



**HAL**  
open science

# Hybridization of deep and prototypical neural network for rare defect classification on aircraft fuselage images acquired by an unmanned aerial vehicle

Julien Miranda, Jannic Veith, Stanislas Larnier, Ariane Herbulot, Michel Devy

## ► To cite this version:

Julien Miranda, Jannic Veith, Stanislas Larnier, Ariane Herbulot, Michel Devy. Hybridization of deep and prototypical neural network for rare defect classification on aircraft fuselage images acquired by an unmanned aerial vehicle. *Journal of Electronic Imaging*, 2020, 29 (04), pp.1-10.1117/1.JEI.29.4.041010 . hal-04931079

**HAL Id: hal-04931079**

**<https://laas.hal.science/hal-04931079v1>**

Submitted on 5 Feb 2025

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Copyright

# 1 Hybridization of deep and prototypical neural network for rare 2 defect classification on aircraft fuselage images acquired by an UAV

3 **Julien Miranda**<sup>a,b,c,\*</sup>, **Jannic Veith**<sup>c,d</sup>, **Stanislas Larnier**<sup>c</sup>, **Ariane Herbulot**<sup>a,b</sup>, **Michel Devy**<sup>a,b</sup>

4 <sup>a</sup>LAAS, CNRS, 7 avenue du colonel Roche, F-31400 Toulouse, France

5 <sup>b</sup>Univ de Toulouse, UPS, LAAS, F-31400 Toulouse, France

6 <sup>c</sup>Donecle, 201 Rue Pierre et Marie Curie, F-31670 Labège, France

7 <sup>d</sup>Swiss Federal Institute of Technology-ETHZ, CH-8092 Zürich, Switzerland

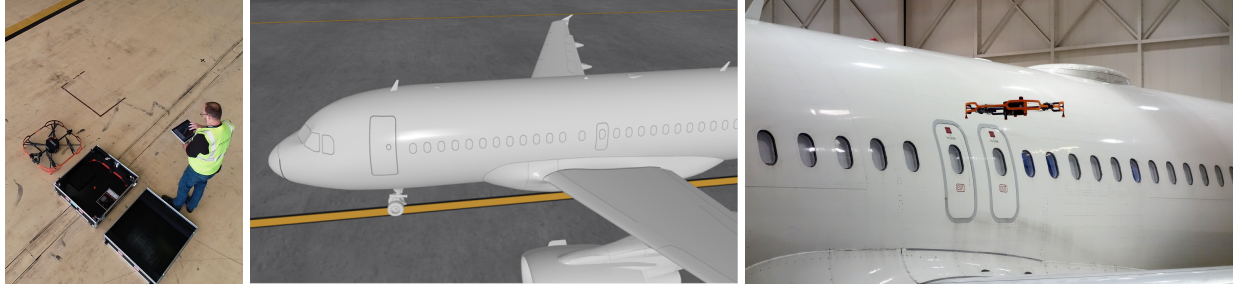
8 **Abstract.** In order to ease visual inspections of exterior aircraft fuselage, new technical approaches have been  
9 recently deployed. Automated UAVs are now acquiring high quality images of the aircraft in order to perform offline  
10 analysis. At first, some acquisitions are annotated by human operators in order to provide a large dataset required  
11 to train machine learning methods, especially for critical defects detection. An intrinsic problem of this dataset is  
12 its extreme imbalance (i.e there is an unequal distribution between classes). The rarest and most valuable samples  
13 represent few elements among thousands of annotated objects. Deep Learning-only based approaches have proven to  
14 be very effective when a sufficient amount of data is available for each desired class, whereas less complex systems  
15 such as Support Vector Machine theoretically need less data, and Few-Shot Learning dedicated methods (Matching  
16 Network, Prototypical Network, etc.) can learn from only few examples. Those approaches are compared on our  
17 applicative case. Preliminary results show the existence of empirical frontiers in term of training dataset volume that  
18 indicate which approach might be promoted. We propose a method to combine different approaches in order to achieve  
19 best performances on defect classification, that is an extension of previous work.<sup>1</sup>

20 **Keywords:** Deep Learning, Few-Shot Learning, Hybrid Model, Defect Detection, Support Vector Machine, Visual  
21 Inspection.

22 \*Julien Miranda: [jmiranda@laas.fr](mailto:jmiranda@laas.fr)

## 23 1 Introduction

24 Visual inspections are one of the most common operations for aircraft maintenance. A major  
25 inspection task, performed by maintenance operators, consists in detecting defects on an aircraft  
26 fuselage. To do this, they must use mobile elevating platforms to reach positions from where they  
27 can properly observe the aircraft skin, looking for those defects. To make those inspections faster,  
28 more effective and less painful for human experts, mobile platforms can be used,<sup>2,3</sup> especially  
29 automated UAV deployed by the French start-up Donecle, for which a localization with respect to  
30 the aircraft can be accurate to a few centimeters, using 3D models as shown in Figure 1.

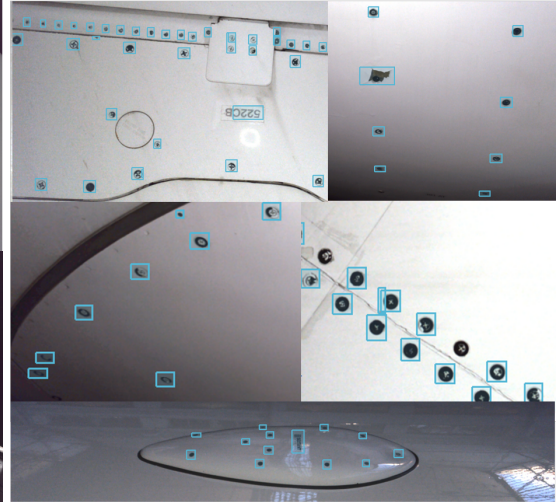


**Fig 1** Automated drone inspection from left to right: drone with a tablet running the analysis software application, 3D model used for autonomous localization, drone inspecting aircraft.

31 Thus using automated drone to acquire images on the whole surface opens new perspectives for  
32 aircraft maintenance traceability and automation.<sup>1</sup> Visual inspection from images can be organized  
33 in two main tasks: object detection (in our context, an "object" corresponds to a salient image  
34 region on the fuselage), and object classification. In this paper we focus on the second task, using  
35 of state-of-the art object detection methods: automatic visual inspection relies on the use of a Deep  
36 Neural Network (DNN) as object detector<sup>4</sup> that performs well enough for our application (some  
37 detections results are displayed in Figure 3). Figure 2 contains images acquired by the drone under  
38 various conditions representative of the variability of the inputs and the difficulties encountered:  
39 **top views** (top, left) show that images often contain several specular areas due to the external  
40 lightning of the hangar. This difficulty is minimized by the image overlapping (about 30%) which  
41 allows to see the areas under reflection on the next acquisitions. **Bottom views** (top, right) show  
42 images that are much less bright. These variations in lightness illumination are mitigated by the  
43 presence of an on-board lightning device (LED rings). In addition, as shown in 3 some scenes  
44 contain almost no objects (bottom, left), while others contain dozens of them (bottom, right).



**Fig 2** Images acquired by drone.



**Fig 3** Detected objects.

45 Growing popularity of Deep Learning (DL) methods has led to great advances in Computer  
46 Vision during past years: more specifically, image classification has become a relatively simple  
47 problem provided there are enough available data to train deep models (they overpass human per-  
48 formances for this task since 2015). ImageNet<sup>5</sup> and CIFAR<sup>6</sup> challenges give to the research com-  
49 munity ways to compare and improve their algorithms on public datasets. However in real world,  
50 accessing those data is often an impassable barrier. For the target application, defects such as light-  
51 ning burns are not frequent enough to envision a classic DL approach, while a lot of other objects  
52 that have to be discriminated are very common. Thus, this paper concerns object classification,  
53 considering very high imbalance ratio (1:5000) between classes.

54 We first establish that with extreme class imbalance ratio, state-of-the-art methods are not suf-  
55 ficient and that there is a need to take advantage from both the power of big data algorithms and  
56 from more specific methods dealing with few data for some classes. To do so, we demonstrate that  
57 with our industrial data, different machine learning approaches are relevant for different volumes  
58 of balanced training set. Then we evaluate those models on extremely imbalanced datasets and



59 propose a method to combine models into an hybrid classifier able to deal with common objects  
60 as well as very rare ones. Finally, this hybrid strategy is also validated and characterized on public  
61 datasets, modified to create different imbalance ratios.

