



1 **Machine learning significantly improves the simulation of hourly-to-**
2 **yearly scale cloud nuclei concentration and radiative forcing in**
3 **polluted atmosphere**

4 Jingye Ren^{1,2}, Songjian Zou³, Honghao Xu³, Guiquan Liu³, Zhe Wang³, Anran Zhang³,
5 Chuanfeng Zhao⁴, Min Hu⁵, Dongjie Shang⁵, Lizi Tang⁵, Ru-Jin Huang¹, Yele Sun⁶,
6 Fang Zhang^{3*}

7 ¹State Key Laboratory of Loess Science, Institute of Earth Environment, Chinese Academy of
8 Sciences, Xi'an, 710061, China

9 ²Xi'an Institute for Innovative Earth Environment Research, Xi'an, 710061, China

10 ³Shenzhen Key Laboratory of Organic Pollution Prevention and Control, School of Eco-
11 Environment, Harbin Institute of Technology Shenzhen, Shenzhen, 518055, China

12 ⁴Department of Atmospheric and Oceanic Sciences, School of Physics, Peking University, Beijing,
13 100871, China

14 ⁵State Key Joint Laboratory of Environmental Simulation and Pollution Control, College of
15 Environmental Sciences and Engineering, Peking University, Beijing, 100871, China

16 ⁶State Key Laboratory of Atmospheric Boundary Layer Physics and Atmospheric Chemistry,
17 Institute of Atmospheric Physics, Chinese Academy of Sciences, Beijing, 100029, China

18

19 Corresponding author: Fang Zhang, zhangfang2021@hit.edu.cn

20

21

22

23

24

25



26 **Abstract**

27 The accurate prediction of cloud condensation nuclei (CCN) number
28 concentration (N_{CCN}) on a large spatiotemporal scale is challenging but critical to
29 evaluate the aerosol cloud interaction effect. Combining multi-source dataset and the
30 N_{CCN} simulated by the Weather Research and Forecasting coupled with Chemistry
31 (WRF-Chem) model, we have developed a Random Forest Regression method (RFRM)
32 model which achieves well prediction of hourly-to-yearly scale N_{CCN} at typical
33 supersaturations in polluted North China Plain (NCP). We show that the prediction bias
34 of N_{CCN} compared to observations is reduced from -59% with the WRF-Chem model
35 to approximately -31% with the RFRM model (the prediction precision is improved by
36 1.6 times accordingly) during the campaigns. The greatest improvement is seen in both
37 very polluted and clean cases. The RFRM model captures well the spatial variation and
38 better describes long-term trends of N_{CCN} . More importantly, the prediction reveals a
39 significant long-term decreasing trend of N_{CCN} in NCP due to a rapid reduction in
40 aerosol concentrations from 2014 to 2018, during which a series of strict emission
41 reduction measures were implemented by the Chinese government. This reflects the
42 climate benefit of pollution control. Our study further illustrates that the RFRM model
43 reduces the uncertainty in simulating cloud radiative forcing from an overestimation of
44 $1.89 \pm 0.78 \text{ W m}^{-2}$ to $0.81 \pm 0.63 \text{ W m}^{-2}$, illustrating the high sensitivity of climate forcing
45 to changes in N_{CCN} . This work offers a new modeling framework that guides the way
46 to simulate CCN in other regions around the world and has the potential to effectively
47 filling the observation gap of CCN concentrations.



48 1. Introduction

49 Aerosol indirect radiative effects caused by aerosol-cloud interactions (ACI) are
50 the largest source of uncertainty when assessing climate change (IPCC 2021). A major
51 issue is the lack of an accurate characterization of cloud condensation nuclei (CCN, or
52 cloud nuclei) number concentrations (N_{CCN}) in global climate models (Sotiropoulou et
53 al., 2007; Fanourgakis et al., 2019). This is largely due to the nonlinear interactions
54 between the aerosol physical and chemical properties and CCN, making the
55 quantification of the N_{CCN} remain highly uncertain. Current models tend to
56 underestimate CCN number concentrations by 20–40% on average, based on
57 comparisons between results from 14 models and observations at sites distributed
58 globally (Fanourgakis et al., 2019). Biases were greater in the Northern Hemisphere
59 due to intensive human activities (Sotiropoulou et al., 2007). It has been proposed that
60 ~10–30% changes in cloud droplet concentrations might be associated with the
61 uncertainties in cloud radiative forcing by $\sim 0.1\text{--}2\text{ W m}^{-2}$ (Charlson et al., 1987;
62 Sotiropoulou et al., 2007; Yu et al., 2022). Towards improving estimates of the ACI
63 effect in modeling, it is critical to obtain accurate spatiotemporal distributions of CCN.

64 Predicting N_{CCN} remains challenging because aerosol CCN activity varies greatly
65 in time and space and involves microphysical and chemical processes. Previous studies
66 have underscored that the main uncertainties in simulating N_{CCN} at regional and global
67 scales are the simplified representations of particle size distribution (Menon & Rotstayn,
68 2006), as well as the lack of detailed treatment of the microphysical and chemical



69 processes in current models (Sha et al., 2019; Yu et al., 2020). Therefore, a considerable
70 number of CCN closure studies have been carried out to predict N_{CCN} in different
71 environments (Moore et al., 2012; Zhang et al., 2014, 2016, 2017, 2019; Xu et al., 2021;
72 Ren et al., 2018, 2023). Although extensive CCN observations and closure studies
73 might help to reconcile this uncertainty, field measurements are relatively sparse and
74 have only been collected during a few campaigns at a few sites (Schmale et al., 2018;
75 Rose et al., 2021) due to the limitations of techniques and cost. Closure studies, however,
76 have mostly focused on investigating the relative importance of the aerosol physical
77 and chemical properties to CCN prediction and have not yet provided a CCN
78 parameterization scheme that is applicable to different regions over the globe. Some
79 studies have attempted to develop a correlation between aerosol optical properties and
80 CCN number concentrations (Rosenfeld et al., 2016). Compared with the measurement
81 of CCN at ground sites, satellite retrieval methods offer global coverage with high
82 spatial and temporal resolutions (Rosenfeld et al., 2016; Liu et al., 2020). However,
83 they are limited to clear-sky conditions. Due to the aerosol swelling effect (Liu et al.,
84 2018), there are typically large deviations of -30% to +90% in the estimation accuracy
85 of N_{CCN} in different environments (Shen et al., 2019).

86 In recent years, machine learning (ML) has been used for the inversion of
87 atmospheric environmental parameters such as tropospheric ozone and particulate
88 matter with a diameter of 2.5 μm or smaller ($\text{PM}_{2.5}$) (Grange et al., 2018; Geng et al.,
89 2021; Wei et al., 2023). To our knowledge, ML-based prediction of CCN properties is
90 few and far between, with studies only focused on analyzing the importance of different



91 variables to estimating CCN spectra at a single field site (Liang et al., 2022) or were
92 conducted in relative clean regions (Nair et al., 2020, 2021). It would be a step forward
93 to use an ML-based method to predict N_{CCN} in polluted areas because it would help to
94 verify the applicability of the method in different regions, but most importantly, it
95 would improve model simulations of the ACI effect in polluted regions where errors in
96 predicting N_{CCN} are greater than in clean regions.

97 In this study, we have developed a ML-based model for predicting N_{CCN} on hourly-
98 to-yearly scale in the heavily polluted North China Plain (NCP) by using a multi-source
99 dataset of atmospheric variables and CCN concentration outputs from the Weather
100 Research and Forecasting coupled with online chemistry (WRF-Chem) model. We have
101 presented and analyzed the relative importance of the different parameters to CCN
102 prediction. Moreover, we have verified the performance of the RFRM model in
103 predicting the N_{CCN} over different temporal and spatial scales in the NCP. Finally, by
104 incorporating the N_{CCN} prediction biases into the evaluation of cloud parameters and
105 radiative forcing, we investigate the sensitivity of aerosols indirect climate forcing to
106 CCN concentrations changes.

107 **2. Methods**

108 **2.1 Study area**

109 In this work, we select the North China Plain (NCP) (32°-40°N and 114°-121°E)
110 as the study area. Being one of the most polluted areas in China, the aerosol particles
111 in NCP are with more complex composition and mixing state, which leads to great



112 challenge in accurate prediction of cloud concentration nuclei (CCN) concentrations.
113 In recent years, emissions of gas pollutants and fine particles have shown a significant
114 downward trend year by year (Wei et al., 2023) due to the implementation of the
115 vigorous emission reduction in China (Zheng et al., 2018). This also makes changes in
116 aerosols CCN activity in the study area from the point of view in assessment of the
117 climate effect of aerosols.

118 **2.2 Model construction and validation**

119 Here we develop the ML-based N_{CCN} prediction model by employing the Random
120 Forest Regression method (RFRM) that has been demonstrated and can solve
121 multivariate and nonlinear regression problems (Nair and Yu, 2020; Liang et al., 2022).
122 The diagram of the model construction and the N_{CCN} prediction is shown in Figure 1.
123 The main dataset that are used to construct the ML-based N_{CCN} prediction model were
124 from six field campaigns at three sites in the NCP (more details see section 2.3.1). The
125 observed CCN number concentration (at supersaturations of 0.2% and 0.4%) is as
126 targeted parameter, and the simultaneously measured atmospheric gaseous precursors,
127 fine particles chemical compositions, and meteorological parameters are as model input
128 parameters. Considering the motivation and main goal of this study is to achieve the
129 prediction of CCN number concentration on a relatively large spatiotemporal scale (e.g.
130 regional scale), especially in polluted areas. Therefore, to better constrain the
131 performance and accuracy of the model in regional-scale predictions, we also
132 incorporated the N_{CCN} simulated by WRF-Chem (more details see section 2.3.2) as an
133 input parameter for model training when constructing the model.



