. Cai, *Digital Discovery*, 2025, DOI: 10.1039/D5DD00504C.



This is an Accepted Manuscript, which has been through the Royal Society of Chemistry peer review process and has been accepted for publication.

Accepted Manuscripts are published online shortly after acceptance, before technical editing, formatting and proof reading. Using this free service, authors can make their results available to the community, in citable form, before we publish the edited article. We will replace this Accepted Manuscript with the edited and formatted Advance Article as soon as it is available.

You can find more information about Accepted Manuscripts in the [Information for Authors](#).

Please note that technical editing may introduce minor changes to the text and/or graphics, which may alter content. The journal's standard [Terms & Conditions](#) and the [Ethical guidelines](#) still apply. In no event shall the Royal Society of Chemistry be held responsible for any errors or omissions in this Accepted Manuscript or any consequences arising from the use of any information it contains.

1 **AI-Driven Antiviral Natural Products Drug Development: A**

2 **Technical Overview**

3 Junxi Song^{1,2†}, Kunhuan Yang^{2†}, Yingcai Xiong^{3†}, Keyu Tao^{2†}, Liangyu Cai^{1*}, Peng
4 Cao^{4*}, Jianjian Ji^{1,2*}

5 †These authors contributed equally: Junxi Song, Kunhuan Yang, Yingcai Xiong, Keyu
6 Tao

7 ¹Wuxi Affiliated Hospital of Nanjing University of Chinese Medicine, Wuxi, China.

8 ²Jiangsu Key Laboratory of Children's Health and Chinese Medicine, The First Clinical
9 College, Nanjing University of Chinese Medicine, Nanjing, China.

10 ³The State Key Laboratory of Pharmaceutical Biotechnology, Chemistry and
11 Biomedicine Innovation Center (ChemBIC), Division of Immunology, Medical School,
12 Nanjing University, Nanjing, China, 210093

13 ⁴State Key Laboratory on Technologies for Chinese Medicine Pharmaceutical Process
14 Control and Intelligent Manufacture, Nanjing University of Chinese Medicine, Nanjing,
15 China.

16 ***Correspondence:**

17 Jianjian Ji (jijj@njucm.edu.cn), Nanjing University of Chinese Medicine, 138 Xianlin
18 Avenue, Nanjing, China, 210023.

19 Peng Cao (cao_peng@njucm.edu.cn), Jiangsu Provincial Medical Innovation Center,
20 Affiliated Hospital of Integrated Traditional Chinese and Western Medicine, Nanjing
21 University of Chinese Medicine, Nanjing, 210028, China.

22 Liangyu Cai (wxzy018@njucm.edu.cn) Wuxi Affiliated Hospital of Nanjing
23 University of Chinese Medicine, Wuxi, China.

24 **Co-authors E-mails**

25 Junxi Song: 039222333@njucm.edu.cn

26 Kunhuan Yang: 039222215@njucm.edu.cn

27 Yingcai Xiong: yingcai_x@163.com



28 Keyu Tao: tky615322123@163.com

View Article Online
DOI: 10.1039/D5DD00504C

29

30 **Abstract**

31 The emergence of viral pandemics and rapid pathogen evolution presents formidable
32 challenges for conventional antiviral development, including prolonged timelines, high
33 costs, and susceptibility to resistance mechanisms. Natural products (NPs) offer
34 promising antiviral potential through structural diversity and multi-target synergism,
35 while their development faces critical bottlenecks in structural characterization, target
36 identification, and synthetic optimization. Given the current situation, artificial
37 intelligence (AI), particularly machine learning (ML) and deep learning (DL), is
38 revolutionizing drug development by transforming data analysis and predictive
39 modeling. This review explores AI applications across the antiviral NPs drug
40 development continuum, providing insights for AI-driven pharmaceutical research.

41

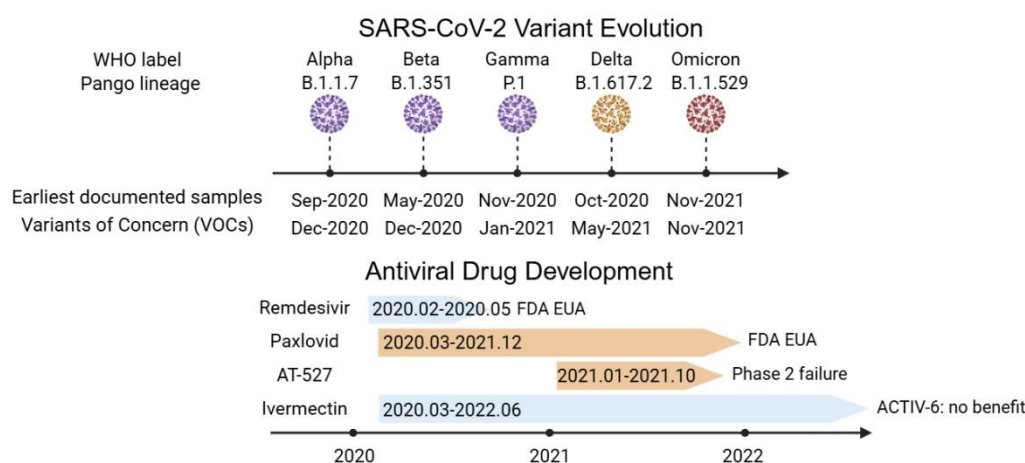
42 **1. Introduction**

43

44 In recent years, many new and re-emerging viral pathogens—such as severe acute
45 respiratory syndrome coronavirus (SARS-CoV), Middle East respiratory syndrome
46 coronavirus (MERS-CoV), and SARS-CoV-2—have continued to pose a risk to global
47 public health(1). Antiviral medicine development is associated with high costs and
48 extended timescales(2). Moreover, the current pace of pharmaceutical research and
49 development critically lags behind the exponentially growing need for rapid therapeutic
50 interventions during emergent pandemic scenarios (**Fig.1**). The remdesivir
51 development timeline exemplifies both the opportunities and persistent limitations in
52 this race against time. Initially explored in 2009 for hepatitis C and later Ebola, it gained
53 emergency use authorization in May 2020 for COVID-19—only half a year after the



54 SARS-CoV-2 outbreak, showcasing a smooth translation from scientific discovery to
 55 emergency response(3). Yet, as depicted in Fig. 1, even this rapid repurposing occurred
 56 amid ongoing viral evolution, with variants like Alpha and Beta emerging before full
 57 clinical implementation, highlighting that such successes depend on pre-existing
 58 scaffolds and may not suffice for de novo threats. For novel viral threats where no such
 59 prior knowledge exists, the de novo drug development process remains substantially
 60 slower than the pace of outbreaks (4). Due to the virus's unique genetic system (lack of
 61 complex genetic information synthesis proofreading system), the virus can rapidly
 62 mutate and evade drug treatment. Consider the H1N1 influenza variants: Ingenious
 63 mutations in the Sb region of the HA protein serve like a molecular camouflage kit,
 64 allowing them to evade vaccine-educated antibodies and leaving vaccination
 65 campaigns in the lurch(5). Such complex pressures require more sophisticated
 66 therapeutic strategies, which demand transformed drug development paradigms.



67
 68 **Fig. 1 Asynchrony Between SARS-CoV-2 Variant Evolution and Antiviral Drug**
 69 **Development.** This figure illustrates the evolutionary dynamics of major SARS-CoV-
 70 2 variants during the COVID-19 pandemic, aligned with the development and clinical
 71 implementation timelines of four representative small-molecule antivirals. The upper
 72 axis delineates SARS-CoV-2 evolutionary dynamics using a generational color scheme:
 73 purple signifies the initial waves of Variants of Concern (Alpha, Beta, and Gamma)
 74 marked by early increases in transmissibility; yellow identifies the Delta variant as a



75 pivotal transition toward significantly higher viral loads and pathogenicity, and red
76 represents the Omicron lineage, which constitutes a fundamental paradigm shift in
77 immune evasion and mutation density. Each variant is annotated with its earliest
78 documented detection date and official WHO Variant of Concern (VOC) designation,
79 with earliest documented detections referring to retrospectively identified sequences
80 rather than real-time discovery, to highlight the inherent delay in global surveillance
81 and response. Antivirals in the lower axis were selected according to two defining
82 dimensions: the development pathway (repurposing of established agents versus de
83 novo structural design) and the clinical resolution (regulatory authorization versus
84 termination due to futility). Blue bands (Remdesivir and Ivermectin) represent the drug
85 repurposing strategy, aimed at immediate deployment based on known safety profiles.
86 In contrast, orange bands (Paxlovid and AT-527) signify de novo discovery programs.
87 Despite such technological acceleration, vertical alignment across the axes reveals that
88 by the time these high-potency agents reached their respective clinical endpoints, the
89 viral landscape had already transitioned through multiple generational cycles,
90 illustrating the persistent structural lag between therapeutic intervention and emergent
91 pandemic needs. This figure highlights the asynchrony between viral evolution and
92 therapeutic development, underscoring the challenges in maintaining efficacy against
93 phylogenetically divergent lineages. Figure created with BioRender.com.

94 Natural products—complex metabolites produced by plants, fungi, animals, and
95 microorganisms—exhibit the highest chemical diversity in nature(6), and are an
96 important source of antiviral medicines. Many licensed treatments and prospects stem
97 directly or indirectly from plants, microbes, and marine animals. For instance,
98 artemisinin's antimalarial potency signified a milestone for natural-product-based anti-
99 infective discovery(7); diammonium glycyrrhizinate from licorice root is authorized in
100 China and Japan as an adjuvant for chronic hepatitis B(8); and numerous antivirals (e.g.,
101 acyclovir, ganciclovir, vidarabine, zidovudine) emerged from natural leads via
102 structural optimization(9). Despite this promise, development confronts obstacles:
103 limited access to bioactive substances (only ~1% of microbial species are



104 culturable(10)); complex metabolite isolation and characterization (bioassay-guided
105 fractionation and structural elucidation are often essential); uncertain pharmacological
106 mechanisms (e.g., the multi-component synergy of Lianhua Qingwen capsules remains
107 incompletely defined(11)); and complex synthesis (e.g., up to 30 enzymatic cascade
108 steps to generate vinblastine(12)).

