

Journal Pre-proof

SeedLLM-Rice: A large language model integrated with rice biological knowledge graph

Fan Yang, Huanjun Kong, Jie Ying, Zihong Chen, Tao Luo, Wanli Jiang, Zhonghang Yuan, Zhefan Wang, Zhaona Ma, Shikuan Wang, Wanfeng Ma, Xiaoyi Wang, Xiaoying Li, Zhengyin Hu, Xiaodong Ma, Minguo Liu, Xi-Qing Wang, Fan Chen, Nanqing Dong

PII: S1674-2052(25)00172-8

DOI: <https://doi.org/10.1016/j.molp.2025.05.013>

Reference: MOLP 1916

To appear in: *MOLECULAR PLANT*

Received Date: 21 February 2025

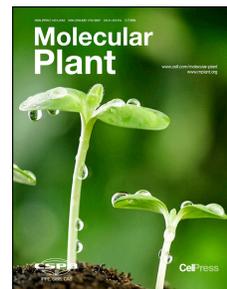
Revised Date: 8 May 2025

Accepted Date: 23 May 2025

Please cite this article as: Yang F., Kong H., Ying J., Chen Z., Luo T., Jiang W., Yuan Z., Wang Z., Ma Z., Wang S., Ma W., Wang X., Li X., Hu Z., Ma X., Liu M., Wang X.-Q., Chen F., and Dong N. (2025). SeedLLM-Rice: A large language model integrated with rice biological knowledge graph. Mol. Plant. doi: <https://doi.org/10.1016/j.molp.2025.05.013>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2025 Published by Elsevier Inc. on behalf of CAS Center for Excellence in Molecular Plant Sciences, Chinese Academy of Sciences, and Chinese Society for Plant Biology.



1
2
3 **SeedLLM-Rice: A large language model integrated with rice biological knowledge graph**
4

5 **Fan Yang^{a,1,2}, Huanjun Kong^b, Jie Ying^b, Zihong Chen^b, Tao Luo^a, Wanli Jiang^b,**
6 **Zhonghang Yuan^b, Zhefan Wang^b, Zhaona Ma^a, Shikuan Wang^a, Wanfeng Ma^a,**
7 **Xiaoyi Wang^a, Xiaoying Li^a, Zhengyin Hu^c, Xiaodong Ma^c, Minguo Liu^f, Xi-Qing Wang^{c,g,1},**
8 **Fan Chen^a, and Nanqing Dong^{b,d,1,2}**
9
10

11
12 ^aYazhouwan National Laboratory, Sanya, 572025, China.

13 ^bShanghai Artificial Intelligence Laboratory, Shanghai, 200232, China.

14 ^cState Key Laboratory of Plant Environmental Resilience, College of Biological Sciences, China
15 Agricultural University, Beijing 100193, China.

16 ^dShanghai Innovation Institute, Shanghai, 200231, China.

17 ^eNational Science Library (Chengdu), Chengdu, 610299, China

18 ^fCollege of Grassland Agriculture, Northwest A & F University, Shaanxi, 712100, China.

19 ^gFrontier Technology Research Institute, China Agricultural University, Shenzhen, 518000,
20 China.

21 ¹Correspondence to: dongnanqing@pjlab.org.cn; yangfan@yzwlab.cn; wangxq21@cau.edu.cn

22 ²These authors contributed equally.
23
24

25 ABSTRACT

26 Rice biology research involves complex decision-making, requiring researchers to navigate a
27 vast and growing body of knowledge that includes extensive literature and multiomics data. The
28 exponential increase in biological data and scientific publications has posed significant
29 challenges in efficiently extracting meaningful insights. While large language models (LLMs)
30 show promise for knowledge retrieval, their application to rice-specific research is hindered by
31 the absence of specialized models and the challenge of synthesizing multimodal data integral to
32 the field. Moreover, the lack of standardized evaluation frameworks for domain-specific tasks
33 impedes the assessment of model performance in this area. To address these challenges, we
34 introduce SeedLLM·Rice (SeedLLM), a 7-billion-parameter model trained using 1.4 million rice-
35 related publications, which represent nearly 98.24% of global rice research. Additionally, we
36 present a novel human evaluation framework designed to assess LLM performance in rice
37 biology tasks. Initial evaluations of rice-specific tasks demonstrate that SeedLLM outperforms
38 general-purpose models such as OpenAI GPT-4o1 and DeepSeek-R1, achieving win rates
39 ranging from 57% to 88%. Furthermore, SeedLLM is integrated with the rice biological
40 knowledge graph (RBKG), which consolidates genome annotations for *Nipponbare* and large-
41 scale synthesis of transcriptomic and proteomic information from over 1,800 studies. This
42 integration enhances the ability of SeedLLM to address complex research questions requiring the
43 fusion of textual and multiomics data. To facilitate global collaboration, we provide free access
44 to SeedLLM and the RBKG via an interactive web portal (<https://seedllm.org.cn/>). SeedLLM
45 represents a transformative tool for rice biology research, facilitating unprecedented discoveries
46 in crop improvement and climate adaptation through its advanced reasoning capabilities and
47 comprehensive data integration.

48

49 KEY WORDS

50 LLM, Knowledge Graph, Multiomics Data Integration, GPT, DeepSeek

51

52 SUMMARY

53 SeedLLM·Rice, a 7-billion-parameter language model trained on 1.4 million rice-related
54 publications, outperforms general-purpose models such as GPT-4o and DeepSeek-R1 in rice-
55 specific tasks. Through integration with a rice biological knowledge graph, it demonstrates

56 superior capacity for multiomics data synthesis, positioning it as a robust tool for AI-enabled
57 crop genomics and systems biology, with continued validation anticipated to broaden its
58 applicability.

Journal Pre-proof

59 INTRODUCTION

60 Paddy rice (*Oryza sativa*) is a crucial staple crop that supports nearly half of the global
61 population (Shi et al., 2023). Research on rice biology is inherently complex, requiring a
62 comprehensive understanding of literature, experimental data, and hypothesis formulation to
63 advance knowledge in the field (Majumdar et al., 2017). This iterative process often challenges
64 or refines established theories, leading to new insights. For example, investigating the function of
65 uncharacterized genes in rice requires a thorough understanding of related biological pathways,
66 as documented in prior studies (Consortium, 2009; Warde-Farley et al., 2010). Similarly, crop
67 breeding efforts rely heavily on genetic and phenotypic data to select optimal parental lines
68 (Huang et al., 2022b). Therefore, a deep understanding of existing knowledge and empirical data
69 is fundamental to advancing rice research.

70 However, the rapid expansion of biological data, coupled with an overwhelming volume
71 of research literature, presents significant challenges (Bornmann and Mutz, 2015). The
72 sequencing of the rice genome in 2002 marked a key turning point, and subsequent
73 advancements in high-throughput technologies have accelerated the generation of vast datasets
74 (Li et al., 2014; Project, 2007). Despite the increasing availability of data, extracting meaningful
75 insights remains a labor-intensive task. Researchers require extensive academic training to
76 navigate this expanding body of knowledge, making the discovery of novel insights increasingly
77 difficult. This highlights the urgent need for advanced tools that are capable of efficiently
78 navigating and extracting relevant information from the growing corpus of rice-related biological
79 data.

80 A promising solution lies in the application of large language models (LLMs) (Naveed et
81 al., 2023). However, several challenges hinder their effective use in rice biology research. First,
82 there is a lack of standardized evaluation frameworks tailored to rice biology, making it difficult
83 to assess the performance of general-purpose LLMs in domain-specific tasks such as
84 comprehension, reasoning, and language generation (Lam et al., 2024). Although benchmarks
85 exist for multilingual tasks or biomedical applications (Li et al., 2020; Singhal et al., 2023), no
86 such tools have been developed for rice biology. Second, the absence of domain-specific training
87 limits the effectiveness of general-purpose models in this domain (Nazi and Peng, 2024), where
88 domain-specific LLMs have been shown outperforming performance over general-purpose
89 LLMs. This gap stems from the limited availability of large-scale, diverse corpora needed to

90 train rice-specific models. Finally, rice biology studies generate complex, multiomics data
91 (Huang et al., 2022a)—such as transcriptomics and genomic sequences—that are difficult to
92 fully represent in textual formats. LLMs, which are predominantly trained using textual data,
93 struggle to synthesize multimodal information, restricting their ability to address complex
94 research questions that require an integrated approach.

95 To address these challenges, we introduce SeedLLM·Rice (hereafter referred to as
96 SeedLLM), a 7-billion-parameter large language model trained using 1.4 million rice-related
97 publications, covering nearly 98.24% of the global literature in the field. To fill this gap,
98 SeedLLM is designed as a domain-specific LLM capable of processing and integrating diverse
99 datasets relevant to rice biology. A novel human evaluation framework is designed to assess the
100 performance of SeedLLM on tasks such as gene function prediction, textual integration of
101 transcriptional and proteomic data, and variety breeding. Initial results from over human
102 evaluations demonstrate that SeedLLM outperforms general-purpose models such as DeepSeek-
103 R1 in rice-specific tasks, with a 88.14% win rate. Furthermore, SeedLLM is integrated with the
104 rice biological knowledge graph (RBKG), a comprehensive resource that includes the latest
105 genome annotations for *Nipponbare* and large-scale synthesis of transcriptional and proteomic
106 information from over 1,800 academic publications. This integration enables SeedLLM to
107 address complex rice biology questions by drawing from both textual data. Our results show that,
108 when augmented with the RBKG, SeedLLM significantly outperforms all general-purpose LLMs
109 in advanced-level rice omics tasks, despite some limitations in reasoning ability. To increase
110 accessibility, we have developed an interactive web portal (<https://seedllm.org.cn/>) that allows
111 researchers worldwide to freely access both SeedLLM and the RBKG, thus accelerating the pace
112 and depth of rice biology research.

