

Network for knowledge Organization (NEKO): An AI knowledge mining workflow for synthetic biology research

Zhengyang Xiao^a, Himadri B. Pakrasi^b, Yixin Chen^{c,*}, Yinjie J. Tang^{a,*}

^a Department of Energy, Environment, and Chemical Engineering, Washington University in St. Louis, St. Louis, MO, 63130, United States

^b Department of Biology, Washington University in St. Louis, St. Louis, MO, 63130, United States

^c Department of Computer Science, Washington University in St. Louis, St. Louis, MO, 63130, United States

ARTICLE INFO

Keywords:

Foundation model
Large language model
Qwen
Retrieval augmented generation
Knowledge graph

ABSTRACT

Large language models (LLMs) can complete general scientific question-and-answer, yet they are constrained by their pretraining cut-off dates and lack the ability to provide specific, cited scientific knowledge. Here, we introduce Network for Knowledge Organization (NEKO), a workflow that uses LLM Qwen to extract knowledge through scientific literature text mining. When user inputs a keyword of interest, NEKO can generate knowledge graphs to link bioinformation entities and produce comprehensive summaries from PubMed search. NEKO significantly enhance LLM ability and has immediate applications in daily academic tasks such as education of young scientists, literature review, paper writing, experiment planning/troubleshooting, and new ideas/hypothesis generation. We exemplified this workflow's applicability through several case studies on yeast fermentation and cyanobacterial biorefinery. NEKO's output is more informative, specific, and actionable than GPT-4's zero-shot Q&A. NEKO offers flexible, lightweight local deployment options. NEKO democratizes artificial intelligence (AI) tools, making scientific foundation model more accessible to researchers without excessive computational power.

1. Introduction

Biomanufacturing has a potential US market of over 30 billion dollars annually (2023 Government Accountability Office report). However, a primary challenge that needs urging solutions is the high cost to develop cellular factories that meet commercially relevant performance. Biological systems are complex and there are many important levers (e. g. genetic regulations, enzyme functions, cellular metabolism, and extracellular conditions) that need to be analyzed and tuned to engineer a desired phenotype (i.e., design-build-test-learn cycles for bioprocess development) (Liao et al., 2022). Therefore, a holistic knowledgebase for bioprocess development is essential. Currently, vast amount of synthetic biology and biomanufacturing literatures have been published. A PubMed search of “synthetic biology or metabolic engineering” queries generated ~125,000 publications, which offers wealthy bioinformation. The pressing need for efficient knowledge integration has coincided with the advent of large language models (LLMs) (Bai et al., 2023; OpenAI, 2023), which now facilitate rapid text information processing. Notably, recent advances in retrieval-augmented generation (RAG) have

made LLMs powerful tools for information processing and knowledge mining from text (Jiang et al., 2023; Lewis et al., 2020; Ni et al., 2024). While general-purpose pretrained LLMs can provide answers to scientific inquiries (Stribling et al., 2024), encapsulating infinite knowledge within LLM's finite parameter space remains an inherent challenge (Martino et al.; Sun et al., 2023). Therefore, LLM needs a database and a RAG pipeline to store and distillate factual knowledge. Knowledge graphs have emerged as a promising solution, offering an intuitive representation and knowledge synthesis that is interpretable by both LLMs and humans (Pan et al., 2024; Yang et al., 2023). The synergy between LLM and knowledge graph construction can help scientists collect key information and make informed decisions on research.

This study introduces NEKO, a Knowledge Graph workflow to enhance generative AI for extracting trustworthy knowledge from the literature. Exemplified through case studies, NEKO demonstrates broad applicability in data collections and information distillations from diverse reliable sources (Fig. 1). First, NEKO helps scholars to quickly inquire into a research topic or scientific question by combining, distillate and organize information from multiple articles. Second,

* Corresponding author.

** Corresponding author.

E-mail addresses: yichen25@wustl.edu (Y. Chen), yinjie.tang@wustl.edu (Y.J. Tang).

<https://doi.org/10.1016/j.ymben.2024.11.006>

Received 8 July 2024; Received in revised form 30 September 2024; Accepted 17 November 2024

Available online 21 November 2024

1096-7176/© 2024 The Author(s). Published by Elsevier Inc. on behalf of International Metabolic Engineering Society. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

NEKO can solve logic problems and generate new ideas or hypothesis to be tested by connecting key concepts and synthesizing holistic knowledge from closely relevant articles. Third, it can help experiment planning and troubleshooting by presenting all past considerations in a systematic way. Therefore, NEKO can not only accelerate the learning and training process of synthetic biology researchers, but also assist them to write literature review, prepare research proposal, and optimize experimental design. NEKO is compatible with any instruction-following LLM, including proprietary models like GPT-4 (from OpenAI) (OpenAI, 2023) and open-source alternatives like Qwen (from Alibaba) (Bai et al., 2023). Recently, the concept of foundation model (Bommasani et al., 2021; Qin et al., 2023) becomes a new research interest due to its specialized knowledge in various fields. NEKO is a lightweight application that anyone can construct synthetic biology foundation model according to specific needs. Developed using beginner-friendly Python code, the entire workflow is readily accessible via GitHub.

2. Methods

Code availability. The codes and examples used in this study are deposited in this GitHub repository: <https://github.com/xiao-zhengyang/NEKO>.