134 When constructing the model, the N_{CCN} simulated by WRF-Chem were processed
135 to match the locations of each observation site and their corresponding measurement
136 periods. The data were split into 7:3 ratio for model training and testing respectively. In
137 order to assure a stronger generalization ability of the N_{CCN} prediction model, the 10-
138 fold cross-validation is adopted (Wei et al., 2023). The optimization parameters of RF
139 model were examined by varying hyperparameters. In addition, cross-validation (CV)
140 is applied to select the hyperparameters during the data preprocessing (Yang et al.,
141 2022).

142 Here the quality metrics for model performance are based on the correlation
143 coefficient (R^2), root mean square error (RMSE), the slope of the RFRM predicted with
144 the observed CCN concentrations and the normalized mean bias (NMB). It showed that
145 the estimated N_{CCN} are highly correlated with the observed values for the testing dataset,
146 with the correlation R^2 of $\sim 0.86-0.95$ and slopes of $0.80-0.88$ at $S=0.2\%$ and $S=0.4\%$
147 (Fig. 2). This suggests that our model works well in estimating N_{CCN} with a high aerosol
148 loading environment.

149 Given that the primary goal of model development is to more accurately predict
150 the relatively large spatial and temporal CCN concentrations, finally, the RFRM model
151 was employed by integrating with multisource-fusion datasets (more details see 2.3.3)
152 to predict regional and yearly scale CCN concentrations. Supplemental Table 1 also
153 provides more details about the multisource-fusion datasets. According to the results of
154 the RFRM model and SHAP (SHapley Additive exPlanations) analysis (Fig. 2), fine
155 particle chemical components such as organic matter (OM), ammonium ions (NH_4^+),



156 and nitrate ions (NO_3^-) are the top-ranked factors with significant impacts on the
157 prediction results. Particularly, OM is the most important factor among all influencing
158 factors for the predictions. Therefore, prior to conducting predictions based on the
159 multi-source dataset, we compared and validated the reliability of the TAP multi-source
160 dataset against the available observational data (Fig. S1). Compared with the
161 observations at the BJ site, mass concentrations for the OM and sulfate from TAP
162 dataset were largely underestimated by approximately 50% (Fig. S1). Therefore, a
163 twofold correction factor was applied to these components when estimating the regional
164 scale and interannual CCN concentrations in NCP. Cartesian coordinates were also
165 added as input due to the spatiotemporal nature of the input data (Yang et al., 2022).

166 **2.3 Data and other details in the model construction**

167 **2.3.1 Ground-based measurements and datasets**

168 Ground measurements of atmospheric gaseous precursors, fine particles chemical
169 compositions, and CCN number concentration (at supersaturations of 0.2% and 0.4%)
170 were collected during six field campaigns at three sites in the NCP, used to construct
171 the ML-based N_{CCN} prediction model. The six campaigns were conducted as follows:
172 at the Beijing (BJ) site from 8–30 November 2014, 20 August to 6 October 2015, 16
173 November to 20 December 2016, and 28 May to 27 June 2017; at the Xingtai (XT) site
174 from 17 May to 14 June 2016; and at the Gucheng (GC) site from 16 November to 15
175 December 2018. Here it should be noted that this study additionally selected six
176 independent cases to assess the performance of the developed ML-based model in



177 predicting N_{CCN} (Fig. S2). They are accordingly named BJ2014_WIN, BJ2015_AUT,
178 BJ2016_WIN, BJ2017_SUM, XT2016_SPR, and GC2018_WIN (Fig. 3a). These cases
179 were excluded from both the training and testing datasets described in Section 2.2.

180 The BJ site (Longitude: 116.37° E; Latitude: 39.97° N) is located at the
181 meteorological tower station of the Institute of Atmospheric Physics, Chinese Academy
182 of Sciences. It is representative of the general emission conditions in urban areas of the
183 northern NCP. The primary pollution sources here are surrounding traffic and
184 residential emissions. The XT site (Longitude: 114.37° E; Latitude: 37.18° N) is
185 situated at a national weather station. It is primarily influenced by emissions from
186 surrounding towns and factories (e.g., coal-fired power plants, coking, steel, cement,
187 and chemical industries) and thus reflects polluted suburban conditions in the southern
188 NCP. The GC site (Longitude: 115.74° E; Latitude: 39.15° N) is located at the
189 Integrated Ecological-Meteorological Observation and Experiment Station of the
190 Chinese Academy of Meteorological Sciences. Surrounded mainly by nearby villages,
191 farmland, and transportation networks, this site represents the regional background
192 pollution in the northern NCP.

193 The CCN number concentrations were measured by using the Droplet Measurement
194 Technologies CCN counter (model CCNC-100, DMT Inc. Lance et al., 2006) at BJ and
195 XT site. The supersaturation (S) levels set for each CCN measurement cycle were 0.1%,
196 0.2%, 0.4%, and 0.8%, respectively. Another measurement at GC site was referred from
197 Tao et al. (2021). In this study, the comparisons between the measured and predicted
198 N_{CCN} were mostly based on the value at $S=0.2\%$ and $S=0.4\%$. The observed N_{CCN} varies



199 from a few hundred to tens of thousands at these sites, and the campaign mean mass
200 concentration of $PM_{2.5}$ ranges from 35.6 to 160 $\mu\text{g m}^{-3}$, indicating that the observations
201 can represent various atmospheric conditions, spanning from clean to polluted in the
202 region. More details about the observations could be found in Fan et al. (2020), Ren et
203 al. (2018), and Zhang et al. (2019). In addition, the long-term measurement of particle
204 number size distribution (PNSD) at a field site in Beijing (Fig. S3, Shang et al., 2022)
205 is also used for deriving the long-term trend of yearly averaged N_{CCN} .

206 **2.3.2 N_{CCN} simulated by WRF-Chem model**

207 The WRF-Chem version 4.1.5 is used to simulate N_{CCN} in this study, which nested
208 a domain in 10 km \times 10 km covering the entire NCP (Fig. S4) and contained 181 \times 170
209 grids. The simulation in WRF-Chem is conducted from 1 January 2014 to 31 December
210 2018 with an hourly resolution. In the WRF-Chem modeling system, the sectional
211 Model for Simulating Aerosol Interactions and Chemistry (MOSAIC), the Morrison
212 two-moment scheme (Morrison et al., 2009) and the Carbon Bond Mechanism Z
213 photochemical mechanism (Zaveri et al., 1999) are employed. We also compared the
214 simulation using the Regional Acid Deposition Model (Stockwell et al., 1990) and the
215 Lin microphysics scheme (Lin et al., 1983). Considering the calculation efficiency and
216 accuracy with the measurements, the CBMZ-MOSAIC and Morrison 2-moment
217 scheme were finally applied to simulate the long-term CCN concentration. More details
218 about the other parameterizations used for WRF-Chem simulation were given in SI.



219 **2.3.3 Multisource-Fusion datasets**

220 The RFRM model was further employed by integrating with multisource-fusion
221 datasets to estimate regional scale and long-term CCN concentrations (Section 3.3 and
222 3.4). Those dataset include the N_{CCN} simulated by WRF-Chem, the chemical
223 components of $PM_{2.5}$ (organic, sulfate, nitrate, ammonium) that were taken from the
224 Tsinghua University Tracking Air Pollution in China dataset (Liu et al., 2022), the gas
225 and particulate pollutants (nitrogen dioxide (NO_2), sulfur dioxide (SO_2), carbon
226 monoxide (CO), ozone (O_3) and $PM_{2.5}$) which were collected from the China National
227 Environmental Monitoring Centre network, meteorological parameters (from the
228 European Centre for Medium-range Weather Forecasts Reanalysis version 5, ERA-5) of
229 temperature (Tem), relative humidity (RH), planetary boundary layer height (BLH),
230 surface pressure (SP) and surface net solar radiation (SNSR).

