# Automatic classification of Furnariidae species from the Paranaense Littoral region using speech-related features and machine learning<sup>\*</sup>

Enrique M. Albornoz<sup>,\*,a</sup>, Leandro D. Vignolo<sup>a</sup>, Juan A. Sarquis<sup>b</sup>, Evelina Leon<sup>b</sup>

<sup>a</sup>Research Institute for Signals, Systems and Computational Intelligence, sinc(i), UNL-CONICET. <sup>b</sup>National Institute of Limnology, INALI, UNL-CONICET

### Abstract

Over the last years, researchers have addressed the automatic classification of calling bird species. This is important for achieving more exhaustive environmental monitoring and for managing natural resources. Vocalisations help to identify new species, their natural history and macro-systematic relations, while computer systems allow the bird recognition process to be sped up and improved. In this study, an approach that uses state-of-the-art features designed for speech and speaker state recognition is presented. A method for voice activity detection was employed previous to feature extraction. Our analysis includes several classification techniques (multilayer perceptrons, support vector machines and random forest) and compares their performance using different configurations to define the best classification method. The experimental results were validated in a cross-validation scheme, using 25 species of the family Furnariidae that inhabit the Paranaense Littoral region of Argentina (South America). The results show that a high classification rate, close to 90%, is obtained for this family in this Furnariidae group using the proposed features and classifiers.

*Key words:* Bird sound classification, computational bioacoustics, machine learning, speech-related features, Furnariidae.

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27

28

29

31

32

33

34

#### 1. Introduction

Vocalisations are often the most noticeable man-2 ifestations of the presence of avian species in dif-3 ferent habitats [1]. Birds have been widely used 4 to indicate biodiversity since they provide critical 5 ecosystem services, respond quickly to changes, are 6 relatively easy to detect and may reflect changes 7 at lower trophic levels (e.g. insects, plants) [2, 3]. 8 Technological tools (such as photographic cameras, video cameras, microphones, mass storage disks, 10 etc.) are useful for collecting data about several 11 patterns of bird populations. However, there are a 12

number of problems associated with them, such as poor sample representation in remote regions, observer bias [4], defective monitoring [5], and high costs of sampling on large spatial and temporal scales, among others.

Bird vocalisations have become an important research field, influencing ethology [6, 7], taxonomy [8, 9, 10] and evolutionary biology [11, 12]. One of the main activities that benefits from vocalisation identification is ecosystems monitoring, where the technological advances allow registering and processing the recordings, and improving the data collection in the field [13]. This makes possible to gather data in large and disjoint areas, which is essential for conducting reliable studies.

Although some works describe vocalisation changes in certain Furnariidae species [14, 15, 16, 17, 18], none of them simultaneously evaluates several vocalisations of Furnariidae species from South America. In this study, vocalisations belonging to 25 Furnariidae species that are distributed in the Paranaense Littoral region (see Figure 1) are January 25, 2017

 $<sup>^{\</sup>diamond}$ Facultad de Ingeniería y Ciencias Hídricas, Ciudad Universitaria, Paraje El Pozo, (S3000) Santa Fe (Argentina). Tel.: +54 (342) 457-5233/39 ext 191; http://fich.unl.edu.ar/sinc/

<sup>\*</sup>Corresponding author.

Email addresses: emalbornoz@sinc.unl.edu.ar

<sup>(</sup>Enrique M. Albornoz), ldvignolo@sinc.unl.edu.ar (Leandro D. Vignolo), juansarquis@conicet.gov.ar (Juan A. Sarquis), evelinaleon@conicet.gov.ar (Evelina Leon) Preprint submitted to Ecological Informatics



Figure 1: Paranaense Littoral region (Argentina).

analysed. This region comprises the Argentinian 35 Mesopotamia (Misiones, Corrientes and Entre Ríos 36 provinces) along with the provinces of Chaco, For-37 mosa and Santa Fe, and it is lapped by great rivers 38 of the Plata basin [19]. Over the last years, this 39 region has become an interesting place for study-40 ing bird vocalisations [16, 17, 20, 21]. In addition, 41 the work of researchers from the National Institute 42 of Limnology (INALI) along with the availability 43 of Furnariidae species would allow us to record and 44 analyse these species in real-life conditions in future 45 studies. Recently, some authors have researched the 46 47 vocalisations and the natural history of Furnariidae. Zimmer et al. [15] used morphometric analysis, 48 behavioural data and vocalisations to analyse the 101 49 Pseudoseisura cristata. The role of several habitats 102 50 as well as natural history, taxonomy, morphology, 103 51 vocalisations and evolution for the Upucerthia sat- 104 52 uratior was studied in [16, 17]. 53

Recognition of species in passeriformes is a chal- 106 54 lenging task because to they produce complex songs 107 55 and can adapt their content over time. It is inter-56 esting to note that the song content can be changed 109 57 depending on the audience, for example, when the 110 58 59 receiver is male or female [22], or in order to match it with that of their neighbours [23]. Furthermore, 60 they can take possession of new songs or syllables 61

during their lifetime [24]. The family Furnariidae produces several songs and some species manifest these as duets. It represents a synchroniation of physiological rhythms in a natural behaviour, which adds more complexity to the analysis. In addition, some species of the same family show similar structures in their songs. These similarities are manifested in introductory syllables or in the trill format, while the complexity of duets within the family makes the analysis and classification of vocalisations more difficult. Previous studies demonstrated that there are differences in tone and note intervals between males and females [25, 17, 16, 15]. For this family, the complexity of vocalisations was proved by means of playback experiments. These showed that the different taxa express dissimilar responses to similar patterns.

It should be noted that environmental conditions (humidity, wind, temperature, etc.) may alter the recording process, modifying the features that are present in the structure of songs and in the calls (e.g. frequency, duration, amplitude, etc.) [26, 27, 28]. Since these conditions may lead to errors and distort subsequent analyses and results, researchers usually use recordings from known databases. Even though these registrations can be also affected by environmental issues, their attributes and labels are validated by the scientific community and consequently, they are more reliable than "homemade" records.

As mentioned in [29], new frontiers have been opened in ecology (besides the analysis performed by expert ecologists) due to the propagation of projects like  $Xeno-canto^1$  and  $EcoGrid^2$ . The access to multimedia data has promoted an interdisciplinary and collaborative science for analysing the environment. Although human experts (who are sufficiently trained) can recognise bioacoustic events with a high performance, this is a laborious and expensive process that would be more efficient if they had the technical support of a semiautomatic tool [30]. Finally, the goal pursued is the development of an automatic classifier that provide a high accuracy and involve the expert only for evaluating the results. Automatic tools allow simultaneous studies to be conducted and diverse bird communities to be monitored in several areas at the same time, in order to identify when and how the species vocalise. In addition, said tools

