# Language models in molecular discovery

Nikita Janakarajan[1], Tim Erdmann[2], Sarath Swaminathan[2], Teodoro Laino[1], and Jannis Born[1*]

[1]IBM Research Europe, Zurich, Switzerland
[2]IBM Research Almaden, San Jose, CA, United States
[*]Corresponding author: jab@zurich.ibm.com

September 29, 2023

## Abstract

The success of language models, especially transformer-based architectures, has trickled into other domains giving rise to "scientific language models" that operate on small molecules, proteins or polymers. In chemistry, language models contribute to accelerating the molecule discovery cycle as evidenced by promising recent findings in early-stage drug discovery. Here, we review the role of language models in molecular discovery, underlining their strength in de novo drug design, property prediction and reaction chemistry. We highlight valuable open-source software assets thus lowering the entry barrier to the field of scientific language modeling. Last, we sketch a vision for future molecular design that combines a chatbot interface with access to computational chemistry tools. Our contribution serves as a valuable resource for researchers, chemists, and AI enthusiasts interested in understanding how language models can and will be used to accelerate chemical discovery.

## 1 Introduction

Despite technological advances constantly reshaping our understanding of biochemical processes, the chemical industry persistently faces escalating resource costs of up to 10 years and 3 billion dollar per new market release [102]. The intricacy of the problem is typically attested by an exorbitant attrition rate in *in vitro* screenings [77], the sheer size of the chemical space [68] and the frequency of serendipity [40].

Language models (LMs) emerged recently and demonstrated an astonishing ability to understand and generate human-like text [65]. Machine learning (ML) in general and LMs in particular hold the potential to profoundly accelerate the molecular discovery cycle (see Figure 1). In this chapter, we explore applications of LMs to chemical design tasks. Although LMs were originally developed for natural language, they have shown compelling results in scientific discovery settings when applied to "scientific languages", e.g., in protein folding [55] or *de novo* design of small molecules [105], peptides [23] or polymers [66]. But what exactly is a language model? By definition, it is any ML model that consumes a sequence of text chunks (so-called tokens) and is capable to reason about the content of the sequence. Since each token is essentially a vector [62], a LM is a pseudo-discrete time series model. Most typically, LMs learn probability distributions over sequences of words thus also facilitating the generation of new text given some input, for example in a language translation task. While all LMs rely on neural networks, contemporary models almost exclusively leverage the Transformer architecture [93]. Now, all of this begs the question – what is the need for LMs in molecular discovery?

First, when applied to serializations of chemical entities (e.g., SMILES [98]), LMs can learn highly structured representations, often even tailored for desired functional properties [36]. This allows to perform smooth and property-driven exploration of the originally deemed discrete protein or molecular space. Another attractive feature of scientific LMs is their ability to seamlessly bridge natural and scientific languages. This can give rise to ChatGPT-style chatbot interfaces that allow chemists to
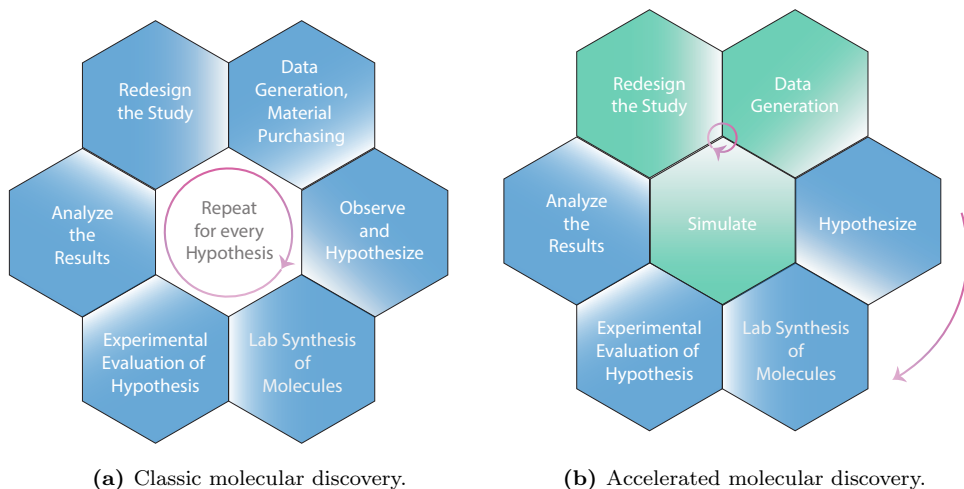
**(a)** Classic molecular discovery.　　　　**(b)** Accelerated molecular discovery.

**Figure 1:** A comparison of molecular discovery workflows: (a) classic approach, where each hypothesis (a.k.a. molecule) requires a new experimental cycle. (b) *Accelerated* molecular discovery cycle with machine-generated hypotheses and assisted validation, enabling simultaneous generation and testing of numerous molecules.

formulate their design objectives through natural language and to iteratively refine their result with an interactive agent thus potentially accomplishing complex chemical tasks more rapidly. Here, we present an overview of the role of LMs toward accelerated molecular discovery. We commence with the conventional scientific discovery method and then discuss how molecular generative models can be coupled with molecular property prediction models. Seeking for practical usability, we then present the reader with selected software tools and libraries for scientific language modeling. We close with a vision for future molecule design that integrates natural language models into the discovery process through chatbots.

## 2　Accelerated molecular discovery

Molecule discovery, intricately linked to optimizing diverse properties in a vast space, challenges conventional scientific methods. In chemistry's Design-Make-Test-Analyze (DMTA) cycle, synthesis costs and time constraints create a bottleneck that hampers hypothesis refinement (cf. Figure 1a). Traditional approaches are largely driven by medicinal chemists who design "molecule hypotheses" which are biased, ad-hoc and non-exhaustive. This hinders progress in addressing global issues, driving crucial necessity for an accelerated process of molecule discovery. Thus, a key challenge lies in improving speed and quality of evaluating such "molecule hypotheses" grounded on laboratory work.

Deep generative models have recently emerged as a promising tool to expedite the hypothesis/design phase in molecular discovery. However, even the most advanced molecular generative models require an efficient method for large-scale virtual screening to test their hypotheses. The *accelerated molecular discovery* cycle adds a validation loop to DMTA, rapidly evaluating numerous hypotheses inexpensively (cf. Figure 1b). This loop enhances the design-phase generative model, ensuring only promising hypotheses advance to the synthesis and physical experimentation stages.

### 2.1　Molecule Representation

Data representation is critical as it determines which information is available for the model. As illustrated in Figure 2, various molecular representations exist. Due to popularity of chemical language models (CLMs), this section focuses on text-representations of molecules. A more focused discussion on CLMs was published by Grisoni [38].

**Figure 2:** An illustration of popular ways of representing a chemical molecule as input to a ML model. The representations may be (a) String-based, such as SMILES, SELFIES, or InChI which use characters to represent different aspects of a molecule, (b) Structure-based, such as Graphs or MolFiles that encode connectivity and atomic position, and (c) Feature-based, such as Morgan Fingerprints, which encode local substructures as bits.

