

A Platform for the Biomedical Application of Large Language Models

Sebastian Lobentanzer

sebastian.lobentanzer@uni-heidelberg.de

Julio Saez-Rodriguez

pub.saez@uni-heidelberg.de

*Heidelberg University, Faculty of Medicine and Heidelberg University Hospital,
Institute for Computational Biomedicine, Heidelberg, Germany*

Abstract

The wealth of knowledge we have amassed in the context of biomedical science has grown exponentially in the last decades. Consequently, understanding and contextualising scientific results has become increasingly difficult for any single individual. In contrast, current Large Language Models (LLMs) can remember an enormous amount of information, but have notable shortcomings, such as a lack of generalised awareness, logical deficits, and a propensity to hallucinate. To improve biomedical analyses, we propose to combine human ingenuity and machine memory by means of an open and modular conversational platform, ChatGSE (<https://chat.biocypher.org>). We safeguard against common LLM shortcomings using general and biomedicine-specific measures and allow automated integration of popular bioinformatics methods. Ultimately, we aim to improve the AI-readiness of biomedicine and make LLMs more useful and trustworthy in research applications.

Main

Despite our technological advances, biology and biomedicine continue to pose incredible challenges (Gallagher, 2023). We measure more and more data points with ever-increasing resolution to such a degree that their analysis and interpretation have become the bottleneck for their exploitation. One reason for this challenge may be the inherent limitation of human knowledge (Marois and Ivanoff, 2005). Even seasoned domain experts cannot know the implications of every molecule, be it metabolite, DNA, RNA, or protein, even in their own domain. In addition, biological events are context-dependent, for instance with respect to a cell type or specific disease.

Large Language Models (LLMs) of the current generation, on the other hand, can access enormous amounts of knowledge, encoded (incomprehensibly) in their billions of parameters (Chowdhery *et al.*, 2022; Thoppilan *et al.*, 2022; OpenAI, 2023). Trained correctly, they can

recall and combine virtually limitless knowledge from their training set. ChatGPT has taken the world by storm, and many biomedical researchers already use LLMs in their daily work, for general as well as bioinformatics-specific tasks (Hou and Ji, 2023; Vert, 2023). However, the current, predominantly manual, way of interacting with LLMs is virtually non-reproducible, and their behaviour can be erratic. For instance, they are known to hallucinate: they make up facts as they go along, and, to make matters worse, are convinced - and convincing - regarding the truth of their hallucinations (Moor *et al.*, 2023; Vert, 2023). While current efforts towards AGI (Artificial General Intelligence) manage to ameliorate some of the shortcomings by ensembling multiple models (LangChain, 2023) with long-term memory stores (Richards, 2023), the current generation of AI does not inspire adequate trust to be applied to biomedical problems without supervision (Moor *et al.*, 2023). Additionally, biomedicine demands greater care in data privacy, licensing, and transparency than most other real-world issues.

A major aim of computational biology is to distil high-dimensional molecular measurements into a humanly digestible form by projecting the measurements into a lower-dimensional space composed of gene programs, pathways, or other functional groupings of biological entities, for example via gene set enrichment analyses. However, even this distilled knowledge requires advanced expertise and thorough literature research to effectively interpret and exploit, and benchmarking the methods' performance is non-trivial (Geistlinger *et al.*, 2021).

To improve and accelerate this interpretation and exploration, we have developed ChatGSE (<https://chat.biocypher.org>), a platform for communicating with LLMs specifically tuned to biomedical research (**Figure 1**). The platform guides the human researcher intuitively through the interaction with the model, while counteracting the problematic behaviours of the LLM. Since the interaction is mainly based on plain text, it can be used by virtually any researcher. We engineer prompts around the queries of the user to improve model performance with regard to biomedicine, and automate the integration of popular bioinformatics methods, such as differential expression and gene set enrichment (**Supplementary Note 1**).

On the model side, we implement several measures in addition to the prompt engineering around the user's queries. For instance, we deploy a second model to safeguard the factual correctness of the primary LLM's responses (**Supplementary Note 2**). These interactions are handled by a pre-programmed conversational "Assistant," which dynamically orchestrates LLM agents with distinct tasks using a Python model chaining framework (LangChain, 2023). Using

vector database approaches, the user's prompts can be further supplemented with information extracted from pertinent, user-provided literature (**Supplementary Note 3**).