62

63

64

65

66

67

68

69

70

71

72

73

74

75

76

77

78

79

80

81

82

83

84

85

86

87

88

89

90

91

92

93

94

95

96

97

98

99

100

105

<sup>&</sup>lt;sup>1</sup>http://www.xeno-canto.org/

<sup>&</sup>lt;sup>2</sup>http://www.aiai.ed.ac.uk/project/ecogrid/

could be used to create complete inventories of bird
communities in unknown or restricted areas, which
are essential for conservation or management plans.
In particular the bird call identification task can

be used in two ways [31]: call retrieval (detec-167 115 tion) and call classification. In the call retrieval 116 168 task, the objective is to identify one or more calls 169 117 in an audio recording, which can contain multiple 170 118 calls of different species overlapped or at different 171 119 times. In the classification task, a set of call classes 172 120 must be defined and the classifier will be trained 173 121 to recognise this fixed set. In this way, every in-174 122 put audio (expected to contain only one call) will 175 123 be classified to one of those classes. A classification 176 124 scheme can be defined as a pipeline of three mod-177 125 ules: preprocessing, feature extraction and classi-126 178 fication. The first one depends strongly on the 179 127 recording process and involves filtering, segmenta- 180 128 tion and enhancement of audio signals. Further-181 129 more, automatic methods for voice activity detec-182 130 tion (VAD) have been recently incorporated [32]. 183 131 Regarding feature extraction, time- and frequency-132 184 based information was employed [30, 1, 33, 34]. In 185 133 addition, characteristics that were originally devel-186 134 oped for speech analysis are used in the context of 187 135 bird call recognition. Some of the features present 188 136 in the literature are mel frequency cepstral coeffi-137 cients (MFCCs) [35], linear frequency cepstral co-138 efficients (LFCCs) [36], and standard functionals 191 139 (mean, standard deviation, kurtosis, etc.) com-140 192 puted over these [32, 37, 38]. Various techniques 193 141 have been applied to bird call classification: Gaus-142 10/ sian mixture model (GMM) [39], gaussian mixture 143 model-universal background model (GMM-UBM) 144 196 [40], support vector machines (SVM) [41], random 145 197 forest (RF) [42], among others. In [32], LFCC fea-146 tures were used along with GMM-UBM to identify 147 148 some subjects from the same bird species.

A similar approach was proposed in [43] for recog-149 nising a single bird species using MFCCs. An in- 199 150 teresting strategy based on the pairwise similar- 200 151 ity measurements, computed on bird-call spectro-201 152 grams, was evaluated in [33], where the authors 153 used different classifiers to recognise four species. In 154 [37], thirty-five species were classified using a SVM 155 classifier and six functionals were obtained from 203 156 each MFCC. A different approach was proposed 204 157 in [44], where a classifier based on hidden Markov 205 158 159 models (HMMs) was used to recognise bird calls 206 through their temporal dynamics. Previous works 207 160 developing full-automatic methods for vocalisation 208 161 recognition can be examined in [45, 46, 47, 48], and 209 162

the current relevance of this topic is shown in some recent works [32, 43]. However, none of these works has addressed the vocalisation recognition of species belonging to the Furnariidae family, which present similar parameters in their vocalisations. Moreover, only a small part of the state-of-the-art speech features have been employed in bird classification tasks. In [49], a large set of state-of-the-art speech features is described, comprising more than 6000 features, and many of these are considered within this task for the first time in this work.

This study proposes the development of a bird call recognition model for dealing with the family Furnariidae from the Paranaense Littoral region of Argentina, which is the first approach for these species. Our model is designed to use stateof-the-art classifiers with speech-related parameterisations, and some feature selection techniques are used to reduce dimensionality while maximising accuracy. As a first step, a method for performing the VAD is included. The model is tested in a crossvalidation scheme in all cases. Furthermore, the best results are discussed, and the confusion matrix is analysed to introduce the misclassification and how some similarities among some species could be addressed in order to improve the performance.

The following section introduces the proposed features and classifiers. Section 3 deals with the experimental setup, presents the implementation details and describes the validation scheme. The results are presented and discussed in Section 4. In addition, the implementation of a web-demo and an android application for testing the model is explained. Finally, conclusions are summarised and future work is commented in the last section.

# 2. Proposed features and classifiers

This section introduces the feature extraction process, two different feature selection techniques and the classifier models.

#### 2.1. Feature extraction

3

As mentioned above, the use of speech-based features is known in bird call analysis, identification and classification. For these tasks, the LFCCs and MFCCs sets (standards in speech recognition) showed good performances [37, 32]. An extended state-of-the-art set of features related to human speech is introduced below.

#### 2.1.1. Speech inspired features 210

In the speech processing area, researchers have 211 made a great effort to find the best set of features 212 213 for speech recognition, speaker recognition, emotion recognition, illness state detection, etc. [50, 51, 52]. 214 In the INTERSPEECH 2013 ComParE Challenge 215 [50], a set of 6373 features was presented which is 216 considered the state-of-the-art in speech processing. 217 The feature set is built from 65 low-level descriptors 218 (LLDs) such as energy, spectral, cepstral (MFCC), 219 voicing-related characteristics ( $F_0$ , shimmer, jitter, 220 etc.), zero crossing rate, logarithmic harmonic-to-221 noise ratio (HNR), spectral harmonicity, psychoa-222 coustic spectral sharpness, and their deltas (i.e. 223 their first temporal derivatives). These features are 224 computed on a time frame basis, using a 60-ms win-225 dow with 10-ms step for  $F_0$  (pitch) and zero crossing 226 rate. The remaining features are computed using 227 a window size of 20 ms and the time contour of 228 each attribute is smoothed by a moving average fil-229 ter. Specific functionals are then computed for each 230 LLD set. These include the arithmetic mean, maxi-231 mum, minimum, standard deviation, skewness, kur-232 tosis, mean of peak distances, among others. Ta-233 bles 1 and 2 provide an exhaustive enumeration of <sup>266</sup> 234 all the LLDs and functionals used to constitute the  $^{267}$ 235 complete feature vector. In addition to the com-236 plete feature set obtained by combining all LLDs 237 and functionals (Full-Set), this work also proposes 238 a subset consisting of the complete set of function-239 als computed only from the MFCCs, which results <sup>272</sup> 240 in a set of 531 attributes (MFCC+Fun). 241

To the best of our knowledge, no suitable baseline 242 models are available for comparing the performance 243 of our proposal. In order to create the baseline, pre-244 vious works [37, 53] were considered to define the 245 classifiers and feature sets for the bird song identi-246 fication task. The first 17 MFCCs, their deltas and 247 acceleration coefficients were computed using over-248 lapped frames. Then, the mean and variance for 249 each feature (over the entire song) were calculated, 250 which resulted in a 102-dimensional vector for each 251 recording. 252

#### 2.1.2. Feature selection 253

Feature selection techniques were defined in or-254 der to reduce the dimensionality of data while keep-255 ing the most relevant information. This allows less 256 complex models to be generated, which reduces the 257 number of parameters to estimate in the model and 258 the computing cost, and provides a similar or even 259

Table 1: Low-level descriptors (LLDs) [51].  $+\Delta$  means that the first derivative is computed and appended, to the feature vector computed for each analysis frame.

Low-level descriptors
Sum of auditory spectrum (loudness) + $\Delta$
Sum of RASTA-style filtered auditory spectrum + $\Delta$
RMS Energy $+\Delta$
harmonic-to-noise ratio (HNR) + $\Delta$
Zero-Crossing Rate + $\Delta$
RASTA-style filtering. Bands 1-26 (0-8 kHz) + $\Delta$
MFCC 1-14 + $\Delta$
Spectral energy 25-650 Hz, 1 k-4 kHz + $\Delta$
Spectral Roll Off Point 0.25, 0.50, 0.75, 0.90 + $\Delta$
Spectral Flux, Entropy, Variance + $\Delta$
Skewness, Kurtosis, Slope + $\Delta$
F0, Probability of voicing $+\Delta$
Jitter (local, delta) + $\Delta$
Shimmer (local) + $\Delta$

improved performance. Feature or attribute selection is commonly carried out by searching the space of feature subsets, and each candidate subset is evaluated according to some criteria [54].

