

Machine learning uncovers aerosol size information from chemistry and meteorology to quantify potential cloud-forming particles

Arshad Nair (✉ aanair@albany.edu)

State University of New York at Albany <https://orcid.org/0000-0003-2530-7757>

Fangqun Yu

University at Albany, State University of New York <https://orcid.org/0000-0001-8862-4835>

Pedro Campuzano Jost

University of Colorado, Boulder <https://orcid.org/0000-0003-3930-010X>

Paul DeMott

Colorado State University <https://orcid.org/0000-0002-3719-1889>

Ezra Levin

Handix Scientific

Jose Jimenez

University of Colorado-Boulder <https://orcid.org/0000-0001-6203-1847>

Jeff Peischl

NOAA <https://orcid.org/0000-0002-9320-7101>

Ilana Pollack

Colorado State University

Carley Fredrickson

University of Washington

Andreas Beyersdorf

California State University

Benjamin Nault

Aerodyne Research Inc <https://orcid.org/0000-0001-9464-4787>

Minsu Park

Yonsei University

Seong Soo Yum

Yonsei University

Brett Palm

University of Colorado Boulder

Lu Xu

California Institute of Technology

Ilann Bourgeois

University of Colorado-Boulder

Bruce Anderson

NASA Langley Research Center

Athanasios Nenes

EPFL

Luke Ziemba

NASA Langley Research Center

Richard H Moore

NASA Langley Research Center <https://orcid.org/0000-0003-2911-4469>

Taehyoung Lee

Hankuk University of Foreign Studies

Taehyun Park

Hankuk University of Foreign Studies

Chelsea Thompson

University of Colorado Boulder

Frank Flocke

National Center for Atmospheric Research

Lewis Huey

Georgia Institute of Technology

Michelle Kim

Everactive, Inc.

Qiaoyun Peng

University of Washington

Article

Keywords: cloud condensation nuclei, climate change prediction, aerosol composition, machine learning/artificial intelligence

Posted Date: February 16th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-244416/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Geophysical Research Letters on November 9th, 2021. See the published version at <https://doi.org/10.1029/2021GL094133>.

1 Machine learning uncovers aerosol size information
2 from chemistry and meteorology to quantify
3 potential cloud-forming particles

4 Arshad Arjunan Nair*, Fangqun Yu*, Pedro Campuzano-Jost,
Paul J. DeMott, Ezra J. T. Levin, Jose L. Jimenez, Jeff Peischl,
Ilana B. Pollack, Carley D. Fredrickson, Andreas J. Beyersdorf,
Benjamin A. Nault, Minsu Park, Seong Soo Yum, Brett B. Palm,
Lu Xu, Ilann Bourgeois, Bruce E. Anderson, Athanasios Nenes,
Luke D. Ziemba, Richard H. Moore, Taehyoung Lee, Taehyun Park,
Chelsea R. Thompson, Frank Flocke, Lewis Gregory Huey,
Michelle J. Kim & Qiaoyun Peng

5 **Abstract**

6 Cloud condensation nuclei (CCN) are mediators of aerosol–cloud interactions, which
7 contribute to the largest uncertainty in climate change prediction. Here, we present
8 a machine learning/artificial intelligence model that quantifies CCN from variables of
9 aerosol composition, atmospheric trace gases, and meteorology. Comprehensive multi-
10 campaign airborne measurements, covering varied physicochemical regimes in the tro-
11 posphere, confirm the validity of and help probe the inner workings of this machine
12 learning model: revealing for the first time that different ranges of atmospheric aerosol
13 composition and mass correspond to distinct aerosol number size distributions. Ma-
14 chine learning extracts this information, important for accurate quantification of CCN,
15 additionally from both chemistry and meteorology. This can provide a physicochemi-
16 cally explainable, computationally efficient, robust machine learning pathway in global
17 climate models that only resolve aerosol composition; potentially mitigating the un-
18 certainty of effective radiative forcing due to aerosol–cloud interactions (ERF_{aci}) and
19 improving confidence in assessment of anthropogenic contributions and climate change
20 projections.

21 Introduction

22 Atmospheric aerosol effects, particularly on cloud radiative forcing, remain the largest source
23 of uncertainty (or model diversity) in climate change prediction¹. Those aerosols capable
24 of condensing water droplets and forming clouds— cloud condensation nuclei (CCN) —
25 contribute to this uncertainty. CCN interactions with water vapor thus impact cloud micro-
26 and macrophysics, and consequently modulate cloud formation, its properties (size, number,
27 and optical), and dynamics (and that of precipitation)^{2–11}. These resultant effects conse-
28 quently impact Earth’s energy budget and influence climate and weather.

29 Obtaining agreement of CCN predictions with observations is crucial towards mitigat-
30 ing the uncertainty associated with aerosol–cloud interactions. Two factors play the largest
31 role in determining CCN (at a given water supersaturation): aerosol particle number size
32 distributions (PNSD) and aerosol chemical composition (speciation)^{12,13}. While the debate
33 continues^{14–17} as to which factor plays a larger role, the more predominant effect is arguably
34 that of PNSD due to the third order dependence on size for the solute effect that permits
35 water vapor condensation as well as the greater variability of PNSD than that of speciation,
36 except in polluted regions. However, most global climate models (GCMs) use simplified
37 prescriptions to estimate aerosol numbers or CCN from speciation while assuming a fixed
38 PNSD^{18–20}. This is due to current computational constraints, which limit the incorporation
39 into GCMs of size-resolved microphysics models with detailed treatment of processes perti-
40 nent to a more accurate representation of PNSD and hence CCN number concentrations.

41 Machine learning (ML) is a subset of Artificial Intelligence (AI) where computers are
42 trained on a large number of scenarios to acquire knowledge by statistical learning and
43 without explicit instructions. While ML has been in use (the humble and ubiquitous linear
44 regression model) for the last several decades^{21,22}, in recent years, novel techniques and rapid
45 advances in this field have led to its emergent applications in the atmospheric sciences^{23–33},
46 especially in grappling with ordinal, non-linear, complex, and massive amounts of data. It is
47 key, however, that these increasingly black-box ML/AI techniques remain grounded in reality
48 for trustworthiness and generalizability; we therefore set out to probe the inner workings of
49 our recently proposed ML model³³ for deriving CCN number concentrations.

50 Here, we present an ML model to derive CCN number concentrations ([CCN]) without ex-
51 plicit PNSD information. Comprehensive multi-campaign airborne measurements over varied
52 physicochemical regimes across the tropospheric extent demonstrate its validity. We show
53 that for accurately quantifying [CCN], both size and chemistry of the aerosol matter. More
54 importantly, we demonstrate for the first time that aerosol speciation (and other commonly
55 available atmospheric variables) contain PNSD information. This information, important for

56 accurately deriving [CCN], is successfully extracted by the ML model. This study provides a
57 robust and explainable AI pathway, without compromising on computational efficiency, for
58 GCMs to incorporate more realistic PNSD information in their simulations of cloud-forming
59 aerosols; this can have potential implications towards reducing uncertainties in the effective
60 radiative forcing of aerosol–cloud interactions.

61 Results and discussion

62 Comprehensive (global scope, tropospheric vertical extent, varied seasons, and high temporal
63 resolution) airborne measurements (detailed in the Methods and in Supplementary Fig. 1
64 and Supplementary Tables 1 and 2) of atmospheric state and composition variables provide
65 an unparalleled opportunity to evaluate the machine learning derivation of number concen-
66 trations of CCN at 0.4% supersaturation ([CCN0.4]). We present three approaches: (1)
67 LinReg: linear regression on the airborne measurements of aerosol speciation for [CCN0.4]
68 as an effective representation for current aerosol mass to number prescriptions in GCMs,
69 (2) RFRM-PM: a random forest regression model, trained on a global model of atmospheric
70 chemical composition with size-resolved microphysics (GEOS-Chem-APM), for [CCN0.4] on
71 aerosol speciation (PM_{10} , NH_4 , SO_4 , NO_3 , and OA (organic aerosol)) as a possible improve-
72 ment on LinReg, and (3) RFRM: a random forest regression model, trained on GEOS-Chem-
73 APM, for [CCN0.4] on aerosol speciation and additional variables of (*Gas-phase chemistry*)
74 [SO_2], [NO_x], and [O_3], and (*Meteorology*) temperature (T) and relative humidity (RH).
75 These models are detailed in the Methods.

76
77 **Machine learning successfully derives CCN number concentrations.** Illustrated
78 in Fig. 1 is the comparison for each of these methods with aggregated airborne campaign
79 measurements (for individual campaign comparisons see Supplementary Figs. 5–7) demon-
80 strating improved machine learning skill from LinReg (Fig. 1a) → RFRM-PM (Fig. 1b) →
81 RFRM (Fig. 1c). Fig. 1d provides the summary statistics quantifying model–observation de-
82 gree of agreement and correlation. In comparison with airborne measurements of [CCN0.4],
83 RFRM-derived values show strong agreement (%-Good, defined in Methods, of $\approx 80\%$)
84 and high correlation ($R_K \approx 0.76$). The highest density is on or around the dotted white
85 line indicating 1:1 model–observation agreement and the majority of the derived values are
86 within the corridor of good-agreement between the dashed light red and dashed light blue
87 lines. While the RFRM is overall robust, we examine the cases where it deviates from
88 airborne measurements. When these model–observation disagreements (absolute Fractional
89 Bias ($|\text{FB}| > 1$) do occur, they are rare (5.9%) and in a regime where their effect on cloud

