# Journal Pre-proof

A new ChatGPT-empowered, easy-to-use machine learning paradigm for environmental science

Haoyuan An, Xiangyu Li, Yuming Huang, Weichao Wang, Yuehan Wu, Lin Liu, Weibo Ling, Wei Li, Hanzhu Zhao, Dawei Lu, Qian Liu, Guibin Jiang

Please cite this article as: H. An, X. Li, Y. Huang, W. Wang, Y. Wu, L. Liu, W. Ling, W. Li, H. Zhao, D. Lu, Q. Liu, G. Jiang, A new ChatGPT-empowered, easy-to-use machine learning paradigm for environmental science, *Eco-Environment & Health*, https://doi.org/10.1016/j.eehl.2024.01.006.

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.

1    PERSPECTIVE

2    **A new ChatGPT-empowered, easy-to-use machine learning**

3    **paradigm for environmental science**

4    Haoyuan An[a,b], Xiangyu Li[a], Yuming Huang[a], Weichao Wang[a], Yuehan Wu[a], Lin Liu[a], Weibo

5    Ling[a], Wei Li[b], Hanzhu Zhao[b], Dawei Lu*[a], Qian Liu[a], Guibin Jiang[a]

6    [a] State Key Laboratory of Environmental Chemistry and Toxicology, Research Center for Eco-

7    Environmental Sciences, Chinese Academy of Sciences, Beijing 100085, China

8    [b] Biomedical Engineering Institute, School of Control Science and Engineering, Shandong University,

9    Jinan 250061, China

10   * Corresponding author. Email: dwlu@rcees.ac.cn

11

## Abstract

The quantity and complexity of environmental data show exponential growth in recent years. High-quality big data analysis is critical for performing a sophisticated characterization of the complex network of environmental pollution. Machine learning (ML) has been employed as a powerful tool for decoupling the complexities of environmental big data based on its remarkable fitting ability. Yet, due to the knowledge gap across different subjects, ML concepts and algorithms have not been well-popularized among researchers in environmental sustainability. In this context, we introduce a new research paradigm—"ChatGPT + ML + Environment", providing an unprecedented chance for environmental researchers to reduce the difficulty of using ML models. For instance, each step involved in applying ML models to environmental sustainability, including data preparation, model selection and construction, model training and evaluation, and hyper-parameter optimization, can be easily performed with guidance from ChatGPT. We also discuss the challenges and limitations of using this research paradigm in the field of environmental sustainability. Furthermore, we highlight the importance of "secondary training" for future application of "ChatGPT + ML + Environment".
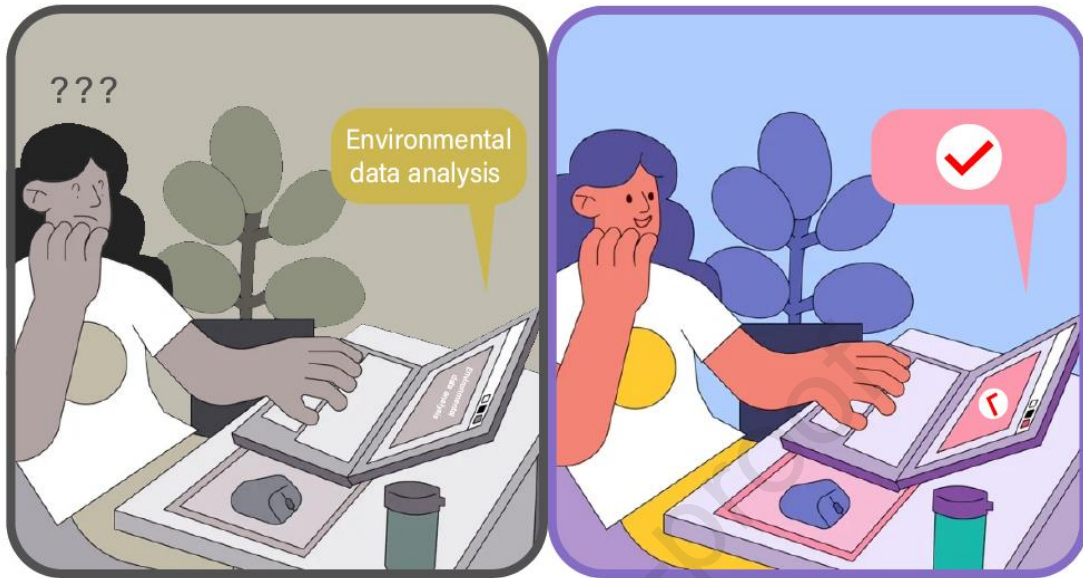
## Keywords

31  **Table of Contents (TOC) Graphic**



"Machine
learning+Environment"                    "ChatGPT+Machine
                                          learning+Environment"

32
33

## 1. Introduction

An environmental issue usually involves multiple substances, factors, and processes, leading to the generation of environmental big data generally characterized by rich sets of input features, e.g., the data of real-time monitoring[1, 2], human activities[3-6], meteorological parameters[7-10], emission inventories[11-14], chemical composition[15, 16], environmental transportation[17, 18], and pollution exposure[19, 20]. In addition to numbers, the input formats of environmental data also include texts, graphs, and images[21]. Hence, environmental big data analysis requires more advanced approaches and powerful tools. In recent years, machine learning (ML), an emerging data mining tool for addressing the multi-dimensional/variety data[22], has triggered a revolutionary development in the field of environmental science[8, 21, 23-28]. ML is defined as "developing a model based on a set of example data, known as 'training data', to generate predictions or decisions without the need for explicit programming"[29]. ML algorithms show an excellent capacity for handling data with various input features and formats, outperforming traditional statistical tools that are often limited to data showing linear relationships with the outcomes[30-32]. It is worth noting that the dataset to be processed can be directly packaged and input into an ML model without prior knowledge of relevant features, and their patterns or trends can be identified or predicted.