231 **3. Results**

232 **3.1 The relative importance of the input parameters to the prediction of N_{CCN}**

233 The SHAP algorithm is employed to interpret the outputs of the RFRM model, as
234 illustrated in Fig. 2. Organic matter emerges as the most crucial indicator with the
235 highest SHAP value. It indicates a strong, monotonic positive effect of organic matter
236 on N_{CCN} , with low-OM contributing negatively to N_{CCN} , while high-OM contributing
237 positively to N_{CCN} . It reflects the synergistic positive effect of organic matter
238 concentration on the variation of CCN number concentration, which differs from the
239 conventional view that inorganic salts contribute more CCN due to their strong



240 hygroscopicity (Petters and Kreidenweis, 2007). However, in fact, previous studies
241 have shown that in the North China region where the proportion and concentration of
242 OM are both high, organic particles affected by strong anthropogenic emission sources
243 was found exhibit strong hygroscopicity, enabling them to serve as more effective CCN
244 (Liu et al., 2021); in addition, the surface tension lowering effect of OM particles in this
245 region can also enhance particle CCN activity (Fan et al., 2024). Therefore, the SHAP
246 analysis results further confirm the conclusions of previous studies. $PM_{2.5}$ concentration
247 is the second most important factor, only after OM. It shows that higher $PM_{2.5}$ mass
248 concentration correspond to positive SHAP contributions, meaning that increase in
249 $PM_{2.5}$ will increase CCN concentrations. Actually, it has been demonstrated that the
250 strong association of $PM_{2.5}$ mass concentration with the CCN concentrations in NCP,
251 due to the synchronously growth in particle size and enhanced hygroscopicity (Zhang
252 et al., 2019). The SO_2 and NO_2 are also among the relatively top-ranked factors, with
253 positive effect on CCN levels according to the SHAP value. This highlights the
254 potential impact of gaseous precursors on CCN activity. Ammonium (NH_4^+) also
255 contributes positively to CCN predictions, with relative larger SHAP values, though
256 occasional negative values might suggest context-dependent effects under certain
257 chemical regimes (Dinar et al., 2008). Nitrate (NO_3^-) exhibited a moderate positive
258 effect in the model, with relatively concentrated SHAP value distributions indicating
259 its stable contribution to the model output. Temperature demonstrated a bidirectional
260 influence, suggesting nonlinear modulation of CCN activity potentially associated with



261 the temperature dependence of nucleation growth and secondary generation of particles
262 (Song et al., 2022).

263 In contrast, other meteorological parameters (RH, SP, PBL, SNSR) showed lower
264 SHAP values, implying marginal contribution to the model output during the study
265 period. Overall, the high SHAP values of organic matter (OM), PM_{2.5}, SO₂ and NO₂
266 underscore the critical role of chemical constituents and gaseous precursors in CCN
267 predictions, which can be well explained by previous known physical mechanisms of
268 the impact of aerosol particles atmospheric processes on CCN activity.

269 **3.2 Performance of the RFRM model in predicting N_{CCN} at field sites in the NCP**

270 To assess the performance of the RFRM model in predicting N_{CCN} , we compare
271 the predicted N_{CCN} with both the simulations by WRF-Chem and the observations in
272 NCP (Fig. 3). Here, the selected independent cases of the six campaigns (BJ2014_WIN,
273 BJ2015_AUT, BJ2016_WIN, BJ2017_SUM, XT2016_SPR, and GC2018_WIN) used
274 for validating the model performance (Fig. 3a). The observed N_{CCN} varies from a few
275 hundred to tens of thousands at these sites, and the campaign mean mass concentration
276 of PM_{2.5} ranges from 35.6 to 160 $\mu\text{g m}^{-3}$ (Fig. 3b), indicating that the observations can
277 represent various atmospheric conditions, spanning from clean to polluted in the region.
278 Fig. 3a shows N_{CCN} at a supersaturation of 0.2% (the typical range of supersaturations
279 in clouds). It shows that, for all the six periods, the time series of N_{CCN} predicted by the
280 RFRM model agrees better with the observed N_{CCN} compared to simulations by the
281 WRF-Chem model. Although both the RFRM model and WRF-Chem exhibit
282 underestimations in observed N_{CCN} , the slope between the predicted and observed N_{CCN}



283 ranged from 0.41 by the WRF-Chem model (with the NMB of -40%, RMSE of 4569
284 cm^{-3}) to 0.68 by the RFRM model (the NMB of -8% and RMSE of 2045 cm^{-3}) (Fig. 3c),
285 corresponding the R^2 increased from 0.33 to 0.86. Compared to WRF-Chem
286 simulations, the RFRM model showed the greatest improvement during the winter
287 campaigns with the NMB decreased from -45% to -14% (eg., slope increased from 0.39
288 to 0.67, R^2 from 0.32 to 0.83) when $\text{PM}_{2.5}$ concentrations were usually higher. For
289 example, during the GC2018_WIN campaign, the observed N_{CCN} is underestimated as
290 large as 71% by the WRF-Chem (Fig. S6), while the underestimation of observed
291 values by the RFRM has been significantly improved, with an enhancement of
292 approximately 36% (Fig. S5). WRF-Chem simulations for warm seasons noticeably
293 improved (the NMB of -29%), e.g., the uncertainty decreased to 25% during the
294 BJ2017_SUM campaign (Fig. S5). Overall, the RFRM model still performs better than
295 the WRF-Chem model and is with the NMB of 4% during summer campaigns (eg.,
296 slope of 0.67 and R^2 of 0.67). Occasionally, the WRF-Chem model overestimated the
297 N_{CCN} apparently, e.g., the episodes of September 24 to 25 during the BJ2015_AUT
298 campaign, and May 7 to 8 during the XT2017_SPR campaign. In addition, we also
299 evaluated the model performance based on another observation at GC site in January
300 (Zhang et al., 2020) (Fig. S6). And Figure S7 shows comparisons of N_{CCN} at $S=0.4\%$.
301 Both exhibit similarly good performance as shown in Fig. 3. Overall, the RFRM model
302 performs well and can accurately capture the observed fluctuations during these
303 episodes. The improvements in RFRM model also demonstrate the effectiveness of the
304 model trained on atmospheric variables to revise the simulation in model.



305 In our case, the underestimation of N_{CCN} by the WRF-Chem model is likely due to
306 the overestimation of the organics and BC mass fraction induced by WRF-Chem (Fig.
307 S8), but the underestimation of the hygroscopic parameter of organics, and the
308 simplified prescriptions in particle size distribution (Fanourgakis et al., 2019). In fact,
309 the much lower and fixed hygroscopicity parameter value in WRF-Chem model does
310 not represent the hygroscopicity of organics throughout the study period (Liu et al.,
311 2021). Neglecting to distinguish between POA and SOA information during the training
312 of the RFRM model may cause the overestimation of N_{CCN} when POA dominates.
313 Uncertainties incurred by the RFRM model could also originate from the lack of
314 physical interpretability in these ML-based models (Wei et al., 2023). Additional input
315 parameters that carry rich and more meaningful information (e.g., particle number size
316 distribution, aerosol sources and other secondary processes) are expected to further
317 improve the predictability of N_{CCN} in future.

318 **3.3 Performance of the RFRM model in predicting hourly-to-yearly-scale N_{CCN}**

319 To further examine the performance of the RFRM model in predicting N_{CCN} at
320 different time scales, we compare the RFRM model-predicted hourly-to-yearly N_{CCN} in
321 Beijing with both WRF-Chem simulations and observed values (Fig. 4). The RFRM
322 model captures well the diurnal cycle (Fig. 4a), while the WRF-Chem model
323 underestimates N_{CCN} , especially at night. Concerning seasonal variations, similarly, the
324 RFRM model performs better with the NMB of $\sim 4\%$ compared to observations. While
325 the mean bias can increase to be $\sim 28\%$ by the WRF-Chem (Fig. 4b). Note that, the bias
326 is much greater in the cold seasons than that in the warm seasons for the WRF-Chem.



327 This is probably due to the higher wintertime CN and CCN concentrations which are
328 more difficult for models to capture and simulate (Fanourgakis et al., 2019).

329 Figure 4c shows the long-term trend of yearly averaged N_{CCN} . Here, the real
330 atmospheric long-term trend of N_{CCN} (denoted as N_{CCN_obs}) is derived using the long-
331 term measurement of particle number size distribution (PNSD) at a field site in Beijing
332 (Fig. S3) (Fig. 4d, Shang et al., 2022) and the κ calculated from the measured chemical
333 compositions based on the κ -Köhler theory (Petters and Kreidenweis, 2007). The results
334 show that the predicted average annual N_{CCN} at $S=0.2\%$ based on the RFRM model
335 agrees well with N_{CCN_obs} in terms of magnitude and long-term trend (Fig. 4c), showing
336 a decreasing trend year by year with the average annual CCN number concentration of
337 about $6216 \pm 3624 \text{ cm}^{-3}$ in 2014 and $3278 \pm 2306 \text{ cm}^{-3}$ in 2018; however, although the
338 WRF-Chem simulations also show a similar decreasing trend year by year, it
339 significantly underestimates the average annual N_{CCN} of all years (with the NMB of
340 43%), resulting in the inter-annual trend lines being parallel but not coincident. The
341 small bias (within $\pm 6\%$) between the RFRM model predictions and the observations
342 may be due to the uncertainty from how N_{CCN_obs} is calculated, i.e., using the Tracking
343 Air Pollution in China (TAP) dataset to calculate κ . A comparison of the values of κ and
344 N_{CCN} between that derived using field observations and the TAP dataset shows little
345 differences (Fig. S9); actually, the long-term change of N_{CCN} is much less sensitive to
346 changes in κ values compared to the PNSD (Fig. S9c). Although absolute
347 concentrations of components in the TAP dataset deviate from observations, their mass
348 fractions are consistent (Fig. S9d), rendering the impact on the calculated κ negligible.



349 In addition, the method to calculate N_{CCN} at $S=0.2\%$ based on κ -Köhler theory would
350 cause an upper-limit uncertainty of 7% (Ren et al., 2018).