113

114 **RESULTS**

115 **Overview of SeedLLM development**

116 To develop a specialized LLM for rice research, we collaborated with experts in rice
117 biology to create RiceCorpus, a rice-specific corpus designed to address the lack of a specialized
118 dataset for rice-domain LLMs (Figure 1A). This rice-specific corpus was developed under the
119 guidance of experts in the field of rice biology research, ensuring comprehensive coverage
120 (Supplemental Figure 1). RiceCorpus integrates both English and Chinese data, reflecting the

121 primary languages used in rice research. The corpus comprises 1.4 million peer-reviewed papers
122 containing keywords such as "rice" and "*Oryza sativa*" published over the past 40 years. This
123 collection represents approximately 98.24% of global rice-related research in these languages
124 during this period, providing a robust foundation for training rice-specific LLMs. Additionally,
125 RiceCorpus includes 1,207 rice-related books. The corpus is exclusively textual and covers a
126 wide range of disciplines in rice research, including molecular biology, plant breeding, and
127 management practices. With a total size of 3,397.49 GB of textual data, RiceCorpus was
128 processed through a multistep quality control pipeline (see the Methods for details), including
129 language detection, content filtering, and deduplication, ensuring that only high-quality data
130 were retained. This meticulous process meets the stringent requirements necessary for training
131 large-scale LLMs.

132

133 **SeedLLM construction and automated evaluation**

134 We selected Qwen2.5-7B, a Transformer-based general-purpose LLM with 7 billion
135 parameters (Qwen et al., 2024), as the base model (Figure 1B). This base model that had been
136 pretrained using a large, multiphase dataset comprising 18 trillion tokens. Previous reports have
137 demonstrated that Qwen excels in language proficiency, comprehension, reasoning, and
138 mathematics. We pretrained the base model using RiceCorpus and GeneralCorpus, the latter
139 being a widely used general-purpose corpus for LLM training (Huang et al., 2024; Penedo et al.,
140 2024). To evaluate its effectiveness, we assessed the model on two rice biology datasets: MCQ-
141 ACC, consisting of 300 single-choice questions, and Gen-QA-ACC, containing 517 short-answer
142 questions (Figure 1C). The pretrained model outperformed the base model that had not been
143 pretrained with RiceCorpus (Figure 1D), demonstrating that pretraining with rice-specific data
144 enabled the LLM to acquire complex rice biology knowledge and domain-specific language
145 patterns. Subsequently, we fine-tuned the pretrained model using RiceQA, a large annotated
146 question–answer dataset in rice biology, along with GeneralQA, a commonly used dataset in
147 general domains (Dong et al., 2024). This process enhanced the model's performance on rice-
148 related tasks where preserving its general-purpose capabilities. The resulting model, SeedLLM
149 (Figure 1B), shares the same architecture as the base model but exhibits improved understanding
150 of rice biology and domain-specific linguistic features due to the pretraining and fine-tuning
151 procedures (see Supplemental Methods for more details).

152 We also conducted an automated evaluation of SeedLLM using the Agri series dataset,
153 which consists of 1,975 question–answer pairs across 10 subdatasets, each with various task
154 types, such as essay-style, summary, language understanding, and multiple-choice questions.
155 SeedLLM outperformed general-purpose LLMs, including Qwen2.5 and Llama3.1 (Grattafiori et
156 al., 2024), across all subdatasets, achieving the highest accuracy, F1, and ROUGE scores in
157 automated evaluations (Figure 1E-G). These results demonstrate the effectiveness of the fine-
158 tuning process in optimizing SeedLLM for rice-specific tasks. Additionally, we assessed the
159 generalizability of SeedLLM by fine-tuning it with GeneralQA, a benchmark dataset for general
160 knowledge widely used to assess LLMs' abilities in general-purpose tasks. Despite being
161 primarily pretrained and fine-tuned for rice-related tasks, SeedLLM achieved accuracy scores
162 comparable to those of general-purpose models in various general knowledge tasks, such as
163 mathematics problem-solving with GSM8K (Figure 1H). This finding suggests that SeedLLM
164 retains the ability to perform general language understanding, reasoning, and mathematics tasks,
165 making it a versatile tool for both domain-specific and general-purpose applications.

166 In summary, we developed SeedLLM, a 7-billion-parameter large language model
167 specifically designed for rice research. By leveraging RiceCorpus, a comprehensive rice-specific
168 corpus, and applying a robust pretraining and posttraining methodology, we created a model
169 capable of outperforming general-purpose LLMs in rice-related tasks, as validated by two rounds
170 of automated evaluations.

171

172 **Human-Centric Evaluation of SeedLLM Performance**

173 LLMs are capable of generating long, coherent, and complex responses. However, they
174 are also prone to factual inaccuracies (Huang et al., 2023), necessitating careful verification by
175 human experts, particularly in specialized fields such as rice biology. To fully evaluate
176 SeedLLM's domain-specific comprehension and knowledge retrieval, we developed a human-
177 centric framework (Figure 2A). This framework enables the assessment of the model's ability to
178 generate complete and accurate answers to real-world research questions.

179 We began by constructing HumanDesignRiceQA, a high-quality, human-designed
180 question-answering benchmark tailored for rice biology. This benchmark enables LLM-
181 generated responses and facilitates subsequent assessment of their accuracy and quality.
182 Developed by rice biology experts, the benchmark comprises 253 questions spanning 6 topics,

183 including rice gene function, multiomics, genome-wide association studies (GWAS), traditional
184 breeding and molecular breeding (Figure 2B). Questions are categorized into three levels on the
185 basis of their complexity: basic, intermediate, and advanced. Basic-level questions can be
186 answered via publicly available information (e.g., research abstracts or online resources) without
187 requiring prior education or experience in rice biology. In contrast, advanced-level questions
188 require individuals to have undergone at least minimal formal training in rice biology and to
189 synthesize information from multiple academic papers, integrating insights from scientific
190 literature and experimental data. The benchmark design reflects the cognitive challenges
191 encountered by individuals with varying levels of expertise when addressing rice biology
192 problems. Additionally, input from a diverse group of experts ensured that the benchmark
193 encompassed a broad spectrum of contemporary rice research topics.

194 Next, we tasked SeedLLM, along with several general-purpose LLMs, with generating
195 responses to the questions in the HumanDesignRiceQA benchmark. To establish a baseline for
196 comparison, undergraduate students specializing in agronomy or crop breeding who had
197 completed relevant coursework also provided responses to the same set of questions. These
198 undergraduate students' answers served as a representative baseline ability of human
199 performance, allowing for a direct comparison between SeedLLM and typical human
200 understanding of rice biology.

201 We assembled a panel of human evaluators to assess whether responses demonstrated
202 correct or incorrect rice-specific reading comprehension and knowledge retrieval. Over 326
203 individuals with academic backgrounds in agronomy, including 83 experts in rice biology and
204 variety development, participated in the evaluation (Figure 2C). Evaluations were conducted in a
205 blinded manner, ensuring that the evaluators were unaware of which responses were generated
206 by SeedLLM. Responses were rated on a scale from 0 to 100 or ranked from best to worst, using
207 answer keys or expert experience as reference.

208 SeedLLM received higher human evaluation scores than all other tested LLMs in the
209 HumanDesignRiceQA benchmark across both rounds (Figure 2D, Supplemental Figure 2).
210 Superior performance of SeedLLM was observed across all question difficulties, as it
211 outperformed all other tested LLMs in basic, intermediate, and advanced-level questions
212 (Supplemental Figure 3). Evaluators also assigned that the highest number of the best responses
213 came to SeedLLM (Figure 2E). These results demonstrate that SeedLLM outperforms general-

214 purpose LLMs in rice-specific question-answering tasks, as validated by human evaluators over
215 time. This suggests that SeedLLM maintains its leading performance, even as general LLMs
216 evolve.

217 In conclusion, the human-centric evaluation confirmed that SeedLLM outperforms both
218 general-purpose LLMs and the human baseline across a broad range of rice biology questions.
219 The model demonstrated state-of-the-art performance in tasks requiring a deep understanding of
220 rice biology. However, SeedLLM achieved an average score of 69.98 in answering advanced-
221 level questions, highlighting areas for improvement, particularly in multistep reasoning and
222 integrating complex biological data into textual information.

223

224 **Construction of the rice biological knowledge graph**

225 Recent studies have shown that incorporating external knowledge graphs into LLMs
226 enhances their reasoning and data fusion capabilities (Pan et al., 2023; Peng et al., 2024).
227 Motivated by these findings, we developed the rice biological knowledge graph (RBKG), a
228 multimodal graph that integrates transcriptional and proteomic data from over 1,879 papers and
229 gene annotation information (Figure 3A and 3B). The construction of the RBKG occurred in
230 three phases: textual integration of transcriptional and proteomic data, integration of rice genome
231 annotation, and comprehensive data fusion.

