



1 **Global Scenario Reference Datasets for Climate Change Integrated Assessment** 2 with Machine Learning 3 Yi-Ming Wei^{1, 2, 3, 4, *}, Li-Jing Liu^{1, 2, 3, 4, *}, Shu-Xin Zhang^{1, 2, 3}, Jia-Xin Liu^{1, 2, 3}, Luo Zhao^{1, 2, 3}, Rong-Gang Cong^{1, 2, 3, 4, *}, Qiao-Mei Liang^{1, 2, 3, 4, *}, Te Han^{1, 2, 3, 4}, Xiao-Chen Yuan^{1, 2, 3, 4}, Yi-Ming Chen^{1, 2, 3}, 4 Yu-Xiang Hu^{1, 2, 3}, Ze-Qi Xu^{1, 2, 3}, Hong-Bo Chen^{1, 2, 3}, Yu-Xuan Xiao^{1, 2, 3}, Peng Wang^{1, 2, 3}, Song-Yang 5 Yan^{1, 2, 3}, Xiao-Ling Huang^{1, 2, 3}, Tian-Yuan Wang^{1, 2, 3}, Xiao-Qi Li^{1, 2, 3}, Hao-Ran Xu^{1, 2, 3}, Wen-Chang 6 Zhao^{1, 2, 3}, Biying Yu^{1, 2, 3, 4}, Baojun Tang^{1, 2, 3, 4}, Lan-Cui Liu⁵, Hua Liao^{1, 2, 3, 4}, Zhi-Qiang Li⁶, Rui Cui⁶ 7 8 9 ¹Center for Energy and Environmental Policy Research, Beijing Institute of Technology, Beijing, China 10 ² School of Management, Beijing Institute of Technology, Beijing, China 11 ³ Beijing Lab for System Engineering of Carbon Neutrality, Beijing, China 12 ⁴ NSFC Basic Science Center for Energy and Climate Change, Beijing, China 13 ⁵ Business School, Beijing Normal University, Beijing, China 14 ⁶ Information Technology Center, Beijing Institute of Technology, Beijing, China 15 16 **Correspondence:** 17 Yi-Ming Wei (wei@bit.edu.cn) 18 Li-Jing Liu (liulijing@bit.edu.cn) 19 Rong-Gang Cong (leocongbit@bit.edu.cn) 20 Qiao-Mei Liang (liangqiaomei@bit.edu.cn) 21 22 Abstract. The deepening of global climate change research and increasingly complex integrated 23 assessment methods generate large amounts of heterogeneous data. The rapid development of 24 artificial intelligence (AI) models, particularly large language models (LLMs) and deep learning 25 techniques, has enhanced the ability to handle vast data, providing new approaches and perspectives 26 for climate analysis. To address the demand for multi-dimensional and comparable scenario design 27 in climate change prediction and policy simulation, this study employs hybrid machine learning 28 techniques to collect and process scenario data from existing literature, developing the Global 29 Climate Scenario Reference datasets (GCSR). The GCSR incorporates data from approximately 30 90,000 articles across multiple temporal and spatial scales and extracts approximately 53,185 scenarios. With its large scale, extensive coverage, and detailed classification, the GCSR provides 31 32 a robust foundation for climate change prediction, risk assessment, mitigation policy, and adaptation 33 strategy planning, supporting scenario design in related fields.

34



35 1. Introduction

| 36 | The global climate crisis poses the greatest threat to both humanity and ecosystems (Zabala, |
|----|--|
| 37 | 2021). Since the Stockholm Declaration was signed in 1972, global attention to climate change has |
| 38 | continued to rise. The international community has successively adopted several important |
| 39 | initiatives, such as the 1997 Kyoto Protocol and the 2015 Paris Agreement, which established |
| 40 | national mitigation targets and emission reduction mechanisms (Masood, 1997; Tollefson and Weiss, |
| 41 | 2015). Assessing climate change and managing climate governance have become crucial challenges |
| 42 | that require urgent attention. |
| 43 | In recent years, Machine Learning (ML) and Artificial Intelligence (AI) have emerged as |
| 44 | powerful tools driving technological advancements. How to more effectively apply these tools to |
| 45 | tackle climate change has garnered widespread attention from both the ML community and climate |
| 46 | experts (Taddeo et al., 2021). Considering the exceptional capabilities of AI, particularly Large |
| 47 | Language Models (LLMs), in handling large-scale heterogeneous data, extracting latent variables, |
| 48 | and managing incomplete and uncertain data (Toetzke et al., 2022; Wu et al., 2024), AI technology |
| 49 | is deemed to have significant advantages in assisting data collection to support decision-making and |
| 50 | guiding future development plans (Rolnick et al., 2022). By extending to climate change governance |
| 51 | goals, AI technology offers researchers new analytical frameworks and tools, enhancing climate |
| 52 | prediction, policy simulation, risk assessment, and adaptation strategy planning. At the same time, |
| 53 | it supports policymakers in formulating more scientific and effective response measures. |
| 54 | Global climate change and greenhouse gas emissions involve both climate and economic |
| 55 | systems, and simulation approaches represented by integrated assessment models have become a |
| 56 | powerful tool for global climate governance (Wei et al., 2023). Specifically, human economic |
| 57 | activities generate carbon emissions, which enter the carbon cycle and alter greenhouse gas |
| 58 | concentration and temperature, thereby affecting the climate system. Changes in the climate system |
| 59 | in turn feed back to the economic system by impacting agricultural production (Ren et al., 2023; |
| 60 | Bousfield et al., 2024), water resource distribution (Padrón et al., 2020; Utsumi and Kim, 2022), |

61 and energy supply and demand (van Ruijven et al., 2019; Perera et al., 2020), etc. In climate change

62 assessment and response research, scenario design and simulation are the most critical means used

63 to describe the uncertainties between the climate and economic systems (O'Neill et al., 2020). The

64 significant rise in climate change research (Kallbekken, 2015) and the escalating complexity of





corresponding simulation methods (Cointe et al., 2019) have resulted in the generation of a largenumber of heterogeneous scenario data.