**Simplified Molecular Input Line-Entry System (SMILES)**   SMILES [98] is a string representation made up of specific characters for atoms, bonds, branches, aromaticity, rings and stereochemistry in molecular structures. The character-level representation enables easy tokenization, making SMILES an ideal input for LMs. SMILES are non-unique, so each molecule can be written as multiple SMILES strings. Hence, SMILES are either canonicalized or, alternatively, their multiplicity is used as data augmentation strategy [8] which has shown performance improvement in molecular property prediction [8, 51, 88] and molecular generation [3, 92]. In generative modeling, a common issue is the invalidity of SMILES strings due to an uneven number of ring opening/closure symbols or bond valence violations. SMILES strings can undergo further processing, such as kekulization or stereoinformation removal but employing canonicalized SMILES remains the most prevalent approach.

   **Tokenization** is the process of splitting a string into vectorizable units. These units are typically a single character, n-gram characters or words. Instead of splitting at the character level, SMILES are typically tokenized at the atom level with regular expressions [79] or by additionally including positional and connectivity information, thereby acknowledging that the same atom can have different encodings based on its location in the molecular structure [91]. SMILES may also be tokenized at the substructure level, as demonstrated by SMILES Pair Encoding (SMILES-PE) [52]. This method, inspired by byte-pair encoding, iteratively counts and merges frequently occurring SMILES token pairs until a given condition is met. Tokenization enables the creation of a vocabulary for SMILES representations.

   **Vocabularies** are dictionaries mapping tokens to vectors thus serving as gateway to LMs. For LMs to learn from SMILES, tokens are vectorized, either via one-hot encodings (where each row in the binary matrix corresponds to a SMILES position and each column signifies a token). However, this discrete method results in sparse, large matrices and thus, an alluring alternative is to learn a continuous embedding for each token during training. This facilitates the learning of semantic relationships between tokens and enhances performance. Since learning good embeddings requires a lot of data, models pre-trained on natural language corpora are a strong option to learn scientific language embeddings through fine-tuning [22].

**Self Referencing Embedded Strings (SELFIES)**   SELFIES [49] were introduced as an alternative to SMILES to counter the problem of generating invalid molecules. Unlike SMILES, SELFIES are generated using derivation rules to enforce valence-bond validity. They store branch length and ring

size to avoid open branches and rings. These supplementary attributes ensure a valid representation during molecule generation. While this strategy guarantees 100% validity, it could produce strings that are too short to be a useful molecule.

**International Chemical Identifier (InChI)** Introduced by the IUPAC, InChI [41] are strings encoding structural information including charge of the molecule in a hierarchical manner. The strings can get long and complex for larger molecules. To counter this, a hash called 'InChiKey' was developed to help with search and retrieval. InChIs are are less commonly used in LMs [39].

## 2.2 Generative Modelling

Generative modeling involves learning the data's underlying distribution with the intent of generating new samples, a technique pivotal in accelerating de novo drug discovery. A generative model may be conditional or unconditional. A conditional generative model utilizes provided data attributes or labels to generate new samples with desired properties, whereas an unconditional model solely provides a way to sample molecules similar to the training data [36]. The DMTA cycle particularly benefits from the conditional generation approach as it facilitates goal-oriented hypothesis design [9]. This section describes a few influential conditional generation models that act on chemical language to generate molecules satisfying user-defined conditions.



**Figure 3:** An illustration of conditional molecule generation using LMs. The process initiates with the collection and processing of multi-modal data, which is then compressed into a fixed-size latent representation. These representations are subsequently passed to a molecular generative model. The generated molecules then undergo in-silico property prediction, which is linked back to the generative model through a feedback loop during training. The in-silico models direct the generative model to produce property- or task-driven molecules using a reward function. In the inference stage, candidate molecules generated by the optimized model continue through the workflow for lab synthesis and subsequent experimental validation to determine their efficacy for the desired task.

### 2.2.1 Recurrent Neural Network (RNN)

The sequential nature of RNNs makes them suitable models for processing chemical languages. Proposed in the 90s, RNNs were the first flavor of CLMs [8, 79, 85]. Their hidden states are continuously updated as new tokens are passed to the network. During the generation process, tokens are produced auto-regressively. RNNs find use in generating molecule libraries [85] which are extensively used in drug development processes like screening. External scoring functions drive the generation of molecules with desired properties. RNNs are also adept at learning complex distributions [31] and generating a higher proportion of unique and valid SMILES [69], even though their inability to count occurrences of ring opening/closing symbols poses a challenge [46, 70].

### 2.2.2 Variational Autoencoder (VAE)

VAEs learn latent distribution parameters of molecules, thus enabling the generation of new molecules by sampling from this distribution. Their unique ability lies in learning a smooth, latent space that facilitates interpolation of samples, even for notoriously discrete entities like molecules [36]. To make it suitable for chemical language models (CLMs), any network compatible with string inputs can function as a VAE's encoder and decoder. Initial works primarily focused on single-modality applications, assessing latent space quality via downstream tasks [36]. This approach remains prevalent and can be used to generate, e.g., catalysts with an RNN-based VAE [78] . Here, a latent space is learned and assessed by predicting the catalyst binding energy. Lim et al. [53] takes it a step further by concatenating a condition vector to the input and the latent embedding generated by the recurrent network-based VAE's encoder. This approach enables the generation of molecules specifically tailored to the given conditions. The scope of VAEs expanded progressively into multi-modal settings for conditional molecule generation, as visualized in Figure 3 and exemplified by Born et al. [11, 12, 13]. These works on task-driven molecule generation incorporate contextual information like gene expression [13] or protein targets [11, 12] or even both [45]. VAEs learn embeddings of context information and primer drugs, which are merged before decoding to produce molecules. A reinforcement-learning-based approach directs the model to produce molecules with desired properties using rewards.

### 2.2.3 Transformer

The self-attention attribute of Transformers [93] have propelled these models to the forefront of NLP. Transformers have an encoder module that relies on this self-attention to learn embeddings of the input and the context associated with this input. The decoder module predicts tokens using the context learnt by the encoder and previously generated tokens through attention. For generative modeling, decoder-only transformer like the Generative Pre-Training Transformer (GPT) [72] have become the dominant approach. This success was translated to the scientific language domain. One of the first models to use the GPT architecture for conditional molecule generation is MolGPT [4]. SMILES tokens concatenated with a condition vector that summarizes the desired properties and scaffolds are passed as input to this model, which is then trained on the next token prediction task to generate molecules. GPT-like models coupled with RL can also be used to optimize molecular properties like pIC50 [61]. In this two-stage approach, embeddings are first learnt from SMILES strings, and the embedding space is then optimized such that the model samples molecules with the desired properties. Going beyond just using GPT-like architectures for molecule generation, Regression Transformer [10] is a seminal work that formulates conditional sequence modeling as a regression problem. This gives rise to a natural multitask model that concurrently performs property prediction and conditional molecular generation. This is achieved by concatenating conventional molecular tokens with property tokens and employing an training scheme that alternates which parts of the sequence are masked.

All these works are testament to the generative capabilities of Transformer-based models. The superior quality of learned embeddings coupled with its ability to handle parallel processing and scalability makes it a top choice for the task of conditional molecule generation, with promising applications in drug discovery and other areas of molecular design [66].