351 According to Fig. 4d-e, the long-term decreasing trend of N_{CCN} at $S=0.2\%$ from
352 2014 to 2018 is mainly attributed to a significant reduction in aerosol particle number
353 concentrations in the atmosphere. In addition, the peak diameter of the PNSD shows a
354 shift toward the left, decreasing slightly from about 70 nm in 2014 to 30 nm in 2018
355 due to the enhanced new particle formation events in recent years (Zhu et al., 2021).
356 This would also result in less aerosol particles serving as CCN. Although the κ_{chem} has
357 a slight upward trend from 2014 to 2018 (Fig. 4e), yielding decreases in activation
358 diameter and thereby more CCN, the aerosol particle hygroscopicity, however, plays
359 less significant role in regulating the long-term total N_{CCN} variations compared to the
360 changes in particle number size distribution during the period.

361 **3.4 Spatial variations of N_{CCN} derived by the RFRM model**

362 We further examine the spatiotemporal changes of N_{CCN} at $S=0.2\%$ in the NCP
363 derived by the RFRM model (Fig. 5). Regionally, the N_{CCN} predicted by the RFRM
364 model is also generally higher than that simulated by WRF-Chem at most of the sites.
365 The N_{CCN} derived by the RFRM model and WRF-Chem both decrease from 2014 to
366 2018 but with different decreasing rates (Fig. 5c-e). On average, N_{CCN} derived by the
367 RFRM model and WRF-Chem decrease from 4415 ± 643 to 2910 ± 789 cm^{-3} and from
368 2834 ± 1366 to 2111 ± 546 cm^{-3} respectively from 2014 to 2018 in the NCP region (Fig.
369 5c), corresponding to annual decreasing rates of approximately ~ -250 $\text{cm}^{-3} \text{ yr}^{-1}$ for
370 the RFRM model and ~ -167 $\text{cm}^{-3} \text{ yr}^{-1}$ for the WRF-Chem model (Fig. 5d-e).



371 Moreover, N_{CCN} and its changes from 2014 to 2018 predicted by the RFRM model show
372 more significant spatial variations than that simulated by the WRF-Chem model.
373 Differences in RFRM -model-predicted N_{CCN} between 2014 and 2018 (2018 minus
374 2014) show negative values at ~90% of the sites, i.e., downward trends in N_{CCN} (Fig.
375 5c1). The sites with apparent N_{CCN} reduction are mainly located in the central and
376 northern of NCP, especially in Beijing-Tianjin-Hebei (BTH) and central Shandong.
377 Sites in southern NCP have slight downward trends in N_{CCN} . The downward trend is
378 consistent with the variations in concentration of gaseous pollutants due to the emission
379 reduction in past years in China (Fig. S10). Interestingly, note a few sites with positive
380 values (upward trends in N_{CCN}) are mainly located along the coast. An increase in the
381 fraction of accumulation-mode particles in coastal areas has been reported contributing
382 more CCN (Zhu et al., 2021). This demonstrates the good performance of the RFRM
383 model in capturing the real-time spatial variations of CCN on a regional scale. While
384 the WRF-Chem model might mask the variations of N_{CCN} among different sites. This
385 will smooth out the true impact of aerosols on weather and climate at regional or local
386 scales, leading to uncertainties in model simulations.

387 **3.5 Sensitivity of the cloud parameters and radiative forcing to CCN prediction**

388 **biases**

389 To evaluate the effects introduced by N_{CCN} prediction biases to the aerosol indirect
390 effects, we further incorporate the deviations between observed N_{CCN} (denoted as
391 CCN_{OBS}) and N_{CCN} predicted by the RFRM model (denoted as CCN_{RFRM}) and the
392 simulated by the WRF-Chem model (denoted as $CCN_{WRF-Chem}$) into calculations of the



393 cloud parameters and radiative forcing, as are shown in Fig. 6 (for $S=0.2\%$) and Fig.
394 S_{11} (for $S=0.4\%$). Typically, aerosol particles serving as CCN could indirectly affect
395 the global climate by the Twomey (Twomey, 1977) and Albrecht effects (Albrecht,
396 1989). According to Wang et al. (2019), two parameters of cloud optical thickness (τ)
397 and the absorption coefficient ($1-\omega_0$) can be used to estimate the Twomey effects. The
398 process of cloud-to-rain conversion, which can be parameterized by the critical radius
399 (r_c) and the cloud-to-rain conversion threshold function (TA), is critical to estimate the
400 Albrecht effect. Therefore, the r_c and TA is also calculated here. Indirect (cloud)
401 radiative forcing (F_c) is also evaluated based on the deviations in CCN number
402 concentration under the assumption of a constant liquid water content (Charlson, 1992;
403 Wang et al., 2008). Section S2 provides details about the methods used to evaluate
404 aerosol indirect effects.

405 In general, the results show that these cloud properties are more sensitive to the
406 changes in N_{CCN} when the models underestimate the CCN number concentrations (Δ
407 $N_{CCN}<0$) compared to the cases with an overestimation (Figure. 6a-d). For example, a
408 $\sim 50\%$ underestimation (overestimation) of N_{CCN} could lead to relative deviations
409 (uncertainties) of -21% (14%) for τ , 27% (-12%) for $(1-\omega_0)$, and -11% (7%) for r_c at
410 $S=0.2\%$. Note that, on average, both the RFRM model and WRF-Chem in this study
411 show underestimations in N_{CCN} within the sensitivity zone of the cloud effect (Fig. 6),
412 It is thus expected to cause large uncertainties in evaluating the cloud radiative forcing,
413 a topic worthy of further attention. Given that the uncertainty in N_{CCN} predicted by the
414 WRF-Chem model is much greater than that of RFRM model, the uncertainties and



415 variation ranges of these cloud parameters from WRF-Chem simulations are also
416 greater. Specifically, the uncertainties of CCN_{RFRM} and $CCN_{WRF-Chem}$ lead to the
417 uncertainties of -33% to +78% and -77% to +92% respectively, for the τ (Fig. 6a and
418 a1), -44% to +50% and -48% to +344% respectively, for the $1-\omega_0$ (Fig. 6b and b1), -
419 18% to +34% and -53% to +38% respectively, for the r_c (Fig. 6c and c1), and -118% to
420 +94% and -258% to +353% respectively, for the TA (Fig. 6d and d1).

421 In addition, the underestimation of CCN would lead to underestimations of cloud
422 optical thickness τ and the critical radius r_c of the automatic cloud/rain transformation,
423 but overestimations of $(1-\omega_0)$ and the threshold function TA of the automatic cloud/rain
424 transformation, all of which depend on their physical mechanisms within the realm of
425 aerosol-cloud interactions (Stier et al., 2024) (Fig. S12). This is also the case at the other
426 supersaturation levels considered (Fig. S11).

427 As a result, we derive that the mean underestimation of $\sim 59 \pm 19\%$ in N_{CCN} at $S=0.2\%$
428 caused by the WRF-Chem leads to underestimations of $26 \pm 11\%$ in the τ , $14 \pm 7\%$ in the
429 r_c , and an overestimation of $35 \pm 191\%$ in the absorption coefficient $(1-\omega_0)$ and $93 \pm 42\%$
430 in the TA . While, the mean uncertainties for all these parameters are largely reduced
431 when the mean underestimation of $\sim 31 \pm 15\%$ in N_{CCN} at $S=0.2\%$ that is caused by
432 RFRM model is applied (Fig. 6e). For example, the underestimation of cloud optical
433 thickness τ decreases to $\sim 12\%$, an improvement compared to the underestimation of
434 about 14% by the WRF-Chem model. Also, the RFRM model reduces the
435 underestimation of the critical radius r_c of the automatic cloud/rain transformation to
436 only $\sim 6\%$. Ultimately, the uncertainty of cloud radiative forcing F_c has been



437 significantly reduced from an overestimation of $1.89 \pm 0.78 \text{ W m}^{-2}$ by the WRF-Chem
438 model to only $0.81 \pm 0.63 \text{ W m}^{-2}$ by the RFRM model, showing the high sensitivity of
439 climate forcing to the uncertainties in CCN number concentrations. Note that a
440 limitation when evaluating the cloud radiative forcing based on the assumption of cloud
441 fraction and the fractional transmission is the approximate analytical expression.
442 Therefore, the results presented here may represent the upper limit, and the sensitivity
443 of the radiative forcing to changes in N_{CCN} would be weaker over continental areas
444 (Wang et al., 2008; Yu et al., 2022).

