



1 **Flow cytometry and machine learning enable identification of allergenic urban tree pollen.**

2 Authors : Sarah Tardif^{1,2}, Maria Raquel Kanieski⁴, Gauthier Lapa^{1,2}, Grégoire Bonnamour^{1,3}, Rita Sousa-
3 Silva⁵, Isabelle Laforest-Lapointe^{2,6}, Alain Paquette^{1,2}

4

5 ¹ Département des sciences biologiques, Université du Québec à Montréal, Montréal, QC, Canada.

6 ² Centre for Forest Research, Université du Québec à Montréal, Montréal, QC, Canada.

7 ³Centre d'excellence en recherche sur les maladies orphelines-Fondation Courtois (CERMO-FC),
8 Université du Québec à Montréal, Montréal, QC, Canada.

9 ⁴ Universidade do Estado de Santa Catarina, Depto. Engenharia Florestal, Lages, SC, Brasil.

10 ⁵ Institute of Environmental Sciences, Department of Environmental Biology, Leiden University, Leiden,
11 The Netherlands.

12 ⁶ Département de Biologie, Université de Sherbrooke, Sherbrooke, QC, Canada.

13 *Correspondence to:* Sarah Tardif (sarahtardif02@gmail.com)

14 **Abstract**

15 Exposure to allergenic pollen is a major public health concern, as it is a key trigger for respiratory allergies,
16 including seasonal allergic rhinitis, which affects approximately 20% of the global population. Monitoring
17 airborne pollen is essential for prevention and clinical management, yet traditional identification methods,
18 such as light microscopy, are time-consuming and often limited to genus- or family-level resolution. Here,
19 we present a high-throughput approach combining flow cytometry with machine learning to identify pollen
20 from urban environments. We collected a reference database of pollen from 97 species across 34 genera,
21 representing the dominant allergenic trees and other common airborne taxa in Montreal, Canada. Using flow
22 cytometry, we measured particle size, granularity, and fluorescence intensity across multiple excitation and
23 emission channels, and applied a Random Forest classifier to distinguish pollen taxa. At the species level,
24 the model achieved a mean F_1 -score of 0.76, while genus-level classification reached 0.90, with
25 misclassifications largely occurring among closely related species. Granularity and fluorescence parameters
26 from the violet and blue lasers were the most distinctive features. Our results demonstrate that flow
27 cytometry combined with machine learning provides an efficient, scalable alternative to microscopy, with
28 potential for large-scale urban pollen monitoring.



29 1 Introduction

30 Exposure to allergenic pollen is a major public health concern, as it is a key risk factor for respiratory
31 allergies. Seasonal allergic rhinitis affects approximately 20 % of the global population (Savouré et al.,
32 2022) and is expected to worsen with climate change, which is projected to lengthen pollen seasons
33 (Anderegg et al., 2021; Mousavi et al., 2024; Zhang and Steiner, 2022; Ziska et al., 2019). Rising
34 temperatures and CO₂ levels stimulate plant growth, increasing pollen levels (Kim et al., 2018; Ladeau and
35 Clark, 2006) and the allergenicity of pollen grains (Ahlholm et al., 1998; Kim et al., 2018). For allergy
36 sufferers and healthcare providers, reliable pollen information, including which plant species and pollen
37 traits contribute to different allergenicity properties, is essential for prevention and effective treatment, but
38 remains scarce (Dunker et al., 2022; Medek et al., 2025; Sousa-Silva et al., 2020).

39 Expanding pollen monitoring networks in urban areas, which host most of the world's population, is
40 increasingly recognized as essential (Tummon et al., 2024), yet this also requires processing a large number
41 of pollen samples and thus highlights a clear need for efficient and accurate identification methods. Over
42 the past decades, several analytical techniques have been developed for pollen detection and classification,
43 each having advantages and limitations. Light microscopy remains the standard method used worldwide for
44 pollen identification, but it is time-consuming and requires highly trained specialists (Brennan et al., 2019;
45 Dunker et al., 2021, 2022; Gierlicka et al., 2022; de Weger et al., 2013). Although pollen morphology,
46 defined by size, shape, apertures, and texture (Ogden et al., 1974; Smith, 1984), supports taxonomic
47 identification, subtle interspecific differences restrict identification to genus or family level in most cases.
48 Automated slide scanning, sometimes coupled with a machine learning algorithm, has improved efficiency
49 but still faces limitations in distinguishing species from the same genus or family (Dunker et al., 2021; Holt
50 and Bennett, 2014). Advanced imaging techniques, such as scanning electron microscopy (SEM),
51 transmission electron microscopy (TEM), and optical diffraction tomography (ODT), provide much higher
52 resolution for detailed analysis of pollen structures, but are costly or impractical for large-scale monitoring
53 (Gierlicka et al., 2022). Molecular biology techniques, particularly metabarcoding and PCR-based methods,
54 have the potential to enable species-level identification yet face challenges such as high costs, the presence
55 of DNA inhibitors that can limit sensitivity and cause false negative, the limitations of taxonomic resolution,
56 and the inability to quantify pollen abundance (Dunker et al., 2021; Gierlicka et al., 2022).

57 More recently, fluorescence spectroscopy and flow cytometry have emerged as promising approaches
58 (Gierlicka et al., 2022; Šaulienė et al., 2019). These methods are based on the size and autofluorescence
59 properties of particles, such as the pollen grains, and when combined with holographic images and machine
60 or deep learning, they can improve classification accuracy and enable automated (Dunker et al., 2022; Erb



et al., 2024; Sikoparija et al., 2024; Swanson et al., 2023) and high-throughput identification (≈ 5000 grains s^{-1}) (Dunker et al., 2021; Gierlicka et al., 2022). Because each species has a specific fluorescence and granularity signature, it is possible to distinguish even morphologically similar taxa (Dunker et al., 2021).

Our study aims to develop a classification model capable of identifying airborne pollen in urban environments. We built a reference collection representing the main tree species found across the city of Montreal, Canada. Unlike previous studies that rely on microscopic or imaging data, our approach relies exclusively on flow cytometry measurements, i.e. fluorescence intensity, particle size, and granularity to characterize pollen. This choice is motivated by the fact that most cytometers routinely used in healthcare and clinical settings are limited to these parameters. Consequently, developing a model based on these features enhances its applicability and ensures compatibility with the most widely implemented cytometry platforms. We then evaluated the performance of the machine-learning classification model trained on these flow cytometry parameters and identified those that contribute most to differentiating pollen species and genera.

2 Methodology

2.1 Pollen collection

To train the machine learning classification model, we created a reference database of pollen grains collected directly from plants of known species (mostly trees). The reference collection included pollen from both common urban tree species as well as widely planted hybrid cultivars.

Tree species were selected based on three criteria: (1) their relative abundance on the Island of Montreal, ensuring representation of the dominant urban taxa; (2) their anemophilous nature, since wind-pollinated species are typically the most allergenic (D'Amato et al., 2007; Falagiani, 1989); and (3) the inclusion of multiple species within each genus, to enable species-level discrimination were possible. Other species such as from the Rosaceae family were also included to increase resolution. For each selected species, pollen was collected from three individual trees from the Montreal Botanical Garden (for ease of identification) or among public trees across the city. At flowering time, ten floral units (flowers, catkins or male cones) were collected per tree, sampling different parts of the crown to capture intra-individual variation among pollen grains. We also included pollen from the *Poaceae* family (grasses) and the genus *Ambrosia* (ragweed), given their well-known allergenic potential (D'Amato et al., 2007; Falagiani, 1989). Their inclusion enabled the model to learn to discriminate tree pollen from other common airborne particle types, as real-world environmental samples typically comprise a heterogeneous mix of tree, grass, and weed pollen, along with various non-pollen particulates. In the laboratory, floral units were placed in pre-labelled paper bags with

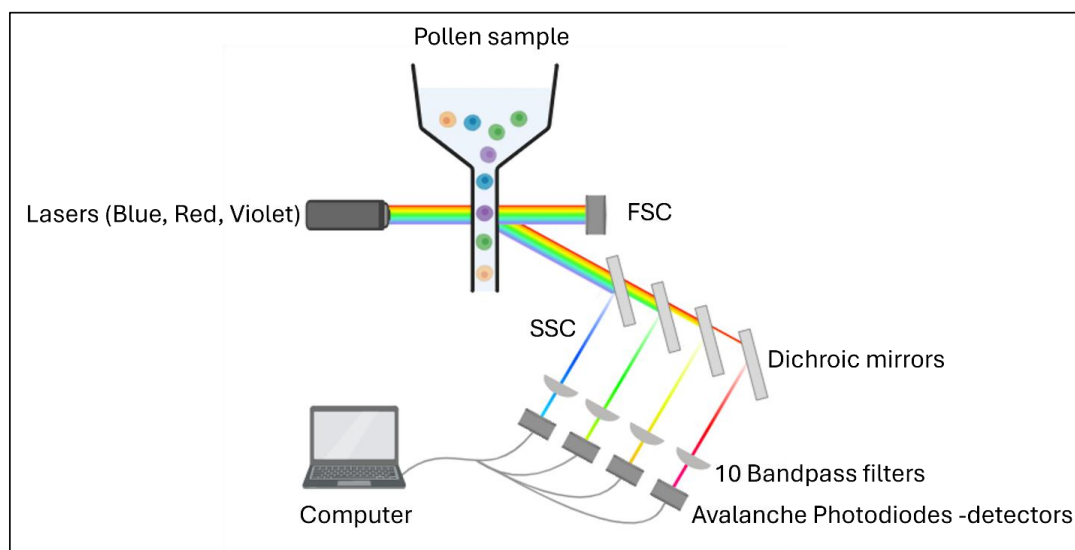


92 desiccant gel. Pollen was extracted from the floral units using a filtration system that retained only particles
93 between 5 and 100 μm in diameter, and the filtrate was suspended in phosphate-buffered saline (PBS)
94 solution to minimize aggregation (see detailed protocol in the supplementary material). A subsample was
95 examined under a light microscope to confirm the presence of pollen grains. If pollen was present, the
96 sample was retained; if not, sampling was repeated, including filtration, and if necessary, additional flowers
97 were collected.