62 Section 2 describes our context: how data are acquired, what are the objects and the classes. In  
63 Section 3 we describe the machine learning approaches that could be evaluated for our application  
64 and justify the choice of Prototypical Network (see Table 2), before showing their limitations  
65 for high imbalance ratio between classes in Section 4. Finally we propose an hybrid method in  
66 Section 5 and compare results with others approaches on an imbalanced dataset. Those sections  
67 refer to our preliminary work published in QCAV conference proceedings.<sup>7</sup> The next ones are  
68 new material and results: based on the previously justified notion of hybridization, we propose to  
69 deepen these methods through general heuristics in Section 6, then by considering multiple views  
70 of the same object in Section 6.1.3. Section 7 is dedicated to the experimentations and analysis  
71 of the results obtained by adding these new material. Finally, Section 8 will discuss possible  
72 improvements for the classification of rare defects and new possibilities offered by the proposed  
73 new system, which go beyond the classification framework.

## 74 **2 Acquisition and dataset**

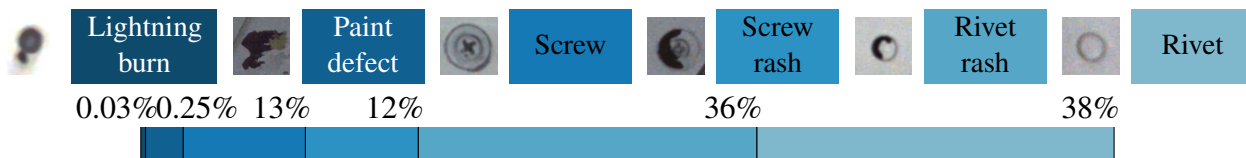
75 During inspection, images are taken by an autonomous UAV in order to cover the entire aircraft  
76 surface. For a typical aircraft, about 1,200 high definition (16 MP) images are required, with an  
77 overlap between two successive acquisitions ensuring that each zone is acquired at least twice.  
78 This reduces the negative impact of specular components and allows to fuse classification results  
79 on the same real object seen on several views. The matching of objects from different points of  
80 view is made possible by using the location data of the drone thus on the on-board camera, but will

81 not be covered in this paper.

82 Acquisitions are sent to a separate laptop or tablet that can process automated analysis using  
 83 GPUs. To create machine learning data sets, maintenance company experts annotate some batches  
 84 of these acquisitions.

85 All objects that can occur on the fuselage have to be registered for various applications com-  
 86 posing a general visual inspection (paint state evaluation, markings analysis, etc.). However using  
 87 such a network with high definition images do not allow very small objects classification (less than  
 88 1 mm<sup>2</sup>) in reasonable time whereas those objects are crucial as they can be critical defects (light-  
 89 ning burns). Thus potential defects (small or ambiguous objects) are gathered for a supplementary  
 90 classification step which is the object of this work.

91 Addressing the defect recognition as an image classification problem allows the use of ad-  
 92 vanced techniques that are much more complex to be introduce into a one-step object recognition  
 93 algorithm in which detection and classification are inseparable,<sup>4,8</sup> that perform the best result in  
 94 the state of the art on object recognition task (using mean average precision as metric).



**Fig 4** Unbalanced dataset description.

Class	Lightning burn	Paint defect	Screw	Screw rash	Rivet rash	Rivet	Total
Samples	11	91	4,758	4,496	13,228	13,959	<b>36,542</b>

**Table 1** Unbalanced dataset composition with number of samples by class.

95 The extremely unbalanced number of samples between classes is a specificity of our data distri-  
 96 bution compared to reference datasets. In this paper, we will consider a dataset with classes given

97 in Table 2 illustrated in Figure 4.

98 This database is not exhaustive and the image quality might not be representative of Donecle  
99 quality of acquisition but still gives an overview of the problem solved. In particular, several other  
100 types of objects are found on the fuselage of an aircraft and will not be mentioned here, for the  
101 sake of clarity and conciseness. Thus, while some classes are numerous enough to envision taking  
102 advantage of DL methods (screws or rivets), others are only represented by few samples (lightning  
103 burns) and may need dedicated Few-Shot Learning (FSL) methods. We splitted the original dataset  
104 into a train set (80%), a validation set (10%) used to tune hyper-parameters, and a test set (10%)  
105 used to obtain the presented results.

### 106 *2.1 Model complexity and required amount of training data*

107 How much data is needed for a given trainable model to perform well on real world data is a cru-  
108 cial question in Machine Learning. Finding an easy and accurate method to determine the required  
109 amount of data to reach a target generalization performance is the topic of many researches: sta-  
110 tistical learning theory has given some clues, introducing capacity measures for such algorithm,  
111 e.g. the Vapnik-Chervonenkis dimension,<sup>9</sup> from which generalization bounds can be applied to  
112 learning algorithms like SVM. However those bounds are vacuous for complex models such as  
113 DNNs.<sup>10</sup> Thus, empirical tests have been also performed to observe performances on classifica-  
114 tion tasks versus volume of training data: a logarithmic relationship seems to exist,<sup>11</sup> but might be  
115 subject to a potential diminishing return on log-scale.<sup>12</sup> Based on those empirical observations the  
116 required number of samples to reach good accuracy is below the number of parameters of a Deep  
117 Network, but it still needs a lot of images to be accurate. Transfer learning is a highly popular way  
118 to train models using representations learned from another task.<sup>13</sup>

## 119 2.2 Learning from imbalanced data

120 Several approaches can be used to cope with imbalanced datasets. **Data-level** method modifies the  
121 data by oversampling, under sampling, transforming or generating training samples. **Algorithm-**  
122 **level** approach tunes existing learning algorithms to adapt them to data with skewed distributions.  
123 **Hybrid methods** combine those two with the possible add of handcrafted rules or another algo-  
124 rithm.

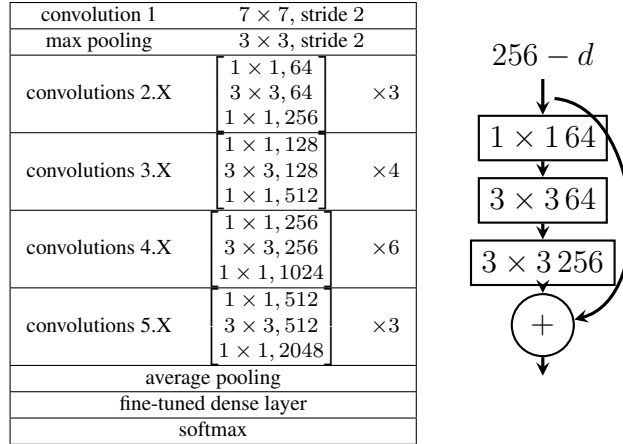
125 While existing works in class imbalance focus on imbalance ratios ranging from 1:4 up to  
126 1:100, classification with extreme imbalance ratio that we are facing, and that can be found as well  
127 in other applications (fraud detection, detection of dangerous behavior, etc.) remains a challenge.<sup>14</sup>  
128 Precisely, the rarest data are often the most valuable ones, like in the present work.

## 129 3 Existing Machine Learning approaches

### 130 3.1 Deep Neural Network

131 DNN are very efficient for the classification task on reference datasets, assuming a sufficient  
132 amount of data. Deep Convolutional Neural Networks (CNNs) have indeed reached high accuracy  
133 rate (exceeding 90%), among the most effective approaches: Wide Residual Network,<sup>15</sup> Fractional  
134 Max Pooling,<sup>16</sup> Dual Path Network<sup>17</sup> or other advanced methods.<sup>18–20</sup>

135 We trained those models and fine-tuned some popular networks (ResNet, Inception, etc.) with  
136 available pre-trained weights. The best accuracy on our validation set was achieved by a fine-  
137 tuned ResNet50 architecture so we used this model, described by Figure 5, as CNN baseline for  
138 our problem. It is composed of residual blocks with skip connections that have proven to be very  
139 efficient.<sup>21</sup>



**Fig 5** ResNet50 architecture (left) and residual block example (right).

140 All those DNNs are sensitive to data imbalance.<sup>22</sup> Common methods to tackle this issue are  
 141 data-based approaches, and consist in creating new samples using transformations on real images,  
 142 or generating random realistic samples with Generative Adversarial Networks (GANs).<sup>23</sup> Those  
 143 methods have proven to be efficient, but they usually do not apply in case of extreme imbalanced  
 144 dataset like ours.