In this study, the performance of two wellknown attribute selection methods is compared: best first (BF) [55] and linear forward selection (LFS) [56]. The BF method performs a greedy hill climbing using backtracking, which means that it can search forward through a specified number of non-improving nodes before the algorithm goes back. This algorithm has proven to guarantee the best global subset without exhaustive enumeration, given that the criterion used satisfies monotonic-The LFS algorithm is an extension of BF, ity which aims to reduce the number of evaluations performed during the search process. The number of attribute expansions is limited in each forward selection step, which drastically improves the runtime performance of the algorithm [56]. Both feature selection methods need a criterion to evaluate each considered subset; therefore the correlation-based feature subset evaluation (CFS) method [54] was applied. This method assesses the predictive ability of each attribute in the subset, and also considers the redundancy among them. Finally, the method picks up the subsets whose attributes are highly correlated within the class and have low intercorrelation among classes. Both feature selection methods were implemented using WEKA library<sup>3</sup>[57].

260

261

262

263

264

265

268

269

270

271

273

274

275

276

277

278

279

280

281

282

283

284

285

286

287

288

<sup>&</sup>lt;sup>3</sup>Software available at http://www.cs.waikato.ac.nz/ ml/weka/

Table 2: Functionals applied to LLDs [51].

Base functionals
quartiles 1-3
3 inter-quartile ranges
1 % percentile ( $\approx$ min), 99 % percentile ( $\approx$ max)
percentile range 1 %-99 $\%$
arithmetic mean, standard deviation
skewness, kurtosis
mean of peak distances
standard deviation of peak distances
mean value of peaks
mean value of peaks - arithmetic mean
linear regression slope and quadratic error
quadratic regression a and b and quadratic error
simple moving average
contour centroid
duration signal is below $25~\%$ range
duration signal is above 90 $\%$ range
duration signal is rising/falling
gain of linear prediction (LP)
Linear Prediction Coefficients 1-5
F0 functionals
percentage of non-zero frames
mean, max, min, std. dev. of segment length
input duration in seconds

#### 2.2. Classifiers 290

Several techniques from machine learning and 303 291 computational intelligence have been used in bird 292 304 call identification [32]. Based on previous studies, 293 the analysis in this work was focused on some of  $_{306}$ 294 the most commonly used classification algorithms. 307 295 The following subsections briefly introduce three 308 296 techniques: multilayer perceptron, random forest 309 297 and support vector machines. WEKA and Scikit- 310 298 Neuralnetwork  $^4$  libraries were employed to apply  $_{\scriptscriptstyle 311}$ 299 these classifiers. 300

#### 2.2.1. Multilayer perceptron 301

A multilayer perceptron (MLP) is a class of artificial neural network that consists of a set of process units (simple perceptrons or neurons) arranged in layers. In the MLP, the nodes are fully connected between layers without connections between units in the same layer (Figure 2). The input of the MLP is the feature vector  $(\mathbf{x})$ , which feeds each of the neurons of the first layer, the outputs of this layer feed into each of the second layer neurons, and so on [58]. The output of a neuron is the weighted sum of its inputs plus the bias term, and its activation



Figure 2: Example of a MLP network model.

is a function (linear or nonlinear) as

$$y = \mathcal{F}\left(\sum_{i=1}^{n} \omega_i x_i + \theta\right). \tag{1}$$

The output of the MLP (i.e. the output of the neurons in the last layer) is decoded to provide the predicted label for a given input example. The backpropagation method [58] is commonly used to obtain the synaptic weights for the connections in the network  $(\omega_i)$ . This method computes the gradient of a loss function, with respect to all network weights. The weights are then updated according to the gradient, with the aim of minimising the loss function (usually the mean square error). Since the method requires a desired output for each training input in order to calculate the error, it is considered as a supervised learning technique.

In this work, three architectures were considered: one hidden layer with the number of neurons set as (Num. of inputs+Num. of outputs)/2 (MLP1),one hidden layer with the number of neurons set to the number of inputs (MLP2), and two hidden layers set as in MLP2 and MLP1, respectively (MLP3).

#### 2.2.2. Random forest

Classification and regression tree (CART) models, the so-called decision trees, are widely known in machine learning and data mining [59]. Some relevant properties include their robustness to different feature transformations, such as scaling, and their ability to discriminate irrelevant information while

302

305

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

<sup>&</sup>lt;sup>4</sup>Software available at http://scikit-neuralnetwork. readthedocs.org

Ecological Informatics, Vol. 38, No., pp. 39 - 49, 2017.



Figure 3: Example of CART using feature vector  $\in \mathbb{R}^3$ .

producing easily analysable models. These models are constructed by recursive partitioning the input space and region-specific models are then defined for the resulting scheme [42]. This can be represented with a tree, where the nodes indicate the decision functions and each leaf stands for a region (Figure 3).

Random forest (RF) is an ensemble learning method whose decision is based on the average of multiple CARTs, which are trained on different parts of the same training set, with the aim of reducing the variance of CART overfitting. The computation can be expressed in terms of the *bagging* technique [59] as

$$f(\mathbf{x}) = \frac{1}{K} \sum_{k=1}^{K} t_k(\mathbf{x})$$
(2)

where  $t_k$  is the k-th tree. Here, the RF was implemented following [42], considering 10 and 100 trees with unlimited depth.

#### 339 2.2.3. Support vector machine

A support vector machine (SVM) is a supervised learning method that is widely used for pattern classification and is supposed to have good generalisation capabilities [60]. Its aim is to find a hyperplane that can separate input patterns in a sufficiently high dimensional space. The distances from the hyperplane to the patterns that are closest to it, on each side, is called a *margin*. This margin needs to be maximised to reach the best generalisation. In the binary case, this is done finding the **w** and  $w_0$  parameters by means of a standard quadratic optimisation [61, 60]:

 $\min \frac{1}{2} \parallel \mathbf{w} \parallel^2 \tag{3}$ 

subject to

$$r^t(\mathbf{w}^T\mathbf{x}^t + w_0) \ge +1, \forall t$$

where  $\{\mathbf{x}^t, r^t\}$  is a pattern with  $r^t = -1$  if  $\mathbf{x}^t$  is class #1, or  $r^t = +1$  in the other case.

It is known that a nonlinear problem could be solved as a linear problem in a new space by making a nonlinear transformation [61]. The new dimensions are then computed using the basis functions by inner product. The *kernel trick* is a method that solves this problem without mapping the features in the new space; therefore, the kernel function is applied to the original space [61]. Some of the more popular kernels used in SVMs are the polynomial of degree q:

$$K(\mathbf{x}^t, \mathbf{x}) = (\mathbf{x}^T \mathbf{x}^t + 1)^q \tag{4}$$

and radial-basis functions:

$$K(\mathbf{x}^{t}, \mathbf{x}) = \exp\left[-\frac{\mathcal{D}(\mathbf{x}^{t}, \mathbf{x})}{2s^{2}}\right]$$
(5)

where  $x^t$  is the centre, s is the radius and  $\mathcal{D}(\mathbf{x}^t, \mathbf{x})$  is a distance function. In our experiments, the SVMs were trained using the sequential minimal optimisation algorithm and considering the polynomial kernel.

### 3. Experiments

342

348

349

350

351

352

353

354

355

356

357

359

360

361

362

363

364

365

366

367

368

369

370

371 6 This section describes the experimental framework used in this study. First, a discussion on why and how the bird species were selected from the known databases. Then, the implementation details of the feature extraction and classifiers are presented. Finally, the validation scheme used to evaluate the models is explained. A general scheme of the whole process for the experiments is shown in Figure 4.

#### 3.1. Study area and target species

The study area islocated between  $22^{\circ}25'S$   $62^{\circ}12'W$  and  $38^{\circ}0'S$   $57^{\circ}26'W$ (Figure 1), and comprises several ecoregions along the Paraná River. These regions are Dry Chaco. Espinal, Pampa, Iberá Wetlands, and Delta and Islands of the Paraná River [2]. The family Furnariidae presents diverse vocalisations and some species can even sing male-female duets. In spite of that, the experts are usually able to identify them, reaching a good performance. The vocalisations obtained from species of this family might be similar and thus difficult to classify. In addition, the vocalisations from one species can change depending on its geographical location.