98

99 **2.2 Flow cytometry**

100 Each pollen sample was analysed using flow cytometry (Fig. 1). Measurements were performed with a
101 CytoFLEX cytometer (Beckam Coulter, Inc.), equipped with three excitation lasers at wavelengths of 405
102 nm (violet), 488 nm (blue), and 640 nm (red). Due to a hydrodynamic flow stream, each pollen grain passes
103 sequentially through each laser, which excites the fluorescent proteins on the surface of the pollen grain's
104 outer wall. Depending on their peptide composition, these proteins absorb light at a certain wavelength and
105 emit light radiation at a different wavelength in return producing a characteristic fluorescence signature that
106 varies among species. For each laser, avalanche photodiode (APD) detectors measure the intensity of light
107 emitted at different wavelengths using ten filters: 450/45, 525/40, 610/20 (violet laser), 525/40, 585/42,
108 690/50, 780/60 (blue laser), 660/10, 712/25, 780/60 nm (red laser). In addition to fluorescence, two scatter
109 parameters were recorded to describe particle morphology: grain size and granularity. The forward scatter
110 (FSC) measures light diffracted by the pollen grain at a flat angle, reflecting the approximate diameter of
111 the grain. The sideways scatter (SSC) measures light diffracted by the pollen grain at a right angle,
112 reflecting its granularity.



113

114 **Figure 1:** Flow cytometry workflow on the CytoFLEX (Beckman Coulter, Inc.). Sample containing pollen
 115 enters at the top, and then is excited by three lasers in the blue ($\lambda=488\text{nm}$), red ($\lambda=640\text{nm}$) and violet
 116 ($\lambda=405\text{nm}$) wavelengths, 10 dichroic mirrors, bandpass filters and detectors in different wavelength ranges
 117 ($\lambda=450/45, 525/40, 610/20, 585/42, 525/40, 690/50, 780/60, 660/10, 712/25$ and $780/60\text{ nm}$). There are two
 118 additional detectors for size and granularity: Forward scatter (FSC) and side scatter (SSC). Created with
 119 BioRender.

120 2.3 Data cleaning

121 Although the samples were filtered to retain only particles within the size range of pollen grains ($5\text{-}100\mu\text{m}$),
 122 some non-pollen particles, such as dust or plant debris, were still present. To distinguish pollen from debris,
 123 we used the recorded size, granularity, and fluorescence parameters for each particle which include one
 124 value for size (FSC), one for granularity (SSC), and ten values for fluorescence, each with two components,
 125 the maximum peak height and the peak area except size which has also a width component. This resulted in
 126 a total of 25 parameter values per particle.

127 Data cleaning was performed using Cytexpert software version 2.4.28 (Beckman Coulter, Inc.). For each
 128 species, pollen grains were manually separated from debris using scatter density plots (size vs. granularity)
 129 and histograms of all fluorescence features. This selection relied primarily on the PB450 and Violet610
 130 fluorescence histograms, while cross-checking against the other recorded parameters to ensure consistency.
 131 Adjustments were made as needed to ensure that only true pollen grains were retained (Fig. A1).



132 The final training dataset included all cleaned pollen data from each species along with a separate category,
133 “OTHER”, which combined all debris data from the cleaning step and the particles from certain species for
134 which it was impossible to distinguish pollen from debris, such as those in the *Thuja* genus. The final
135 reference database used to train the model comprised 97 species from 34 different genera. A detailed list of
136 species is presented in Table A1 and the complete training datasets are available on Figshare
137 (<https://doi.org/10.6084/m9.figshare.30870641>).

138 2.4 Machine learning algorithm

139 Four supervised classification algorithms were initially tested: *Random Forest*, *Gradient Boosting*, *Extreme*
140 *Gradient Boosting* and *Neuronal Network*. Among these, the Random Forest algorithm showed the best
141 performance and was therefore selected for subsequent analysis. In our training dataset, the number of pollen
142 grains varies across taxa (min=306; max=35307). This caused the model to more frequently predict taxa
143 with more training examples (Chawla, 2010). To address this class imbalance, we used the synthetic
144 minority over-sampling technique (Chawla, 2010), resulting in a balanced dataset with 1,000 pollen grains
145 per species for the species-level classification model and 10,000 pollen grains per genus for the genus-level
146 classification model. Each dataset was randomly split into two subsets: 70% for training and 30% for
147 validation. Models were trained using the *train()* function from the *caret* package in R software (version
148 4.4.0), calling the *rf()* function for the random forest model. Model robustness was assessed using 10-fold
149 cross-validation implemented via the *trainControl()* function with the “cv” method (nine repetitions for
150 training and one for validation). We trained the models using the default value of 500 trees. The parameter
151 *mtry*, representing the number of variables randomly selected at each node split, was set to 5, based on prior
152 testing across values from 1 to 10. We assessed the models’ performance using the F_1 -score: $F_1 =$
153 $(2 * \text{precision} * \text{recall}) / (\text{precision} + \text{recall})$. Precision is the proportion of correctly predicted positives out of all
154 predicted positives and recall is the proportion of correctly predicted positives out of all actual positives
155 (Grandini et al., 2020). Variable importance was assessed using the mean decrease in Gini coefficient, which
156 quantifies each variable’s contribution to reducing classification error by decreasing node impurity during
157 tree construction. The trained models are available on Figshare.

158 3 Results

159 3.1 Classification performance

160 At the species level, the model achieved a mean F_1 -score of 0.76 (n=97 species; Fig. 2a). The lowest F_1 -
161 scores were obtained for *Quercus rubra* (0.44), *Salix x pendulina* f. *tristis*. (*Salix alba tristis* hereafter) (0.43)
162 and *Ulmus minor* (0.44). Several other species also showed reduced accuracy, with F_1 -scores ranging



163 between 0.5 and 0.65. These included *Acer x freemanii*, *Acer ukurunduense*, *Fagus grandifolia*, *Fraxinus*
164 *nigra*, *Pinus banksiana*, and *Syringa villosa*, as well as several species of the Betulaceae family (*Betula*
165 *papyrifera*, *Carpinus caroliniana*, and *Corylus colurna*), the Juglandaceae family (*Carya ovata*, *Juglans*
166 *nigra*, and *Juglans virginiana*), and the Ulmus genus (*Ulmus davidiana*, *Ulmus propinqua*, and *Ulmus*
167 *pumila*) (Fig. 2a).

168 When trained at the genus level, model performance improved across the 34 genera, reaching a mean F_1 -
169 score of 0.90 (Fig. 2b). The only notable exception was *Juglans*, with an F_1 -score of 0.73. All other genera
170 achieved F_1 -scores close to or above 0.8. Taxa with relatively lower accuracy at the species level, such as
171 those in the genera *Betula*, *Quercus* and *Ulmus*, showed marked improvement at the genus level. Most
172 misclassifications occurred between species within the same genus, as is evident for species from the genus
173 *Ulmus* (see confusion matrices in Appendix B and in supplement material Table S1 and Table S2).

174 3.2 Variables contribution

175 The ranking of predictors using the Gini index shows that the most important variables for distinguishing
176 pollen grains among taxa were granularity (SSC), two fluorescence variables from the violet laser (PB450
177 and Violet610) and one from the blue laser (FITC). These variables exhibited the highest mean decrease in
178 Gini, indicating a major contribution to the homogeneity of nodes and consequently, to overall classification
179 accuracy in the Random Forest model (Fig. 3).

180 Analysis of the variables contributing most to pollen differentiation revealed that size (FSC) and granularity
181 (SSC) varied more among genera than among species within a given genus, whereas fluorescence
182 parameters primarily accounted for the variation observed among species within genera (Fig.4 and
183 Appendix C). Figure 4 illustrates the distributions for six genera known to be allergenic (see Appendix C
184 for more details). Pollen grains from the *Pinus* genus were larger than those from other genera and also had
185 a specific granularity pattern. For these two parameters, FSC and SSC, intra-genus variation for all genera
186 was very small or absent. In contrast, fluorescence parameters showed more pronounced differences among
187 species within the same genus. For example, *Alnus* species presented distinct values across all three
188 fluorescence channels (FITC, Violet610, PB450), while *Corylus* species differed mainly in the Violet610
189 channel. For other genera, only certain species, such as *Betula nigra*, *Quercus macrocarpa*, and *Salix spp.*,
190 showed distinct fluorescence profiles (Fig. 4).