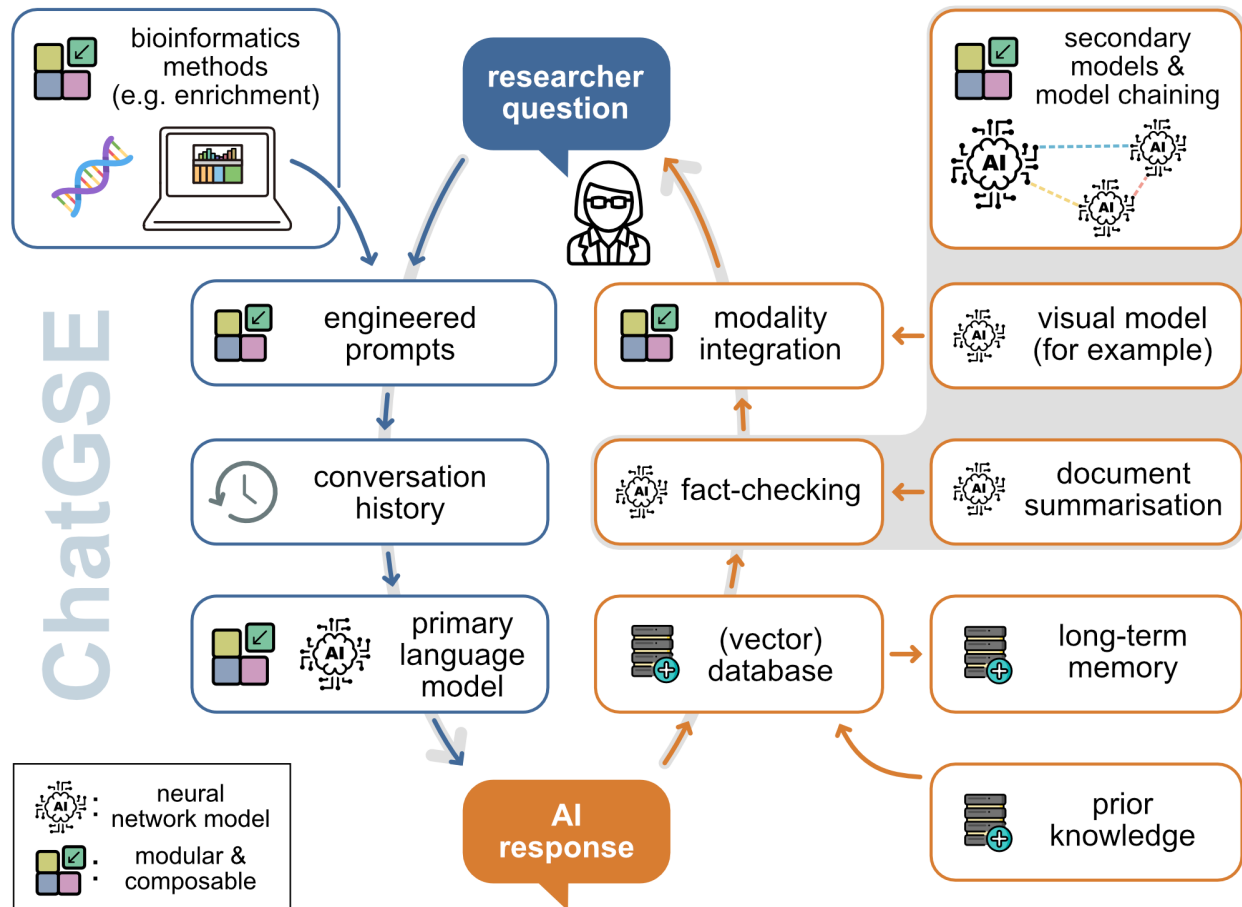


Figure 1: The ChatGSE composable platform architecture (simplified). The user submits a question about a topic of interest (e.g., an experiment) along with the low-dimensional results of a bioinformatics analysis (top left). The platform's main response circuit (blue) composes a number of specifically engineered prompts and passes them (and a conversation history) to the primary LLM, which generates a response for the user based on all inputs. This response is simultaneously used to prompt the secondary circuit (orange), which fulfils auxiliary tasks to complement the primary response. In particular, using search, the secondary circuit queries a database as prior knowledge repository and compares annotations to the primary response. The knowledge graph can also serve as long-term memory extension of the model. Further, an independent LLM receives the primary response for fact-checking, which can be supplemented with context-specific information by a document summarisation model. If this "second opinion" differs from the primary response, a warning is issued. The platform is composable in all aspects, in principle allowing arbitrary extensions to other, specialised models for additional tasks orchestrated by the primary LLM.