In recent years, several reviews have summarized the current state of ML applications in environmental research. In 2021, Zhong et al. reported the working principles of ML algorithms and presented their specific applications in environmental pollution research, including predicting the pollution trends of atmospheric fine particulate matter ($PM_{2.5}$), predicting the future water availability, data processing from different water facilities, predicting sludge bulking in wastewater treatment plants, and identifying the Endocrine Disrupting Chemicals (EDCs)[21]. In 2022, Liu et al. summarized the new gains in using ML algorithms to study environmental issues, and highlighted their applications in estimating the health outcome of exposure[22]. Furthermore, they illustrated the importance of balancing the performance and interpretability of ML models in environmental research. Since 2022, the environmental scenarios of applying ML algorithms have been further expanded. For instance, ML algorithms have been widely used for improving the efficiency of environmental monitoring and policy-

63    making[27], accounting carbon budget[33, 34], decoupling the meteorological impact on air

64    pollution[9, 35], screening the new pollutants from a tremendous number of chemicals[36],

65    predicting the health benefits through reducing pollution[37-42], identifying the impactors

66    affecting the food chain or ecosystem[43, 44], etc. Example ML algorithms used in

67    environmental research include recurrent neural network (RNN)[45], convolutional neural

68    network (CNN)[46], decision tree[47], support vector machine (SVM)[48, 49], random forest

69    (RF)[8, 10], and artificial/deep neural network[22]. Most of these ML models used in

70    environmental research are well-developed, and their concepts, principles, and example codes

71    are publicly shared. Despite that, environmental researchers with less experience in AI

72    techniques still face challenges in appropriate applications of ML algorithms, e.g., misuse of

73    cross-validation to the entire data set[21], or confusion between the validation set and test

74    set[50]. Hence, they usually seek collaborations with researchers in the field of computing,

75    ensuring a correct application of ML algorithms. Yet, some critical parameters for proper ML

76    application, e.g., feature description and hyper-parameter tuning, should be drawn upon

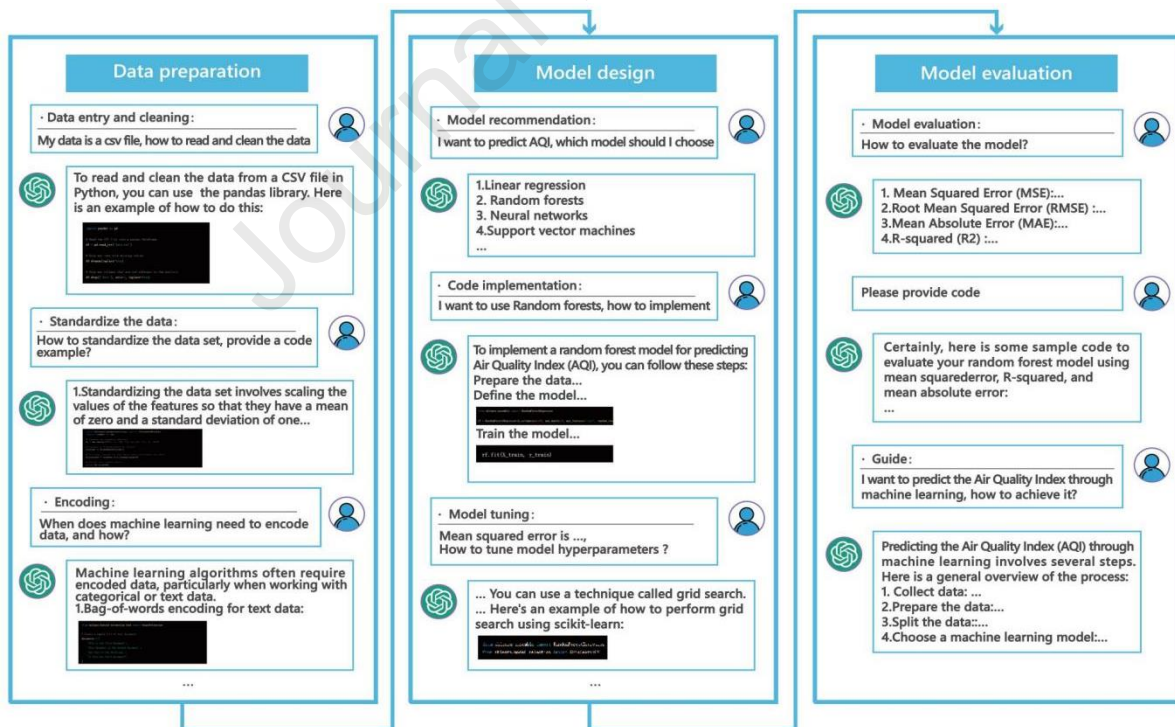77    domain expertise, rather than only AI techniques[21].

78    ChatGPT, as a state-of-the-art version of the dialogue-based model, was launched in

79    November 2022 and will probably simplify ML usage in environmental research[51].

80    Specifically, ChatGPT has been trained on a large corpus of billions of text data, and is

81    embedded with human feedback reinforcement learning and manually supervised fine-

82    tuning[52-55]. This enables it to naturally understand and generate the text like a human[56].