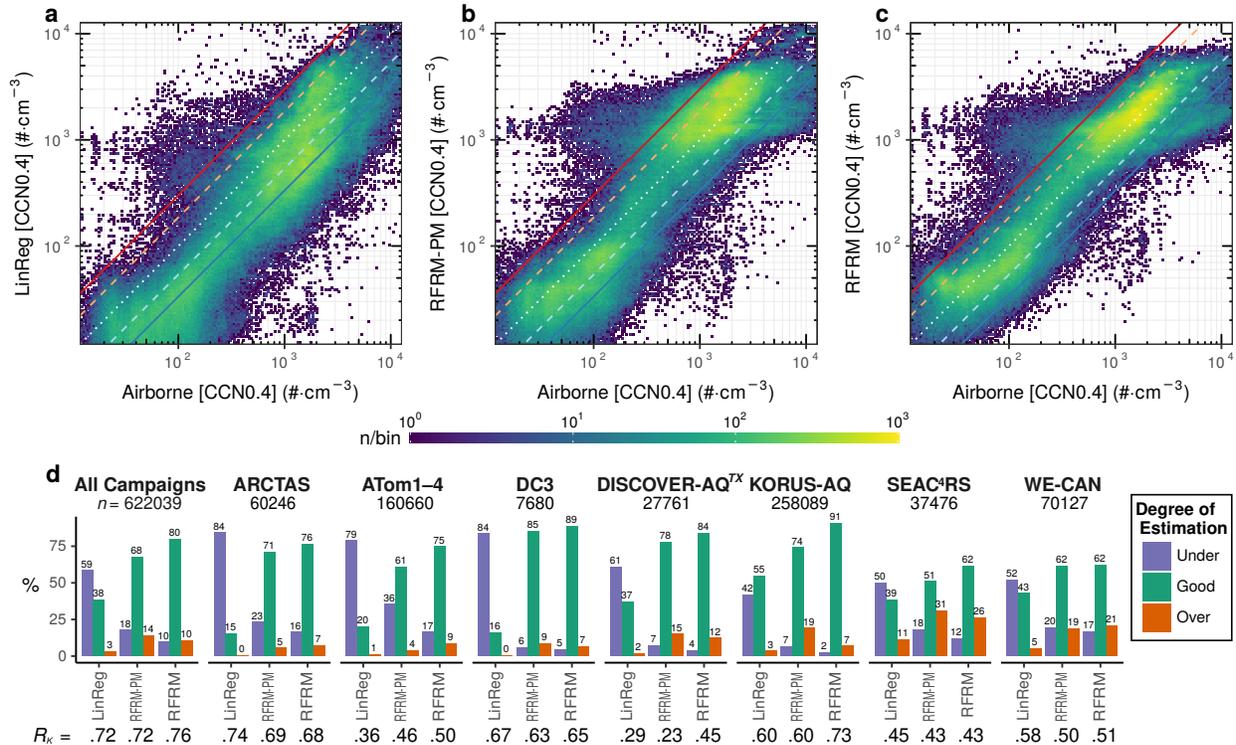


Figure 1: **Comparison of machine learning derived versus airborne measurements of [CCN0.4].** Binned scatter plot for data at the 1Hz resolution from all campaigns. For **a** Linear Regression (LinReg), **b** RFRM-PM, and **c** RFRM. Central 99% range of the airborne-measured [CCN0.4] shown for a zoomed-in view. The lines, in the order of decreasing y-intercept, indicate fractional bias (FB) of (solid red) +1, (dashed light red) +0.6, (dotted white) 0 or 1 : 1 agreement, (dashed light blue) -0.6, and (solid blue) -1, respectively. Logscale colorbar shows the count per bin. Bin-width is 0.02 (arbitrary) on the log-scale. **d** Summary statistics for the degree of model-observation agreement and correlation, as defined in the Methods, for each aircraft campaign.

90 properties will be smallest^{34,35}. For high ($> 3 \times 10^3 \text{ cm}^{-3}$) measured [CCN0.4] RFRM low
 91 bias ($\text{FB} < -1$) is largely associated with the wildfire plume measurements during the ARC-
 92 TAS and WE-CAN campaigns. It must be noted here that the low likelihood of the RFRM
 93 being exposed to these scenarios of high [CCN0.4] and predictor values in its training (on the
 94 GEOS-Chem-APM global simulations) may contribute to this observed low bias. Ultimately,
 95 however, this scenario is infrequent: ARCTAS (8.7% of its measurements), WE-CAN (8.3%),
 96 SEAC⁴RS (2.7%), and other campaigns ($\ll 0.5\%$). The high bias ($\text{FB} > +1$) of RFRM-
 97 derived [CCN0.4] occurs mainly during SEAC⁴RS (14%) and WE-CAN (7.1%). While the
 98 reason for this remains to be determined, there may be measurement uncertainties; for in-
 99 stance, in Supplementary Fig. 4a, [CCN0.4] measured directly and inferred separately are in
 100 large disagreement for SEAC⁴RS during these instances of apparent RFRM-high-bias.

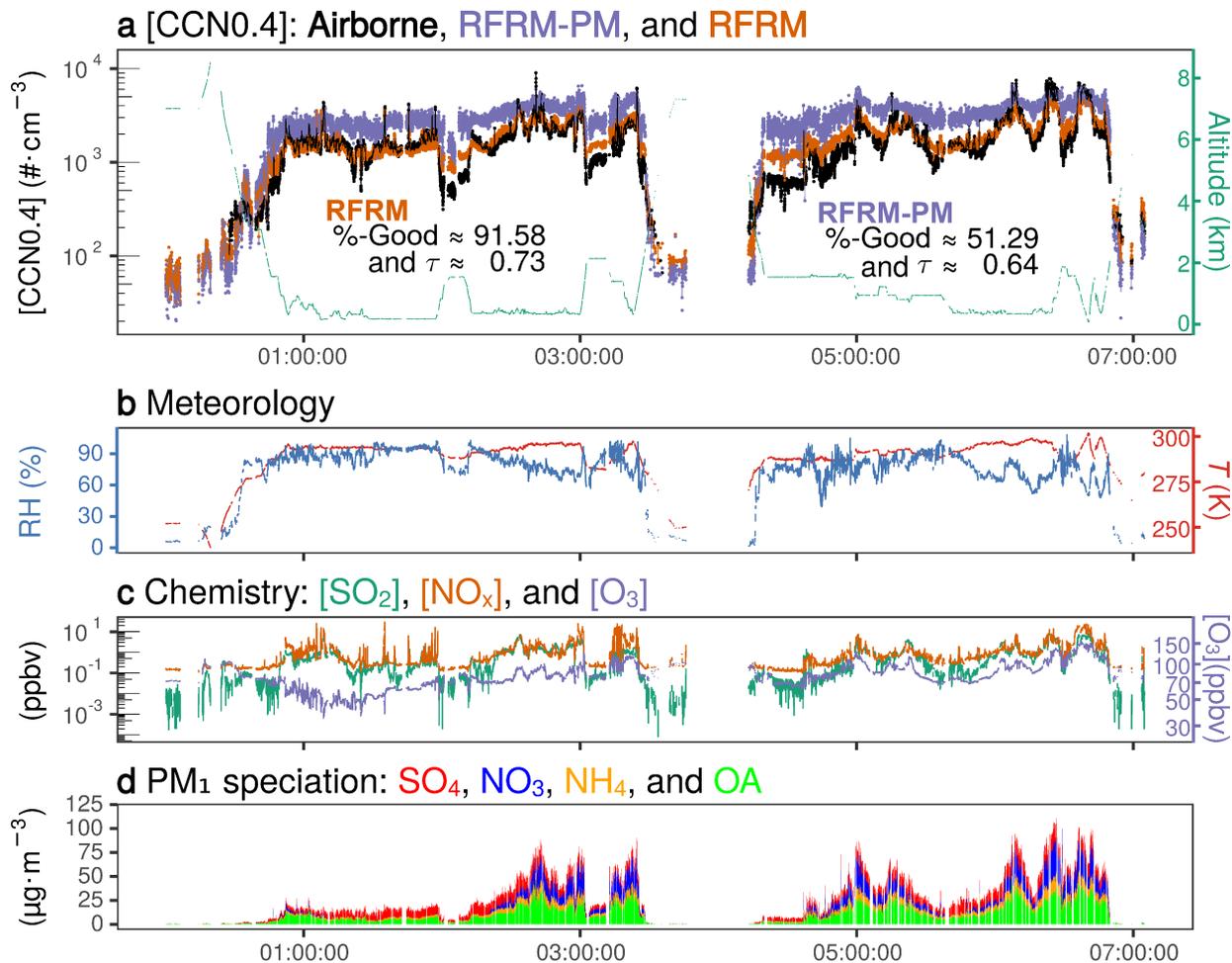


Figure 2: Time series of [CCN0.4] and variables of atmospheric state and composition shown for a selected campaign day (KORUS-AQ: 10 June 2016). **a** [CCN0.4]: (black) Airborne-measurement, (purple) RFRM-PM-derived, and (orange) RFRM-derived; and (green) altitude. **b** Meteorology: (red) temperature (T) and (blue) relative humidity (RH). **c** Chemistry: (green) [SO₂], (orange) [NO_x], and (blue) [O₃]. **d** PM₁ speciated masses of (red) SO₄, (blue) NO₃, (orange) NH₄, and (green) OA. Data is shown at the 1 Hz resolution. Solid lines associated with [CCN0.4] are 5 s rolling means.

101 While the random forest regression models demonstrate a high degree of predictive per-
 102 formance overall, we examine their performance in higher detail, leveraging the high tem-
 103 poral resolution of airborne measurements, in Fig. 2. For illustration, we select a day (10
 104 June 2016 from the KORUS-AQ campaign) with large variability in altitude (surface–8.5
 105 km) as well as the 9 predictors. Shown is the time series of the measurements of these
 106 variables during this day. Shown in Fig. 2a is the measured [CCN0.4], which varies over
 107 a range of roughly 3 orders of magnitude from 10^1 – 10^4 cm^{-3} . The 9 simultaneously mea-
 108 sured predictors (Fig. 2b–d) are used as input predictors for the RFRMs to derive [CCN0.4].

109 RFRM-PM-derived (purple) and RFRM-derived (orange) [CCN0.4] are shown in Fig. 2a.
110 Even down to the 1 Hz resolution, RFRM is able to capture [CCN0.4] variations with high
111 skill (%-Good $\approx 92\%$ and $R_K \approx 0.73$). During periods (1st, 3rd, and 6th hours) of aircraft
112 ascent and descent and the corresponding large change in magnitude of [CCN0.4], the RFRM
113 demonstrates its robustness in varying physicochemical environments. The consistency of
114 the RFRM performance across the vertical extent of the troposphere is illustrated for each
115 campaign in Supplementary Figs. 9&10. For WE-CAN (4–6 km) and ARCTAS (1–3 km),
116 the earlier noted tendency of the RFRM to underpredict [CCN0.4] is seen in the splitting
117 and skewing left of the violin distribution (Supplementary Figs. 9&10). Examining this in
118 further detail, for observations with $\text{PM}_{1\text{OA}} > 40 \mu\text{g}\cdot\text{m}^{-3}$, mean fractional bias (MFB) for
119 ARCTAS(WE-CAN) is $-1.3(-0.6)$ as compared to $-0.03(+0.2)$ when otherwise ($\text{PM}_{1\text{OA}}$
120 $\leq 40 \mu\text{g}\cdot\text{m}^{-3}$). This suggests that the RFRM-underestimation is due mostly to the high
121 organic mass (likely in biomass burning plumes) not experienced by the RFRM during its
122 training or the underestimation of the potential contribution of organic aerosol to CCN num-
123 bers in current models or a combination of these factors.

124

125 **Aerosol mass speciation contains size distribution information as revealed by ma-**
126 **chine learning.** In GCMs that do not resolve particle size distributions, proxies for aerosol
127 numbers or cloud droplet numbers are obtained from aerosol mass speciation alone, assum-
128 ing a fixed aerosol number size distribution. We have demonstrated that LinReg, directly
129 obtained from the comprehensive airborne measurements and by virtue of this overfitting is
130 an effective representation of current aerosol mass to number prescriptions in GCMs, can be
131 inadequate. A potential improvement, RFRM-PM, which employs one of the most accurate
132 machine learning approaches for regression appreciably (%-Good: $38 \rightarrow 68\%$) improves the
133 degree of agreement with CCN measurements. The importance of considering 19 predictor
134 variables of atmospheric state and composition (not limited to aerosol mass speciation) for
135 accurate RFRM-derivation of [CCN0.4] has been demonstrated³³. Considering observational
136 limitations, a maximum reduction to 9 most important predictors including T , RH, $[\text{SO}_2]$,
137 $[\text{NO}_x]$, and $[\text{O}_3]$ is possible without significant deterioration of model performance. RFRM,
138 which considers these variables in addition to only aerosol speciated mass, is in agreement
139 with measured [CCN0.4] to a much greater degree (%-Good: $38 \rightarrow 68 \rightarrow 80\%$; Figs. 1, 2, and
140 Supplementary Figs. 9&10). With the significant amount of measurement data that these
141 airborne campaigns provide, we probe into the reasons for why consideration of predictors
142 beyond $\text{PM}_{1\text{OA}}$ speciation helps improve the machine-learning model derivation of [CCN0.4].