To increase data-awareness of the AI agents, we introduce connectivity to databases, which can extend the long-term memory of the models, semantically ground the biological entities with respect to suitable ontologies, and compare the model's responses to prior knowledge ground truth. By integrating a flexible knowledge graph creation framework (Lobentanzer *et al.*, 2022), we allow versatile use cases across the entire research spectrum. For example, connecting to a knowledge graph of cell markers based on Cell Ontology (Diehl *et al.*, 2016), the task of annotating single cell data sets can be automated and made more reproducible (**Supplementary Note 4**) by abstracting the pioneering efforts of manually executed studies (Hou and Ji, 2023).

Currently, the most powerful conversational AI platform, ChatGPT (OpenAI), is surrounded by data privacy concerns (Sterling, 2023). ChatGSE addresses this issue in two ways. Firstly, we provide access to the OpenAI models through the API (Application Programming Interface), which is subject to different, more stringent data protection than the web interface¹. Secondly, we aim to preferentially support open-source LLMs to facilitate more transparency in their application and increase data privacy by being able to run a model locally (Spirling, 2023). Our orchestration tool supports dozens of LLM providers (LangChain, 2023), such as the Hugging Face API, which can be used to query the recently released open-source ChatGPT-alternative HuggingChat or any other of the more than 100 000 open-source models on Hugging Face Hub (Hugging Face, 2023). Although OpenAI's models currently vastly outperform any alternatives in terms of both LLM performance and API convenience, we expect many open-source developments in this area in the future. Therefore, we support plug-and-play exchange of models to enhance biomedical AI-readiness.

In the future, we aim to integrate biological prior knowledge representation with LLM reasoning. Using the emergent strategies of *in-context learning*, *instruction learning*, and *chain-of-thought prompting* (Shen *et al.*, 2023), this can enable causal inference on relationships between biological entities, for instance via protein-protein interactions (**Supplementary Note 5**), and automated validation of literature references provided by the LLM (**Supplementary Note 6**). While the current models do not yet appear suited for unsupervised reasoning in the biomedical space, they can already save much time otherwise spent on web and literature searches. Additionally, the ChatGSE platform provides a reproducible environment for benchmarking of models and engineered prompts to gauge their biomedical reliability. The ability to chain arbitrary types of models enables advanced

¹ <https://openai.com/policies/terms-of-use>

applications, for instance connecting to visual modalities such as spatial omics. We provide further details and application scenarios in our Supplementary Notes.

While we focus on the biomedical field, the concept of the tool can easily be extended to other scientific domains by adjusting domain-specific prompts and data inputs, which in our framework are accessible in a composable and user-friendly manner. We openly develop the project on GitHub (<https://github.com/biocypher/ChatGSE>) under the permissive MIT licence and encourage contributions and suggestions from the community with regard to the addition of bioinformatics tool inputs, prompt engineering, safeguarding mechanisms, and any other feature.

Author Contributions

SL conceptualised and developed the platform and wrote the manuscript. JSR supervised the project, revised the manuscript, and acquired funding.

Acknowledgements

We thank Hanna Schumacher, Daniel Dimitrov, Pau Badia i Mompel, and Aurelien Dugourd for feedback on the original draft of the manuscript and the software.

Conflict of Interest

JSR reports funding from GSK, Pfizer and Sanofi and fees from Traverre Therapeutics and Astex Pharmaceuticals.

Bibliography

Chowdhery, A. *et al.* (2022) "PaLM: Scaling Language Modeling with Pathways," *arXiv* [Preprint]. doi:10.48550/arxiv.2204.02311.

Diehl, A.D. *et al.* (2016) "The Cell Ontology 2016: enhanced content, modularization, and ontology interoperability.," *Journal of biomedical semantics*, 7(1), p. 44. doi:10.1186/s13326-016-0088-7.

Gallagher, J. (2023) "Study reveals cancer's 'infinite' ability to evolve," *BBC News*, 12 April. Available at: <https://www.bbc.com/news/health-65252510> (Accessed: May 4, 2023).

Geistlinger, L. *et al.* (2021) "Toward a gold standard for benchmarking gene set enrichment analysis.," *Briefings in Bioinformatics*, 22(1), pp. 545–556. doi:10.1093/bib/bbz158.

Hou, W. and Ji, Z. (2023) "Reference-free and cost-effective automated cell type annotation with

GPT-4 in single-cell RNA-seq analysis.,” *BioRxiv* [Preprint]. doi:10.1101/2023.04.16.537094.

