# Automatic identification of species with neural networks

A new automatic identification system using photographic images has been designed to recognize fish, plant, and butterfly species from Europe and South America. The automatic classification system integrates multiple image processing tools to extract the geometry, morphology, and texture of the images. Artificial neural networks (ANNs) were used as the pattern recognition method. We tested a data set that included 740 species and 11,198 individuals. Our results show that the system performed with high accuracy, reaching 91.65% of true positive fish identifications, 92.87% of plants and 93.25% of butterflies. Our results highlight how the neural networks are complementary to species identification, which is useful in today´s taxonomic crisis.

1 **Automatic identification of species with neural networks.**

2 **Andrés Hernández-Serna [a, b *] and Luz Fernanda Jiménez-Segura [a]**

3 [a] Grupo de Ictiología; Instituto de Biología, Universidad de Antioquia, Medellín, Colombia

4 [b] Department of Biology, University of Puerto Rico-Río Piedras, San Juan, Puerto Rico.

5 *Corresponding autor: Andres Hernández-Serna, PO Box 23360, Department of Biology,

6 University of Puerto Rico, San Juan, PR 00931-3360, USA

7  e-mail: andres137@gmail.com

8 **ABSTRACT**

9 A new automatic identification system using photographic images has been designed to

10 recognize fish, plant, and butterfly species from Europe and South America. The automatic

11 classification system integrates multiple image processing tools to extract the geometry,

12 morphology, and texture of the images. Artificial neural networks (ANNs) were used as the

13 pattern recognition method. We tested a data set that included 740 species and 11,198

14 individuals. Our results show that the system performed with high accuracy, reaching 91.65% of

15 true positive fish identifications, 92.87% of plants and 93.25% of butterflies. Our results

16 highlight how the neural networks are complementary to species identification, which is useful in

17 today´s taxonomic crisis.

18 **Keywords:** Fish, plant, butterflies, neural network, feature extraction, digital image, and species

19 **INTRODUCTION**

20  Currently, species identification is a taxonomic challenge and an integral process of all biological

21  research, which generates important information for biodiversity conservation.  Difficulty

22  identifying species and ambiguity in the species concept, are seriously affecting our ability to

23  estimate levels of biodiversity (Gaston & O'Neill, 2004). The Global Taxonomy Initiative

24  highlights the knowledge gaps in our taxonomic system due to the shortage of trained

25  taxonomists and curators; these deficiencies reduce our ability to understand, use, and conserve

26  biological diversity. High levels of global biodiversity and a limited number of taxonomists

27  represents significant challenges to the future of biological study and conservation.  The main

28  problem is that almost all taxonomic information exists in languages and formats not easily

29  understood or shared without a high level of specialized knowledge and vocabularies.  Thus,

30  taxonomic knowledge is localized within limited geographical areas and among a limited number

31  of taxonomists.  This lack of accessibility of taxonomic knowledge to the general public has been

32  termed the "taxonomic crisis" (Dayrat, 2005).


33  Recently, taxonomists have been searching for more efficient methods to meet species

34  identification requirements, such as developing digital image processing and pattern recognition

35  techniques.  These methods automatically identify species based on extracting unique image

36  shape information that distinguishes them by taxonomic groups.  Researchers currently have

37  recognition techniques for  insects, plants, spiders, and plankton (Gaston & O'Neill, 2004). This

38  approach can be extended even further to field-based identification of organisms such as fish

39  (Strachan, Nesvadba & Allen, 1990; Storbeck & Daan, 2001; White, Svellingen & Strachan,

40  2006; Zion, Alchanatis, Ostrovsky, Barki & Karplus, 2007; Hu, Li, Duan, Han, Chen & Si,

41  2012), insects (Mayo & Watson, 2007; O'Neill, 2007 ; Kang, Song & Lee, 2012), zooplankton

42  (Grosjean, Picheral, Warembourg & Gorsky, 2004) and plants (Novotny & Suk, 2013). These

43 methods are helpful in alleviating the "taxonomy crisis". In this research, we present a new

44 methodology for the identification of different taxonomic groups to the species level for fish,

45 plants, and butterflies.