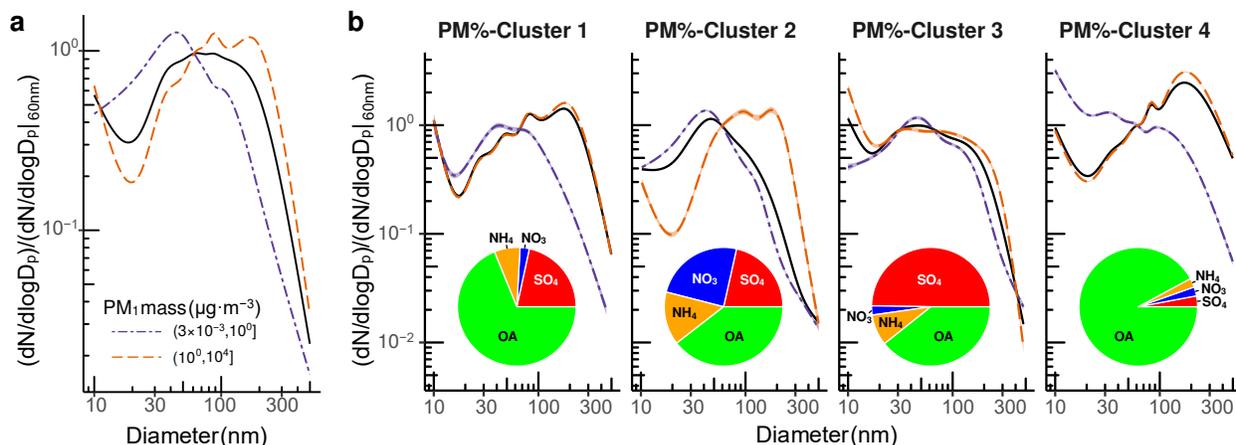


Figure 3: Aerosol mass and composition carry its number size distribution information. Average (generalized additive model) airborne measured aerosol number size distributions (PNSD) normalized to ≈ 60 nm. For (purple, dot-dashed) $PM_{tot} \leq 1 \mu\text{g cm}^{-3}$, and (orange, dashed) $PM_{tot} > 1 \mu\text{g cm}^{-3}$. Solid black curve in **a** is for all data. **b** For each cluster: Cluster 1 (SO_4 : 19–24%, OA: 66–71%, NO_3 : 0–5%, and NH_4 : 4.5–9.5%), Cluster 2 (SO_4 : 19–24%, OA: 37–42%, NO_3 : 22–27%, and NH_4 : 12–17%), Cluster 3 (SO_4 : 47.5–52.5%, OA: 37–42%, NO_3 : 0–5%, and NH_4 : 6–11%), Cluster 4 (SO_4 : 0.5–5.5%, OA: 91–96%, NO_3 : 0–5%, and NH_4 : 0–5%), and (black) respective cluster-wise average. Typical aerosol composition for each cluster is illustrated by the inset pie charts.

143 The RFRM-PM performs better than LinReg for deriving [CCN0.4] when only the PM_1
 144 speciated masses are used as input (Fig. 1). To examine the reason for this, Fig. 3 shows how
 145 the PM_1 mass contains information about the aerosol number size distribution (PNSD; P:
 146 particle/aerosol) that the random forest approach can leverage. The average normalized (to
 147 ≈ 60 nm³⁶: the rough cut-off size for CCN0.4) airborne measured PNSD is shown in Fig. 3.
 148 Fig. 3a shows that for two different total PM_1 mass ranges the PNSD profile varies. While
 149 the linear regression implicitly assumes a fixed average PNSD (black curve), the RFRM
 150 derives [CCN0.4] using decisions in the subspace corresponding to the PM_1 total mass,
 151 which defines more representative variations of PNSD. In addition, Fig. 3b demonstrates
 152 that the aerosol composition (speciated mass fractions of aerosol mass) also carries PNSD
 153 information. The four panels correspond to distinct clusters of aerosol composition, and
 154 each cluster with speciated composition of the total PM_1 mass within a range of $\pm 2.5\%$
 155 to ensure in-cluster homogeneity as well as each cluster spanning the entire range of PM_1
 156 total mass. The clusters are determined with the aid of an unsupervised machine learning
 157 technique (k -means clustering), described in the Methods and illustrated in Supplementary
 158 Figs. 12&13. Thus aerosol mass and composition confer to the RFRM-PM the ability to
 159 implicitly consider the PNSDs pertinent to PM_1 mass and speciation in its derivation of
 160 [CCN0.4] and enhance its skill compared to linear regression with an assumed mean PNSD.

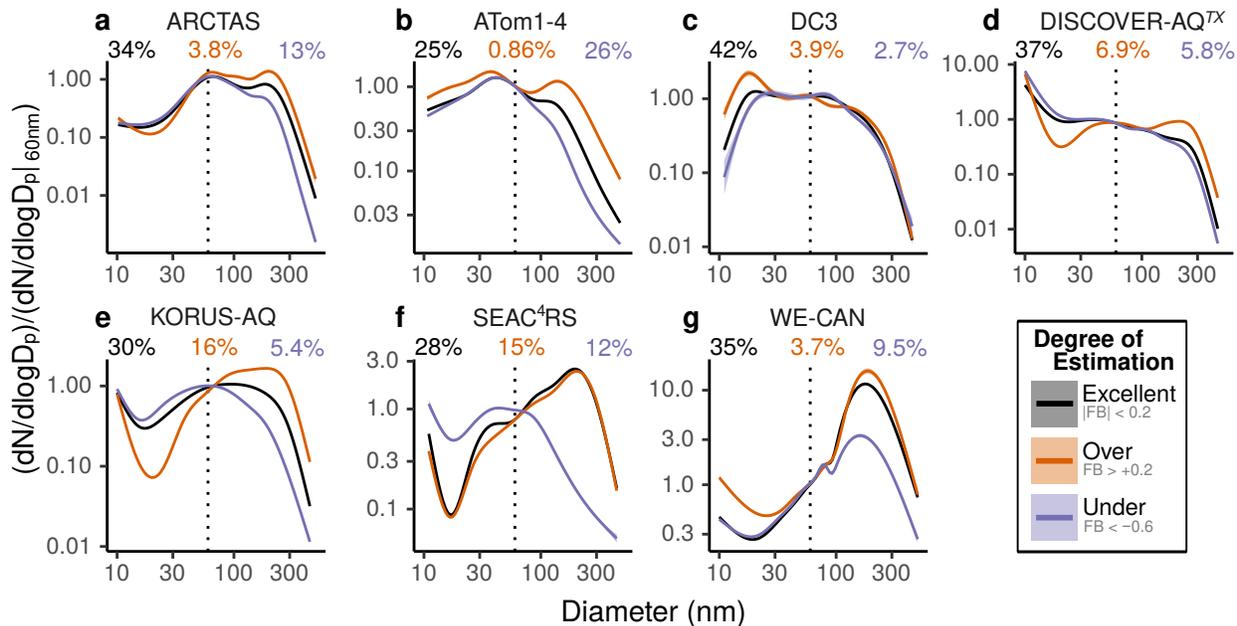


Figure 4: **Machine learning can extract aerosol number size information from chemistry and meteorology.** Average (generalized additive model) aerosol number size distributions (PNSD) normalized to ≈ 60 nm for each campaign: **a** ARCTAS, **b** ATom1–4, **c** DC3, **d** DISCOVER-AQ^{TX}, **e** KORUS-AQ, **f** SEAC⁴RS, and **g** WE-CAN. Data shown for the subset of RFRM in good-agreement and where RFRM-PM (orange) overestimates, (purple) underestimates, or is in (black) excellent agreement with airborne measurements of [CCN0.4]. Percentage of the number of observations in each class of degree of agreement shown with respectively colored text in panel sub-headings.

161 **Further size information can be machine-learned from additional chemistry and**
 162 **meteorology.** To examine why RFRM is more robust than RFRM-PM in its derivation
 163 of [CCN0.4], we consider the subset of the data where RFRM-derived [CCN0.4] is in good-
 164 agreement with airborne measurements. Counterintuitively, RFRM-PM overestimates (FB
 165 > 0.6) mostly (83.6%) when higher [CCN0.4] is measured and underestimates (FB < -0.6)
 166 mostly (82.4%) when lower [CCN0.4] is measured. This is indicative that rather than a
 167 general bias in the RFRM, it is the non-consideration of the predictors other than PM
 168 speciation contributing to the RFRM-PM bias. In Fig. 4, RFRM-PM-derived [CCN0.4]
 169 is classified into excellent-agreement ($|FB| < 0.2$; roughly 22% deviation from airborne
 170 measurement of [CCN0.4]; black), overestimation (orange), and underestimation (purple).
 171 The percentage corresponding to these classes are noted in each campaign’s panel. Illustrated
 172 are the typical PNSD normalized to the ~ 60 nm diameter, corresponding roughly to the cut-
 173 off size of CCN0.4. Across all campaigns, differences in these size distributions with respect
 174 to the degree of estimation remain consistent. More detailed differences in PNSD across

175 the vertical extent of the troposphere are also illustrated in Supplementary Fig. 11. In the
176 scenario of a more typical PNSD, with high Aitken and low accumulation mode, both RFRM
177 and RFRM-PM are in agreement with measurements. When the accumulation mode is much
178 higher and Aitken mode is much lower than average, RFRM is in agreement but RFRM-PM
179 overestimates. This is because the aerosol mass distribution towards the larger diameters
180 results in less numerous particles than a mean size distribution would suggest. When the
181 Aitken mode is much higher and the accumulation mode much lower than average, the
182 corollary follows. The additional consideration of chemical species of SO_2 , NO_x , and O_3 and
183 meteorology (T and RH), which are important for chemistry and gas-to-particle conversion
184 (including new particle formation and growth) and hence PNSD, enables RFRM to contain
185 more discerning subspaces for its decision making than RFRM-PM. With regards to the
186 PNSD, these additional predictors carry rich information about the air mass history, sources
187 of primary aerosols, and occurrence of atmospheric new particle formation and growth and
188 photochemical processing towards secondary aerosol formation. Future investigations will
189 focus on comprehensive assessment of individual contributions of each predictor variable,
190 consideration of all variables in the full-RFRM pertinent towards improved reflection of
191 the ambient PNSD, and delineation of the physicochemical processes that determine CCN
192 (spectrum) number concentrations.