Hugging Face (2023) *Hugging Face Hub Documentation*. Available at: <https://huggingface.co/docs/hub/index> (Accessed: May 4, 2023).

LangChain (2023) *LangChain Documentation*. Available at: <https://python.langchain.com> (Accessed: May 4, 2023).

Lobentanzer, S. *et al.* (2022) “Democratising Knowledge Representation with BioCypher,” *arXiv* [Preprint]. doi:10.48550/arxiv.2212.13543.

Marois, R. and Ivanoff, J. (2005) “Capacity limits of information processing in the brain.,” *Trends in Cognitive Sciences*, 9(6), pp. 296–305. doi:10.1016/j.tics.2005.04.010.

Moor, M. *et al.* (2023) “Foundation models for generalist medical artificial intelligence.,” *Nature*, 616(7956), pp. 259–265. doi:10.1038/s41586-023-05881-4.

OpenAI (2023) “GPT-4 Technical Report,” *arXiv* [Preprint]. doi:10.48550/arxiv.2303.08774.

Richards, T.B. (2023) *AutoGPT, AutoGPT*. Available at: <https://autogpt.net/> (Accessed: May 4, 2023).

Shen, Y. *et al.* (2023) “HuggingGPT: Solving AI Tasks with ChatGPT and its Friends in HuggingFace,” *arXiv* [Preprint]. doi:10.48550/arxiv.2303.17580.

Spirling, A. (2023) “Why open-source generative AI models are an ethical way forward for science.,” *Nature*, 616(7957), p. 413. doi:10.1038/d41586-023-01295-4.

Sterling, T. (2023) “European privacy watchdog creates ChatGPT task force,” *Reuters*, 14 April.

Thoppilan, R. *et al.* (2022) “LaMDA: Language Models for Dialog Applications,” *arXiv* [Preprint]. doi:10.48550/arxiv.2201.08239.

Vert, J.-P. (2023) “How will generative AI disrupt data science in drug discovery?,” *Nature Biotechnology* [Preprint]. doi:10.1038/s41587-023-01789-6.

Supplementary Materials

In our Supplementary Notes, we explain the functions of our platform in more detail. Please note that several of the mentioned features, particularly more advanced ones, are in early developmental stages. For an up-to-date overview of current functionality of the platform, please visit <https://chat.biocypher.org> and <https://github.com/biocypher/ChatGSE>. ChatGSE is developed in Python (version 3.10), according to modern standards of software development (Lin *et al.*, 2020). We use Streamlit (version 1.21.0, <https://streamlit.io>) for the web application user interface. We include a code of conduct and contributor guidelines to offer accessibility and inclusivity to all that are interested in contributing to the framework.

Supplementary Note 1: Prompt Engineering

Recent experience with Large Language Models (LLMs) shows that the clever engineering of model prompts can yield drastic performance increases. For example, when performing logical inference, simply adding “Let’s think step by step” to the end of a question prompt increased LLM performance from ~20% to ~80% (Kojima *et al.*, 2022). As such, we highly prioritise the identification and adjustment of prompts tuned exactly to the requirements of the specific task to be performed, respecting “general” rules of interacting with LLMs as well as biomedicine-specific issues.

We designed the backend of ChatGSE to be completely flexible with regard to the application and rearrangement of prompts and prompt templates (which allow the insertion of variables such as the user input question or data). For templating, we use the corresponding generic functionality provided by LangChain (LangChain, 2023), while adding our biomedicine-specific layer on top. We provide general prompts to the primary model, setting it up to be helpful and concise in its responses, and individual prompts for each tool we want the primary model to “understand.” This includes explanation of the method itself as well as structural information about the data file containing the results.

To make this functionality available to all users, we provide a “Prompt Engineering” tab in the ChatGSE application, which allows the modification (or removal) of existing prompts, as well as the addition of new ones. We also provide functionality to import and export these prompts in JSON format to facilitate reproducible and shareable biomedical prompt engineering. To discuss and share examples of useful or ineffective prompts, we encourage all users to join the

#chatgse stream in our freely accessible Zulip channel at <https://biocypher.zulipchat.com>, or the GitHub discussion thread at <https://github.com/biocypher/ChatGSE/discussions/11>.

Supplementary Note 2: Correcting Agent

The propensity of LLMs to hallucinate untrue facts necessitates control mechanisms and guardrails for their application in research, particularly in high-stakes fields such as biomedicine (Vert, 2023). Some corrective incentive can be included in the instructions to the primary model, for instance by including a “Criticism” directive such as “Constructively self-criticise your big-picture behaviour constantly.” However, this does not guarantee to prevent hallucinations completely, as the same model that hallucinates also is responsible for correction.