46 We designed a simple and effective algorithm (preprocess solution) and defined a range of new

47 features that use pattern recognition with artificial neural network designs (ANN). Our

48 experiments are outlined, discussed, and important conclusions on automatic species image

49 identification are summarized.

50 **MATERIALS AND METHODS**

51 **Images**

52 Image data in this study was taking from two sources: natural history museum records, and

53 online databases. Analyses from each collection were done with respect to country. Ichthyology

54 collections from Colombia were compiled from the Instituto de Investigaciones Marinas y

55 Costeras (INVEMAR), the Colección de Referencia Biología Marina Universidad del Valle

56 (CRBMUV), and the Coleccion Ictiologica Universidad de Antioquia (CIUA). Ichthyology

57 collections from Brazil were found in the Museu de Zoologia da USP (MZUSP), the Instituto

58 Nacional de Pesquisas da Amazônia Manaus (INPA), and the Museu Nacional Rio de Janeiro

59 (MNRJ). Image data from Spain came from the Museo Nacional de Ciencias Naturales Madrid

60 (MNCN). We tested a data set that included a total of 740 species and 11,198 individuals of fish,

61 plants, and butterflies. Fish specimen images were taken using a Cannon EOS 6dD one-use

62 camera with a 1280 x 960 pixel resolution. 697 total fish species, previously identified by

63 experts, were photographed (see Fig. 1 for a subset of photographed species). Images of 32 plant

64 species were downloaded from the Flavia database (http://flavia.sourceforge.net/) (see Fig.2).

65  Image data for 11 species of butterflies were downloaded from the MorphBank database

66  (http://www.morphbank.net/) (see Fig. 3.).

67  **System development**

68  Based on pattern recognition theory  (Marqués de Sá, 2001) and basic computer-processing

69  pathways used in typical automated species identification systems (Gaston & O'Neill, 2004), we

70  designed a system for automatic individual identification at the species level (Fig. 4). In a novel

71  way, our system shares preprocess and extraction components with both the training and

72  recognition processes.  Features of training images are used to build a model of the classification

73  progress pattern after feature extraction. These features and the trained model are then recorded

74  in the database and incorporated in the analysis of subsequent photos. This process uses two

75  types of data to model features of recognition files and results in better species identification

76  results. The following sections provide implementation details for each step in Fig. 4.  Due to its

77  size, a list of features could not be included in this manuscript, but is alternatively available upon

78  request.

79  **Image preprocessing**

80  Image heterogeneity in terms of orientation, size, brightness, and illumination was common (Fig.

81  5.1). Image background was removed with Grabcut's algorithm (Rother, Kolmogorov & Blake,

82  2004) (Fig. 5.2) and  converted to grayscale (Fig. 5.3). Different filters were applied to improve

83  the image by removing image noise; the filters used were smooth and median (Fig. 5.4 and 5.5),

84  and the image was then reduced to one of two possible levels, 0 or 1 (Fig. 5.6).  **Next,** the

85  processed image was brought to a contour (Fig. 5.7) and then a skeleton (Fig. 5.8). All of these

86  processes were performed for each taxonomic group using the image processing in MATLAB

87  R2009b.

88 **Feature extraction**

89 Feature extraction greatly influences species identification from image processing. Features

90 should represent taxonomic information and be easily acquired from data images. A series of

91 geometrical, morphological, and texture features, unique to species, are used in our automatic

92 identification system; these features can be efficiently extracted with image processing. Fifteen

93 intuitive features were used in the system and are described below:

94 **Geometrical**

95 Geometric features contain information about form, position, size, and orientation of the region.

96 The following are some geometric features that are commonly used in pattern recognition.

97 1- *Area* is the total number of pixels of the study area, and is defined as:

$$A(s) = \int_x \int_y I(x, y) \, dy \, dx$$

98 I (x, y) depends on the limits of the shape (see figure 5.7).

99 2- *Perimeter.* The number of pixels that belong to the edge of the region (see figure 5.8). In other

100 words, it is the curve that encloses a region S, defined as:

$$P(s) = \int_t \sqrt{x^2(t) + y^2(t)} \, dt$$

101 3- *Diameter.* Value representing the diameter of a circle with the same area as the region.