193 This work demonstrates, using comprehensive airborne multi-campaign measurements
194 encompassing the varied physicochemical conditions across the troposphere, the overall suc-
195 cess of machine learning in deriving CCN number concentrations. Importantly, machine
196 learning can extract aerosol size information from aerosol composition and additionally from
197 atmospheric chemical and meteorological variables; this demonstrates that the statistical
198 learning of ML/AI algorithms is emergent from the underlying physical (and chemical) laws.
199 This physicochemically explainable and robust machine learning model can provide a com-
200 putationally efficient pathway for a more accurate representation of CCN in global climate
201 models and potentially reduce the uncertainties associated with aerosol–cloud interactions
202 in assessing anthropogenic forcing and climate change projection.

203 Methods

204 **Multi-campaign airborne measurements.** Data were collected from seven airborne
205 campaigns with measurements of the 9 predictors as well as [CCN0.4] identified in Sup-
206 plementary Table 1 and with their spatial domain shown in Supplementary Fig. 1. The 9
207 predictor variables required by the RFRM for its derivation of [CCN0.4] are measured in
208 these campaigns using instrumentation detailed in Supplementary Table 2. PNSD, here, are
209 limited to ~ 1000 nm, by the size of which aerosol numbers sharply taper off and negligibly
210 contribute to [CCN0.4]. For the ATom1–4 campaign, PNSD is measured using the Aerosol
211 Microphysical Properties (AMP) package³⁷ and for the other campaigns using a Scanning
212 Mobility Particle Sizer (SMPS; and nano-SMPS for WE-CAN) and either an Ultra-High
213 Sensitivity Aerosol Spectrometer (UHSAS: ARCTAS, DISCOVER-AQ^{TX}, and WE-CAN)
214 or a Laser Aerosol Spectrometer (LAS: DC3, KORUS-AQ, and SEAC⁴RS).

215
216 **[CCN0.4] when $ss \neq 0.4\%$.** [CCN] measurements during the ARCTAS, DC3, DISCOVER-
217 AQ^{TX}, SEAC⁴RS, and WE-CAN campaigns were typically made over a range of supersatura-
218 tions, either by changing the temperature gradient or the flow rate in the growth chamber³⁸,
219 and not necessarily for constant (0.4%) supersaturation. For these campaigns, the DMT
220 CCN-100 instrument data is reported for a range of supersaturations from 0.08–0.86%, with
221 $> 99\%$ within the range 0.1–0.7%. For the KORUS-AQ campaign, the DMT CCN-100 in-
222 strument data is reported for 0.6% supersaturation. For a 1:1 comparison of RFRM-derived
223 [CCN0.4] to measurements, we convert the [CCN] measured at supersaturations other than
224 0.4% using an empirical fitting function. This is to increase the data coverage during the
225 campaign periods. This function is determined from [CCN] measurements at the US De-
226 partment of Energy’s (DOE) Atmospheric Radiation Measurement (ARM) Southern Great
227 Plains (SGP) site located in Lamont, Oklahoma (36°36’18” N, 9°29’6” W, 318 m). ARM
228 SGP [CCN] measurements made using a DMT CCN-200^{39–41}, which is the two-column ver-
229 sion of the instrument deployed during the airborne campaigns, were used. The measurement
230 data^{42,43} is made publicly available by ARM (<https://adc.arm.gov/discovery>). The two
231 columns provide the opportunity to compare coterminous [CCN0.4] and [CCN ss], where
232 ss is some supersaturation other than 0.4%, and obtain an empirical fit for [CCN0.4] as
233 a function of [CCN ss] and ss . Supersaturations are rounded to the nearest 0.05 and the
234 ratios of [CCN ss]:[CCN0.4] (where $ss \in \{n \times 0.05\% : n \in \mathbb{Z}^+ \text{ where } 1 \leq n \leq 27\}$) are
235 plotted against ss in Supplementary Fig. 2. A polynomial fit is obtained on the median
236 values corresponding to the supersaturation ratios. Within the range of supersaturations
237 of the measurements here (0.08–0.86%), $\mathcal{O}3$ and higher terms have a negligible contribu-

238 tion. The function is equivalent to the Taylor expansion to $\mathcal{O}2$ of the functional relationship
 239 $\text{CCN}(ss)$, with no assumptions of its form (such as Twomey’s approach⁴⁴). This empirical
 240 fit function is used to approximate $[\text{CCN}0.4]$ from airborne measurements of $[\text{CCN}ss]$ where
 241 $ss \neq 0.4$. An example from the WECAN campaign where $[\text{CCN}]$ was measured across a
 242 range of supersaturations ($\{0.2,0.3,0.4,0.6\}\%$) is shown in Supplementary Fig. 3. We note
 243 that although this fit is on the median values in Supplementary Fig. 3 and the median
 244 absolute deviations shown with the error bars are not large, there exists the possibility of de-
 245 viation from the polynomial fit. It must also be noted that the empirical fit is obtained from
 246 the ground-based ARM SGP site, which may be in an atmospheric environment drastically
 247 different from where the airborne measurements are made. However, compositional effects
 248 are ameliorated to an extent by using the ratio of simultaneous $[\text{CCN}ss]:[\text{CCN}0.4]$ rather
 249 than $[\text{CCN}ss]$ alone for this fit. In the absence of sufficient direct measurement of $[\text{CCN}0.4]$,
 250 this approach, while not ideal, is presently the best possible for evaluating the performance
 251 of the machine learning models. While rigorously demonstrating the effectiveness of this
 252 approximation is currently impossible due to the airborne observational constraints, we note
 253 that there are no significant discontinuities in approximated $[\text{CCN}0.4]$ values occurring over
 254 the change in instrument supersaturation; additionally, the approximation is consistent with
 255 RFRM derived values. For ATom1–4, direct measurement of $[\text{CCN}]$ was unavailable and
 256 $[\text{CCN}0.4]$ is estimated as the number of particles with diameter greater than 60 nm based
 257 on Köhler theory³⁶. For illustration, Supplementary Fig. 4a shows this estimation applied
 258 to the other campaigns’ data. Supplementary Fig. 4b illustrates the comparison, in similar
 259 vein as Fig. 1c, of RFRM-derived $[\text{CCN}0.4]$ to these airborne-PNSD-derived $[\text{CCN}0.4]$.

260

261 **Data imputation.** To increase data coverage, if a measurement was missing and if there
 262 were measurements one second prior and/or after, it was imputed with their mean value.
 263 For DC3 $[\text{SO}_2]$ (0.1 Hz) and WE-CAN HR-ToF-AMS (0.2 Hz), measurements were assumed
 264 constant for 10 s and 5 s, respectively.

265

266 **Statistical Estimators to quantify RFRM performance.** In the present study, we
 267 use the following statistical estimators for model–observation comparison: Kendall rank
 268 correlation coefficient (R_K) to quantify correlation and **%-Good** to quantify agreement.
 269 The rationale and advantages of using these statistical metrics to evaluate model–observation
 270 comparisons are described in detail elsewhere⁴⁵. These estimators are defined as follows:

$$R_K = \tau = \frac{\sum_{i=2}^n (\text{sign}(C_i^m - C_{i-1}^m)) (\text{sign}(C_i^o - C_{i-1}^o))}{\sqrt{\binom{n}{2} - \frac{1}{2} \sum_{i=1}^n t_i^m (t_i^m - 1)} \sqrt{\binom{n}{2} - \frac{1}{2} \sum_{i=1}^n t_i^o (t_i^o - 1)}} \quad (1)$$

$$\mathbf{FB} = \frac{(C_i^m - C_i^o)}{\left(\frac{C_i^m + C_i^o}{2}\right)}; \quad \mathbf{MFB} = \frac{1}{n} \sum_{i=1}^n \frac{(C_i^m - C_i^o)}{\left(\frac{C_i^m + C_i^o}{2}\right)} \quad (2)$$

$$\% \text{-Good} = 100 \times \frac{1}{n} \sum_{i=1}^n ((|\mathbf{FB}(i)| \leq 0.6) \rightarrow 1) \quad (3)$$

271 **%-Good** is defined on the basis of the Fractional Bias (**FB**). It is the percentage of
 272 RFRM-derived [CCN0.4] with fractional bias in the range $[-0.6, +0.6]$ with respect to mea-
 273 sured [CCN0.4]. Correspondingly, **%-Over**: $\mathbf{FB} > +0.6$ and **%-Under**: $\mathbf{FB} < -0.6$.

274

275 **Machine learning models.** A Random Forest Regression Model (RFRM) to derive the
 276 number concentrations of cloud condensation nuclei at 0.4% supersaturation ([CCN0.4]) from
 277 atmospheric state and composition variables has been developed and is described in detail
 278 elsewhere³³. The present analysis focuses on [CCN0.4] for the purpose of demonstration and
 279 in future work will be extensible for the full CCN spectrum. A Random Forest⁴⁶ is a machine
 280 learning technique that can be used for regression analysis and understanding the dependence
 281 of an outcome on other variables (its predictors). This is an ensemble (to reduce overfitting)
 282 of several decision trees⁴⁷, each obtained on random subsets⁴⁸ of the training data. Here,
 283 the RFRM is trained on 30-year simulations by GEOS-Chem-APM: a state-of-the-science
 284 chemical transport model with detailed size-resolved microphysics. The present study uses
 285 the RFRM-ShortVars configuration³³, a fast implementation⁴⁹ of Random Forest models⁵⁰
 286 in the statistical computing language R⁵¹, retrained using $\text{PM}_{1.0}$ speciation variables as pre-
 287 dictors in the absence of airborne measurements of $\text{PM}_{2.5}$ speciation. Henceforth referred to
 288 as RFRM, this model derives [CCN0.4] from the following 9 commonly measured variables of
 289 atmospheric state and composition as input predictors: (*Meteorology*) temperature (T) and
 290 relative humidity (RH), (*Gas-phase chemistry*) SO_2 , NO_x , and O_3 , and (*Aerosol composition*
 291 *and mass*) NH_4 , SO_4 , NO_3 , and OA (organic aerosol).

292 To identify clusters of similar aerosol composition, k -means clustering^{52,53} is applied to
 293 partition the four dimensional $\text{PM}_{1.0}(\text{SO}_4, \text{NO}_3, \text{NH}_4, \text{OA})\%$ space into optimally separate
 294 Voronoi cells by minimization of variance within each cluster and maximization of silhouet-
 295 ting⁵⁴. To ensure more homogeneity of composition within each cluster, a tolerance of $\pm 2.5\%$
 296 around the highest density is applied. This clustering method is schematically illustrated in
 297 Supplementary Figs. 12&13. In obtaining these clusters, no cognitive bias is introduced, yet
 298 the statistical basis of clustering is in line with a chemical basis, reflecting varied atmospheric
 299 environments (biogenic/urban/biomass-burning/remote).

300 A Linear Regression model with minimization of least squares using the fast column-

301 pivoted QR decomposition method⁵⁵ is also developed on the airborne dataset. This in-
302 tentional overfitting is to obtain an effective representation of linear regression-like current
303 GCMs’ prescriptions of aerosol mass to number.