83    Moreover, the human-like text ability makes it an indispensable tool for handling a variety of

84    language-based tasks, e.g., providing exampled codes of ML models and connecting up-/down-

85    stream sections in the full-chain study mentioned above. Thus, for environmental researchers

86    with less knowledge of ML algorithms, ChatGPT might reduce the threshold of using ML for

87    environmental big data analysis.

88    Here, we present a novel research paradigm—"ChatGPT + ML + Environment" and

89    highlight its potential in popularizing ML in the field of environmental science. We also discuss

90    the challenges and limitations remaining in this technique. Considering the current version of

91    ChatGPT-3.5 is mainly performed based on a general database, we give our perspectives on its

92    performance improvement by "secondary training" with some professional databases.

93 Furthermore, we also discuss the possibility of coupling ChatGPT with other AI techniques,

94 e.g., intelligent robots and console algorithms. This training provides a chance for generating

95 an integration solution in the full-chain study of environmental sustainability.

## 2. A new paradigm of "ChatGPT + ML + Environment"

97 The workflow of ML models used in environmental research can generally be

98 decomposed into data preparation, model selection and construction, model training and

99 evaluation, hyper-parameter optimization, and output[57]. Note: hyper-parameter optimization

100 means improving the performance and accuracy of the model by adjusting the hyper-

101 parameters (parameters that cannot be learned by the model itself and require to be manually

102 set) in the algorithm[57]. As shown in Fig. 1 and Supplementary discussion, the specific

103 concepts, common errors, features, and example codes of solutions can be obtained by

104 consulting ChatGPT. Therefore, the paradigm of "ChatGPT+ ML + Environment" is a

105 promising tool that provides an unprecedented chance for inexperienced environmental

106 researchers to address complex data analysis.



107

108 **Figure 1.** Schematic overview of "ChatGPT + ML + Environment". The workflow of using

109 ML in environmental research can be roughly decomposed into data preparation, model design,

110 and model evaluation. The dialog boxes show examples of how ChatGPT makes ML

111     algorithms to be easy used in environmental research.

112

113     **2.1 Data preparation**

114       The raw data of environmental analysis and monitoring usually contain a large amount of

115     "noise" and irrelevant information, as well as incorrect, missing, or duplicate results. Moreover,

116     some types of environmental data cannot be read by the ML model. Although some data can

117     be directly inputted into the model, their uneven distribution also leads to unstable model

118     training and slow model convergence. Therefore, to ensure the smooth running of ML models

119     in environmental research, the first step is to perform data preparation of environmental big

120     data by using some algorithms, e.g., Python's Pandas library and Scikit-Learn library[57].

121     Specifically, we can inquire with ChatGPT about the data preparation methods and their

122     functions, and choose an appropriate one according to the specific formats and features of raw

123     data (Fig. 1). Alternatively, we can also enter ChatGPT with our available data storage formats,

124     and then guide it to provide appropriate data preparation methods (Fig. 1). Furthermore,

125     ChatGPT can also generate the code examples for operating data preparation.

126       To further test the reliability of this method, we performed an example procedure of data

127     preparation in Air Quality Index (AQI) prediction[58]. Specifically, we inputted "My data is a

128     csv file, the columns are 'date, $PM_{2.5}$, $PM_{10}$, $SO_2$, CO, $NO_2$, $O_3$, AQI', the date column does not

129     need to be entered into the model, the remaining columns may be partially missing, how to

130     read the file, perform data cleaning and divide it into a training set and a validation set?" into

131     the ChatGPT. As shown in Supplementary discussion, ChatGPT directly provided annotated

132     codes and their description. However, ChatGPT seemed to ignore that "the date column does

133     not need to be entered into the model." Then, a further instruction, "I don't need the data in the

134     date column," was entered into the ChatGPT, which provided a complete set of code and

135     explanation. Hence, ChatGPT can help inexperienced environmental researchers achieve data

136     preparation of complex environmental data.

137     **2.2 Model selection and construction**

138       As aforementioned, ML models have been widely used for environmental big data

139     analysis, including classification, data fitting, clustering analysis, association analysis, and

140     anomaly detection[21]. Theoretically, there are multiple ML models available that can be used

141  to resolve the same type of task in data analysis. Yet, the model capacity, training speed, and

142  functional focus of these ML models are different. Thus, a sophisticated analysis of the

143  fundamentals and functional differences of the numerous models is essential for model

144  selection. ChatGPT provides an effective solution for selecting an appropriate ML model.

145  Specifically, we can learn about the patterns, basics and fundamentals, functional focuses,

146  advantages, and disadvantages of the intentional-required models by inquiring with ChatGPT.

147  It is worth noting that using ChatGPT to select an ML model only requires a few short

148  conversations, saving considerable time compared with manual research and investigation.

149  Considering that different ML models have their own frameworks, the data to be

150  processed should be optimized to achieve the requirements of the selected ML's framework.

151  For example, if a convolutional neural network (CNN) is chosen to perform AQI prediction

152  (Supplementary discussion), bootstrap instructions can be given to ChatGPT, such as "I want

153  to achieve AQI prediction through a one-dimensional convolutional neural network based on

154  the pytorch framework". Then, ChatGPT would present guidelines for converting the pending

155  data into a readable format for Data Loader. Moreover, a complete set of "sample code" for the

156  selected model construction can also be provided by ChatGPT (Supplementary discussion).