102 4- *Compatibility.* The efficiency of the contour or perimeter P(s) that encloses an area A(s)

$$C(s) = \frac{4\pi \, A(s)}{P^2(s)}$$

103   5- *Compactness*. The efficiency with which area A(s) encloses an object is determined by P(s)

$$Co(s) = \frac{P^2(s)}{4\pi \, A(s)}$$

104   6- *Solidity*. The scalar specifying the proportion of the pixels in the convex hull that are also in

105   the region. This property is supported only for 2-D input label matrices.

106   *Solidity*. The number of pixels, specified in terms of area/scalar.

107   **Texture**

108   Textures are important visual patterns for homogeneous description of regions. Intuitive

109   measures provide properties such as smoothing, roughness, and regularity (Glasbey, 1996).

110   Textures depend on the resolution of the image and can follow two approaches: statistical and

111   frequency. We use the statistical approximation in which statistical values are analyzed first order

112   (on the histogram) and second order (on the co-occurrence matrix).

113   *Statistical first order* is obtained from the gray level histogram of the image. Each value is

114   divided by the total number of pixels (area) and has a new histogram representing the probability

115   that a determined gray level is displayed in the region of interest.

116   Obtained properties:

117   7- *Median*

$$\mu = \sum_{x=1}^{n} x h(x)$$

118   8-*Variance*

$$\delta^2 = \sum_{x=1}^{n}\left(x - \mu^2\right)h(x)$$

119    *Statistical second order* is the matrix spatial dependence of gray levels or co-occurrence matrices.

120    Given a vector of polar coordinates, $\delta = (r, \theta),$ one can calculate the conditional probability

121    that two properties appear separated by a given distance $\delta, P_{\delta}$ using an angle $\theta$ of -45 and

122    a distance *r* equal to one pixel. The features that are extracted from this matrix are:

123    9- *Uniformity*

$$\sum_{x=1}^{n}\sum_{y=1}^{n}P_{\delta}(x,y)^2$$

124    10- *Entropy co-occurrence*

$$-\sum_{x=1}^{n}\sum_{y=1}^{n}P_{\delta}(x,y)\log P_{\delta}(x,y)$$

125    11- *Homogeneity*

$$\frac{\sum_{x=1}^{n}\sum_{y=1}^{n}P_{\delta}(x,y)}{1+|x-y|}$$

126    12- *Inertia*

$$\sum_{x=1}^{n}\sum_{y=1}^{n}P_{\delta}(x,y)(x-y)^2$$

127    **Morphological**

128    The morphological features are those that concentrate on the organization of pixels. They

129    perform a comprehensive description of the region of interest. They fall into two categories: two-

130    dimensional Cartesian moments and normalized central moments.

131    *Two-dimensional Cartesian moments* are variable at minor order, and initiate at zero at higher

132    orders. The moment of order p and q of a function I (x, y) is defined as:

$$m_{pq} = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} x^p y^q \, I(x,y) \, dx \, dy$$

133    The parameters $p$ and $q$ denote the order of the moment. When $p = 0$ and $q = 0$, which determines

134    the center of mass or gravity of the overall function in binary images, the center of mass or

135    gravity of the region under study is:

$$\bar{x} = \frac{m_{10}}{m_{00}} \quad \bar{y} = \frac{m_{01}}{m_{00}}$$

136    The center of mass or gravity can define the central moments that are invariant to displacement

137    or translation of the image's region of interest defined as:

$$u_{pq} = \sum_x \sum_y (x - \bar{x})^p (y - \bar{y})^q \, I(x,y) \Delta A$$

138    Where $\Delta A$ is the area of a pixel.

139    *Normalized central moments* are invariant to scale which is defined as:

$$n_{pq} = \frac{u_{pq}}{u^{\gamma}_{00}}$$

$$\gamma = \frac{p+q}{2} \quad \forall \, p+q \geq 2$$

140  Where

141  The above equations can be defined by seven moments that are invariant to rotation, translation,

142  and scale changes, known as the Hu invariant set of moments (Hu, 1962). In this study, we used