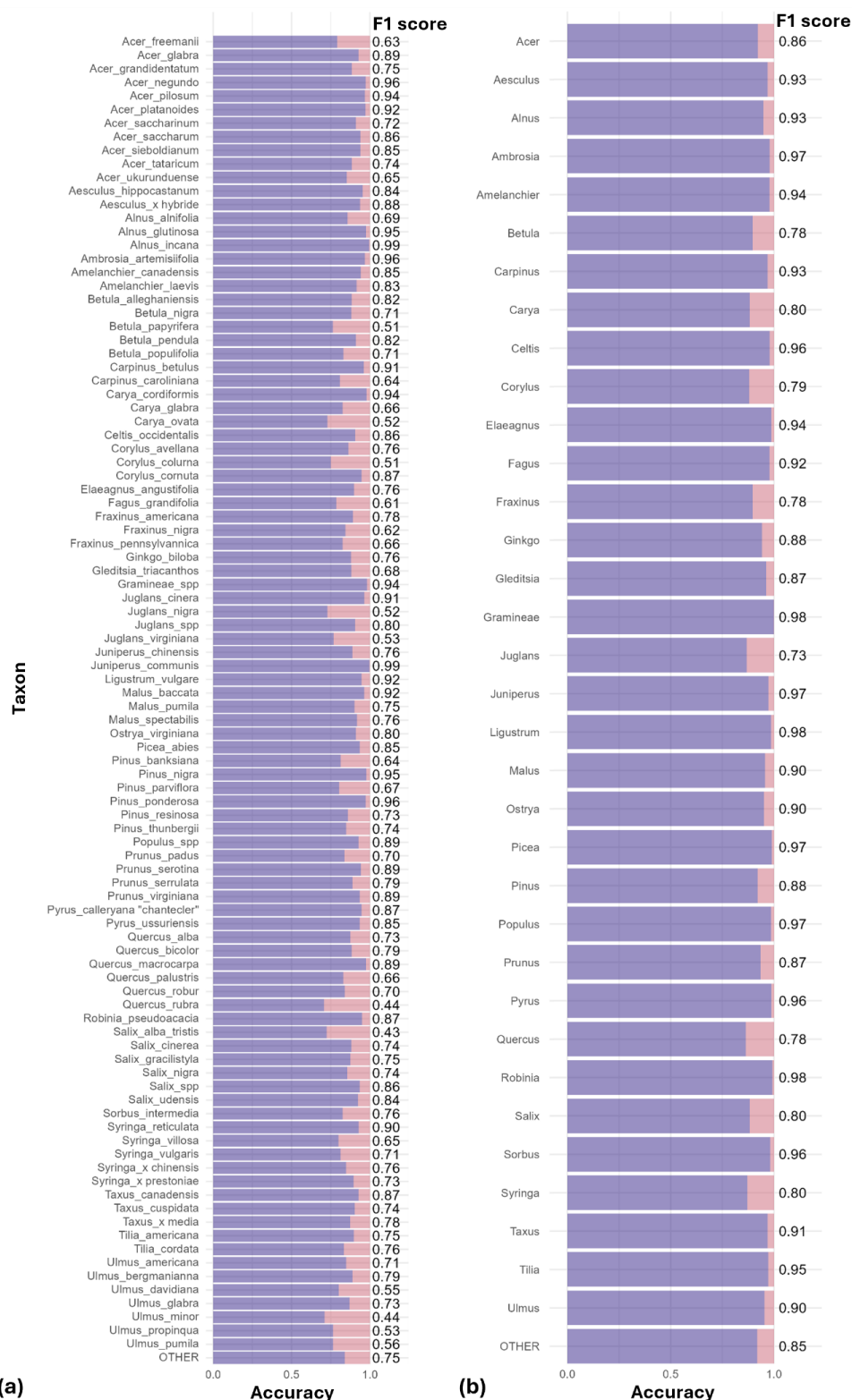




Figure 2: Performance of the classification models at the species (a) and genus levels (b). For each taxon, purple bars represent correct classifications (accuracy) and pink represents misclassifications (1-accuracy). F₁-scores are shown as labels to the right of each bar. Mean F₁-scores were 0.76 for the species-level model and 0.90 for the genus-level model.

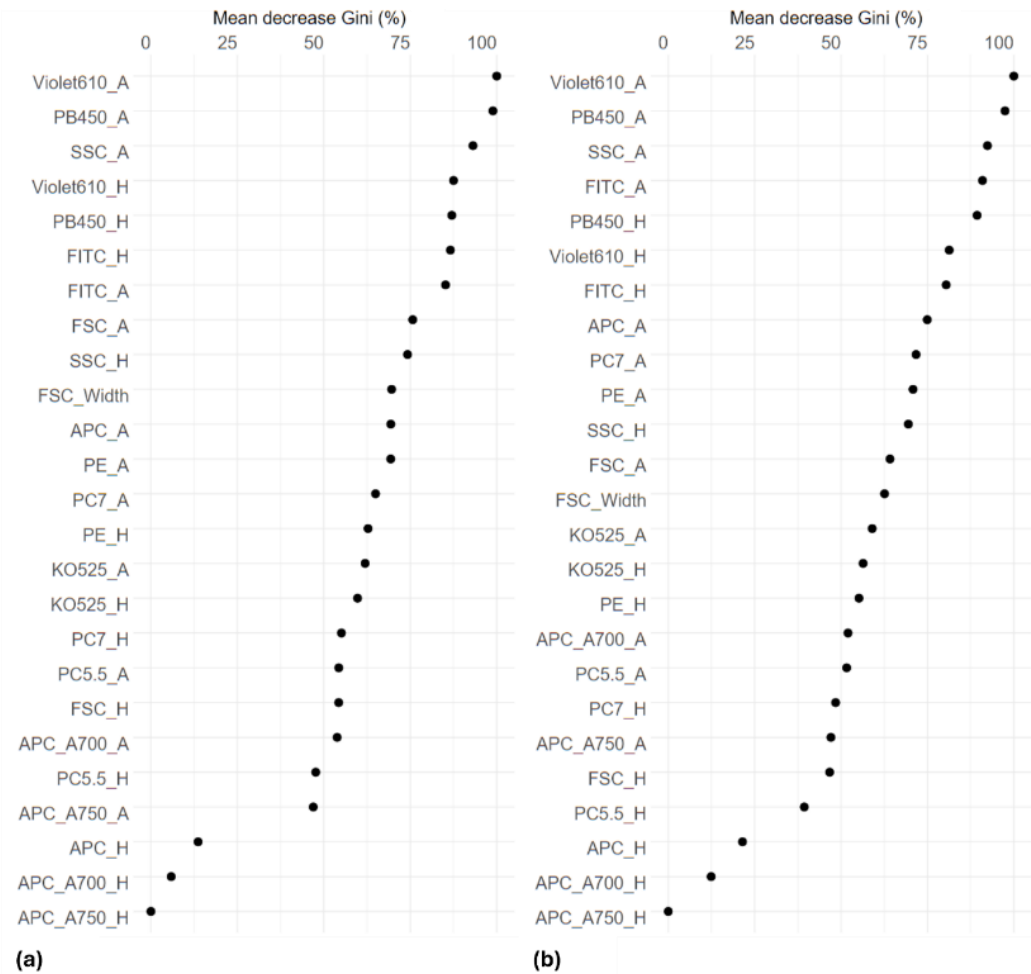
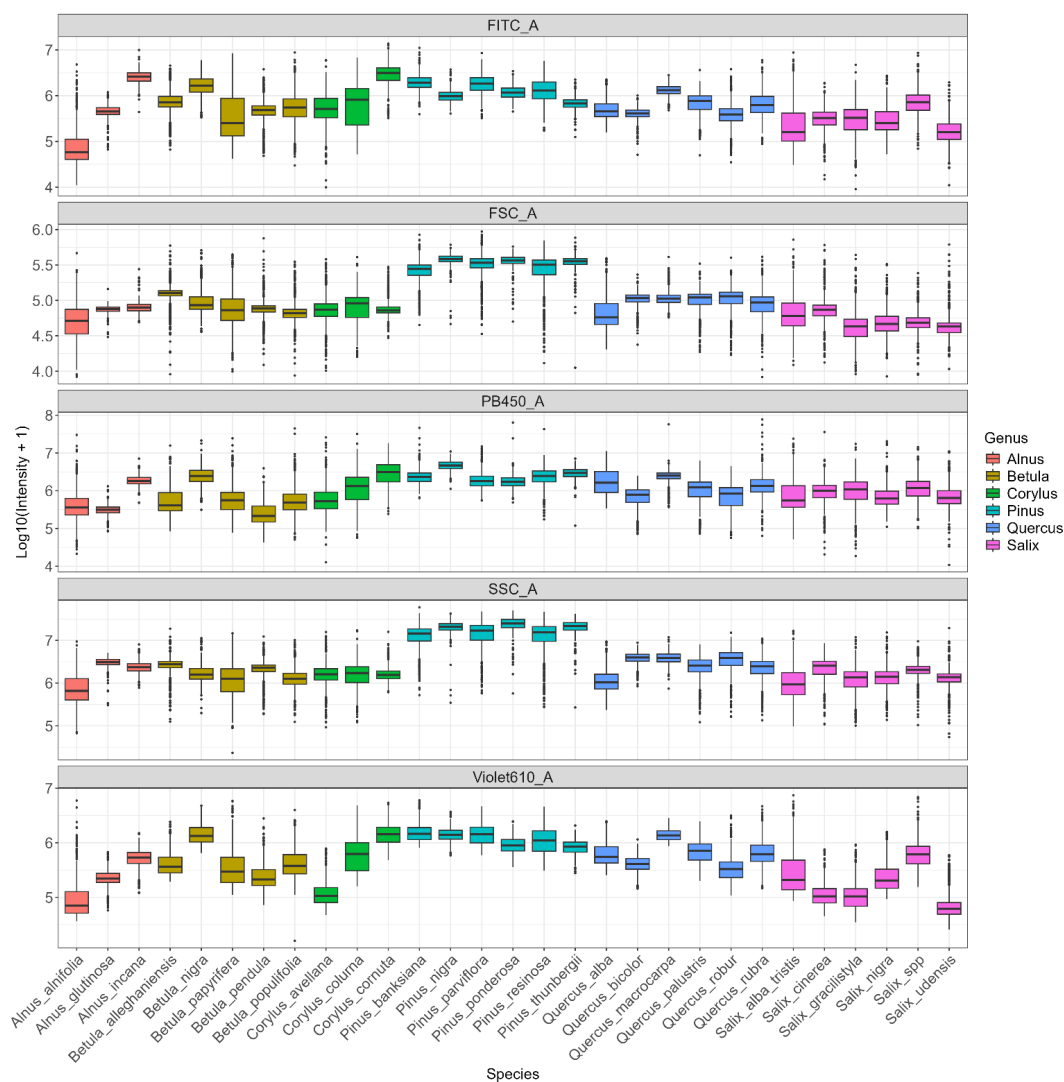


Figure 3: Variable contributions to node and leaf purity in the random forest classification models, measured by mean decrease in Gini index. Higher values indicate greater importance. Results are shown for species-level (a) and genus-level (b) models. Each variable includes two metrics: maximum peak height (H) and peak area (A). Explanation of variable names in Table A2.