445 **4. Summary**

446 **4.1 Discussion and conclusions**

447 In this study, using a multisource dataset of ground-based observations,
448 atmospheric variables, the N_{CCN} simulations by the WRF-Chem model, we have
449 developed a new machine-learning-based model that predicts well the regional-scale
450 N_{CCN} in the polluted NCP region. The results show that the prediction bias of N_{CCN}
451 compared to observations is approximately -31% from the RFRM model. Good
452 accuracy has also been achieved during heavy pollution periods or cold seasons. In
453 general, the RFRM model better captures the spatial differences of N_{CCN} than the WRF-
454 Chem model. In addition, the prediction reveals a long-term downward trend of N_{CCN}
455 coincident with the observed trend for the period of 2014–2018. By further
456 incorporating the N_{CCN} prediction biases into the evaluation of cloud parameters and
457 radiative forcing, we found that the cloud properties are more sensitive to the changes



458 in N_{CCN} when the models underestimated the CCN number concentrations compared to
459 the cases when the models overestimated N_{CCN} . As a result, the simulated uncertainty
460 of cloud radiative forcing F_c could be significantly reduced from an overestimation of
461 $1.89 \pm 0.78 \text{ W m}^{-2}$ by the WRF-Chem model to $0.81 \pm 0.63 \text{ W m}^{-2}$ by the RFRM model.
462 Given the simplified setting in current climate models, this work emphasizes the
463 necessity and urgency to obtain the precise N_{CCN} values, offering a new framework for
464 predicting CCN concentrations based on machine learning algorithms and effectively
465 filling the observation gap of CCN concentrations.

466 **4.2 Limitations and outlook**

467 Although the machine learning-based model established in this study has been
468 shown to significantly improve the prediction of CCN number concentration, the
469 prediction results still exhibit considerable deviations, which may be attributed to
470 several factors. First, it is related to the inherent limitation of the random forest method
471 in describing time series with extreme values and large short-term fluctuations. In the
472 future, other advanced machine learning algorithms (e.g., Long Short-Term Memory,
473 Transformer) can be integrated to optimize and improve the results. Second, in this
474 study, the observational data only from six campaigns at three sites are analyzed.
475 Validating the simulated N_{CCN} through comparisons with observations at more ground
476 sites is warranted in future. Also, it is crucial to obtain comprehensive monitoring data
477 of CCN and other key aerosol properties (e.g., particle size distribution, chemical
478 compositions) in different environments. Overall, the RFRM framework presented here
479 relies on available atmospheric state variables (e. g., chemical compositions, gas



480 pollutants, and meteorology elements) and significantly improves the accuracy of N_{CCN}
481 prediction, thereby helping to bridge observational gaps. Our modeling framework
482 could then be used to simulate ground-level CCN data in other regions around the world
483 and even on a global scale. This new modeling framework could also guide the
484 development of machine learning based models to predict CCN vertical profiles, which
485 are critical for the accurate evaluation of the ACI effect.

486 **Code and Data availability**

487 The data and code are publicly accessible at <https://zenodo.org/records/18932004> (Ren
488 et al., 2026). This includes the WRF-Chem model version 4.1.5 used in this study, the
489 machine learning code, the corresponding training, testing datasets and the CCN
490 observation datasets, the emissions inventory and scripts used in WRF-Chem and the
491 scripts used for plotting, supporting the findings of this study. The release version of
492 WRF-Chem is also open-access and can be publicly available at NCAR
493 https://www2.mmm.ucar.edu/wrf/users/download/get_source.html (Skamarock et al.,
494 2019, last access: 10 May 2025). The initial meteorological variables are from the
495 National Center for Environmental Prediction's Final Operational Global (NCEP/FNL)
496 and available at <https://doi.org/10.5065/D6M043C6> (NCEP, 2000).

497 **Supplement**

498 The Supplement contains the information of additional descriptions of the WRF-Chem
499 simulation (parameterization scheme, emission inventory and initial and boundary



519 valuable comments and suggestions, which have greatly improved the quality of this
520 paper.

521 **References**

- 522 Albrecht, B.: Aerosols, cloud microphysics, and fractional cloudiness, *Science*,
523 245(4923), 1227–1230, <https://doi.org/10.1126/science.245.4923.1227>, 1989.
- 524 Charlson, R., Lovelock, J., Andreae, M., Warren, S.: Oceanic phytoplankton,
525 atmospheric sulphur, cloud albedo and climate, *Nature*, 326(6114), 655–661,
526 <https://doi.org/10.1038/326655a0>, 1987.
- 527 Chen, L., Zhang, F., Zhang, D. et al.: Measurement report: Hygroscopic growth of
528 ambient fine particles measured at five sites in China, *Atmos. Chem. Phys.*, 22(10),
529 6773–6786, <https://doi.org/10.5194/acp-22-6773-2022>, 2022.
- 530 Charlson, R. et al.: Climate Forcing by Anthropogenic Aerosols, *Science*, 255, 423–430,
531 <https://doi.org/10.1126/science.255.5043.423>, 1992.
- 532 Dinar, E., Anttila, T., Rudich, Y.: CCN activity and hygroscopic growth of organic
533 aerosols following reactive uptake of ammonia, *Environmental Science &*
534 *Technology*, 42(3), 912–918, <https://doi.org/10.1021/es071874p>, 2008.
- 535 Fanourgakis, G., Kanakidou, M., Nenes, A. et al.: Evaluation of global simulations of
536 aerosol particle and cloud condensation nuclei number, with implications for cloud
537 droplet formation, *Atmos. Chem. Phys.*, 19(13), 8591–8617, <https://doi.org/10.5194/acp-19-8591-2019>, 2019.
- 539 Fan, X., Liu, J., Zhang, F. et al.: Contrasting size-resolved hygroscopicity of fine
540 particles derived by HTDMA and HR-ToF-AMS measurements between summer
541 and winter in Beijing: the impacts of aerosol aging and local emissions, *Atmos.*
542 *Chem. Phys.*, 20, 915–929, <https://doi.org/10.5194/acp-20-915-2020>, 2020.
- 543 Fan, T., Ren, J., Liu, C., Li, Z., Liu, J., Sun, Y., et al.: Evidence of surface - tension
544 lowering of atmospheric aerosols by organics from field observations in an urban
545 atmosphere: Relation to particle size and chemical composition, *Environmental*



- 546 Science & Technology, 58(26), 11363 – 11375.
547 <https://doi.org/10.1021/acs.est.4c03141>, 2024.
- 548 Grange, S., Carslaw, D., Lewis, A. et al.: Random Forest meteorological normalisation
549 models for Swiss PM₁₀ trend analysis, *Atmos. Chem. Phys.*, 18(9), 6223–6239,
550 <https://doi.org/10.5194/acp-18-6223-2018>, 2018.
- 551 Geng, G., Xiao, Q., Liu, S. et al.: Tracking air pollution in China: near real-time PM_{2.5}
552 retrievals from multisource data fusion, *Environmental Science & Technology*,
553 55(17), 12106–12115, <https://doi.org/10.1021/acs.est.1c01863>, 2021.
- 554 Guo, J., Miao, Y., Zhang, Y. et al.: The climatology of planetary boundary layer height
555 in China derived from radiosonde and reanalysis data, *Atmos. Chem. Phys.*, 16(20),
556 13309–13319, <https://doi.org/10.5194/acp-16-13309-2016>, 2016.
- 557 IPCC. Summary for Policymakers. In *Climate Change 2021: The Physical Science*
558 *Basis. Contribution of Working Group I to the Sixth Assessment Report of the*
559 *Intergovernmental Panel on Climate Change*; Cambridge University Press:
560 Cambridge, United Kingdom and New York, NY, USA. (2021).
- 561 Liu, S., Geng, G., Xiao, Q., Zheng, Y., Liu, X., Cheng, J., & Zhang, Q.: Tracking daily
562 concentrations of PM_{2.5} chemical composition in China since 2000, *Environ Sci*
563 *Technol*, 56, 16517–16527, <https://doi.org/10.1021/acs.est.2c06510>, 2022.
- 564 Lance, S., Nenes, A., Medina, J. et al.: Mapping the operation of the DMT continuous
565 flow CCN counter, *Aerosol Science and Technology*, 40(4), 242–254,
566 <https://doi.org/10.1080/02786820500543290>, 2006.
- 567 Liu, C., Wang, T., Rosenfeld, D. et al.: Anthropogenic effects on cloud condensation
568 nuclei distribution and rain initiation in East Asia, *Geophysical Research Letters*,
569 47, <https://doi.org/10.1029/2019GL086184>, 2020.
- 570 Liu, J., Li, Z.: Significant underestimation in the optically based estimation of the
571 aerosol first indirect effect induced by the aerosol swelling effect, *Geophysical*
572 *Research Letters*, 45(11), 5690–5699, <https://doi.org/10.1029/2018GL077679>,
573 2018.