157  After a slight optimization, we can easily build the selected ML model. Hyper-parameters

158  selection, an important factor for proper model building, directly affects the capacity,

159  convergence speed, and performance of the ML model. Particularly, some hyper-parameters

160  (e.g., the depth of trees in the RF model) are not fixed options, which should be set with a

161  comprehensive account of the number of input data features, data volume, data distribution,

162  and application scenario, etc[21]. Considering that hyper-parameters selection is a dilemma

163  that involves the knowledge of AI and environmental science, inexperienced environmental

164  researchers can seek solutions with the support of ChatGPT. Although ChatGPT might not

165  provide optimum parameter settings, it can provide the detailed meaning of each hyper-

166  parameter and advanced methods (e.g., grid search) for proper selection. Thus, ChatGPT can

167  guide the ML model building in the field of environmental science.

168  To illustrate how to select the most appropriate ML mode, we performed an exampled

169  case of the Shannon index (a critical indicator for measuring biodiversity) prediction with the

170  parameters of nanoparticles (e.g., type, shape, size, potential) and relevant environmental

171 factors (e.g., temperature, pH, soil depth). For instance, we performed an original prediction

172 with linear regression based on this ChatGPT-empowered system. Then, "Can any other model

173 be used to achieve this prediction? Output the performance of each model and select the best

174 one." was inputted into the ChatGPT-empowered system. As shown in Supplementary

175 discussion, the ChatGPT-empowered system provided the codes of linear regression, random

176 forest, and xgb tree models, and output the name and RMSE (Root Mean Square Error) of the

177 most suitable model. Moreover, the ChatGPT-empowered system can provide codes of cross-

178 validation to evaluate the performance of these models. It can also search the most suitable

179 parameters on the internet automatically. For the whole process, we merely provided the output

180 and error message from the last step for ChatGPT, which then generated the subsequent codes

181 of correction and implementation automatically.

182 **2.3 Model training, performance, and hyper-parameter optimization**

183     ChatGPT can further guide the training, performance evaluation, and hyper-parameter

184 optimization of the ML models used in environmental research. For traditional ML models like

185 RF and SVM, most of their codes used for model training are with fixed structures[21, 22]. The

186 corresponding statements and structures can usually be found by ChatGPT in the database of

187 code examples. For instance, the training procedure of the RF model for air quality (AQI)

188 prediction from emissions was smoothly performed with guidance from ChatGPT

189 (Supplementary discussion). With regard to deep learning models, to reduce running problems

190 (e.g., convergence difficulties and declining model generalization ability), the parameters,

191 including learning rate, optimizer, and learning rate decay, are required to be set prior to

192 training[22]. Taking an example of AQI prediction by using CNN (Supplementary discussion),

193 the parameters including adam optimizer, learning rate (0.001), and mean squared error loss

194 were successfully set guided by ChatGPT. Moreover, to further optimize the training process,

195 the procedures of gradient descent and backpropagation, and the codes for learning rate decay

196 were also provided by ChatGPT.

197     Model performance is critical for ML applications, determining the reliability of

198 prediction[57]. Although there are many ways to evaluate an ML model's performance, some

199 evaluation parameters involve computer terminology and are difficult to understand for

200 environmental researchers. ChatGPT can provide formulas, meanings, and examples of

201  application scenarios of the various evaluation parameters for users to understand and select

202  appropriate evaluation methods. Specifically, we can obtain the "Mean Squared Error," "Root

203  Mean Squared Error," "Mean Absolute Error," and "R-squared" of the models used in AQI

204  predictions via inquiring with ChatGPT (Fig.1, Supplementary discussion). More importantly,

205  the implementation codes for model evaluation can be accessed directly from the package

206  provided by ChatGPT. Furthermore, tuning hyper-parameters is usually required to further

207  improve the model performance. Similar to hyper-parameters selection (*see section 2.2*), we

208  can obtain specific tuning codes of the selected model, and find the optimum hyper-parameters

209  by ChatGPT.

210      The aforementioned applications mainly tend to directly use or make slight modifications

211  to the existing code structures. In these applications, ChatGPT can provide clear and concise

212  code examples, preventing us from spending tremendous time studying the user manual of

213  various ML models. This is of extreme importance for those with less knowledge in ML

214  programming, as it can greatly reduce the interference and misdirection caused by complex

215  codes. Additionally, ChatGPT can provide code interpretation and error-checking assistance,

216  enabling us to quickly grasp the logical framework of a code segment and apply it to

217  environmental studies. To facilitate understanding, the whole process of application examples

218  based on the paradigm of "ChatGPT + ML + Environment" has been successfully performed,

219  as detailed in Supplementary discussion.

## 3 Advancement and challenges

221      In addition to the aforementioned text data processing, the ChatGPT-empowered system

222  also shows advantages in processing complex data. For instance, it can be used to predict the

223  toxicology of chemicals based on their physical-chemical properties dataset (see

224  Supplementary discussion). The used dataset consists of 210 features, including a series of

225  specific chemical descriptors (e.g., molecular structure, chemical name, source, and CAS

226  number), a range of refined molecular properties (e.g., polar surface area, adsorption properties,

227  the quantity, state, and size of atoms and functional groups), and some important

228  physicochemical properties (e.g., solubility, lipophilicity, and surface area). Considering that