143  the first Hu moment defined as:

144  13-*Hu1*

$$\varphi_1 = m_{20} + m_{02}$$

145  Normalized central moments can be generated by related moment invariants "AMI" (Flusser &

146  Suk, 1993), based on the theory of algebraic invariants and invariants under general affine

147  transformation. We used two of the four invariants associated with discriminant character

148  moments defined as:

149  14-*Ami1*

$$I_1 = \frac{u_{20}\,u_{02} - u_{11}^2}{u_{00}^4}$$

150  15- *Ami2*

$$I_2 = \frac{u_{30}^2 u_{03}^2 - 6u_{30}u_{21}u_{12} + 4u_{30}u_{12}^3 + 4u_{21}^3 u_{03} - 3u_{21}^2 u_{12}^2}{u_{00}^{10}}$$

151 These moments enable a high degree of insensitivity to noise that is not altered by rotation,

152 translation, or staggering.

153 The use of the above 15 features (Table 1) has two advantages. First, the features can express the

154 structure of the individual's body, which is important for the identification at species level.

155 Second, our features were elaborately chosen to avoid using feature optimization methods like

156 adapted fuzzy reasoning (Lancieri & Boubchir, 2007). We designed and realized automatic

157 extraction algorithms to compute the values of these features so that all variables and features

158 could be calculated automatically.

159 **Neural Network**

160 A neural network is defined as a parallel computer model composed of a large number of

161 adaptive processing (neural) units which communicate via interconnections with variables. A

162 multiple layer network has one or more layers (neurons)  that enable the learning of complex

163 tasks by progressively extracting more meaningful features from the input image patterns (Wu,

164 1997). Compared to other machine learning methods, neural networks learn slower but predict

165 faster and have very good models presenting nonlinear data. The simple perceptron is assigned

166 multiple inputs but generates a single output, similar to different linear combinations that depend

167 on input weights and generate a linear activation function (Rosenblatt, 1958 ). Mathematically,

168 the neural network can be described with the following equation:

$$y = \varphi \left( \sum_{i=1}^{n} w_i * x_i + b \right)$$

169 $W_i$ : weight vector, $X_i$ : input vector, $b$ : bias activation function.

170 A multilayer perceptron consists of a set of source nodes containing one or more input layer and

171 a set of hidden-node outputs. The input signal propagates through the network layer by layer

172 (Zhang, Patuwo & Hu, 1998). Fig. 6 presents a diagram of the multilayer neural network.

173 The neural network structure is composed of N inputs N = [ $N_1$ , $N_2$ , ..., $N_n$], a hidden layer h and

174 an output vector $S$ = [ $S_1$ , $S_2$ , ..., $S_m$ ].  Each $S_i$ is assessed by a single step that transforms the

175 vector S binary signal [0,1]. A supervised training phase, or sigmoid activation,  is based on the

176 back propagation algorithm in which the weights and biases are updated in the direction of the

177 negative gradient of the performance and then updated in the opposite direction( Werbos, 1974;

178 Rumelhart, Hinton & Williams, 1986; Parker, 1987; Smith & Brier, 1996;). The sigmoid

179 activation function for the hidden layer and output layer is determined by the following equation:

$$f(x) = \frac{1}{1 + e^{-x}}$$

180 In this study, the number of input neurons is determined by the number of descriptors that are

181 available in each pattern, which in this case is N=15 (see variables section). The number of

182 neurons in the hidden layer, *h,* has been experimentally determined from the error set data

183 searching for the general training date of the ANN. The number of output neurons is determined

184 by the number of species classified in each database.

185 **RESULTS AND DISCUSSION**

186 All features were extracted from images and defined according to the above mentioned methods.

187 We tested different species from various taxonomic groups, using the developed neural network

188 systems. The results of the main tests with different test species are listed below.

189　Experiments were divided into two groups: 1) images from the training group were used for

190　building the classifications of the model; 2) images from the test group were used for the

191　reorganization and testing of the developed model.