Thus, we implement a modular combination of primary and corrective model, where the corrective agent – a “second opinion” – is set up with instructions tuned exactly to the task of fact-checking the responses of the primary model (“You are a fact-checker. Please judge the following statement for its factual correctness. [etc]”). The performance of the corrective agent can be further increased by using a more powerful model (e.g., moving from gpt-3.5-turbo to gpt-4), by pre-processing the primary model’s response (e.g., splitting the response into single sentences and fact-checking each sentence individually), and comparing the primary model’s statements to prior knowledge stored in a knowledge graph connected to the chat platform. As model development advances and other parties create models as powerful as OpenAI’s, we foresee it will be useful to combine models from different suppliers to increase diversity between primary and corrective models. Using ChatGSE’s modular framework, arbitrary numbers of corrective models can be added to the LLM chain.

To allow all users to interact with correcting functionality, we include a dedicated “Corrective Agent” tab for adjusting settings of the corrective model independent of the settings for the primary model. We also facilitate the testing of corrective agents and their prompts by providing a free-text field for sending false information to the corrective agent. This is necessary since some models, particularly those from OpenAI, are heavily regulated to not purposely provide false information, even for testing purposes.

The comparative power of the LLM agents can be further increased by connecting to a vector database containing embeddings of the contents of specific user-supplied documents. For more information, see the following Supplementary Note.

Supplementary Note 3: Document Summarisation / Vector Database Integration / In-context Learning

While the general knowledge of current LLMs is extensive, they may not know how to prioritise very specific scientific results, or they may not have had access to some research articles in their training data (e.g. due to their recency or licensing issues). To bridge this gap, we can provide additional information from relevant publications to the model via the prompt. However, we cannot add entire publications to the prompt, since the input length of current models still is restricted; we need to isolate the information that is specifically relevant to the question given by the user. To find this information, we perform a similarity search between the user's question and the contents of user-provided scientific articles (or other texts). The most efficient way to do this mapping is by using a vector database.

The contextual background information provided by the user (e.g. by uploading a scientific article of prior work related to the experiment to be interpreted) is split into pieces suitable to be digested by the LLM, which are individually embedded by the model. These embeddings (represented by vectors) are used to store the text fragments in a vector database; the storage as vectors allows fast and efficient retrieval of similar entities via the comparison of individual vectors. For example, the two sentences "Amyloid beta levels are associated with Alzheimer's Disease stage." and "One of the most important clinical markers of AD progression is the amount of deposited A-beta 42." would be closely associated in a vector database (given the embedding model is of sufficient quality, i.e., similar to GPT-3 or better), while traditional text-based similarity metrics probably would not identify them as highly similar.

By comparing the user's question to prior knowledge in the vector database, we can extract the relevant pieces of information from the entire background. These pieces (for instance, single sentences directly related to the topic of the question) are then sufficiently small to be directly added to the prompt. In this way, the model can learn from additional context without the need for retraining or fine-tuning. This method is sometimes described as *in-context learning* (Shen *et al.*, 2023).

To provide access to this functionality in ChatGSE, we add a "Document Summarisation" tab to the platform that allows the upload of text documents to be added to a vector database, which then can be queried to add contextual information to the prompt sent to the primary model. This contextual information is transparently displayed in the main chat window. Since this functionality requires a connection to a vector database system, we provide modular

connectivity to several standard vector database providers, such as Pinecone, Weaviate, or Milvus.

Supplementary Note 4: Cell Type Annotation

A common repetitive task in bioinformatics is to annotate single-cell datasets with cell type labels. This task is usually performed by a human expert, who will look at the expression of marker genes and assign a cell type label based on their knowledge of the cell types present in the tissue of interest. LLMs have been shown to be able to perform this task with high accuracy, and can be used to automate cell type annotation with minimal human input (Hou and Ji, 2023).

We propose to combine LLM inference on cell types with the storage solutions provided by ChatGSE: using a vector database to store embeddings of cells for a more streamlined workflow, and using a traditional database on the basis of BioCypher (Lobentanzer *et al.*, 2022) to inject prior knowledge into the process, as well as store the intermediate decisions of the model/user.