304 While, typically, machine learning models require large amounts of data for their train-
305 ing, we explore the possibility of using the comprehensive airborne measurements to develop
306 an observation-based RFRM (Obs-RFRM). The airborne measurements are split 7 : 3 for
307 training : testing. Illustrated in Supplementary Fig. 8a&b is the performance of the trained
308 Obs-RFRM-PM & Obs-RFRM, respectively, using the airborne measurement data set aside
309 for testing. Obs-RFRM using airborne measurements of the 9 predictors as input derives
310 [CCN0.4] in excellent agreement ($R_K \approx 0.92$ and %-Good $\approx 96\%$) with its airborne mea-
311 surements. This is however only an academic exercise for the sake of completeness, i.e.,
312 for equivalence of panels b&c to the overfit panel a in Fig. 1 as well as for demonstration
313 that the most of the variability of [CCN0.4] can be explained by the 9 predictors consid-
314 ered here. Supplementary Fig. 8c shows this Obs-RFRM applied to a testing dataset for
315 GEOS-Chem-APM (which RFRM and RFRM-PM have not been exposed to during their
316 training, as well). Inset is the RFRM performance with respect to GEOS-Chem-APM. This
317 demonstrates that empirical observations are, at present, still limited in number to develop
318 a practicable and generalizable machine learning model. However, a synergistic approach of
319 measurements–modeling–machine-learning has the potential to provide revealing insights.

320 Data availability

321 All datasets used in this paper are publicly available as detailed below.

322 [CCN0.18–0.86], PNSD, PM₁ composition and mass, [SO₂], [NO_x], [O₃], T , and RH mea-
323 sured during the ARCTAS⁵⁶ campaign are publicly available⁵⁷ ([https://www-air.larc.
324 nasa.gov/cgi-bin/ArcView/arctas](https://www-air.larc.nasa.gov/cgi-bin/ArcView/arctas)).

325 PNSD, PM₁ composition and mass, [SO₂], [NO_x], [O₃], T , and RH measured during the
326 ATOm1–4³⁷ campaigns are publicly available^{58–61} ([https://espo.nasa.gov/atom/archive/
327 browse/atom/DC8](https://espo.nasa.gov/atom/archive/browse/atom/DC8)).

328 [CCN0.13–0.68], PNSD, PM₁ composition and mass, [SO₂], [NO_x], [O₃], T , and RH
329 measured during the DC3⁶² campaign are publicly available⁶³ ([https://www-air.larc.
330 nasa.gov/missions/dc3-seac4rs/](https://www-air.larc.nasa.gov/missions/dc3-seac4rs/)).

331 [CCN0.14–0.60], PNSD, PM composition and mass, [SO₂], [NO_x], [O₃], T , and RH mea-
332 sured during the DISCOVER-AQ^{TX} campaign are publicly available⁶⁴ ([https://www-air.
333 larc.nasa.gov/cgi-bin/ArcView/discover-aq.tx-2013](https://www-air.larc.nasa.gov/cgi-bin/ArcView/discover-aq.tx-2013)).

334 [CCN0.6], PNSD, PM₁ composition and mass, [SO₂], [NO_x], [O₃], T , and RH measured

335 during the KORUS-AQ⁶⁵ campaign are publicly available⁶⁶ ([https://www-air.larc.nasa.](https://www-air.larc.nasa.gov/missions/korus-aq/)
336 [gov/missions/korus-aq/](https://www-air.larc.nasa.gov/missions/korus-aq/)).

337 [CCN0.09–0.56], PNSD, PM₁ composition and mass, [SO₂], [NO_x], [O₃], *T*, and RH
338 measured during the SEAC⁴RS⁶⁷ campaign are publicly available⁶⁸ ([https://www-air.](https://www-air.larc.nasa.gov/cgi-bin/ArcView/seac4rs)
339 [larc.nasa.gov/cgi-bin/ArcView/seac4rs](https://www-air.larc.nasa.gov/cgi-bin/ArcView/seac4rs)).

340 [CCN0.079–0.73], PNSD, PM₁ composition and mass, [SO₂], [NO_x], [O₃], *T*, and RH
341 measured during the WE-CAN campaign are publicly available⁶⁹ ([https://www-air.larc.](https://www-air.larc.nasa.gov/cgi-bin/ArcView/firexaq)
342 [nasa.gov/cgi-bin/ArcView/firexaq](https://www-air.larc.nasa.gov/cgi-bin/ArcView/firexaq)).

343 Dual column CCNc measurement^{42,43} data were obtained from the Atmospheric Radi-
344 ation Measurement (ARM) user facility, a U.S. Department of Energy (DOE) Office of Sci-
345 ence User Facility managed by the Biological and Environmental Research program, which
346 is publicly available at the ARM Discovery Data Portal ([https://www.archive.arm.gov/](https://www.archive.arm.gov/discovery/)
347 [discovery/](https://www.archive.arm.gov/discovery/)), last accessed on October 4 2020.

348 Data to reproduce the figures and analysis in this paper will be made publicly available
349 on Figshare upon publication.

350 Code availability

351 R scripts to reproduce the figures and analysis in this paper will be made publicly available
352 on Figshare upon publication.

353 References

- 354 1. IPCC. *Climate Change 2013: The Physical Science Basis. Contribution of Working*
355 *Group I to the Fifth Assessment Report of the Intergovernmental Panel on Climate*
356 *Change* (Cambridge University Press, Cambridge, United Kingdom and New York, NY,
357 USA, 2013).
- 358 2. Twomey, S. A. Pollution and the planetary albedo. *Atmospheric Environment (1967)*
359 **8**, 1251–1256 (1974).
- 360 3. Twomey, S. A. The influence of pollution on the shortwave albedo of clouds. *Journal of*
361 *the Atmospheric Sciences* **34**, 1149–1152 (1977).
- 362 4. Twomey, S. A., Piepgrass, M. & Wolfe, T. An assessment of the impact of pollution on
363 global cloud albedo. *Tellus B* **36B**, 356–366 (1984).

- 364 5. Albrecht, B. A. Aerosols, Cloud Microphysics, and Fractional Cloudiness. *Science* **245**,
365 1227–1230 (1989).
- 366 6. Liou, K.-N. & Ou, S.-C. The role of cloud microphysical processes in climate: An assess-
367 ment from a one-dimensional perspective. *Journal of Geophysical Research: Atmospheres*
368 **94**, 8599–8607 (1989).
- 369 7. Pincus, R. & Baker, M. B. Effect of precipitation on the albedo susceptibility of clouds
370 in the marine boundary layer. *Nature* **372**, 250–252 (1994).
- 371 8. Hansen, J., Sato, M. & Ruedy, R. Radiative forcing and climate response. *Journal of*
372 *Geophysical Research: Atmospheres* **102**, 6831–6864 (1997).
- 373 9. Ackerman, A. S. *et al.* Reduction of Tropical Cloudiness by Soot. *Science* **288**, 1042–
374 1047 (2000).
- 375 10. Ferek, R. J. *et al.* Drizzle suppression in ship tracks. *Journal of the Atmospheric Sciences*
376 **57**, 2707–2728 (2000).
- 377 11. Rosenfeld, D. Suppression of rain and snow by urban and industrial air pollution. *Science*
378 **287**, 1793–1796 (2000).
- 379 12. Junge, C. & McLaren, E. Relationship of Cloud Nuclei Spectra to Aerosol Size Distri-
380 bution and Composition. *Journal of the Atmospheric Sciences* **28**, 382–390 (1971).
- 381 13. Fitzgerald, J. W. Dependence of the supersaturation spectrum of CCN on aerosol size
382 distribution and composition. *Journal of the Atmospheric Sciences* **30**, 628–634 (1973).
- 383 14. Dusek, U. *et al.* Size matters more than chemistry for cloud-nucleating ability of aerosol
384 particles. *Science* **312**, 1375–1378 (2006).
- 385 15. Hudson, J. G. Variability of the relationship between particle size and cloud-nucleating
386 ability. *Geophysical Research Letters* **34** (2007).
- 387 16. Twohy, C. H. & Anderson, J. R. Droplet nuclei in non-precipitating clouds: composition
388 and size matter. *Environmental Research Letters* **3**, 045002 (2008).
- 389 17. Crosbie, E. *et al.* On the competition among aerosol number, size and composition in
390 predicting CCN variability: a multi-annual field study in an urbanized desert. *Atmo-*
391 *spheric Chemistry and Physics* **15**, 6943–6958 (2015).
- 392 18. Boucher, O. & Lohmann, U. The sulfate-CCN-cloud albedo effect. *Tellus B: Chemical*
393 *and Physical Meteorology* **47**, 281–300 (1995).

- 394 19. Menon, S., Genio, A. D. D., Koch, D. & Tselioudis, G. GCM simulations of the aerosol
395 indirect effect: Sensitivity to cloud parameterization and aerosol burden. *Journal of the*
396 *Atmospheric Sciences* **59**, 692–713 (2002).
- 397 20. Menon, S. & Rotstayn, L. The radiative influence of aerosol effects on liquid-phase
398 cumulus and stratiform clouds based on sensitivity studies with two climate models.
399 *Climate Dynamics* **27**, 345–356 (2006).
- 400 21. Reichstein, M. *et al.* Deep learning and process understanding for data-driven earth
401 system science. *Nature* **566**, 195–204 (2019).
- 402 22. Dramsch, J. S. 70 years of machine learning in geoscience in review. In *Machine Learning*
403 *in Geosciences*, 1–55 (Elsevier, 2020).
- 404 23. Joutsensaari, J. *et al.* Identification of new particle formation events with deep learning.
405 *Atmospheric Chemistry and Physics* **18**, 9597–9615 (2018).
- 406 24. Zaidan, M. A. *et al.* Predicting atmospheric particle formation days by Bayesian clas-
407 sification of the time series features. *Tellus B: Chemical and Physical Meteorology* **70**,
408 1–10 (2018).
- 409 25. Christopoulos, C. D., Garimella, S., Zawadowicz, M. A., Möhler, O. & Cziczo, D. J. A
410 machine learning approach to aerosol classification for single-particle mass spectrometry.
411 *Atmospheric Measurement Techniques* **11**, 5687–5699 (2018).
- 412 26. Hughes, M., Kodros, J., Pierce, J., West, M. & Riemer, N. Machine learning to predict
413 the global distribution of aerosol mixing state metrics. *Atmosphere* **9**, 15 (2018).
- 414 27. Grange, S. K., Carslaw, D. C., Lewis, A. C., Boleti, E. & Hueglin, C. Random forest me-
415 teorological normalisation models for swiss PM₁₀ trend analysis. *Atmospheric Chemistry*
416 *and Physics* **18**, 6223–6239 (2018).
- 417 28. Fuchs, J., Cermak, J. & Andersen, H. Building a cloud in the southeast Atlantic:
418 understanding low-cloud controls based on satellite observations with machine learning.
419 *Atmospheric Chemistry and Physics* **18**, 16537–16552 (2018).
- 420 29. Mauceri, S., Kindel, B., Massie, S. & Pilewskie, P. Neural network for aerosol re-
421 trieval from hyperspectral imagery. *Atmospheric Measurement Techniques* **12**, 6017–
422 6036 (2019).