232 We first identified scientific papers reporting rice transcriptome and proteome data.
233 Through a comprehensive literature search, we identified 1,879 papers provided raw or
234 preanalyzed transcriptomic and proteomic data. However, inconsistencies in experimental
235 protocols and analysis methods present challenges to data standardization. To address this
236 problem, we structured the transcriptomic and proteomic data by representing each gene's
237 transcriptional event as a sentence, using an approach similar to that of CellAnnoation (Fang et
238 al., 2024). This structure included gene expression levels, protein abundance, experimental
239 attributes (e.g., rice variety, genetic background, tissue type, developmental stage), and
240 conditions specified in the respective studies. These data were then modeled within a knowledge
241 graph framework, with nodes representing gene IDs, transcriptional events, and experimental
242 attributes and edges denoting relationships among them (Figure 3A). This approach facilitated
243 the standardization and harmonization of diverse rice omics data, resulting in a cohesive rice
244 omics knowledge graph. Next, we integrated gene annotation data, including Gene Ontology

245 (GO) terms and Kyoto Encyclopedia of Genes and Genomes (KEGG) pathway annotations, for
246 all genes identified in the transcriptomic and proteomic datasets. Annotation nodes were created
247 and linked to the corresponding gene IDs, ensuring that functional annotations were directly
248 connected to the underlying transcriptomic and proteomic data.

249 Finally, we fused the individual knowledge graphs—transcriptomic, proteomic, and gene
250 annotation—into a unified RBKG using graph matching algorithms to reconcile discrepancies
251 and resolve conflicts from overlapping data sources. The resulting RBKG comprises 401,094
252 nodes and 1,573,258 edges, representing transcriptional, translational, and genomic data for
253 33,599 rice genes (Figure 3B). To our knowledge, the RBKG is the first knowledge graph that
254 integrates both transcriptional and proteomic data for rice, providing a comprehensive and
255 structured framework for rice multiomics retrieval and reasoning.

256

257 **SeedLLM integrates knowledge graphs for complex biological queries**

258 To investigate whether the RBKG could enhance the ability of SeedLLM to answer
259 complex biological questions related to rice, we developed a framework that integrates the
260 RBKG into the response generation process of SeedLLM (Figure 3C), hereafter referred to as
261 SeedLLM-KG. This framework comprises three key steps: query decomposition, entity
262 grouping, and knowledge augmentation. In the first step, SeedLLM-KG decomposes user queries
263 into essential entities and attributes, identifying critical components for subsequent processing.
264 The second step involves constructing entity groups by linking the queried entity with relevant
265 concepts from the knowledge graph. Finally, in the knowledge augmentation phase, SeedLLM-
266 KG uses its parametric knowledge base to establish connections between the queried entity and
267 pertinent concepts, thereby enriching the context of the query.

268 For example, given the following query: “Does the rice gene AGIS_Os06g035130
269 respond to various environmental conditions?” SeedLLM-KG initially fails to locate
270 AGIS_Os06g035130 in its textual knowledge base (Figure 4A and 4B). It then queries the
271 RBKG for the node representing AGIS_Os06g035130, expanding the search to identify
272 connected nodes that include gene annotations, functional descriptions, and multiomics data. The
273 system retrieves all relevant nodes and edges corresponding to transcriptomic and translational
274 data, synthesizing this structured knowledge into a coherent, human-readable response. Notably,
275 SeedLLM-KG consistently provided correct responses when the AGIS ID of the rice gene was

276 converted to the corresponding MSU ID and RAP-DB ID in the queries (Supplemental Figures 4
277 and 5).

278

279 **SeedLLM-KG integrates multiomic and literature data for complex rice biology questions**

280 To evaluate the performance of SeedLLM-KG in addressing complex rice biology
281 questions, we conducted human-centric assessments via advanced-level questions from the
282 HumanDesignRiceQA dataset. Human experts evaluated the correctness of responses from
283 SeedLLM-KG, awarding an median score of 85, which was significantly higher than the score of
284 67 awarded to SeedLLM alone (Supplemental Figure 6). This improvement stems from
285 SeedLLM-KG's ability to retrieve and synthesize transcriptional, proteomic and specialized rice
286 research data, which are absent from general-purpose text corpora used to train conventional
287 LLMs (Figure 4A; Supplemental Figures 4 and 5). We also included state-of-the-art general-
288 purpose LLMs, such as DeepSeek-R1 (DeepSeek-AI et al., 2025), DeepSeek-V3 (DeepSeek-AI
289 et al., 2024), GPT-4o1 (OpenAI et al., 2023), and GPT-4o3-mini in our evaluation. SeedLLM-
290 KG demonstrated a unique ability to integrate heterogeneous biological data sources,
291 outperforming all tested general-purpose LLMs (Figure 3D), with win rates ranging from
292 57.63% to 88.14% (Figure 3E). Notably, it surpassed DeepSeek-R1 in rice-specific task, a
293 leading LLM at the time of writing.

294 Interestingly, human evaluators assigned similar scores in assessing SeedLLM's
295 reasoning ability, which were not statistically different from other LLMs (Figure 3F). This result
296 was unexpected, as other LLMs, such as DeepSeek-V3 with its 671 billion parameters, were
297 expected to exhibit superior reasoning capabilities, especially considering that SeedLLM-KG is
298 based on SeedLLM, which has only a 7-billion-parameter architecture and lacks reinforcement
299 learning in the post-training stage. The results suggest that human evaluators considered
300 SeedLLM-KG to exhibit similar reasoning abilities compared to LLMs with larger architectures
301 in rice-specific advanced-level tasks.

302

303 **DISCUSSION**

304 LLMs have revolutionized the field of AI, particularly in general-purpose content
305 generation. However, their application in specialized fields such as rice biology remains limited
306 because of a lack of domain-specific training data. Here, we presented SeedLLM, a 7-billion-

307 parameter model developed from scratch using 1.4 million rice-related research publications,
308 representing nearly 98.24% of the global literature on rice biology. To our knowledge, SeedLLM
309 is the first LLM specifically designed for this domain. Its potential to advance both rice biology
310 and broader plant biology research is substantial, addressing critical gaps in the current research
311 landscape.

312 While general-purpose LLMs have garnered considerable interest, their application in
313 specialized fields such as rice biology raises an important question: can these models effectively
314 contribute to scientific discovery in this domain? An ideal model for rice biology must be
315 capable of mastering a vast body of knowledge, enabling efficient information retrieval and
316 fostering scientific breakthroughs. However, a comprehensive framework for evaluating LLMs
317 in rice biology is lacking. To address this gap, we propose a novel evaluation framework
318 featuring a robust question-answer dataset of rice biology, comprising 253 human-designed
319 questions and 1,975 additional automatically generated questions. Our evaluations—both
320 subjective human-led assessments and objective automated evaluations—demonstrate that
321 general-purpose LLMs underperform relative to SeedLLM across multiple rice biology tasks.
322 These findings underscore the necessity for domain-specific models to increase research
323 efficiency and accuracy. The performance of SeedLLM is significantly enhanced by the
324 integration of a knowledge graph, which consolidates transcriptional, proteomic data from over
325 1,800 research papers—data typically inaccessible to general-purpose models. This knowledge
326 graph is a key factor driving the superior performance of the SeedLLM in rice biology tasks,
327 highlighting the importance of specialized knowledge sources in optimizing model output.

328 Despite these advances, several limitations must be addressed. First, although SeedLLM
329 is specifically designed to improve rice biology knowledge, it still suffers from hallucinations—
330 incorrect or fabricated responses—particularly in specialized tasks. During our evaluation,
331 distinguishing between hallucinations and information retrieval failures proved challenging, as
332 both can lead to inaccurate outputs. This issue is exacerbated by the absence of a more
333 comprehensive knowledge graph, which could serve as a structured repository of accurate data to
334 guide model outputs. Expanding and refining this domain-specific knowledge graph in future
335 iterations of SeedLLM could reduce hallucinations and enhance model reliability.

336 Second, while scaling laws suggest that larger models exhibit improved performance with
337 increased training data and model complexity (Kaplan et al., 2020), the current version of

338 SeedLLM—comprising 7 billion parameters—faces limitations owing to available computational
339 resources and human expertise. Increasing the model size—such as for 14B, 32B, or even 72B
340 parameters—may lead to performance gains. However, it remains uncertain whether such
341 increases significantly impact the discovery of novel phenomena in rice biology that the current
342 configuration of SeedLLM has not yet captured. As such, we propose that future studies should
343 prioritize improving data diversity and integrating cutting-edge rice biology datasets rather than
344 solely focusing on scaling the model.

345 Third, the RiceQA dataset used for supervised fine-tuning, while instrumental to
346 SeedLLM development, has limitations in terms of data quality and diversity. Although RiceQA
347 is comprehensive, it could better represent underexplored areas in rice biology and more diverse
348 experimental conditions. Expanding the dataset to cover a broader range of topics (Han et al.,
349 2023) and ensuring higher accuracy in question design would likely improve the performance of
350 SeedLLM, especially in addressing more complex or nuanced queries. Additionally, a higher-
351 quality dataset could help mitigate issues such as hallucinations, making model outputs more
352 reliable.

353 Rice biology research has traditionally been labor intensive, requiring scientists to
354 manually process vast amounts of textual and biological data. SeedLLM represents a substantial
355 advancement in research efficiency, enabling scientists to interact with the model via natural
356 language queries and obtain critical information much faster than traditional search engines and
357 databases. This is exemplified by our recent integration of the high-quality RiceData database
358 (www.ricedata.cn) into SeedLLM's response pipeline (Supplemental Figure 7). This capability is
359 expected to significantly accelerate the pace of discovery and analysis in rice biology.