- 574 Liang, M., Tao, J., Ma, N. et al.: Prediction of CCN spectra parameters in the North
575 China Plain using a random forest model, *Atmospheric Environment*, 289, 119323,
576 <https://doi.org/10.1016/j.atmosenv.2022.119323>, 2022.
- 577 Li, S., Zhang, F., Jin, X. et al.: Characterizing the ratio of nitrate to sulfate in ambient
578 fine particles of urban Beijing during 2018-2019, *Atmospheric Environment*, 237,
579 <https://doi.org/10.1016/j.atmosenv.2020.117662>, 2020.
- 580 Liu, J., Zhang, F., Xu, W. et al.: Hygroscopicity of organic aerosols linked to formation
581 mechanisms, *Geophysical Research Letters*, 48,
582 <https://doi.org/10.1029/2020GL091683>, 2021.
- 583 Lin, Y., Farley, R., Orville, H.: Bulk parameterization of the snow field in a cloud model,
584 *Journal of Applied Meteorology and Climatology*, 22(6): 1065–1092,
585 [https://doi.org/10.1175/1520-0450\(1983\)022<1065:BPOTSF>2.0.CO;2](https://doi.org/10.1175/1520-0450(1983)022<1065:BPOTSF>2.0.CO;2), 1983.
- 586 Menon & Rotstayn.: The radiative influence of aerosol effects on liquid-phase cumulus
587 and stratiform clouds based on sensitivity studies with two climate models,
588 *Climate Dynamics*, 27, 4, 345–356, <https://doi.org/10.1007/s00382-006-0139-3>,
589 2006.
- 590 Moore, R., Cerully, K., Bahreini, R. et al.: Hygroscopicity and composition of
591 California CCN during summer 2010, *Journal of Geophysical Research:
592 Atmospheres*, 117, <https://doi.org/10.1029/2011JD017352>, 2012.
- 593 Morrison, H., Thompson, G., Tatarskii, V.: Impact of cloud microphysics on the
594 development of trailing stratiform precipitation in a simulated squall line:
595 Comparison of one-and two-moment schemes, *Monthly Weather Review*, 137(3),
596 991–1007, <https://doi.org/10.1175/2008MWR2556.1>, 2009.
- 597 Nair, A. A., Yu, F.: Using machine learning to derive cloud condensation nuclei number
598 concentrations from commonly available measurements, *Atmos. Chem. Phys.*,
599 20(21), 12853–12869, <https://doi.org/10.5194/acp-20-12853-2020>, 2020.
- 600 Nair, A. A., Yu, F., Campuzano-Jost, P. et al.: Machine learning uncovers aerosol size
601 information from chemistry and meteorology to quantify potential cloud-forming



- 602 particles, *Geophysical Research Letters*, 48, [https://doi.org/](https://doi.org/10.1029/2021GL094133)
603 10.1029/2021GL094133, 2021.
- 604 NCEP: NCEP FNL Operational Model Global Tropospheric Analyses, continuing from
605 July 1999, National Centers for Environmental Prediction [Data set],
606 <https://doi.org/10.5065/D6M043C6>, 2000 (last access: 10 May 2025).
- 607 Petters, M. D., Kreidenweis, S. M.: A single parameter representation of hygroscopic
608 growth and cloud condensation nucleus activity, *Atmos. Chem. Phys.*, 7(8), 1961–
609 1971, <https://doi.org/10.5194/acp-7-1961-2007>, 2007.
- 610 Ren, J., Zhang, F., Wang, Y. et al.: Using different assumptions of aerosol mixing state
611 and chemical composition to predict CCN concentrations based on field
612 measurements in urban Beijing, *Atmos. Chem. Phys.*, 18(9), 6907–6921,
613 <https://doi.org/10.5194/acp-18-6907-2018>, 2018.
- 614 Ren, J., Chen, L., Liu, J., and Zhang, F.: The density of ambient black carbon retrieved
615 by a new method: implications for cloud condensation nuclei prediction, *Atmos.*
616 *Chem. Phys.*, 23, 4327–4342, <https://doi.org/10.5194/acp-23-4327-2023>, 2023.
- 617 Rose, C., Collaud Coen, M. et al.: Seasonality of the particle number concentration and
618 size distribution: a global analysis retrieved from the network of Global
619 Atmosphere Watch (GAW) near-surface observatories, *Atmos. Chem. Phys.*, 21,
620 17185–17223, <https://doi.org/10.5194/acp-21-17185-2021>, 2021.
- 621 Rosenfeld, D., Zheng, Y., Hashimshoni, E. et al.: Satellite retrieval of cloud
622 condensation nuclei concentrations by using clouds as CCN chambers,
623 *Proceedings of the National Academy of Sciences*, 113(21), 5828–5834,
624 <https://doi.org/10.1073/pnas.1514044113>, 2016.
- 625 Ren, J., Zou, S., Xu, H., Liu, G., Wang, Z., Zhang, A., Zhao, C., Hu, M., Shang, D.,
626 Tang, L., Huang, R.-J., Sun, Y., & Zhang, F.: Machine learning significantly
627 improves the simulation of hourly-to-yearly scale cloud nuclei concentration and
628 radiative forcing in polluted atmosphere [Data set]. Zenodo.
629 <https://zenodo.org/records/18932004>, 2026.



- 630 Stockwell, W. R., Middleton, P., Chang, J. S., et al.: The second generation regional
631 acid deposition model chemical mechanism for regional air quality modeling,
632 Journal of Geophysical Research: Atmospheres, 95(D10): 16343–16367,
633 <https://doi.org/10.1029/JD095iD10p16343>, 1990.
- 634 Sotiropoulou, R. E. P., Nenes, A., Adams, P. J. et al.: Cloud condensation nuclei
635 prediction error from application of Köhler theory: Importance for the aerosol
636 indirect effect, Journal of Geophysical Research: Atmospheres, 112(D12)
637 <https://doi.org/10.1029/2006JD007834>, 2007.
- 638 Shen, Y., Virkkula, A., Ding, A. et al.: Estimating cloud condensation nuclei number
639 concentrations using aerosol optical properties: role of particle number size
640 distribution and parameterization, Atmos. Chem. Phys., 19(24), 15483–15502,
641 <https://doi.org/10.5194/acp-19-15483-2019>, 2019.
- 642 Sha, T., Ma, X.Y., Jia, H. L. et al.: Exploring the influence of two inventories on
643 simulated air pollutants during winter over the Yangtze River Delta, Atmospheric
644 Environment, 206, 170–182, <https://doi.org/10.1016/j.atmosenv.2019.03.006>,
645 2019.
- 646 Schmale, J., Henning, S., Decesari, S. et al.: Long-term cloud condensation nuclei
647 number concentration, particle number size distribution and chemical composition
648 measurements at regionally representative observatories, Atmos. Chem. Phys.,
649 18(4), 2853–2881, <https://doi.org/10.5194/acp-18-2853-2018>, 2018.
- 650 Song, C., Becagli, S., Beddows, D. C. S. et al.: Understanding Sources and Drivers of
651 Size-Resolved Aerosol in the High Arctic Islands of Svalbard Using a Receptor
652 Model Coupled with Machine Learning, Environmental Science & Technology,
653 56(16), 11189–11198, <https://doi.org/10.1021/acs.est.1c07796>, 2022.
- 654 Shang, D., Tang, L., Fang, X. et al.: Variations in source contributions of particle
655 number concentration under long-term emission control in winter of urban Beijing,
656 Environmental Pollution, 304, 119072,
657 <https://doi.org/10.1016/j.envpol.2022.119072>, 2022.



- 658 Stier, P., van, den, Heever, S. C., Christensen, M. W. et al.: Multifaceted aerosol effects
659 on precipitation, *Nature Geoscience*, 17(8), 719–732,
660 <https://doi.org/10.1038/s41561-024-01482-6>, 2024.
- 661 Skamarock, W., Klemp, J., Dudhia, J., Gill, D. O., Liu, Z., Berner, J., Wang, W., Powers,
662 J. G., Duda, M. G., Barker, D., and Huang, X.-Y.: A Description of the Advanced
663 Research WRF Model Version 4.1, UCAR/NCAR, [https://doi.org/10.5065/1dfh-](https://doi.org/10.5065/1dfh-6p97)
664 [6p97](https://doi.org/10.5065/1dfh-6p97), 2019 (code available at [https://www2.mmm.ucar.edu/wrf/users/download/](https://www2.mmm.ucar.edu/wrf/users/download/get_source.html)
665 [get_source.html](https://www2.mmm.ucar.edu/wrf/users/download/get_source.html), last access: 10 May 2025).
- 666 Tao, J., Kuang, Y., Ma, N., Hong, J., Sun, Y., Xu, W., Zhang, Y., He, Y., Luo, Q., Xie,
667 L., Su, H., and Cheng, Y.: Secondary aerosol formation alters CCN activity in the
668 North China Plain, *Atmos. Chem. Phys.*, 21, 7409–7427,
669 <https://doi.org/10.5194/acp-21-7409-2021>, 2021.
- 670 Twomey, S.: The influence of pollution on the shortwave albedo of clouds, *Journal of*
671 *the atmospheric sciences*, 34(7), 1149–1152, [https://doi.org/10.1175/1520-](https://doi.org/10.1175/1520-0469(1977)034<1149:TIOPOP>2.0.CO;2)
672 [0469\(1977\)034<1149:TIOPOP>2.0.CO;2](https://doi.org/10.1175/1520-0469(1977)034<1149:TIOPOP>2.0.CO;2), 1977.
- 673 Wei, J., Li, Z., Wang, J. et al.: Ground-level gaseous pollutants (NO₂, SO₂, and CO) in
674 China: daily seamless mapping and spatiotemporal variations, *Atmos. Chem.*
675 *Phys.*, 23, 1511–1532, <https://doi.org/10.5194/acp-23-1511-2023>, 2023.
- 676 Wang, Y., Niu, S., Lv, J. et al.: A new method for distinguishing unactivated particles in
677 cloud condensation nuclei measurements: implications for aerosol indirect effect
678 evaluation, *Geophysical Research Letters*, 46, 14,185–14,194,
679 <https://doi.org/10.1029/2019GL085379>, 2019.
- 680 Wang, J., Lee, Y.-N., Daum, P. H., Jayne, J., and Alexander, M. L.: Effects of aerosol
681 organics on cloud condensation nucleus (CCN) concentration and first indirect
682 aerosol effect, *Atmos. Chem. Phys.*, 8, 6325–6339, [https://doi.org/10.5194/acp-8-](https://doi.org/10.5194/acp-8-6325-2008)
683 [6325-2008](https://doi.org/10.5194/acp-8-6325-2008), 2008.
- 684 Xu, W., Fossum, K. N., Ovadnevaite, J. et al.: The impact of aerosol size-dependent
685 hygroscopicity and mixing state on the cloud condensation nuclei potential over