## 2.3 Property Prediction

Whether a discovery is novel or not, property prediction is a key step in validating the molecules for a given use case. The success of a molecule depends on a myriad of factors, including how it interacts with its environment. The MoleculeNet datasets [103] are a commonly used benchmark for property prediction. It is curated from public datasets and comprises over 700,000 compounds tested on various properties. Born et al. [15] uses a multiscale convolutional attention model to predict toxicity from SMILES. The model has three kernel sizes for the convolutional network and uses a a Bahdanau attention mechanism [5]. The model shows a superior performance overall on various MoleculeNet tasks compared to all other SMILES-based models. A recent trend is to use transformer-encoders to learn embeddings for molecules and then apply a multilayer perceptron (MLP) on the embeddings for property prediction. MolBERT [29] and ChemBERTA [20]) are two such examples.

These transformer-based models use a BERT backbone to learn molecular embeddings from SMILES and predict properties. Similarly, Molformer [75] uses a transformer-encoder with linear attention and relative positional encoding to learn compressed molecular representations which are then fine-tuned on chemical property prediction benchmarks. To equip transformers with better inductive biases to handle molecules, adaptations of the attention mechanism were proposed. The molecule attention transformer (MAT) incorporates inter-atomic distances and graph structure into the attention mechanism [58]. An improvement over this model is the *relative*-MAT which fuses the distance embedding, bond embedding and neighbourhood embedding and achieves competitive performances on a range of property prediction tasks [59].

# 3    Software tools for scientific language modeling

The paradigm shift towards open-sourcing software has exerted a profound influence in chemistry. Commonly listed implications of open-sourcing in the context of drug discovery include catalyzation of methodological development, fostering of collaboration and ease of scientific reproducibility [35]. In this section we present several software assets (e.g., Python packages or cloud-based web apps) that are key to enable molecular discovery.

## 3.1    Natural language models

The success story of the Transformer [93] as most widely adopted neural network architecture goes hand in hand with the rise of the `transformers` library [101], developed since 2019 by HuggingFace. Initially intended for NLP applications, Transformers were adopted interdisciplinarily, e.g in computer vision [25], reinforcement learning [19], protein folding [47] and, of course, chemistry [84]. *HuggingFace* provides the largest public hub of language models and it offers implementations of all recent models as well as a diverse collection of pretrained models available for fine-tuning or inference. While most of their models focus on NLP, selected models are designed for life science applications, in particular molecular property prediction (e.g., *ChemBerta* [20]), molecular captioning (e.g., *MolT5* [26]), text-based molecular generation (e.g., *MolT5* [26]) but also unsupervised protein language models (e.g., *ProtBert*, *ProtAlbert*, *ProtXLNet* or *ProtT5* [27]). Moreover, some available models like *Multimodal Text and Chemistry T5* [22] are prompt-based multitasker that besides the above mentioned tasks also perform additional tasks such as forward/backward reaction prediction.

## 3.2    GT4SD – Generative modeling toolkits

Python libraries like `GT4SD` (the Generative Toolkit for Scientific Discovery [57]), `TdC` (Therapeutics Data Commons [43]) or `deepchem` [73] were developed primarily for molecular discovery applications, but especially `GT4SD` offers ample support of language models (LMs). `GT4SD` is designed to enable researchers and developers to use, train, fine-tune and distribute state-of-the-art generative models for sciences with a focus on the design of organic materials. It is compatible and inter-operable with many existing libraries and, beyond `transformers`, it also gives access to diffusion models (`diffusers` [96]) or graph generative models (`TorchDrug` [106]). Next to established molecular generation benchmark like `Moses` [69] and `GuacaMol` [16] that include VAEs, generative adversarial networks (GANs), genetic algorithms, and many evaluation metrics for molecular design, `gt4sd` also supports very contemporary models like the *Regression Transformer* for concurrent sequence regression and property-driven molecular design [10], *GFlowNets* for highly diverse candidate generation [6] or *MoLeR* for motif-constrained molecule generation [60]. `GT4SD` ships with a harmonized interface and a set of command line tools that access a registry of generative models to run or train any model with a few lines of code. Trained models can be shared to a cloud-hosted model hub and the library is build to facilitate consumption by containerization or distributed computing systems. To date, it includes $\sim 50$ property prediction endpoints for small molecules, proteins and crystals and overall hosts $\sim 30$ pre-trained algorithms for material design, 20 free webapps [2] and many Jupyter/Colab notebooks.

## 3.3 RXN for Chemistry: Reaction and synthesis language models

Once a molecule has been selected for experimental validation, a tangible synthesis route has to be identified. Since the most important tasks in chemical reaction modeling can be framed as sequence conversion problems, the methodology developed for natural language translation can be seamlessly translated to chemistry [84]. In this analogy, atoms are characters, molecules are words, reactions are sentences and precursors are translated into a product or vice versa.

The most mature and flexible library for reaction modeling with LMs is the package `rxn4chemistry` [32]. It wraps the API of the *IBM RXN for Chemistry* platform, a freely accessible web application that gives access to a rich set of language models for different tasks in reaction chemistry. The flagship architecture has been the *Molecular Transformer* (MT), an autoregressive encoder-decoder model, originally applied to predict outcomes of chemical reactions in organic chemistry [80]. Notably, the MT uses a purely data-driven, template-free approach that, unlike many graph-based models, can directly represent stereochemistry and thus also exhibits excellent performance on regio- and stereoselective reactions [67]. The MT was applied to single-step retrosynthesis [90] and became the linchpin of a multi-step retrosynthesis model with a hypergraph exploration strategy [81]. This approach was later generalized to enzymatic reactions with a tokenization scheme based on enzyme classes which facilitated biocatalyzed synthesis planning and paved the road towards more sustainable and green chemistry [71]. Derivatives of the MT helped to enhance diversity in single-step retrosynthesis [90] and a prompt-based disconnection scheme proposed by Thakkar et al. [89] significantly improved controllability by allowing the user to mark a disconnection side in the reactant. Interestingly, an encoder-only derivative of the MT (that replaced the autoregressive decoder with a classification head and leveraged BERT-style [24] self-supervised pretraining on reactions) excelled in predicting reaction classes [83]. The hidden representations of such a model were found to encode reaction types and thus allowing to map reaction atlases and to perform reaction similarity search. This gave rise to the `rxnfp` package for chemical reaction fingerprinting. Strikingly, masked language modeling also led later to the discovery that the learned attention weights of the Transformer are "secretly" performing atom mapping between products and reactions [82]. The epiphany that CLMs accomplish atom mapping without supervision or human labeling bridged the gap between rule-based and data-driven approaches in reaction modeling, making this once tedious experimental task more efficient.

In the quest for automation in organic chemistry, once the precursors for a molecule's synthesis route are identified, the subsequent crucial phase involves seeking an actionable, stepwise synthesis protocol that is ideally amenable for autonomous execution on a robotic platform, such as *IBM RoboRXN*. In two seminal works Vaucher et al. demonstrated that encoder-decoder Transformers can extract chemical synthesis actions, first from experimental procedures described in patents [94] and later predict them directly from the reaction SMILES [95]. Notable, all the aforementioned models are available via the *IBM RXN for Chemistry* platform which even allows to control and monitor the robotic platform directly from the web interface. For the daunting task of multistep retrosynthesis planning, *RXN* also includes non-transformer based models like *AiZynthFinder* [34], a Monte Carlo Tree Search approach build on top of a RNN. Most of the *RXN* models can be executed also via the `rxn4chemistry` Python package.