### 145 3.2 Support Vector Machine

146 Another widely used learning approach, popular until the rise of DL hegemony in Computer Vision  
 147 is Support Vector Machine (SVM). It is a statistical learning approach that needs image represen-  
 148 tation as inputs. It was first used with hand-crafted descriptors, such as Histogram Of Gradient  
 149 (HOG).

150 We also tested SVM with representations learned from unsupervised learning, using a GAN  
 151 trained to generate realistic images. Recent works have shown that replacing the softmax layer of  
 152 a DNN by a SVM can give significant gain on classification datasets.<sup>24</sup> Moreover, data imbalance  
 153 can be integrated into the algorithm using class weights or cost-sensitive learning. SVMs are  
 154 sensitive to lack of data and to imbalanced classes as well as CNNs and can benefit from the same

155 data augmentation techniques.

156 On our test set SVM with HOG performs poorly (maximum accuracy is 0.76 while using the  
157 full training dataset), and provides good results on medium datasets (100 – 1000) when combined  
158 with learned representations (from pre-trained models and GAN). However it never outperformed  
159 our CNN baseline, so we did not include SVM on further comparisons.

### 160 3.3 Few-shot learning: algorithm-level approaches to face the lack of data

161 FSL algorithms are usually performing  $n$  shot,  $k$ -way learning, with  $n$  being the number of needed  
162 learning samples and  $k$  is the number of possible classes for a new sample during inference. Two  
163 main public datasets dedicated to FSL algorithm training and test have been created. Omniglot<sup>25</sup>  
164 is a dataset composed of handwritten characters from different alphabets and MiniImageNet is a  
165 subset of ImageNet. Samples are shown in Figures 6 and 7. Next, we describe some of most  
166 popular FSL algorithms.

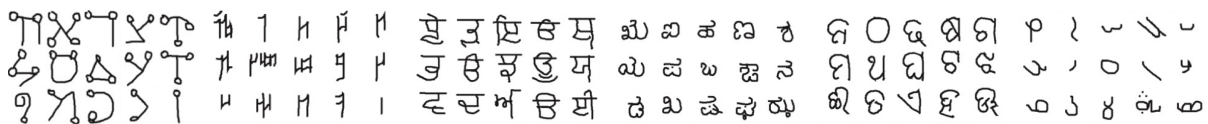


Fig 6 Omniglot samples from 6 alphabets.

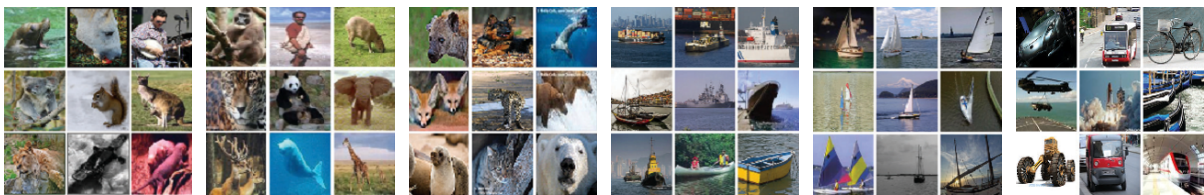


Fig 7 MiniImageNet samples from 6 classes.

167 To learn from few data, algorithm-level solutions have been proposed. Some of them are con-  
168 sidering meta-learning processes<sup>26–28</sup> while others focus on the ability to learn metric.<sup>29</sup>



169 **Siamese Networks.** This approach combines multiple networks.<sup>30</sup> To be called 'siamese', two  
170 networks have to share the same architecture (same layers with the same parameters) and to share  
171 learning loss and weights of the junction layer. The choice of the loss function is of crucial impor-  
172 tance. Some known examples are:

173 **Usual cross-entropy:**  $L = -y \log(p) + (1 - y) \log(1 - p)$

174 **Triplet-loss:**<sup>31</sup>  $L = \max(d(a, p) - d(a, n) + m, 0)$

With the triplet-loss formulation,  $d$  is the  $L2$  loss (or another distance function),  $a$  is the input sample from the dataset,  $p$  is a sample from the target class (randomly picked),  $n$  is a sample from another class.  $m$  is an hyper-parameter (margin). With this method, it is possible to perform metric learning, then to use a classical nearest-neighbor classification algorithm to separate classes. This last operation is not part of an end-to-end training process and so cannot be called optimized for the task.

**Matching Networks**<sup>32</sup> introduced few-shot networks based on the idea of making nearest neighbor algorithm learnable during training process by using it in a differentiable form. This allows to perform end-to-end fully optimized learning and is achieved by embedding a sample into a representation space and then performs a nearest-neighbor-like algorithm with the equation:

$$\hat{y} = \sum_{i=1}^k a(\hat{x}, x_i) y_i$$

175 where  $\hat{y}$  is the model prediction,  $x_i$  the support features,  $y_i$  the support labels,  $\hat{x}$  the query sample  
176 features and  $a$  is a similarity function. This approach is the first end-to-end few-shot dedicated  
177 model, and outperforms Siamese network for this task, see Table 2.

178 **Prototypical Network**<sup>29</sup> are built within the assumption that a single prototype per class can be

179 used to compute distance in the representation space. If the used distance is the Euclidean distance,  
 180 then the best prototype (in the sense that it minimizes the distance between a prototype and supports  
 181 points) is the mean of the support set representations.

182 **Model Agnostic Meta Learning**<sup>28</sup> (MAML) is a very different approach for FSL that uses a clever  
 183 initialization for deep models.<sup>18,33</sup> MAML is a way to learn how to initialize weights by optimizing  
 184 the generalization of the model. We took Prototypical Networks as few-shot algorithm baseline,  
 185 as they obtain the best performances on reference datasets and can dynamically perform k-way  
 186 classification.

Omniglot	5-way Accuracy		20-way Accuracy		MiniImageNet	5-way Accuracy	
	1-shot	5-shot	1-shot	5-shot		1-shot	5-shot
Siamese	97.3%	98.4%	88.2%	97.0%	Matching	44.2%	57.0%
Matching	98.1%	98.9%	93.8%	98.5%	Prototypical	<b>49.4%</b>	<b>68.2%</b>
Prototypical	<b>98.8%</b>	99.7%	<b>96.0%</b>	<b>98.9%</b>	MAML	48.07%	63.15%
MAML	98.7%	<b>99.9%</b>	95.8%	98.6%			

Table 2 Recent methods on FSL task.

187 We took the results from<sup>28</sup> and original papers, gathered them in Table 2.

#### 188 4 State of the art methods evaluation

189 Based on the state of the art, we tested the described approaches on our datasets : original training  
 190 dataset is sub-sampled to obtain training datasets of increasing size from 1 sample per class to  
 191 6,000 samples per class, first to evaluate which one performs the best for a given number of  
 192 training samples when there is no extreme imbalance. Results displayed on Figure 8 show that  
 193 these methods have good performances when the classes are limited ( $< 4$ ). With a large number  
 194 ( $\geq 4$ ), we found DL methods trained with large dataset (6,000 per classes) are more accurate. We  
 195 have also implemented classifiers built on a SVM fed by HOG descriptor as well as on a SVM  
 196 fed by ResNet features, that performed poorly on our test. We used prototypical networks, which

197 regarding the current state of the art metrics, see Table 2 appears to be the best alternative in terms  
198 of precision and recall for both the simple (Omniglot) and the hardest (minImageNet) datasets.

#### 199 4.1 Experiment

200 To evaluate model relevancy regions in term of real training dataset size, we used object classes  
201 that we have in sufficient number to train deep models, then create decreasing subsets from the  
202 original dataset to observe the effect on top-1, top-2 and top-3 accuracies, with top- $k$  accuracy  
203 being the proportion of samples for which the ground truth label is one of the  $k$  most probable  
204 predicted classes. We used common data augmentation techniques (horizontal and vertical flip,  
205 rotation, crop) to obtain presented results. For more general application of our work we also  
206 tested our method on public available datasets that where truncated in order to artificially create  
207 extremely imbalance datasets that are subsets of well known sets. For those dataset, we observe  
208 the performances on the randomly under-sampled class, which were randomly chosen.