Web-based literature search. This workflow starts with online literature search. This step is compatible with any web-based databases with application programming interface (API) services. In this study, we used PubMed and arXiv as examples. User of NEKO input a search keyword into corresponding API, and the API would return the article

title and abstract. The search result was saved in an article list excel file. This workflow also applies to PDF files. After downloading research articles in PDF, the files were read and divided into 1000 words segments. Then these texts were passed to the next step.

LLM text processing. LLMs were used to process the text from API search or from downloaded PDFs. LLMs with strong retrieval augmented generation (RAG) capabilities were recommended. In this study, we demonstrated with Qwen1.5 (developed by Alibaba, China), but we also provided codes compatible with GPT-4. Qwen1.5-72B-chat were downloaded from Hugging Face and deployed on a high-performance computing cluster with Nvidia A100 GPUs. The text from previous step were sequentially input into LLM with the following system prompt:

“You are specialized for analyzing scientific paper abstracts, focusing on identifying entities causal relationships related to biological studies, such as performance, species, genes, methods of genetic engineering, enzymes, proteins, and bioprocess conditions (e.g., growth conditions). You output the identified causal relationships between entities in combination pairs. The output strictly follows the format: (Entity A, Entity B), (Entity C, Entity D) ... with no additional text.”

User can customize this prompt based on specific needs. This prompt is a general prompt for synthetic biology studies. We also included a specific prompt example focused on gene expression/deletion and cell response in our GitHub page. The response from LLM was processed by the Word2Vec method (Rong, 2014). A word embedding from sentence-transformers (Reimers and Gurevych, 2019) was used to identify and combine entities with same meanings. The processed lists of

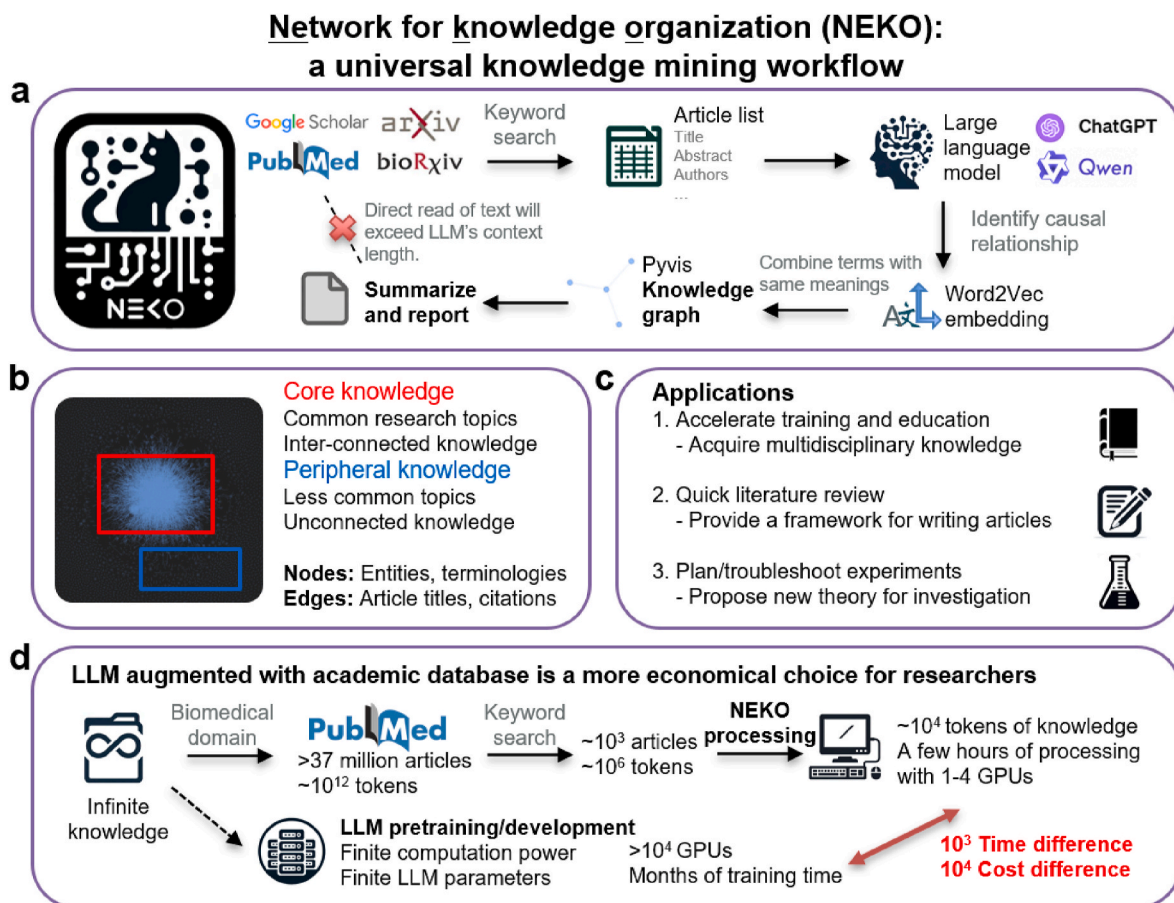


Fig. 1. NEKO: a universal knowledge mining workflow. (a) NEKO's logo and workflow illustration. (b) Knowledge graph has two regions, core and peripheral. (c) Applications in academic tasks. (d) LLM augmented with academic database is a more economical choice for researchers. Pretraining LLM requires thousands of GPUs and months of training time, while our RAG work flow requires minimal GPU resource.

entities were saved and passed to the next step.