229  the dataset is a mixture of both useful and irrelevant information, including numerical and
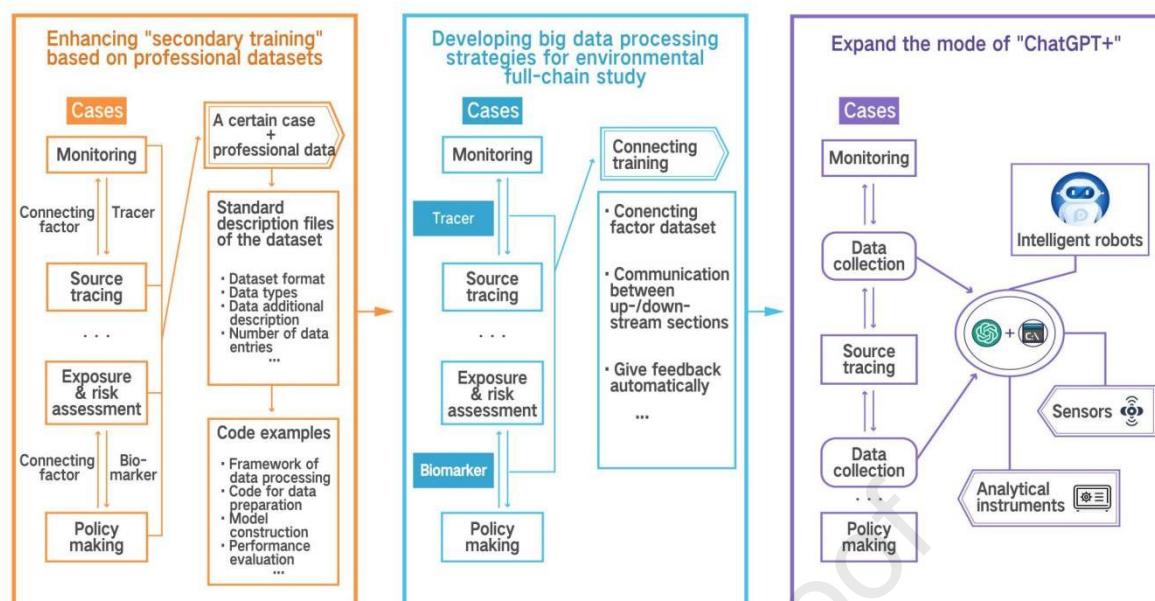
230  character-based data, we initially used the ChatGPT-3.5 to generate the code of a random forest

231  model, yielding an RMSE of 1.39. To address the possible limitations of ChatGPT-3.5 missing

232  some contextual information in complex datasets, we further performed this prediction by using

233  the ChatGPT4.0-empowered system. As shown in Supplementary discussion, the RMSE is

234  0.67 with an R-squared ($R^2$) of 0.57, which demonstrates the potential of the ChatGPT-

235  empowered system in addressing complex ML tasks.

236    However, ChatGPT, one of the first human-like language models, still faces challenges

237  and limitations in environmental applications. For instance, 1) Honest use. Most of ChatGPT's

238  output is difficult to distinguish from the text written by humans. Recently, ChatGPT was

239  directly listed as the author of several publications, which has triggered a widespread

240  discussion among the academic community[53-55]. Indeed, the use of ChatGPT must strictly

241  adhere to academic ethics and standards. To popularize the applications of public-shared tools

242  (i.e., ML) in the field of environmental science, the details of ChatGPT usage should be clearly

243  disclosed in the publications. Furthermore, for better regulation, the usage record can be

244  documented accurately with the time stamp in blockchain technique. 2) Model development.

245  The training of ChatGPT is still based on a large amount of existing data. Therefore, ChatGPT

246  can provide code examples for the well-developed ML models used in environmental research

247  but fails to develop new models. As shown in Supplementary discussion, the ChatGPT-

248  empowered system can perform almost all ML tasks in environmental science. Yet, it is still a

249  probability-based AI model[51]. Its responses are the results of analyzing a large amount of

250  training data, lacking thought of the context and background information. Therefore, it may not

251  understand why we perform these analyses, and hence the whole data processing strategy

252  should be designed by the researchers. Moreover, ChatGPT would be unaware of the parameter

253  errors existing in its generated codes, which can only be found when the codes are actually

254  executed. 3) Professional database. The current ChatGPT database is limited to general data

255  prior to 2021[51, 53], lacking a professional dataset of environmental sustainability. This may

256  result in suboptimal performance in solving environmental problems. Therefore, the ChatGPT-

257  empowered plug-in can be embedded into the professional system of environmental research

258  to promptly provide ML applications. Additionally, to obtain high-quality big data analysis,

259  some environmental data are encouraged to be open to the public.

## 4. Discussion

Although ML is a powerful tool for addressing complex environmental problems, it can be a challenging task for environmental scientists without AI research backgrounds. Integrating ChatGPT can provide effective solutions, including the concepts, principles and exampled codes, for ML applications. For environmental researchers with no prior knowledge, it can help them to perform ML analysis smoothly; for scientists with some AI knowledge, this process will improve their efficiency by saving their time to edit the codes. Notably, almost all programming tools or languages like Python and R can be used to build the ChatGPT-based process. In addition to environmental science, this process will extend ML application to other fields, e.g., industrial, biology, and geochemistry. Furthermore, it is noted that other Generative Pre-trained Transformer-based tools like Claude and Bard have similar effects as the ChatGPT[51], reducing the threshold of environmental application of ML. With the development of generative models and AI technologies, the application of the "ChatGPT + ML + Environment" research paradigm will be further expanded. For instance, the processed data will not be limited to text, and graphic data might be understood and processed as the ChatGPT evolves [53]. In the future, these techniques, used correctly in accordance with academic ethics and usage guidelines, would provide excitement for solving complex environmental problems:

1) Enhancing "secondary training" based on professional datasets. As shown in Fig. 2, the first step involves choosing a certain type of environmental case (e.g., environmental monitoring, source tracing, and policy making) and introducing a specific professional dataset. Moreover, a standard description file of the professional dataset, including dataset format, data types, additional data description, number of data entries, and dataset content description, should be set for the system of "ChatGPT + ML + Environment." This step will help ChatGPT to learn about the overview of the dataset. Afterward, a "secondary training" model, including the framework of data processing, the code for data preparation, model construction, and performance evaluation, would be built for the professional dataset. The detailed implementation procedures are similar to that mentioned in *Section 2*. Through further training or optimization, the "secondary training" model would show a capacity to provide effective and quick solutions for such environmental problems, especially for some emergency events.