## 3.4 Specialized libraries

**Molecular property prediction.** `HuggingMolecules` is a library solely devoted to aggregate, standardize and distribute molecular property prediction LMs [33]. It contains many encoder-only CLMs, some of them with geometrical and structure-aware inductive biases (e.g., the MAT [58] or its successor, the R-MAT [59]) while others being pure BERT-based models that were trained on SMILES (e.g,. *MolBERT* [29] or *ChemBERTA* [20]).

**Data processing.** RDKit [50] is a library for manipulating molecules in Python. For narrower applications like ML data preparation several tools exist. First, `rxn-chemutils` is a library with chemistry-related utilities from RXN for Chemistry. It includes functionalities for standardizing SMILES (e.g., canonicalization or sanitization) but also conversions to other representations (e.g., InChI). It harmonizes reaction SMILES and prepares them for consumption by CLMs, including also SMILES aug-

mentation (by traversing the molecular graph in a non-canonical order) and tokenization. Another library with a similar focus is `pytoda` [12, 13]. It does not support reaction SMILES but implements richer preprocessing utilities, allowing to chain >10 SMILES transformations (e.g., kekulization [15]). It supports different languages (e.g., SELFIES [49] or BigSMILES [54]) and tokenization schemes (e.g., SMILES-PE [52]). Similar functionalities are available for proteins including different languages (IU-PAC, UniRep or Blosum62) and protein sequence augmentation strategies [14]. For small molecules, proteins, and polymers, dedicated language classes facilitate the integration with LMs by storing vocabularies, performing online transformations and feeding to custom datasets. Datasets exist for predicting molecular properties, drug sensitivity, protein-ligand affinity or for self-supervision on small molecules, proteins or polymers.

## 3.5   General purpose platforms

Several general-purpose platforms for molecular discovery have been launched recently, sometimes even preserving privacy through federated learning (i.e., decentralized, distributed training). For example, MELLODDY [42] is a collaborative effort aimed at cross-pharma federated learning of 2.6 billion confidential activity data points. Similarly, VirtualFlow [37] is an open-source platform facilitating large-scale virtual screening that was shown to identify potent KEAP1 inhibitors. With a focus on *de novo* drug design, Chemistry42 [44] is a proprietary platform integrating AI with computational and medicinal chemistry techniques.

# 4   Future of molecular discovery

A few years ago, the idea of querying an AI model – like one would a search engine – to not only extract scientific knowledge but also perform computational analyses was an overly ambitious feat. Scientific thinking comes from the ability to reason, and AI models cannot reason like humans, yet. However, these models can **learn** from humans. Our propensity to document everything has enabled us to train Large Language Models (LLMs), like ChatGPT [64] and GitHub Copilot [1], to mimic human responses. When brought into the context of computational science, this could equip non-experts to confidently conduct computational analyses through well-designed prompts. With human-in-the-loop, a synergistic effect could be created where the scientist provides feedback to the model on its output, thus aiding in better model optimization (a strategy called reinforcement learning from human feedback (RLHF) that has been proven critical for ChatGPT [21]). These applications also reduce the barrier for individuals from non-scientific backgrounds to gain a more hands-on experience in conducting scientific analyses without having to go through formal training in computational analysis.

This section provides a sneak peak into what's next for molecular discovery. Riding the LLM wave, the future holds a place for chatbot-like interfaces that may take care of all things computational in molecular discovery. This includes, for example, generating and iteratively improving design ideas, synthesis planning, material purchasing, performing routine safety checks, and validating experiments.

**The rise of foundation models in chemistry**

Conventionally, neural networks are trained for a single given task to achieve maximum performance. This essentially renders the models useless for other tasks, thus requiring a new model for every new task, even when the training domain is the same, which in turn imposes a constraint on the rate of our technological advancements. Over the last few years, this conventional approach has been challenged by Large Language Models (LLMs). It has been found that scaling up LLMs leads to astonishing performances in few-shot [17] and even zero-shot task generalization [76]. Referred to as "foundation models" [30, 63], these models, with typically billions of parameters, can perform multiple tasks despite being trained on one large dataset. Essentially, this multi-task learning is achieved by prompting LLMs with task instructions along with the actual query text which has been found to induce exceptional performance in natural language inference and sentence completion [76]. These findings have kicked off new research directions, such as prompt engineering [97] and in-context learning [17], in NLP.

The foundation model paradigm also finds an increasing adoption in chemistry. There is an increase in task-specific models integrating natural and chemical languages [26, 94, 95, 104]. Concurrently, multi-tasking in pure CLMs has also been advancing through models that combined tasks such as property prediction, reaction prediction and molecule generation either with small task-specific heads (e.g., T5Chem [56]) or via mask infilling (e.g., Regression Transformer [10]). Christofidellis et al. [22] were the first to bridge the gap and develop a fully prompt-based multi-task chemical and natural language model. Despite only 250M parameters, the *Multitask Text and Chemistry T5* was shown to outperform ChatGPT [64] and Galactica [87] on a contrived discovery workflow for re-discovering a common herbicide (natural text → new molecule → synthesis route → synthesis execution protocol).

## 4.1   The coalescence of chatbots with chemistry tools

Given the aforementioned strong task generalization performances of LLMs, building chatbot interfaces around it was a natural next step and thus next to ChatGPT [64], many similar tools were launched. Such tools were found to perform well on simplistic chemistry tasks [18, 99], opening potential to



**Figure 4:** Screenshot of the LLM-powered chatbot application `ChemChat`. Embedding the capabilities of existing resources such as PubChem [48], RDKit [50] or GT4SD [57] enables the assistant to execute programming routines in the background and thus answer highly subject-matter specific user requests without the user needing programming skills.

reshape how chemists interact with chemical data, enabling intuitive access to complex concepts and make valuable suggestions for diverse chemical tasks. Furthermore, AI models specifically developed by computer scientists for e.g. drug discovery or material science can be made available through applications powered by LLMs, such as chatbots. This minimizes the access barrier for subject matter experts who would otherwise require the respective programming skills to utilize these AI models. The power of such chatbots is reached through the coalscence of LLMs and existing chemistry software tools like PubChem [48], RDKit [50] or GT4SD [57]. Together, such applications can unleash the full

potential and value of these models by the strongly enhanced usage. An example of how the interaction with such a tool could look like is shown in Figure 4.

In this example, a user provides a molecule (either as SMILES string or via a molecule sketcher) and asks to identify the molecule. The chatbot relies on prompt-engineering in order to inform the LLM about all its available tools. The user input is first sent to the LLM which recognizes that one of its supported tools, in this case PubChem, can answer the question. The chatbot then sends a request to the PubChem API and returns a concise description of the molecule. The user subsequently asks to compute the logP partition coefficient [100] and the quantitative estimate of drug-likeness (QED) [7]. Calculation of both properties is enabled through the GT4SD tool [57] allowing the chatbot to answer the request with certainty. This will trigger a programming routine to accurately format the API request for GT4SD, i.e., composing the SMILES string with the logP or QED endpoint. The computation is then performed asynchronously and a separate call to the post-processing routine formats the LLM-generated string reply and composes the response object for the frontend.