360 We look forward to further enhancing the performance of SeedLLM. Insights gained
361 from this study—particularly regarding the importance of knowledge graphs—will guide the
362 development of an expanded, more sophisticated rice knowledge graph. This graph integrates
363 recent advances in rice biology, including single-cell sequencing and spatial transcriptomics, and
364 transforms these complex, high-dimensional datasets into a format that SeedLLM can process
365 more effectively. Traditional academic literature often fails to capture the intricacies of such
366 cutting-edge datasets, creating a challenge for researchers. By enabling SeedLLM to incorporate
367 and analyze these advanced datasets, we aim to create a more comprehensive and accurate model
368 of rice biology. Additionally, recent studies indicate that LLM agents, which utilize domain-

369 specific external tools to autonomously execute tasks, provide significant advantages over
370 traditional chatbots (Kapoor et al., 2024). We propose integrating tools like RiceNavi (Wei et al.,
371 2021) with SeedLLM to streamline breeding tasks such as target and parental line selection,
372 thereby significantly enhancing the practical utility of SeedLLM. Although SeedLLM currently
373 answers queries by retrieving rice epistatic QTL pairs (Wei et al., 2024) that serve as input data
374 for RiceNavi pipeline, this integration would further elevate its capabilities (Supplemental Figure
375 8).

376 Our long-term vision for SeedLLM is to evolve into a global, comprehensive knowledge
377 atlas that will provide researchers with unprecedented access to insights previously hidden owing
378 to reliance on isolated data points. This vision aligns with the broader trend of AI-driven
379 knowledge synthesis, where models such as SeedLLM will enable new discoveries by
380 integrating diverse and complex data sources. Additionally, as we develop specialized LLMs for
381 other crops, we anticipate that cross-species knowledge reasoning will become increasingly
382 feasible, empowering researchers with more powerful tools for AI-assisted seed design in the
383 future.

384

385 **METHODS**

386 **Comprehensive search and retrieval of rice-related publications.** The dataset used in this
387 study comprises scientific publications in both English and Chinese. English-language
388 publications were retrieved from the Web of Science using a search query that included various
389 rice-related terms, resulting in 1,148,299 publications. Chinese-language publications were
390 sourced from the China National Knowledge Infrastructure using a corresponding set of rice-
391 related terms, yielding 232,445 publications. All publications used for the construction of
392 RiceCorpus are up to December 31, 2024.

393

394 **Construction of RiceCorpus.** To generate a high-quality, rice-specific corpus, we developed a
395 reusable, high-granularity data cleaning pipeline consisting of four primary stages. First, raw
396 PDF documents of rice publication were converted to text using the MinerU tool (Wang et al.,
397 2024a), with the Layout model accurately recognizing document sections (e.g., titles, abstracts)
398 to ensure semantic coherence. Text extraction, formula recognition, and table conversion were
399 performed, while PaddleOCR assisted in optical character recognition (Du et al., 2020). Post-

400 processing with regular expressions optimized the identification of rice-related terms, gene
401 names, numbers, and punctuation. Second, heuristic cleaning was applied to address redundancy
402 and irrelevant content in the resulting TXT files. Statistical analysis of rice literature informed
403 the development of regular expression-based rules to filter low-quality text, retaining 72.82%
404 high-quality documents, thereby enhancing the corpus' rice knowledge density. To further
405 improve model training efficiency and reduce overfitting, we performed sentence-level
406 deduplication using MinHash (Broder, 1997), which calculates n-gram similarity between text
407 pairs. After replacing MinHash's tokenizer with SeedLLM's, deduplication was conducted at the
408 sentence level, removing 1,834,317 sentences (25% of the data). Lastly, model-based filtering
409 was used to eliminate non-rice content by applying the IndustryCorpus2_Classifier (Wang et al.,
410 2024b), which classified text into 31 domain categories, retaining only agricultural, biological,
411 and chemical content. The CCI3-HQ-Classifer provided quality scores (Wang et al., 2024b),
412 filtering out segments with scores below 2. These combined methods resulted in a defined corpus
413 termed as RiceCorpus consisting of 1.1 billion tokens, ensuring high relevance and quality for
414 subsequent model training.

415

416 **Model Pretraining.** To enhance rice domain capabilities without compromising generalization,
417 we pretrained base model Qwen-2.5-7B with RiceCorpus along with GeneralCorpus. The
418 GeneralCorpus included downsampled Fineweb-Edu (English) (Penedo *et al.*, 2024), Fineweb-
419 Edu-Chinese-V2.1 (Chinese) (Huang *et al.*, 2024), as well as code-mathematical corpus such as
420 Opc-Fineweb-Math-Corpus and Opc-Fineweb-Code-Corpus. All model-related hyperparameters
421 of pre-training matching those of Qwen-2.5-7B (Qwen *et al.*, 2024). The AdamW optimizer was
422 used with $\beta_1=0.9$, $\beta_2=0.999$, weight decay=0.0, and a maximum context length of 4K tokens. The
423 learning rate followed a linear increase from $8e-10$ to $6e-07$ for the first 10% of training steps,
424 then decayed to $2e-14$ following a cosine curve. Pretraining was conducted on 16 NVIDIA H100
425 GPUs with a batch size of 128K tokens. An ablation experiment using only the RiceQA
426 demonstrated the model's ability to avoid catastrophic forgetting, enhancing both cross-domain
427 generalization and rice-domain capabilities, as evidenced by a substantial improvement in the
428 BBH metric compared to the baseline.

429

430 **Model Post-training.** GraphGen generates high-quality synthetic data for LLM fine-tuning by
431 leveraging knowledge graphs. The process begins with knowledge construction, where text is
432 segmented into semantically coherent chunks, and a synthesizer model M_{synth} extracts entities and
433 relationships, which are merged to form a structured knowledge graph (Ibrahim et al., 2024).
434 This enables effective long-text processing and reduces content hallucination. In the
435 comprehension evaluation phase, M_{synth} generates paraphrased statements and negations to assess
436 the model’s understanding, with a comprehension loss computed based on the model’s
437 confidence scores. Graph organization follows, where subgraphs are extracted using methods
438 like k-hop and selection strategies to balance complexity and relevance. Finally, question-answer
439 pairs generation is performed for various scenarios, including atomic, aggregated, and multi-hop
440 question-answer pairs, based on the subgraphs. For the SFT phase, RiceQA dataset is created by
441 categorizing the Infinity Instruct dataset into six categories, selecting the top 300k question-
442 answer pairs based on vector similarity, and augmenting the data with the synthetic pairs via
443 AutoIF (Dong *et al.*, 2024). The final dataset of question-answer pairs integrates domain-specific
444 knowledge with general capabilities while maintaining data quality through expert curation and
445 overlap removal.

446
447 **Training Configuration.** During the training phase, we utilized XTuner as the training
448 framework, based on the Transformer architecture and optimized using the AdamW optimizer.
449 The learning rate followed a linear schedule with a warm-up phase, and gradient clipping was
450 applied to stabilize training. The model training employed several key parameters to optimize
451 performance. The maximum sequence length was set to 2048, defining the maximum input
452 sequence size. A learning rate of $2e-5$ was used to control the step size for weight updates, while
453 a weight decay of 0.1 helped mitigate overfitting by penalizing large weights. Gradient clipping
454 was applied with a threshold of 1 to stabilize training and prevent gradient explosion. A batch
455 size of 64 (16×4) was selected, determining the number of samples processed per optimization
456 step. The AdamW optimizer, a variant of Adam with decoupled weight decay, was used for
457 improved generalization, with β_1 and β_2 values set at 0.9 and 0.999, respectively, for exponential
458 decay of moment estimates. To gradually increase the learning rate, a warm-up ratio of 0.03 was
459 applied for the initial fraction of the total training steps, and the model was trained for 2 epochs,
460 completing two full passes over the dataset.

461

462 **Automated evaluations of LLM performance.** To assess model performance, we use accuracy
463 for classification tasks and perplexity (PPL) for language modeling (Hu et al., 2024). The pre-
464 trained model was evaluated on PPL-MCQ-ACC and Gen-QA-ACC datasets. The supervised
465 fine-tuned model was evaluated on several general-purpose benchmarks CMMLU (Li et al.,
466 2023), GSM8K (Cobbe et al., 2021), BBH (Srivastava et al., 2022), MMLU (Hendrycks et al.,
467 2020) and Agri series dataset, which a domain-specific rice dataset from SeedBench (Ying et al.,
468 2025), focusing on accuracy for multiple-choice tasks and PPL for fill-in-the-blank tasks.

469

470 **Human-mediated evaluation of LLM performance.** To evaluate LLM within the domain of
471 rice biology, we adapted a structured human evaluation framework inspired by the methodology
472 proposed by Petrov et al. (2025). With domain-specific modifications, we developed
473 HumanDesignRiceQA, a curated benchmark comprising 253 expert-authored questions spanning
474 six major topics: gene function, multi-omics, genome-wide association studies (GWAS),
475 traditional breeding, molecular breeding, and gene editing. Each question was classified into one
476 of three complexity tiers—basic, intermediate, and advanced—based on the depth of biological
477 knowledge and reasoning required. Reference answers were derived from peer-reviewed
478 literature, and evaluation rubrics were constructed by biological science experts with Master's-
479 level training. Responses generated by LLMs, including SeedLLM, were assessed alongside
480 those written by students through a blinded review process conducted by 326 human evaluators,
481 of whom 83 were domain experts in rice biology. Each response was independently scored by
482 three evaluators using a 0–100 scale and ranked relative to other answers based on predefined
483 criteria encompassing factual accuracy, logical structure, and clarity. This evaluation framework
484 constitutes the first domain-adapted, human-mediated assessment pipeline in plant science,
485 establishing a rigorous benchmark for comparing LLM outputs in rice biology. The results
486 demonstrate that SeedLLM consistently outperforms peer models across all levels of question
487 complexity.

