

---

# PI@ntNet-300K: a new plant image dataset for the evaluation of set-valued classification (supplementary material)

---

Camille Garcin<sup>\*1</sup>, Alexis Joly<sup>†2</sup>, Pierre Bonnet<sup>‡3</sup>, Antoine Affouard<sup>§2,3</sup>, Jean-Christophe Lombardo<sup>¶2,3</sup>, Maximilien Servajean<sup>||4</sup>, and Joseph Salmon<sup>\*\*5</sup>

<sup>1</sup>IMAG, Univ Montpellier, Inria, CNRS, Montpellier, France

<sup>2</sup>Inria, LIRMM, Univ Montpellier, CNRS, Montpellier, France

<sup>3</sup>CIRAD, AMAP

<sup>4</sup>LIRMM, AMIS, UPVM, Univ Montpellier, CNRS, Montpellier

<sup>5</sup>IMAG, Univ Montpellier, CNRS, Montpellier, France

## 1 URL to download the dataset

<https://doi.org/10.5281/zenodo.4726653>

## 2 Motivation

- For what purpose was the dataset created? Was there a specific task in mind? Was there a specific gap that needed to be filled? Please provide a description.  
**PI@ntNet-300k dataset was created to evaluate set-valued classification, in particular for plant identification. Unlike previous datasets, PI@ntNet-300k is designed so as to preserve the high level of ambiguity across classes of the initial real-world dataset (PI@ntNet) as well as its long tail distribution.**
- Who created the dataset (e.g., which team, research group) and on behalf of which entity (e.g., company, institution, organization)?  
**The dataset was created by PI@ntNet team, PI@ntNet being a consortium composed of four French research organisms (Inria, INRAE, CIRAD and IRD).**
- Who funded the creation of the dataset? If there is an associated grant, please provide the name of the grantor and the grant name and number.  
**The creation of the dataset was funded by the European Union's Horizon 2020 research and innovation program under grant agreement No 863463 (Cos4Cloud project) and by the French national research agency under the grant agreement ANR-20-CHIA-0001-01 (CaMeLOt project). PI@ntNet has also received the support of Agropolis Fondation for the platform creation.**

---

\*camille.garcin@inria.fr

†alexis.joly@inria.fr

‡pierre.bonnet@cirad.fr

§antoine.affouard@cirad.fr

¶jean-christophe.lombardo@inria.fr

||servajean@lirmm.fr

\*\*joseph.salmon@umontpellier.fr

### 21 3 Composition

- 22 • What do the instances that comprise the dataset represent (e.g., documents, photos, people,  
23 countries)? Are there multiple types of instances (e.g., movies, users, and ratings; people  
24 and interactions between them; nodes and edges)? Please provide a description.  
25 **The dataset is composed of pictures of plants. We are in the multi-class classification**  
26 **setting: there is a single plant species per image.**
- 27 • How many instances are there in total (of each type, if appropriate)?  
28 **There are 306,146 plant images : 243,916 in the training set, 31,118 in the validation**  
29 **set and 31,112 in the test set.**
- 30 • Does the dataset contain all possible instances or is it a sample (not necessarily random) of  
31 instances from a larger set? If the dataset is a sample, then what is the larger set? Is the  
32 sample representative of the larger set (e.g., geographic coverage)? If so, please describe  
33 how this representativeness was validated/verified. If it is not representative of the larger  
34 set, please describe why not (e.g., to cover a more diverse range of instances, because in-  
35 stances were withheld or unavailable).  
36 **The dataset is sampled from a larger set such that two particular features are pre-**  
37 **served. These features are inherent to the way the images are acquired and to the**  
38 **intrinsic diversity of plants morphology: i) The dataset exhibits a strong class imbal-**  
39 **ance, meaning that a few species represent most of the images. ii) Many species are**  
40 **visually similar, making identification difficult even for the expert eye. More details**  
41 **about these properties are available in Section 3.**
- 42 • What data does each instance consist of? “Raw” data (e.g., unprocessed text or images) or  
43 features? In either case, please provide a description.  
44 **Each instance is an image of a single plant.**
- 45 • Is there a label or target associated with each instance? If so, please provide a description.  
46 **Each instance is associated to its species.**
- 47 • Is any information missing from individual instances? If so, please provide a description,  
48 explaining why this information is missing (e.g., because it was unavailable). This does not  
49 include intentionally removed information, but might include, e.g., redacted text.  
50 **There is no missing information.**
- 51 • Are relationships between individual instances made explicit (e.g., users’ movie ratings,  
52 social network links)? If so, please describe how these relationships are made explicit.  
53 **There is no particular relationships between our instances.**
- 54 • Are there recommended data splits (e.g., training, development/validation, testing)? If so,  
55 please provide a description of these splits, explaining the rationale behind them.  
56 **The dataset already provides a train/validation/test. For more detail see section 3.1.**
- 57 • Are there any errors, sources of noise, or redundancies in the dataset? If so, please provide  
58 a description.  
59 **Only PI@ntNet observations with a valid species name were included the dataset. The**  
60 **species name validation is based on a Bayesian inference taking as input the names**  
61 **proposed by PI@ntNet users with a principle of adaptive weights depending on the**  
62 **user’s expertise.**
- 63 • Is the dataset self-contained, or does it link to or otherwise rely on external resources (e.g.,  
64 websites, tweets, other datasets)? If it links to or relies on external resources, a) are there  
65 guarantees that they will exist, and remain constant, over time; b) are there official archival  
66 versions of the complete dataset (i.e., including the external resources as they existed at the  
67 time the dataset was created); c) are there any restrictions (e.g., licenses, fees) associated  
68 with any of the external resources that might apply to a future user? Please provide descrip-  
69 tions of all external resources and any restrictions associated with them, as well as links or  
70 other access points, as appropriate.  
71 **The dataset is self contained.**
- 72 • Does the dataset contain data that might be considered confidential (e.g., data that is pro-  
73 tected by legal privilege or by doctorpatient confidentiality, data that includes the content  
74 of individuals’ non-public communications)? If so, please provide a description.  
75 **No protected data are available in the paper.**