Knowledge graph visualization. Each pair of entities indicates their causal relationship. An open-source Pyvis package (Perrone et al., 2020) was used to construct knowledge graph. The nodes represent entities, while edges (lines connecting nodes) are labeled as the article title as citation/source.

Keyword search and summarization. After obtaining the

knowledge graph, it is recommended to search for a keyword and filter the knowledge graph for cleaner visualization. The entities/nodes related to the search keyword were extracted and input into LLM for summary report generation. LLMs with good reasoning ability, such as GPT-4 and Qwen1.5–32B/72B, are recommended.

Problem solving and hypothesis generations. After obtaining summary reports about relevant concepts, these reports were combined.



These are the terms related to *Y. lipolytica* β -carotene production, categorize them and write a summary report.



NEKO knowledge mining from >200 PDF files

Sample report: β -carotene production in *Y. lipolytica*

The terms provided are related to the field of metabolic engineering and synthetic biology, specifically focusing on the optimization and enhancement of β -carotene production in microorganisms, particularly *Yarrowia lipolytica*. β -carotene is a precursor to vitamin A and has significant nutritional and pharmaceutical importance. Here's a categorized summary of the key concepts:

- 1. Microbial Hosts:** Mainly *Yarrowia lipolytica*, but also includes other species like *Xanthophyllomyces dendrorhous*, *Escherichia coli*, *Saccharomyces cerevisiae*, and *Blakeslea trispora*, which are used as platforms for β -carotene biosynthesis.
- 2. Genetic Modifications:** Techniques such as gene deletion (Δ CLA4, Δ MHY1), overexpression (ERG13, ERG9, HMG, CarB, CarRP), promoter optimization, integration of multiple copies of pathway genes, and introduction of heterologous genes (from *Mucor circinelloides*, cyanobacteria) are employed to enhance β -carotene production.
- 3. Pathway Genes and Enzymes:** Key genes and enzymes involved in β -carotene synthesis include HMG, GGPS (Geranylgeranyl pyrophosphate synthase), crtI, crtYB, crtZ, phytoene dehydrogenase (carB), and β -carotene ketolase (crtW).
- 4. Metabolic Engineering Strategies:** These include nitrogen limitation, deletion of competing pathways (beta-oxidation), overexpression of central metabolic modules, and manipulation of MVA (mevalonate) and MEP (methylerythritol phosphate) pathways.
- 5. Fermentation Conditions:** Factors like high glucose feeding rate, fed-batch fermentation, carbon-to-nitrogen ratio, temperature regulation, and use of different media (YPD, synthetic medium, glucose, xylose) are manipulated to optimize β -carotene production.
- 6. Yield and Titer:** Achieved improvements are mentioned, such as a 49% increase in β -carotene titers, 272% increase in β -carotene production, and β -carotene yields reaching 6.5 g/L.
- 7. Recombinant Strains and Transformation:** Genetic engineering tools are used to create recombinant strains, often through iterative integration or homologous recombination, to achieve higher β -carotene content.
- 8. Bioreactors and Analytical Techniques:** Flask cultures, fermentation broth, and β -carotene extraction using solvents like acetone are mentioned, along with techniques like photometric measurement, 2,3-oxidosqualene analysis, and carotenoid profiling.
- 9. Lipid and Metabolite Accumulation:** β -carotene sequestration in lipid bodies, lipid overproduction, and squalene overproduction are noted as strategies to enhance carotenoid accumulation.
- 10. Synthetic Biology Applications:** Synthetic biology approaches are employed for pathway engineering, promoter shuffling, and the creation of carotenoid platform strains.

Overall, these terms represent a comprehensive exploration of the strategies and techniques used to improve β -carotene production in microorganisms, highlighting the importance of metabolic engineering and synthetic biology in optimizing bioconversion processes. The research focuses on strain optimization, pathway manipulation, and fermentation conditions to maximize β -carotene yields while considering factors like cellular stress response and morphological engineering.

Highlights: Quantifiable knowledge, Gene targets, Strain engineering strategies, Cultivation conditions, Reliable sources and citations

Fig. 2. Compare NEKO and ChatGPT zero-shot Q&A. The topic is to write a summary of production of β -carotene in *Y. lipolytica*.



Write a report about optimizing β -carotene production in *Y. lipolytica*.



ChatGPT Zero-shot direct Q&A

Producing β -carotene in *Yarrowia lipolytica* represents a compelling example of the application of metabolic engineering and synthetic biology in microorganisms to produce high-value compounds. This report will delve into various aspects of β -carotene production in *Yarrowia lipolytica*, including its significance, metabolic pathways involved, genetic engineering strategies, cultivation conditions, and challenges faced in scaling up production.

1. Introduction to β -Carotene and *Yarrowia lipolytica*

β -Carotene is a highly valued carotenoid with widespread applications in the food, pharmaceutical, and cosmetic industries due to its antioxidant properties and role as a vitamin A precursor. *Yarrowia lipolytica* is a non-conventional yeast that has garnered attention as an efficient host for the production of various compounds, including lipids, citric acid, and high-value carotenoids like β -carotene. Its ability to utilize a wide range of substrates, robustness in different cultivation conditions, and well-characterized genetic background make it an attractive platform for metabolic engineering.