#### 209 4.2 Results

210 Considering ResNet model, we found that the assertion made in Sun *et al.*<sup>11</sup> that "performance  
211 increases logarithmically based on volume of training data" fits our results in accuracy for top- $k$   
212 accuracy score in classification ( $k \in \mathbb{N}, k < N_c$ ) with  $N_c$  the number of classes, however we did  
213 not test this assumption for datasets larger than 6,000 elements per class. In Figure 8, dotted-lines  
214 represent log-regression on the obtained results.

215 We also observed the same relationship for FSL baseline for 2-way to  $N_c$ -way task. As ex-  
216 pected it appears than under a certain amount of available data, DL methods are less accurate than  
217 prototypical network: top-1 accuracy of ResNet is lower than 7-way prototypical classifier: for

218 those classes it would be preferable to use a FSL dedicated approach. Results are illustrated in  
 219 Figure 8: a frontier between two approaches is empirically set around 100 samples per class. The  
 220 problem that stands next is simple: we cannot know in advance which model will be the most  
 221 relevant for a given sample.

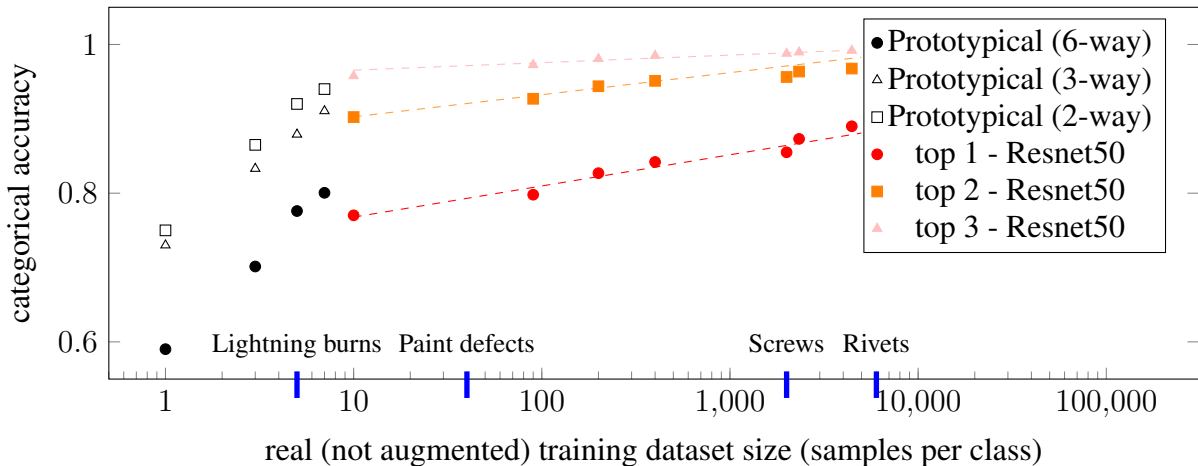


Fig 8 Different approaches performances regarding training dataset size.

## 222 5 Hybridization with Random Support Sampling (HRSS)

### 223 5.1 Description

224 Since CNN baseline top-3 accuracy giving good results, and 3-way Prototypical Network perform-  
 225 ing significantly better than 6-way (see Figure 8), we propose to combine classic DL model and  
 226 FSL approaches in an hybrid system to combine their advantages.

227 The idea of our approach is to first estimate which classes are the most likely ones for a given  
 228 input image, using a DL model with data augmentation. Then, in case one of the possible classes  
 229 is known to be highly under-represented, few-shot dedicated model is applied, and a combination  
 230 of the results from those two model gives the final output. If all the possible classes are known to  
 231 be well represented, the output is the CNN baseline output. We tested our approach and compared

232 it with other methods on our test dataset. To do so, we took a defect class with 10 training sam-  
 233 ples while other classes are trained with 3,000 samples before any augmentation. We compared  
 234 ResNet50 with data augmentation, Prototypical Network with (10-shot, 6-way), and our hybrid  
 235 method with top-3-way linking and two hybridization rules (few-shot wins and averaging outputs  
 236 of deep model and few-shot model).

## 237 5.2 Results

238 We observed global categorical accuracy, the precision, recall and average precision on the imbal-  
 239 anced defect class (varying the algorithm defect output probability), results are shown in Table 3.  
 240 We can see that hybrid method improves precision, recall and average precision (AP) for rare  
 defect class. The best performances were obtained using average output hybridization method.

Truncated dataset	Rare defect classification performances			
	Categorical accuracy	Precision	Recall	AP
<b>Hybrid Method (combined)</b>	<b>0.877</b>	<b>0.97</b>	<b>0.77</b>	<b>0.79</b>
Hybrid Method (FS)	0.867	0.94	0.71	0.71
Prototypical network	0.821	0.84	0.75	0.75
ResNet50	0.863	0.90	0.75	0.77

**Table 3** Classification results on imbalanced dataset.

241

242 The biggest gain is visible on precision score, which is very important for the addressed indus-  
 243 trial application with big and imbalanced dataset, because operators can only handle a reasonable  
 244 amount of false positive alerts. Nevertheless, the recall score is not high enough to constitute a  
 245 truly reliable aid to the operator. Nevertheless, the recall score is not high enough to constitute a  
 246 truly reliable aid to the operator. This is why, despite these encouraging initial results with regard  
 247 to the hybridization of methods, this naïve approach seems to have considerable room for improve-  
 248 ment. Indeed, the choice of the support element vectors of the highly represented classes was not  
 249 questioned, and the interfacing of the two models is naive.

## 250 6 Heuristics for class subsampling

251 In order to overcome the imbalance of classes in the context of a classification problem by ma-  
252 chine learning, some methods have been proposed in the literature<sup>34</sup> for oversampling as well as  
253 undersampling operations. The oversampling corresponds to the data augmentation process al-  
254 ready mentioned in the previous sections. We describe here some undersampling heuristics that  
255 can be applied to the most well represented class sets.

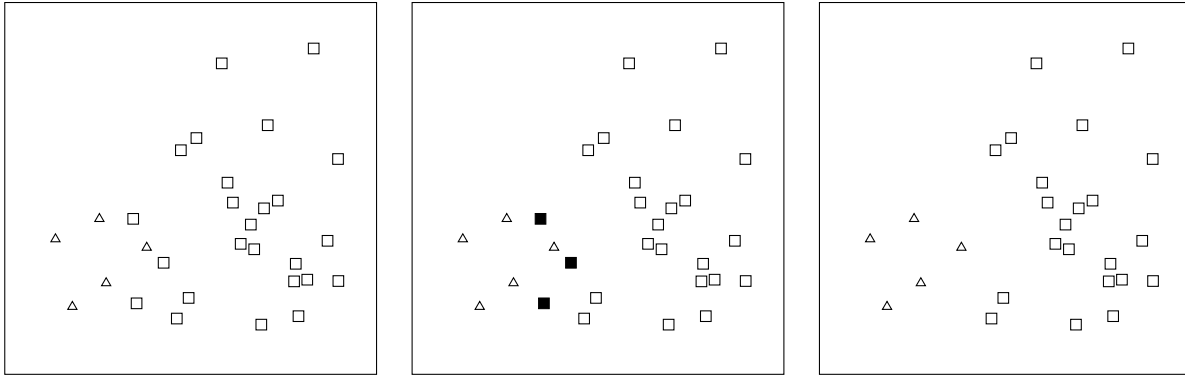
### 256 6.0.1 Tomek's links

257 A sub-sampling approach is to try to remove unrepresentative or too noisy samples. The notion  
258 of Tomek's link<sup>35</sup> can be used for this purpose. Given a metric space, two examples  $E_i$  and  $E_j$  of  
259 two different are a Tomek link if there is not an example  $E_k$ , such that  $d(E_i, E_j) > d(E_i, E_k)$  or  
260  $d(E_i, E_j) > d(E_j, E_k)$ . This concept can be used via two purposes:

- 261 • Undersampling: the elements of the majority classes that belong to a Tomek link are deleted.
- 262 • Cleaning: all elements that belong to a Tomek link are deleted.

263 The use of this technique is possible in our case by considering the representation space learned  
264 by the prototypical network provided with the Euclidean distance. Thus we can use this method to  
265 sub-sample over-represented classes.





**Fig 9** Undersampling iteration using Tomek's link: Two classes (left), sample of biggest class belonging to Tomek's link (center) and resulting datasets (right).