**Figure 2.** The conceptual mode of "ChatGPT + ML + Environment" in future environmental research. The left box shows the secondary training by introducing environmental professional dataset. The middle box mainly shows the potential in connecting the up-/down-stream tasks of data analysis in the full chain study of environmental sustainability. The right box mainly gives a perspective on coupling data processing with data collection via using an integration of ChatGPT, control algorithms, ML, and robots, etc.

2) Developing big data processing strategies for full-chain environmental study. An environmental event usually involves the coupling of multiple substances, factors, and processes across various scales, requiring a comprehensive research route covering "monitoring—source tracing—environmental behavior and transformation—exposure and risk assessment—policy making." Each of them can generate different datasets (Fig. 2). These datasets might have become "data islands" due to a lack of proper data analysis techniques, hampering the proposal of a systematic solution for real environmental problems[22]. Identifying the connection factors and developing an intelligent data processing system is critical for achieving full-chain environmental study. For instance, we would first establish a dataset composed of connection factors (Fig. 2), e.g., tracers, transformation reactions, biomarkers, and policy implementation date. The specific communication instructions for connecting up-/down-stream sections would be well-trained by ChatGPT with its human-like text ability[54]. In this way, the ML-based data processing in a down-stream section can be

310 operated automatically after receiving the output from the up-stream section. Alternatively,

311 they can provide feedback of the output to the up-stream section, guiding its optimization. Thus,

312 the integration of ChatGPT and ML algorithms is a promising tool for future full-chain

313 environmental research.

314      3) Expanding the application mode of "ChatGPT +". The integration of ChatGPT and ML

315 significantly improves the processing capacity of environmental big data, promoting the rapid

316 development of environmental science. For instance, the current environmental monitoring

317 system is capable of continuously collecting real-time environmental data and outputting brief

318 reports[48, 58]. Such operations are tasks consisting of specific sequences of steps, where the

319 execution of each task is based on previously normalized instructions. However, these tasks

320 pose challenges in terms of generating predictions, making decision, and developing smart

321 feedback to optimize the next step of data collection. In the future, the "ChatGPT + ML" mode

322 can be further expanded by combining with other intelligent techniques like intelligent robots

323 and control algorithms. Specifically, multiple environmental data collection devices (e.g.,

324 intelligent robots, sensors, and analytical instruments) and their carriers would be connected

325 by the "ChatGPT + ML" system integrated with computer control algorithms (Fig. 2). This will

326 integrate static environmental big data processing with dynamic environmental analysis,

327 providing a novel tool for future environmental research, especially for some environmental

328 monitoring under extreme conditions.

## Declaration of competing interests

330 The authors declare no competing interests.

331

## Acknowledgments

## References

[1] G. Geng, Q. Xiao, S. Liu, X. Liu, J. Cheng, Y. Zheng, et al., Tracking air pollution in China: Near real-time $PM_{2.5}$ retrievals from multisource data fusion, Environ. Sci. Technol. 55 (2021) 12106.

[2] H. Messer, A. Zinevich, P. Alpert, Environmental monitoring by wireless communication networks, Science 312 (2006) 713.

[3] D.G. Streets, H.M. Horowitz, D.J. Jacob, Z. Lu, L. Levin, A.F.H. ter Schure, et al., Total mercury released to the environment by human activities, Environ. Sci. Technol. 51 (2017) 5969.

[4] A. Pruden, M. Arabi, H.N. Storteboom, Correlation between upstream human activities and riverine antibiotic resistance genes, Environ. Sci. Technol. 46 (2012) 11541.

[5] D.G. Streets, M.K. Devane, Z. Lu, T.C. Bond, E.M. Sunderland, D.J. Jacob, All-time releases of mercury to the atmosphere from human activities, Environ. Sci. Technol. 45 (2011) 10485.

[6] A.R. Ferro, R.J. Kopperud, L.M. Hildemann, Source strengths for indoor human activities that resuspend particulate matter, Environ. Sci. Technol. 38 (2004) 1759.

[7] J. Klánová, P. Èupr, J. Kohoutek, T. Harner, Assessing the influence of meteorological parameters on the performance of polyurethane foam-based passive air samplers, Environ. Sci. Technol. 42 (2008) 550.

[8] Z. Shi, C. Song, B. Liu, G. Lu, J. Xu, T. Van Vu, et al., Abrupt but smaller than expected changes in surface air quality attributable to COVID-19 lockdowns, Sci. Adv. 7 (2021) eabd6696.

[9] P. Zuo, Y. Huang, P. Liu, J. Zhang, H. Yang, L. Liu, et al., Stable iron isotopic signature reveals multiple sources of magnetic particulate matter in the 2021 Beijing sandstorms, Environ. Sci. Technol. Lett. 9 (2022) 299.

[10] P. Zuo, Z. Zong, B. Zheng, J. Bi, Q. Zhang, W. Li, et al., New insights into unexpected severe $PM_{2.5}$ pollution during the SARS and COVID-19 pandemic periods in Beijing, Environ. Sci. Technol. 56 (2022) 155.