109 AI provides transformative solutions for these challenges. Advanced algorithms—
110 including Transformer architectures and graph neural networks (GNNs)—have made
111 great strides in accuracy for predicting drug-target interaction(13). In addition, the
112 family of biomedical data is exploding (e.g. the ChEMBL compound library
113 contains >20.3 million bioactivity measurements and >2.4 million unique
114 compounds(14). Furthermore, the Traditional Chinese Medicine Systems
115 Pharmacology Database (TCMSP) includes 29,384 components, 3,311 with targets,
116 and 837 linking to diseases(15), thus simplifying model training. On the other hand,
117 GPU clusters speedup computational tasks, improving by several orders of magnitude
118 the speed of training large-scale AI models, such as DL architectures for drug-target
119 interaction prediction, when compared with CPUs(16, 17).

120 The recent emergence of AI-driven platforms has already yielded significant
121 breakthroughs in antiviral discovery. For instance, a sophisticated AI pipeline recently
122 identified established antiretrovirals, such as bictegravir and etravirine, as potent broad-
123 spectrum inhibitors against monkeypox virus and related poxviruses(18). While such
124 successes underscore the transformative potential of AI in accelerating drug
125 repurposing and novel application discovery, there remains a need for a more granular
126 technical overview focused specifically on the end-to-end integration of AI across the
127 entire antiviral natural product development continuum.

128 In this Review we summarize some of the critical applications of AI in upstream
129 (resource mining and target identification), midstream (drug candidate screening and
130 optimization), and downstream (preclinical and clinical stages). We subsequently
131 explore the current technological limitations and emerging possibilities. Finally, we
132 discuss the outline future directions for the field. We hope to highlight a new era of



133 technology, efficiency and precision in drug development that is expected to speed
 134 delivery of new and improved medicines to patients.

135

136 2. AI applications in antiviral natural products drug development

137

138 AI technologies are increasingly integrated into various stages of natural product-based
 139 antiviral drug development, forming a systematic and intelligent pipeline from resource
 140 mining and target identification to preclinical and clinical applications (Fig. 2).

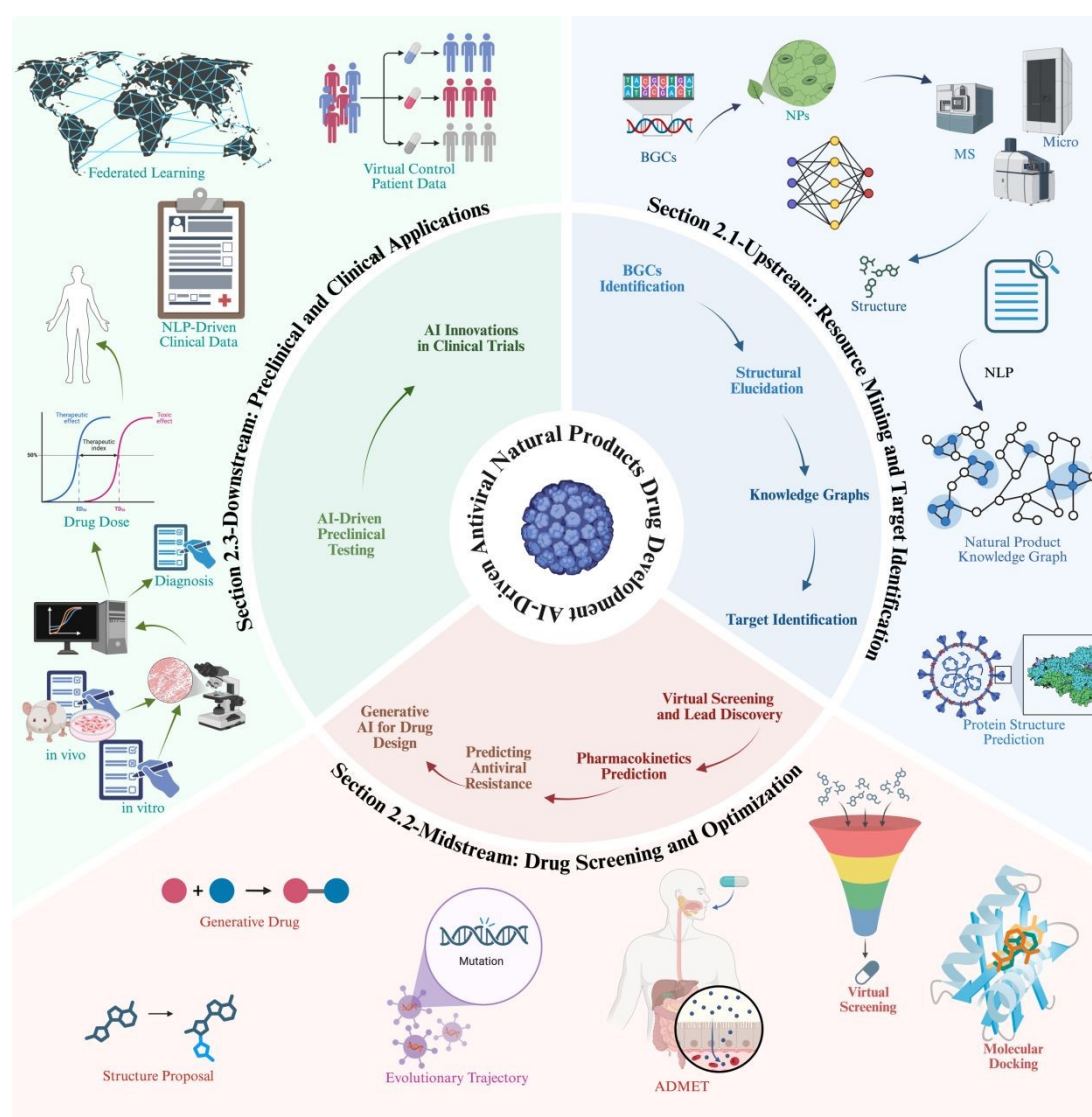


Fig. 2 AI-Driven Strategies for Antiviral Natural Products Drug Development Across the Full Pipeline. This figure illustrates the comprehensive application of



AI throughout the full pipeline of antiviral natural products drug development. It is divided into three main stages: Upstream (Resource Mining and Target Identification), Midstream (Drug Screening and Optimization), and Downstream (Preclinical and Clinical Applications). In the upstream stage, AI facilitates biosynthetic gene clusters (BGCs) identification, structural elucidation of NPs, construction of knowledge graphs, protein structure prediction, and target identification. In the midstream stage, AI supports virtual screening, molecular docking, ADMET prediction, evolutionary trajectory analysis, antiviral resistance prediction, and generative drug design. In the downstream stage, AI enhances preclinical testing, dose optimization, diagnostic assistance, NLP-driven clinical data mining, federated learning, virtual control cohort construction, and innovations in clinical trials. Different color-coded sections visualize key processes across stages, highlighting the integrative role of AI in advancing natural product-based antiviral drug discovery. Figure created with BioRender.com.

141

142 **2.1 Upstream: Resource Mining and Target Identification**

143 **2.1.1 AI for Genome Mining and Biosynthetic Gene Clusters (BGCs)**

144 Multiple antiviral NPs are the products of secondary metabolites prescribed in
145 microbial(19, 20) and plant genomes(21-23). AI is enhancing the search for these
146 biosynthetic gene clusters. Standard genome-mining tools, such as anti-SMASH,
147 employ rule-based pattern matching to identify known classes of BGCs, and this results
148 in under-representative identification of clusters that are “atypical” clusters and in
149 significant false-negative rates(24). Recognizing small sequence motifs beyond human-
150 defined elements, DL approaches can be used to discover novel BGCs. For example,
151 the RNN-based DeepBGC found novel BGCs in Streptomyces genomes(25). Such AI
152 models leverage gene context, conserved domains, and amino acid properties
153 generation to illuminate potentially hidden BGCs that may yield new antiviral drugs.



154 Bridging genetic and metabolomic data, these growing algorithms can correlate
155 putative gene clusters with actual molecules, completing the link from genes to the
156 NPs they produce. Genome-to-metabolite prediction is essential for discovering new
157 antiviral chemicals that are stored in nature's genomic libraries.

158 2.1.2 AI-Assisted Structural Elucidation

159 The unambiguous structure determination of NPs is typically a labor-intensive, time-
160 consuming process(26). AI enhances understanding of complex spectrometric data,
161 including mass spectrometry (MS) and liquid chromatography–mass spectrometry (LC-
162 MS)(27, 28). Methods for AI-assisted structural elucidation can be broadly classified
163 into three categories: (1) ML/DL-based MS/MS spectrum annotation and substructure
164 prediction; (2) predictive modeling for LC-MS retention time, peak feature extraction,
165 and spectrum grouping; and (3) integrated pipelines for advanced techniques like
166 microcrystal electron diffraction (MicroED). This classification reflects the shift from
167 rule-based manual analysis to data-driven automation(29).

168 Recent advancements further strengthen this capability: MZmine 3, a scalable
169 open-source platform, supports integrative processing of multimodal MS data
170 (including LC-MS and ion mobility), enabling efficient feature detection, visualization,
171 and annotation tailored to natural product workflows(30). Similarly, studies employing
172 LC-MS/MS and molecular networking have identified marine-derived secondary
173 metabolites with anti-SARS-CoV-2 activity, such as homofascaplysin A and aureol,
174 though structural confirmation remains labor-intensive due to stereochemical
175 complexity(31). MicroED with streamlined AI pipelines has enabled 3D structure
176 determination of macrocyclic NPs with antiviral potential, overcoming the need for
177 large single crystals by utilizing microcrystals—yet sample preparation is challenging,
178 as NPs often form amorphous or poorly diffracting microcrystals, limiting resolution
179 and throughput(32). ML-based retention-time prediction further enhances confidence
180 by avoiding re-isolation of known compounds, but struggles with NPs' high chemical



181 diversity and batch variability(33-35). In summary, AI-facilitated spectrometric data
182 analysis has begun to streamline structural elucidation in natural product research,
183 thereby facilitating the identification of potential antiviral candidates. Nevertheless,
184 overcoming NPs-inherent hurdles—such as mixture complexity, data scarcity for rare
185 scaffolds, and the substantial validation gap—will be essential to translate these
186 computational advances into robust, clinically relevant antiviral natural products.