- 423 30. Okamura, R., Iwabuchi, H. & Schmidt, K. S. Feasibility study of multi-pixel retrieval
424 of optical thickness and droplet effective radius of inhomogeneous clouds using deep
425 learning. *Atmospheric Measurement Techniques* **10**, 4747–4759 (2017).
- 426 31. Dou, X. & Yang, Y. Comprehensive evaluation of machine learning techniques for esti-
427 mating the responses of carbon fluxes to climatic forces in different terrestrial ecosystems.
428 *Atmosphere* **9**, 83 (2018).
- 429 32. Jin, J., Lin, H. X., Segers, A., Xie, Y. & Heemink, A. Machine learning for observation
430 bias correction with application to dust storm data assimilation. *Atmospheric Chemistry
431 and Physics* **19**, 10009–10026 (2019).
- 432 33. Nair, A. A. & Yu, F. Using machine learning to derive cloud condensation nuclei number
433 concentrations from commonly available measurements. *Atmospheric Chemistry and
434 Physics* **20**, 12853–12869 (2020).
- 435 34. Martin, G. M., Johnson, D. W. & Spice, A. The measurement and parameterization of
436 effective radius of droplets in warm stratocumulus clouds. *Journal of the Atmospheric
437 Sciences* **51**, 1823–1842 (1994).
- 438 35. Ramanathan, V. Aerosols, climate, and the hydrological cycle. *Science* **294**, 2119–2124
439 (2001).
- 440 36. Williamson, C. J. *et al.* A large source of cloud condensation nuclei from new particle
441 formation in the tropics. *Nature* **574**, 399–403 (2019).
- 442 37. Brock, C. A. *et al.* Aerosol size distributions during the Atmospheric Tomography
443 Mission (ATom): methods, uncertainties, and data products. *Atmospheric Measurement
444 Techniques* **12**, 3081–3099 (2019).
- 445 38. Moore, R. H. & Nenes, A. Scanning Flow CCN Analysis—A Method for Fast Measure-
446 ments of CCN Spectra. *Aerosol Science and Technology* **43**, 1192–1207 (2009).
- 447 39. Roberts, G. C. & Nenes, A. A Continuous-Flow Streamwise Thermal-Gradient CCN
448 Chamber for Atmospheric Measurements. *Aerosol Science and Technology* **39**, 206–221
449 (2005).
- 450 40. Uin, J. Cloud Condensation Nuclei Particle Counter (CCN) Instrument Handbook.
451 Tech. Rep., DOE Office of Science Atmospheric Radiation Measurement (ARM) Program
452 (2016).

- 453 41. Uin, J. & Smith, S. Southern Great Plains (SGP) Aerosol Observing System (AOS)
454 Instrument Handbook (2021).
- 455 42. Uin, J., Salwen, C. & Senum, G. Cloud Condensation Nuclei Particle Counter
456 (AOSCCN2COLA). 2017-04-12 to 2020-08-11, Southern Great Plains (SGP) Lamont,
457 OK (Extended and Co-located with C1) (E13) (2017). URL [https://adc.arm.gov/
458 discovery/#/results/datastream::sgpaosccn2colaE13.b1](https://adc.arm.gov/discovery/#/results/datastream::sgpaosccn2colaE13.b1).
- 459 43. Uin, J., Salwen, C. & Senum, G. Cloud Condensation Nuclei Particle Counter
460 (AOSCCN2COLB). 2017-04-12 to 2020-08-11, Southern Great Plains (SGP) Lamont,
461 OK (Extended and Co-located with C1) (E13) (2017). URL [https://adc.arm.gov/
462 discovery/#/results/datastream::sgpaosccn2colbE13.b1](https://adc.arm.gov/discovery/#/results/datastream::sgpaosccn2colbE13.b1).
- 463 44. Twomey, S. The nuclei of natural cloud formation part II: The supersaturation in natural
464 clouds and the variation of cloud droplet concentration. *Geofisica Pura e Applicata* **43**,
465 243–249 (1959).
- 466 45. Nair, A. A., Yu, F. & Luo, G. Spatioseasonal Variations of Atmospheric Ammonia
467 Concentrations Over the United States: Comprehensive Model-Observation Comparison.
468 *Journal of Geophysical Research: Atmospheres* **124**, 6571–6582 (2019).
- 469 46. Breiman, L. Random forests. *Machine Learning* **45**, 5–32 (2001).
- 470 47. Breiman, L., Friedman, J. H., Olshen, R. A. & Stone, C. J. *Classification And Regression*
471 *Trees* (Routledge, 1984).
- 472 48. Breiman, L. Bagging predictors. *Machine Learning* **24**, 123–140 (1996).
- 473 49. Wright, M. N. & Ziegler, A. ranger: A fast implementation of random forests for high
474 dimensional data in c++ and r. *Journal of Statistical Software* **77** (2017).
- 475 50. Breiman, L. *Manual for Setting Up, Using, and Understanding Random Forest V4.0*
476 (2003). [https://www.stat.berkeley.edu/~breiman/Using_random_forests_v4.0.
477 pdf](https://www.stat.berkeley.edu/~breiman/Using_random_forests_v4.0.pdf).
- 478 51. R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation
479 for Statistical Computing, Vienna, Austria (2020). URL <https://www.R-project.org/>.
- 480 52. Lloyd, S. Least squares quantization in PCM. *IEEE Transactions on Information Theory*
481 **28**, 129–137 (1982).

- 482 53. Hartigan, J. A. & Wong, M. A. Algorithm AS 136: A k-means clustering algorithm.
483 *Applied Statistics* **28**, 100 (1979).
- 484 54. Kaufman, L. & Rousseeuw, P. J. (eds.) *Finding Groups in Data* (John Wiley & Sons,
485 Inc., 1990).
- 486 55. Bates, D. & Eddelbuettel, D. Fast and Elegant Numerical Linear Algebra Using the
487 RcppEigenPackage. *Journal of Statistical Software* **52** (2013).
- 488 56. Jacob, D. J. *et al.* The Arctic Research of the Composition of the Troposphere from Air-
489 craft and Satellites (ARCTAS) mission: design, execution, and first results. *Atmospheric*
490 *Chemistry and Physics* **10**, 5191–5212 (2010).
- 491 57. ARCTAS Team. ARCTAS Field Campaign Data (2020). URL [https://www-air.larc.](https://www-air.larc.nasa.gov/cgi-bin/ArcView/arctas)
492 [nasa.gov/cgi-bin/ArcView/arctas](https://www-air.larc.nasa.gov/cgi-bin/ArcView/arctas).
- 493 58. Allen, H. M., Crouse, J. D., Kim, M. J., Teng, A. P. & Wennberg, P. O. ATom: L2
494 In Situ Data from Caltech Chemical Ionization Mass Spectrometer (CIT-CIMS) (2019).
495 URL <https://espo.nasa.gov/atom/archive/browse/atom/DC8/CIT-S02>.
- 496 59. Ryerson, T. B., Thompson, C. R., Peischl, J. & Bourgeois, I. ATom: L2 In Situ Mea-
497 surements from NOAA Nitrogen Oxides and Ozone (NOyO3) Instrument (2019). URL
498 <https://espo.nasa.gov/atom/archive/browse/atom/DC8/>.
- 499 60. Jimenez, J. L. *et al.* ATom: L2 Measurements from CU High-Resolution Aerosol
500 Mass Spectrometer (HR-AMS) (2019). URL [https://espo.nasa.gov/atom/archive/](https://espo.nasa.gov/atom/archive/browse/atom/DC8/AMS)
501 [browse/atom/DC8/AMS](https://espo.nasa.gov/atom/archive/browse/atom/DC8/AMS).
- 502 61. Brock, C. A. *et al.* ATom: L2 In Situ Measurements of Aerosol Microphysical Proper-
503 ties (AMP) (2019). URL [https://espo.nasa.gov/atom/archive/browse/atom/DC8/](https://espo.nasa.gov/atom/archive/browse/atom/DC8/SDAerosol)
504 [SDAerosol](https://espo.nasa.gov/atom/archive/browse/atom/DC8/SDAerosol).
- 505 62. Barth, M. C. *et al.* The deep convective clouds and chemistry (DC3) field campaign.
506 *Bulletin of the American Meteorological Society* **96**, 1281–1309 (2015).
- 507 63. DC3 Team. DC3 Field Campaign Data (2013). URL [https://www-air.larc.nasa.](https://www-air.larc.nasa.gov/missions/dc3-seac4rs/)
508 [gov/missions/dc3-seac4rs/](https://www-air.larc.nasa.gov/missions/dc3-seac4rs/).
- 509 64. DISCOVER-AQ Team. DISCOVER-AQ Field Campaign Data (2014). URL [https:](https://www-air.larc.nasa.gov/cgi-bin/ArcView/discover-aq.tx-2013)
510 [//www-air.larc.nasa.gov/cgi-bin/ArcView/discover-aq.tx-2013](https://www-air.larc.nasa.gov/cgi-bin/ArcView/discover-aq.tx-2013).

- 511 65. Jordan, C. E. *et al.* Investigation of factors controlling PM_{2.5} variability across the south
512 korean peninsula during KORUS-AQ. *Elementa: Science of the Anthropocene* **8** (2020).
- 513 66. KORUS-AQ Team. KORUS-AQ Field Campaign Data (2018). URL [https://www-air.
514 larc.nasa.gov/missions/korus-aq/](https://www-air.larc.nasa.gov/missions/korus-aq/).
- 515 67. Toon, O. B. *et al.* Planning, implementation, and scientific goals of the studies of
516 emissions and atmospheric composition, clouds and climate coupling by regional surveys
517 (SEAC4RS) field mission. *Journal of Geophysical Research: Atmospheres* **121**, 4967–
518 5009 (2016).
- 519 68. SEAC4RS Team. SEAC4RS Field Campaign Data (2014). URL [https://www-air.
520 larc.nasa.gov/cgi-bin/ArcView/seac4rs](https://www-air.larc.nasa.gov/cgi-bin/ArcView/seac4rs).
- 521 69. WE-CAN Team. WE-CAN Field Campaign Data (2019). URL [https://www-air.
522 larc.nasa.gov/cgi-bin/ArcView/firexaq?MERGE=1](https://www-air.larc.nasa.gov/cgi-bin/ArcView/firexaq?MERGE=1).

523 Acknowledgments

524 We are grateful to Chuck Brock (NOAA) for the in situ measurements of aerosol microphysi-
525 cal properties during the ATom1–4 campaigns, Andrew J. Weinheimer, Denise D. Montzka &
526 David J. Knapp (ARCTAS, DISCOVER-AQ^{TX}, KORUS-AQ, and WE-CAN) and Thomas
527 B. Ryerson (ATom1–4, DC3, and SEAC⁴RS) for [NO_x] and [O₃] measurements, Paul O.
528 Wennberg, John D. Crouse & Hannah M. Allen for [SO₂] measurements during the ARC-
529 TAS and ATom1–4 campaigns, Kevin R. Barry (funded by NSF grant no. AGS-1660486;
530 WE-CAN: CCNc and SMPS), Sonia M. Kreidenweis (WE-CAN: CCNc and HR-ToF-AMS),
531 Kathryn A. Moore (funded by the NSF Graduate Research Fellowship grant no. 006784;
532 WE-CAN: CCNc), Darin W. Toohey & Michael Reeves (WE-CAN: UHSAS), Lauren A.
533 Garofalo & Delphine K. Farmer (funded by the NOAA grant no. NA17OAR4310010; WE-
534 CAN: HR-ToF-AMS), and Joel A. Thornton (WE-CAN: CIMS [SO₂] measurements). We are
535 thankful to Michael Shook & Gao Chen at the NASA Langley Research Center Airborne Sci-
536 ence Data for Atmospheric Composition (<https://www-air.larc.nasa.gov/index.html>)
537 for data curation. We also thank the DOE ARM SGP Research Facility data teams for the
538 operation and maintenance of instruments, quality checks, and making their measurement
539 data publicly available.