Climate scenarios not only provide core hypotheses for exploring future development pathways 67 68 and assessing the impacts of current decisions (Rogelj et al., 2019), but also facilitate the 69 identification of blind spots and the generation of creative solutions (Finch et al., 2024). The 70 adequacy of current scenario settings and their ability to accurately reflect the volatile realities of 71 the future have become one of the widely focused issues in the academic community (Weber et al., 72 2018; Anderson and Jewell, 2019; O'Neill et al., 2020; Guivarch et al., 2022). Institutions and 73 researchers from diverse backgrounds (such as knowledge systems, disciplines, and experiences) may develop different scenarios (Cornell et al., 2013; Trutnevyte et al., 2016; Carrington and 74 75 Stephenson, 2018), leading to significant variations in research results (Kriegler et al., 2016; van 76 Vuuren et al., 2020). Therefore, it is essential to develop a more objective, detailed, and 77 comprehensive climate change scenario reference dataset. This dataset focuses on the climate 78 change-related literature that carries out simulation analysis and extracts relevant descriptive 79 information for scenario design. The scenario descriptions of this dataset help to enhance scholars' 80 understanding of the simulation of climate and economic-related variables and provide a reference 81 for designing appropriate scenarios for research objectives, thus it can serve as a complement to 82 existing scenario databases that focus on simulation results under different scenarios (such as the 83 IPCC AR6 database).

84 Here, this study leverages AI's capabilities in processing large-scale heterogeneous data and 85 efficiently extracting value information from complex and diverse data (Toetzke et al., 2022; Wu et 86 al., 2024), attempts to employ hybrid ML methods to collect and process scenario data from existing 87 literature, and develops the Global Climate Scenario Reference datasets (GCSR). Approximately 88 90,000 climate change-related studies, spanning different temporal and spatial scales, have been 89 collected. The content covers multiple key dimensions: cause analysis, impact assessment, 90 prediction methods, and governance strategies of climate change. Drawing on these extensive 91 literature resources, this study uses an LLM to extract and organize around 53,185 scenario data points. Employing ML techniques such as keyword recognition and topic extraction, this study 92 93 categorizes the scenarios more precisely. Researchers can then use issue-specific keywords to 94 swiftly identify related scenarios in the GCSR, thereby guiding their exploration of uncertainties in 95 scenario design research. Policymakers, after locating relevant scenarios, can further filter related





- 96 literature and scenario results to inform policy design. It provides powerful data support for scenario
- 97 design in research areas such as global climate change prediction, risk assessment, mitigation policy
- 98 simulation, and adaptation strategy planning.

99 2. Methods

100 This study provides a complete process from literature identification, literature collection, 101 scenario extraction, scenario keyword recognition, and topic classification, to the development of 102 GCSR (see Figure 1). Based on a thematic analysis of the IPCC reports (IPCC, 2022), the major 103 issues related to climate change can be grouped into four key categories: its underlying causes, 104 associated impacts and risks, future emission and temperature projections, and governance and 105 response strategies. Accordingly, the literature can be classified into four main categories: climate 106 change causes, impacts, predictions, and governance. The Web of Science core database is used to 107 check and supplement the search formula of various literature, ultimately finalizing the search 108 formula and collecting the relevant studies. Subsequently, this study utilizes the LLM to extract 109 scenarios and reclassifies each scenario into four categories. To facilitate easy access and use by 110 researchers, natural language processing techniques are applied to the initial scenario dataset. This 111 is followed by high-frequency word analysis, keyword extraction, and topic classification using 112 BERTopic semantic recognition, allowing for more specific subcategories of scenarios. Finally, the 113 scenario data is organized and stored according to these refined categories.







114

115 Figure 1: Overview of the methodological framework for developing the GCSR.

116 **2.1 Literature classification and collection**

117 The subcategories of the four issues are shown in Figure 2(a). Climate change causes can be 118 attributed to both natural and human factors. In terms of human factors, economic activities, 119 lifestyles, and energy use generate greenhouse gas emissions; changes in land use patterns affect 120 carbon sequestration. The impacts of climate change are observed across both natural and human 121 systems. In natural systems, climate change alters the structure and function of ecosystems and 122 exacerbates extreme weather events. In the human system, climate change threatens water and food 123 security and human health and has cascade effects on key economic sectors. Climate change 124 predictions mainly involve estimating future trends in key climate elements such as temperature rise 125 and precipitation. Climate change governance mainly covers two core strategies: mitigation and 126 adaptation. Mitigation, beyond the application of emission reduction technologies, relies on a series 127 of emission reduction policies, including mainstream policy measures such as international 128 cooperative emission reduction mechanisms, carbon pricing mechanisms, and power market





129 reforms.

130 This study first identifies initial search formulas for various types of literature. On August 6, 131 2024, a search was performed using the Web of Science Core Collection database, which 132 encompasses the Science Citation Index Expanded (SCI-E) and Social Science Citation Index 133 (SSCI). During the search process, keywords from highly relevant literature are reviewed to refine 134 initial search formulas. Subsequently, search formulas are revised, and the search is conducted again. 135 After multiple iterations of verification and modification, the study finalizes search formulas for 136 each type of literature (see Supporting Information 1). In total, approximately 90,000 articles are 137 collected.



Figure 2: Datasets structure and number of scenarios. (a) Literature classification standards. (b) Thestructure of the datasets.