187 **2.1.3 Knowledge Graphs and natural language processing (NLP) for Natural** 188 **Product Data**

189 New tools, including AI, are also increasingly utilized in organizing and mining this
190 vast knowledge repository of NPs and traditional medicine. Knowledge graphs are
191 structured networks that integrate heterogeneous data sources, such as chemical
192 structures, biological targets, biosynthetic pathways, and literature references, enabling
193 cross-domain analysis(36). NLP methods, on the other hand, extract structured
194 information from unstructured texts, particularly historical herb manuals and medical
195 literature. knowledge graphs offer high connectivity and query efficiency but face
196 challenges in entity resolution and data integration due to the heterogeneity of NP
197 sources (e.g., varying nomenclature and incomplete annotations). NLP excels at text
198 mining but struggles with ancient language ambiguity, OCR errors, and domain-
199 specific terminology in traditional medicine texts.

200 In applications, knowledge graphs have demonstrated value. For example, a
201 knowledge graph on natural products connects segments of tandem MS to predicted
202 metabolites and those predicted metabolites to potential generating genes, as
203 implemented in frameworks like the Experimental Natural Products Knowledge Graph
204 (ENPKG), which integrates multimodal data for plant-derived compounds(37). Recent
205 efforts further demonstrate this by leveraging AI to associate MS fragmentation patterns
206 with biosynthetic gene clusters through substructure discovery and BGC-metabolite
207 mapping(38). This graph can emulate an expert chemist's intuition, mining



208 correlations that yield novel antivirals. In the textual side, NLP methods are being
209 applied to the vast literature on medicinal plants and traditional therapies. One of them
210 constructed TCMBank, one of the largest integrative databases associating TCM with
211 multi-omics data, based on text-mining historical herb manuals and medical texts. The
212 system, employing advanced NLP techniques initially based on bidirectional LSTM
213 networks and conditional random fields with subsequent enhancements incorporating
214 Transformer-based models, amassed structured data on herbs, chemicals and known
215 effects from dozens of ancient manuscripts(39). This “knowledge reconstruction” re-
216 works past empirical material into a machine-readable database enabling the AI to
217 rapidly sift through potential antiviral medicines in conventional literature and then
218 match them against prevailing biomedical information.

219 By transforming fragmented knowledge into connected, machine-readable data
220 through knowledge graphs and curated databases, AI enhances upstream discovery of
221 natural antivirals. Researchers can now query these systems to identify promising
222 compounds and targets far more efficiently than manual curation. Nevertheless, the
223 field must address NPs-specific hurdles—such as textual ambiguity in ancient sources,
224 data heterogeneity, and the persistent validation gap—to ensure reliable, translational
225 impacts in antiviral drug development.

226 **2.1.4 AI-Driven Target Identification**

227 Another essential upstream step is targeting molecular targets (viral or host) for natural
228 antivirals. AI accelerates this process through a combination of data-driven and
229 structure-based approaches, enabling more efficient prediction and validation(40). To
230 achieve comprehensive coverage, we organize AI-driven target identification into a
231 trinity framework: Chemical-centric, Systems-centric, and Physics-centric. This
232 structure addresses ligand-based, network-based, and structure-based paradigms in
233 pharmacology, while incorporating emerging techniques to fill critical gaps such as
234 phenotypic profiling, metabolomic interference, and dynamic simulations.



235 (i) Chemical-centric strategies focus on ligand-chemical space using deep transfer
236 learning and matrix factorization to associate "orphan ligands" with known targets,
237 including phenotypic/image-based AI (e.g., Cell Painting assays) that enables forward
238 pharmacology by inferring pathways from cellular morphological fingerprints without
239 prior ligand-target data(41). For instance, the STarFish platform, a stacked ensemble
240 model, identifies potential targets for NPs by leveraging known ligand-target data,
241 achieving high accuracy in multi-target predictions on benchmark datasets(42).
242 Similarly, DeepPurpose employs DL for drug-target interaction prediction, facilitating
243 virtual screening of NPs with improved hit rates(43).

244 (ii) Systems-centric strategies integrate knowledge graphs and multimodal GNNs
245 to reveal hidden nodes in virus-host interaction networks, integrated with AI-driven
246 metabolomic analysis (e.g., flux balance analysis in integrated metabolic models) to
247 predict how NPs alter host metabolic environments, uncovering indirect antiviral
248 targets(44). For example, graph neural networks have been applied to construct and
249 analyze SARS-CoV-2 knowledge graphs based on virus-host interactions, pathways,
250 and drug associations, identifying potential host genes and biological processes for
251 antiviral drug repurposing(45). Tools like TCMBank leverage NLP and knowledge
252 graphs for traditional medicine data mining, enabling synergistic target discovery(39).
253 This dimension also extends to genomic surveillance for pathogen evolution, enhancing
254 target relevance in dynamic viral contexts(46).

255 (iii) Physics-centric strategies employ generative AI for biophysical predictions
256 and dynamic simulations. AlphaFold 3 exemplifies this by providing ~50% improved
257 accuracy in predicting protein-ligand and nucleic acid interactions compared to
258 physics-based docking(47) and RoseTTAFold All-Atom for all-atom modeling of
259 protein-ligand dynamics(48).

260 Overall, these AI technologies—from knowledge graphs to predictive models—
261 hold significant potential to streamline upstream discovery by connecting chemical



262 leads with biological targets, paving the way for next-generation antiviral drugs.
263 Notably, polypharmacology represents a core advantage of NPs: their multi-target
264 effects enable synergistic therapeutic outcomes, and AI shows emerging potential to
265 actively design and quantify these synergies or antagonisms for optimized efficacy(49).
266 However, challenges persist, including data bias in training sets (e.g.,
267 underrepresentation of rare interactions leading to high false-positive rates in a data-
268 starved setting) and the experimental-computational gap, necessitating rigorous
269 validation and diverse datasets to mitigate biases and improve generalization.

270 **2.2 Midstream: Drug Screening and Optimization**

271 **2.2.1 Virtual Screening and Lead Discovery**

272 In the screening and lead identification phase, AI has significantly enhanced processes
273 by enabling ultra-efficient virtual screening of enormous chemical libraries. While
274 traditional high-throughput wet lab screens are often resource-intensive, AI in silico
275 methods can evaluate billions of compounds more rapidly. However, the effectiveness
276 of these in silico screenings heavily relies on the quality and maturity of the underlying
277 machine learning models, such as the accuracy of training data and model
278 generalization to avoid common pitfalls like overfitting or data biases(50, 51).
279 Quantitative structure–activity relationship (QSAR) models play a central role in this,
280 with performance that is nuanced and contingent upon the available data regime. In
281 high-data regimes, modern approaches utilizing GNNs and transformers excel at
282 learning complex molecular features from structural information, eliminating the need
283 for manual descriptor selection and improving prediction accuracy for large
284 datasets(52). For example, DL-QSAR models integrating molecular fingerprints and
285 GNN-derived features can accelerate antiviral activity prediction and prioritize natural
286 product analogs(53). In contrast, in low-data regimes—common in natural product
287 antiviral discovery—classic techniques such as tree-based models (e.g., random forests)
288 combined with circular fingerprints often perform comparably or better, offering



289 robustness, simplicity, and superior generalization with fewer samples(51, 54, 55).
290 Hybrid strategies blending classic and advanced methods may yield optimal results
291 across diverse scenarios, balancing computational efficiency and reliability.

292 In structure-based virtual screening, traditional physics-based molecular docking
293 faces prohibitive computational costs when traversing billion-scale chemical spaces,
294 necessitating AI-driven strategies that optimize both accuracy and throughput while
295 addressing limitations in generalization, especially for structurally complex NPs. Key
296 facets of this paradigm shift encompass enhancing evaluation precision through DL-
297 driven affinity estimation—exemplified by curvature-based GNNs such as CurvAGN,
298 which capture intricate 3D molecular geometries and multi-scale interactions to
299 mitigate systematic biases in classical empirical scoring functions(56). To surmount
300 scalability constraints of exhaustive docking, Deep Docking-type paradigms deploy DL
301 surrogate models trained on representative subsets to forecast docking scores for the
302 remainder of the library, thereby excluding over 99% of non-binders without explicit
303 physics-based computations—a capability indispensable for trillion-scale ultra-large
304 chemical library screening(57). This acceleration is further augmented by Active
305 Learning (AL) strategies, which reconfigure virtual screening into a dynamic
306 'sampling-docking-training-prediction' iterative loop(58). Through iterative selection
307 of the most informative compounds, AL frameworks identify top-tier leads while
308 requiring docking of less than 1% of the library, drastically mitigating resource
309 demands. Integrating these AI accelerators with robust error-control mechanisms, such
310 as the Conformal Prediction (CP) methodology, assures high sensitivity and reliability
311 in ligand discovery for challenging targets like G protein-coupled receptors
312 (GPCRs)(59). Although these innovations markedly expedite antiviral natural product
313 screening, persistent challenges include limited model generalization to novel scaffolds
314 and substantial computational infrastructure requirements.

315 AI also permits multi-target screening, interrogating interactions between
316 compounds and multiple viral proteins (i.e., polymerase and protease inhibitors), and



317 the imposition of broad-spectrum antivirals with tuned polypharmacology(60). The
318 synergistic use of generative models, CP-guided docking, and multi objective
319 optimization allows AI-driven virtual screening to accelerate hit discovery while
320 increasing chemical diversity when compared to brute-force approaches in terms of
321 both speed and lead quality.

322 2.2.2 AI-Enhanced Pharmacokinetics and Toxicity Prediction

323 The optimization of pharmacokinetic (PK) and toxicity profiles has traditionally
324 represented a high-attrition, late-stage bottleneck in drug development; however, AI is
325 increasingly transforming this process into an early-stage, parallel Multi-Parameter
326 Optimization (MPO) paradigm. Rather than treating ADMET (Absorption, Distribution,
327 Metabolism, Excretion, and Toxicity) as a sequential experimental filter, multi-task
328 learning algorithms now enable the simultaneous prediction of diverse drug-likeness
329 endpoints from a single molecular representation (Fig. 3)(61, 62). This holistic
330 evaluation facilitates the navigation of complex trade-offs—such as balancing oral
331 bioavailability against cardiotoxicity risks—during the hit-to-lead transition(63, 64).
332 Furthermore, the integration of transfer learning addresses data sparsity challenges in
333 natural product research by adapting animal-derived datasets to human-specific
334 predictions, while interpretability tools like SHAP (SHapley Additive exPlanations)
335 offer mechanistic insights by identifying toxicophoric substructures for targeted
336 medicinal chemistry modifications(65, 66).