[11] J. Hao, H. Tian, Y. Lu, Emission inventories of NOx from commercial energy consumption in China, 1995−1998, Environ. Sci. Technol. 36 (2002) 552.

[12] K. Breivik, V. Vestreng, O. Rozovskaya, J.M. Pacyna, Atmospheric emissions of some POPs in Europe: a discussion of existing inventories and data needs, Environ. Sci. Policy 9 (2006) 663.

[13] X. Lu, X. Ye, M. Zhou, Y. Zhao, H. Weng, H. Kong, et al., The underappreciated role of agricultural soil nitrogen oxide emissions in ozone pollution regulation in North China, Nat. Commun. 12 (2021) 5021.

[14] M. Li, H. Liu, G. Geng, C. Hong, F. Liu, Y. Song, et al., Anthropogenic emission inventories in China: a review, Natl. Sci. Rev. 4 (2017) 834.

[15] X. Feng, H. Sun, X. Liu, B. Zhu, W. Liang, T. Ruan, et al., Occurrence and ecological impact of

370    chemical mixtures in a semiclosed sea by suspect screening analysis, Environ. Sci. Technol. 56 (2022) 10681.

371    [16] S. Breinlinger, T.J. Phillips, B.N. Haram, J. Mareš, J.A. Martínez Yerena, P. Hrouzek, et al., Hunting

372    the eagle killer: A cyanobacterial neurotoxin causes vacuolar myelinopathy, Science 371 (2021) eaax9050.

373    [17] Z. Tian, H. Zhao, K.T. Peter, M. Gonzalez, J. Wetzel, C. Wu, et al., A ubiquitous tire rubber-derived

374    chemical induces acute mortality in coho salmon, Science 371 (2021) 185.

375    [18] Y. Yin, Y. Li, C. Tai, Y. Cai, G. Jiang, Fumigant methyl iodide can methylate inorganic mercury species

376    in natural waters, Nat. Commun. 5 (2014) 4633.

377    [19] D. Lu, Q. Luo, R. Chen, Y. Zhuansun, J. Jiang, W. Wang, et al., Chemical multi-fingerprinting of

378    exogenous ultrafine particles in human serum and pleural effusion, Nat. Commun. 11 (2020) 2567.

379    [20] R. Vermeulen, E.L. Schymanski, A.-L. Barabási, G.W. Miller, The exposome and health: Where

380    chemistry meets biology, Science 367 (2020) 392.

381    [21] S. Zhong, K. Zhang, M. Bagheri, J.G. Burken, A. Gu, B. Li, et al., Machine learning: New ideas and

382    tools in environmental science and engineering, Environ. Sci. Technol. 55 (2021) 12741.

383    [22] X. Liu, D. Lu, A. Zhang, Q. Liu, G. Jiang, Data-driven machine learning in environmental pollution:

384    Gains and problems, Environ. Sci. Technol. 56 (2022) 2124.

385    [23] Z. Cao, J. Zhou, M. Li, J. Huang, D. Dou, Urbanites' mental health undermined by air pollution, Nat.

386    Sustain. (2023).

387    [24] W. Li, W.-Y. Guo, M. Pasgaard, Z. Niu, L. Wang, F. Chen, et al., Human fingerprint on structural density

388    of forests globally, Nat. Sustain. (2023).

389    [25] M. Toetzke, N. Banholzer, S. Feuerriegel, Monitoring global development aid with machine learning,

390    Nat. Sustain. 5 (2022) 533.

391    [26] Z. Mehrabi, M.J. McDowell, V. Ricciardi, C. Levers, J.D. Martinez, N. Mehrabi, et al., The global

392    divide in data-driven farming, Nat. Sustain. 4 (2021) 154.

393    [27] M. Hino, E. Benami, N. Brooks, Machine learning for environmental monitoring, Nat. Sustain. 1 (2018)

394    583.

395    [28] M. Callaghan, C.-F. Schleussner, S. Nath, Q. Lejeune, T.R. Knutson, M. Reichstein, et al., Machine-

396    learning-based evidence and attribution mapping of 100,000 climate impact studies, Nat. Clim. Change 11

397    (2021) 966.

398    [29] J.R. Koza, F.H. Bennett, D. Andre, M.A. Keane, in: J.S. Gero and F. Sudweeks (Eds.), Automated

399    Design of Both the Topology and Sizing of Analog Electrical Circuits Using Genetic Programming, Springer

400    Netherlands. Dordrecht, 1996, pp. 151.

401    [30] D. Seng, Q. Zhang, X. Zhang, G. Chen, X. Chen, Spatiotemporal prediction of air quality based on

402    LSTM neural network, Alex. Eng. J. 60 (2021) 2021.

403    [31] Y. Zhao, L. Wang, J. Luo, T. Huang, S. Tao, J. Liu, et al., Deep learning prediction of polycyclic

404    aromatic hydrocarbons in the high arctic, Environ. Sci. Technol. 53 (2019) 13238.

405    [32] R. Janarthanan, P. Partheeban, K. Somasundaram, P. Navin Elamparithi, A deep learning approach for

406    prediction of air quality index in a metropolitan city, Sustain. Cities Soc. 67 (2021) 102720.

407    [33] M. Mugabowindekwe, M. Brandt, J. Chave, F. Reiner, D.L. Skole, A. Kariryaa, et al., Nation-wide

408    mapping of tree-level aboveground carbon stocks in Rwanda, Nat. Clim. Change 13 (2023) 91.