141 2.2 Scenario extraction using LLM

138

This study utilizes the excellent natural language deep semantic understanding ability of the LLM to efficiently identify and extract key information from complex texts (Ray, 2023; Qu and Wang, 2024). After multiple rounds of testing and evaluation, this study selects the DeepSeek-V2 model for scenario extraction from the literature, as it outperforms others in terms of accuracy and operational efficiency. DeepSeek-V2 supports processing context information up to 128,000 tokens, making it excellent at handling long documents and performing complex scenario extraction tasks with high accuracy in capturing key scenario-related information (DeepSeek-AI, 2024a; DeepSeek-





- 149 AI, 2024b). To precisely identify scenarios, this study optimizes the design of input prompts for
- 150 LLM, developing an optimal prompt structure capable of guiding LLM to accurately capture
- 151 scenarios. The specific scenario extraction process includes the following four main steps:
- 152 (1) Processing of Input Literature Data
- 153 The literature data undergoes text preprocessing and segmentation, transforming it into a
- 154 sequence of tokens $X = [x_1, x_2, ..., x_n]$. *n* represents the number of tokens in the literature.
- 155 (2) Generation of Latent Vectors
- 156 The Multi-head Latent Attention (MLA) mechanism is employed to transform the input tokens
- 157 into compact latent vectors, denoted as z_i , and can be expressed as Eq. (1).
- 158 $MLA(X) = Concat(Head_1, Head_2, \dots, Head_h)W^0$ (1)
- 159 Where $Head_i$ represents the output of the *i-th* attention head. The final multi-head attention 160 result is obtained by concatenating the outputs of all heads and combining them with the output 161 weight matrix W^o . Each attention head $Head_i$ is computed as Eq. (2):

$$Head_{i} = softmax \left(\frac{QW_{i}^{Q}\left(KW_{i}^{K}\right)^{T}}{\sqrt{d_{k}}}\right) \left(VW_{i}^{V}\right)$$

$$(2)$$

162

- 163 Where Q is the matrix of Queries, K is the matrix of Keys, V is the matrix of Values, 164 W_i^Q, W_i^K, W_i^V are learnable weight matrices for the *i-th* head that project the input queries, keys, 165 and values into appropriate spaces. The *softmax* operation normalizes the attention scores. 166 Through this mechanism, the model captures diverse information via multi-head parallel processing, 167 enabling it to integrate the most representative features and generate refined latent vectors z_i (Han 168 et al., 2024). These refined vectors provide an efficient input representation for subsequent scenario 169 recognition and extraction tasks.
- 170 (3) Scenario Recognition and Extraction

171 The core of scenario extraction lies in accurately recognizing and outputting scenarios, which 172 is achieved through the Mixture-of-Experts (MoE). The latent vectors z_i are input into multiple 173 experts, and the activation weights of each expert are computed based on a soft selection mechanism. 174 The experts relevant to the scenario extraction task are selectively activated for computation, and 175 the corresponding scenario representations are generated, as shown in Eq.(3).

$$z_i^* = \sum_{e=1}^k \alpha_{i,e} \cdot E_e(z_i)$$
(3)





177 Where z_i^* denotes the final latent vector, including scenario information extracted from 178 different experts. $\alpha_{i,e}$ is the selection weight of the *e*-th expert for the latent vector z_i , indicating 179 the expert's contribution to the input. $E_e(z_i)$ is the processing result of the *e*-th expert on the latent 180 vector z_i , typically a transformation or mapping function used to generate scenario-related 181 representations. Through the MoE mechanism, DeepSeek-V2 can selectively activate appropriate 182 experts, ensuring that only the scenario-relevant contents are processed, thereby improving the 183 efficiency and accuracy of scenario extraction. 184 (4) Output Generation 185 The latent vectors z_i processed by the expert network are further analyzed and aggregated to 186 generate the final scenario data \mathcal{Y} . The model outputs the extracted scenario information as 187 structured text, as shown in Eq.(3). $y = DeepSeekMoE(z_1^*, z_2^*, \dots, z_n^*)$

188 $y = DeepSeekhoE(z_1, z_2, ..., z_n)$ (4) 189 Where $(z_1^*, z_2^*, ..., z_n^*)$ denotes the processed latent vector, and y represents the extracted 190 scenario data. The design of this model perfectly aligns with the requirements of scenario extraction 191 tasks, efficiently handling multimodal data and diverse corpora. These output results consist of 192 paragraphs, sentences, or key information about the literature's scenarios, ultimately summarizing 193 the scenario descriptions and hypotheses.

In total, 136,754 scenarios are extracted from the literature. To improve the accuracy of the categorization, this study reclassifies scenarios using the LLM to avoid misalignment with their assigned categories or the omission of scenarios that exhibit characteristics of multiple categories. Additionally, a small sample of scenarios from each category is manually verified to assess the accuracy of the model's classification results. This verified subset is then used as a training set to reclassify the remaining scenarios. Ultimately, based on the four main categories, this study further subdivides the scenarios into 12 specific subcategories, as shown in Figure 2(a).

201 2.3 Semantic check

The initial scenarios contain issues such as unrecognizable special characters, misrecognized letters, and semantic confusion. To resolve these problems, this study systematically screens and corrects the data, addressing character recognition errors and ensuring the integrity and readability of the scenario datasets. To tackle the issues of semantic confusion, this study utilizes the Pythonbased LanguageTool (Jumanov and Karshiev, 2020), which combines predefined grammatical rules, 207





208 Following the automated error detection and suggestions, a manual review is conducted to 209 ensure detailed modifications, ensuring that each scenario description is semantically accurate and 210 logically coherent. Moreover, the extracted data source may not have real climate scenarios. The 211 scenario mentioned in the original paper may have other meanings. So, we use a supervised learning 212 method to judge the extracted data, and finally, only retain the data with real climate scenarios. 213 2.4 Scenario expansion using BERTopic 214 Based on the initially obtained scenario datasets, this study applies a series of methods, 215 including scenario coding, word frequency statistics, similarity calculation, and topic extraction, to

language patterns, and statistical models to achieve efficient grammatical checking.

216 conduct in-depth processing on each scenario dataset separately (Berio Fortini et al., 2023). These 217 steps remove duplicate scenarios and assign clear thematic tags to each group of scenarios, thereby 218 enhancing the structural clarity and operability of the scenario datasets. The optimization of the 219 scenario datasets can provide a clearer and more efficient retrieval environment for subsequent 220 scholars' research, which helps to promote in-depth research in related fields.