266 An iteration of sub-sampling over represented classes dataset using Tomek's link is illustrated  
 267 in Figure 9: The elements of over-represented classes are highlighted (in black) and then deleted.

### 268 6.0.2 Condensed Nearest Neighbour Rule

269 Another iterative approach is based on the progressive construction of sets, selecting only those  
 270 examples that cannot be explained by a simple classification based on elements already known.  
 271 The idea is that the information provided by the elements that are correctly classified by this sim-  
 272 ple approach is already largely contained in the previously selected examples. Condensed Nearest  
 273 Neighbour Rule proposed by Hart<sup>36</sup> and described in 1 is an algorithm reflecting this reasoning us-  
 274 ing neighbour algorithm as the naive classifier. It does not guarantee to find the smallest consistent  
 275 subset.

### 276 6.0.3 One-sided selection (OSS)

277 It is possible to combine the two previous approaches by first applying the Tomek rule (under  
 278 sampling) to eliminate noise and edges, then use the Condensed Nearest Neighbour algorithm.<sup>37</sup>

---

**Algorithm 1** CNN rule.

---

**Require:** Let  $E$  be the whole considered set  $E_{c1}$  be the under-represented class and  $E_{c2}$  be the over-represented class .

$\hat{E} \leftarrow E_{c1}$

Randomly select  $x \in E$

$\hat{E} \leftarrow \hat{E} \cup \{x\}$

Classify each remaining examples in  $E$  using 1- $k$ -NN. All misclassified elements are noted  $X_f$ .

**for all**  $x_f \in X_f$  **do**

$\hat{E} \leftarrow \hat{E} \cup \{x_f\}$

**end for**

---

279 As it was proven to be more efficient than only one of these heuristics,<sup>34</sup> we will consider this  
280 approach.

### 281 *6.1 Proposed methods*

282 The described above methods allow to under-sample the data under consistency consideration.  
283 However, two major pitfalls remain. First, they are likely to delete ambiguous examples, that  
284 might be noisy samples or outliers, but can also be interesting samples in crucial feature space  
285 zones. Second, there is still a significant amount of randomness in these algorithms. We propose  
286 a method to select the best samples that tackle those issues. The idea is to clusterize the potential  
287 support samples in the feature space provided by the prototypical neural network, then to keep  
288 only the closest elements to the prototype of each cluster. Using the prototype could be an option,  
289 but as it does not correspond to a real sample, this will increase the dependency on the encoding  
290 function (which is learnt), and the operation will not be under-sampling any more.

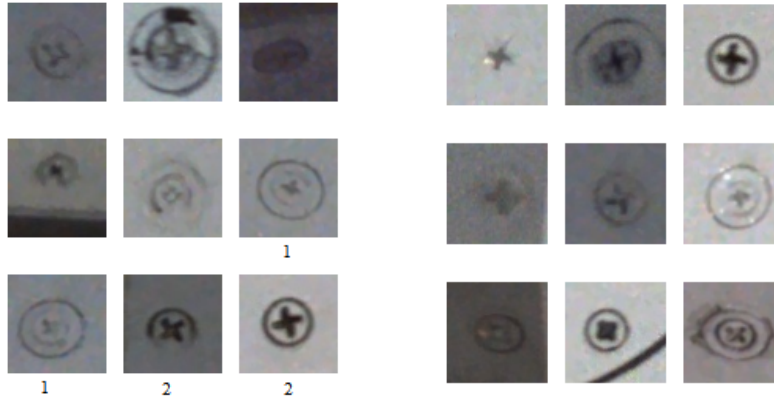
#### 291 *6.1.1 Meta-Data Filters on support dataset (MDF)*

292 Since Deep Learning methods require thousands of data to be accurate, it is necessary to gather  
293 learning examples from a variety of sources (different airline companies, different aircraft models,

294 indoor and outdoor acquisitions, etc.). On the contrary, FSL is based on a handful of examples  
295 and it is therefore possible to avoid these mixtures to automatically consider problem-specific  
296 databases. We introduce meta data filters to create those adapted subsets for all categories. To  
297 improve the synergy between the deep model and the FSL method, we can also select examples to  
298 use as supports between two classes among those that could not be separated by the deep model  
299 during the learning phase. This is done using the learning confusion matrix. This makes it possible  
300 to obtain a matrix of the support examples to be selected as a priority for the separation of two  
301 given classes.

### 302 6.1.2 Cluster-based Medoid Prototypes (CMP)

303 To prepare  $N$ -shot learning task, we need to create prototypes from support samples. As we ob-  
304 served that classification accuracy growth with  $N$ , we set  $N$  to be equal to  $N_{c1}$  the cardinal of the  
305 under-represented class. We want to select the  $N$  best samples among  $N_{c2}$ , with  $N_{c2} \gg N_{c1}$ .  
306 Prototypical Neural Network provides a feature space where euclidean distance can be used. We  
307 propose to create clusters in this space with Density-Based Spatial Clustering of Applications with  
308 Noise (DBSCAN) that is a very popular method,<sup>11, 38</sup> Once the clusters are created, we simply take  
309 the acquired sample that is the closest to the each cluster centroid (the medoid) as support sample  
310 for the class. The Figure 10 illustrates the difference between 9 randomly selected supports for  
311 the 'screw' class (left) where some screw models (numbered 1 and 2) can be represented more  
312 than once in the support set and medoids selection after clustering (right) where each support is a  
313 different type of screw.



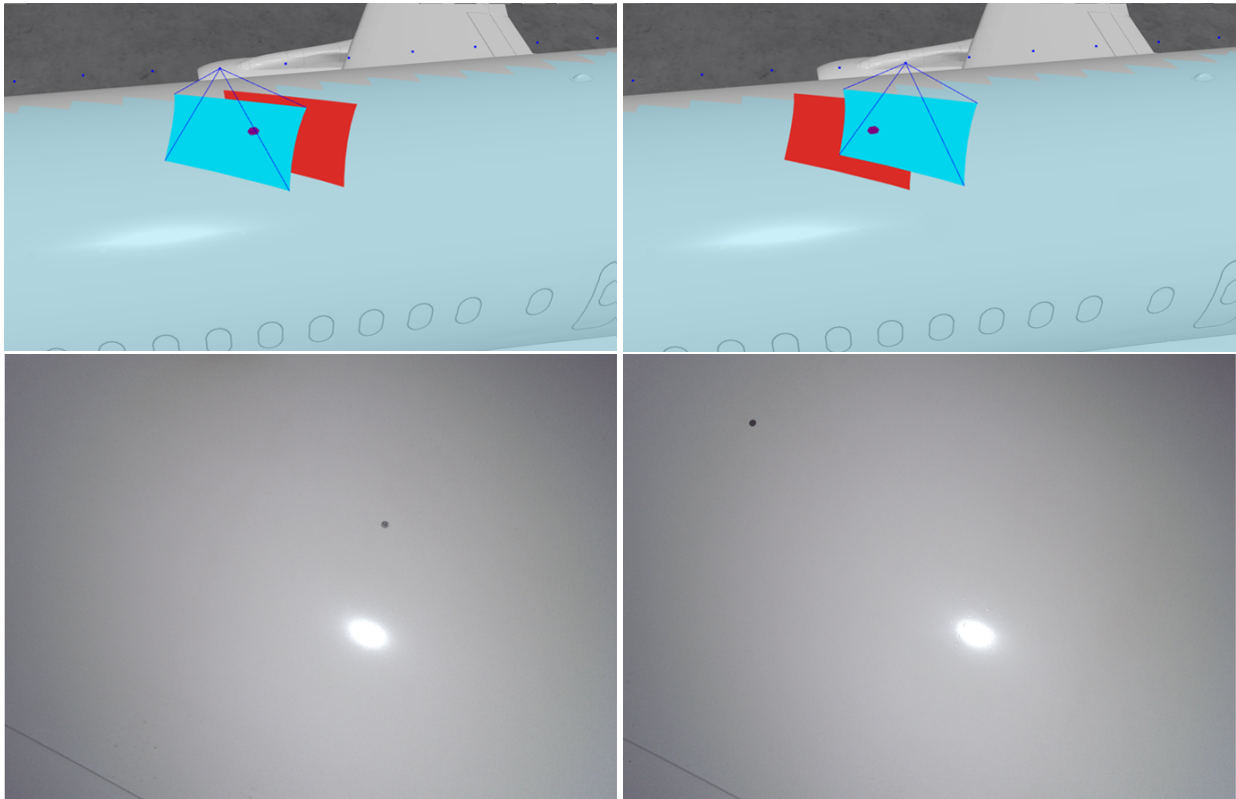
**Fig 10** 9 random screw supports (left) and 9 screw medoids selected after clustering (right).