201



202

Figure 4: Distribution of log-transformed values for the five variables that contributed the most to distinguish taxa. Fluorescence channels: FITC_A (excitation: 488 nm/emission: 525 nm), PB450_A (excitation: 405 nm/ emission: 450 nm), Violet610_A (excitation: 405 nm/ emission: 610 nm); scatter parameters: SSC_A (granularity) and FSC_A (size). The suffix _A indicates that we consider the signal's peak area. Only species from six known allergenic genera (*Alnus*, *Betula*, *Corylus*, *Pinus*, *Quercus*, *Salix*) are shown. For more species see Appendix C. Colors indicate genus.

209



210 4 Discussion

211 Our results demonstrate that flow cytometry combined with machine learning can reliably identify pollen
212 across a wide range of taxa. The models achieved high classification performance ($F_1=0.76$ at the species
213 level and 0.90 at the genus level) highlighting the potential of this approach as a scalable alternative to
214 traditional microscopy for pollen identification. This represents a significant improvement over
215 conventional methods, such as microscopy, which typically only resolve pollen to the genus or family level.
216 The improved performance of the genus-level model over the species-level model most likely reflects
217 biological and structural similarities among species within the same genus. This was particularly evident for
218 species in the *Betulaceae* family, which are wind-pollinated and considered highly allergenic (D'Amato et
219 al., 2007; Falagiani, 1989), but also for other genera especially abundant in Montreal, such as *Acer*, *Syringa*,
220 and *Ulmus*.

221 The advantage of flow cytometry coupled with machine learning lies not only in its performance in
222 classifying at the genus or species level, but especially in its ability to enable automated, high-throughput
223 identification (≈ 5000 grains $\cdot s^{-1}$) while avoiding the lengthy and costly training required for human
224 specialists. Accurate monitoring is clinically important, as even low pollen concentrations (10–50 grains
225 per cubic meter) can trigger allergic symptoms (Steckling-Muschack et al., 2021). From a public health
226 perspective, the genus-level model is therefore appropriate, as it provides higher accuracy for the taxa most
227 relevant to allergy monitoring.

228 The fluorescence variables that contributed most to pollen classification were associated with blue and violet
229 excitation lasers, with emission detected in the blue (PB450), red-orange (Violet610), and green (FITC)
230 channels. This pattern is consistent with the known autofluorescence properties of sporopollenin, the main
231 biopolymer in the pollen exine, which emits strongly near 475 nm (Pöhlker et al., 2013). Additional
232 emissions likely originate from secondary compounds such as flavonoids, carotenoids, and terpenes located
233 in the exine or pollenkit coating (Donaldson, 2020; Pöhlker et al., 2013). The distribution of the most
234 discriminative variables indicates that size and granularity primarily differentiate genera, while blue, red-
235 orange and green fluorescence channels capture species-level differences within genera. This pattern
236 explains the model's higher accuracy at the genus-level and its partial success in distinguishing closely
237 related species. The misclassifications at species-level likely stem from the high similarity in pollen size
238 and fluorescence spectra among closely related species, which makes them harder to distinguish. In addition,
239 because our classification relied on size and fluorescence alone, without complementary morphological data
240 such as holography images (Erb et al., 2024; Gierlicka et al., 2022; Zhang and Abdulla, 2023), the model's
241 performance may have been constrained by limited representation of some taxa in the reference dataset.



242 Increasing both the number of pollen grains per species and the diversity of species within each genus would
243 help train more robust models. Future research should prioritize expanding reference datasets, ideally
244 through the creation of a global database of pollen fluorescence signatures, which represent the emission
245 spectrum for given excitation wavelengths. Such a resource, similar to *The Global Pollen Project*, for
246 microscopic images (Martin and Harvey, 2017), would provide a valuable foundation for machine learning
247 and deep learning applications in aerobiology, but also ecology, palynology, paleoecology, and other pollen
248 related fields.

249 Another factor that may explain the reduced model accuracy is that some species in our reference collection
250 could not be included in the model's training dataset due to the impossibility to distinguish pollen from
251 debris during the data cleaning, even though we had visually confirmed the presence of pollen grains in our
252 samples. These data were included in the training dataset under the category "OTHER" rather than assigned
253 to individual taxa. Such was the case for *Thuja*, a genus abundant in Montreal (Paquette et al., submitted),
254 likely due to the small size of its pollen grains, which can easily mix with debris or because pollen grains
255 included in our dataset may have been limited in quantity or had not fully matured. Indeed, distinguishing
256 male from female *Thuja* cones and assessing the phenological stage to collect mature pollen is difficult, and
257 the small size of the cones is another challenge for pollen extraction. Improving collection and extraction
258 protocols for this genus could help reduce debris contamination in future sampling.

259 A crucial next step is to adapt these models for use on complex airborne samples collected in urban
260 environments. Such samples often contain large amounts of debris as during atmospheric transport, pollen
261 grains may remain airborne for days or weeks, during which they can fold, crack, or adhere to air pollutants
262 (De Weger et al., 2024). They are also exposed to ultraviolet radiation and humidity fluctuations that can
263 alter fluorescence properties. These factors complicate the discrimination of true pollen grains from other
264 particles and represent a major challenge for operational implementation.

265 Because small pollen grains, folded grains and debris can have overlapping size distributions,
266 misclassification remains a possibility, with pollen occasionally identified as debris, and vice versa. Future
267 research could therefore explore multidimensional hierarchical classification frameworks, especially when
268 complementarity data such as holographic images are available for validation. For example, when
269 classification confidence is high, the model could assign a species-level label, but default to a broader
270 taxonomic category such as genus or family when uncertainty is greater (Hernández et al., 2014). This
271 flexibility would prevent incorrect fine-level classifications and improve overall reliability under complex
272 environmental conditions.



273 Another limitation of flow cytometry-based models concerns their device dependency, as fluorescence
274 intensity values are typically linked to the specific cytometer used during model training, which limits model
275 transferability across instruments. Standardization procedures, such as calibrating cytometers using
276 Rainbow beads and Quality Control beads could help ensure consistent signal outputs across different
277 instruments (Solly et al., 2013). The present work was carried out using a conventional cytometer with three
278 lasers and ten filters; using equipment with more lasers and detectors could refine the detection of
279 fluorescent signatures and detect more of them. Spectral cytometry also opens up new possibilities for
280 analyzing fluorescent signatures on a larger scale (Konecny et al., 2024), which could enable even better
281 characterization of pollen based on its fluorescence.

282 The combination of flow cytometry and a Random Forest classification model proves to be a highly
283 promising approach for the identification of airborne pollen in urban environments. By relying exclusively
284 on routinely measured cytometric parameters, rather than images, this method ensures broad applicability
285 and compatibility with standard healthcare and clinical cytometers. Integrating this approach into existing
286 aerobiological monitoring networks could enable near-real-time identification and quantification of
287 allergenic pollen. We also built an extensive reference pollen collection comprising 97 species across 34
288 genera. For each species, we have several floral units (flower, catkins, cones) containing pollen, microscopic
289 slides, and flow cytometry data for all pollen grains. This reference collection could be reused for different
290 purposes such as future model training.

291 **5 Conclusion**

292 This study demonstrates a significant advancement in pollen identification by combining flow cytometry
293 with a random forest classification model. This approach achieved high accuracy at both the genus ($F_1 =$
294 0.90) and species levels ($F_1 = 0.76$), surpassing several limitations of traditional microscopy. While species-
295 level classification remains challenging for certain taxa, the results highlight the method's robustness and
296 potential for large-scale implementation. With continued refinement and standardization, this approach
297 could enable near-real-time, cheap, high throughput pollen identification and broaden its applications in
298 aerobiological monitoring, while supporting public health applications and advancing research in pollen
299 ecology worldwide.

300 **6 Code availability**

301 The code is available on the public Github repository SarahTardif/Pollen-classification-model.

302 **7 Data availability**



303 Training datasets and trained models are available on a Figshare repository
304 (<https://doi.org/10.6084/m9.figshare.30870641>). More data can be provided upon request.

305 **8 Author contribution**

306 ST: conceptualization, data collection, analyses, writing – original draft; AP, IL and RSS:
307 conceptualization, funding acquisition, supervision, validation, support, writing – review and editing; GB:
308 Methodology (cytometry), writing – review and editing; MRK: Methodology (lab protocols), writing –
309 review and editing; GL: Methodology (initial algorithm for the machine learning model), writing – review
310 and editing

311 **9 Competing interests**

312 The authors declare that they have no conflict of interest.

313 **10 Acknowledgements**

314 We thank the CERMO-UQAM Imaging Platform and the Aerobiology Research Laboratories for their
315 technical support. We are grateful for the precious help of Kira Safranova, Emily Ducharme, Maya Héon,
316 and Kim Florentin in sampling, filtering, and running fresh pollen through the cytometer. We thank the
317 Montreal Botanical Garden for permitting pollen collection from tree flowers. Model training was
318 performed on supercomputers managed by Calcul Québec and the Digital Research Alliance of Canada.