This fusion of LLMs with existing tools gives rise to a chatbot assistant for material science and data visualization that can perform simple programming routines without requiring the user to know programming or have access to compute resources. A continuation of the conversation involving more complex user queries is shown in Figure 5. Having identified the initial molecule as theobromine with
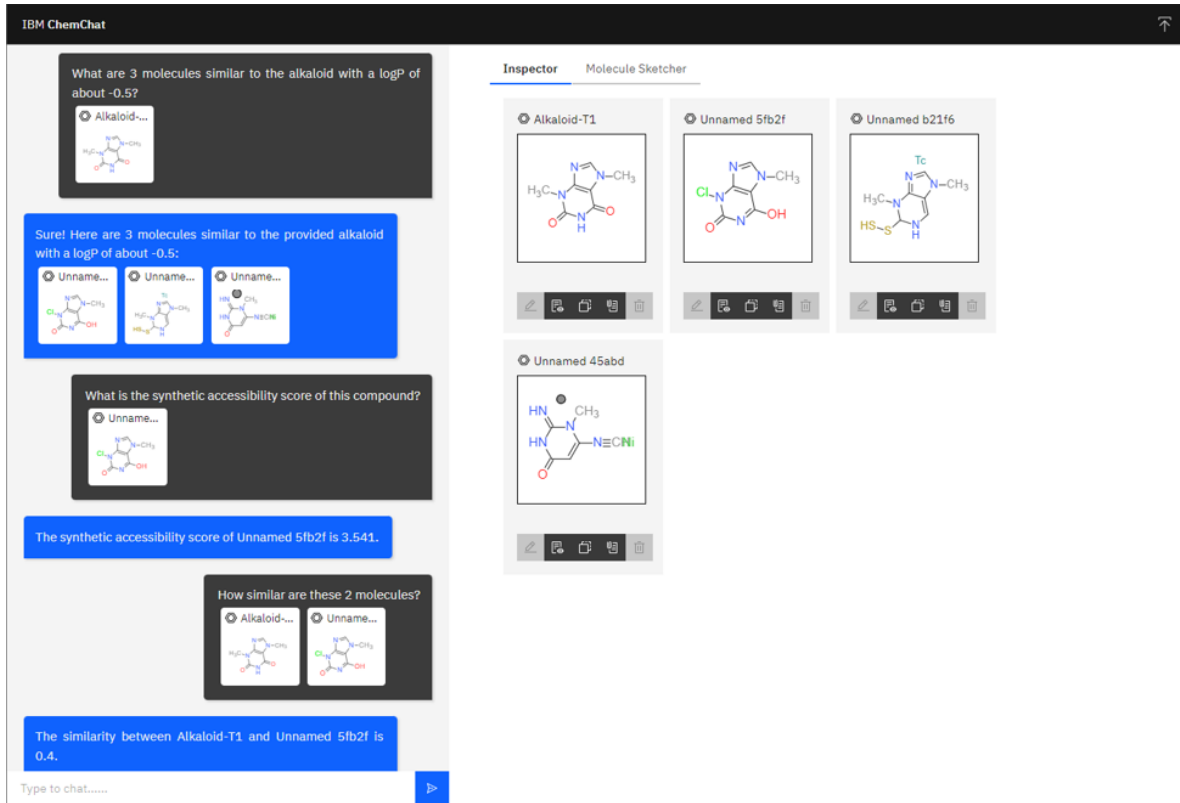


**Figure 5:** Screenshot of the LLM-powered chatbot application `ChemChat` showing the continuation of the conversation involving generative tasks through GT4SD's Regression Transformer [10] as well as property [28] and similarity calculation [74, 86].

a logP of -1.04, the user requests three similar molecules with a slightly increased logP of -0.5. Here, `ChemChat` identifies the Regression Transformer [10] as the available tool to perform substructure-constrained, property-driven molecule design. Once the routine has been executed and the three candidate SMILES are collected, the text result is post-processed to add more response data objects such as molecule visualizations, datasets or Vega Lite specs for interactive visualizations.

In conclusion, chatbots can facilitate the integration of essentially all major cheminformatics software in a truly harmonized and seamless manner. While LLMs are not intrinsically capable to perform

complex routines, at least not with high precision and in a trustworthy manner, the synergy between their natural language abilities with existing chemistry tools has the potential to transform the way chemistry is performed.

# References

[1] Github copilot. https://copilot.github.com/, 2021. Accessed: August 8, 2023.

[2] Abubakar Abid, Ali Abdalla, Ali Abid, Dawood Khan, Abdulrahman Alfozan, and James Zou. Gradio: Hassle-free sharing and testing of ml models in the wild. *arXiv preprint arXiv:1906.02569*, 2019.

[3] Josep Arús-Pous, Simon Viet Johansson, Oleksii Prykhodko, Esben Jannik Bjerrum, Christian Tyrchan, Jean-Louis Reymond, Hongming Chen, and Ola Engkvist. Randomized smiles strings improve the quality of molecular generative models. *Journal of cheminformatics*, 11(1):1–13, 2019.

[4] Viraj Bagal, Rishal Aggarwal, PK Vinod, and U Deva Priyakumar. Molgpt: molecular generation using a transformer-decoder model. *Journal of Chemical Information and Modeling*, 62(9):2064–2076, 2021.

[5] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.

[6] Yoshua Bengio, Tristan Deleu, Edward J Hu, Salem Lahlou, Mo Tiwari, and Emmanuel Bengio. Gflownet foundations. *Preprint at https://arxiv.org/abs/2111.09266*, 2021.

[7] G Richard Bickerton, Gaia V Paolini, Jérémy Besnard, Sorel Muresan, and Andrew L Hopkins. Quantifying the chemical beauty of drugs. *Nat. Chem.*, 4(2):90–98, 2012.

[8] Esben Jannik Bjerrum. Smiles enumeration as data augmentation for neural network modeling of molecules. *arXiv preprint arXiv:1703.07076*, 2017.

[9] Jannis Born and Matteo Manica. Trends in deep learning for property-driven drug design. *Current medicinal chemistry*, 28(38):7862–7886, 2021.

[10] Jannis Born and Matteo Manica. Regression transformer enables concurrent sequence regression and generation for molecular language modelling. *Nature Machine Intelligence*, 5(4):432–444, 2023.

[11] Jannis Born, Tien Huynh, Astrid Stroobants, Wendy D Cornell, and Matteo Manica. Active site sequence representations of human kinases outperform full sequence representations for affinity prediction and inhibitor generation: 3d effects in a 1d model. *Journal of Chemical Information and Modeling*, 62(2):240–257, 2021.

[12] Jannis Born, Matteo Manica, Joris Cadow, Greta Markert, Nil Adell Mill, Modestas Filipavicius, Nikita Janakarajan, Antonio Cardinale, Teodoro Laino, and María Rodríguez Martínez. Data-driven molecular design for discovery and synthesis of novel ligands: a case study on sars-cov-2. *Mach. Learn.: Sci. Technol.*, 2(2):025024, 2021.

[13] Jannis Born, Matteo Manica, Ali Oskooei, Joris Cadow, Greta Markert, and María Rodríguez Martínez. PaccMann$^{RL}$: De novo generation of hit-like anticancer molecules from transcriptomic data via reinforcement learning. *iScience*, 24(4):102269, 2021.

[14] Jannis Born, Yoel Shoshan, Tien Huynh, Wendy D Cornell, Eric J Martin, and Matteo Manica. On the choice of active site sequences for kinase-ligand affinity prediction. *Journal of chemical information and modeling*, 62(18):4295–4299, 2022.