540 **Funding**

541 This research has been supported by the National Aeronautics and Space Administration
542 (NASA grant no. NNX17AG35G) and the National Science Foundation (NSF grant no. AGS-
543 1550816). B.A.N., P.C.J. & J.L.J. were supported by NASA (grant nos. NNX15AJ23G,
544 NNX15AH33A, 80NSSC19K0124, and 80NSSC18K0630). P.J.D. & E.J.T.L. acknowledge
545 support from the NSF (grant no. AGS-1650786). Any opinions, findings, and conclusions or
546 recommendations expressed in this material are those of the authors and do not necessarily
547 reflect the views of the NSF. M.P. & S.S.Y. acknowledge the support from the National
548 Research Foundation of Korea (NRF) funded by the Korean government (MSIT) (grant
549 no. NRF-2018R1A2B2006965). M.J.K. was supported by an NSF Atmospheric and Geospace
550 Sciences Postdoctoral Research Fellowship (AGS-PRF; grant no. 1524860). C.D.F., B.B.P. &
551 Q.P. acknowledge support from the NSF (grant no. AGS-1652688) and the National Oceanic
552 and Atmospheric Administration (NOAA grant no. NA17OAR4310012).

553 **Author information**

554 **Benjamin A. Nault**

555 Present address: Center for Aerosols and Cloud Chemistry, Aerodyne Research, Inc., Biller-
556 ica, Massachusetts 01821, USA

558 **Ezra J. T. Levin**

559 Present address: Handix Scientific, Boulder, Colorado 80301, USA

560 **Affiliations**

561 **Atmospheric Sciences Research Center, State University of New York, Albany,**
562 **New York 12203, USA**

563 Arshad Arjunan Nair & Fangqun Yu

564 **Cooperative Institute for Research in Environmental Sciences (CIRES), Univer-**
565 **sity of Colorado Boulder, Boulder, Colorado 80309, USA**

566 Pedro Campuzano-Jost, Jose L. Jimenez, Jeff Peischl, Benjamin A. Nault, Chelsea R.
567 Thompson & Ilann Bourgeois

568 **Department of Chemistry, University of Colorado Boulder, Boulder, Colorado**
569 **80309, USA**

570 Pedro Campuzano-Jost, Jose L. Jimenez & Benjamin A. Nault

571 **Department of Atmospheric Science, Colorado State University, Fort Collins,**
572 **Colorado 80523, USA**
573 Paul J. DeMott, Ezra J. T. Levin & Ilana B. Pollack

574 **NOAA Chemical Science Laboratory, Boulder, Colorado 80305, USA**
575 Jeff Peischl, Ilann Bourgeois, Chelsea R. Thompson, Thomas B. Ryerson

576 **Department of Atmospheric Sciences, University of Washington, Seattle, Wash-**
577 **ington 98195, USA**
578 Carley D. Fredrickson, Brett B. Palm & Qiaoyun Peng

579 **NASA Langley Research Center, Hampton, Virginia 23666, USA**
580 Andreas J. Beyersdorf, Bruce E. Anderson, Luke D. Ziemba & Richard H. Moore

581 **California State University, San Bernardino, California 92407, USA**
582 Andreas J. Beyersdorf

583 **Department of Atmospheric Sciences, Yonsei University, Seoul 03722, South**
584 **Korea**
585 Minsu Park & Seong Soo Yum

586 **Division of Geological and Planetary Sciences, California Institute of Technology,**
587 **Pasadena, California 91125, USA**
588 Lu Xu, & Michelle J. Kim

589 **Institute of Chemical Engineering Sciences, Foundation for Research & Technology-**
590 **Hellas, Patras 26504, Greece**
591 Athanasios Nenes

592 **School of Architecture, Civil and Environmental Engineering, Swiss Federal In-**
593 **stitute of Technology Lausanne, Lausanne 1015, Switzerland**
594 Athanasios Nenes

595 **School of Earth and Atmospheric Sciences and School of Chemical and Biomolec-**
596 **ular Engineering, Georgia Institute of Technology, Atlanta, Georgia 30332, USA**
597 Athanasios Nenes

598 **Department of Environmental Science, Hankuk University of Foreign Studies,**
599 **Yongin 17035, South Korea**
600 Taehyoung Lee & Taehyun Park

601 **Atmospheric Chemistry Observations and Modeling, National Center for Atmo-**
602 **spheric Research, Boulder, Colorado 80307, USA**
603 Frank Flocke

604 **School of Earth and Atmospheric Sciences, Georgia Tech, Atlanta, Georgia**
605 **30332, USA**

606 Lewis Gregory Huey

607 **Contributions**

608 A.A.N. and F.Y. conceptualized this study. B.E.A., A.N., A.J.B., L.D.Z., R.H.M., M.P.,
609 S.S.Y., E.J.T.L., P.J.D., P.C.J., B.A.N., B.B.P., J.L.J., T.L., T.P., C.R.T., F.F., I.B.P., I.B.,
610 J.P., L.X., M.J.K., L.G.H., C.D.F., and Q.P. carried out, analyzed, and archived airborne
611 measurements. A.A.N. led the analysis and writing with major input from F.Y. and further
612 input from P.C.J., P.J.D., E.J.T.L., J.L.J., J.P., I.B.P., C.D.F., A.J.B., B.A.N., M.P., S.S.Y.,
613 B.B.P., L.X., and I.B..

614 **Corresponding authors**

615 Correspondence to Arshad Arjunan Nair (aanair@albany.edu) and Fangqun Yu (fyu@albany.edu).

616 **Ethics declarations**

617 **Competing interests**

618 The authors declare no competing interests.

Figures

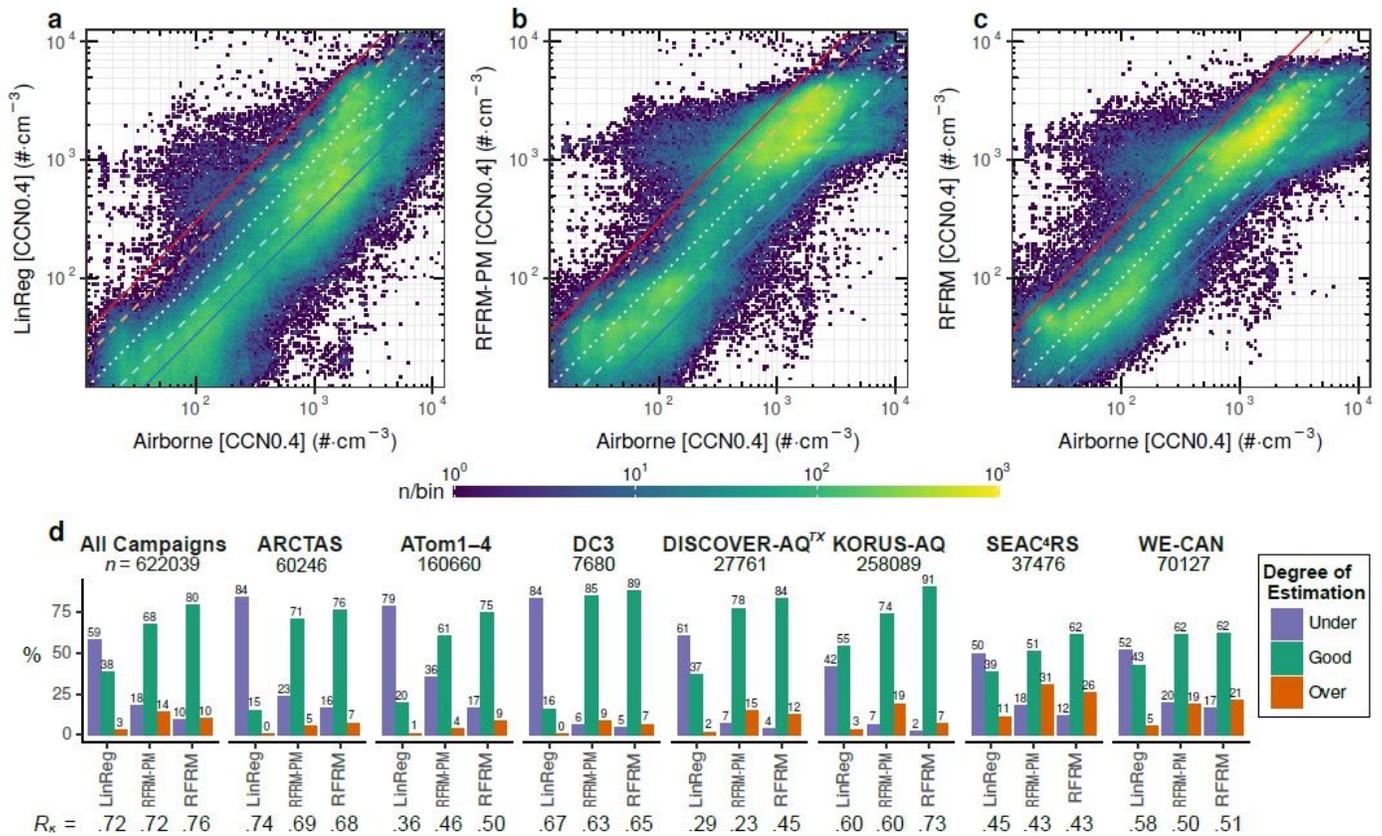


Figure 1

Comparison of machine learning derived versus airborne measurements of [CCN0.4]. Binned scatter plot for data at the 1Hz resolution from all campaigns. For a Linear Regression (LinReg), b RFRM-PM, and c RFRM. Central 99% range of the airborne measured [CCN0.4] shown for a zoomed-in view. The lines, in the order of decreasing yintercept, indicate fractional bias (FB) of (solid red) +1, (dashed light red) +0:6, (dotted white) 0 or 1 : 1 agreement, (dashed light blue) -0:6, and (solid blue) -1, respectively. d Summary statistics for the degree of model/observation agreement and correlation, as defined in the Methods, for each aircraft campaign.