- 76 • Does the dataset contain data that, if viewed directly, might be offensive, insulting, threat-  
77 ening, or might otherwise cause anxiety? If so, please describe why.  
78 **No, the dataset only contains plant pictures.**
- 79 • Does the dataset relate to people? If not, you may skip the remaining questions in this  
80 section.  
81 **No.**
- 82 • Does the dataset identify any subpopulations (e.g., by age, gender)? If so, please describe  
83 how these subpopulations are identified and provide a description of their respective distri-  
84 butions within the dataset.  
85 **Irrelevant.**
- 86 • Is it possible to identify individuals (i.e., one or more natural persons), either directly or  
87 indirectly (i.e., in combination with other data) from the dataset? If so, please describe  
88 how.  
89 **Irrelevant.**
- 90 • Does the dataset contain data that might be considered sensitive in any way (e.g., data that  
91 reveals racial or ethnic origins, sexual orientations, religious beliefs, political opinions or  
92 union memberships, or locations; financial or health data; biometric or genetic data; forms  
93 of government identification, such as social security numbers; criminal history)? If so,  
94 please provide a description.  
95 **Irrelevant.**
- 96 • Any other comments?

#### 97 4 Collection Process

- 98 • How was the data associated with each instance acquired?  
99  
100 Was the data directly observable (e.g., raw text, movie ratings), reported by subjects  
101 (e.g., survey responses), or indirectly inferred/derived from other data (e.g., part-of-speech  
102 tags, model-based guesses for age or language)? **Each image comes from the picture of**  
103 **a plant taken by a user of the PI@ntNet application**  
104
- 105 If data was reported by subjects or indirectly inferred/derived from other data, was the data  
106 validated/verified? If so, please describe how.  
107 **Only PI@ntNet observations with a valid species name were included the dataset. The**  
108 **species name validation is based on a Bayesian inference taking as input the names**  
109 **proposed by PI@ntNet users with a principle of adaptive weights depending on the**  
110 **user's expertise.**
- 111 • What mechanisms or procedures were used to collect the data (e.g., hardware apparatus or  
112 sensor, manual human curation, software program, software API)?  
113 **The data was collected through PI@ntNet mobile application and curated through**  
114 **crowdsourcing (by PI@ntNet users) in addition to the automated filtering (CNN-**  
115 **based) of unappropriated or irrelevant content (faces, humans, animals, buildings,**  
116 **etc.).**  
117
- 118 How were these mechanisms or procedures validated?  
119 **The mechanisms were validated by PI@ntNet curators (expert botanists) and by the**  
120 **scientific and technical committee of PI@ntNet.**
- 121 • If the dataset is a sample from a larger set, what was the sampling strategy (e.g., deter-  
122 ministic, probabilistic with specific sampling probabilities)? **The sampling is done at the**  
123 **genus level : 10% of the genus are randomly sampled, and all images that belong to**  
124 **these genera are kept. As a last step, we only retained species with at least 4 images.**
- 125 • Who was involved in the data collection process (e.g., students, crowdworkers, contractors)  
126 and how were they compensated (e.g., how much were crowdworkers paid)? **The data**  
127 **is collected by users of the PI@ntNet application (which has more than 10 millions**  
128 **users). PI@ntNet users are citizen scientist who gracefully participate to the project.**