[15] Jannis Born, Greta Markert, Nikita Janakarajan, Talia B Kimber, Andrea Volkamer, María Rodríguez Martínez, and Matteo Manica. Chemical representation learning for toxicity prediction. *Digital Discovery*, 2023.

[16] Nathan Brown, Marco Fiscato, Marwin HS Segler, and Alain C Vaucher. Guacamol: benchmarking models for de novo molecular design. *J. Chem. Inf. Model.*, 59(3):1096–1108, 2019.

[17] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

[18] Cayque Monteiro Castro Nascimento and André Silva Pimentel. Do large language models understand chemistry? a conversation with chatgpt. *Journal of Chemical Information and Modeling*, 63(6):1649–1655, 2023.

[19] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097, 2021.

[20] Seyone Chithrananda, Gabriel Grand, and Bharath Ramsundar. Chemberta: large-scale self-supervised pretraining for molecular property prediction. *arXiv preprint arXiv:2010.09885*, 2020.

[21] Paul F Christiano, Jan Leike, Tom Brown, Miljan Martic, Shane Legg, and Dario Amodei. Deep reinforcement learning from human preferences. *Advances in neural information processing systems*, 30, 2017.

[22] Dimitrios Christofidellis, Giorgio Giannone, Jannis Born, Ole Winther, Teodoro Laino, and Matteo Manica. Unifying molecular and textual representations via multi-task language modelling. In *International Conference on Machine Learning*, 2023.

[23] Payel Das, Tom Sercu, Kahini Wadhawan, Inkit Padhi, Sebastian Gehrmann, Flaviu Cipcigan, Vijil Chenthamarakshan, Hendrik Strobelt, Cicero Dos Santos, Pin-Yu Chen, et al. Accelerated antimicrobial discovery via deep generative models and molecular dynamics simulations. *Nat. Biomed. Eng.*, 5(6):613–623, 2021.

[24] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.

[25] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*, 2021.

[26] Carl Edwards, Tuan Lai, Kevin Ros, Garrett Honke, Kyunghyun Cho, and Heng Ji. Translation between molecules and natural language. In *2022 Conference on Empirical Methods in Natural Language Processing, EMNLP 2022*, 2022.

[27] Ahmed Elnaggar, Michael Heinzinger, Christian Dallago, Ghalia Rehawi, Wang Yu, Llion Jones, Tom Gibbs, Tamas Feher, Christoph Angerer, Martin Steinegger, Debsindhu Bhowmik, and Burkhard Rost. Prottrans: Towards cracking the language of life's code through self-supervised deep learning and high-performance computing. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pages 1–1, 2021. doi: 10.1109/TPAMI.2021.3095381.

[28] Peter Ertl and Ansgar Schuffenhauer. Estimation of synthetic accessibility score of drug-like molecules based on molecular complexity and fragment contributions. *Journal of cheminformatics*, 1:1–11, 2009.

[29] Benedek Fabian, Thomas Edlich, Héléna Gaspar, Marwin Segler, Joshua Meyers, Marco Fiscato, and Mohamed Ahmed. Molecular representation learning with language models and domain-relevant auxiliary tasks. *arXiv preprint arXiv:2011.13230*, 2020.

[30] Nanyi Fei, Zhiwu Lu, Yizhao Gao, Guoxing Yang, Yuqi Huo, Jingyuan Wen, Haoyu Lu, Ruihua Song, Xin Gao, Tao Xiang, et al. Towards artificial general intelligence via a multimodal foundation model. *Nature Communications*, 13(1):3094, 2022.

[31] Daniel Flam-Shepherd, Kevin Zhu, and Alán Aspuru-Guzik. Language models can learn complex molecular distributions. *Nature Communications*, 13(1):3293, 2022.

[32] IBM RXN for Chemistry team. rxn4chemistry: Python wrapper for the IBM RXN for Chemistry API. https://github.com/rxn4chemistry/rxn4chemistry, 2023.

[33] Piotr Gaiński, Łukasz Maziarka, Tomasz Danel, and Stanisław Jastrzebski. Huggingmolecules: An open-source library for transformer-based molecular property prediction (student abstract). In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 12949–12950, 2022.

[34] Samuel Genheden, Amol Thakkar, Veronika Chadimová, Jean-Louis Reymond, Ola Engkvist, and Esben Bjerrum. Aizynthfinder: a fast, robust and flexible open-source software for retrosynthetic planning. *Journal of cheminformatics*, 12(1):70, 2020.

[35] J Daniel Gezelter. Open source and open data should be standard practices, 2015.

[36] Rafael Gómez-Bombarelli, Jennifer N Wei, David Duvenaud, José Miguel Hernández-Lobato, Benjamín Sánchez-Lengeling, Dennis Sheberla, Jorge Aguilera-Iparraguirre, Timothy D Hirzel, Ryan P Adams, and Alán Aspuru-Guzik. Automatic chemical design using a data-driven continuous representation of molecules. *ACS central science*, 4(2):268–276, 2018.

[37] Christoph Gorgulla, Andras Boeszoermenyi, Zi-Fu Wang, Patrick D Fischer, Paul W Coote, Krishna M Padmanabha Das, Yehor S Malets, Dmytro S Radchenko, Yurii S Moroz, David A Scott, et al. An open-source drug discovery platform enables ultra-large virtual screens. *Nature*, 580(7805):663–668, 2020.

[38] Francesca Grisoni. Chemical language models for de novo drug design: Challenges and opportunities. *Current Opinion in Structural Biology*, 79:102527, 2023.

[39] Jennifer Handsel, Brian Matthews, Nicola J Knight, and Simon J Coles. Translating the inchi: adapting neural machine translation to predict iupac names from a chemical identifier. *Journal of cheminformatics*, 13(1):1–11, 2021.

[40] Emily Hargrave-Thomas, Bo Yu, and Jóhannes Reynisson. Serendipity in anticancer drug discovery. *World journal of clinical oncology*, 3(1):1, 2012.

[41] Stephen R Heller, Alan McNaught, Igor Pletnev, Stephen Stein, and Dmitrii Tchekhovskoi. Inchi, the iupac international chemical identifier. *Journal of cheminformatics*, 7(1):1–34, 2015.

[42] Wouter Heyndrickx, Lewis Mervin, Tobias Morawietz, Noé Sturm, Lukas Friedrich, Adam Zalewski, Anastasia Pentina, Lina Humbeck, Martijn Oldenhof, Ritsuya Niwayama, et al. Melloddy: cross pharma federated learning at unprecedented scale unlocks benefits in qsar without compromising proprietary information. 2022.

[43] Kexin Huang, Tianfan Fu, Wenhao Gao, Yue Zhao, Yusuf Roohani, Jure Leskovec, Coley Connor W, Cao Xiao, Jimeng Sun, and Marinka Zitnik. Therapeutics data commons: Machine learning datasets and tasks for drug discovery and development. *Advances in Neural Information Processing System*, 35, 2021.

[44] Yan A Ivanenkov, Daniil Polykovskiy, Dmitry Bezrukov, Bogdan Zagribelnyy, Vladimir Aladinskiy, Petrina Kamya, Alex Aliper, Feng Ren, and Alex Zhavoronkov. Chemistry42: an ai-driven platform for molecular design and optimization. *Journal of Chemical Information and Modeling*, 63(3):695–701, 2023.