221 (1) Duplication removal

To enhance the retrieval efficiency of the scenario datasets, this study removes duplicate scenarios, effectively reducing data redundancy. Through this processing, scenario datasets become more concise and clearer, significantly improving the efficiency of users during retrieval and usage, allowing them to quickly and accurately locate the required key information.

226 (2) Text cleaning

This study tokenizes the text data, divides each document into independent lexical units, and converts all words to lowercase, eliminating the impact of case differences on the statistical results. Furthermore, to more accurately focus on the core content of the scenarios, this study removes punctuation, special characters, and stop words that appear frequently but contribute little to actual semantics, such as "the", "and", etc.

- 232 (3) High-frequency word statistics
- To provide basic data support for subsequent analysis of the scenario datasets, this study uses
 Natural Language Toolkit (NLTK) (Bird and Loper, 2004) to count high-frequency words in the text





- 235 data. Specifically, the program counts the word frequency in each scenario and selects the most
- 236 representative high-frequency words. Finally, this study outputs the top five high-frequency words
- 237 in each scenario.
- 238 (4) Keyword extraction
- To identify the core keywords in the scenario description, this study uses the Best Matching 25 (BM25) algorithm. Compared to traditional methods, BM25 can more accurately identify keywords that are representative of scenario texts, and it avoids the common issue of term frequency bias in long documents, thereby improving the accuracy and effectiveness of keyword extraction (Amati and Rijsbergen, 2002). The weight calculation for each term in the scenario y is shown in Eq.(5).

244
$$BM 25(t, y) = \sum_{i=1}^{n} IDF(t_i) \cdot \frac{\mathrm{TF}(t_i, y) \cdot (k_i + 1)}{\mathrm{TF}(t_i, y) + k_1 \cdot \left(1 - b + b \cdot \frac{|y|}{\mathrm{avgdl}}\right)}$$
(5)

245
$$IDF(t) = log\left(\frac{N - df(t) + 0.5}{df(t) + 0.5} + 1\right)$$
(6)

Inverse Document Frequency (IDF) is used to measure the scarcity of a term t_i across the 246 entire scenario dataset. N is the total number of scenarios in the datasets, and df(t) is the 247 248 number of scenarios containing the term. $f(t_i, y)$ represents the term frequency (TF) of t_i in the scenario y. k_1 and b are adjustment parameters to balance the influence of term frequency 249 250 and document length on the score. |y| is the length of the scenario y, and avgdl is the average 251 length of all scenario texts. The BM25 algorithm not only optimizes the recognition of common 252 words but also comprehensively considers TF and IDF, making the extracted keywords more aligned 253 with the core content of the scenario.

254 (5) Topic extraction and classification

To effectively extract and classify the groups of scenarios, this study adopts BERTopic, a topic 255 modeling method based on Transformer. Traditional topic classification methods typically focus on 256 the surface information of words, neglecting the contextual semantics and deep relationships 257 between words. When dealing with semantically complex or syntactically flexible texts, this 258 259 limitation becomes particularly prominent, often resulting in information loss or misjudgment. However, BERTopic is a topic modeling method based on Transformer, which combines deep 260 261 learning and clustering techniques. It is capable of handling large-scale text, automatically and 262 efficiently identifying topics within the text (Grootendorst, 2022). Furthermore, BERTopic supports





| 263 | dynamic updates and extensions, allowing it to adapt to the continuous evolution of scenario datasets. |
|-----|--|
| 264 | $\mathbf{v} = \mathrm{ST}(y), \mathbf{v}' = UMAP(v), Cluster_i = \{\mathbf{v}' Density(\mathbf{v}') \ge Threshold \} $ (7) |
| 265 | The BERTopic model consists of three main steps. Firstly, it uses a pre-trained model, sentence- |
| 266 | transformers (all-MiniLM-L6-v2), to encode the scenario text into a high-dimensional vector v , |
| 267 | allowing the semantic relationships between texts to be effectively reflected. $ST(y)$ represents the |
| 268 | process of encoding scenario text using sentence-transformers. Furthermore, the Uniform Manifold |
| 269 | Approximation and Projection (UMAP) algorithm is used to reduce the dimension of the semantic |
| 270 | vector, while retaining the location information between the documents (McInnes et al., 2018). |
| 271 | Finally, the Hierarchical Density-Based Spatial Clustering of Applications with Noise (HDBSCAN) |
| 272 | algorithm is employed to cluster the dimensionality-reduced scenario text vectors. Density (\dot{v}) |
| 273 | denotes the final scenario classifications obtained through HDBSCAN clustering. Threshold |
| 274 | serves as the minimum density threshold to filter out text clusters with higher density. |
| 275 | Through this process, this study can automatically and efficiently identify topics in scenarios |

and output keywords for each topic category through Class-based TF-IDF, providing effective support for the efficient categorization and dynamic evolution of scenarios. This ensures the longterm adaptability and sustainability of the model.

279 3. Results and discussion

280 Figure 3 illustrates the number of scenarios for four categories across different years, as well 281 as the proportion of scenarios in each sub-category. Overall, after 2010, there was a significant 282 increase in the number of scenarios related to the climate change causes and impacts. In contrast, 283 the number of scenarios concerning prediction and governance did not show a slow growth trend 284 until 2015. Regarding the climate change causes, approximately 70% - 80% of the scenarios focused 285 on greenhouse gas emissions, around 10% - 20% concentrated on land use patterns, and the remaining scenarios (less than 10%) paid attention to natural factors. In terms of the climate change 286 287 impacts, the number of scenarios centered on ecological structure and function impacts ranked first. 288 However, over time, the quantity of such scenarios exhibited a decreasing trend. Meanwhile, the 289 focus on the impacts on economic sectors gradually increased, with their proportion approaching 290 25% in recent years. With respect to climate change governance, mitigation policies received the 291 most attention, with the proportion of scenarios in this category remaining at around 60%. It was 292 followed by scenarios related to mitigation technology (about 30%), and the remaining 10% focused









299 Figure 4 displays the word frequency statistics results of top words under different scenario categories. The statistics indicate that in scenarios focusing on the climate change causes, words 300 301 such as "carbon", "emission", "CO2" appear with relatively high frequency. This reflects that 302 relevant studies primarily concentrate on the core factor that emission activities trigger climate 303 change. In scenarios centered around climate change prediction, words like "future", "temperature", 304 and "model" occur frequently. This demonstrates a research trend of utilizing various models to 305 predict indicators such as future temperature rise. In addition, studies focusing on the climate change impacts tend to emphasize areas like "energy", "temperature", and "water". As for research on 306 307 climate change governance, scenarios are often designed around elements such as "tax", "electricity", 308 and "price".