192　To determine the optimal number of neurons given a data image, the relationship between the

193　identification success rate and the number of neurons was explored. Fig. 7 shows this

194　relationship for the different configurations considered. We display values for neuron number for

195　each species in the database in Table 2.

196　Table 3 shows the performance average of the artificial neural networks using image data and the

197　15 analyzing features. The data set was randomly divided into 60-70-80-90% training images,

198　resulting in 40-30-20-10% test images. The results with the highest average accuracy for species

199　identification were networks using 80-90% training and 20-10% test images. For these tests, the

200　declared success rate was related to the number of species. Recognition became more difficult

201　with increased species number, as observed in the fish result collections from MZUSP, INPA,

202　INVEMAR, CRBMUV, and MNCN which averaged below 90% recognition.

203　Similar to previous findings (Strachan et al., 1990; Storbeck & Daan, 2001; White et al., 2006;

204　Zion et al., 2007; Novotny & Suk, 2013), the neural network used classified species from image

205　data.  However, most other studies only employ databases with low levels of species richness

206　usually spanning many different orders and families and are easily classified due to distinct

207　differences in morphological characteristics. Our neural network builds on the work of these

208　networks, and requires low operator expertise, costs, and response time, but also offers high

209　reproducibility, species identification accuracy, and usability. The ANN algorithm is optimized

210　for testing datasets with high levels of species richness, in this case 740 species (11,198

211　individuals) of fishes, plants and butterflies.

212  The predictive ability of the ANNs was affected by the high phenotypic similarity between

213  species in the analysis, for example small fish species such as those from the family Characidos

214  (Annex 1, Fig. 8). The magnitude of this error comes from low phenotypic differences of some

215  species that vary only in minor details, like teeth or fin radii, which hinders classification.

216  However, the error obtained on the neural network model has been low in other taxonomic

217  families (Table 3).  Overall performance of the system achieved high accuracy and precision,

218  with 91.65% true positive fish identifications, 92.87% plant identifications, and 93.25% butterfly

219  identifications. The evaluation of results appears simple at first glance: the comparison of success

220  rates appears sufficient, however upon closer examination, the success rates in tests on closed

221  data sets strongly depend on the number of species and the ratio of test to training image

222  samples. The data sets with a lower species number have higher success rates, possibly explained

223  by species with very distinct morphological characteristics.

224  The strength of this research is in its applicability to combat the "taxonomic crisis". In the past

225  three decades, many promising techniques for fish identification have emerged. Many of them

226  are based on genetics, interactive computer software, image recognition, hydro-acoustics, and

227  morphometrics (Fischer, 2013). In our study, neural networks were tested as a possible method

228  for species identification. However, taking advantage of the fast performance of the ANNs and

229  the speed of modern PCs, further research should explore the applications of the ANN

230  methodology to automate biomass estimation and real-time species classifications.  This could

231  produce useful tools for both scientific and commercial use.  Fischer (2013) concludes that the

232  image recognition methods are useful but their transferability and resolution are poor because

233  species differ between geographic regions.  This is a clear obstacle to future ANN development

234  and network identification success. Our advances in this field in relation to species identification

235  should be developed for specific geographic regions and translated into user-friendly

236  applications.  We support the development of species identification methods that are globally

237  interchangeable but also tailored to regional biodiversity composition.

248  **REFERENCES**

249  **Dayrat B. 2005.** Towards integrative taxonomy. *Biological Journal of the Linnean Society* **85:**

250       407-415.

251  **Erickson G, Jörgensen P, Jörgensen C, Riccardi G, Ronquist F, van Engelen R. 2007.**

252       Morphbank: Web Image Database Technology for Comparative Morphology and

253       Biodiversity Research *http://www.morphbank.net/* (accessed on 18 October 2013)

254  **Fischer J. 2013.** Fish identification tools for biodiversity and fisheries assessments: review

255       and guidance for decision-makers. *FAO, Rome ed: FAO Fisheries and Aquaculture* **107**.