314 The previously described heuristic can be used in many cases of extreme imbalance dataset.  
 315 Another aspect, taking advantage of the specificity of our acquisitions, which consists of multiple  
 316 shots of each fuselage area, is the linking of images corresponding to the same part of the aircraft.  
 317 The following section discusses this approach.

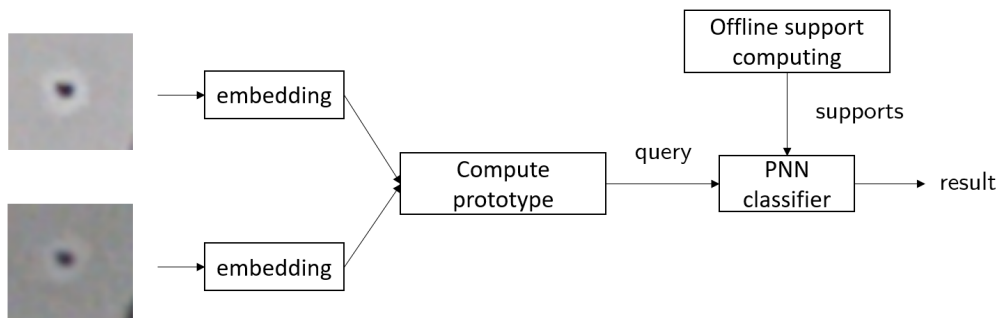
### 318 6.1.3 Sequence-Wise Prototypical Network (SWPN)

319 In this section we discuss the possibility to use several views of the same object to create a proto-  
 320 type in the prototypical feature space for further classification. It should be remembered that the  
 321 images acquired and all defects detected are geolocated relative to a 3D model of the aircraft: an  
 322 example of such a model is given in Figure 2. In addition, a controlled overlay allows each of  
 323 them to be visible on several of those. We propose to match these elements based on the position  
 324 in space and the type of object predicted by the Deep Neural Network (referred to is previous sec-  
 325 tions). We can thus form sequence of images of the same object. From this point we embed all  
 326 the elements of the sequence into the prototypical feature space, then create a query prototype to  
 327 be classified by dedicated FSL classifier. Two views of the same objects are shown in Figure 11:  
 328 here the objects are lightning strikes: in itself measures only a few square millimetres, but as it

329 was previously detected during a manual inspection, a black circular pad (much more visible in the  
 330 image) with an identification number was added to the fuselage. The sequence embedding process  
 331 is illustrated in Figure 12 with two extracted images of the same detected lighting strike.



**Fig 11** Lightning strike from two points of view.



**Fig 12** Sequential embedding for classification using prototypical neural network.

## 332 **7 Experiment and analysis**

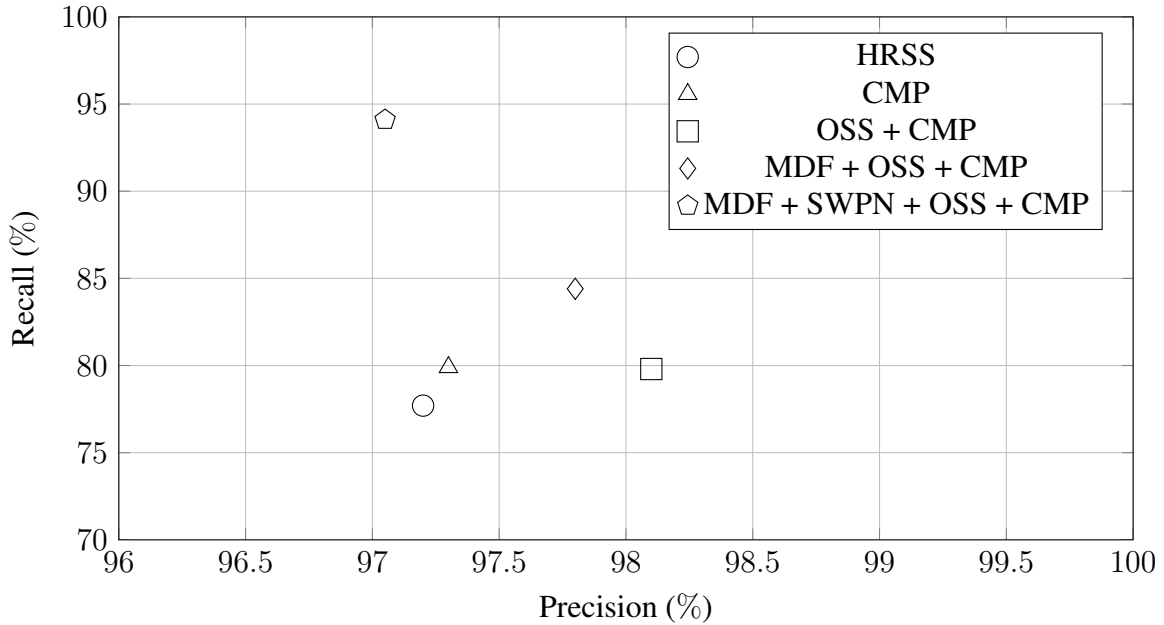
333 In this section, we compare the results obtained by new materials combinations under one consis-  
334 tent dataset: Hybridization with Random Support Sampling (HRSS), Cluster-based Medoid Pro-  
335 totypes (CMP), One-Sided pre-processing for Support dataset cleaning (OSS), Meta-Data Filters  
336 (MDF) on support dataset and Sequence-Wise Prototypical Network (SWPN).

337 We sequentially apply the described algorithms to combine them, taking care of the order  
338 of operations. For example, the combination MDF + SWPN + OSS + CMP is obtained by the  
339 following steps:

- 340 1. Filter the elements using the available meta data: here we use the airline and aircraft model  
341 that correspond to the acquisitions after noticing that many elements are specific to them.
- 342 2. Associate views of the same objects using the SWPN.
- 343 3. Use OSS to reduce the size of over-represented class sets.
- 344 4. Select support prototypes using CMP.

345 This allows to evaluate performance gain of each of the proposed hybridization scheme. We  
346 compare the Deep Learning and FSL baselines with several hybrid models on both CIFAR-10  
347 truncated dataset (we truncated each class then consider the mean average precision over those  
348 classes) and industrial dataset.





**Fig 13** Precision and Recall for different methods.

349 In this figure we can see the significant improvements made by the different additive modules.  
 350 In particular, the recall score is increased. It should be noted that the methods most specific to our  
 351 application, namely the use of meta-data or sequencing of requests for the prototypical network,  
 352 provide the most significant benefits. Nevertheless heuristics methods of sub-sampling which are  
 353 model-agnostic are also effective. We have therefore presented general and effective methods as  
 354 well as methods that are more specific to the application in question and even more effective.

## 355 **8 Conclusion and prospect**

356 In this paper, we analysed different machine learning approaches that give incomplete answer to  
 357 a practical use-case of FSL due to extreme data imbalanced. As in many applications,<sup>39,40</sup> the  
 358 rarest data are the most critical, so there is a need to be specifically accurate on those classes. We  
 359 compared those methods with different training dataset sizes with balanced classes, then with a

360 training dataset that contains one few-sampled defect class. We proposed an hybrid method that  
361 performs better on our dataset for the under represented defect classification.

362 Some drawbacks of our approach might lead to future works on the following aspects:

363 **Required data for training hybrid method.** By now we needed to train multiple models with  
364 sub-sampled training sets to define boundaries (in term of training samples) that delimits which  
365 model is appropriate for which class. Using the observed logarithmic relationship between train-  
366 ing set size and accuracy an analytical estimation of those boundaries might help reducing the  
367 number of needed training phases.

368 **Few shot dataset representativeness.** For FSL methods, the representativeness of the support  
369 set is not granted. Performances expectations are conditioned by the likelihood of not having an  
370 unrepresentative support set and the quality of these tiny data sets can be easily altered by the sim-  
371 plest noisy or outlier example.