129 **Their reward is the acclaimed performance of the application which enables them to**  
130 **identify plant species.**

- 131 • Over what timeframe was the data collected? Does this timeframe match the creation  
132 timeframe of the data associated with the instances (e.g., recent crawl of old news articles)?  
133 If not, please describe the timeframe in which the data associated with the instances was  
134 created. **The dataset was created with images collected by the Plantnet application**  
135 **from 2011 up to November 2020.**
- 136 • Were any ethical review processes conducted (e.g., by an institutional review board)? If  
137 so, please provide a description of these review processes, including the outcomes, as well  
138 as a link or other access point to any supporting documentation. **An ethical review was**  
139 **processed by CIRAD’s institutional review board. The main outcome was the terms**  
140 **of use of PI@ntNet application ([https://api.plantnet.org/views/terms\\_of\\_use?lang=en](https://api.plantnet.org/views/terms_of_use?lang=en))**
- 141 • Does the dataset relate to people? **No.** If not, you may skip the remainder of the questions  
142 in this section.
- 143 • Did you collect the data from the individuals in question directly, or obtain it via third  
144 parties or other sources (e.g., websites)? **not applicable**
- 145 • Were the individuals in question notified about the data collection? If so, please describe  
146 (or show with screenshots or other information) how notice was provided, and provide a  
147 link or other access point to, or otherwise reproduce, the exact language of the notification  
148 itself. **not applicable**
- 149 • Did the individuals in question consent to the collection and use of their data? If so, please  
150 describe (or show with screenshots or other information) how consent was requested and  
151 provided, and provide a link or other access point to, or otherwise reproduce, the exact  
152 language to which the individuals consented. **not applicable**
- 153 • If consent was obtained, were the consenting individuals provided with a mechanism to  
154 revoke their consent in the future or for certain uses? If so, please provide a description, as  
155 well as a link or other access point to the mechanism (if appropriate). **not applicable**
- 156 • Has an analysis of the potential impact of the dataset and its use on data subjects (e.g., a  
157 data protection impact analysis) been conducted? If so, please provide a description of this  
158 analysis, including the outcomes, as well as a link or other access point to any supporting  
159 documentation. **not applicable**

## 160 **5 Preprocessing/cleaning/labeling**

- 161 • Was any preprocessing/cleaning/labeling of the data done (e.g., discretization or bucketing,  
162 tokenization, part-of-speech tagging, SIFT feature extraction, removal of instances, pro-  
163 cessing of missing values)? If so, please provide a description. If not, you may skip the  
164 remainder of the questions in this section.  
165 **No pre-preprocessing was applied (apart from the curation process, see previous**  
166 **section).**
- 167 • Was the “raw” data saved in addition to the preprocessed/cleaned/labeled data (e.g., to  
168 support unanticipated future uses)? If so, please provide a link or other access point to the  
169 “raw” data.  
170 **not applicable**
- 171 • Is the software used to preprocess/clean/label the instances available? If so, please provide  
172 a link or other access point.  
173 **No.**

## 174 **6 Uses**

- 175 • Has the dataset been used for any tasks already? **Not this specific PI@ntNet subset.** If so,  
176 please provide a description.
- 177 • Is there a repository that links to any or all papers or systems that use the dataset? If so,  
178 please provide a link or other access point. **The list all or some papers that use our**

179 **dataset will be displayed and updated at the following address: [https://github.com/](https://github.com/plantnet/PlantNet-300K/)**  
180 **[plantnet/PlantNet-300K/](https://github.com/plantnet/PlantNet-300K/)**