319 **11 Financial support**

320 This work was funded by NSERC-Alliance ALLRP 554373 – 21 and *Fonds vert dans le cadre du Plan*
321 *d'action 2013-2020 sur les changements climatiques du gouvernement québécois* awarded to AP.
322 ST also received funding from the Urban forestry program NSERC-CREATE -543300-20.



323 **Appendix A: Reference pollen collection**

324 **Table A1:** Species in the reference pollen collection

Family	Genus	Species (scientific name)
Asteraceae	Ambrosia	<i>Ambrosia artemisiifolia</i>
Betulaceae	Alnus	<i>Alnus alnifolia</i>
		<i>Alnus glutinosa</i>
		<i>Alnus incana</i>
	Betula	<i>Betula alleghaniensis</i>
		<i>Betula nigra</i>
		<i>Betula papyrifera</i>
		<i>Betula pendula</i>
		<i>Betula populifolia</i>
	Carpinus	<i>Carpinus betulus</i>
		<i>Carpinus caroliniana</i>
	Corylus	<i>Corylus avellana</i>
		<i>Corylus colurna</i>
		<i>Corylus cornuta</i>
	Ostrya	<i>Ostrya virginiana</i>
Cannabacées	Celtis	<i>Celtis occidentalis</i>
Cupressaceae	Juniperus	<i>Juniperus chinensis</i>
		<i>Juniperus communis</i>
Elaeagnaceae	Elaeagnus	<i>Elaeagnus angustifolia</i>
Fabaceae	Gleditsia	<i>Gleditsia triacanthos</i>
	Robinia	<i>Robinia pseudoacacia</i>
Fagaceae	Fagus	<i>Fagus grandifolia</i>
	Quercus	<i>Quercus alba</i>
		<i>Quercus bicolor</i>
		<i>Quercus macrocarpa</i>
		<i>Quercus palustris</i>
		<i>Quercus robur</i>
		<i>Quercus rubra</i>
Ginkgoaceae	Ginkgo	<i>Ginkgo biloba</i>
Gramineae	-	<i>Gramineae spp</i>
Juglandaceae	Carya	<i>Carya cordiformis</i>
		<i>Carya glabra</i>
		<i>Carya ovata</i>
	Juglans	<i>Juglans cinera</i>
		<i>Juglans nigra</i>
		<i>Juglans spp</i>
		<i>Juglans virginiana</i>



325

Family	Genus	Species (scientific name)
Oleaceae	Fraxinus	<i>Fraxinus americana</i>
		<i>Fraxinus nigra</i>
		<i>Fraxinus pennsylvannica</i>
	Ligustrum	<i>Ligustrum vulgare</i>
	Syringa	<i>Syringa reticulata</i>
		<i>Syringa villosa</i>
		<i>Syringa vulgaris</i>
		<i>Syringa x chinensis</i>
		<i>Syringa x prestoniae</i>
Pinaceae	Picea	<i>Picea abies</i>
	Pinus	<i>Pinus banksiana</i>
		<i>Pinus nigra</i>
		<i>Pinus parviflora</i>
		<i>Pinus ponderosa</i>
		<i>Pinus resinosa</i>
		<i>Pinus thunbergii</i>
Rosaceae	Amelanchier	<i>Amelanchier canadensis</i>
		<i>Amelanchier laevis</i>
	Malus	<i>Malus baccata</i>
		<i>Malus pumila</i>
		<i>Malus spectabilis</i>
	Prunus	<i>Prunus padus</i>
		<i>Prunus serotina</i>
		<i>Prunus serrulata</i>
		<i>Prunus virginiana</i>
	Pyrus	<i>Pyrus calleryana</i> "chantecler"
		<i>Pyrus ussuriensis</i>
	Sorbus	<i>Sorbus intermedia</i>
Salicaceae	Populus	<i>Populus spp</i>
	Salix	<i>Salix alba tristis</i>
		<i>Salix cinerea</i>
		<i>Salix gracilistyla</i>
		<i>Salix nigra</i>
		<i>Salix spp</i>
		<i>Salix udensis</i>

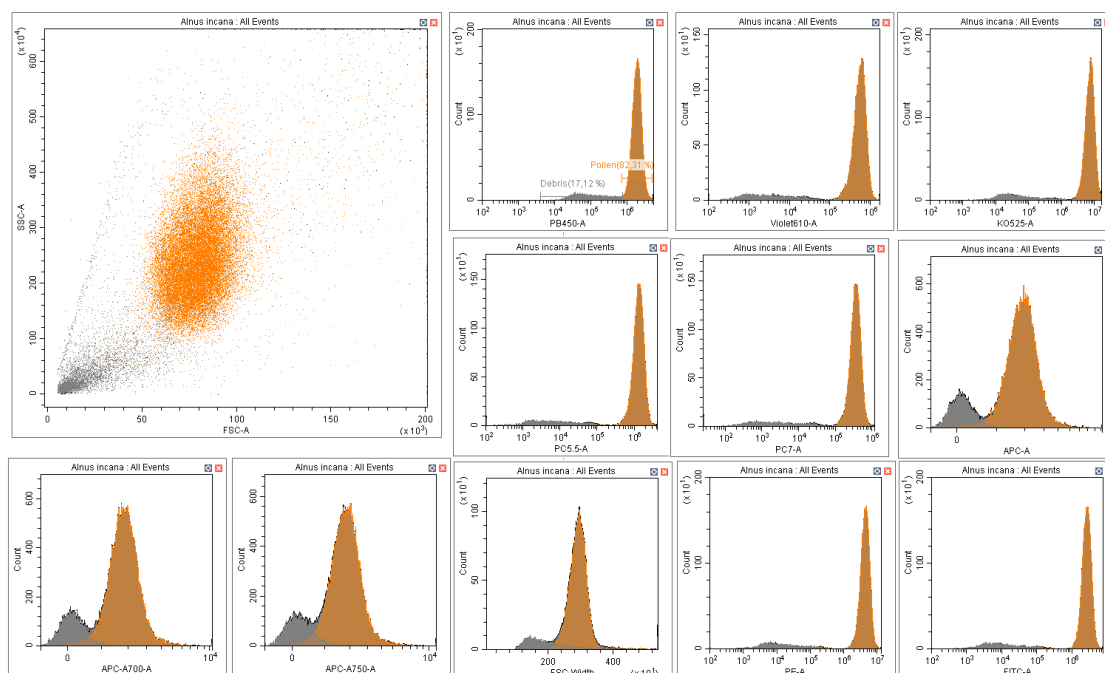
326



Family	Genus	Species (scientific name)
Sapindaceae	Acer	<i>Acer freemanii</i>
		<i>Acer glabra</i>
		<i>Acer grandidentatum</i>
		<i>Acer negundo</i>
		<i>Acer pilosum</i>
		<i>Acer platanoides</i>
		<i>Acer saccharinum</i>
		<i>Acer saccharum</i>
		<i>Acer sieboldianum</i>
		<i>Acer tataricum</i>
		<i>Acer ukurunduense</i>
	Aesculus	<i>Aesculus hippocastanum</i>
		<i>Aesculus x hybride</i>
Taxaceae	Taxus	<i>Taxus canadensis</i>
		<i>Taxus cuspidata</i>
		<i>Taxus x media</i>
Tiliaceae	Tilia	<i>Tilia americana</i>
		<i>Tilia cordata</i>
Ulmaceae	Ulmus	<i>Ulmus americana</i>
		<i>Ulmus bergmanianna</i>
		<i>Ulmus davidiana</i>
		<i>Ulmus glabra</i>
		<i>Ulmus minor</i>
		<i>Ulmus propinqua</i>
		<i>Ulmus pumila</i>

Table A2: Explanation of cytometry variable names, showing the respective excitation lasers and emission detectors with their wavelengths and associated colors.

Detector name	PB 450	KO 525	Violet 610	PE	FITC	PC5,5	PC7	APC	APC-A700	APC-A750	Granularity	Size
											SSC	FSC
Excitation (nm)	405			488				640			-	-
Emission (nm)	450	525	610	585	525	690	780	660	712	780	-	-



331

332 **Figure A1:** Distinction pollen (orange) versus debris (grey) on CytExpert software: Example of *Alnus*
 333 *incana*