488

489 **Construction of Rice Biological Knowledge Graph.** We developed a Python pipeline to
490 identify publications on rice transcriptomics and proteomics by searching for relevant keywords
491 in titles and abstracts. After filtering potential papers, each was manually reviewed to confirm

492 focus on rice transcriptomics or proteomics, and whether raw or processed data were available.
493 From these studies, we curated lists of the most upregulated and downregulated genes and
494 proteins, standardizing them to rice AGIS IDs (Shang et al., 2023). Experimental metadata,
495 including genotype, tissue or organ, growth stage, and treatments, were extracted and used to
496 generate structured annotations (e.g., “Genotype_X under treatment_Y shows differential
497 expression of gene AGIS_ID_1 in organ_Z at growth stage_W”). These annotations were
498 converted into triples (subject, relation, object) to represent transcriptional and translational
499 events, as well as experimental conditions. Each AGIS_ID was cross-referenced with databases
500 like RAP-DB (Sakai et al., 2013) and Gramene (Jaiswal, 2011) to obtain functional annotations
501 and subcellular localization information. These data were incorporated into the knowledge graph,
502 linking them to the corresponding gene or protein nodes.

503

504 **Visualization of Rice Biological Knowledge Graph.** For basic visualization, we used the
505 networkx library, assigning node and edge styles based on entity types (e.g., proteins in blue,
506 growth stages in green). The spring layout algorithm optimized node positioning for clarity. For
507 advanced visualizations, the graph was exported in GraphML format and imported into Gephi
508 and Cytoscape. These tools enabled customization, such as adjusting node size by protein
509 expression magnitude and edge thickness by interaction strength, allowing for a more detailed
510 exploration of the data.

511

512 **Construction of HumanDesignRiceQA.** We developed a benchmark for rice-specific
513 knowledge based on academic papers by creating three question levels—basic, intermediate, and
514 advanced—covering five major research areas: gene function, transcriptomics, proteomics,
515 traditional breeding, and molecular breeding. The levels are distinguished by the complexity of
516 reasoning and knowledge integration, not by specific topics. Basic questions rely on readily
517 accessible information, such as abstracts or general knowledge from search engines. Intermediate
518 questions require a deeper understanding of rice biology, focusing on the paper’s results with
519 experimental details. Advanced questions demand specialized expertise in rice research, with the
520 ability to analyze biological data within a broader biological context. The levels differ in analysis
521 depth, knowledge integration, and reasoning complexity. Academic papers from each research

522 area were randomly selected, reviewed, and used to design the questions and corresponding
523 answer sheets.

524
525 **Graph-based Retrieval-Augmented Generation (GraphRAG).** To enhance the performance
526 of SeedLLM, we developed a framework that integrates external knowledge through a structured
527 retrieval process to guide SeedLLM’s response generation. This approach combines graph
528 structures with dense indexing methods to represent relationships between knowledge fragments,
529 facilitating the retrieval of relevant information for LLM-generated responses (Peng et al., 2024).
530 The GraphRAG framework consists of three main components graph-guided indexing, retrieval,
531 and text generation. In the graph-guided indexing phase, data is preprocessed into manageable
532 chunks, followed by entity and relation extraction to form “entity, relation, description” tuples,
533 which are organized into a knowledge graph. Persistence ensures that these relationships are
534 embedded and stored in a database for efficient access. In the graph-guided retrieval phase, logic
535 form method and dual-level method is conducted for rice-related queries and for general queries,
536 respectively. The logic form method decomposes the query into operators and parameters,
537 generating sub-queries whose results are merged to form a retrieval context. Pre- and post-checks
538 verify whether the context sufficiently supports the LLM’s response. If the logic form method is
539 insufficient, the dual-level method decomposes the query into high-level semantic
540 representations and low-level entity-based components, using fuzzy matching to identify relevant
541 nodes and relationships. The results are merged to create a comprehensive retrieval context.
542 Finally, in graph-guided text generation, the retrieved context and the original query are input
543 into the LLM, which generates the final output based on the enriched context. This integrated
544 approach enables LLMs to leverage structured external knowledge, leading to more accurate and
545 contextually relevant responses.

546
547 **Code Availability.** SeedLLM will continue to grow and improve through version control.
548 Currently, SeedLLM-Rice, version 0.6a is available via an interactive web portal
549 <https://seedllm.org.cn/>.

550
551 **SUPPLEMENTAL INFORMATION.** Supplemental Information is available at
552 Molecular Plant Online.

553

554 **FUNDING.** SeedLLM research and development was supported by Yazhouwan National
555 Laboratory Project (grant no. 2310CF01) and Shanghai Artificial Intelligence Laboratory.
556 Corpus preparation was supported by the Hainan Yazhou Bay Seed Laboratory Project (grant no.
557 B21HJ0001). Human evaluation of LLM performance was supported by Biological Breeding-
558 National Science and Technology Major Project (2023ZD04076).

559

560 **AUTHOR CONTRIBUTIONS.** F.Y. and N.Q.D. conceptualized the idea, designed the
561 experiments, wrote the manuscript, supervised the project, and acquired fundings. F.Y., N.Q.D.,
562 and H.J.K designed the methodology. H.J.K, J.Y, Z.H.C,W.L.J, Z.H.Y, and Z.F.W trained the
563 SeedLLM·Rice and SeedLLM·Rice-KG models, conducted automated model evaluations, and
564 developed the web portal for user access to SeedLLM·Rice-KG. T.L., Z.N.M., S.K.W., X.Y.W.,
565 W.F.M., and X.Y.L. prepared the RiceCorpus and HumanDesignRiceQA datasets and performed
566 the data analysis in human evaluations. X.D.M., M.G.L., and X.Q.W. organized human-centric
567 model evaluations. Z.Y.H. conducted comprehensive searching for rice-related publications.

568

569 **ACKNOWLEDGEMENTS.** We thank members of the Innovation Platform for Seed Design
570 Team and Hongqing Ling in YNL fruitful discussions. We thank Xinquan Yang for inviting
571 human evaluators. We thank YZBSTCACC for support with GPU resources and Pengcheng
572 Cloud Brain II for data storage support. We thank SiliconFlow for support with GPU resources
573 for SeedLLM deployment on the cloud. We thank Gaojun Fan in Pengcheng Laboratory for
574 OpenI Community support. No conflict of interest is declared.

575

576 **REFERENCES**

- 577 Bornmann, L., and Mutz, R. (2015). Growth rates of modern science: A bibliometric analysis
578 based on the number of publications and cited references. *Journal of the Association for*
579 *Information Science and Technology* 66:2215-2222. <https://doi.org/10.1002/asi.23329>.
- 580 Broder, A.Z. (1997). On the resemblance and containment of documents. *Proceedings.*
581 *Compression and Complexity of SEQUENCES 1997 (Cat. No.97TB100171)*.
- 582 Cobbe, K., Kosaraju, V., Bavarian, M., Chen, M., Jun, H., Kaiser, L., Plappert, M., Tworek, J.,
583 Hilton, J., Nakano, R., et al. (2021). Training Verifiers to Solve Math Word Problems.
584 Consortium, T.U. (2009). The Universal Protein Resource (UniProt) in 2010. *Nucleic Acids*
585 *Research* 38:D142-D148. 10.1093/nar/gkp846.
- 586 DeepSeek-AI, Liu, A., Feng, B., Xue, B., Wang, B., Wu, B., Lu, C., Zhao, C., Deng, C., Zhang,
587 C., et al. (2024). DeepSeek-V3 Technical Report.
- 588 DeepSeek-AI, Guo, D., Yang, D., Zhang, H., Song, J., Zhang, R., Xu, R., Zhu, Q., Ma, S., Wang,
589 P., et al. (2025). DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement
590 Learning.
- 591 Dong, G., Lu, K., Li, C., Xia, T., Yu, B., Zhou, C., and Zhou, J. (2024). Self-play with Execution
592 Feedback: Improving Instruction-following Capabilities of Large Language Models.
- 593 Du, Y., Li, C., Guo, R., Yin, X., Liu, W., Zhou, J., Bai, Y., Yu, Z., Yang, Y., Dang, Q., et al.
594 (2020). PP-OCR: A Practical Ultra Lightweight OCR System.
- 595 Fang, C., Wang, Y., Song, Y., Long, Q., Lu, W., Chen, L., Wang, P., Feng, G., Zhou, Y., and Li,
596 X. (2024). How do Large Language Models understand Genes and Cells.
597 *bioRxiv:2024.2003.2023.586383*. 10.1101/2024.03.23.586383.
- 598 Grattafiori, A., Dubey, A., Jauhri, A., Pandey, A., Kadian, A., Al-Dahle, A., Letman, A., Mathur,
599 A., Schelten, A., Vaughan, A., et al. (2024). The Llama 3 Herd of Models.
- 600 Han, L., Zhong, W., Qian, J., Jin, M., Tian, P., Zhu, W., Zhang, H., Sun, Y., Feng, J.-W., Liu, X.,
601 et al. (2023). A multi-omics integrative network map of maize. *Nature Genetics* 55:144-153.
602 10.1038/s41588-022-01262-1.
- 603 Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. (2020).
604 *Measuring Massive Multitask Language Understanding*.
- 605 Hu, Y., Huang, Q., Tao, M., Zhang, C., and Feng, Y. (2024). Can Perplexity Reflect Large
606 Language Model's Ability in Long Text Understanding?