[45] Nikita Janakarajan, Jannis Born, and Matteo Manica. A fully differentiable set autoencoder. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pages 3061–3071, 2022.

[46] Armand Joulin and Tomas Mikolov. Inferring algorithmic patterns with stack-augmented recurrent nets. *Advances in neural information processing systems*, 28, 2015.

[47] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Olaf Ronneberger, Kathryn Tunyasuvunakool, Russ Bates, Augustin Žídek, Anna Potapenko, et al. Highly accurate protein structure prediction with alphafold. *Nature*, 596(7873):583–589, 2021.

[48] Sunghwan Kim, Jie Chen, Tiejun Cheng, Asta Gindulyte, Jia He, Siqian He, Qingliang Li, Benjamin A Shoemaker, Paul A Thiessen, Bo Yu, et al. Pubchem 2019 update: improved access to chemical data. *Nucleic acids research*, 47(D1):D1102–D1109, 2019.

[49] Mario Krenn, Florian Häse, AkshatKumar Nigam, Pascal Friederich, and Alan Aspuru-Guzik. Self-referencing embedded strings (selfies): A 100% robust molecular string representation. *Machine Learning: Science and Technology*, 1(4):045024, 2020.

[50] Greg Landrum. Rdkit documentation. *Release*, 1(1-79):4, 2013.

[51] Xinhao Li and Denis Fourches. Inductive transfer learning for molecular activity prediction: Next-gen qsar models with molpmofit. *Journal of Cheminformatics*, 12(1):1–15, 2020.

[52] Xinhao Li and Denis Fourches. Smiles pair encoding: a data-driven substructure tokenization algorithm for deep learning. *Journal of chemical information and modeling*, 61(4):1560–1569, 2021.

[53] Jaechang Lim, Seongok Ryu, Jin Woo Kim, and Woo Youn Kim. Molecular generative model based on conditional variational autoencoder for de novo molecular design. *Journal of cheminformatics*, 10(1):1–9, 2018.

[54] Tzyy-Shyang Lin, Connor W Coley, Hidenobu Mochigase, Haley K Beech, Wencong Wang, Zi Wang, Eliot Woods, Stephen L Craig, Jeremiah A Johnson, Julia A Kalow, et al. Bigsmiles: a structurally-based line notation for describing macromolecules. *ACS central science*, 5(9): 1523–1531, 2019.

[55] Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin, Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, 2023.

[56] Jieyu Lu and Yingkai Zhang. Unified deep learning model for multitask reaction predictions with explanation. *Journal of Chemical Information and Modeling*, 62(6):1376–1387, 2022.

[57] Matteo Manica, Jannis Born, Joris Cadow, Dimitrios Christofidellis, Ashish Dave, Dean Clarke, Yves Gaetan Nana Teukam, Giorgio Giannone, Samuel C Hoffman, Matthew Buchan, et al. Accelerating material design with the generative toolkit for scientific discovery. *npj Computational Materials*, 9(1):69, 2023.

[58] Łukasz Maziarka, Tomasz Danel, Sławomir Mucha, Krzysztof Rataj, Jacek Tabor, and S Jastrzkebski. Molecule-augmented attention transformer. In *Workshop on Graph Representation Learning, Neural Information Processing Systems*, 2019.

[59] Łukasz Maziarka, Dawid Majchrowski, Tomasz Danel, Piotr Gaiński, Jacek Tabor, Igor Podolak, Paweł Morkisz, and Stanisław Jastrzkebski. Relative molecule self-attention transformer. *arXiv preprint arXiv:2110.05841*, 2021.

[60] Krzysztof Maziarz, Henry Jackson-Flux, Pashmina Cameron, Finton Sirockin, Nadine Schneider, Nikolaus Stiefl, Marwin Segler, and Marc Brockschmidt. Learning to extend molecular scaffolds with structural motif. In *The Tenth International Conference on Learning Representations, ICLR*, 2022.

[61] Eyal Mazuz, Guy Shtar, Bracha Shapira, and Lior Rokach. Molecule generation using transformers and policy gradient reinforcement learning. *Scientific Reports*, 13(1):8799, 2023.

[62] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.

[63] Michael Moor, Oishi Banerjee, Zahra Shakeri Hossein Abad, Harlan M Krumholz, Jure Leskovec, Eric J Topol, and Pranav Rajpurkar. Foundation models for generalist medical artificial intelligence. *Nature*, 616(7956):259–265, 2023.

[64] OpenAI. Chatgpt. https://chat.openai.com/chat, 2023. Accessed: August 8, 2023.

[65] OpenAI. Gpt-4 technical report, 2023.

[66] Nathaniel H Park, Matteo Manica, Jannis Born, James L Hedrick, Tim Erdmann, Dmitry Yu Zubarev, Nil Adell-Mill, and Pedro L Arrechea. Artificial intelligence driven design of catalysts and materials for ring opening polymerization using a domain-specific language. *Nature Communications*, 14(1):3686, 2023.

[67] Giorgio Pesciullesi, Philippe Schwaller, Teodoro Laino, and Jean-Louis Reymond. Transfer learning enables the molecular transformer to predict regio-and stereoselective reactions on carbohydrates. *Nature communications*, 11(1):4874, 2020.

[68] Pavel G Polishchuk, Timur I Madzhidov, and Alexandre Varnek. Estimation of the size of drug-like chemical space based on gdb-17 data. *J. Comput. Aid. Mol. Des.*, 27(8):675–679, 2013.

[69] Daniil Polykovskiy, Alexander Zhebrak, Benjamin Sanchez-Lengeling, Sergey Golovanov, Oktai Tatanov, Stanislav Belyaev, Rauf Kurbanov, Aleksey Artamonov, Vladimir Aladinskiy, Mark Veselov, et al. Molecular sets (moses): a benchmarking platform for molecular generation models. *Front. Pharmacol.*, 11:1931, 2020.

[70] Mariya Popova, Olexandr Isayev, and Alexander Tropsha. Deep reinforcement learning for de novo drug design. *Science advances*, 4(7):eaap7885, 2018.

[71] Daniel Probst, Matteo Manica, Yves Gaetan Nana Teukam, Alessandro Castrogiovanni, Federico Paratore, and Teodoro Laino. Biocatalysed synthesis planning using data-driven learning. *Nature communications*, 13(1):964, 2022.

[72] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.

[73] Bharath Ramsundar, Peter Eastman, Patrick Walters, Vijay Pande, Karl Leswing, and Zhenqin Wu. *Deep Learning for the Life Sciences*. O'Reilly Media, 2019. https://www.amazon.com/Deep-Learning-Life-Sciences-Microscopy/dp/1492039837.

[74] David Rogers and Mathew Hahn. Extended-connectivity fingerprints. *Journal of chemical information and modeling*, 50(5):742–754, 2010.

[75] Jerret Ross, Brian Belgodere, Vijil Chenthamarakshan, Inkit Padhi, Youssef Mroueh, and Payel Das. Large-scale chemical language representations capture molecular structure and properties. *Nature Machine Intelligence*, 4(12):1256–1264, 2022.

[76] Victor Sanh, Albert Webson, Colin Raffel, Stephen H Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud Stiegler, Teven Le Scao, Arun Raja, et al. Multitask prompted training enables zero-shot task generalization. In *ICLR 2022-Tenth International Conference on Learning Representations*, 2022.