Figure 4: Conceptual flowchart of the general whole process for the experiments.

408

409

410

411

412 7

The family Furnariidae includes 68 genera com-372 posed of 302 species [62]. Being distributed in 373 South America and in a region of Central Amer-374 ica [63], it is one of the most impressive examples 375 of continental adaptive radiation. This family has 376 probably the highest morpho-ecological diversity in 377 birds, living in diverse habitats such as desert or 378 arid regions, rocky coasts, ravines, swamps, grass-379 lands and forests [64, 65]. The characteristics de-380 scribed above plus the large number of studies 381 about its taxonomy, the biological and natural his-382 tory [64, 66, 67, 68, 65, 17] and our own experi-383 ence at INALI make the family Furnariidae an in-384 teresting and open challenge to study. Figure 5 385 shows the tree structure of the 25 studied Furnari-386 idae species/genera. 387

### 388 3.2. Bird call corpus

To obtain a suitable number of vocalisations for 389 training the classifiers and evaluating the perfor-390 mance, records from two well-known databases were 391 selected, obtaining a total of 206 recordings. From 392 these, 90 recordings were selected from the Xeno-393  $canto^5$  database [69, 1, 70] and 116 recordings were 394 taken from the Birds of Argentina & Uruguay: A 395 Field Guide Total Edition corpus [71, 21, 72]. This 396 combination of different data sources involves an 397 additional complexity that the model should be 398 able to handle<sup>6</sup>. 399

### 400 3.3. Feature extraction

As mentioned earlier, the step prior to feature
extraction is usually the preprocessing and it is carried out to standardise the audio signals. A Wienerbased noise filter [73] was applied to the audio signals to reduce noise in the recordings. As all of the



Figure 5: Tree structure of the 25 studied Furnariidae species.

utterances have an initial silence, the noise could be modelled.

The acoustic activity detection (where the information is contained) is an active area of research [74]. In this work, the endpoints of acoustic activity were computed using a voice activity detector (VAD) based on Rabiner and Schafer's method [75].

<sup>&</sup>lt;sup>5</sup>http://www.xeno-canto.org/

<sup>&</sup>lt;sup>6</sup>The list of audio files used in this work was included as supplementary material.

sinc(i) Research Center for Signals, Systems and Computational Intelligence (fich.unl.edu.ar/sinc)

The openSMILE toolkit [76] was used to extract 441 413 the state-of-the-art features [50] mentioned in the 414 previous section. This is a feature extraction tool 442 415 that allows a large set of audio features to be ex- 443 416 tracted, and it is distributed free of charge for re-444 417 search and personal use<sup>7</sup>. 418 445

#### 3.4. Validation 419

Coefficients in vectors were normalised using the 420 maximum and minimum values (for each dimen-421 sion) in the training set as follows: 422

$$C_{i,j}^{norm} = \frac{(C_{i,j} - C_{min,j})}{(C_{max,j} - C_{min,j})},$$
(6)

where  $C_{i,j}^{norm}$  is the normalised coefficient j from 455 423 recording  $i, C_{i,j}$  represents the original value, while 456 424  $C_{min,j}$  and  $C_{max,j}$  represent the minimum and 457 425 maximum values of coefficient j from all the train-426 ing recordings. 427

The recognition rate estimation may be biased 460 428 if only one training partition and one test parti-461 429 tion are used. To avoid these estimation biases, 430 a cross-validation was performed with the k-fold 463 431 method [77]. For each experiment the classification <sup>464</sup> 432 results by 10-fold stratified cross-validation (SCV) 465 433 were computed, where each fold was composed of 466 434 90% of data for training and the remaining 10% was  $_{467}$ 435 used for testing. Finally, the results were computed 468 436 and averaged over the 10 test sets. 437

Several classification measures were computed 470 for accurately visualising the performance of the 471 models. The weighted average recall or accuracy 472 (ACC) is the number of correctly classified in- 473 stances divided by the total number of instances. 474 Although this measure is widely used, it can be <sup>475</sup> biased when the classes are not balanced. If the 476 classes (species) are unbalanced, the unweighted av- 477 erage recall (UAR) gives a more accurate estima- 478 tion of the performance [78]. The UAR was com-<sup>479</sup> puted as the average of all class accuracies as:

$$\text{UAR} = \frac{1}{K} \sum_{i=1}^{K} \frac{A_{ii}}{\sum_{j=1}^{K} A_{ij}},$$
 (7)

where K is the number of classes and  $A_{ij}$  is the 438 number of instances belonging to class i that are 439 classified as i. 440

# 4. Results and discussion

446

447

448

449

450

451

452

453

454

459

The baseline feature set and the proposed feature sets were evaluated using all the classifiers described in Section 2.2, considering the normalised attributes explained in Section 3.4. Also, LFS and BF feature selection methods were used to reduce the size of the Full-Set (6373 features), maximising accuracy while keeping the most relevant information. Tables 3 and 4 present the results obtained in terms of accuracy and UAR, respectively<sup>6</sup>. Table 3 shows that the baseline set (102 features) provides high accuracy rates while the proposed sets improve these results, and the best results are close to 90%. However, the performance is lower when the Full-Set is used because the models cannot be properly trained. This means that the complexity of the classifiers is increased due to the high number of inputs (especially in the case of MLP), and the small amount of data available is not enough for appropriately training them, which causes poor performance.

In order to assess how the imbalance of classes affects the results, the UAR values should be analysed, taking into account the hit rates for each class (Table 4). This table presents similar results, where the proposed feature sets improve the baseline performance. The MFCC+Fun set (531 features) performs better than the baseline for almost all classifiers, whereas both feature selection methods applied over the Full-set achieve the best performances. It is interesting to note that MLPs and SVMs produce better results than RF for all the feature sets. Finally, one can be conclude that the best performance is obtained using the multilayer perceptron (MLP1) and applying the LFS method over the Full-Set.

The dimension of the best feature set is 153, thus the system has kept a very low dimensionality in addition to achieving the best rates. The retained features include mostly spectral and cepstral coefficients as described next. Thirty-six features were computed based on the MFCC coefficients and some functionals (quartiles, percentiles and mean, among others). Eleven features obtained from the first derivative of MFCC (delta MFCC) [79] and the same functionals. Twenty-four spectral features were selected, including roll-off (percentile of the power spectral distribution), slope (which describes how rapidly the amplitudes of successive component change), harmonicity (which evaluates the total strength of harmonic structure) and flux

480

481

482

483

484

485

486

487

488

489

<sup>&</sup>lt;sup>7</sup>Software available at http://www.audeering.com/ research/opensmile/

Table 3: Weighted average recall (accuracy) [%].

	-		- 1	. ,		
Feature vector	MLP1	MLP2	MLP3	RF10	RF100	SVM
Baseline	85.92	86.89	78.64	68.45	80.10	84.95
MFCC+Fun	89.32	88.83	79.61	69.42	83.01	85.92
Full-Set	74.27	65.05	08.25	68.93	80.10	83.50
Full-Set + LFS	89.32	86.89	80.58	76.70	86.41	87.38
$\mathrm{Full}\mathrm{-}\mathrm{Set}+\mathrm{BF}$	89.32	89.32	80.58	76.70	86.41	87.38

Table 4: Unweighted average recall (UAR) [%]

Feature vector	MLP1	MLP2	MLP3	RF10	RF100	SVM			
Baseline	77.24	79.21	72.06	58.25	67.00	74.07			
MFCC+Fun	79.96	80.85	69.16	58.08	70.43	75.18			
Full-Set	61.90	53.55	05.06	55.74	65.24	72.46			
Full-Set + LFS	82.21	78.74	68.65	64.20	73.65	77.82			
Full-Set + BF	82.10	80.25	70.42	64.20	73.35	77.82			

Table 5: Confusion matrix for the MLP1 and Full-Set+LFS. References for the classes are included in the supplementary material.