337 The industrial applicability of these AI-driven workflows is supported by
338 documented improvements in throughput and predictive accuracy. For instance,
339 platforms such as ADMETLab 3.0, which employ directed message-passing neural
340 networks, have shown the ability to evaluate over 119 endpoints with AUROC values
341 up to 0.94, contributing to enhanced efficiency in computational screening compared
342 to traditional empirical models(67, 68). Regulatory horizon-scanning reports, such as
343 those from the European Medicines Agency (EMA), underscore the deployment of



344 tools like Toxometris.ai and the Deep-PK framework in toxicity and pharmacokinetic
 345 predictions and preclinical study designs within pharmaceutical pipelines(59, 69).
 346 These implementations suggest that AI can enhance predictive capabilities and
 347 potentially de-risk the development of natural antivirals by supporting a more
 348 streamlined transition from hit identification to clinical candidate nomination.
 349 Nonetheless, ongoing challenges in model generalization to novel scaffolds highlight
 350 the need for rigorous validation and hybrid approaches to ensure translational
 351 reliability(59, 61).

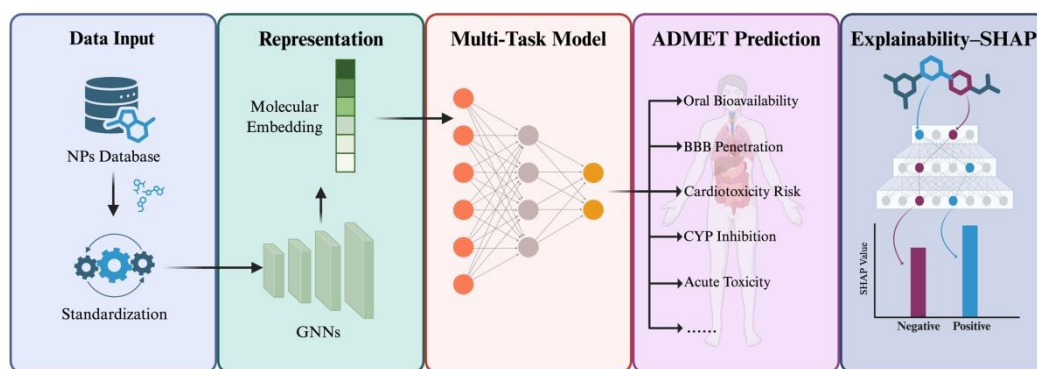


Fig. 3 Schematic workflow of AI-driven multi-task learning for ADMET prediction in natural product drug development. The workflow starts with data input from NPs databases. Data standardization follows, normalizing formats (e.g., SMILES) and curating biases. GNNs then generate molecular embeddings by encoding compounds as graphs, capturing NP-specific features like macrocycles. A multi-task model predicts interrelated endpoints simultaneously, such as oral bioavailability, blood-brain barrier (BBB) penetration (e.g., via permeability coefficients), cardiotoxicity risk (e.g., hERG channel inhibition), cytochrome P450 (CYP) enzyme inhibition (e.g., CYP3A4 isoform specificity), and acute toxicity (e.g., LD50 estimates), using shared representations for efficiency. SHAP, a game-theoretic interpretability framework based on Shapley values from cooperative game theory, quantifies feature contributions (positive or negative) to each prediction,



visualizing substructure impacts (e.g., highlighting aromatic rings contributing to hepatotoxicity) to guide targeted structural modifications in lead optimization. Figure created with BioRender.com.

View Article Online
DOI: 10.1039/D5DD000504C

352 2.2.3 Predicting and Mitigating Antiviral Resistance

353 Viruses evolve rapidly, posing a persistent challenge to antiviral drug efficacy as
354 resistance mechanisms emerge, often rendering treatments obsolete within months or
355 years. Traditional approaches rely on reactive surveillance and empirical testing, but AI
356 introduces a proactive paradigm by forecasting evolutionary trajectories, identifying
357 resistance signatures, and guiding resilient inhibitor design.

358 AI models leverage sequence data and structural predictions to anticipate viral
359 mutations at binding sites, enabling the development of inhibitors that maintain potency
360 against future variants. For example, EVEscape computationally produced multi-
361 mutant SARS-CoV-2 spikes to replicate immune escape, with experimental
362 validation(70). Predicting mutations at binding sites can guide design of inhibitors
363 resilient to future variants. AI may also search sequence databases for resistance
364 signatures and offer chemical modifications to bypass common mechanisms(71). These
365 methods also facilitate chemical modifications, using generative AI to suggest scaffold
366 alterations that bypass common resistance pathways, shifting from post-resistance
367 response to preemptive antiviral engineering.

368 Despite these advances, significant dilemmas arise from the interplay of viral
369 biology and technological limitations. Viruses' high mutation rates create out-of-
370 distribution challenges for AI models, where training on historical variants may fail to
371 generalize to novel, phylogenetically divergent strains, leading to inaccurate
372 predictions and false confidence in drug resilience(46). Computationally, scaling
373 simulations for multi-mutant landscapes demands immense resources, often exceeding
374 available infrastructure and introducing biases from incomplete datasets. NPs' structural
375 diversity offers a vast pool for discovering new scaffolds less prone to resistance, but



376 integrating this with AI remains underdeveloped due to data scarcity in NP-specific
377 resistance profiles.

378 **2.2.4 Generative AI for Novel Drug Design and Synthesis Planning**

379 Generative AI serves as an impressive midstream approach in antiviral NPs drug
380 development, supporting the de novo generation of molecular structures inspired by
381 NPs' structural diversity and multi-target synergies(72). These models process chemical
382 patterns from large datasets, such as general chemical databases like ChEMBL (which
383 includes NPs) or NP-focused libraries like TCMSP and COCONUT, to suggest
384 candidates that could address viral evolution and resistance issues(73). Generative
385 frameworks commonly include variational autoencoders (VAEs), generative
386 adversarial networks (GANs), and diffusion models, often implemented with sequence-
387 based backbones (e.g., recurrent neural networks (RNNs) or Transformers) or graph-
388 based structures, allowing exploration of chemical spaces beyond conventional
389 screening and potentially aiding in timeline efficiency(74). These major generative AI
390 paradigms can be summarized as follows (**Fig. 4**). They efficiently generate molecular
391 structures inspired by natural products through mechanisms such as sequence modeling,
392 latent space exploration, diffusion processes, or evolutionary optimization. Notably,
393 generative AI generally produces de novo molecules that are NP analogs or mimics,
394 often struggling to replicate the intricate structural features of authentic NPs, such as
395 high chirality centers or complex ring systems, yet it can integrate NP-like motifs (e.g.,
396 macrocycles or polyketide scaffolds) to approximate bioactivity benefits like
397 polypharmacology, as discussed in recent AI-NP literature(75). This method aims to
398 align natural-inspired designs with synthetic feasibility, facilitating multi-objective
399 considerations for aspects like binding affinity, ADMET properties, and accessibility
400 in antiviral settings.

401 Examining individual architectures, VAEs project molecular structures into latent
402 spaces for sampling NP-inspired analogs, supporting the development of structurally



403 complex antiviral candidates, though accurately mirroring full NP polypharmacology
 404 can be difficult(76). GANs utilize adversarial training to yield plausible NP-like
 405 molecules, such as analogs of artemisinin for possible broad-spectrum antiviral
 406 exploration, via generator-discriminator refinement(75). Diffusion models construct
 407 full 3D molecular conformations by progressively denoising atomic coordinates and
 408 atom types in an equivariant manner, offering advantages in generating geometrically
 409 accurate NP-mimetic structures that might counter viral mutations(77). RNNs and
 410 Transformers, as sequence-based backbones, manage formats like SMILES strings to
 411 produce varied analogs, enabling investigation of NP-adjacent chemical spaces for
 412 antiviral purposes(78). GAs, meanwhile, simulate evolution via selection and mutation
 413 of known structures, embedding synthetic accessibility measures to optimize NP
 414 analogs(74). By linking these architectures with retrosynthesis planning tools like
 415 AiZynthFinder, generative AI can propose candidate structures alongside practical
 416 synthesis routes, encouraging an iterative design-make-test process that broadens
 417 chemical diversity in line with NP bioactivity concepts, although validation against
 418 diverse viral variants remains a key limitation(79).

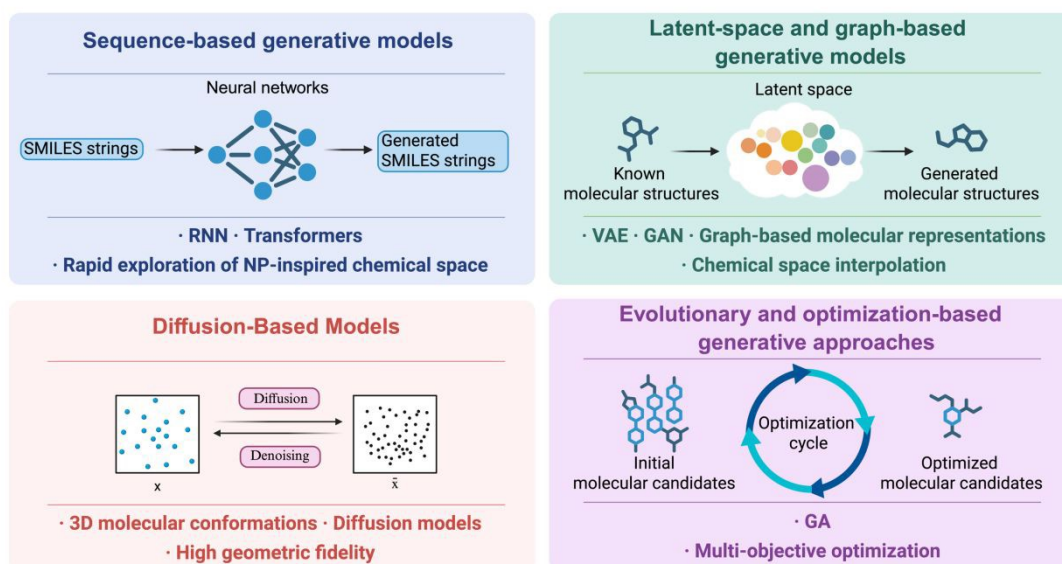


Fig. 4 Overview of major generative AI paradigms applied to natural product–inspired molecular design. Sequence-based generative models treat molecular



generation as a language modeling problem, in which molecular structures are encoded as SMILES strings and processed by sequence-based neural networks to rapidly explore NP-inspired chemical space. Latent-variable and graph-based generative models embed known molecular structures into a continuous latent space using latent representations and graph-based molecular encodings, enabling interpolation and sampling to generate novel molecular structures. Diffusion-based generative models generate molecular structures by progressively denoising stochastic representations, allowing the reconstruction of high-fidelity three-dimensional molecular conformations. Evolutionary and optimization-based generative approaches iteratively refine molecular candidates through mutation and selection under multi-objective optimization criteria, explicitly incorporating considerations such as chemical properties and synthetic feasibility. Figure created with BioRender.com.