607 Huang, F., Jiang, Y., Chen, T., Li, H., Fu, M., Wang, Y., Xu, Y., Li, Y., Zhou, Z., Jia, L., et al.
608 (2022a). New Data and New Features of the FunRiceGenes (Functionally Characterized Rice
609 Genes) Database: 2021 Update. *Rice* 15:23. 10.1186/s12284-022-00569-1.

610 Huang, L., Yu, W., Ma, W., Zhong, W., Feng, Z., Wang, H., Chen, Q., Peng, W., Feng, X., Qin,
611 B., et al. (2023). A Survey on Hallucination in Large Language Models: Principles, Taxonomy,
612 Challenges, and Open Questions.

613 Huang, S., Cheng, T., Liu, J.K., Hao, J., Song, L., Xu, Y., Yang, J., Liu, J.H., Zhang, C., Chai, L.,
614 et al. (2024). OpenCoder: The Open Cookbook for Top-Tier Code Large Language Models.

615 Huang, X., Huang, S., Han, B., and Li, J. (2022b). The integrated genomics of crop
616 domestication and breeding. *Cell* 185:2828-2839. 10.1016/j.cell.2022.04.036.

617 Ibrahim, N., Aboulela, S., Ibrahim, A., and Kashef, R. (2024). A survey on augmenting
618 knowledge graphs (KGs) with large language models (LLMs): models, evaluation metrics,
619 benchmarks, and challenges. *Discover Artificial Intelligence* 4:76. 10.1007/s44163-024-00175-8.

620 Jaiswal, P. (2011). Gramene database: a hub for comparative plant genomics. *Methods Mol Biol*
621 678:247-275. 10.1007/978-1-60761-682-5_18.

622 Kaplan, J., McCandlish, S., Henighan, T., Brown, T.B., Chess, B., Child, R., Gray, S., Radford,
623 A., Wu, J., and Amodei, D. (2020). Scaling Laws for Neural Language Models.

624 Kapoor, S., Stroebel, B., Siegel, Z.S., Nadgir, N., and Narayanan, A. (2024). AI Agents That
625 Matter.

626 Lam, H.Y.I., Ong, X.E., and Mutwil, M. (2024). Large language models in plant biology. *Trends*
627 *in Plant Science* 29:1145-1155. <https://doi.org/10.1016/j.tplants.2024.04.013>.

628 Li, H., Arora, A., Chen, S., Gupta, A., Gupta, S., and Mehdad, Y. (2020). MTOP: A
629 Comprehensive Multilingual Task-Oriented Semantic Parsing Benchmark.

630 Li, H., Zhang, Y., Koto, F., Yang, Y., Zhao, H., Gong, Y., Duan, N., and Baldwin, T. (2023).
631 CMMLU: Measuring massive multitask language understanding in Chinese.

632 Li, J.-Y., Wang, J., and Zeigler, R.S. (2014). The 3,000 rice genomes project: new opportunities
633 and challenges for future rice research. *GigaScience* 310.1186/2047-217x-3-8.

634 Majumdar, J., Naraseeyappa, S., and Ankalaki, S. (2017). Analysis of agriculture data using data
635 mining techniques: application of big data. *Journal of Big Data* 4:20. 10.1186/s40537-017-0077-
636 4.

637 Naveed, H., Ullah Khan, A., Qiu, S., Saqib, M., Anwar, S., Usman, M., Akhtar, N., Barnes, N.,
638 and Mian, A. (2023). A Comprehensive Overview of Large Language Models.

639 Nazi, Z.A., and Peng, W. (2024). Large Language Models in Healthcare and Medical Domain: A
640 Review. *Informatics* 11:57.

641 OpenAI, Achiam, J., Adler, S., Agarwal, S., Ahmad, L., Akkaya, I., Leoni Aleman, F., Almeida,
642 D., Altenschmidt, J., Altman, S., et al. (2023). GPT-4 Technical Report.

643 Pan, S., Luo, L., Wang, Y., Chen, C., Wang, J., and Wu, X. (2023). Unifying Large Language
644 Models and Knowledge Graphs: A Roadmap.

645 Penedo, G., Kydlíček, H., Ben allal, L., Lozhkov, A., Mitchell, M., Raffel, C., Von Werra, L., and
646 Wolf, T. (2024). The FineWeb Datasets: Decanting the Web for the Finest Text Data at Scale.

647 Peng, B., Zhu, Y., Liu, Y., Bo, X., Shi, H., Hong, C., Zhang, Y., and Tang, S. (2024). Graph
648 Retrieval-Augmented Generation: A Survey.

649 Project, R.A. (2007). The Rice Annotation Project Database (RAP-DB): 2008 update *. *Nucleic
650 Acids Research* 36:D1028-D1033. 10.1093/nar/gkm978.

651 Qwen, Yang, A., Yang, B., Zhang, B., Hui, B., Zheng, B., Yu, B., Li, C., Liu, D., Huang, F., et al.
652 (2024). Qwen2.5 Technical Report.

653 Sakai, H., Lee, S.S., Tanaka, T., Numa, H., Kim, J., Kawahara, Y., Wakimoto, H., Yang, C.C.,
654 Iwamoto, M., Abe, T., et al. (2013). Rice Annotation Project Database (RAP-DB): an integrative
655 and interactive database for rice genomics. *Plant Cell Physiol* 54:e6. 10.1093/pcp/pcs183.

656 Shang, L., He, W., Wang, T., Yang, Y., Xu, Q., Zhao, X., Yang, L., Zhang, H., Li, X., Lv, Y., et al.
657 (2023). A complete assembly of the rice Nipponbare reference genome. *Molecular Plant*
658 16:1232-1236. 10.1016/j.molp.2023.08.003.

659 Shi, J., An, G., Weber, A.P.M., and Zhang, D. (2023). Prospects for rice in 2050. *Plant, Cell &
660 Environment* 46:1037-1045. <https://doi.org/10.1111/pce.14565>.

661 Singhal, K., Azizi, S., Tu, T., Mahdavi, S.S., Wei, J., Chung, H.W., Scales, N., Tanwani, A., Cole-
662 Lewis, H., Pfohl, S., et al. (2023). Large language models encode clinical knowledge. *Nature*
663 620:172-180. 10.1038/s41586-023-06291-2.

664 Srivastava, A., Rastogi, A., Rao, A., Shoeb, A.A.M., Abid, A., Fisch, A., Brown, A.R., Santoro,
665 A., Gupta, A., Garriga-Alonso, A., et al. (2022). Beyond the Imitation Game: Quantifying and
666 extrapolating the capabilities of language models.

- 667 Wang, B., Xu, C., Zhao, X., Ouyang, L., Wu, F., Zhao, Z., Xu, R., Liu, K., Qu, Y., Shang, F., et
668 al. (2024a). MinerU: An Open-Source Solution for Precise Document Content Extraction.
- 669 Wang, L., Zhang, B.-W., Wu, C., Zhao, H., Shi, X., Gu, S., Li, J., Ma, Q., Pan, T., and Liu, G.
670 (2024b). CCI3.0-HQ: a large-scale Chinese dataset of high quality designed for pre-training large
671 language models.
- 672 Warde-Farley, D., Donaldson, S.L., Comes, O., Zuberi, K., Badrawi, R., Chao, P., Franz, M.,
673 Grouios, C., Kazi, F., Lopes, C.T., et al. (2010). The GeneMANIA prediction server: biological
674 network integration for gene prioritization and predicting gene function. *Nucleic Acids Research*
675 38:W214-W220. 10.1093/nar/gkq537.
- 676 Wei, X., Qiu, J., Yong, K., Fan, J., Zhang, Q., Hua, H., Liu, J., Wang, Q., Olsen, K.M., Han, B.,
677 et al. (2021). A quantitative genomics map of rice provides genetic insights and guides breeding.
678 *Nat Genet* 53:243-253. 10.1038/s41588-020-00769-9.
- 679 Wei, X., Chen, M., Zhang, Q., Gong, J., Liu, J., Yong, K., Wang, Q., Fan, J., Chen, S., Hua, H., et
680 al. (2024). Genomic investigation of 18,421 lines reveals the genetic architecture of rice. *Science*
681 385:eadm8762. doi:10.1126/science.adm8762.
- 682 Ying, J., Chen, Z., Wang, Z., Jiang, W., Wang, C., Yuan, Z., Su, H., Kong, H., Yang, F., and
683 Dong, N. (2025). SeedBench: A Multi-task Benchmark for Evaluating Large Language Models
684 in Seed Science.
- 685

686
687 **Figure 1. Overview of SeedLLM development and automated evaluation. (A)** The
688 RiceCorpus is a comprehensive dataset of scientific publications and books related to rice,
689 encompassing 1.38 million academic papers in both Chinese and English. **(B)** The RiceCorpus
690 was used to train the base model, Qwen2.5, which was fine-tuned with the riceQA dataset to
691 specialize SeedLLM for rice-specific questions. GeneralCorpus, consisting of datasets like
692 Curated FineWeb-Edu and Curated OpenCoder, was used for general language model training.
693 RiceQA, designed for rice biology, includes Key-info QA (extracting key data from
694 RiceCorpus), Bad-case QA (addressing difficult scenarios), and Graph QA (using a knowledge
695 graph for relevant questions). GeneralQA, which includes Curated Infinity Instruct and Curated
696 AutoIF, further fine-tunes the model for instruction-following and automated inference tasks. **(C)**
697 The pre-trained model was evaluated on two tasks: Gen-QA-ACC for open-ended questions and
698 PPL-MCQ-ACC for multiple-choice questions. SeedLLM's performance was assessed using
699 rice-specific datasets from the Agr series, which included single-choice, multiple-choice, and fill-
700 in-the-blank questions. The number of questions across these datasets was also recorded. **(D)** The
701 pre-trained model outperforms the baseline on both Gen-QA-ACC and PPL-MCQ-ACC,
702 confirming the effectiveness of the pretraining process. **(E-G)** SeedLLM demonstrates its
703 superiority over other LLMs in rice-specific tasks in terms of accuracy, F1 score, and ROUGE.
704 Accuracy measures the proportion of correct predictions, F1 score balances precision and recall,
705 and ROUGE evaluates the overlap between model-generated outputs and reference texts. **(H)**
706 SeedLLM exhibited robust performance on general-purpose tasks across multiple datasets,
707 including CMMLU (Chinese multitask understanding), GSM8K (grade school math), BBH
708 (beyond current model capabilities), and MMLU (general multitask language understanding). All
709 evaluations in (D-H) utilized an automated pipeline to extract model responses and compare
710 them to a reference key for correctness.