372 **Hybrid learning by confusion transfer.** We introduced some methods to improve classification  
373 performance of an hybrid algorithm, however those methods are focusing on the choice of the sup-  
374 port vector for prototypical network regardless of the existing confusion of tit could be possible to  
375 use samples not correctly learnt by the DL model as support vector for the FSL method, using the  
376 learning confusion matrix.

377

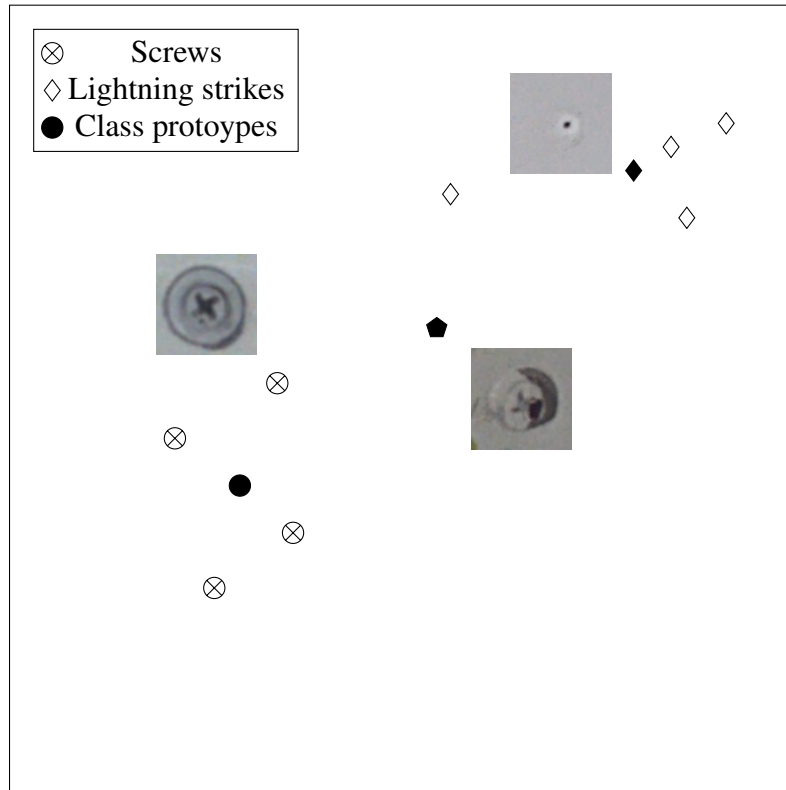
378 Overall, we discussed several machine approaches and proposed a data-aware hybridization  
379 method that applies to rare defects detection on aircraft fuselage and could be extended to many  
380 other fields with extreme imbalanced classes as it relies on the well-known effectiveness of DL  
381 approaches using transfer learning and simplicity of Prototypical Networks. The chosen method  
382 not only allows a significant improvement in the measured empirical risk, but also provides a real

383 mean of adapting the classification model through comprehensible actions in a semantic space.  
384 Opens the way to new prospects.

385 **Online-adaptive classification model** In the same vein, the presented work paves the way for  
386 online learning to complement the offline learning of the core model. Indeed, corrections made  
387 by a human operator can be considered as excellent candidates to support a prototype class with  
388 regard to the knowledge of effective confusions.

### 389 **Using zero-shot learning for non-exclusive classification**

390 Up to now we have discussed the case of classification of mutually exclusive classes. However,  
391 this assumption is not generally valid. Indeed, two distinct classes can be present in the same area  
392 of interest, in particular defects that tend to appear on protruding objects such as lightning strikes  
393 on screws. Various methods to tackle this issue can be found in the Machine Learning literature:  
394 it is possible to learn exclusive classes with a classic neural network classifier and modify the  
395 output interpretation function taking not only the class with the maximum probability but applying  
396 a threshold on the probabilities of all classes instead. We propose to use zero-shot learning to  
397 avoid those modifications: the idea is to create artificial prototypes that are the combination of  
398 N-shot embedding of real classes. By this we can produce all the class combinations automatically  
399 and classify them. This technique can be applied for refining classification results. In particular,  
400 defects that may occur preferentially on structural elements (such as lightning strikes on screws  
401 or rivets) can be categorized in this way. Figure 14 illustrates the creation of such a prototype:  
402 black filled circle and diamond are the medoids for respectively screws and lightning strikes and  
403 the black filled pentagon is the mean of those medoid and represent the theoretical prototype for  
404 screws struck by lightning.



**Fig 14** Zero-shot learning for class combination.

405 *References*

- 406 1 J. Miranda, S. Larnier, and M. Claybrough, “Caractérisation d’objets sur des images acquises  
 407 par drone,” in *Conférence Reconnaissance des Formes, Image, Apprentissage et Perception*  
 408 (*RFIAP*), 2018.
- 409 2 I. Jovančević, S. Larnier, J.-J. Orteu, and T. Sentenac, “Automated exterior inspection of  
 410 an aircraft with a pan-tilt-zoom camera mounted on a mobile robot,” *Journal of Electronic*  
 411 *Imaging* **24**(6), p. 061110, 2015.
- 412 3 J. Miranda, S. Larnier, A. Herbulot, and M. Devy, “UAV-based inspection of airplane exterior  
 413 screws with computer vision,” in *14th International Joint Conference on Computer Vision,*  
 414 *Imaging and Computer Graphics Theory and Applications (VISIGRAPP)*, 2019.

- 415 4 J. Redmon and A. Farhadi, “Yolo v3: An incremental improvement,” 2018.  
416 arXiv:1804.02767.
- 417 5 J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and F. F. Li, “ImageNet: A Large-Scale Hier-  
418 archical Image Database,” in *IEEE Conference on Computer Vision and Pattern Recognition*  
419 (*CVPR*), pp. 248–255, 2009.
- 420 6 A. Krizhevsky and G. Hinton, “Learning multiple layers of features from tiny images,” tech.  
421 rep., Citeseer, 2009.
- 422 7 J. Miranda, J. Veith, S. Larnier, A. Herbulot, and M. Devy, “Machine learning approaches for  
423 defect classification on aircraft fuselage images acquired by an UAV,” in *Fourteenth Interna-*  
424 *tional Conference on Quality Control by Artificial Vision*, **11172**, pp. 49–56, SPIE, 2019.
- 425 8 W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, “SSD: Single  
426 shot multibox detector,” in *European Conference on Computer Vision (ECCV)*, pp. 21–37,  
427 Springer, 2016.
- 428 9 A. Blumer, A. Ehrenfeucht, D. Haussler, and M. K. Warmuth, “Learnability and the  
429 Vapnik-Chervonenkis dimension,” *Journal of the Association for Computing Machinery*  
430 (*JACM*) **36**(4), pp. 929–965, 1989.
- 431 10 G. K. Dziugaite and D. M. Roy, “Computing nonvacuous generalization bounds for  
432 deep (stochastic) neural networks with many more parameters than training data,” 2017.  
433 arXiv:1703.11008.
- 434 11 C. Sun, A. Shrivastava, S. Singh, and A. Gupta, “Revisiting unreasonable effectiveness of  
435 data in deep learning era,” in *IEEE International Conference on Computer Vision (ICCV)*,  
436 pp. 843–852, 2017.