419 **2.3 Downstream: Preclinical and Clinical Applications**

420 **2.3.1 AI-Driven Preclinical Testing**

421 Before advancing to clinical trials, antiviral natural products must undergo rigorous
422 preclinical evaluation, including *in vitro* assays, animal studies, and toxicity
423 assessments. AI has been integrated into these stages to enhance efficiency, but its
424 application in NPs—characterized by structural complexity and multi-target
425 interactions—presents unique challenges. This subsection provides a systematic
426 overview of AI methods in preclinical testing for antiviral NPs, focusing on
427 experimental optimization, automated data analysis, and predictive modeling, while
428 critically discussing limitations and real-world impacts.

429 First, AI facilitates experimental design through optimization algorithms, such as
430 Bayesian optimization (BO) and active learning. BO iteratively selects experimental
431 parameters (e.g., compound dosages or combinations) based on prior data to maximize



432 information gain with minimal trials. In the context of antiviral NPs, BO has been
433 applied to refine in vitro assays for plant-derived compounds like artemisinin analogs
434 against broad-spectrum antivirals including coronaviruses, achieving comparable
435 enrichment factors while reducing the experimental footprint by approximately 40-50%
436 compared to traditional methods(80). Similarly, active learning frameworks combine
437 machine learning with high-throughput screening to prioritize NPs from microbial
438 sources, such as polyketides discovered through genome mining with anti-viral
439 potential(81). These approaches accelerate preclinical workflows by intelligently
440 navigating the vast chemical space of NPs.

441 Second, AI enhances automated analysis of preclinical data, particularly in imaging
442 and histopathology. DL models, such as CNNs, automate the scoring of pathology
443 slides from animal models. For instance, the CSGO (Cell Segmentation with Globally
444 Optimized boundaries) pipeline has recently been validated for high-throughput,
445 whole-cell segmentation in Hematoxylin-and-Eosin (H&E)-stained tissues, enabling
446 precise quantification of inflammatory cell infiltration in lung injury models(82). Such
447 methodologies provide objective, high-resolution readouts for evaluating the
448 therapeutic efficacy of NP candidates, such as quercetin, in alleviating virus-induced
449 pulmonary inflammation. This automation significantly reduces inter-observer
450 variability and ensures the reproducibility of toxicity and efficacy assessments for
451 natural antivirals in preclinical animal trials.

452 Third, AI-driven predictive modeling, including physiologically based PBPK
453 models augmented with ML, bridges in vitro data to in vivo outcomes. These models
454 predict pharmacokinetics (e.g., clearance, half-life) for NPs by integrating
455 physicochemical properties and multi-omics data. A notable application involves
456 forecasting human systemic exposure for plant-derived antivirals like berberine. For
457 instance, mechanistic PBPK models have been refined to capture complex interactions
458 between berberine and multiple transporters (e.g., P-gp and OCTs), providing a
459 quantitative framework to evaluate its therapeutic potential against viral infections,



460 particularly in the context of drug-drug interactions (DDI)(83). Furthermore, ML-
461 enhanced PBPK has demonstrated significant superiority, reducing prediction errors
462 (e.g., RMSE) by 20–30% compared to empirical scaling methods, as shown in small-
463 molecule validations(84). Although primarily validated in targeted therapeutics like
464 oligonucleotides(85), the transfer learning strategies utilized for cross-species
465 translation are increasingly being adapted to NPs, facilitating more accurate animal-to-
466 human extrapolations for complex natural compounds and offering methodological
467 insights for antiviral drug development.

468 However, despite these advancements, AI in preclinical testing for antiviral NPs
469 faces significant failure modes. Recent high-profile blind challenges, such as the
470 CACHE series and the ASAP-Polaris initiatives, have exposed a stark reality: the
471 majority of AI-prioritized compounds fail to exhibit reproducible activity in wet-lab
472 assays, with failure rates exceeding 90% in many cases(86, 87). During the COVID-19
473 pandemic, the “rush to screen” led to an influx of low-quality in silico studies where
474 herbal compounds like quercetin were frequently identified as hits. Retrospective
475 analyses now confirm these were largely false positives—often due to the molecules
476 being Pan-Assay Interference Compounds (PAINS) or the AI overestimating binding
477 affinities by neglecting complex solvation effects and structural plasticity(87). These
478 failures underscore a systemic translational gap, necessitating a shift toward more
479 rigorous, physics-informed AI models for natural products.

480 2.3.2 AI Innovations in Clinical Trials

481 Antiviral natural products drug, like other therapeutics, undergo expensive and lengthy
482 clinical trial phases. Artificial intelligence is beginning to offer targeted improvements
483 in efficiency and informativeness, though many applications remain in early stages or
484 face substantial limitations in real-world translation. Several AI-based methodologies
485 are under exploration in Phase I–III studies for antiviral agents:



486 (i) Federated learning (FL) for multi-center data integration — FL enables
487 collaborative model training across institutions without centralizing sensitive patient
488 data, thereby preserving privacy while leveraging diverse cohorts. A landmark example
489 is its application to predict clinical outcomes in COVID-19 patients from multiple
490 hospitals, demonstrating improved generalizability across heterogeneous
491 populations(88). Although this approach has not yet been widely reported for antiviral
492 NP trials, it holds methodological promise for modeling disease trajectories or
493 treatment responses in multi-ethnic or multi-risk-group settings, provided data
494 harmonization and model robustness challenges are addressed.

495 (ii) NLP for unstructured clinical data — NLP algorithms can extract relevant
496 outcomes, adverse events, or symptom patterns from electronic health records,
497 physician notes, and free-text entries in near real-time. Systematic reviews indicate that
498 NLP enhances signal detection in clinical decision support and could support more
499 responsive monitoring in trials(89). In the context of antiviral NPs, such tools may
500 facilitate earlier identification of efficacy or safety signals in adaptive trial designs,
501 although current implementations are largely limited to general medical contexts rather
502 than NP-specific endpoints.

503 (iii) Synthetic control arms and generative models — Generative AI methodologies,
504 particularly GANs, is being explored to create synthetic patient data (SPD) to augment
505 or partially replace traditional control arms. Early benchmarks using tabular clinical
506 datasets demonstrate that GAN-based frameworks, such as GANerAid, can synthesize
507 patient-level records that preserve the complex statistical correlations and longitudinal
508 trajectories of actual trial participants(90). This approach is particularly promising for
509 establishing synthetic control arms in rare disease or oncology trials, where simulating
510 survival data (e.g., progression-free survival) can reduce the reliance on large placebo
511 cohorts while maintaining statistical power(91). However, applications to antiviral NPs
512 remain nascent, with significant challenges in ensuring biological plausibility—the risk
513 that GANs may generate pharmacologically impossible trajectories for multi-



514 component natural extracts. Furthermore, regulatory acceptance requires rigorous
515 validation against "hallucinated" correlations and the avoidance of bias amplification
516 in complex patient profiles.

517 **3. Bottlenecks In Ai-Enabled Antiviral NPs Development**

518 **3.1 Data constraints: NPs complexity and viral dynamics**

519 Data scarcity and heterogeneity are the key obstacles. NPs data are multi-modal,
520 complex, and unevenly standardized across sources and formats(92). High-quality,
521 empirically validated annotations for structures, biosynthetic processes, antiviral
522 activity (including negatives), and toxicity are limited. Multi-target/network
523 pharmacology mechanisms are tough to capture adequately, hindering mechanistic
524 modeling. Knowledge integrated in traditional medicine literature is valuable; however,
525 it is unstructured, archaic, and philosophically complex, challenging NLP and
526 knowledge-graph development(93).

527 Viral evolution data are limited in sample size and dynamic. For new strains or
528 unexpected resistance mutations, establishing robust models is challenging(94). Rapid
529 genomic change needs frequent model updating(46). Many viruses replicate within
530 liquid–liquid phase-separated (LLPS) condensates, which lack well-defined structural
531 targets, complicating design; changing sequences may affect condensate
532 physicochemical features and dynamics(95). Host-target inhibition may generate
533 compensating responses that current AI rarely can foresee(49). Addressing such
534 “moving targets” requires models that spot dynamic or allosteric sites or even intervene
535 in phase behavior, potentially via system-level “digital twins”—well beyond existing
536 capabilities.

537 **3.2 Model limitations: from pattern recognition to biological understanding**



538 AI's advancements primarily rely on pattern recognition, but deep biological and
539 chemical understanding for intricate NPs structures and polypharmacology remains
540 inadequate. NPs feature intricate ring structures, multiple chiral centers, and flexible
541 conformations; capturing fine-grained conformational dynamics and flexible protein–
542 ligand interactions is difficult(76). Even employing AlphaFold 3, difficulties exist in
543 modeling dynamics, allostery, and binding of unusual ligands, which can affect docking
544 accuracy(96). NPs frequently act via several viral and host targets, forming complex
545 networks; current AI struggles to interpret synergy or antagonism or predict system-
546 level ramifications. Multi-component traditional formulations confront extra “black-
547 box” issues that limit optimization and sensible combination design(97). For swiftly
548 growing viruses, generalization to out-of-distribution variations is limited; even single-
549 point mutations might have structural ramifications that models fail to capture(98).

550 3.3 Synthetic accessibility

551 Bridging computational hits to physical molecules remains a formidable hurdle for NPs.
552 While NPs benefit from innate bioactivity and initial accessibility through extraction or
553 fermentation, scaling their production to meet clinical or industrial demands reveals
554 significant bottlenecks. Direct large-scale extraction from natural sources is often
555 unsustainable due to resource scarcity, seasonal variability, environmental degradation
556 (e.g., overharvesting), and extremely low yields (frequently below 0.01% dry
557 weight)(6). Similarly, native microbial fermentation is hampered by the low
558 productivity of wild-type strains and difficult-to-control fermentation conditions,
559 though modern strain improvement and optimization techniques (e.g., fed-batch or
560 continuous fermentation) can enhance productivity(99).