2. Metabolic Pathways for β -Carotene Synthesis

The biosynthesis of β -carotene in *Yarrowia lipolytica* involves the mevalonate pathway for the synthesis of isoprenoid precursors, which are then converted into β -carotene through several enzymatic steps. Key enzymes involved include geranylgeranyl diphosphate synthase (GGPPS), phytoene synthase (PSY), phytoene desaturase (PDS), and lycopene β -cyclase (LCY). Enhancing the flux through these pathways and balancing the expression levels of these enzymes are crucial for optimizing β -carotene production.

3. Genetic Engineering Strategies

Several genetic engineering strategies have been employed to increase β -carotene production in *Yarrowia lipolytica*:

Overexpression of Biosynthetic Genes: Targeted overexpression of genes involved in the β -carotene biosynthesis pathway, such as PSY, PDS, and LCY, has been shown to enhance production.

CRISPR/Cas9 Mediated Genome Editing: The CRISPR/Cas9 system has been utilized for precise genomic modifications, including knockout of competing pathways and integration of biosynthetic gene clusters to increase β -carotene yield.

Promoter Engineering: Utilizing strong, constitutive, or inducible promoters to drive the expression of biosynthetic genes has proven effective in optimizing production.

Metabolic Flux Analysis: Computational models to analyze and redirect metabolic flux towards β -carotene synthesis have been developed, identifying bottlenecks and targets for genetic engineering.

4. Cultivation Conditions

Optimizing the cultivation conditions, such as the carbon source, pH, temperature, and aeration, is essential for maximizing β -carotene production. *Yarrowia lipolytica* can utilize a variety of carbon sources, including glucose, glycerol, and industrial waste streams. Controlled batch, fed-batch, or continuous cultures under optimized conditions have been developed to enhance yield and productivity.

5. Challenges and Future Perspectives

Scaling up β -carotene production in *Yarrowia lipolytica* to industrial levels poses several challenges:

Product Inhibition and Toxicity: High levels of β -carotene can be toxic to the cells or inhibit growth, necessitating strategies to balance production and cell viability.

Economic Feasibility: The economic viability of the process depends on substrate cost, product yield, and the efficiency of downstream processing. Utilizing low-cost substrates and improving extraction methods are key areas of research.

Regulatory and Safety Considerations: Meeting food and pharmaceutical grade standards for β -carotene produced in *Yarrowia lipolytica* requires thorough safety and regulatory assessment.

6. Conclusion

The production of β -carotene in *Yarrowia lipolytica* exemplifies the potential of metabolic engineering in creating sustainable and economically viable bioprocesses for high-value compounds. Ongoing research and technological advancements in genetic engineering, bioprocess optimization, and systems biology continue to unlock the potential of *Yarrowia lipolytica* as a versatile platform for industrial biotechnology applications.

Fig. 2. (continued).

NEKO leverages the inference capabilities of LLMs to propose novel research plans and hypotheses. In the cyanobacterial nitrogen fixation example, concepts of *Anabaena* sp. nitrogen fixation and common bio-fertilizer formulation were input to Qwen1.5-72B. The LLM inference gives a brief description of potential research directions and hypotheses, which are evaluated by human researchers.

3. Results

NEKO (Fig. 1a) is a generative AI tool which can be a companion for synthetic biology researchers. Using *Y. lipolytica* PubMed abstract knowledge mining as a demonstration (Fig. 1b), we plotted all information on a single knowledge graph. We identified two main regions: the core knowledge region at the center, which encapsulates common

research topics and interconnected knowledge from various articles, critical for tasks such as process optimization and scale-up in bio-production applications. The peripheral knowledge region contains fewer common topics. Researchers are encouraged to integrate this peripheral knowledge with the core by conducting experiments. We demonstrated NEKO's applicability in academic tasks through several case studies (Fig. 1C). This method of augmenting LLM with existing academic infrastructure is a more economical choice. Users can locally deploy NEKO and quantized LLM even on consumer-level desktops, with as little as 24 GB GPU RAM (Nvidia RTX 3090 or 4090). Compared to pretraining science LLM, NEKO saves tremendous time and monetary costs (Fig. 1d).

We choose Qwen as an example LLM in this study due to its exceptional Retrieval Augmented Generation (RAG) ability. With 14B parameters, Qwen can achieve comparable RAG score as GPT-4 (QwenTeam, 2024). The use of local LLM is more economical and safer for copyright materials. It can bypass the request rate limitation of OpenAI API. Interestingly, during our practical deployment of NEKO, we found that Qwen1.5 is as good as GPT-4 in terms of instruction following, scientific reasoning, and the comprehensiveness of responses. To date, there was few Qwen application journal articles, leading us to conclude that Qwen1.5 is an underappreciated yet reliable LLM for practical applications. The field of LLM/AI is evolving very fast. During publishing our article, Qwen2.5 was released with better RAG capabilities compared to its predecessor (Yang et al., 2024). User can easily substitute new LLM using the same Hugging Face transformers framework.