711

712

713 **Figure 2. Human-centric evaluation confirms SeedLLM's superior performance compared**

714 **to general-purpose LLMs in HumanDesignRiceQA. (A)** Overview of human-centric

715 evaluation. HumanDesignRiceQA, a dataset comprising question-answering pairs with plant

716 biology expertise derived from academic publications, was constructed. SeedLLM, other

717 general-purpose LLMs, and undergraduate responses were tasked with answering questions from

718 this dataset. A panel of evaluators, experts in rice biology, ranked the quality of responses for

719 each question from best to worst. Evaluators also assigned grades based either on the provided

720 answer key or their own expertise. **(B)** Distribution of questions in HumanDesignRiceQA is

721 categorized into three difficulty levels basic, intermediate, and advanced. **(C)** Distribution of

722 evaluator educational backgrounds. All evaluators possess academic training in agronomy, with

723 degrees ranging from Bachelor's to PhD. Experts are defined as individuals who have published

724 research, filed patents, or contributed to rice variety development in the past five years. **(D)**

725 Human evaluation scores evaluated in the HumanDesignRiceQA dataset. **(E)** Human-assigned

726 rankings from top1 to top5 for responses generated by SeedLLM and other LLMs.

727

728
729 **Figure 3. Encoding multiomics data into the rice biological knowledge graph for SeedLLM-**
730 **based agricultural query responses. (A)** Schematic overview for encoding transcriptomic,
731 proteomic, and genome annotation data as graph structures. Transcriptional and proteomic events
732 for each gene are represented as sentences and converted into triples. These triples are
733 transformed into nodes (representing rice gene IDs, transcriptional and translational events, and
734 experimental metadata) and edges (depicting relationships). Gene ID nodes are linked to genome
735 annotation attributes such GO terms and KEGG. **(B)** Overview of the Rice Biological
736 Knowledge Graph. Node types are color-coded to represent different data categories. The graph
737 includes 35,599 unique rice AGIS IDs identified from rice transcriptomic and proteomic studies,
738 shown with their relationships as edges. **(C)** SeedLLM-KG working mechanism. SeedLLM
739 processes agricultural queries by decomposing them into sub-queries, retrieving answers from
740 either the LLM's database or the knowledge graph. These sub-answers are integrated to form a
741 comprehensive response. Non-agricultural queries are handled directly by SeedLLM. **(D)**
742 Correctness of SeedLLM-KG and other LLMs were assessed by human evaluators using
743 advanced-level questions of the HumanDesignRiceQA dataset. Asterisks indicate statistical
744 significance between comparisons (t-test, $P < 0.001$). **(E)** Win rate of SeedLLM-KG against
745 various LLMs. SeedLLM-KG is considered to win if it receives a higher human evaluation
746 scores than the other model on the same questions from HumanDesignRiceQA dataset. The win
747 rate represents the percentage of questions where SeedLLM-KG outperforms each model, as
748 labeled in the plot. **(F)** Reasoning ability of SeedLLM-KG and other LLMs were assessed by
749 human evaluators using advanced-level questions of the HumanDesignRiceQA dataset. n.s.
750 indicates no significant difference between comparisons (t-test, $P < 0.05$).
751

752

753 **Figure 4. Comparison of response quality between SeedLLM and other LLMs. (A)** Models

754 were tasked with generating responses to the query: "Does the rice gene AGIS_Os06g035130

755 respond to various environmental conditions?" The models tested include SeedLLM-KG,

756 DeepSeek-R1, and OpenAI GPT-4o1. Response quality was assessed by verifying content

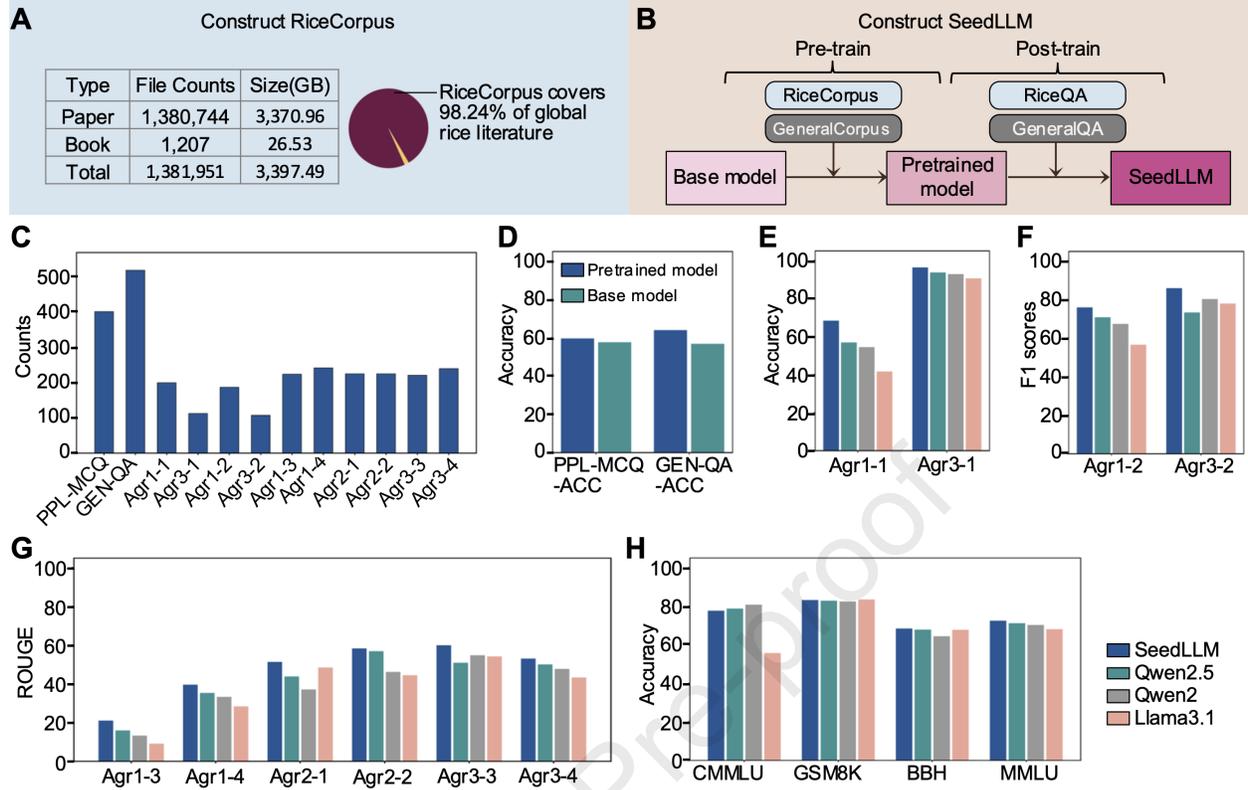
757 against RiceCorpus literature. Correct content is highlighted in blue, model reasoning in green,

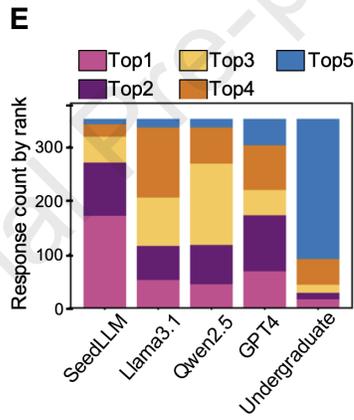
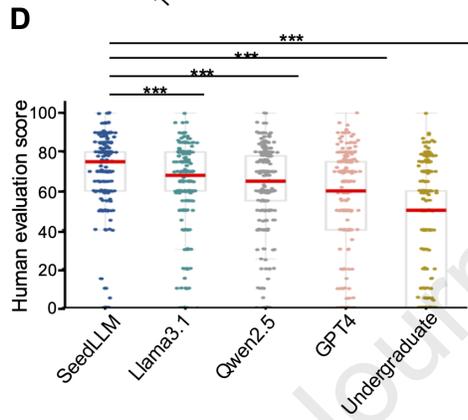
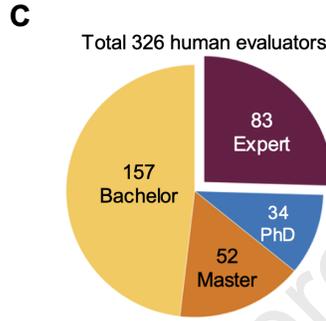
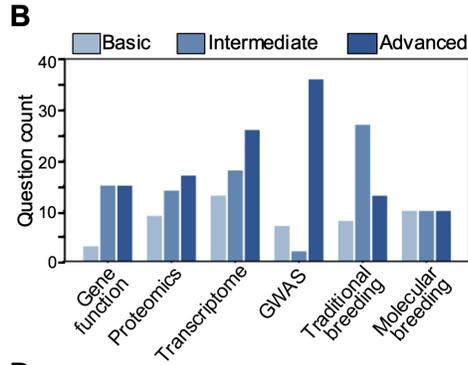
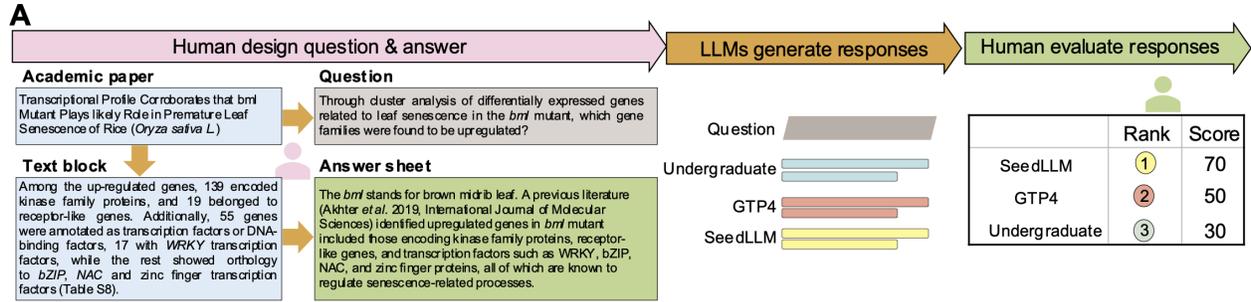
758 and inability to answer in red. **(B)** A list of references identified in the literature search that