309

Figure 4: Word cloud of top words for different scenario categories. (a) scenarios in climate change
causes. (b) scenarios in climate change prediction. (c) scenarios in climate change impacts. (d) scenarios

312 in climate change governance.

313 GCSR stores a wealth of scenario design cases. Its detailed classification and labeling enable 314 researchers and policymakers to efficiently identify and utilize scenario information relevant to 315 specific issues. For instance, if researchers are studying the impact of climate change on air pollution, they can use keywords like "climate change impact," "human systems," and "health and wellbeing" 316 317 to find scenarios related to "air", thereby easily obtaining closely related scenario cases. These cases 318 contain detailed scenario descriptions, serving as invaluable reference materials for researchers. 319 Policymakers who have screened the GCSR for scenarios relevant to their policy issues can further 320 utilize the literature information provided by the GCSR to directly access the original literature. By 321 comparing simulation results under different scenarios, policymakers can assess the potential 322 impacts and feasibility of various policy options. Hence, the significance of the GCSR dataset lies 323 in the following: ①Enhancing Research Efficiency: The GCSR significantly reduces the time 324 researchers and policymakers spend searching for relevant cases and literature through efficient 325 keyword searching and scenario matching functions, thereby improving their work efficiency and 326 research quality. 2 Promoting Knowledge Sharing: The scenario cases in the GCSR originate from 327 global researchers and practitioners, fostering knowledge sharing and exchange within the field of 328 scenario design research and contributing to a more comprehensive and in-depth understanding. (3)





| 329 | Enh | ancing Scientific Decision-Making: Based on scenario cases in the GCSR, policymakers can |
|-----|-------|--|
| 330 | conc | duct more scientific and objective decision-making analysis, reducing the blindness and |
| 331 | unce | ertainty in decision-making and enhancing the relevance and effectiveness of policies. $\textcircled{4}$ |
| 332 | Driv | ing Disciplinary Development: The establishment and use of the GCSR also promote the |
| 333 | cont | inuous improvement and development of scenario design methods, providing new perspectives |
| 334 | and | ideas for research in this field. |
| 335 | 4. | Data availability |
| 336 | | To facilitate the use by researchers, decision-makers, industry professionals, and other relevant |
| 337 | pers | onnel, the GCSR provides two forms of data storage: MySQL and EXCEL. The final results are |
| 338 | uplo | aded through Zenodo: https://doi.org/10.5281/zenodo.15536298 (Wei et al., 2025). The |
| 339 | struc | cture of the GCSR core data is shown in Figure 2 (b). |
| 340 | (1) | Literature: Stores information related to the source literature of scenarios, including the title of |
| 341 | | the literature, the name of the journal published, the year of publication, the DOI number, and |
| 342 | | keywords. This information provides comprehensive background support for the source of |
| 343 | | scenarios. |
| 344 | (2) | Scenario: Contains extracted and summarized scenario descriptions and assumptions. Each |
| 345 | | document may contain multiple scenarios, with each scenario stored as a separate piece of |

extracted by the Best Matching 25 (BM25), and the scenario topics based on BERTopicclustering.

content. In addition, it also includes the high-frequency words of the scenario, keywords

- 349 (3) Topic: Includes scenario topic categories generated by BERTopic clustering, the name of each
 350 topic category, and its core keywords, facilitating further organization and classification of
 351 scenario content.
- (4) Classification: Provides scenario classification based on LLM categorization, and hierarchy
 for labeling scenarios (including the first, second, and third-level categories shown in Figure
 2(a)), aiding in structured analysis from broad categories to specific scenarios.
- 355 5. Technical Validation

346

356 As machine learning increasingly relies on large datasets, there is a growing acknowledgment





of the importance of data quality (Gero et al., 2023). The data quality of the GCSR is affected by 357 358 several factors. The primary factor lies in whether the data sources are comprehensive and extensive, 359 capable of covering a majority of climate change scenarios, which directly impacts the training 360 capability of the model. Secondly, the effectiveness of scenario extraction is also crucial. It requires 361 the effective extraction of large-scale and heterogeneous scenarios while preserving the original 362 meaning of the scenarios, which helps enhance the model's understanding capability. Furthermore, 363 the accuracy of scenario classification is a key step in ensuring the overall reliability of the scenario 364 dataset, directly influencing the efficient organization and retrieval of information within the dataset. 365 Therefore, to improve data quality, this study adopts the following specific measures.

366 5.1 Ensure more comprehensive literature sources

367 In this study, the search query about climate change research was repeatedly modified and 368 validated according to the identified relevant literature topics and keywords, and the final search 369 query ensured comprehensive coverage of the literature. Additionally, literature not automatically 370 marked as successfully downloaded by the program was manually retrieved as a supplement.

371 5.2 Check the extracted scenarios' semantics

372 On the one hand, the scenarios extracted by the original LLM may contain issues such as 373 incorrect letter recognition and semantic confusion. This study utilizes a language tool to help 374 identify potential grammatical errors and unclear expressions in the initial scenarios. On the other 375 hand, the extracted data might not represent real climate scenarios. The term "scenario" mentioned 376 in the original papers may carry other meanings. To address this, we employ a supervised learning approach to evaluate the extracted data. By carefully reading the original texts, we classify 1,862 377 378 data samples, assigning a label of 1 to those that truly contain climate scenarios and 0 to those that 379 do not. Distinct from manual assessment, there might contain some generalization errors of such a 380 big data analysis. However, after conducting five-fold stratified cross-validation based on these 381 samples, the accuracy rate of the model reaches over 90% (see Table 1).