3.1. Case study 1: knowledge acquisition and distillation on β -carotene production in *Y. lipolytica*

NEKO can help researchers quickly acquire up-to-date knowledge about one research topic. To illustrate, NEKO rapidly read more than 200 research articles in PDF files and produce a knowledge base for the oleaginous yeast *Y. lipolytica*. When searching for β -carotene production in *Y. lipolytica*, NEKO's analysis pinpointed relevant nodes, with the summary presented in Fig. 2. NEKO also produced a knowledge graph illustrating the connections between knowledge entities (Fig. S1). Users can click one line connecting two nodes and view the title of the article containing the knowledge connection. Compared to ChatGPT-4 zero-shot Q&A, NEKO gave 200% more gene targets, 200% more strain engineering strategies, and 57% more bioprocess cultivation conditions, with knowledge from 37 reliable peer-reviewed sources (Table 1). We also did a negative control by using base LLM Qwen1.5 to write a summary about β -carotene production in *Y. lipolytica* (Fig. S2). Qwen1.5 only gave general optimization suggestions like GPT-4 did. This demonstrates that without NEKO knowledge mining workflow, the LLM does not have specific knowledge about β -carotene production. NEKO augments the LLM by providing a light-weight knowledge database.

Another information-intensive academic task is experiment planning and troubleshooting. We illustrate this through NEKO's ability to synthesize experiment procedures and actionable insights. For example, when aiming to genetically engineer *Y. lipolytica*, user needs to obtain the procedures for strain transformation. NEKO provided a detailed, step-by-step methodological guidance (Fig. S3), including specific DNA amount used (2 μ g), procedure name (LiAc or electroporation), selection

markers (hygromycin resistance or *ura3d4* defective allele), etc. Take another example, assume the user is having trouble expressing GGPPS, an essential enzyme in β -carotene synthesis pathway. NEKO can search for information regarding to GGPPS and compile a report with suggestions to troubleshoot this scenario (Fig. S4). We noticed that GPT-4's suggestions were heavily templated. When given another β -carotene production gene *CarRP* (Fig. S5), GPT-4's suggestions were very similar to GGPPS's case. Contrasting NEKO's outputs with GPT-4's templated responses highlighted NEKO's ability to generate pertinent suggestions based on mined knowledge.

3.2. Case study 2: literature review of non-model species *Rhodospiridium toruloides*

Literature review is one of the time-consuming academic tasks, and we demonstrated NEKO's quick application in review paper writing in this case study. *Rhodospiridium toruloides* is gaining research attention recently for its high lipid content and native carotenoid production. We applied NEKO to *R. toruloides* article abstracts (total 392 articles) on PubMed. By visualizing the knowledge graph and comparing with *Y. lipolytica* (Fig. S6), *R. toruloides* literature is still at its early stage of research. The knowledge from studies is disconnected from each other. After LLM summarization (Fig. 3 and Fig. S7), NEKO identified trending research areas, including genetic engineering, metabolic pathway, adaptive evolution experiments, and biochemical production. Moreover, NEKO enables users to refine their searches by specifying target products. For instance, focusing on lipid production yielded a detailed breakdown of the research topic (Fig. S8). NEKO excels at providing details like strain engineering gene targets and fermentation conditions. Overall, NEKO provides a framework for users to organize the knowledge and structure their review paper writing. By further inspecting the article list and entity list, users can produce rich and precise literature study.

3.3. Case 3: hypotheses generation and new theory for exploration

Being able to summarize literature information is part of a competent synthetic biology foundation model. In this case study, we demonstrate that LLM can do inference and propose new hypotheses given relevant concepts. In essence, relevant concepts were "prompts" for LLM to generate new experiment plan. For example, we are interested in engineering a nitrogen-fixing cyanobacteria *Anabaena* sp. to produce guanidine and urea as bio-fertilizers. We combined relevant concepts of *Anabaena* sp. nitrogen fixation and common bio-fertilizer formulations. The LLM Qwen1.5-72B can do inference and write a brief research proposal and experiment plan (Fig. 4, Fig. S9). Interestingly, same as our research plan, this workflow correctly identified gene targets as nitrogenase, arginine decarboxylase (ethylene-forming enzyme), and urease. Compared to GPT-4 zero-shot Q&A in Fig. S10, NEKO pointed out that urea and guanidine are more stable than ammonia, and they are suitable for long term storage and transportation. The proposed experiments for cyanobacterial production of guanidine are detailed, incorporating tools like CRISPR-Cas9, metabolic engineering, and dynamic expression control systems. In terms of project deliverables, NEKO's report had a thorough plan that includes *in vitro* and *in vivo* characterization, field trials, and an economic/environmental evaluation, offering a more holistic project design than GPT-4's zero-shot Q&A.

4. Discussion

NEKO can compile massive literature reports, fill knowledge gap, remove redundant data, and connect information streams, which can be used to collect both features and targets from literature for developing standardized datasets (Xiao et al., 2023). NEKO can be widely used for Synthetic biology research. It distinguishes itself from other knowledge-based Q&A platforms by overcoming the context length

Table 1
Compare NEKO and ChatGPT zero-shot Q&A.