759 corroborate the responses. Note that none of the LLMs generated references post-response.

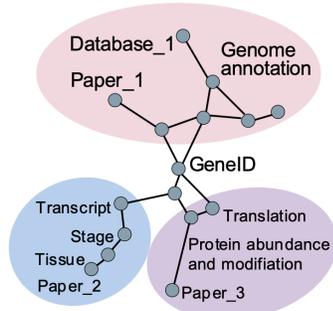
760

Journal Pre-proof

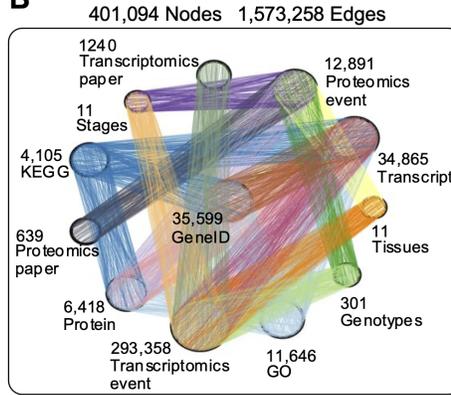




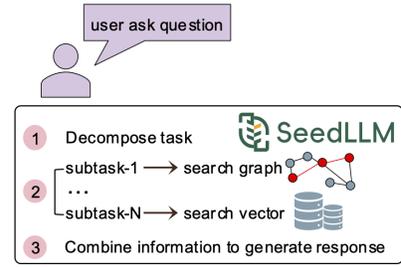
A



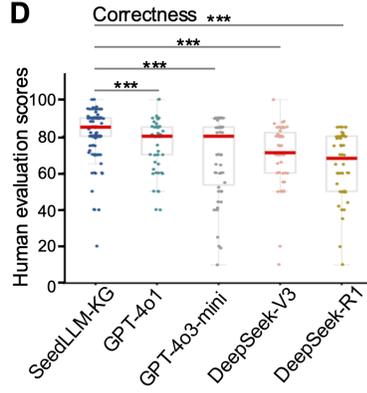
B



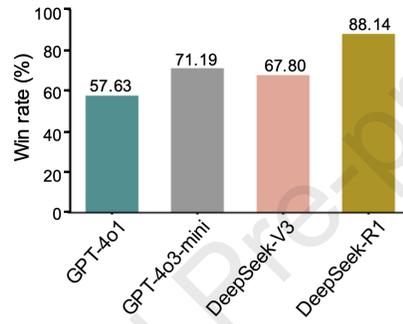
C



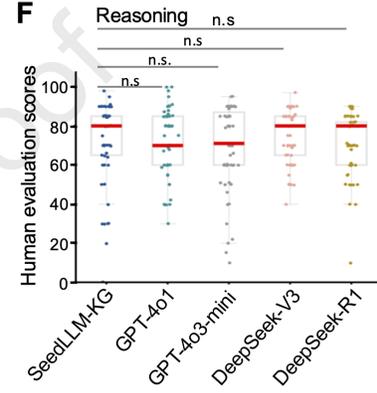
D



E



F



A

Does the rice gene AGIS_Os06g035130 respond to various environmental conditions?

Correct Answer

Reasoning

Unable to Answer

 SeedLLM

Yes, the rice gene AGIS_Os06g035130 does respond to various environmental conditions. It shows changes in transcriptional levels in response to different environmental stimuli. This gene's expression is measured under a wide range of experimental conditions, including **drought stress**, **anaerobic germination**, **early chilling stress**, and **osmotic stress**. Specifically, AGIS_Os06g035130's transcript expression is **upregulated** during the **seedling stage** under certain conditions, such as **osmotic stress**, and **in the panicle** during specific developmental stages under **drought stress**. Additionally, its expression is studied in different rice varieties and growth stages, highlighting its dynamic regulatory mechanisms in response to various stress conditions.

 deepseek

Okay, **so I need to figure out if** the rice gene AGIS_Os06g035130 responds to various environmental conditions and whether its transcriptional or translational levels change under different environmental stimuli. **Let me start by breaking down the problem.** **First, I need to confirm** the correct identifier for the gene. The user provided AGIS_Os06g035130. Wait, sometimes gene IDs can have different formats depending on the database. For example, Rice Genome Annotation Project (RGAP) uses LOC_Os06g035130, where "LOC" is a common prefix. The user's ID starts with AGIS, which might be a typo or from a different database. **Let me check** that first **because if** the ID is incorrect, any further steps **could be wrong**. **Assuming** AGIS_Os06g035130 is correct, **but maybe it's supposed to be** LOC_Os06g035130. **Maybe** "AGIS" is part of another naming convention? **Alternatively, perhaps it's from the AgBase database?** Hmm. **Alternatively, maybe** it's a typo and should be RAP-DB's ID, which **might be** Os06g0135200 or similar. Wait, RGAP and RAP-DB are two major databases for rice genes. Let me cross-reference this. The format Os06g035130 would correspond to chromosome 6, locus 035130. But RGAP uses LOC_Os06g035130. For example, LOC_Os06g035130 is a known gene encoding a protein kinase. **If** AGIS is part of another system, maybe **I need to check** that. **Alternatively, maybe** the user intended to write "LOC" but wrote "AGIS". Let me proceed with LOC_Os06g035130 as a possible correct ID and see what information I can find. **Now**, moving on to the main question: does this gene respond to environmental conditions? Changes in transcription or translation levels? I know that plant genes often respond to environmental stimuli such as drought, salinity, cold, heat, pathogen attack, etc. Transcription factors or kinases might be involved in stress responses. Looking up LOC_Os06g035130: According to RGAP (Rice Genome Annotation Project), this gene is annotated as a receptor-like protein kinase. Protein kinases are often involved in signal transduction pathways in response to environmental stresses. **To check expression data under different conditions:** There are databases like RiceXPro (Rice Expression Profile Database), which provide expression profiles under various conditions. Also, CREP (Database of Composit).

 OpenAI

The rice gene AGIS_Os06g035130, identified in the Gene Expression Omnibus (GEO) dataset GSE73609, is involved in the plant's response to various environmental conditions. This dataset examined global gene expression patterns in rice under different cultivation methods and climatic conditions, highlighting how environmental factors influence gene expression. **Additionally**, research on Environmental Gene Regulatory Influence Networks (EGRINs) in rice has shown that genes like AGIS_Os06g035130 are part of complex networks that coordinate responses to environmental stimuli such as water deficit and high temperature. These networks involve interactions between transcription factors and target genes, modulating gene expression in response to environmental changes. **Furthermore**, studies on stress-induced gene expression in rice have demonstrated that environmental stresses can lead to changes in DNA methylation patterns near stress-responsive genes. This suggests that genes like AGIS_Os06g035130 may undergo transcriptional and epigenetic modifications in response to environmental stimuli. In summary, AGIS_Os06g035130 is responsive to various environmental conditions, with its expression levels and epigenetic modifications adjusting to environmental stimuli.

B

1. Kaur, S., Seem, K., Duhan, N., Kumar, S., Kaundal, R., & Mohapatra, T. (2023). [Transcriptome and Physio-Biochemical Profiling Reveals Differential Responses of Rice Cultivars at Reproductive-Stage Drought Stress](https://doi.org/10.3390/ijms24021002). *International Journal of Molecular Sciences*, 24(2), 1002. <https://doi.org/10.3390/ijms24021002>
2. Baldoni, E., Bagnaresi, P., Locatelli, F., Mattana, M., & Genga, A. (2016). [Comparative leaf and root transcriptomic analysis of two rice japonica cultivars reveals major differences in the root early response to osmotic stress](https://doi.org/10.1186/s12284-016-0098-1). *Rice*, 9(1), 23. <https://doi.org/10.1186/s12284-016-0098-1>
3. Han, B., Ma, X., Cui, D., Wang, Y., Geng, L., Cao, G., Zhang, H., & Han, L. (2019). [Comprehensive evaluation and analysis of the mechanism of cold tolerance based on the transcriptome](#)