561 Heterologous biosynthesis has emerged as a transformative strategy, but metabolic
562 engineering still faces a "scale-up gap," where titer optimization often stalls at the mg/L
563 level, failing to reach the g/L threshold required for commercial viability(100). For
564 molecules with extreme structural complexity, semi-synthesis offers a middle ground,



565 yet it remains constrained by limited reagent-accessible space and high operation
566 costs(101).

567 **4. Outlook And Future Directions**

568 **4.1 High-quality, multi-modal benchmark datasets:**

569 To overcome data scarcity and fragmented evaluation, the field needs FAIR (Findable,
570 Accessible, Interoperable, Reusable), high-quality, multi-modal benchmark datasets.
571 These should integrate NPs structures, biosynthetic pathways, multi-dimensional
572 bioactivity profiles (including negative results), and ADMET properties, along with
573 standardized metrics for fair model comparison and industrial translation.

574

575 **4.2 Model and algorithmic innovation**

576 Beyond correlational prediction, future AI should prioritize interpretability (e.g.,
577 explainable AI, XAI) and robust out-of-distribution generalization to novel biological
578 systems (e.g., new variants). Incorporating causal inference and neuro-symbolic
579 techniques may facilitate a move from pattern recognition to scientific discovery,
580 yielding better mechanistic understanding.

581

582 **4.3 Stronger experimental closed loops and high-throughput validation**

583 To match AI's high-throughput hypothesis development, experimental validation must
584 accelerate. Automated platforms (self-driving laboratories) and active learning can
585 connect in silico forecasts tightly with wet-lab feedback during DBTL cycles,
586 eliminating the compute-experiment gap.

587

588 **4.4 Cross-disciplinary collaboration and ecosystem building**

589 Antiviral NPs discovery needs fundamental integration across chemistry, biology,
590 medicine, and data science. Training cross-domain competence and promoting
591 academia-industry-international partnerships will enable secure data exchange (e.g.,



592 privacy-preserving federated learning). Building an automated and intelligent end-to-
593 end pipeline—from data collection and modeling to candidate generation and
594 validation—will require continual investment but is essential to realize AI's full
595 promise in antiviral NPs discovery.

596 **5. Conclusion**

597 Research on antiviral drugs today faces significant difficulties: expensive, time-
598 consuming, and limited against rapidly evolving viruses. NPs are still hindered in their
599 analysis, target detection, and synthesis because of their structural diversity and multi-
600 target synergy. AI is a powerful new paradigm that is changing how the drug discovery
601 process is conducted and through the integration of data and algorithmic learning.

602 In the upstream phase, AI helps to refine genome mining and the modeling of
603 biosynthetic gene clusters, while knowledge graphs and natural language processing
604 aid in extracting insight from conventional medical databases. AI-driven virtual
605 screening and generative design during the midstream phase significantly expand
606 chemical space exploration, yielding multi-target antiviral candidates with optimal
607 pharmacokinetics and resilience against resistance. On the downstream, AI optimizes
608 preclinical models and clinical trial efficiency using techniques such as federated
609 learning and synthetic control arms, which accelerates translational outcomes.

610 There are many critical issues, including data scarcity, limited model
611 interpretability, and synthetic accessibility challenges, necessitating integrated
612 solutions to advance the field. The complexity of NPs and viral dynamics demands
613 FAIR, multi-modal benchmark datasets that unify structural, biosynthetic, and
614 bioactivity data, while interpretable AI and causal inference are critical to move beyond
615 pattern recognition toward mechanistic understanding of polypharmacology and
616 dynamic targets like liquid-liquid phase-separated condensates. Automated platforms,
617 such as self-driving laboratories, must bridge the compute-experiment gap through



618 high-throughput Design-Build-Test-Learn cycles, embedding synthetic feasibility early
619 in the design process. By fostering cross-disciplinary collaboration and privacy-
620 preserving data-sharing ecosystems, the field can build an intelligent, end-to-end
621 pipeline, transforming reactive antiviral NP discovery into a proactive, resilient strategy
622 against evolving viral threats.

623 Although not a cure all, AI represents the best catalyst for converting NP-based
624 antiviral discovery from empirical serendipity to predictive engineering. If adopted
625 more widespread—as long as there's solid validation and international data-sharing
626 agreements—the virus's arsenal of antiviral responses may be unleashed. We are on the
627 cusp of the fourth drug discovery revolution enabled by symbiotic human-AI
628 collaboration, where the speed of machines meets the smarts of chemical evolution to
629 move beyond the capacity for viral adaptation.

630 **Acknowledgments**

631 We would like to express our profound gratitude to all participants for their invaluable
632 contributions to this research. Additionally, we appreciate the assistance of
633 BioRender.com in creating the figure for this study.

634 **Author Contributions**

635 **Conceptualization:** Jianjian Ji, Peng Cao, Liangyu Cai, Junxi Song, Kunhuan Yang,
636 Yingcai Xiong, and Keyu Tao.

637 **Methodology:** Junxi Song, Kunhuan Yang, Yingcai Xiong, and Keyu Tao.

638 **Investigation:** Junxi Song, Kunhuan Yang, Yingcai Xiong, and Keyu Tao.

639 **Writing—Original Draft Preparation:** Jianjian Ji, Peng Cao, Liangyu Cai, Junxi
640 Song, Kunhuan Yang, Yingcai Xiong, and Keyu Tao.

641 **Visualization:** Junxi Song, Kunhuan Yang, Yingcai Xiong, and Keyu Tao.

642 **Writing—Review & Editing:** Jianjian Ji, Peng Cao, and Liangyu Cai.

643 **Funding Acquisition:** Jianjian Ji, Peng Cao, and Liangyu Cai.



644 **Supervision:** Jianjian Ji, Peng Cao, and Liangyu Cai.

645 **Project Administration:** Jianjian Ji, Peng Cao, and Liangyu Cai.

646 **References**

- 647 1. Abdelrahman Z, Li M, Wang X. Comparative Review of SARS-CoV-2, SARS-CoV, MERS-CoV,
648 and Influenza A Respiratory Viruses. *Front Immunol.* 2020;11:552909.
- 649 2. DiMasi JA, Grabowski HG, Hansen RW. Innovation in the pharmaceutical industry: New
650 estimates of R&D costs. *J Health Econ.* 2016;47:20-33.
- 651 3. Eastman RT, Roth JS, Brimacombe KR, Simeonov A, Shen M, Patnaik S, et al. Remdesivir:
652 A Review of Its Discovery and Development Leading to Emergency Use Authorization for
653 Treatment of COVID-19. *ACS Cent Sci.* 2020;6(5):672-83.
- 654 4. Cihlar T, Mackman RL. Journey of remdesivir from the inhibition of hepatitis C virus to the
655 treatment of COVID-19. *Antivir Ther.* 2022;27(2):13596535221082773.
- 656 5. Guarnaccia T, Carolan LA, Maurer-Stroh S, Lee RT, Job E, Reading PC, et al. Antigenic drift
657 of the pandemic 2009 A(H1N1) influenza virus in A ferret model. *PLoS Pathog.*
658 2013;9(5):e1003354.
- 659 6. Newman DJ, Cragg GM. Natural Products as Sources of New Drugs over the Nearly Four
660 Decades from 01/1981 to 09/2019. *J Nat Prod.* 2020;83(3):770-803.
- 661 7. Miller LH, Su X. Artemisinin: discovery from the Chinese herbal garden. *Cell.*
662 2011;146(6):855-8.
- 663 8. Li JY, Cao HY, Liu P, Cheng GH, Sun MY. Glycyrrhizic acid in the treatment of liver diseases:
664 literature review. *Biomed Res Int.* 2014;2014:872139.
- 665 9. De Clercq E. Antiviral drugs in current clinical use. *J Clin Virol.* 2004;30(2):115-33.
- 666 10. Rappé MS, Giovannoni SJ. The uncultured microbial majority. *Annu Rev Microbiol.*
667 2003;57:369-94.
- 668 11. Liu T, Lin S. Comprehensive characterization of the chemical constituents of Lianhua
669 Qingwen capsule by ultra high performance liquid chromatography coupled with Fourier
670 transform ion cyclotron resonance mass spectrometry. *Heliyon.* 2024;10(6).
- 671 12. Zhang J, Hansen LG, Gudich O, Viehrig K, Lassen LMM, Schruebbers L, et al. A microbial
672 supply chain for production of the anti-cancer drug vinblastine. *Nature.* 2022;609(7926):341-
673 +.