- 686 the north-east Atlantic, *Atmos. Chem. Phys.*, 21, 8655–8675,
687 <https://doi.org/10.5194/acp-21-8655-2021>, 2021.
- 688 Yang, N., Shi, H., Tang, H. et al.: Geographical and temporal encoding for improving
689 the estimation of PM_{2.5} concentrations in China using end-to-end gradient boosting,
690 *Remote Sensing of Environment*, 269, 112828,
691 <https://doi.org/10.1016/j.rse.2021.112828>, 2022.
- 692 Yu, F., Luo, G., Nair, A. A., Tsigaridis, K., & Bauer, S. E.: Use of machine learning to
693 reduce uncertainties in particle number concentration and aerosol indirect
694 radiative forcing predicted by climate models, *Geophysical Research Letters*, 49,
695 <https://doi.org/10.1029/2022GL098551>, 2022.
- 696 Yu, F., Luo, G., Nair, A. A. et al.: Wintertime new particle formation and its contribution
697 to cloud condensation nuclei in the Northeastern United States, *Atmos. Chem.*
698 *Phys.*, 20(4), 2591–2601, <https://doi.org/10.5194/acp-20-2591-2020>, 2020.
- 699 Zhu, Y., Shen, Y., Li, K. et al.: Investigation of particle number concentrations and new
700 particle formation with largely reduced air pollutant emissions at a coastal semi-
701 urban site in northern China, *Journal of Geophysical Research: Atmospheres*, 126,
702 e2021JD035419, <https://doi.org/10.1029/2021JD035419>, 2021.
- 703 Zhang, Y., Tao, J., Ma, N. et al.: Predicting cloud condensation nuclei number
704 concentration based on conventional measurements of aerosol properties in the
705 North China Plain, *Science of The Total Environment*, 719, 137473,
706 <https://doi.org/10.1016/j.scitotenv.2020.137473>, 2020.
- 707 Zhang, F., Li, Y., Li, Z. et al.: Aerosol hygroscopicity and cloud condensation nuclei
708 activity during the AC3Exp campaign: Implications for cloud condensation nuclei
709 parameterization, *Atmos. Chem. Phys.*, 14(24), 13423–13437,
710 <https://doi.org/10.5194/acp-14-13423-2014>, 2014.
- 711 Zhang, F., Li, Z., Li, Y. et al.: Impacts of organic aerosols and its oxidation level on
712 CCN activity from measurement at a suburban site in China, *Atmos. Chem. Phys.*,
713 16(8), 5413–5425, <https://doi.org/10.5194/acp-16-5413-2016>, 2016.
- 714 Zhang, F., Ren, J., Fan, T. et al.: Significantly enhanced aerosol CCN activity and
715 number concentrations by nucleation-initiated haze events: A case study in urban



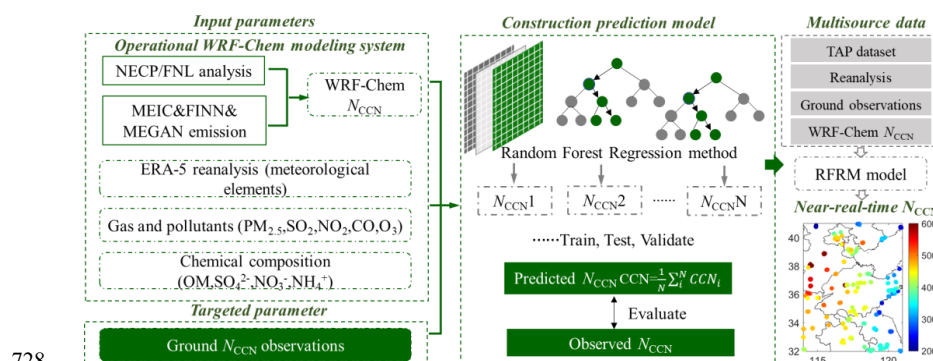
716 Beijing, Journal of Geophysical Research: Atmospheres, 124(24), 14102–14113,
 717 <https://doi.org/10.1029/2019JD031457>, 2019.

718 Zhang, F., Wang, Y., Peng, J. et al.: Uncertainty in predicting CCN activity of aged and
 719 primary aerosols, Journal of Geophysical Research: Atmospheres, 122(21),
 720 11723–11736, <https://doi.org/10.1002/2017JD027058>, 2017.

721 Zaveri, R. A., Peters, L. K.: A new lumped structure photochemical mechanism for
 722 large-scale applications, Journal of Geophysical Research: Atmospheres, 104,
 723 30387–30415, <https://doi.org/10.1029/1999JD900876>, 1999.

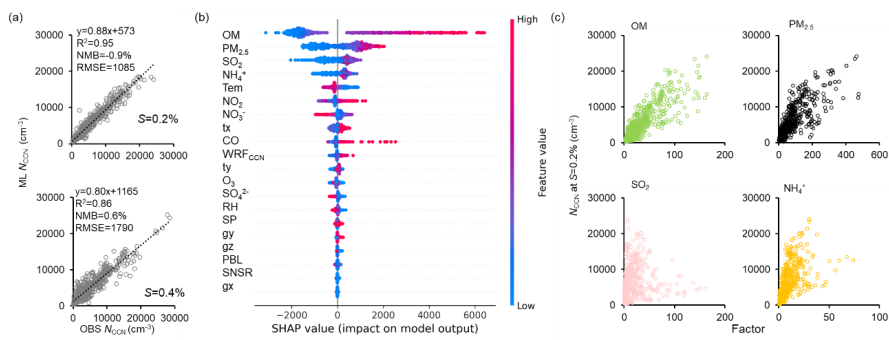
724 Zheng, B., Tong, D., Li, M., Liu, F. et al.: Trends in China’s anthropogenic emissions
 725 since 2010 as the consequence of clean air actions, Atmos. Chem. Phys., 18,
 726 14095–14111, <https://doi.org/10.5194/acp-18-14095-2018>, 2018.

727 **Figures**



729 **Fig. 1** Methodological framework of CCN number concentration prediction.

730



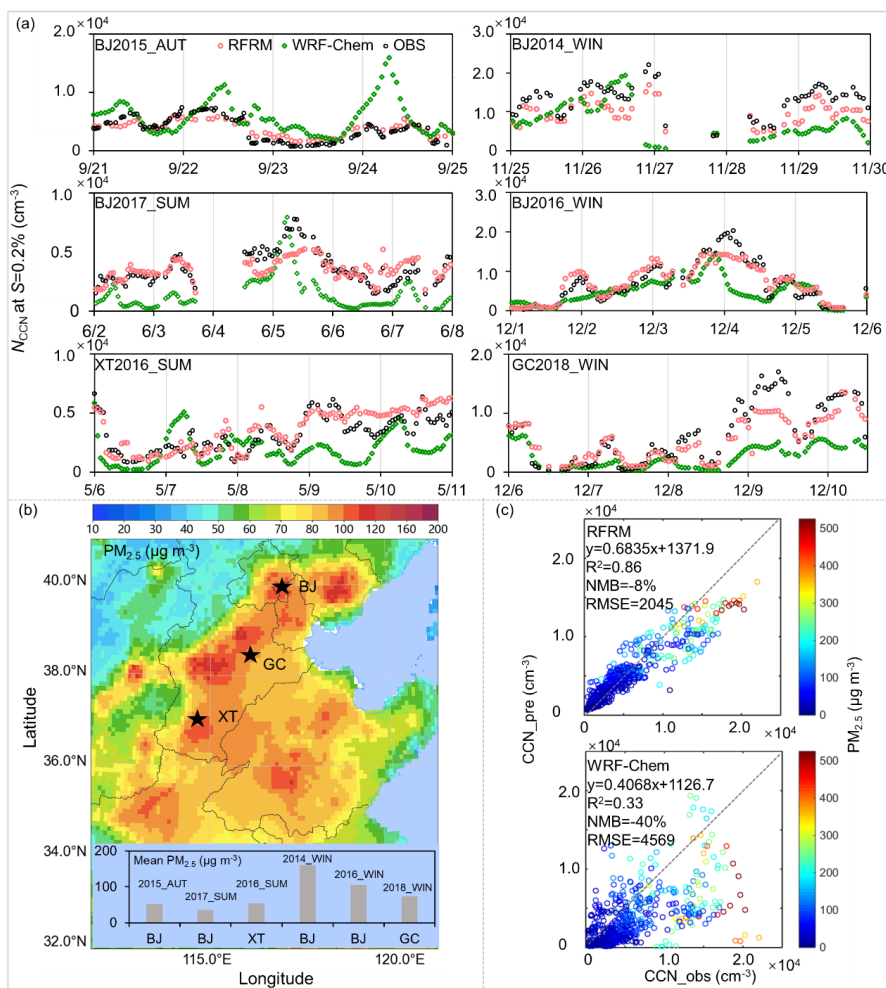
731

732 **Fig. 2** Comparison of RFRM model retrieval and observed N_{CCN} at $S=0.2\%$ and 0.4%

733 for the testing dataset (a), The SHAP value of the input parameters to the prediction of

734 N_{CCN} (b), Scatter plots of the input parameters (OM, PM_{2.5}, SO₂, NH₄⁺) with CCN

735 number concentration at $S=0.2\%$ (c).