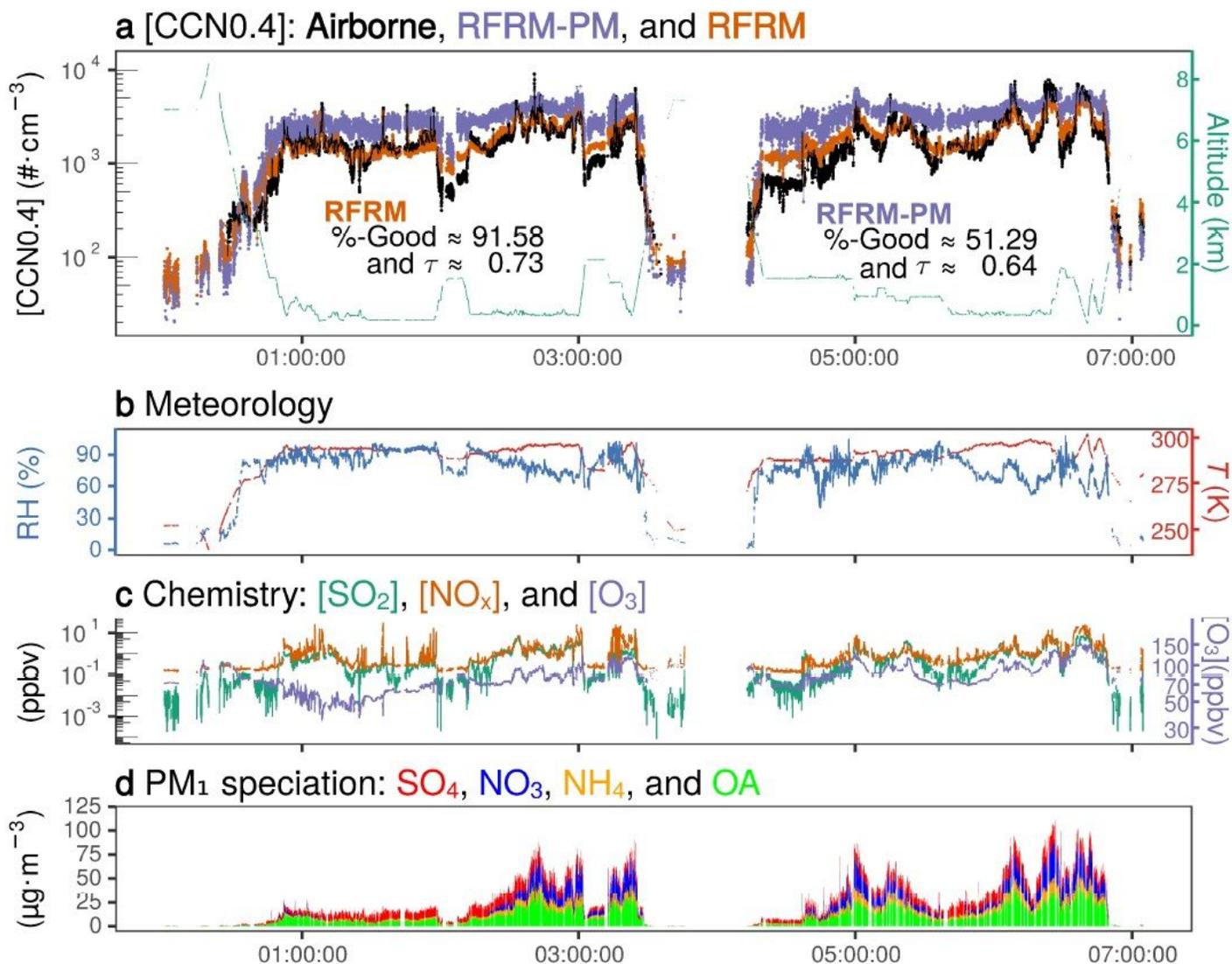


Figure 2

Time series of [CCN0.4] and variables of atmospheric state and composition shown for a selected campaign day (KORUS-AQ: 10 June 2016). a [CCN0.4]: (black) Airborne-measurement, (purple) RFRM-PM-derived, and (orange) RFRM-derived; and (green) altitude. b Meteorology: (red) temperature (T) and (blue) relative humidity (RH). c Chemistry: (green) [SO₂], (orange) [NO_x], and (blue) [O₃]. d PM₁ speciated masses of (red) SO₄, (blue) NO₃, (orange) NH₄, and (green) OA. Data is shown at the 1 Hz resolution. Solid lines associated with [CCN0.4] are 5 s rolling means.

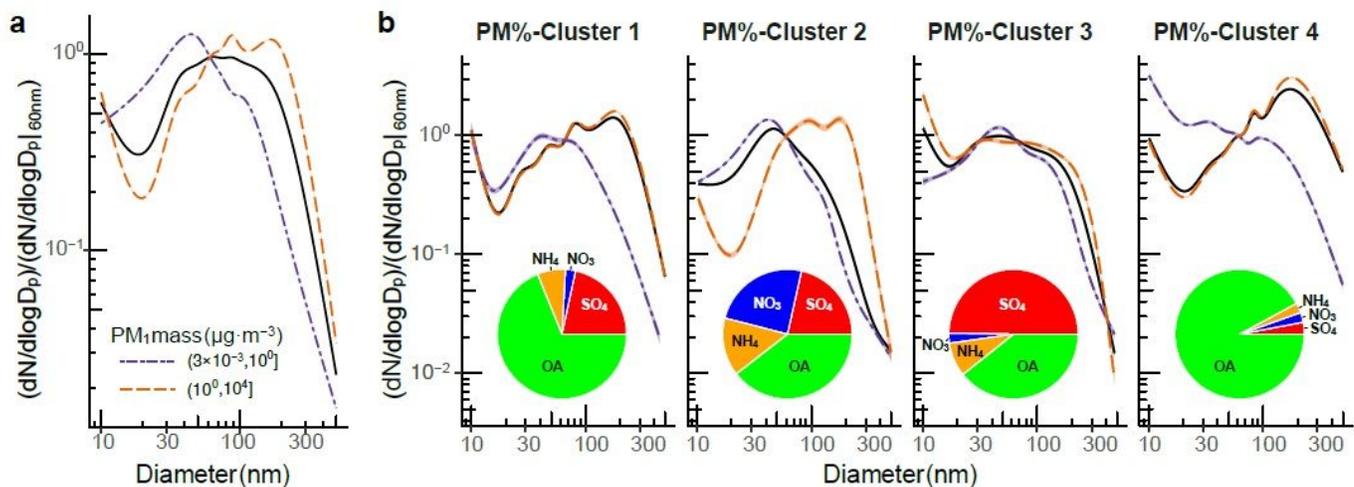


Figure 3

Aerosol mass and composition carry its number size distribution information. Average (generalized additive model) airborne measured aerosol number size distributions (PNSD) normalized to -60 nm. For (purple, dot-dashed) PM_{tot} < 1 $\mu\text{g}\cdot\text{cm}^{-3}$, and (orange, dashed) PM_{tot} > 1 $\mu\text{g}\cdot\text{cm}^{-3}$. Solid black curve in a is for all data. b For each cluster: Cluster 1 (SO₄: 19{24%, OA: 66{71%, NO₃: 0{5%, and NH₄: 4.5{9.5%), Cluster 2 (SO₄: 19{24%, OA: 37{42%, NO₃: 22{27%, and NH₄: 12{17%), Cluster 3 (SO₄: 47.5{ 52.5%, OA: 37{42%, NO₃: 0{5%, and NH₄: 6{11%), Cluster 4 (SO₄: 0.5{5.5%, OA: 91{96%, NO₃: 0{5%, and NH₄: 0{5%), and (black) respective cluster-wise average. Typical aerosol composition for each cluster is illustrated by the inset pie charts.

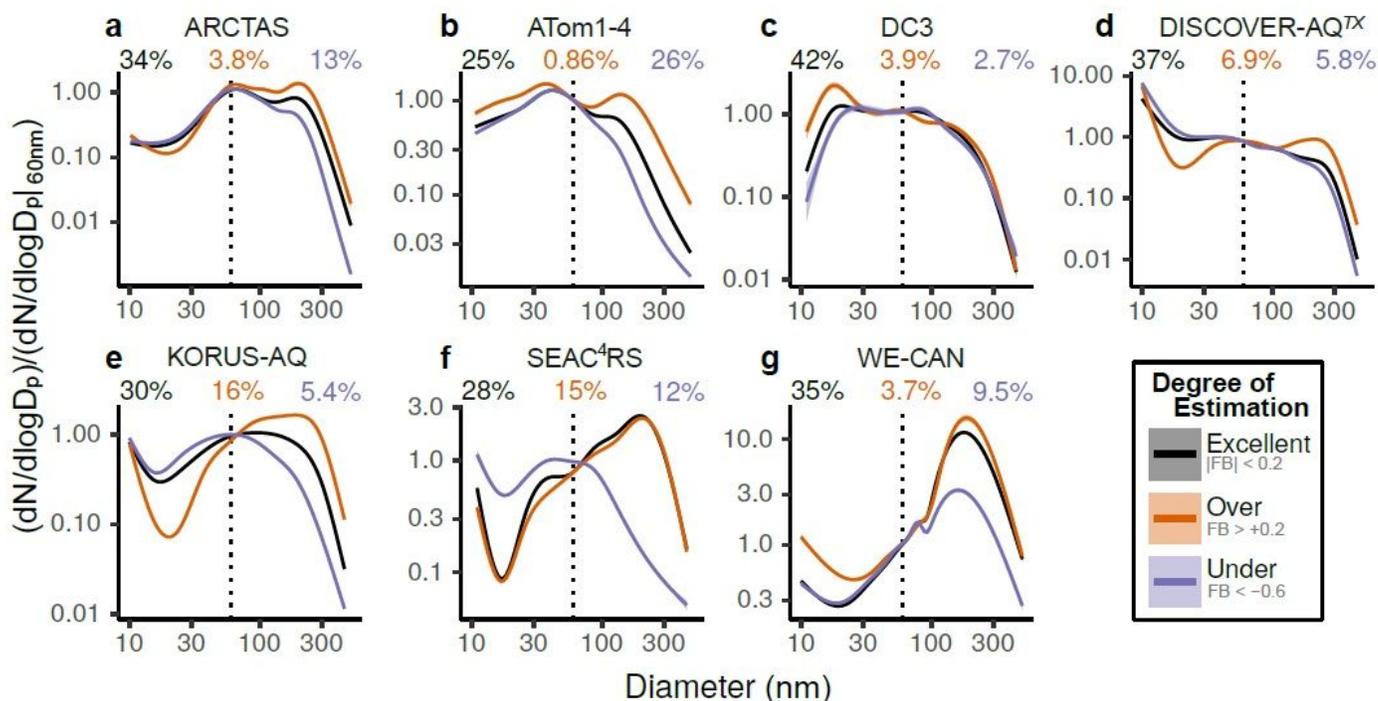


Figure 4

Machine learning can extract aerosol number size information from chemistry and meteorology. Average (generalized additive model) aerosol number size distributions (PNSD) normalized to ~60 nm for each campaign: a ARCTAS, b ATom1{4, c DC3, d DISCOVER-AQTX, e KORUS-AQ, f SEAC4RS, and g WE-CAN. Data shown for the subset of RFRM in good-agreement and where RFRM-PM (orange) overestimates, (purple) underestimates, or is in (black) excellent agreement with airborne measurements of [CCN0.4]. Percentage of the number of observations in each class of degree of agreement shown with respectively colored text in panel sub-headings.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [supplementary.pdf](#)