(a measure that indicates how quickly the power 504 492 spectrum of a signal is changing) [76]. Twelve fea- 505 493 tures computed as functionals from frequency band 506 494 energies, particularly in bands of 250-650Hz and 507 495 1000-4000Hz. Forty-four features obtained by ap- 508 496 plying functionals to 26 spectral bands filtered with 497 509 RASTA (RASTA uses bandpass filtering in the log 498 spectral domain to remove slow channel variations) 510 499 [80]. Eleven features computed from the auditory 511 500 512 spectrum, which is inspired by psychoacoustic stud-501 513 ies on human primary auditory cortex and produces 502 514 a time-frequency representation. Five features com-503

puted as functionals from the auditory spectrum filtered with RASTA [80]. Twelve features computed from the root mean square energy, voicing, harmonic-to-voice ratio, jitter and zero crossing rate.

As the performance obtained is highly satisfactory (close to 90%) and the amount of data is limited, a test of statistical significance like the paired T-test [81] is not relevant. However, our results suggest that 5 samples per species are required for properly training the model (see Table 5). Evidently, patterns from the same species present some

differences, therefore analyses where only one sam- 565 516 ple is used to represent the species (as in [37]) could 566 517 be not very reliable. Furthermore, confusions may 567 518 be explained by certain similarities in vocalisations, 568 519 such as waveform shapes, harmonic content, place-569 520 ment and separation of syllables, among others. 521 570 These should be deeply explored in future analy-522 571 ses and modelled in order to improve the results. 523

Since the limited amount of data might make the 573 524 result obtained through 10-fold cross validation un- 574 525 stable, the performance using leave-one-out cross 575 526 validation (LOOCV) was also evaluated. LOOCV 576 527 was performed for the alternative with the best 577 528 performance (Full-Set + LFS features with MLP1 578 529 classifier) and the baseline with best performance 579 530 (baseline features with MLP2 classifier). As a re-531 580 sult, UARs of 85.09% and 80.18% were obtained <sub>581</sub> 532 for the proposed features and the baseline, respec- 582 533 tively. The accuracy achieved was 91.75% and  $_{583}$ 534 88.35% for the proposed features and the baseline, <sub>584</sub> 535 respectively. Therefore, the results obtained with 585 536 LOOCV show an even better improvement (almost 537 586 5% for UAR) of the proposed approach over the 587 538 baseline. Moreover, the performances for both al-539 ternatives were improved comparing the results ob- 589 540 tained with LOOCV and 10-fold cross validation. 590 541 Given the small amount of data available, it is rea-542 sonable that the higher number of training exam-543 ples used in each LOOCV iteration<sup>8</sup> helps the clas-544 sifier to provide a better performance. These results 545 594 suggest that the overall performance could be fur-546 ther improved if more data was available for train-547 596 ing the classifiers. 548

The results can be further analysed by using con- 598 549 fusion matrices. Confusion matrices give a good 599 550 representation of the results per each class, which 551 allows making a detailed analysis of performance 552 and finding the main classification errors. The con-553 fusion matrix (adding all partitions) of our best 554 model (MLP1 and Full-Set + LFS) is shown in  $_{601}$ 555 Table 5. The rows correspond to the actual class  $_{602}$ 556 labels, the columns show the predicted labels of 603 557 bird species, and the main diagonal indicates the 604 558 species that are correctly recognised. In this ma- 605 559 trix there are no-major errors and the unbalance 606 560 between the number of examples per species can 607 561 be noticed. Some confusions (underlined numbers) 608 562 might be due to the small amount of available pat-563 609 terns for these species when the model is trained 564 610

(see FuR, GeC, PhS and PhSt in Table 5). The remaining confusions may be explained by the acoustic likeness between species. By contrast, species of the same genus are not confused. Nevertheless, a deeper acoustic analysis would be very useful to define these "similarities". The acoustic similarities could be exploited to define groups of species without taking into account information from the traditional taxonomy of the bird family. Therefore, a hierarchical classification scheme could be defined [82, 83], which allows the mistakes to be addressed more efficiently, classifying these groups at a first stage and then, the more confusing species within the groups.

Figure 6 shows spectrograms of vocalisation segments from species Limnoctites rectirostris (LiR), Phleocryptes melanops (PhM), Upucerthia dumetaria (UpD), and Phacellodomus sibilatrix (PhS). Examples from these species were selected because they are highly confused by the model, as as presented in Table 5. The spectral characteristics of all the four vocalisations are very similar. For example, they show successive high energy peaks, which are regular in time and centred around 5000 Hz. Similarly, all the spectrograms present some weaker energy peaks around 10kHz, which are also regular in time. Since most of the features we considered are based on the spectrum, the auditory spectrum and the spectrogram, it is reasonable that these species be misclassified. Therefore, in order to obtain high a performance for these four species. it would probably be appropriate to include some features based on temporal dynamics of the vocalisations, or to consider a dynamics models for the classification, like hidden Markov models [84].

# 5. Conclusions and future work

The identification of bird species is of increasing importance for ecologists in order to monitor the terrestrial environment, as it reflects important ecosystem processes and human activities. This study explores the bird call classification using speech-related features, and compares the performance using different classification techniques and configurations. Species from the family Furnariidae in the Paranaense Littoral region were analysed, which are well-known in the community but were never studied considering a big group. In addition, our work was motivated by the hypothesis that an extended state-of-the-art feature set, de-

611

612

572

588

595

597

<sup>&</sup>lt;sup>8</sup>It is compared to the number of training examples in each fold for 10-fold cross validation.



Figure 6: Spectrograms of vocalisation segments from species *Limnoctites rectirostris* (LiR), *Phleocryptes melanops* (PhM), *Upucerthia dumetaria* (UpD) and *Phacellodomus sibilatrix* (PhS).

fined for speech-related tasks, would obtain a better 638 614 performance than the feature set used at present. 639 615 The research demonstrated that the baseline re-640 616 sults can be improved using additional LLDs, keep-641 617 ing low-dimensional data. The results were poorer 642 618 when the Full-Set was used, which is expectable due 643 619 to the high dimensionality of data and the num- 644 620 ber of samples used to train the multi-class models. 645 621 This means that the large number of inputs makes 646 622 the model more complex, and the scarce number of 623 647 examples available is not enough for appropriately 648 624 training it. Finally, the best performances (ACC 649 625 and UAR) were obtained, keeping a low dimension- 650 626 ality, when feature selection techniques were used. 651 627 This indicates that said techniques are appropriate 652 628 for extracting the more discriminative information 653 629 from the full set of features, and exhibit a good 654 630 behaviour with unbalanced data. Particularly, the 655 631 best result is reached using a MLP classifier and the 656 632 LFS technique. From an ecological monitoring and 657 633 management point of view, our approach would be 658 634 useful for developing autonomous tools that allow 635 ornithologists to know which species are present in 660 636 particular areas. Specifically, it could reduce the 661 637

effort of manually reviewing recordings of Furnariidae species for labelling. Moreover, it would enable ornithologists to perform remote and simultaneous monitoring in different areas.