- 437 12 J. Krause, B. Sapp, A. Howard, H. Zhou, A. Toshev, T. Duerig, J. Philbin, and L. Fei-Fei,  
438 “The unreasonable effectiveness of noisy data for fine-grained recognition,” in *European*  
439 *Conference on Computer Vision*, pp. 301–320, Springer, 2016.
- 440 13 H.-C. Shin, H. R. Roth, M. Gao, L. Lu, Z. Xu, I. Nogues, J. Yao, D. Mollura, and R. M. Sum-  
441 mers, “Deep convolutional neural networks for computer-aided detection: CNN architectures,  
442 dataset characteristics and transfer learning,” *IEEE Transactions on Medical Imaging* **35**(5),  
443 pp. 1285–1298, 2016.
- 444 14 B. Krawczyk, “Learning from imbalanced data: open challenges and future directions,”  
445 *Progress in Artificial Intelligence* **5**(4), pp. 221–232, 2016.
- 446 15 S. Zagoruyko and N. Komodakis, “Wide residual networks,” 2016. arXiv:1605.07146.
- 447 16 B. Graham, “Fractional max-pooling,” 2014. arXiv:1412.6071.
- 448 17 Y. Chen, J. Li, H. Xiao, X. Jin, S. Yan, and J. Feng, “Dual path networks,” in *Advances in*  
449 *Neural Information Processing Systems*, pp. 4467–4475, 2017.
- 450 18 D. Mishkin and J. Matas, “All you need is a good init,” 2015. arXiv:1511.06422.
- 451 19 C.-Y. Lee, P. W. Gallagher, and Z. Tu, “Generalizing pooling functions in convolutional neu-  
452 ral networks: Mixed, gated, and tree,” in *Artificial intelligence and statistics*, pp. 464–472,  
453 2016.
- 454 20 J. Snoek, O. Rippel, K. Swersky, R. Kiros, N. Satish, N. Sundaram, M. Patwary, M. Prabhat,  
455 and R. Adams, “Scalable bayesian optimization using deep neural networks,” in *International*  
456 *Conference on Machine Learning (ICML)*, pp. 2171–2180, 2015.
- 457 21 C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, “Inception-v4, inception-resnet and



- 458 the impact of residual connections on learning,” in *31st AAAI Conference on Artificial Intel-*  
459 *ligence*, 2017.
- 460 22 C. Huang, Y. Li, C. Change Loy, and X. Tang, “Learning deep representation for imbal-  
461 anced classification,” in *The IEEE Conference on Computer Vision and Pattern Recognition*  
462 *(CVPR)*, June 2016.
- 463 23 F. Konidaris, T. Tagaris, M. Sdraka, and A. Stafylopatis, “Generative adversarial networks as  
464 an advanced data augmentation technique for mri data,” in *14h International Joint Conference*  
465 *on Computer Vision, Imaging and Computer Graphics Theory and Applications.*, 2019.
- 466 24 Y. Tang, “Deep learning using linear support vector machines,” 2013. arXiv:1306.0239.
- 467 25 B. M. Lake, R. Salakhutdinov, and J. B. Tenenbaum, “Human-level concept learning through  
468 probabilistic program induction,” *Science* **350**(6266), pp. 1332–1338, 2015.
- 469 26 S. Ravi and H. Larochelle, “Optimization as a model for few-shot learning,” in *5th Interna-*  
470 *tional Conference on Learning Representations (ICLR)*, 2017.
- 471 27 N. Mishra, M. Rohaninejad, X. Chen, and P. Abbeel, “A simple neural attentive meta-learner,”  
472 2017. arXiv:1707.03141.
- 473 28 C. Finn, P. Abbeel, and S. Levine, “Model-agnostic meta-learning for fast adaptation of deep  
474 networks,” in *34th International Conference on Machine Learning-Volume 70*, pp. 1126–  
475 1135, 2017.
- 476 29 J. Snell, K. Swersky, and R. Zemel, “Prototypical networks for few-shot learning,” in *Ad-*  
477 *vances in Neural Information Processing Systems*, pp. 4077–4087, 2017.
- 478 30 J. Bromley, I. Guyon, Y. LeCun, E. Säcinger, and R. Shah, “Signature verification using a”

- 479 siamese” time delay neural network,” in *Advances in neural information processing systems*,  
480 pp. 737–744, 1994.
- 481 31 F. Schroff, D. Kalenichenko, and J. Philbin, “Facenet: A unified embedding for face recog-  
482 nition and clustering,” in *IEEE conference on Computer Vision and Pattern Recognition*  
483 *(CVPR)*, pp. 815–823, 2015.
- 484 32 O. Vinyals, C. Blundell, T. Lillicrap, D. Wierstra, *et al.*, “Matching networks for one shot  
485 learning,” in *Advances in Neural Information Processing Systems (NIPS)*, pp. 3630–3638,  
486 2016.
- 487 33 X. Glorot and Y. Bengio, “Understanding the difficulty of training deep feedforward neu-  
488 ral networks,” in *thirteenth International Conference on Artificial Intelligence and Statistics*  
489 *(AISTATS)*, pp. 249–256, 2010.
- 490 34 G. E. Batista, R. C. Prati, and M. C. Monard, “A study of the behavior of several methods  
491 for balancing machine learning training data,” *ACM SIGKDD explorations newsletter* **6**(1),  
492 pp. 20–29, 2004.
- 493 35 I. Tomek, “Two modifications of cnn,” *IEEE Transactions on Systems, Man, and Cybernet-*  
494 *ics SMC-6*, pp. 769–772, Nov 1976.
- 495 36 P. Hart, “The condensed nearest neighbor rule,” *IEEE transactions on information the-*  
496 *ory* **14**(3), pp. 515–516, 1968.
- 497 37 M. Kubat, S. Matwin, *et al.*, “Addressing the curse of imbalanced training sets: one-sided  
498 selection,” in *International Conference on Machine Learning*, **97**, pp. 179–186, Nashville,  
499 USA, 1997.

- 500 38 S. Chakraborty, N. Nagwani, and L. Dey, “Performance comparison of incremental k-means  
501 and incremental dbSCAN algorithms,” 2014. arXiv:1406.4751.
- 502 39 B. Krawczyk, M. Galar, Ł. Jeleń, and F. Herrera, “Evolutionary undersampling boosting for  
503 imbalanced classification of breast cancer malignancy,” *Applied Soft Computing* **38**, pp. 714–  
504 726, 2016.
- 505 40 M. J. Siers and M. Z. Islam, “Software defect prediction using a cost sensitive decision forest  
506 and voting, and a potential solution to the class imbalance problem,” *Information Systems* **51**,  
507 pp. 62–71, 2015.

508 **Julien Miranda** is a Ph.D student at LAAS -CNRS laboratory. He is working in partnership with  
509 the start-up Donecle on the subject of defect recognition on aircraft fuselages based on images ac-  
510 quired by a drone. He is interested in machine learning theory and in particular the expressiveness  
511 of models, decisional robustness and generalization theory.

512 **Stanislas Larnier** obtained a PhD degree in Applied Mathematics from Université Paul Sabatier,  
513 Toulouse, France, in 2011. Then he worked on Image and Video Processing for INRIA, LAAS-  
514 CNRS then AKKA Research. Now, he is the Computer Vision R&D Manager in Donecle. His field  
515 of expertise is scene analysis. He worked on collaborative projects with applications in various  
516 fields: coastal engineering, microbiology, visual grading, and robotics.

517 **Ariane Herbulot** is an Associate Professor from University Paul Sabatier (Toulouse, France). She  
518 obtained her PhD in 2007 at University of Nice-Sophia Antipolis (France) on image segmenta-  
519 tion by active contours. She is interested on detection and object tracking for video-surveillance  
520 and robotics applications. Her research concerns segmentation, detection, object recognition and

521 motion estimation and tracking.

522 **Michel Devy** Michel DEVY is CNRS Research Director at /Laboratory of Analysis and Architec-  
523 ture of Systems/ at Toulouse, France. For more than 35 years, he has participated to the Robotics  
524 department, now member (as team leader until 2013) of the RAP team (/Robotics, Action and Per-  
525 ception)/. His research has been devoted to the application of computer vision in Automation and  
526 Robotics. He was PhD advisor or co-advisor for about 60 PhD students and co author of about 200  
527 scientific communications.

## 528 **List of Figures**

- 529 1 Automated drone inspection from left to right: drone with a tablet running the  
530 analysis software application, 3D model used for autonomous localization, drone  
531 inspecting aircraft.
- 532 2 Images acquired by drone.
- 533 3 Detected objects.
- 534 4 Unbalanced dataset description.
- 535 5 ResNet50 architecture (left) and residual block example (right).
- 536 6 Omniglot samples from 6 alphabets.
- 537 7 MiniImageNet samples from 6 classes.
- 538 8 Different approaches performances regarding training dataset size.
- 539 9 Undersampling iteration using Tomek's link: Two classes (left), sample of biggest  
540 class belonging to Tomek's link (center) and resulting datasets (right).
- 541 10 9 random screw supports (left) and 9 screw medoids selected after clustering (right).

- 542 11 Lightning strike from two points of view.
- 543 12 Sequential embedding for classification using prototypical neural network.
- 544 13 Precision and Recall for different methods.
- 545 14 Zero-shot learning for class combination.

## 546 **List of Tables**

- 547 1 Unbalanced dataset composition with number of samples by class.
- 548 2 Recent methods on FSL task.
- 549 3 Classification results on imbalanced dataset.