	NEKO workflow	GPT-4 zero-shot Q&A
Quantifiable knowledge	3	0
Gene targets	12	4
Strain engineering strategies	12	4
Cultivation conditions	11	7
Reliable sources and citations	37	0

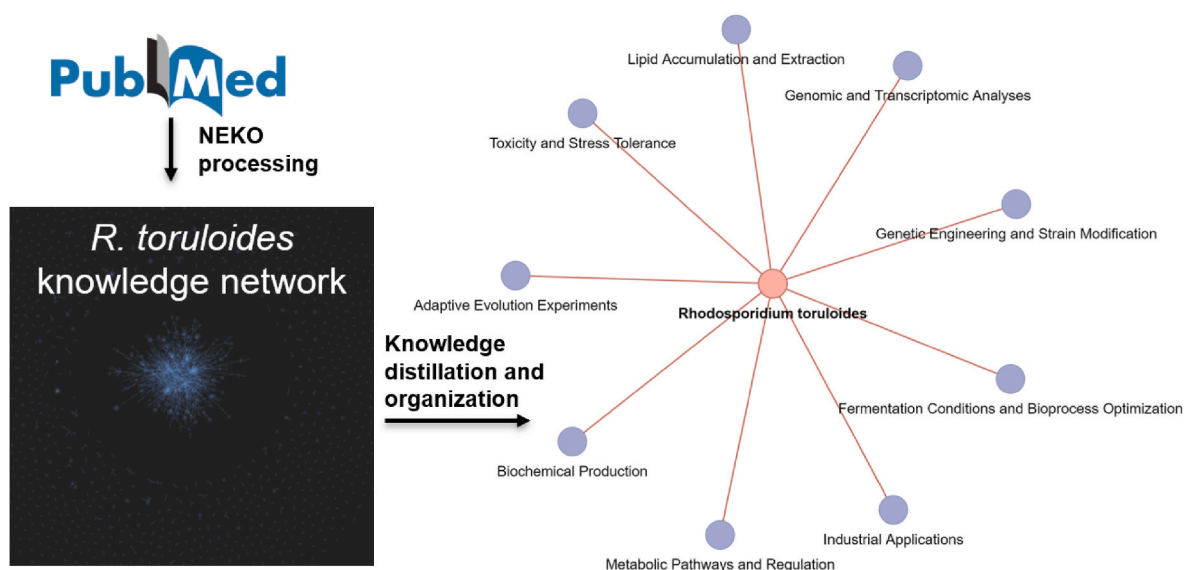


Fig. 3. NEKO processes *R. toruloides* publications on PubMed and produced a quick research topics review.

Topic: How to engineer *Anabaena* sp. to produce biofertilizer as urea and guanidine?

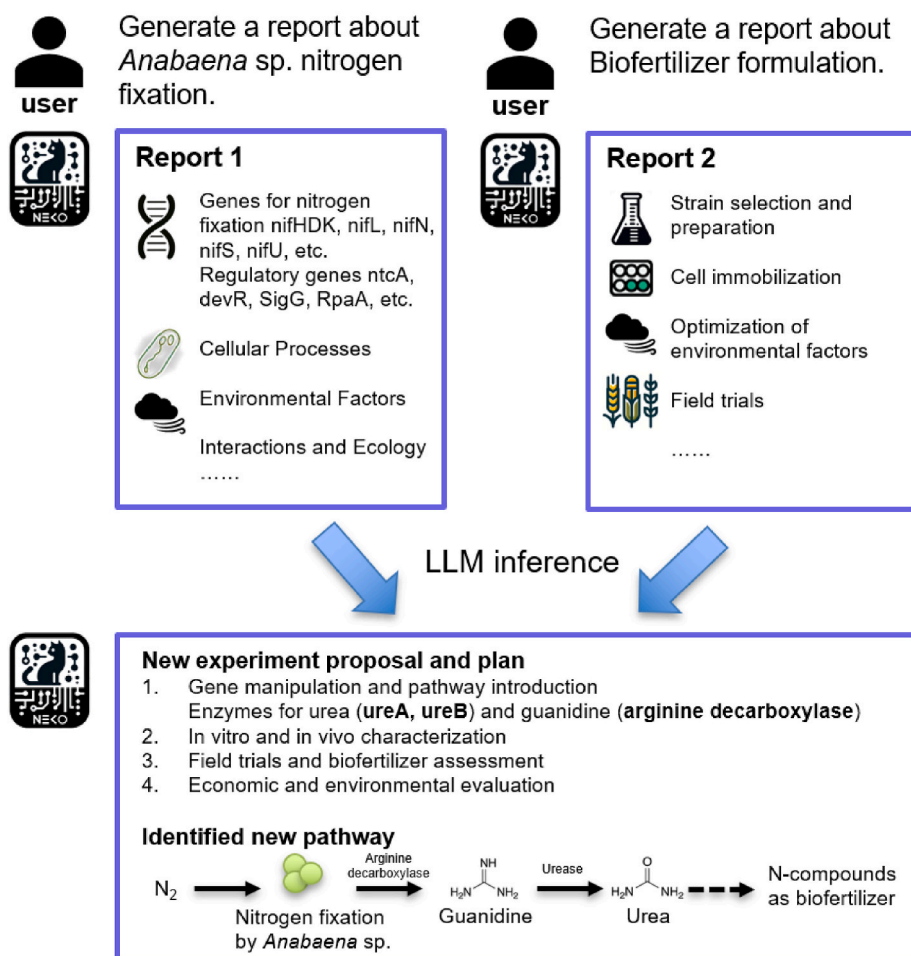


Fig. 4. Based on the knowledge mining of NEKO, LLM can perform inference and give new research plan about new compounds production pathway.