In future research, the model will be improved to detect more than one species in each audio file, performing a dynamic analysis of the vocalisations, i.e. frame by frame instead of using static (averaged) features. This could be achieved by matching every frame with short "templates" [85] that should be first obtained for the species. Said matching could be done in terms of cross correlation [86] or dynamic time warping [87]. Then, a "dictionary" should be built including several templates that capture the characteristics of each species. In addition, it would be interesting to extend this research to perform the classification considering a large number of families with all the genus and species included. A hierarchical classification scheme could also be used, in which the first step would classify bird families, the second step would classify genus and the last step would determine the species. This means that the first classifier would focus on families only. The second step would consist of a set of different clas-

sifiers, each of which would be trained to recognize 710 662 the genus of a particular family, which would be de-711 663 712 termined in the previous step. Finally, the last step 664 713 would consist of a classifier for each of the genus un-665 714 der study, which would determine the species given 666 715 the genus predicted in the previous step. The pos-716 667 717 sibility of developing a semi-automatic tool to pro-668 718 vide a list of the most probable species could be 669 719 also considered. Ornithologists could then select 720 670 721 the correct species from the list provided, based on 671 722 their expertise. 672 723

# 673 6. Web-demo for reproducible research

728 A web interface was implemented using the web-674 729 demo tool [88] in order to obtain further details and 675 730 test our proposal with some experimental setups. 676 731 This web interface is available at http://fich.unl. 677 732 edu.ar/sinc/blog/web-demo/furnariidae/. Also, an 733 678 android application with the same functionalities was 734 679 735 680 developed, which can be downloaded from the men-736 tioned web page. The system can be tested using an ex-681 737 ample register or uploading a register. The preprocess-682 738 ing can be set to use or not to use Wiener-based filter 683 739 and acoustic activity detector. Then, after the feature 684 740 extraction process, the sample is classified by the best 741 685 742 model trained using all the reported data. Moreover, 686 743 the graphical results of the audio file preprocessing, the 687 744 features file (arff format), the trained model and the 688 745 recognised species are freely available for download. 689 746

### 690 7. Acknowledgements

The authors would like to thank the National 691 Agency for Scientific and Technological Promotion 692 (ANPCyT)(with PICT #2014-1442 and PICT-2015-693 977), Universidad Nacional del Litoral (with PACT 694 2011 #58, CAI+D 2011 #58-511, CAI+D 2011 #58-695 525), Herpetology Laboratory from INALI (CON-696 ICET), and the National Scientific and Technical Re-697 search Council (CONICET), from Argentina, for their 698 699 support.

# 700 References

701

702

703

704

705

706

707

708

709

- I. Potamitis, Unsupervised dictionary extraction of bird vocalisations and new tools on assessing and visualising bird activity, Ecological Informatics 26, Part 3 (2015)
   6–17.
- [2] R. Burkart, N. Bárbaro, R. Sánchez, D. Gómez, Eco-Regiones de la Argentina, Administración de Parques Nacionales (APN). Secretaría de Recursos Naturales y Desarrollo Sostenible, Presidencia de la Nación Argentina, 1999.

- [3] M. Louette, L. Bijnens, D. Upoki Agenong'a, R. Fotso, The utility of birds as bioindicators: case studies in equatorial africa, Belgian Journal of Zoology 125 (1) (1995) 157–165.
- [4] R. Laje, G. B. Mindlin, Highly structured duets in the song of the south american hornero, Physical review letters 91 (25) (2003) 258104.
- [5] M. Betts, D. Mitchell, A. Diamond, J. Bêty, Uneven rates of landscape change as a source of bias in roadside wildlife surveys, The Journal of Wildlife Management 71 (7) (2007) 2266–2273.
- [6] N. Seddon, J. A. Tobias, Character displacement from the receiver's perspective: species and mate recognition despite convergent signals in suboscine birds, Proceedings of the Royal Society of London B: Biological Sciences (2010) 1–9.
- [7] N. Hesler, R. Mundry, T. Dabelsteen, Does song repertoire size in common blackbirds play a role in an intrasexual context?, Journal of Ornithology 152 (3) (2011) 591–601.
- [8] B. Schottler, Canary islands blue tits (parus caeruleus ssp.)-differences and variation in territorial song: preliminary results.
- H. Bergmann, B. Schottler, Tenerife robin erithacus (rubecula) superbus-a species of its own, Dutch Birding 23 (2001) 140–146.
- [10] M. A. Raposo, E. Höfling, Overestimation of vocal characters in suboscine taxonomy (aves: Passeriformes: Tyranni): causes and implications, Lundiana 4 (1) (2003) 35–42.
- [11] A. Lynch, A. J. Baker, A population memetics approach to cultural evolution in chaffinch song: differentiation among populations, Evolution (1994) 351–359.
- [12] M. Päckert, J. Martens, J. Kosuch, A. A. Nazarenko, M. Veith, Phylogenetic signal in the song of crests and kinglets (aves: Regulus), Evolution 57 (3) (2003) 616– 629.
- [13] M. Towsey, J. Wimmer, I. Williamson, P. Roe, The use of acoustic indices to determine avian species richness in audio-recordings of the environment, Ecological Informatics 21 (2014) 110–119.
- [14] D. I. MacKenzie, J. D. Nichols, G. B. Lachman, S. Droege, J. Andrew Royle, C. A. Langtimm, Estimating site occupancy rates when detection probabilities are less than one, Ecology 83 (8) (2002) 2248–2255.
- [15] K. J. Zimmer, A. Whittaker, The rufous cacholote (furnariidae: Pseudoseisura) is two species, The Condor 102 (2) (2000) 409–422.
- [16] J. I. Areta, M. Pearman, Natural history, morphology, evolution, and taxonomic status of the earthcreeper upucerthia saturatior (furnariidae) from the patagonian forests of south america, The Condor 111 (1) (2009) 135–149.
- [17] J. I. Areta, M. Pearman, Species limits and clinal variation in a widespread high andean furnariid: The buffbreasted earthcreeper (upucerthia validirostris), The Condor 115 (1) (2013) 131–142.
- [18] I. Potamitis, S. Ntalampiras, O. Jahn, K. Riede, Automatic bird sound detection in long real-field recordings: Applications and tools, Applied Acoustics 80 (2014) 1– 9.
- [19] V. Arzamendia, A. R. Giraudo, Influence of large south american rivers of the plata basin on distributional patterns of tropical snakes: a panbiogeographical analysis, Journal of Biogeography 36 (9) (2009) 1739–1749.