Using the graphical user interface of ChatGSE, we can provide recommendations of inferred cell types to the human user, such that the ultimate decision about a cell type annotation remains with the domain expert, while the tedious aspects of annotation are reduced significantly. If the model reaches a threshold of perfect annotation for any cell type (e.g., a >99% success rate in more than 50 cells), the user can decide to trust the model in instances of this cell type and fully automate the annotation in these instances. Similarly, confidence metrics can be used to trigger user input only at cell type inferences that are not straightforward.

Supplementary Note 5: Causal Inference

More recent models have shown considerable capacity for common sense reasoning, in particular, GPT-4 (OpenAI, 2023). There is an unmet need to compare the reasoning outcomes of LLMs to other, more traditional modes of inference, such as the do-calculus (Pearl, 2009), logic models (Le Novère, 2015), and semantic reasoning approaches (Hu *et al.*, 2019). With ChatGSE, we provide a Python platform for the side-by-side application of reasoning algorithms, in the “Causal Inference” tab.

The ability to connect to flexible databases created by BioCypher enables the representation of the same basic knowledge, but tailored to each individual reasoning mode (for instance, a

labelled property graph as input for the logic model, an RDF graph for the semantic reasoner, and a vector database for the LLM). This increases the ease-of-use as well as the reproducibility of benchmarking the reasoning abilities of different algorithms.

Supplementary Note 6: Literature Reference Database

LLMs are very good feature extractors. In addition to the correcting agent described in Supplementary Note 2, a “literature agent” could additionally be used to ameliorate the occasionally untruthful statements of the primary model. Given a response from said model, the literature agent is tasked with extracting references to academic papers from the response (if it contains such), and validating the existence of the claimed reference as well as its attributes (such as title and digital object identifier). The validation of extracted article authors, titles, and identifiers can be performed by a simple search in a connected database of scientific publications.

Supplementary Note 7: Experimental Design

Experimental design is a crucial step in any biological experiment. However, it can be a subtle and complex task, requiring a deep understanding of the biological system under study as well as statistical and computational expertise. LLMs can potentially fill the gaps that exist in most research groups, which traditionally focus on either the biological or the statistical aspects of experimental design.

Similar to the general chat functionality, we provide facilities to upload experimental design plans and ask about their feasibility in the biological or statistical context in an “Experimental Design” tab. Coupled with a suitable knowledge graph (for instance containing methods literature) and/or a document summarisation of relevant articles, ChatGSE can be a simple and effective tool to consider the design of an experiment from multiple perspectives.

Bibliography

Hou, W. and Ji, Z. (2023) “Reference-free and cost-effective automated cell type annotation with GPT-4 in single-cell RNA-seq analysis.,” *BioRxiv* [Preprint]. doi:10.1101/2023.04.16.537094.

Hu, P. *et al.* (2019) “Datalog Reasoning over Compressed RDF Knowledge Bases,” in *Proceedings of the 28th ACM International Conference on Information and Knowledge Management - CIKM '19. the 28th ACM International Conference*, New York, New York, USA: ACM Press, pp. 2065–2068. doi:10.1145/3357384.3358147.

Kojima, T. *et al.* (2022) “Large Language Models are Zero-Shot Reasoners,” *arXiv* [Preprint]. doi:10.48550/arxiv.2205.11916.

LangChain (2023) *LangChain Documentation*. Available at: <https://python.langchain.com> (Accessed: May 4, 2023).

Le Novère, N. (2015) “Quantitative and logic modelling of molecular and gene networks.,” *Nature Reviews. Genetics*, 16(3), pp. 146–158. doi:10.1038/nrg3885.

Lin, D. *et al.* (2020) “The TRUST Principles for digital repositories.,” *Scientific data*, 7(1), p. 144. doi:10.1038/s41597-020-0486-7.

Lobentanzer, S. *et al.* (2022) “Democratising Knowledge Representation with BioCypher,” *arXiv* [Preprint]. doi:10.48550/arxiv.2212.13543.

OpenAI (2023) “GPT-4 Technical Report,” *arXiv* [Preprint]. doi:10.48550/arxiv.2303.08774.

Pearl, J. (2009) *Causality: Models, reasoning, and inference*. 2nd ed. Cambridge, U.K: Cambridge University Press.

Shen, Y. *et al.* (2023) “HuggingGPT: Solving AI Tasks with ChatGPT and its Friends in HuggingFace,” *arXiv* [Preprint]. doi:10.48550/arxiv.2303.17580.

Vert, J.-P. (2023) “How will generative AI disrupt data science in drug discovery?,” *Nature Biotechnology* [Preprint]. doi:10.1038/s41587-023-01789-6.