limitations of LLMs without directly reading article texts. All NLP tools used in this study are off-the-shelf, without the need for pretraining. NEKO is compatible with any LLM, although the effectiveness is influenced by the LLM's ability to follow instructions and its RAG capabilities. Generally, LLMs with a higher number of parameters tend to perform better in identifying causal relationships and generating logical summary reports. When using Qwen1.5-72B running on 4*A100 GPUs, NEKO can process about 300 PubMed abstracts in 1 h. It seems that GPU memory (GPU RAM) and CPU computation are the limiting factors. If using GPT-4, user's API request rate is subject to OpenAI policies. Our preliminary findings indicate that NEKO typically offers more factual information compared to GPT-4. However, users are advised to verify important information. We also applied this workflow to other disciplines such as computer science. By modifying the knowledge mining prompt, NEKO was capable of produce a mini review of recent knowledge retrieval studies on arXiv (Fig. S11). In general, NEKO works with any online article database with API services. Users need to modify the knowledge mining prompt based on their specific needs.

As academic research grows in complexity and increasingly spans multiple disciplines, researchers are investing substantial time in literature review and knowledge acquisition. NEKO addresses this challenge by streamlining literature study tasks, enabling knowledge synthesis from literature (Kastner et al., 2012; Whittemore et al., 2014). This trend towards continuously interacting research with LLM represents a new common. We envision a roadmap for three tiers of science AI and LLM knowledge synthesis (Fig. 5). Currently most pretrained LLMs are at tier 2, and they can answer general scientific questions based on pretraining text corpus. NEKO is at tier 1. It can output summary reports based on literature search and give reliable cited knowledge. To some extent, it can help guide experiments in the physical world. It serves as a “human prompter” by brainstorming potential experiment plans. The top tier 0 is one of the ultimate development goals for science foundation models, which can provide quantitative feasibility score based on human evaluation and experiment feedback. However, there are challenges associated with this goal. First, data generated by experiments are multi-modal, and it is difficult to integrate multi-omics datasets by text LLM only (Bi et al., 2024; Xu et al., 2024). Multi-modal LLM should be able to process images, videos, and data plots, minimizing manual efforts for data curation. Second, the inherent variability of experiment (Gilman et al., 2021) requires standardized workflow to facilitate effective feedback loops between practical research and LLM inferences. Third, biological systems are chaotic and subject to the influence of experiment initial conditions (Elowitz et al., 2002; Paulsson, 2005; Raj and Van Oudenaarden, 2008), so the notion of creating an omniscient AI Laplace's demon is not feasible. A pragmatic approach acknowledges these limitations, focusing on the integration of AI within a defined set of constraints and assumptions. The key is to integrate knowledge mining workflow with physical world, and NEKO provides a lightweight deployment example for integrating LLM with routine education and academic tasks. During the drafting of our article, Microsoft released GraphRAG, a state-of-the-art knowledge graph generation tool. While NEKO and GraphRAG share similar concepts, we emphasize NEKO's simplicity and its advantage for local deployment. NEKO's simplicity brings trustworthiness and all steps in the workflow is fully traceable. More importantly, GraphRAG or OpenAI are not easily accessible in certain countries (e.g., China). In contrast, NEKO uses an open source LLM Qwen, which may serve broader synthetic biology community with low costs.

5. Conclusion

NEKO is a low-cost and in-house pipeline for knowledge graph construction. It improves LLM performance by extracting synthetic biology information and connecting relevant knowledge entities from the literature. NEKO also shows potentials to design synthetic biology experiments and generate new research hypotheses. In the future, we

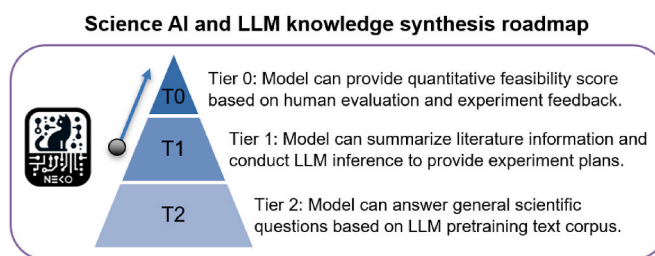


Fig. 5. Proposed science AI and LLM knowledge synthesis roadmap.

will test this knowledge mining workflow coupled with experiment feedback, facilitating design-build-test-learn (DBTL) cycles and offering a foundation model for synthetic biology applications.

CRedit authorship contribution statement

Zhengyang Xiao: Writing – original draft, Visualization, Software, Methodology, Investigation, Formal analysis, Data curation. **Himadri B. Pakrasi:** Writing – review & editing, Supervision, Project administration, Funding acquisition, Conceptualization. **Yixin Chen:** Writing – review & editing, Project administration, Funding acquisition, Conceptualization. **Yinjie J. Tang:** Writing – review & editing, Supervision, Resources, Project administration, Funding acquisition, Conceptualization.

Conflict of interest statement

The authors declare no conflict of interest.

Acknowledgements

This study is funded by United States NSF award number 2225809 and DOE Energy Earthshots award number DE-SC0024702.

Appendix A. Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.ymben.2024.11.006>.

Data availability

Data will be made available on request.