724

725

726 727

747 748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

776

777

778

779

780

781

782

783

784

785

786

787

788

789

790

791

792

793 794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

- [20] E. León, A. Beltzer, M. Quiroga, El jilguero dorado 840 (sicalis flaveola) modifica la estructura de sus vocaliza-841 ciones para adaptarse a hábitats urbanos [the saffron 842 finch (sicalis flaveola) modifies its vocalizations to adapt 843 to urban habitats], Revista mexicana de biodiversidad 844 85 (2) (2014) 546-552. 845
- [21] E. J. Leon, A. H. Beltzer, P. F. Olguin, C. F. Reales, 846 G. V. Urich, V. Alessio, C. G. Cacciabué, M. A. 847 Quiroga, Song structure of the golden-billed saltator 848 (saltator aurantiirostris) in the middle parana river 849 floodplain, Bioacoustics 24 (2) (2015) 145-152. 850
- B. E. Byers, Geographic variation of song form within [22]851 and among chestnut-sided warbler populations, The 852 Auk (1996) 288-299. 853
- R. B. Payne, Song traditions in indigo buntings: ori- 854 [23]gin, improvisation, dispersal, and extinction in cultural 855 evolution, Ecology and evolution of acoustic communi-856 cation in birds (1996) 198-220. 857
- P. Marler, Three models of song learning: Evidence [24]858 from behavior, Journal of Neurobiology 33 (5) (1997) 859 501 - 516860
- [25]J. F. Pacheco, L. P. Gonzaga, A new species of synal-861 laxis of the ruficapilla/infuscata complex from eastern 862 brazil (passeriformes: Furnariidae), Revista Brasileira 863 de Ornitologia-Brazilian Journal of Ornithology 3 (3) 864 (2013) 10.
- [26]C. M. Harris, Absorption of sound in air in the audio-866 frequency range, The Journal of the Acoustical Society 867 of America 35 (1) (1963) 11–17. 868
- [27]C. M. Harris, Absorption of sound in air versus hu-869 midity and temperature. The Journal of the Acoustical 870 Society of America 40 (1) (1966) 148-159. 871
- S. A. Zollinger, H. Brumm, Why birds sing loud songs [28]872 and why they sometimes don't, Animal Behaviour 105 873 (2015) 289-295.
- [29] C. Spampinato, V. Mezaris, B. Huet, J. van Ossen-875 bruggen, Editorial - special issue on multimedia in ecol-876 ogy, Ecological Informatics 23 (2014) 1 – 2, special Issue 877 on Multimedia in Ecology and Environment. 878
- A. Truskinger, M. Towsey, P. Roe, Decision support for [30]879 the efficient annotation of bioacoustic events, Ecological 880 Informatics 25 (2015) 14-21. 881
- [31]X. Dong, M. Towsey, A. Truskinger, M. Cottman-Fields, J. Zhang, P. Roe, Similarity-based birdcall re-883 trieval from environmental audio, Ecological Informat-884 ics 29, Part 1 (2015) 66-76.
- [32] L. Ptacek, L. Machlica, P. Linhart, P. Jaska, L. Muller, 886 Automatic recognition of bird individuals on an open 887 set using as-is recordings, Bioacoustics 25 (1) (2015) 888 1 - 19
- [33] S. Keen, J. C. Ross, E. T. Griffiths, M. Lanzone, 890 A. Farnsworth, A comparison of similarity-based ap-891 proaches in the classification of flight calls of four 892 species of north american wood-warblers (parulidae), 893 Ecological Informatics 21 (2014) 25-33. 894
- E. R. Cramer, Measuring consistency: spectrogram [34]895 cross-correlation versus targeted acoustic parameters, 896 Bioacoustics 22 (3) (2013) 247-257.
- S. Molau, M. Pitz, R. Schluter, H. Ney, Computing mel-[35]898 frequency cepstral coefficients on the power spectrum, 899 in: Acoustics, Speech, and Signal Processing, 2001. Pro-900 ceedings.(ICASSP'01). 2001 IEEE International Con-901 ference on, Vol. 1, IEEE, 2001, pp. 73–76.
- 838 [36]X. Zhou, D. Garcia-Romero, R. Duraiswami, C. Espy-903 Wilson, S. Shamma, Linear versus mel frequency cep-839 904

stral coefficients for speaker recognition, in: Automatic Speech Recognition and Understanding (ASRU), 2011 IEEE Workshop on, IEEE, 2011, pp. 559-564.

- [37]O. Dufour, T. Artieres, H. Glotin, P. Giraudet, Soundscape Semiotics - Localization and Categorization, In-Tech Open Book, 2014, Ch. Clusterized Mel Filter Cepstral Coefficients and Support Vector Machines for Bird Song Identification
- F. Briggs, B. Lakshminarayanan, L. Neal, X. Z. Fern, [38] R. Raich, S. J. Hadley, A. S. Hadley, M. G. Betts, Acoustic classification of multiple simultaneous bird species: A multi-instance multi-label approach, The Journal of the Acoustical Society of America 131 (6) (2012) 4640-4650.
- [39] M. A. Roch, M. S. Soldevilla, J. C. Burtenshaw, E. E. Henderson, J. A. Hildebrand, Gaussian mixture model classification of odontocetes in the southern california bight and the gulf of california, The Journal of the Acoustical Society of America 121 (3) (2007) 1737-1748
- [40] Z. Xiong, T. F. Zheng, Z. Song, F. Soong, W. Wu, A tree-based kernel selection approach to efficient gaussian mixture model-universal background model based speaker identification, Speech communication 48 (10) (2006) 1273–1282.
- M. A. Hearst, S. T. Dumais, E. Osman, J. Platt, [41]B. Scholkopf, Support vector machines, IEEE Intelligent Systems and their Applications 13 (4) (1998) 18-28.
- [42]L. Breiman, Random forests, Machine Learning 45 (1) (2001) 5-32.
- T. D. Ganchev, O. Jahn, M. I. Marques, J. M. de [43]Figueiredo, K.-L. Schuchmann, Automated acoustic detection of vanellus chilensis lampronotus, Expert Systems with Applications 42 (15-16) (2015) 6098 -6111.
- T. M. Ventura, A. G. de Oliveira, T. D. Ganchev, J. M. [44]de Figueiredo, O. Jahn, M. I. Marques, K.-L. Schuchmann, Audio parameterization with robust frame selection for improved bird identification. Expert Systems with Applications 42 (22) (2015) 8463-8471.
- [45] ICML int. Conf., Proc. 1st workshop on Machine Learning for Bioacoustics - ICML4B, http://sabiod.univtln.fr.
- D. Stowell, M. D. Plumbley, Feature design for multil-[46]abel bird song classification in noise, in: NIPS4B 2013 Bird Challenge, 2013.
- [47]D. Stowell, M. D. Plumbley, Segregating event streams and noise with a markov renewal process model, Journal of Machine Learning Research 14 (2013) 1891-1916.
- [48]D. Giannoulis, E. Benetos, D. Stowell, M. Rossignol, M. Lagrange, M. D. Plumbley, Detection and classification of acoustic scenes and events: an ieee aasp challenge, in: Proceedings of the Workshop on Applications of Signal Processing to Audio and Acoustics (WAS-PAA).
- [49]B. Schuller, S. Steidl, A. Batliner, J. Epps, F. Eyben, F. Ringeval, E. Marchi, , Y. Zhang, The INTER-SPEECH 2014 Computational Paralinguistics Challenge: Cognitive & Physical Load, Proc. Interspeech, ISCA (2014) 427-431.
- [50]B. Schuller, S. Steidl, A. Batliner, A. Vinciarelli, K. Scherer, F. Ringeval, M. Chetouani, F. Weninger, F. Eyben, E. Marchi, M. Mortillaro, H. Salamin, A. Polychroniou, F. Valente, S. Kim, The INTER-SPEECH 2013 Computational Paralinguistics Chal-