409    [34] Z. Ban, X. Hu, J. Li, Tipping points of marine phytoplankton to multiple environmental stressors, Nat.

410    Clim. Change 12 (2022) 1045.

411    [35] Z. Zhang, B. Xu, W. Xu, F. Wang, J. Gao, Y. Li, et al., Machine learning combined with the PMF model

412    reveal the synergistic effects of sources and meteorological factors on $PM_{2.5}$ pollution, Environ. Res. 212

413    (2022) 113322.

[36] D. Xia, J. Chen, Z. Fu, T. Xu, Z. Wang, W. Liu, et al., Potential application of machine-learning-based quantum chemical methods in environmental chemistry, Environ. Sci. Technol. 56 (2022) 2115.

[37] J. Jeong, J. Choi, Artificial intelligence-based toxicity prediction of environmental chemicals: Future directions for chemical management applications, Environ. Sci. Technol. 56 (2022) 7532.

[38] L. Conibear, C.L. Reddington, B.J. Silver, Y. Chen, C. Knote, S.R. Arnold, et al., Sensitivity of air pollution exposure and disease burden to emission changes in China using machine learning emulation, GeoHealth 6 (2022) e2021GH000570.

[39] E. Isaev, B. Ajikeev, U. Shamyrkanov, K.-u. Kalnur, K. Maisalbek, R.C. Sidle, Impact of climate change and air pollution forecasting using machine learning techniques in Bishkek, Aerosol Air Qual. Res. 22 (2022) 210336.

[40] L. Zhang, X. Li, H. Chen, Z. Wu, M. Hu, M. Yao, Haze air pollution health impacts of breath-borne VOCs, Environ. Sci. Technol. 56 (2022) 8541.

[41] G.D. Thurston, L.C. Chen, M. Campen, Particle toxicity's role in air pollution, Science 375 (2022) 506.

[42] H. Tan, J. Wu, R. Zhang, C. Zhang, W. Li, Q. Chen, et al., Development, validation, and application of a human reproductive toxicity prediction model based on adverse outcome pathway, Environ. Sci. Technol. 56 (2022) 12391.

[43] C. Zhan, H. Matsumoto, Y. Liu, M. Wang, Pathways to engineering the phyllosphere microbiome for sustainable crop production, Nat. Food 3 (2022) 997.

[44] H. Meyer, E. Pebesma, Machine learning-based global maps of ecological variables and the challenge of assessing them, Nat. Commun. 13 (2022) 2208.

[45] G. Kurnaz, A.S. Demir, Prediction of $SO_2$ and $PM_{10}$ air pollutants using a deep learning-based recurrent neural network: Case of industrial city Sakarya, Urban Clim. 41 (2022) 101051.

[46] V. Nikolopoulou, R. Aalizadeh, M.-C. Nika, N.S. Thomaidis, TrendProbe: Time profile analysis of emerging contaminants by LC-HRMS non-target screening and deep learning convolutional neural network, J. Hazard. Mater. 428 (2022) 128194.

[47] A. Coors, A.R. Brown, S.K. Maynard, A. Nimrod Perkins, S. Owen, C.R. Tyler, Minimizing experimental testing on ffish for legacy pharmaceuticals, Environ. Sci. Technol. 57 (2023) 1721.

[48] F. Camastra, V. Capone, A. Ciaramella, A. Riccio, A. Staiano, Prediction of environmental missing data time series by Support Vector Machine Regression and Correlation Dimension estimation, Environ. Modell. Softw. 150 (2022) 105343.

[49] X.-C. Song, N. Dreolin, E. Canellas, J. Goshawk, C. Nerin, Prediction of collision cross-section values for extractables and leachables from plastic products, Environ. Sci. Technol. 56 (2022) 9463.

[50] M. Lastra-Mejias, A. Villa-Martinez, M. Izquierdo, R. Aroca-Santos, J.C. Cancilla, J.S. Torrecilla, Combination of LEDs and cognitive modeling to quantify sheep cheese whey in watercourses, Talanta 203 (2019) 290.

[51] ChatGPT: optimizing language models for dialogue, https://openai.com/blog/chatgpt, (2022).

[52] Much to discuss in AI ethics, Nat. Mach. Intell. 4 (2022) 1055.

[53] C. Stokel-Walker, AI bot ChatGPT writes smart essays — should professors worry?, (2022) https://doi.org/10.1038/d41586.

[54] M. Hutson, Could AI help you to write your next paper?, Nature 611 (2022) 192.

[55] J.B. Eva A. M. van Dis, Willem Zuidema, Robert van Rooij, Claudi L. Bockting, ChatGPT: five priorities for research, Nature 614 (2023) 224.

[56] The AI writing on the wall, Nat. Mach. Intell. 5 (2023) 1.

[57] L. Bottou, Stochastic gradient descent tricks in Neural Networks: Tricks Trade, Berlin,

458    Germany:Springer, 7700 (2012).

459    [58] PM$_{2.5}$ Prediction Based on Random Forest Algorithm., https://github.com/StephenZheng0315/PM2.5-

460    Prediction-Based-on-Random-Forest-Algorithm. (2023).

**Highlights**

- A new paradigm of "ChatGPT + Machine learning (ML) + Environment" is presented.

- The novelty and knowledge gaps of ML for decoupling the complexity of environmental big data are discussed.

- The new paradigm guided by GPT reduces the threshold of using Machine Learning in environmental research.

- The importance of "secondary training" for using "ChatGPT + ML + Environment" in the future is highlighted.