References

- Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., Fan, Y., Ge, W., Han, Y., Huang, F., 2023. Qwen Technical Report. arXiv Preprint arXiv:2309.16609.
- Bi, Z., Dip, S.A., Hajialigol, D., Kommu, S., Liu, H., Lu, M., Wang, X., 2024. AI for biomedicine in the era of large language models. arXiv preprint arXiv:2403.15673.
- Bommasani, R., Hudson, D.A., Adeli, E., Altman, R., Arora, S., von Arx, S., Bernstein, M. S., Bohg, J., Bosselut, A., Brunskill, E., 2021. On the opportunities and risks of foundation models. arXiv preprint arXiv:2108.07258.
- Elowitz, M.B., Levine, A.J., Siggia, E.D., Swain, P.S., 2002. Stochastic gene expression in a single cell. *Science* 297, 1183–1186.
- Gilman, J., Walls, L., Bandiera, L., Menolascina, F., 2021. Statistical design of experiments for synthetic biology. *ACS Synth. Biol.* 10, 1–18.
- Jiang, Z., Xu, F.F., Gao, L., Sun, Z., Liu, Q., Dwivedi-Yu, J., Yang, Y., Callan, J., Neubig, G., 2023. Active Retrieval Augmented Generation arXiv preprint arXiv:2305.06983.
- Kastner, M., Tricco, A.C., Soobiah, C., Lillie, E., Perrier, L., Horsley, T., Welch, V., Cogo, E., Antony, J., Straus, S.E., 2012. What is the most appropriate knowledge synthesis method to conduct a review? Protocol for a scoping review. *BMC Med. Res. Methodol.* 12, 1–10.
- Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W.-t., Rocktäschel, T., 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Adv. Neural Inf. Process. Syst.* 33, 9459–9474.
- Liao, X., Ma, H., Tang, Y.J., 2022. Artificial intelligence: a solution to involution of design-build-test-learn cycle. *Curr. Opin. Biotechnol.* 75, 102712.
- Martino, A., Iannelli, M., Truong, C., Knowledge Injection to Counter Large Language Model (LLM) Hallucination. Springer, pp. 182–185.

- Ni, S., Bi, K., Guo, J., Cheng, X., 2024. When do LLMs need retrieval augmentation? Mitigating LLMs' overconfidence helps retrieval augmentation. arXiv preprint arXiv:2402.11457.
- OpenAI, 2023. GPT-4 Technical report. arXiv preprint. 2303.08774.
- Pan, S., Luo, L., Wang, Y., Chen, C., Wang, J., Wu, X., 2024. Unifying large language models and knowledge graphs: a roadmap. IEEE Trans. Knowl. Data Eng.
- Paulsson, J., 2005. Models of stochastic gene expression. Phys. Life Rev. 2, 157–175.
- Perrone, G., Unpingco, J., Lu, H.-m., 2020. Network Visualizations with Pyvis and VisJS arXiv preprint arXiv:2006.04951.
- Qin, Y., Hu, S., Lin, Y., Chen, W., Ding, N., Cui, G., Zeng, Z., Huang, Y., Xiao, C., Han, C., 2023. Tool Learning with Foundation Models arXiv preprint arXiv:2304.08354.
- QwenTeam, 2024. Introducing Qwen1.5. <https://qwenlm.github.io/blog/qwen1.5/>.
- Raj, A., Van Oudenaarden, A., 2008. Nature, nurture, or chance: stochastic gene expression and its consequences. Cell 135, 216–226.
- Reimers, N., Gurevych, I., 2019. Sentence-bert: Sentence Embeddings Using Siamese Bert-Networks arXiv preprint arXiv:1908.10084.
- Rong, X., 2014. word2vec parameter learning explained. arXiv preprint arXiv:1411.2738.
- Stribling, D., Xia, Y., Amer, M.K., Graim, K.S., Mulligan, C.J., Renne, R., 2024. The model student: GPT-4 performance on graduate biomedical science exams. Sci. Rep. 14, 5670.
- Sun, K., Xu, Y.E., Zha, H., Liu, Y., Dong, X.L., 2023. Head-to-tail: How knowledgeable are large language models (llm)? AKA will llms replace knowledge graphs? arXiv preprint arXiv:2308.10168.
- Whittemore, R., Chao, A., Jang, M., Minges, K.E., Park, C., 2014. Methods for knowledge synthesis: an overview. Heart Lung 43, 453–461.
- Xiao, Z., Li, W., Moon, H., Roell, G.W., Chen, Y., Tang, Y.J., 2023. Generative artificial intelligence GPT-4 accelerates knowledge mining and Machine learning for synthetic biology. ACS Synth. Biol. 12, 2973–2982.
- Xu, X., Li, J., Zhu, Z., Zhao, L., Wang, H., Song, C., Chen, Y., Zhao, Q., Yang, J., Pei, Y., 2024. A comprehensive review on synergy of multi-modal data and AI Technologies in Medical Diagnosis. Bioengineering 11, 219.
- Yang, A., Yang, B., Hui, B., Zheng, B., Yu, B., Zhou, C., Li, C., Li, C., Liu, D., Huang, F., 2024. Qwen2 Technical Report arXiv preprint arXiv:2407.10671.
- Yang, L., Chen, H., Li, Z., Ding, X., Wu, X., 2023. Chatgpt is not enough: Enhancing large language models with knowledge graphs for fact-aware language modeling. arXiv preprint arXiv:2306.11489.