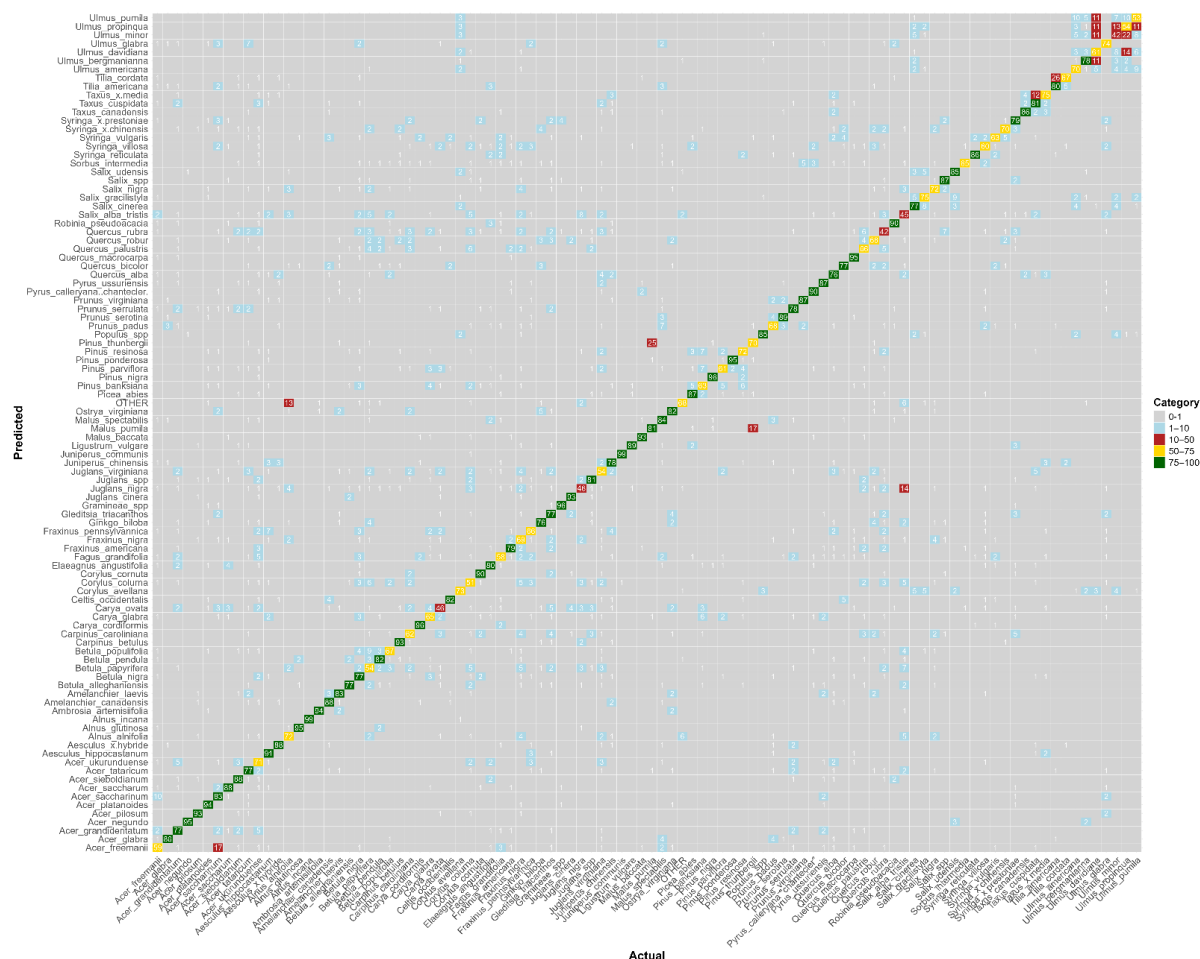
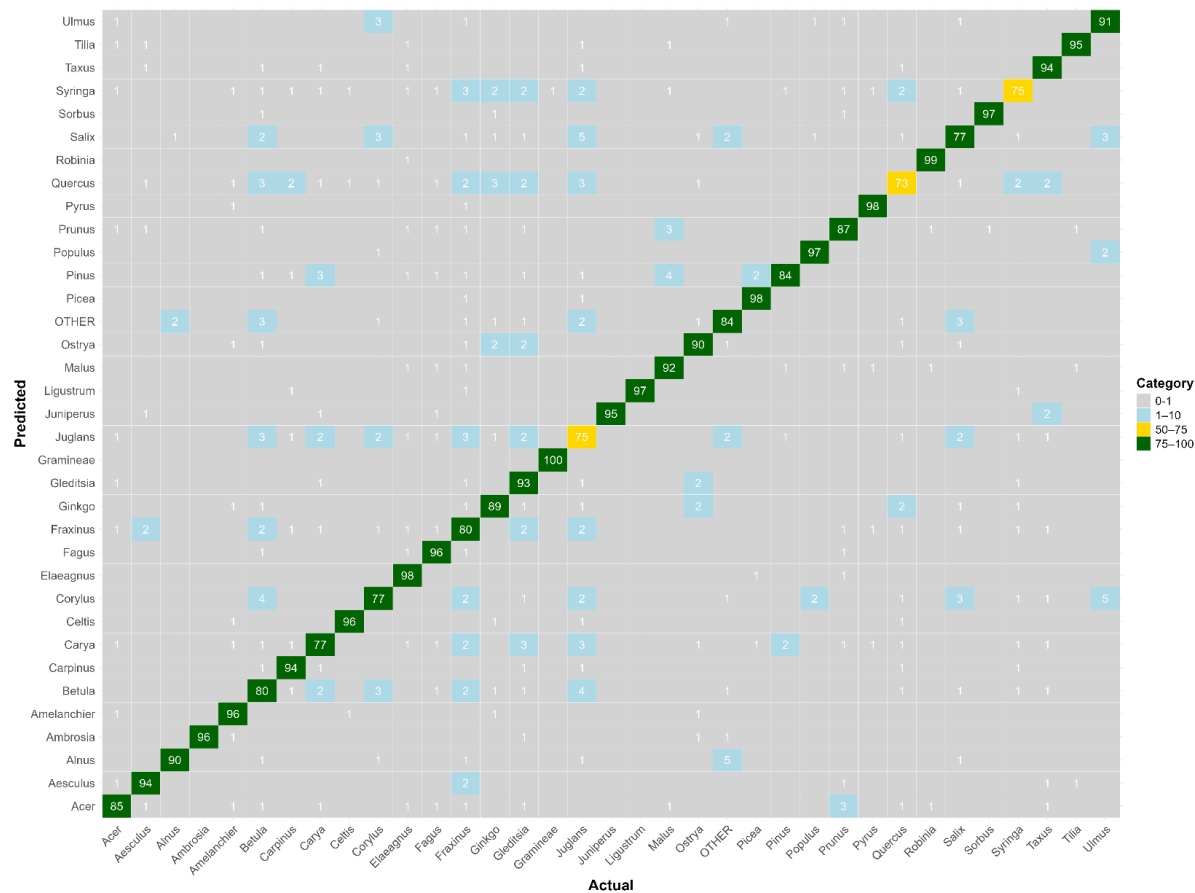
334 **Appendix B: Confusion matrices**

Figure B1: Confusion matrix for the species-level model. The values represent, for each species, the percentage of pollen grains correctly classified (on the diagonal) and misclassified with the actual corresponding species (on the x-axis). Colors correspond to categories (0–1% in gray, 1–10% in blue, 10–50% in red, 50–75% in yellow, and 75–100% in green). Raw data are provided in supplement material Table S1.

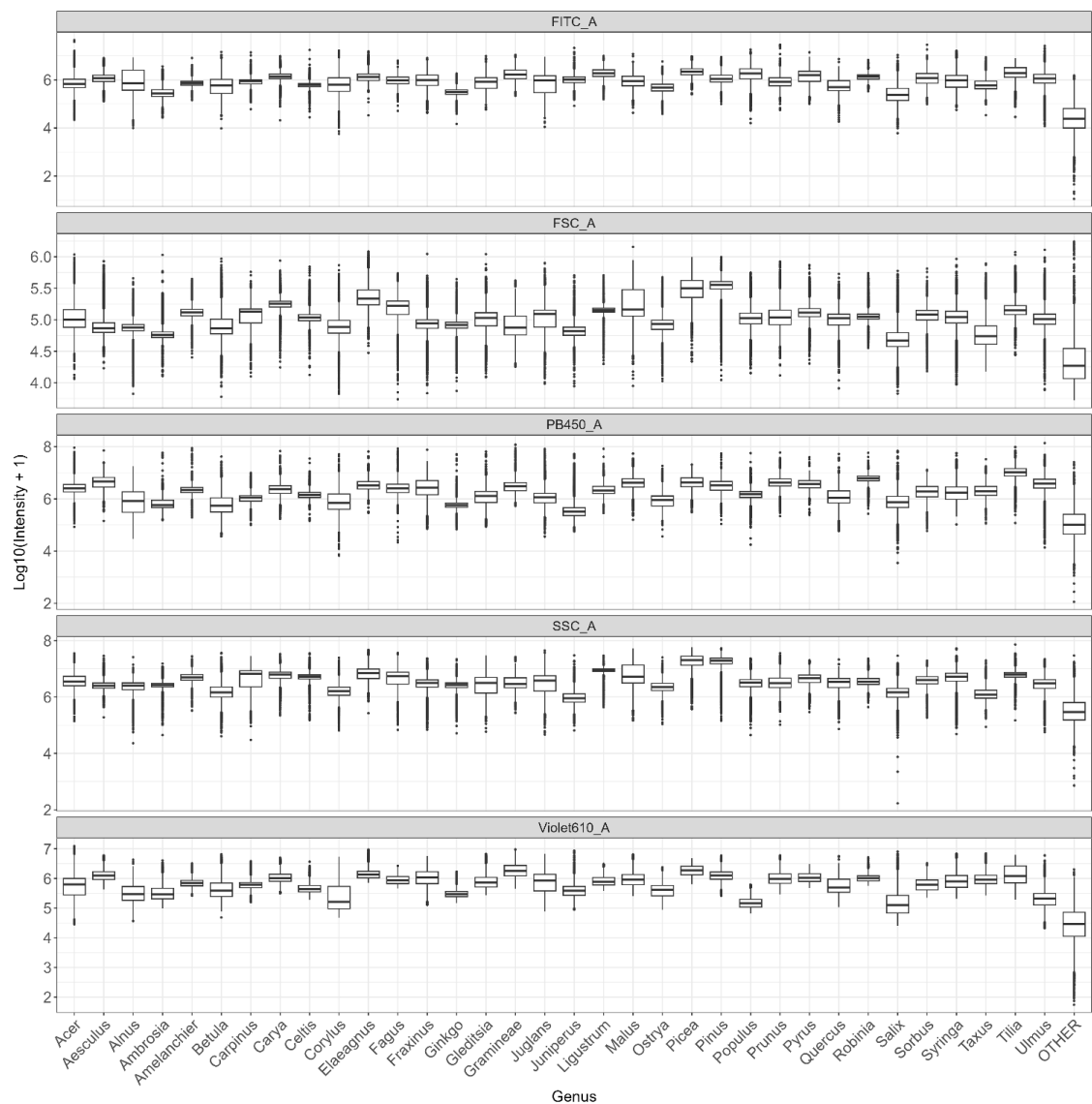


341

342 **Figure B2:** Confusion matrix for the genus-level model. The values represent, for each genus, the
343 percentage of pollen grains correctly classified (on the diagonal) and misclassified with the actual
344 corresponding genus (on the x-axis). Colors correspond to categories (0–1% in gray, 1–10% in blue,
345 10–50% in red, 50–75% in yellow, and 75–100% in green). Raw data are provided in supplement material
346 Table S2.