382 Table 1 Comparison between supervised learning and manual verification of scenario results

| Fold | Precision | Recall | F1 |
|------|-------------|-------------|-------------|
| 1 | 0.961111111 | 0.940217391 | 0.950549451 |
| 2 | 0.941489362 | 0.961956522 | 0.951612903 |
| 3 | 0.972826087 | 0.978142077 | 0.975476839 |





| 4 | 0.928205128 | 0.989071038 | 0.957671958 |
|---|-------------|-------------|-------------|
| 5 | 0.977142857 | 0.929347826 | 0.95264624 |

383 **5.3 Verify the scenario classification's training set**

To improve the classification accuracy, a small subset of scenarios is first classified using the LLM, and these preliminary classifications are then further verified by assessing the rationality of the features, content, and assigned categories. Inconsistencies are corrected to develop a highquality training set, which is subsequently used to classify additional scenarios. Additionally, to account for scenarios outside the four predefined categories, the overall recall rate is calculated by comparing the number of categorized scenarios to the total (see Table 2). The results show that the quality of the GCSR is reliable.

391 **Table 2** Comparison between LLM and manual verification of classification results.

| Classification | LLM classification vs. manual check | | | |
|----------------------------|--|-----------|--------|------|
| | | Precision | Recall | F1 |
| Climate change | Climate change caused by natural factors | 0.99 | 0.94 | 0.96 |
| causes | Greenhouse gas emissions | 0.92 | | 0.93 |
| | Land use patterns | 0.94 | | 0.94 |
| Climate change | Ecosystem structure and function | 0.80 | | 0.86 |
| impacts | Extreme weather events | 0.98 | | 0.96 |
| | Food production | 0.97 | | 0.96 |
| | Health and wellbeing | 0.92 | | 0.93 |
| | Economic sectors | 0.92 | | 0.93 |
| Climate change predictions | | 0.90 | | 0.92 |
| Climate policies | Mitigation policy | 0.92 | | 0.93 |
| and governance | Mitigation technology | 0.93 | | 0.93 |
| | Adaption measures | 0.93 | | 0.93 |

392 **6. Code availability**

393 Code, dataset, and some intermediate results are presented in Supporting Information 2.





Reference

| 395 | Amati, G. and Rijsbergen, C. J. V.: Probabilistic models of information retrieval based on measuring the |
|-----|--|
| 396 | divergence from randomness, ACM Trans. Inf. Syst., 20, 357-389, |
| 397 | https://doi.org/10.1145/582415.582416, 2002. |
| 398 | Anderson, K. and Jewell, J.: Debating the bedrock of climate-change mitigation scenarios, Nature, 573, |
| 399 | 348–349, https://doi.org/10.1038/d41586-019-02744-9, 2019. |
| 400 | Berio Fortini, L., Kaiser, L. R., Frazier, A. G., and Giambelluca, T. W.: Examining current bias and future |
| 401 | projection consistency of globally downscaled climate projections commonly used in climate |
| 402 | impact studies, Climatic Change, 176, 169, https://doi.org/10.1007/s10584-023-03623-z, 2023. |
| 403 | Bird, S. and Loper, E.: NLTK: The Natural Language Toolkit, In Proceedings of the ACL Interactive |
| 404 | Poster and Demonstration Sessions, pages 214–217. Barcelona, Spain, Association for |
| 405 | Computational Linguistics, 2004. |
| 406 | Bousfield, C. G., Morton, O., and Edwards, D. P.: Climate change will exacerbate land conflict between |
| 407 | agriculture and timber production. Nature Climate Change. https://doi.org/10.1038/s41558-024- |
| 408 | 02113-z. 2024. |
| 409 | Carrington G and Stephenson I. The politics of energy scenarios: Are International Energy Agency |
| 410 | and other conservative projections hampering the renewable energy transition? Energy Research & |
| 411 | Social Science 46 103-113 https://doi.org/10.1016/j.erss.2018.07.011.2018 |
| 412 | Cointe B. Cassen, C. and Nadaï A.: Organising Policy-Relevant Knowledge for Climate Action: |
| 412 | Integrated Assessment Modelling the IPCC and the Emergence of a Collective Expertise on |
| 414 | Sociacoonomia Emission Sociación Science & Tachnelogy Studios 22, 26,57 |
| 415 | https://doi.org/10.2087/at.65021.2010 |
| 415 | https://doi.org/10.2596//sts.05051,2019. Consell S. Davidsout E. Tvinster W. Takara I. D. Esser I. Chakay I. da Wit D. Langleis P. Milla |
| 410 | D. Mall D. Otta I. M. Dataraan A. Dahl C. and van Karkhoff I. C. Onaning vin Institution |
| 41/ | D., Moh, F., Otto, I. M., Petersen, A., Pont, C., and van Kerknon, L.: Opening up knowledge |
| 418 | systems for better responses to global environmental change, Environmental Science & Policy, 28, |
| 419 | 00-70, https://doi.org/10.1010/j.envsci.2012.11.006, 2015. |
| 420 | Deepseek-A: Deepseek-V2: A Strong, Economical, and Efficient Mixture-of-Experts Language Model, |
| 421 | https://doi.org/10.48550/arXiv.2405.04434, 2024a. |
| 422 | DeepSeek-AI: DeepSeek LLM Scaling Open-Source Language Models with Longtermism, |
| 423 | nttps://doi.org/10.48550/arXiv.2401.02954, 20246. |
| 424 | Finch, M., Older, M., Mahon, M., and Robertson, D.: Climate action and the vantage point of imagined |
| 425 | futures: a scenario-based conversation, npj Climate Action, 3, 45, https://doi.org/10.1038/s44168- |
| 426 | 024-00123-3, 2024. |
| 427 | Gero, K. I., Das, P., Dognin, P., Padhi, I., Sattigeri, P., and Varshney, K. R.: The incentive gap in data |
| 428 | work in the era of large models, Nature Machine Intelligence, 5, 565-567, |
| 429 | https://doi.org/10.1038/s42256-023-00673-x, 2023 . |
| 430 | Grootendorst, M.: BERTopic: Neural topic modeling with a class-based TF-IDF procedure, |
| 431 | https://doi.org/10.48550/arXiv.2203.05794, 2022. |
| 432 | Guivarch, C., Le Gallic, T., Bauer, N., Fragkos, P., Huppmann, D., Jaxa-Rozen, M., Keppo, I., Kriegler, |
| 433 | E., Krisztin, T., Marangoni, G., Pye, S., Riahi, K., Schaeffer, R., Tavoni, M., Trutnevyte, E., van |
| 434 | Vuuren, D., and Wagner, F.: Using large ensembles of climate change mitigation scenarios for robust |
| 435 | insights, Nature Climate Change, 12, 428-435, https://doi.org/10.1038/s41558-022-01349-x, 2022. |
| 436 | Han, T., Cong, RG., Yu, B., Tang, B., and Wei, YM.: Integrating local knowledge with ChatGPT-like |
| 437 | large-scale language models for enhanced societal comprehension of carbon neutrality, Energy and |
| 438 | AI, 18, 100440, https://doi.org/10.1016/j.egyai.2024.100440, 2024. |
| 439 | IPCC: Climate Change 2022: Impacts, Adaptation and Vulnerability. Contribution of Working Group II |
| 440 | to the Sixth Assessment Report of the Intergovernmental Panel on Climate Change [HO. Pörtner, |
| 441 | D.C. Roberts, M. Tignor, E.S. Poloczanska, K. Mintenbeck, A. Alegría, M. Craig, S. Langsdorf, S. |
| 442 | Löschke, V. Möller, A. Okem, B. Rama (eds.)], Cambridge University Press. Cambridge University |
| 443 | Press, Cambridge, UK and New York, NY, USA, 3056, https://doi.org/10.1017/9781009325844, |
| 444 | 2022. |
| 445 | Jumanov, I. I. and Karshiev, K. B.: Mechanisms for optimization of detection and correction of text errors |
| 446 | based on combining multilevel morphological analysis with n-gram models, Journal of Physics: |
| 447 | Conference Series, 1546, 012082, https://doi.org/10.1088/1742-6596/1546/1/012082, 2020. |
| 448 | Kallbekken, S.: National climate ambition must match international targets, Nature, |