[77] Jack W Scannell, Alex Blanckley, Helen Boldon, and Brian Warrington. Diagnosing the decline in pharmaceutical r&d efficiency. *Nat. Rev. Drug Discov.*, 11(3):191–200, 2012.

[78] Oliver Schilter, Alain Vaucher, Philippe Schwaller, and Teodoro Laino. Designing catalysts with deep generative models and computational data. a case study for suzuki cross coupling reactions. *Digital Discovery*, 2(3):728–735, 2023.

[79] Philippe Schwaller, Theophile Gaudin, David Lanyi, Costas Bekas, and Teodoro Laino. "found in translation": predicting outcomes of complex organic chemistry reactions using neural sequence-to-sequence models. *Chemical science*, 9(28):6091–6098, 2018.

[80] Philippe Schwaller, Teodoro Laino, Théophile Gaudin, Peter Bolgar, Christopher A Hunter, Costas Bekas, and Alpha A Lee. Molecular transformer: a model for uncertainty-calibrated chemical reaction prediction. *ACS central science*, 5(9):1572–1583, 2019.

[81] Philippe Schwaller, Riccardo Petraglia, Valerio Zullo, Vishnu H Nair, Rico Andreas Haeuselmann, Riccardo Pisoni, Costas Bekas, Anna Iuliano, and Teodoro Laino. Predicting retrosynthetic pathways using transformer-based models and a hyper-graph exploration strategy. *Chemical science*, 11(12):3316–3325, 2020.

[82] Philippe Schwaller, Benjamin Hoover, Jean-Louis Reymond, Hendrik Strobelt, and Teodoro Laino. Extraction of organic chemistry grammar from unsupervised learning of chemical reactions. *Science Advances*, 7(15):eabe4166, 2021.

[83] Philippe Schwaller, Daniel Probst, Alain C Vaucher, Vishnu H Nair, David Kreutter, Teodoro Laino, and Jean-Louis Reymond. Mapping the space of chemical reactions using attention-based neural networks. *Nature machine intelligence*, 3(2):144–152, 2021.

[84] Philippe Schwaller, Alain C Vaucher, Ruben Laplaza, Charlotte Bunne, Andreas Krause, Clemence Corminboeuf, and Teodoro Laino. Machine intelligence for chemical reaction space. *Wiley Interdisciplinary Reviews: Computational Molecular Science*, 12(5):e1604, 2022.

[85] Marwin HS Segler, Thierry Kogej, Christian Tyrchan, and Mark P Waller. Generating focused molecule libraries for drug discovery with recurrent neural networks. *ACS central science*, 4(1):120–131, 2018.

[86] Taffee T Tanimoto. Ibm internal report. *Nov*, 17:1957, 1957.

[87] Ross Taylor, Marcin Kardas, Guillem Cucurull, Thomas Scialom, Anthony Hartshorn, Elvis Saravia, Andrew Poulton, Viktor Kerkez, and Robert Stojnic. Galactica: A large language model for science. *arXiv preprint arXiv:2211.09085*, 2022.

[88] Igor V Tetko, Pavel Karpov, Eric Bruno, Talia B Kimber, and Guillaume Godin. Augmentation is what you need! In *International Conference on Artificial Neural Networks*, pages 831–835. Springer, 2019.

[89] Amol Thakkar, Alain C Vaucher, Andrea Byekwaso, Philippe Schwaller, Alessandra Toniato, and Teodoro Laino. Unbiasing retrosynthesis language models with disconnection prompts. *ACS Central Science*, 2023.

[90] Alessandra Toniato, Philippe Schwaller, Antonio Cardinale, Joppe Geluykens, and Teodoro Laino. Unassisted noise reduction of chemical reaction datasets. *Nature Machine Intelligence*, 3(6):485–494, 2021.

[91] Umit V Ucak, Islambek Ashyrmamatov, and Juyong Lee. Improving the quality of chemical language model outcomes with atom-in-smiles tokenization. *Journal of Cheminformatics*, 15(1): 55, 2023.

[92] Ruud van Deursen, Peter Ertl, Igor V Tetko, and Guillaume Godin. Gen: highly efficient smiles explorer using autodidactic generative examination networks. *Journal of Cheminformatics*, 12 (1):1–14, 2020.

[93] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[94] Alain C Vaucher, Federico Zipoli, Joppe Geluykens, Vishnu H Nair, Philippe Schwaller, and Teodoro Laino. Automated extraction of chemical synthesis actions from experimental procedures. *Nature communications*, 11(1):3601, 2020.

[95] Alain C Vaucher, Philippe Schwaller, Joppe Geluykens, Vishnu H Nair, Anna Iuliano, and Teodoro Laino. Inferring experimental procedures from text-based representations of chemical reactions. *Nature communications*, 12(1):2573, 2021.

[96] Patrick von Platen, Suraj Patil, Anton Lozhkov, Pedro Cuenca, Nathan Lambert, Kashif Rasul, Mishig Davaadorj, and Thomas Wolf. Diffusers: State-of-the-art diffusion models, 10 2022. URL https://github.com/huggingface/diffusers.

[97] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in Neural Information Processing Systems*, 35:24824–24837, 2022.

[98] David Weininger. Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. *J. Chem. Inf. Comp. Sci.*, 28(1):31–36, 1988.

[99] Andrew D White, Glen M Hocky, Heta A Gandhi, Mehrad Ansari, Sam Cox, Geemi P Wellawatte, Subarna Sasmal, Ziyue Yang, Kangxin Liu, Yuvraj Singh, et al. Assessment of chemistry knowledge in large language models that generate code. *Digital Discovery*, 2(2):368–376, 2023.

[100] Scott A Wildman and Gordon M Crippen. Prediction of physicochemical parameters by atomic contributions. *Journal of chemical information and computer sciences*, 39(5):868–873, 1999.

[101] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38–45, 2020.

[102] Olivier J Wouters, Martin McKee, and Jeroen Luyten. Estimated research and development investment needed to bring a new medicine to market, 2009-2018. *Jama*, 323(9):844–853, 2020.

[103] Zhenqin Wu, Bharath Ramsundar, Evan N Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S Pappu, Karl Leswing, and Vijay Pande. Moleculenet: a benchmark for molecular machine learning. *Chemical science*, 9(2):513–530, 2018.

[104] Zheni Zeng, Yuan Yao, Zhiyuan Liu, and Maosong Sun. A deep-learning system bridging molecule structure and biomedical text with comprehension comparable to human professionals. *Nature communications*, 13(1):862, 2022.

[105] Alex Zhavoronkov, Yan A Ivanenkov, Alex Aliper, Mark S Veselov, Vladimir A Aladinskiy, Anastasiya V Aladinskaya, Victor A Terentiev, Daniil A Polykovskiy, Maksim D Kuznetsov, Arip Asadulaev, et al. Deep learning enables rapid identification of potent ddr1 kinase inhibitors. *Nat. Biotechnol.*, 37(9):1038–1040, 2019.

[106] Zhaocheng Zhu, Chence Shi, Zuobai Zhang, Shengchao Liu, Minghao Xu, Xinyu Yuan, Yangtian Zhang, Junkun Chen, Huiyu Cai, Jiarui Lu, et al. Torchdrug: A powerful and flexible machine learning platform for drug discovery. *Preprint at https://arxiv.org/abs/2202.08320*, 2022.