- 674 13. Zhu ZQ, Zheng X, Qi GQ, Gong YF, Li YY, Mazur N, et al. Drug-target binding affinity
675 prediction model based on multi-scale diffusion and interactive learning. *Expert Systems with*
676 *Applications*. 2024;255.
- 677 14. Zdrzil B, Felix E, Hunter F, Manners EJ, Blackshaw J, Corbett S, et al. The ChEMBL
678 Database in 2023: a drug discovery platform spanning multiple bioactivity data types and time
679 periods. *Nucleic Acids Research*. 2023;52(D1):D1180-D92.
- 680 15. Ru J, Li P, Wang J, Zhou W, Li B, Huang C, et al. TCMSP: a database of systems
681 pharmacology for drug discovery from herbal medicines. *Journal of Cheminformatics*. 2014;6.
- 682 16. Darmawan RD, Kusuma WA, Rahmawan H. Deep learning optimization for drug-target
683 interaction prediction in COVID-19 using graphic processing unit. *International Journal of*
684 *Electrical and Computer Engineering (IJECE)*. 2023;13(3):3111-23.
- 685 17. Chen Y, Luo D, Xue W. Deep Learning for Drug-Target Interaction Prediction: A
686 Comprehensive Review. *Chemical Biology & Drug Design*. 2025;106(4):e70183.
- 687 18. Wang Y, Ünlü A, Wang X, Çevrim E, Offermans DM, Flesseman MP, et al. AI-driven
688 discovery of antiretroviral drug bictegavir and etravirine as inhibitors against monkeypox and
689 related poxviruses. *Communications Biology*. 2025;8(1):1734.
- 690 19. Mishra AK, Sudalaimuthasari N, Hazzouri KM, Saeed EE, Shah I, Amiri KMA. Tapping into
691 Plant-Microbiome Interactions through the Lens of Multi-O mics Techniques.
692 *Cells*.11(20):3254.
- 693 20. Zhang J, Li B, Qin Y, Karthik L, Zhu G, Hou C, et al. A new abyssomicin polyketide with anti-
694 influenza A virus activity from a marine-derived *Verrucosipora* sp. MS100137. *Appl Microbiol*
695 *Biotechnol*.104(4):1533-43.
- 696 21. Li H, Yang L, Liu F-F, Ma X-N, He P-L, Tang W, et al. Overview of therapeutic drug research
697 for COVID-19 in China. *Acta Pharmacol Sin*.41(9):1133-40.
- 698 22. Ge F, Yang Y, Bai Z, Si L, Wang X, Yu J, et al. The role of Traditional Chinese medicine in
699 anti-HBV: background, progress, and challenges. *Chin Med*.18(1):159.
- 700 23. Ma L, Ji L, Wang T, Zhai Z, Su P, Zhang Y, et al. Research progress on the mechanism of
701 traditional Chinese medicine regulating intestinal microbiota to combat influenza a virus
702 infection. *Virol J*.20(1):260.
- 703 24. Medema MH, Blin K, Cimermancic P, de Jager V, Zakrzewski P, Fischbach MA, et al.
704 antiSMASH: rapid identification, annotation and analysis of secondary metabolite
705 biosynthesis gene clusters in bacterial and fungal genome sequences. *Nucleic Acids*
706 *Res*.39(Web Server issue):W339-46.



- 707 25. Hannigan GD, Prihoda D, Palicka A, Soukup J, Klempir O, Rampula L, et al. A deep learning
708 genome-mining strategy for biosynthetic gene cluster prediction. *Nucleic Acids Research*.
709 2019;47(18).
- 710 26. Zhang Y, Zhang J, Li M, Qiao Y, Wang W, Ma L, et al. Target discovery of bioactive natural
711 products with native-compound-coupled CNBr-activated Sepharose 4B beads (NCCB):
712 Applications, mechanisms and outlooks. *Bioorg Med Chem*.96:117483.
- 713 27. Xu Y, Cao L, Chen Y, Zhang Z, Liu W, Li H, et al. Integrating Machine Learning in
714 Metabolomics: A Path to Enhanced Diagnostics and Data Interpretation. *Small Methods*.
715 2024;8(12):e2400305.
- 716 28. Zhang Z, Yang H, Wang Y, Zhang L, Lin SH. QuanFormer: A Transformer-Based Precise
717 Peak Detection and Quantification Tool in LC-MS-Based Metabolomics. *Anal Chem*.
718 2025;97(5):2698-706.
- 719 29. Hu G, Qiu M. Machine learning-assisted structure annotation of natural products based
720 on MS and NMR data. *Natural Product Reports*. 2023;40(11):1735-53.
- 721 30. Schmid R, Heuckeroth S, Korf A, Smirnov A, Myers O, Dyrland TS, et al. Integrative analysis
722 of multimodal mass spectrometry data in MZmine 3. *Nature biotechnology*. 2023;41(4):447-
723 9.
- 724 31. Chhetri BK, Tedbury PR, Sweeney-Jones AM, Mani L, Soapi K, Manfredi C, et al. Marine
725 natural products as leads against SARS-CoV-2 infection. *Journal of natural products*.
726 2022;85(3):657-65.
- 727 32. Song JG, Ye WC, Wang Y. Advanced crystallography for structure determination of natural
728 products. *Nat Prod Rep*. 2025;42(3):429-42.
- 729 33. Dymura SA, Viniichuk OO, Melnykov KP, Radchenko DS, Grygorenko OO. Machine
730 Learning-Based Retention Time Prediction Tool for Routine LC-MS Data Analysis. *J Chem Inf*
731 *Model*. 2025;65(14):7415-25.
- 732 34. Liu Y, Yoshizawa AC, Ling Y, Okuda S. Insights into predicting small molecule retention
733 times in liquid chromatography using deep learning. *J Cheminform*. 2024;16(1):113.
- 734 35. Song D, Tang T, Wang R, Liu H, Xie D, Zhao B, et al. Enhancing compound confidence in
735 suspect and non-target screening through machine learning-based retention time prediction.
736 *Environ Pollut*. 2024;347:123763.
- 737 36. Peng C, Xia F, Naseriparsa M, Osborne F. Knowledge Graphs: Opportunities and
738 Challenges. *Artif Intell Rev*.1-32.

View Article Online
DOI: 10.1039/D5DD00504C



- 739 37. Gaudry A, Pagni M, Mehl F, Moretti S, Quiros-Guerrero L-M, Cappelletti L, et al. A sample
740 centric and knowledge-driven computational framework for natural products drug discovery.
741 ACS Publications; 2024.
- 742 38. Leão TF, Wang M, da Silva R, Gurevich A, Bauermeister A, Gomes PWP, et al. NPOmix: a
743 machine learning classifier to connect mass spectrometry fragmentation data to biosynthetic
744 gene clusters. PNAS nexus. 2022;1(5):pgac257.
- 745 39. Lv Q, Chen G, He H, Yang Z, Zhao L, Chen H-Y, et al. TCMBank: bridges between the largest
746 herbal medicines, chemical ingredients, target proteins, and associated diseases with
747 intelligence text mining. Chemical science. 2023;14(39):10684-701.
- 748 40. Galati S, Di Stefano M, Martinelli E, Poli G, Tuccinardi T. Recent Advances in In Silico Target
749 Fishing. Molecules. 2021;26(17).
- 750 41. Sivanandan S, Leitmann B, Lubeck E, Sultan MM, Stanitsas P, Ranu N, et al. A pooled cell
751 painting CRISPR screening platform enables de novo inference of gene function by self-
752 supervised deep learning. Nature Communications. 2025.
- 753 42. Cockroft NT, Cheng X, Fuchs JR. STarFish: A Stacked Ensemble Target Fishing Approach
754 and its Application to Natural Products. J Chem Inf Model. 2019;59(11):4906-20.
- 755 43. Huang K, Fu T, Glass LM, Zitnik M, Xiao C, Sun J. DeepPurpose: a deep learning library for
756 drug-target interaction prediction. Bioinformatics. 2020;36(22-23):5545-7.
- 757 44. Kundu P, Beura S, Mondal S, Das AK, Ghosh A. Machine learning for the advancement of
758 genome-scale metabolic modeling. Biotechnology Advances. 2024;74:108400.
- 759 45. Hsieh K, Wang Y, Chen L, Zhao Z, Savitz S, Jiang X, et al. Drug repurposing for COVID-19
760 using graph neural network with genetic, mechanistic, and epidemiological validation.
761 Research square. 2020:rs. 3. rs-114758.
- 762 46. Lytras S, Lamb KD, Ito J, Grove J, Yuan K, Sato K, et al. Pathogen genomic surveillance and
763 the AI revolution. Journal of Virology. 2025;99(2):e01601-24.
- 764 47. Abramson J, Adler J, Dunger J, Evans R, Green T, Pritzel A, et al. Accurate structure
765 prediction of biomolecular interactions with AlphaFold 3. Nature. 2024;630(8016):493-500.
- 766 48. Krishna R, Wang J, Ahern W, Sturmfels P, Venkatesh P, Kalvet I, et al. Generalized
767 biomolecular modeling and design with RoseTTAFold All-Atom. Science.
768 2024;384(6693):eadl2528.
- 769 49. Abdelsayed M. AI-Driven Polypharmacology in Small-Molecule Drug Discovery. Int J Mol
770 Sci. 2025;26(14).



- 771 50. Bender A, Cortes-Ciriano I. Artificial intelligence in drug discovery: what is realistic, what
772 are illusions? Part 2: a discussion of chemical and biological data. *Drug Discov Today*.
773 2021;26(4):1040-52.
- 774 51. Deng J, Yang Z, Wang H, Ojima I, Samaras D, Wang F. A systematic study of key elements
775 underlying molecular property prediction. *Nature Communications*. 2023;14(1):6395.
- 776 52. Wen T, Cai X, Li J. Graph Neural Networks vs. Traditional QSAR: A Comprehensive
777 Comparison for Multi-Label Molecular Odor Prediction. *Molecules*. 2025;30(23):4605.
- 778 53. Cai H, Zhang H, Zhao D, Wu J, Wang L. FP-GNN: a versatile deep learning architecture for
779 enhanced molecular property prediction. *Brief Bioinform*.23(6):bbac408.
- 780 54. Praski M, Adamczyk J, Czech W. Benchmarking pretrained molecular embedding models
781 for molecular representation learning. *arXiv preprint arXiv:250806199*. 2025.
- 782 55. Dobbelaere MR, Lengyel I, Stevens CV, Van Geem KM. Geometric deep learning for
783 molecular property predictions with chemical accuracy across chemical space. *Journal of*
784 *Cheminformatics*. 2024;16(1):99.
- 785 56. Wu J, Chen H, Cheng M, Xiong H. CurvAGN: Curvature-based Adaptive Graph Neural
786 Networks for Predicting Protein-Ligand Binding Affinity. *BMC Bioinformatics*. 2023;24(1):378.
- 787 57. Gentile F, Agrawal V, Hsing M, Ton A-T, Ban F, Norinder U, et al. Deep docking: a deep
788 learning platform for augmentation of structure based drug discovery. *ACS central science*.
789 2020;6(6):939-49.
- 790 58. Graff DE, Shakhnovich EI, Coley CW. Accelerating high-throughput virtual screening
791 through molecular pool-based active learning. *Chemical science*. 2021;12(22):7866-81.
- 792 59. Luttens A, Cabeza de Vaca I, Sparring L, Brea J, Martínez AL, Kahlous NA, et al. Rapid
793 traversal of vast chemical space using machine learning-guided docking screens. *Nature*
794 *Computational Science*. 2025:1-12.
- 795 60. Zhang L-C, Zhao H-L, Liu J, He L, Yu R-L, Kang C-M. Design of SARS-CoV-2 Mpro, PLpro
796 dual-target inhibitors based on deep reinforcement learning and virtual screening. *Future*
797 *Med Chem*.14(6):393-405.
- 798 61. European Medicines Agency. Review of artificial intelligence and machine learning
799 applications in medicines lifecycle (2024): Horizon Scanning Short Report. In: Agency EM,
800 editor. Amsterdam, Netherlands: European Medicines Agency; 2025.
- 801 62. Zhang J, Li H, Zhang Y, Huang J, Ren L, Zhang C, et al. Computational toxicology in drug
802 discovery: applications of artificial intelligence in ADMET and toxicity prediction. *Briefings in*
803 *Bioinformatics*. 2025;26(5):bbaf533.