| 449 | https://doi.org/10.1038/nature.2015.19077, 2015. |
|-----|---|
| 450 | Kriegler, E., Mouratiadou, I., Luderer, G., Bauer, N., Brecha, R. J., Calvin, K., De Cian, E., Edmonds, J., |
| 451 | Jiang, K., Tavoni, M., and Edenhofer, O.: Will economic growth and fossil fuel scarcity help or |
| 452 | hinder climate stabilization?, Climatic Change, 136, 7-22, https://doi.org/10.1007/s10584-016- |
| 453 | 1668-3, 2016. |
| 454 | Masood, E.: Kvoto agreement creates new agenda for climate research, Nature, 390, 649-650. |
| 455 | https://doi.org/10.1038/37686. 1997. |
| 456 | McInnes, L., Healy, J., and Melville, J.: UMAP: Uniform Manifold Approximation and Projection for |
| 457 | Dimension Reduction, https://doi.org/10.48550/arXiv.1802.03426, 2018. |
| 458 | O'Neill, B. C., Carter, T. R., Ebi, K., Harrison, P. A., Kemp-Benedict, E., Kok, K., Kriegler, E., Preston, |
| 459 | B. L., Riahi, K., Sillmann, J., van Ruijven, B. J., van Vuuren, D., Carlisle, D., Conde, C., Fuglestvedt, |
| 460 | J., Green, C., Hasegawa, T., Leininger, J., Monteith, S., and Pichs-Madruga, R.: Achievements and |
| 461 | needs for the climate change scenario framework, Nature Climate Change, 10, 1074-1084, |
| 462 | https://doi.org/10.1038/s41558-020-00952-0, 2020. |
| 463 | Padrón, R. S., Gudmundsson, L., Decharme, B., Ducharne, A., Lawrence, D. M., Mao, J., Peano, D., |
| 464 | Krinner, G., Kim, H., and Seneviratne, S. I.: Observed changes in dry-season water availability |
| 465 | attributed to human-induced climate change, Nature Geoscience, 13, 477-481, |
| 466 | https://doi.org/10.1038/s41561-020-0594-1, 2020. |
| 467 | Perera, A. T. D., Nik, V. M., Chen, D., Scartezzini, JL., and Hong, T.: Quantifying the impacts of climate |
| 468 | change and extreme climate events on energy systems, Nature Energy, 5, 150-159, |
| 469 | https://doi.org/10.1038/s41560-020-0558-0, 2020. |
| 470 | Qu, Y. and Wang, J.: Performance and biases of Large Language Models in public opinion simulation, |
| 471 | Humanities and Social Sciences Communications, 11, 1095, https://doi.org/10.1057/s41599-024- |
| 472 | 03609-x, 2024. |
| 473 | Ray, P. P.: ChatGPT: A comprehensive review on background, applications, key challenges, bias, ethics, |
| 474 | limitations and future scope, Internet of Things and Cyber-Physical Systems, 3, 121-154, |
| 475 | https://doi.org/10.1016/j.iotcps.2023.04.003, 2023. |
| 476 | Ren, C., Zhang, X., Reis, S., Wang, S., Jin, J., Xu, J., and Gu, B.: Climate change unequally affects |
| 477 | nitrogen use and losses in global croplands, Nature Food, 4, 294-304, |
| 478 | https://doi.org/10.1038/s43016-023-00730-z, 2023. |
| 479 | Rogelj, J., Huppmann, D., Krey, V., Riahi, K., Clarke, L., Gidden, M., Nicholls, Z., and Meinshausen, |
| 480 | M.: A new scenario logic for the Paris Agreement long-term temperature goal, Nature, 573, 357- |
| 481 | 363, https://doi.org/10.1038/s41586-019-1541-4, 2019. |
| 482 | Rolnick, D., Donti, P. L., Kaack, L. H., Kochanski, K., Lacoste, A., Sankaran, K., Ross, A. S., Milojevic- |
| 483 | Dupont, N., Jaques, N., Waldman-Brown, A., Luccioni, A. S., Maharaj, T., Sherwin, E. D., |
| 484 | Mukkavilli, S. K., Kording, K. P., Gomes, C. P., Ng, A. Y., Hassabis, D., Platt, J. C., Creutzig, F., |
| 485 | Chayes, J., and Bengio, Y.: Tackling Climate Change with Machine Learning, ACM Comput. Surv., |
| 486 | 55, 42, https://doi.org/10.1145/3485128, 2022. |
| 487 | Taddeo, M., Tsamados, A., Cowls, J., and Floridi, L.: Artificial intelligence and the climate emergency: |
| 488 | Opportunities, challenges, and recommendations, One Earth, 4, 776-779, |
| 489 | https://doi.org/10.1016/j.oneear.2021.05.018, 2021. |
| 490 | Toetzke, M., Banholzer, N., and Feuerriegel, S.: Monitoring global development aid with machine |
| 491 | learning, Nature Sustainability, 5, 533-541, https://doi.org/10.1038/s41893-022-00874-z, 2022. |
| 492 | Tollefson, J. and Weiss, K. R.: Nations approve historic global climate accord, Nature, 528, 315-316, |
| 493 | https://doi.org/10.1038/528315a, 2015. |
| 494 | Trutnevyte, E., McDowall, W., Tomei, J., and Keppo, I.: Energy scenario choices: Insights from a |
| 495 | retrospective review of UK energy futures, Renewable and Sustainable Energy Reviews, 55, 326- |
| 496 | 337, https://doi.org/10.1016/j.rser.2015.10.067, 2016. |
| 497 | Utsumi, N. and Kim, H.: Observed influence of anthropogenic climate change on tropical cyclone heavy |
| 498 | rainfall, Nature Climate Change, 12, 436-440, https://doi.org/10.1038/s41558-022-01344-2, 2022. |
| 499 | van Ruijven, B. J., De Cian, E., and Sue Wing, I.: Amplification of future energy demand growth due to |
| 500 | climate change, Nature Communications, 10, 2762, 10.1038/s41467-019-10399-3, 2019. |
| 501 | van Vuuren, D. P., van der Wijst, KI., Marsman, S., van den Berg, M., Hof, A. F., and Jones, C. D.: The |
| 502 | costs of achieving climate targets and the sources of uncertainty, Nature Climate Change, 10, 329- |
| 503 | 334, https://doi.org/10.1038/s41558-020-0732-1, 2020. |
| 50/ | Weber C. McCollum D. L. Edmonds, L. Forio, P. Dyanet A. Pogeli, I. Tayoni, M. Thoma, J. and |