902

865

874

882

885

889

906

907

908 909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

- lenge: Social Signals, Conflict, Emotion, Autism, Proc. 970 Interspeech, ISCA (2013) 148–152. 971
- [51] B. Schuller, S. Steidl, A. Batliner, F. Schiel, J. Kra-972 jewski, The INTERSPEECH 2011 Speaker State Chal-973 lenge, Proc. Interspeech, ISCA (2011) 3201-3204. 974
- B. Schuller, S. Steidl, A. Batliner, The INTERSPEECH 975 [52]2009 emotion challenge, Proc. Interspeech, ISCA (2009) 976 312 - 315977
- [53] S. Fagerlund, Bird species recognition using support 978 vector machines, EURASIP Journal on Applied Signal 979 Processing 2007 (1) (2007) 64–64. 980
- M. A. Hall, Correlation-based feature subset selec-[54]981 tion for machine learning, Ph.D. thesis, University of 982 Waikato, Hamilton, New Zealand (1998). 983
- [55] L. Xu, P. Yan, T. Chang, Best first strategy for feature 984 selection, in: 9th International Conference on Pattern 985 Recognition, Vol. 2, 1988, pp. 706-708. 986
- [56]M. Gütlein, E. Frank, M. Hall, A. Karwath, Large-scale 987 attribute selection using wrappers, in: Computational 988 Intelligence and Data Mining, 2009. CIDM'09. IEEE 989 Symposium on, IEEE, 2009, pp. 332–339. 990
- [57] M. Hall, E. Frank, G. Holmes, B. Pfahringer, P. Reute-991 mann, I. H. Witten, The WEKA Data Mining Software: 992 An Update, SIGKDD Explorations 11 (1) (2009) 10–18. 993 S. Haykin, Neural Networks: A Comprehensive Foun-[58]994
- dation, 2nd Edition, Prentice Hall, 1998. 995 K. P. Murphy, Machine learning: a probabilistic per-[59]996
- spective, MIT press, 2012. [60]V. Vapnik, C. Cortes, Support-vector networks, Ma-
- 998 chine learning 20 (3) (1995) 273-297. 999
- E. Alpaydin, Introduction to Machine Learning, 2nd 1000 [61]Edition, The MIT Press, 2010. 1001
- J. Clements, T. Schulenberg, M. Iliff, D. Rober- 1002 [62]son, T. Fredericks, B. Sullivan, C. Wood, The 1003 ebird/clements checklist of birds of the world (2015). 1004 URL www.birds.cornell.edu/clementschecklist/ 1005
- [63] J. I. Noriega, Un nuevo género de Furnariidae (ave: 1006 Passeriformes) del pleistoceno inferior-medio de la 1007 provincia de Buenos Aires, Argentina, Ameghiniana 28 1008 (1991) 317-323. 1009
- [64] F. Vuilleumier, C. Vaurie, Taxonomy and geographi- 1010 cal distribution of the furnariidae (aves, passeriformes), 1011 Bulletin of the American Museum of Natural History 1012 166(1980)1-357.1013
- [65] M. Irestedt, J. Fjeldså, L. Dalén, P. G. Ericson, Conver- 1014 gent evolution, habitat shifts and variable diversifica- 1015 tion rates in the ovenbird-woodcreeper family (furnari- 1016 idae), BMC evolutionary biology 9 (1) (2009) 1. 1017
- [66] J. Garciá-Moreno, P. Arctander, J. Fjeldså, A case of 1018 rapid diversification in the neotropics: phylogenetic re- 1019 lationships among cranioleuca spinetails (aves, furnari- 1020 idae), Molecular phylogenetics and evolution 12 (3) 1021 (1999) 273–281. 1022
- J. Fjeldså, M. Irestedt, P. G. Ericson, Molecular data 1023 [67]reveal some major adaptational shifts in the early evo- 1024 lution of the most diverse avian family, the furnariidae, 1025 Journal of Ornithology 146 (1) (2005) 1–13. 1026
- [68] S. L. Olson, M. Irestedt, P. G. Ericson, J. Fjeldså, In- 1027 dependent evolution of two darwinian marsh-dwelling 1028 ovenbirds (furnariidae: Limnornis, limnoctites), Orni- 1029 tologia Neotropical 16 (2005) 347-359. 1030
- [69] B. Planqué, W.-P. Vellinga, Xeno-cano.org, accessed: 1031 2015-07-10. 1032 1033
- URL http://www.xeno-canto.org
- [70] A. Joly, H. Goëau, H. Glotin, C. Spampinato, P. Bon- 1034

net, W.-P. Vellinga, R. Planque, A. Rauber, R. Fisher, H. Müller, Lifeclef 2014: Multimedia life species identification challenges, in: E. Kanoulas, M. Lupu, P. Clough, M. Sanderson, M. Hall, A. Hanbury, E. Toms (Eds.), Information Access Evaluation. Multilinguality, Multimodality, and Interaction, Vol. 8685 of Lecture Notes in Computer Science, Springer International Publishing, 2014, pp. 229-249.

- [71]T. Narosky, D. Yzurieta, Aves de Argentina y Uruguay-Birds of Argentina & Uruguay: Guía de Identificación Edición Total-A Field Guide Total Edition, 16th Edition, Buenos Aires, 2010.
- [72]J. R. Contreras, F. Agnolin, Y. E. Davies, I. Godoy, A. Giacchino, E. E. Ríos., Atlas ornitogeográfico de la provincia de Formosa, Vazquez Mazzini, 2014.
- C. Plapous, C. Marro, P. Scalart, Improved Signal-to-[73]Noise Ratio Estimation for Speech Enhancement, IEEE Transactions on Audio, Speech, and Language Processing 14 (6) (2006) 2098–2108.
- [74]A. G. de Oliveira, T. M. Ventura, T. D. Ganchev, J. M. de Figueiredo, O. Jahn, M. I. Marques, K.-L. Schuchmann, Bird acoustic activity detection based on morphological filtering of the spectrogram, Applied Acoustics 98 (2015) 34-42.
- T. Giannakopoulos, A. Pikrakis, Introduction to Audio [75]Analysis: A MATLAB® Approach, 1st Edition, Academic Press, Oxford, 2014.
- [76] F. Eyben, F. Weninger, F. Gross, B. Schuller, Recent Developments in openSMILE, the Munich Open-Source Multimedia Feature Extractor, in: 21st ACM International Conference on Multimedia, Barcelona, Spain, 2013, pp. 835-838.
- D. Michie, D. Spiegelhalter, C. Taylor, Machine Learn-[77]ing, Neural and Statistical Classification, Ellis Horwood, University College, London, 1994.
- [78] A. Rosenberg, Classifying skewed data: Importance weighting to optimize average recall, in: INTER-SPEECH 2012, Portland, USA, 2012.
- F. Zheng, G. Zhang, Z. Song, Comparison of different [79]implementations of mfcc, Journal of Computer Science and Technology 16 (6) (2001) 582-589.
- [80]H. Hermansky, N. Morgan, Rasta processing of speech, IEEE transactions on speech and audio processing 2 (4) (1994) 578–589.
- [81]J. Demšar, Statistical comparisons of classifiers over multiple data sets, Journal of Machine learning research 7 (Jan) (2006) 1-30.
- E. M. Albornoz, D. H. Milone, H. L. Rufiner, Spoken [82] emotion recognition using hierarchical classifiers, Computer Speech & Language 25 (3) (2011) 556–570.
- C. Lee, E. Mower, C. Busso, S. Lee, S. Narayanan, Emo-[83] tion recognition using a hierarchical binary decision tree approach, Proc. Interspeech, ISCA (2009) 320-323.
- [84]Y. Ephraim, Hidden markov models, Encyclopedia of Operations Research and Management Science (2013) 704 - 708.
- [85]D. Wachter, M. Matton, K. Demuynck, M. P. Wambacq, R. Cools, D. V. Compernolle, Templatebased continuous speech recognition, IEEE Transactions on Audio, Speech, and Language Processing 15 (4) (2007) 1377-1390. doi:10.1109/TASL.2007.894524.
- [86]D.-M. Tsai, C.-T. Lin, Fast normalized cross correlation for defect detection, Pattern Recognition Letters 24 (15) (2003) 26252631.doi:http://dx.doi.org/10.1016/S0167-8655(03)00106-5.

- 1035URLhttp://www.sciencedirect.com/science/1036article/pii/S0167865503001065
- [87] M. Müller, Dynamic time warping, Information retrieval for music and motion (2007) 69–84.
  [88] G. Stegmayer, M. Pividori, D. H. Milone, A very
- [88] G. Stegmayer, M. Pividori, D. H. Milone, A very
  simple and fast way to access and validate algorithms
  in reproducible research., Briefings in Bioinformatics.
  17 (1) (2016) 180–183.
- 1043URLhttp://fich.unl.edu.ar/sinc/1044sinc-publications/2016/SPM16