- 804 63. Lee H, Kim J, Kim J-W, Lee Y. Recent advances in AI-based toxicity prediction for drug
805 discovery. *Frontiers in Chemistry*. 2025;13:1632046.
- 806 64. Lavecchia A. Explainable artificial intelligence in drug discovery: bridging predictive power
807 and mechanistic insight. *Wiley Interdisciplinary Reviews: Computational Molecular Science*.
808 2025;15(5):e70049.
- 809 65. Pathan I, Raza A, Sahu A, Joshi M, Sahu Y, Patil Y, et al. Revolutionizing pharmacology: AI-
810 powered approaches in molecular modeling and ADMET prediction. *Medicine in Drug*
811 *Discovery*. 2025:100223.
- 812 66. Wang Q, Sun B, Yi Y, Velkov T, Shen J, Dai C, et al. Progress of AI-driven drug–target
813 interaction prediction and lead optimization. *International Journal of Molecular Sciences*.
814 2025;26(20):10037.
- 815 67. Ferreira FJ, Carneiro AS. AI-Driven Drug Discovery: A Comprehensive Review. *ACS omega*.
816 2025.
- 817 68. Fu L, Shi S, Yi J, Wang N, He Y, Wu Z, et al. ADMETlab 3.0: an updated comprehensive
818 online ADMET prediction platform enhanced with broader coverage, improved performance,
819 API functionality and decision support. *Nucleic acids research*. 2024;52(W1):W422-W31.
- 820 69. Myung Y, de Sá AG, Ascher DB. Deep-PK: deep learning for small molecule
821 pharmacokinetic and toxicity prediction. *Nucleic acids research*. 2024;52(W1):W469-W75.
- 822 70. Youssef N, Gurev S, Ghantous F, Brock KP, Jaimes JA, Thadani NN, et al. Computationally
823 designed proteins mimic antibody immune evasion in viral evolution. *Immunity*.
824 2025;58(6):1411-21.e6.
- 825 71. Sahibzada KI, Shahid S, Akhter M, Abid R, Azhar M, Hu Y, et al. HIV OctaScanner: A
826 Machine Learning Approach to Unveil Proteolytic Cleavage Dynamics in HIV-1 Protease
827 Substrates. *J Chem Inf Model*. 2025;65(2):640-8.
- 828 72. Bilodeau C, Jin W, Jaakkola T, Barzilay R, Jensen KF. Generative models for molecular
829 discovery: Recent advances and challenges. *Wiley Interdisciplinary Reviews: Computational*
830 *Molecular Science*. 2022;12(5):e1608.
- 831 73. Atanasov AG, Zotchev SB, Dirsch VM, Supuran CT. Natural products in drug discovery:
832 advances and opportunities. *Nature reviews Drug discovery*. 2021;20(3):200-16.
- 833 74. Nigam A, Pollice R, Aspuru-Guzik A. Parallel tempered genetic algorithm guided by deep
834 neural networks for inverse molecular design. *Digital Discovery*. 2022;1(4):390-404.
- 835 75. Meyers J, Fabian B, Brown N. De novo molecular design and generative models. *Drug*
836 *discovery today*. 2021;26(11):2707-15.



- 837 76. Ochiai T, Inukai T, Akiyama M, Furui K, Ohue M, Matsumori N, et al. Variational
838 autoencoder-based chemical latent space for large molecular structures with 3D complexity.
839 Communications Chemistry. 2023;6(1):249. [View Article Online
DOI: 10.1039/D5DD00504C](#)
- 840 77. Hoogeboom E, Satorras VG, Vignac C, Welling M, editors. Equivariant diffusion for
841 molecule generation in 3d. International conference on machine learning; 2022: PMLR.
- 842 78. Parrot M, Tajmouati H, da Silva VBR, Atwood BR, Fourcade R, Gaston-Mathé Y, et al.
843 Integrating synthetic accessibility with AI-based generative drug design. Journal of
844 Cheminformatics. 2023;15(1):83.
- 845 79. Westerlund AM, Saigiridharan L, Genheden S. Human-guided synthesis planning via
846 prompting. Chemical Science. 2025;16(32):14655-67.
- 847 80. Long T-Z, Jiang D-J, Shi S-H, Deng Y-C, Wang W-X, Cao D-S. Enhancing Multi-species Liver
848 Microsomal Stability Prediction through Artificial Intelligence. Journal of Chemical Information
849 and Modeling. 2024;64(8):3222-36.
- 850 81. Yan D, Zhou M, Adduri A, Zhuang Y, Guler M, Liu S, et al. Discovering type I cis-AT
851 polyketides through computational mass spectrometry and genome mining with Seq2PKS.
852 Nature Communications. 2024;15(1):5356.
- 853 82. Gu Z, Wang S, Rong R, Zhao Z, Wu F, Zhou Q, et al. Cell segmentation with globally
854 optimized boundaries (csgo): A deep learning pipeline for whole-cell segmentation in
855 hematoxylin-and-eosin-stained tissues. Laboratory Investigation. 2025;105(2):102184.
- 856 83. Adiwidjaja J, Boddy AV, McLachlan AJ. Physiologically based pharmacokinetic model
857 predictions of natural product-drug interactions between goldenseal, berberine, imatinib and
858 bosutinib. European Journal of Clinical Pharmacology. 2022;78(4):597-611.
- 859 84. Li Y, Wang Z, Li Y, Du J, Gao X, Li Y, et al. A combination of machine learning and PBPK
860 modeling approach for pharmacokinetics prediction of small molecules in humans.
861 Pharmaceutical Research. 2024;41(7):1369-79.
- 862 85. Derbalah A, Stader F, Liu C, Zyla A, Abdulla T, Wu Q, et al. Cross-species translational
863 modelling of targeted therapeutic oligonucleotides using physiologically based
864 pharmacokinetics. Journal of Pharmacokinetics and Pharmacodynamics. 2025;52(4):35.
- 865 86. CACHE WI, CHALLENGE JA. CRITICAL ASSESSMENT OF COMPUTATIONAL HIT-FINDING
866 EXPERIMENTS.
- 867 87. MacDermott-Opeskin H, Scheen J, Wognum C, Horton JT, West D, Payne AM, et al. A
868 computational community blind challenge on pan-coronavirus drug discovery data. 2025.



- 869 88. Dayan I, Roth HR, Zhong A, Harouni A, Gentili A, Abidin AZ, et al. Federated learning for
870 predicting clinical outcomes in patients with COVID-19. *Nature medicine*. 2021;27(10):1735-
871 43.
- 872 89. Eguia H, Sánchez-Bocanegra CL, Vinciarelli F, Alvarez-Lopez F, Saigí-Rubió F. Clinical
873 decision support and natural language processing in medicine: systematic literature review.
874 *Journal of Medical Internet Research*. 2024;26:e55315.
- 875 90. Krenmayr L, Frank R, Drobig C, Braungart M, Seidel J, Schaudt D, et al. GANerAid: realistic
876 synthetic patient data for clinical trials. *Informatics in Medicine Unlocked*. 2022;35:101118.
- 877 91. Akiya I, Ishihara T, Yamamoto K. Comparison of synthetic data generation techniques for
878 control group survival data in oncology clinical trials: simulation study. *JMIR medical*
879 *informatics*. 2024;12(1):e55118.
- 880 92. Mallowney MW, Duncan KR, Elsayed SS, Garg N, van der Hooft JJJ, Martin NI, et al.
881 Artificial intelligence for natural product drug discovery. *Nat Rev Drug Discov*.
882 2023;22(11):895-916.
- 883 93. Zhou L, Liu S, Li C, Sun Y, Zhang Y, Li Y, et al. Natural Language Processing Algorithms for
884 Normalizing Expressions of Synonymous Symptoms in Traditional Chinese Medicine. *Evid*
885 *Based Complement Alternat Med*. 2021;2021:6676607.
- 886 94. Mallapaty S. What will viruses do next? AI is helping scientists predict their evolution.
887 *Nature*. 2025;637(8046):527-8.
- 888 95. Galloux M, Longhi S. Unraveling Liquid-Liquid Phase Separation (LLPS) in Viral Infections
889 to Understand and Treat Viral Diseases. *Int J Mol Sci*. 2024;25(13).
- 890 96. Shen SY, Li JR, Wang YS, Li SN, Xu HE, He XH. An update for AlphaFold3 versus
891 experimental structures: assessing the precision of small molecule binding in GPCRs. *Acta*
892 *Pharmacol Sin*. 2025.
- 893 97. Yang L, Wang H, Zhu Z, Yang Y, Xiong Y, Cui X, et al. Network Pharmacology-Driven
894 Sustainability: AI and Multi-Omics Synergy for Drug Discovery in Traditional Chinese Medicine.
895 *Pharmaceuticals (Basel)*. 2025;18(7).
- 896 98. Wee J, Wei GW. Rapid response to fast viral evolution using AlphaFold 3-assisted
897 topological deep learning. *Virus Evol*. 2025;11(1):veaf026.
- 898 99. Liu X, Ding W, Jiang H. Engineering microbial cell factories for the production of plant
899 natural products: from design principles to industrial-scale production. *Microbial cell factories*.
900 2017;16(1):125.



901 100. Zhou H, Eun H, Lee SY. Systems metabolic engineering for the production of pharmaceutical natural products. *Current Opinion in Systems Biology*. 2024;37:100491. [View Article Online](#)
902 [DOI: 10.1039/D5DD00504C](#)

903 101. Shenvi RA. Natural product synthesis in the 21st century: Beyond the mountain top. *ACS*
904 *Central Science*. 2024;10(3):519-28.

905



Data Availability Statement

View Article Online
DOI: 10.1039/D5DD00504C

This article is a review, and as such, does not report any new primary data. All data and information discussed are available in the original publications and sources cited within the manuscript's references section