Weber, C., McCollum, D. L., Edmonds, J., Faria, P., Pyanet, A., Rogelj, J., Tavoni, M., Thoma, J., and
 Kriegler, E.: Mitigation scenarios must cater to new users, Nature Climate Change, 8, 845-848,





- 506 https://doi.org/10.1038/s41558-018-0293-8, 2018.
- Wei, Y.-M., Liu, L.-J., Zhang, S.-X., and Liu, J.-X.: Global Climate Scenario Reference Dataset, [Data set]. Zenodo. https://doi.org/10.5281/zenodo.15536298, 2025.
- Wei, Y. M., Liang, Q. M., Yu, B. Y., and Liao, H. 2023. Climate Change Integrated Assessment Model
 and Its Applications, Science Press.
- Wu, L., Huang, Z., Zhang, X., and Wang, Y.: Harmonizing existing climate change mitigation policy
 datasets with a hybrid machine learning approach, Scientific Data, 11, 580,
 https://doi.org/10.1038/s41597-024-03411-z, 2024.
- 514 Zabala, A.: The right to a sound environment, Nature Sustainability, 4, 750-751, 515 https://doi.org/10.1038/s41893-021-00733-3, 2021.

516 Acknowledgements

- 517 We gratefully acknowledge the financial support of the National Natural Science Foundation of
- 518 China [grant number 72293600, 72293605, 72488101, 72474024 and 72474021]. We thank our
- 519 colleagues for their support and acknowledge help from CEEP-BIT and Information Technology
- 520 Center of BIT.

521 Author contributions

- 522 Y.-M. W., L.-J. L., Q.-M. L. conceptualized the paper. S.-X. Z., J.-X. L., L. Z., Y.-M. C., Y.-X. H.,
- 523 Z.-Q. X., H.-B. C., Y.-X. X., P. W., S.-Y. Y., X.-L. H., T.-Y. W., X.-Q. L., H.-R. X., W.-C. Z., Z.-Q.
- 524 L., and R. C. acquired the data. J.-X. L. and L. Z. carried out the model. L.-J. L., S.-X. Z., J.-X. L.,
- 525 Q.-M. L., and R.-G. C. contributed to the interpretation of the results. S.-X. Z., J.-X. L., Y.-M. C.,
- 526 Y.-X. H. and P. W. implemented the visualization. J.-X. L. and S.-X. Z. contributed to the
- 527 Supplementary Information. Y.-M. W., L.-J. L., S.-X. Z., J.-X. L., Q.-M. L., R.-G. C., T. H., X.-C.
- 528 Y., B. Y., B. T., L.-C. L. and H. L. wrote and reviewed the main manuscript.

529 Competing interests

- 530 The authors declare no competing interests.
- 531