- 181 • What (other) tasks could the dataset be used for?  
182 **The dataset can be used for any supervised or unsupervised classification tasks.**
- 183 • Is there anything about the composition of the dataset or the way it was collected and pre-  
184 processed/cleaned/labeled that might impact future uses? For example, is there anything  
185 that a future user might need to know to avoid uses that could result in unfair treatment  
186 of individuals or groups (e.g., stereotyping, quality of service issues) or other undesirable  
187 harms (e.g., financial harms, legal risks) If so, please provide a description. Is there any-  
188 thing a future user could do to mitigate these undesirable harms? **No**
- 189 • Are there tasks for which the dataset should not be used? If so, please provide a description.  
190 **No**

## 191 7 Distribution

- 192 • Will the dataset be distributed to third parties outside of the entity (e.g., company, insti-  
193 tution, organization) on behalf of which the dataset was created? If so, please provide a  
194 description. **the dataset will be publicly available**
- 195 • How will the dataset will be distributed (e.g., tarball on website, API, GitHub)? Does the  
196 dataset have a digital object identifier (DOI)? **the dataset will be distributed through**  
197 **zenodo under doi: <https://doi.org/10.5281/zenodo.4726653>**
- 198 • When will the dataset be distributed? **the dataset will be distributed after acceptance of**  
199 **the paper**
- 200 • Will the dataset be distributed under a copyright or other intellectual property (IP) license,  
201 and/or under applicable terms of use (ToU)? If so, please describe this license and/or ToU,  
202 and provide a link or other access point to, or otherwise reproduce, any relevant licensing  
203 terms or ToU, as well as any fees associated with these restrictions.  
204 **The dataset and all images composing it will be distributed under Creative-Common**  
205 **Attribution-ShareAlike 2.0 license.**
- 206 • Have any third parties imposed IP-based or other restrictions on the data associated with the  
207 instances? If so, please describe these restrictions, and provide a link or other access point  
208 to, or otherwise reproduce, any relevant licensing terms, as well as any fees associated with  
209 these restrictions.  
210 **No.**
- 211 • Do any export controls or other regulatory restrictions apply to the dataset or to individual  
212 instances? If so, please describe these restrictions, and provide a link or other access point  
213 to, or otherwise reproduce, any supporting documentation.  
214 **No.**

## 215 8 Maintenance

- 216 • Who is supporting/hosting/maintaining the dataset?  
217 **The Pl@ntnet team will maintain the dataset and provide support. The dataset is**  
218 **hosted by <http://zenodo.org>.**
- 219 • How can the owner/curator/manager of the dataset be contacted (e.g., email address)?  
220 **The owner/manager of the dataset can be contacted by mail at [plantnet-300k@inria.](mailto:plantnet-300k@inria.fr)**  
221 **fr.**
- 222 • Is there an erratum? If so, please provide a link or other access point.  
223 **Zenodo will provide a versioning of any correction of the dataset. We will keep the**  
224 **users informed at <https://github.com/plantnet/PlantNet-300K/>**
- 225 • Will the dataset be updated (e.g., to correct labeling errors, add new instances, delete ins-  
226 tances)? If so, please describe how often, by whom, and how updates will be communi-  
227 cated to users (e.g., mailing list, GitHub)?  
228 **The dataset will be updated if errors are spotted. The update will be performed by the**

- 229 **Plantnet team, and these modifications will be listed at [https://github.com/plantnet/](https://github.com/plantnet/PlantNet-300K/)**  
230 **[PlantNet-300K/](https://github.com/plantnet/PlantNet-300K/)**
- 231 • If the dataset relates to people, are there applicable limits on the retention of the data as-  
232 sociated with the instances (e.g., were individuals in question told that their data would be  
233 retained for a fixed period of time and then deleted)? If so, please describe these limits and  
234 explain how they will be enforced.  
235 **Irrelevant.**
  - 236 • Will older versions of the dataset continue to be supported/hosted/maintained? If so, please  
237 describe how. If not, please describe how its obsolescence will be communicated to users.  
238 **Zenodo will provide a versioning of any correction of the dataset.**
  - 239 • If others want to extend/augment/build on/contribute to the dataset, is there a mechanism  
240 for them to do so? If so, please provide a description. Will these contributions be vali-  
241 dated/verified? If so, please describe how. If not, why not? Is there a process for commu-  
242 nicating/distributing these contributions to other users? If so, please provide a description.  
243 **No.**
  - 244 • Any other comments?

## 245 **9 Author statement**

246 The authors confirm that all data in Pl@ntNet-300K dataset are under a Creative-Common  
247 Attribution-ShareAlike 2.0 license (see terms of use here) and bear responsibility in case of vio-  
248 lation of copyrights.