736

737 **Fig. 3** Performance of the RFRM model in predicting N_{CCN} at field sites in NCP. (a)

738 Time series of the observed and predicted CCN number concentrations at $S=0.2\%$ for

739 the six periods (BJ2015_AUT, BJ2017_SUM, XT2016_SUM, BJ2014_WIN,

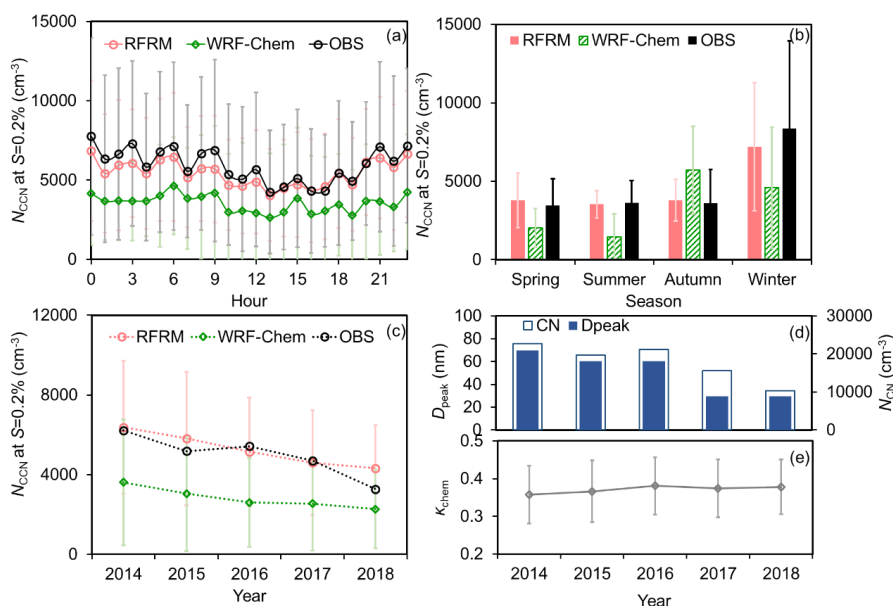
740 BJ2016_WIN, GC2018_WIN) in the North China Plain; (b) Map for average mass

741 concentration of $\text{PM}_{2.5}$ of 2014 from TAP dataset in NCP (<http://tapdata.org.cn/>) and

742 field observed average mass concentration of $\text{PM}_{2.5}$ during the six field campaigns (see

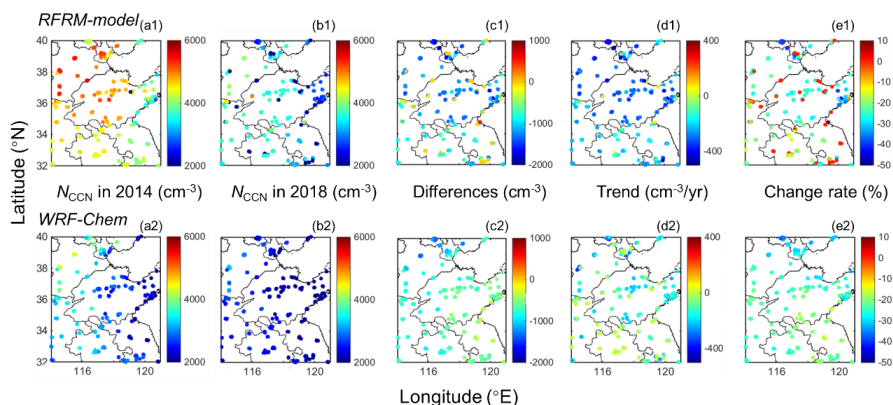


743 embedded histogram); (c) Scatter plots of the observed N_{CCN} at $S=0.2\%$ with the RFRM
 744 model predicted (top) and WRF-Chem simulated (bottom) respectively.



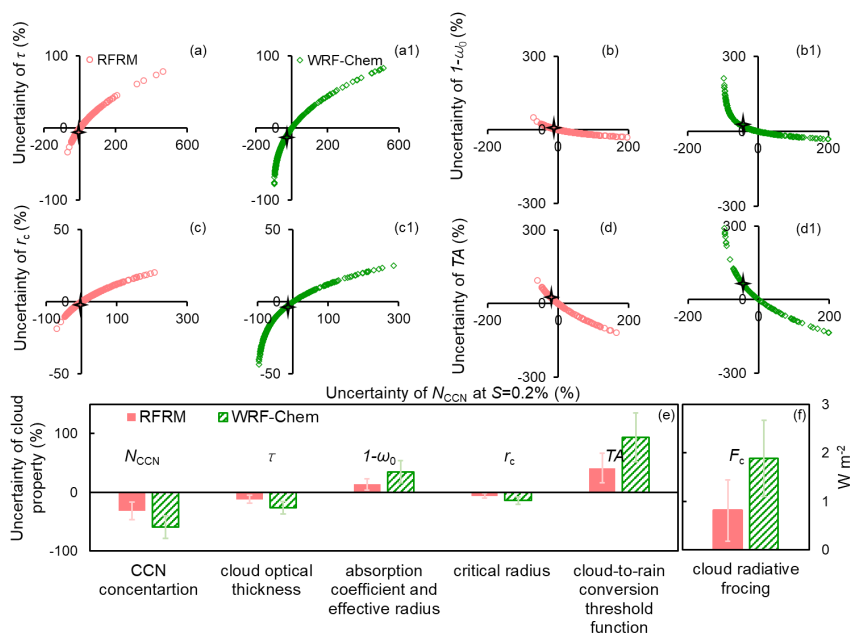
745

746 **Fig. 4** Performance of the RFRM model in predicting hourly-to-yearly scale N_{CCN} . (a)
 747 Diurnal variations of N_{CCN} at $S=0.2\%$ derived from the RFRM model, the WRF-Chem,
 748 and the observations from the field campaigns; (b) Seasonal variations, here the
 749 comparison in spring, summer, autumn and winter were conducted using the campaign
 750 averages of XT2016_SPR, BJ2017_SUM, BJ2015_AUT, and
 751 BJ2014_WIN&BJ2016_WIN respectively with the RFRM model and WRF-Chem
 752 predictions at corresponding periods; (c) Trends of annual mean N_{CCN} from 2014 to
 753 2018; (d) Trends of annual mean particle number concentration and peak diameter; (e)
 754 Trends of annual mean of the hygroscopic parameter κ_{chem} calculated from TAP dataset
 755 in Beijing.



756

757 **Fig. 5** Spatial variations of N_{CCN} derived by the RFRM model (top) and WRF-Chem
 758 (bottom) at the sites in the studied region. (a1) Average N_{CCN} at $S=0.2\%$ in 2014
 759 predicted by the RFRM model; (a2) Average N_{CCN} at $S=0.2\%$ in 2014 by the WRF-
 760 Chem; (b1 and b2) Same as a1 and a2 but in 2018; (c1) Differences in N_{CCN} at $S=0.2\%$
 761 between the year of 2014 and 2018 predicted by the RFRM model; (c2) Same as (c1)
 762 but by the WRF-Chem; (d1) Trends of N_{CCN} at $S=0.2\%$ from 2014 to 2018 predicted by
 763 the RFRM model ; (d2) Same as (d1) but by the WRF-Chem; (e1) Change rates of N_{CCN}
 764 at $S=0.2\%$ from 2014 to 2018 predicted by the RFRM model; (e2) Same as (e1) but by
 765 the WRF-Chem.



766

767 **Fig. 6** Sensitivity of the cloud parameters and radiative forcing to CCN prediction

768 biases. (a) Dependence of the uncertainty of the cloud optical thickness (τ) on the

769 uncertainty of N_{CCN} at $S=0.2\%$ with the RFRM model; (a1) Same as (a) but by the WRF-

770 Chem; (b) Dependence of the uncertainty of the absorption coefficient ($1-\omega_0$) on the

771 uncertainty of N_{CCN} at $S=0.2\%$ with the RFRM model; (b1) Same as (b) but by the

772 the WRF-Chem; (c) Dependence of the uncertainty of the critical radius (r_c) on the

773 uncertainty of N_{CCN} at $S=0.2\%$ with the RFRM model; (c1) Same as (c) but by the WRF-

774 Chem; (d) Dependence of the uncertainty of the cloud-to-rain conversion threshold

775 function (TA) on the uncertainty of N_{CCN} at $S=0.2\%$ with the RFRM model; (d1) Same

776 as (d) but by the WRF-Chem; (e) Mean uncertainty in simulating the cloud properties

777 and (f) radiative forcing (F_c) by the RFRM model and the WRF-Chem; Black star

778 shows the mean value for the observation.