347 **Appendix C: Distributions of values for the main discriminant variables.**
348



349
350 **Figure C1:** Distribution of log-transformed values for the variables that contribute the most to distinguish
351 taxa (FITC,FSC,SSC,Violet610, PB450) across all genera.

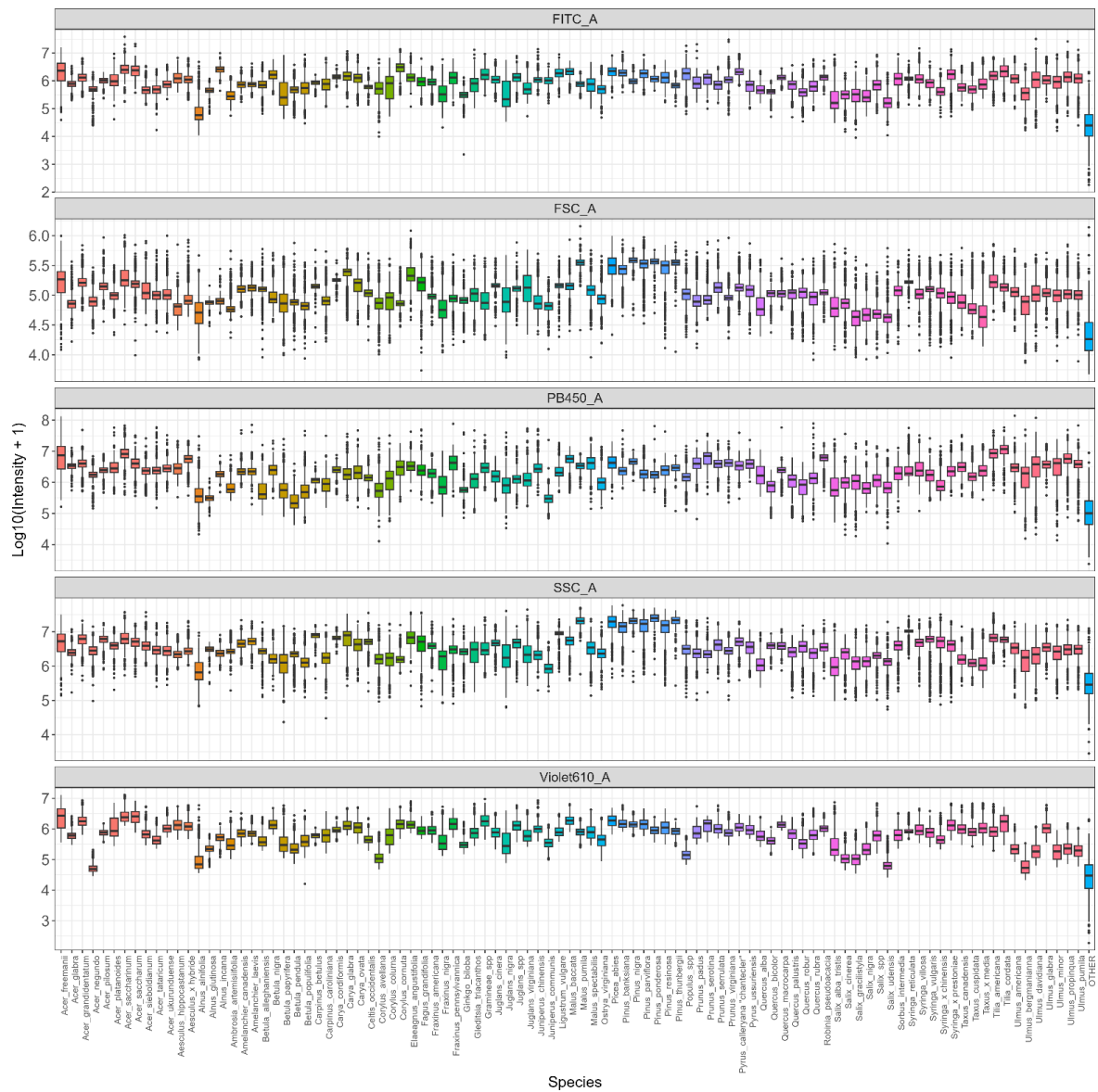


Figure C2: Distribution of log-transformed values for the variables that contribute the most to distinguish taxa (FITC,FSC,SSC,Violet610, PB450) across all species. Colors indicate genus.



References

- Ahlholm, Helander, and Savolainen: Genetic and environmental factors affecting the allergenicity of birch (*Betula pubescens* ssp. *czerepanovii* [Orl.] Hämet-Ahti) pollen, *Clinical & Experimental Allergy*, 28, 1384–1388, <https://doi.org/10.1046/j.1365-2222.1998.00404.x>, 1998.
- Anderegg, W. R. L., Abatzoglou, J. T., Anderegg, L. D. L., Bielory, L., Kinney, P. L., and Ziska, L.: Anthropogenic climate change is worsening North American pollen seasons, *Proc. Natl. Acad. Sci. U.S.A.*, 118, e2013284118, <https://doi.org/10.1073/pnas.2013284118>, 2021.
- Brennan, G. L., Potter, C., de Vere, N., Griffith, G. W., Skjøth, C. A., Osborne, N. J., Wheeler, B. W., McInnes, R. N., Clewlow, Y., Barber, A., Hanlon, H. M., Hegarty, M., Jones, L., Kurganskiy, A., Rowney, F. M., Armitage, C., Adams-Groom, B., Ford, C. R., Petch, G. M., and Creer, S.: Temperate airborne grass pollen defined by spatio-temporal shifts in community composition, *Nat Ecol Evol*, 3, 750–754, <https://doi.org/10.1038/s41559-019-0849-7>, 2019.
- Chawla, N. V.: Data mining for imbalanced datasets: An overview, in: *Data mining and knowledge discovery handbook*, 875–886, 2010.
- D’Amato, G., Cecchi, L., Bonini, S., Nunes, C., Annesi-Maesano, I., Behrendt, H., Liccardi, G., Popov, T., and van Cauwenberge, P.: Allergenic pollen and pollen allergy in Europe, *Allergy*, 62, 976–990, <https://doi.org/10.1111/j.1398-9995.2007.01393.x>, 2007.
- De Weger, L. A., Verbeek, C., Markey, E., O’Connor, D. J., and Gosling, W. D.: Greater difference between airborne and flower pollen chemistry, than between pollen collected across a pollution gradient in the Netherlands, *Science of The Total Environment*, 934, 172963, <https://doi.org/10.1016/j.scitotenv.2024.172963>, 2024.
- Donaldson, L.: Autofluorescence in Plants, *Molecules*, 25, 2393, <https://doi.org/10.3390/molecules25102393>, 2020.
- Dunker, S., Motivans, E., Rakosy, D., Boho, D., Mäder, P., Hornick, T., and Knight, T. M.: Pollen analysis using multispectral imaging flow cytometry and deep learning, *New Phytologist*, 229, 593–606, <https://doi.org/10.1111/nph.16882>, 2021.
- Dunker, S., Boyd, M., Durka, W., Erler, S., Harpole, W. S., Henning, S., Herzsuh, U., Hornick, T., Knight, T., Lips, S., Mäder, P., Švara, E. M., Mozarowski, S., Rakosy, D., Römermann, C., Schmitt-Jansen, M., Stoof-Leichsenring, K., Stratmann, F., Treudler, R., Virtanen, R., Wendt-Potthoff, K., and Wilhelm, C.: The potential of multispectral imaging flow cytometry for environmental monitoring, *Cytometry Pt A*, 101, 782–799, <https://doi.org/10.1002/cyto.a.24658>, 2022.
- Erb, S., Graf, E., Zeder, Y., Lionetti, S., Berne, A., Clot, B., Lieberherr, G., Tummon, F., Wullschlegel, P., and Crouzy, B.: Real-time pollen identification using holographic imaging and fluorescence measurements, *Atmos. Meas. Tech.*, 17, 441–451, <https://doi.org/10.5194/amt-17-441-2024>, 2024.
- Falagiani, P.: *Pollinosis*, CRC Press, 288 pp., 1989.
- Gierlicka, I., Kasprzyk, I., and Wnuk, M.: Imaging Flow Cytometry as a Quick and Effective Identification Technique of Pollen Grains from Betulaceae, Oleaceae, Urticaceae and Asteraceae, *Cells*, 11, 598, <https://doi.org/10.3390/cells11040598>, 2022.



- 393 Grandini, M., Bagli, E., and Visani, G.: Metrics for Multi-Class Classification: an Overview,
394 <https://doi.org/10.48550/arXiv.2008.05756>, 13 August 2020.
- 395 Hernández, J., Sucar, L. E., and Morales, E. F.: Multidimensional hierarchical classification, Expert
396 Systems with Applications, 41, 7671–7677, <https://doi.org/10.1016/j.eswa.2014.05.054>, 2014.
- 397 Holt, K. A. and Bennett, K. D.: Principles and methods for automated palynology, New Phytologist, 203,
398 735–742, <https://doi.org/10.1111/nph.12848>, 2014.
- 399 Kim, K. R., Oh, J.-W., Woo, S.-Y., Seo, Y. A., Choi, Y.-J., Kim, H. S., Lee, W. Y., and Kim, B.-J.: Does
400 the increase in ambient CO₂ concentration elevate allergy risks posed by oak pollen?, Int J Biometeorol,
401 62, 1587–1594, <https://doi.org/10.1007/s00484-018-1558-7>, 2018.
- 402 Konecny, A. J., Mage, P. L., Tyznik, A. J., Prlic, M., and Mair, F.: OMIP-102: 50-color phenotyping of
403 the human immune system with in-depth assessment of T cells and dendritic cells, Cytometry Part A, 105,
404 430–436, <https://doi.org/10.1002/cyto.a.24841>, 2024.
- 405 Ladeau, S. L. and Clark, J. S.: Pollen production by *Pinus taeda* growing in elevated atmospheric CO₂,
406 Functional Ecology, 20, 541–547, <https://doi.org/10.1111/j.1365-2435.2006.01133.x>, 2006.
- 407 Martin, A. C. and Harvey, W. J.: The Global Pollen Project: a new tool for pollen identification and the
408 dissemination of physical reference collections, Methods Ecol Evol, 8, 892–897,
409 <https://doi.org/10.1111/2041-210X.12752>, 2017.
- 410 Medek, D. E., Katelaris, C. H., Milic, A., Beggs, P. J., Lampugnani, E. R., Vicendese, D., Erbas, B., and
411 Davies, J. M.: Aerobiology matters: Why people in the community access pollen information and how
412 they use it, Clinical & Translational All, 15, e70031, <https://doi.org/10.1002/cla2.70031>, 2025.
- 413 Mousavi, F., Oteros, J., Shahali, Y., and Carinanos, P.: Impacts of climate change on allergenic pollen
414 production: A systematic review and meta-analysis, Agricultural and Forest Meteorology, 349, 109948,
415 <https://doi.org/10.1016/j.agrformet.2024.109948>, 2024.
- 416 Ogden, E. C., Museum, N. Y. S., Service, S., and Commission, U. S. A. E.: Manual for Sampling
417 Airborne Pollen, Hafner Press, 1974.
- 418 Paquette, A., Sousa-Silva, R., Fernandez, M., Faticov, M., Schillé, L., Bacon, E., Cameron, E., Fraysse, J.,
419 gagnon Koudji, E., Poirier, S., Rondeau-Leclaire, J., Tardif, S., Handa, T., Laforest-Lapointe, I., Puric-
420 Mladenovic, D., and Ziter, C.: Montreal Urban Observatory: research platform to monitor urban forest
421 ecosystems for global change adaptation and health, submitted.
- 422 Pöhlker, C., Huffman, J. A., Förster, J.-D., and Pöschl, U.: Autofluorescence of atmospheric bioaerosols:
423 spectral fingerprints and taxonomic trends of pollen, Atmos. Meas. Tech., 6, 3369–3392,
424 <https://doi.org/10.5194/amt-6-3369-2013>, 2013.
- 425 Šaulienė, I., Šukienė, L., Daunys, G., Valiulis, G., Vaitkevičius, L., Matavulj, P., Brdar, S., Panic, M.,
426 Sikoparija, B., Clot, B., Crouzy, B., and Sofiev, M.: Automatic pollen recognition with the Rapid-E
427 particle counter: the first-level procedure, experience and next steps, Atmos. Meas. Tech., 12, 3435–3452,
428 <https://doi.org/10.5194/amt-12-3435-2019>, 2019.
- 429 Sauvageat, E., Zeder, Y., Auderset, K., Calpini, B., Clot, B., Crouzy, B., Konzelmann, T., Lieberherr, G.,
430 Tummon, F., and Vasilatou, K.: Real-time pollen monitoring using digital holography, Atmos. Meas.
431 Tech., 13, 1539–1550, <https://doi.org/10.5194/amt-13-1539-2020>, 2020.



- 432 Savouré, M., Bousquet, J., Jaakkola, J. J. K., Jaakkola, M. S., Jacquemin, B., and Nadif, R.: Worldwide
433 prevalence of rhinitis in adults: A review of definitions and temporal evolution, *Clinical & Translational*
434 *All*, 12, e12130, <https://doi.org/10.1002/clin.12130>, 2022.
- 435 Sikoparija, B., Matavulj, P., Simovic, I., Radisic, P., Brdar, S., Minic, V., Tesendic, D., Kadantsev, E.,
436 Palamarchuk, J., and Sofiev, M.: Classification accuracy and compatibility across devices of a new Rapid-
437 E+ flow cytometer, <https://doi.org/10.5194/egusphere-2024-187>, 2 April 2024.
- 438 Smith, E. G.: Sampling and identifying allergenic pollens and molds. An illustrated manual for physicians
439 and lab technicians., *Sampling and identifying allergenic pollens and molds. An illustrated manual for*
440 *physicians and lab technicians.*, 1984.
- 441 Solly, F., Rigollet, L., Baseggio, L., Guy, J., Borgeot, J., Guérin, E., Debliquis, A., Drenou, B., Campos,
442 L., Lacombe, F., and Béné, M. C.: Comparable flow cytometry data can be obtained with two types of
443 instruments, Canto II, and Navios. A GEIL study, *Cytometry A*, 83, 1066–1072,
444 <https://doi.org/10.1002/cyto.a.22404>, 2013.
- 445 Sousa-Silva, R., Smargiassi, A., Paquette, A., Kaiser, D., and Kneeshaw, D.: Exactly what do we know
446 about tree pollen allergenicity?, *The Lancet Respiratory Medicine*, 8, e10, [https://doi.org/10.1016/S2213-](https://doi.org/10.1016/S2213-2600(19)30472-2)
447 [2600\(19\)30472-2](https://doi.org/10.1016/S2213-2600(19)30472-2), 2020.
- 448 Steckling-Muschack, N., Mertes, H., Mittermeier, I., Schutzmeier, P., Becker, J., Bergmann, K.-C., Böse-
449 O'Reilly, S., Buters, J., Damialis, A., Heinrich, J., Kabesch, M., Nowak, D., Walser-Reichenbach, S.,
450 Weinberger, A., Zamfir, M., Herr, C., Kutzora, S., and Heinze, S.: A systematic review of threshold
451 values of pollen concentrations for symptoms of allergy, *Aerobiologia*, 37, 395–424,
452 <https://doi.org/10.1007/s10453-021-09709-4>, 2021.
- 453 Swanson, B., Freeman, M., Rezgui, S., and Huffman, J. A.: Pollen classification using a single particle
454 fluorescence spectroscopy technique, *Aerosol Science and Technology*, 57, 112–133,
455 <https://doi.org/10.1080/02786826.2022.2142510>, 2023.
- 456 Tummon, F., Adams-Groom, B., Antunes, C. M., Bruffaerts, N., Buters, J., Cariñanos, P., Celenk, S.,
457 Choël, M., Clot, B., Cristofori, A., Crouzy, B., Damialis, A., Fernández, A. R., González, D. F., Galán, C.,
458 Gedda, B., Gehrig, R., Gonzalez-Alonso, M., Gottardini, E., Gros-Daillon, J., Hajkova, L., O'Connor, D.,
459 Östensson, P., Oteros, J., Pauling, A., Pérez-Badia, R., Rodinkova, V., Rodríguez-Rajo, F. J., Ribeiro, H.,
460 Sauliene, I., Sikoparija, B., Skjøth, C. A., Spanu, A., Sofiev, M., Sozinova, O., Srncic, L., Visez, N., and
461 De Weger, L. A.: The role of automatic pollen and fungal spore monitoring across major end-user
462 domains, *Aerobiologia*, 40, 57–75, <https://doi.org/10.1007/s10453-024-09820-2>, 2024.
- 463 de Weger, L. A., Bergmann, K. Ch., Rantio-Lehtimäki, A., Dahl, A., Buters, J., Déchamp, C., Belmonte,
464 J., Thibaudon, M., Cecchi, L., Besancenot, J.-P., Galán, C., and Waisel, Y.: Impact of Pollen, in:
465 *Allergenic Pollen: A Review of the Production, Release, Distribution and Health Impacts*, 161,203, 2013.
- 466 Zhang, G. and Abdulla, W.: Identifying Pollen Species Using Multispectral Imaging Flow Cytometry and
467 Neural Networks, <https://doi.org/10.2139/ssrn.4375939>, 2023.
- 468 Zhang, Y. and Steiner, A. L.: Projected climate-driven changes in pollen emission season length and
469 magnitude over the continental United States, *Nat Commun*, 13, 1234, [https://doi.org/10.1038/s41467-](https://doi.org/10.1038/s41467-022-28764-0)
470 [022-28764-0](https://doi.org/10.1038/s41467-022-28764-0), 2022.



471 Ziska, L. H., Makra, L., Harry, S. K., Bruffaerts, N., Hendrickx, M., Coates, F., Saarto, A., Thibaudon,
472 M., Oliver, G., Damialis, A., Charalampopoulos, A., Vokou, D., Heidmarsson, S., Gudjohnsen, E., Bonini,
473 M., Oh, J.-W., Sullivan, K., Ford, L., Brooks, G. D., Myszkowska, D., Severova, E., Gehrig, R., Ramón,
474 G. D., Beggs, P. J., Knowlton, K., and Crimmins, A. R.: Temperature-related changes in airborne
475 allergenic pollen abundance and seasonality across the northern hemisphere: a retrospective data analysis,
476 The Lancet Planetary Health, 3, e124–e131, [https://doi.org/10.1016/S2542-5196\(19\)30015-4](https://doi.org/10.1016/S2542-5196(19)30015-4), 2019.