256  **Flavia. (2009).**  Flavia at a glance. *http://flavia.sourceforge.net/* (accessed on 28 October 2013)

257  **Flusser J, Suk T. 1993.** Pattern-recognition by affine moment invariants. *Pattern Recognition*

258       **26:** 167-174.

259   **Gaston KJ, O'Neill MA. 2004.** Automated species identification: why not? *Royal Society*

260       *Philosophical Transactions Biological Sciences* **359:** 655-667.

261   **Glasbey C. 1996.** Handbook of pattern recognition and computer vision - Chen,CH, Pau,LF,

262       Wang,PSP. *Journal of Classification* **13:** 350-352.

263   **Grosjean P, Picheral M, Warembourg C, Gorsky G. 2004.** Enumeration, measurement, and

264       identification of net zooplankton samples using the ZOOSCAN digital imaging system.

265       *Ices Journal of Marine Science* **61:** 518-525.

266   **Hu J, Li D, Duan Q, Han Y, Chen G, Si X. 2012.** Fish species classification by color, texture

267       and multi-class support vector machine using computer vision. *Computers and*

268       *Electronics in Agriculture* **88:** 133-140.

269   **Hu M-K. 1962.** Visual pattern recognition by moment invariants: Information Theory, IRE

270       Transactions. 179 - 187.

271   **Kang S-H, Song S-H, Lee S-H. 2012.** Identification of butterfly species with a single neural

272       network system. *Journal of Asia-Pacific Entomology* **15:** 431-435.

273   **Lancieri L, Boubchir L. 2007.** Using multiple uncertain examples and adaptive fuzzy

274       reasoning to optimize image characterization. *Knowledge-Based Systems* **20:** 266-276.

275   **Marqués de Sá J. 2001.** *Pattern Recognition: Concepts, Methods and Applications*. Springer:

276       Springer.

277   **Mayo M, Watson AT. 2007.** Automatic species identification of live moths. *Knowledge-Based*

278       *Systems* **20:** 195-202.

279   **Novotny P, Suk T. 2013.** Leaf recognition of woody species in Central Europe. *Biosystems*

280       *Engineering* **115:** 444-452.

281   **O'Neill MA. 2007.** DAISY: a practical computed-based tool for semi-automated species

282       identification. *Systematics Association Special Volume Series* **74:** 101-114.

283  **Parker DB. 1987.** Optimal Algorithms for Adaptive Networks: Second order Back Propagation,

284  Second Order Direct Propagation, and Second Order Hebbian Learning: Proceedings of

285  the IEEE First International Conference on Neural Networks \rm San Diego, CA. 593-

286  600.

287  **Rosenblatt. F.** *1958* The Perceptron: A Probabilistic Model for Information Storage and

288  Organization in the Brain. *Psychological Review* **65:** *386-408*

289  **Rother C, Kolmogorov V, Blake A. 2004.** "GrabCut" - Interactive foreground extraction using

290  iterated graph cuts. *Acm Transactions on Graphics* **23:** 309-314.

291  **Rumelhart DE, Hinton GE, Williams RJ. 1986.** Learning representations by back-propagating

292  errors. *Nature* **323:** 533-536.

293  **Smith BP, Brier ME. 1996.** Statistical approach to neural network model building for

294  gentamicin peak predictions. *Journal of Pharmaceutical Sciences* **85:** 65-69.

295  **Storbeck F, Daan B. 2001.** Fish species recognition using computer vision and a neural network.

296  *Fisheries Research* **51:** 11-15.

297  **Strachan NJC, Nesvadba P, Allen AR. 1990.** Fish species recognition by shape-analysis of

298  images. *Pattern Recognition* **23:** 539-544.

299  **Werbos PJ. 1974.** Beyond Regression: New Tools for Prediction and Analysis in the Behavioral

300  Sciences: *Harvard University*.

301  **White DJ, Svellingen C, Strachan NJC. 2006.** Automated measurement of species and length

302  of fish by computer vision. *Fisheries Research* **80:** 203-210.

303  **Wu CH. 1997.** Artificial neural networks for molecular sequence analysis. *Computers &*

304  *Chemistry* **21:** 237-256.

305  **Zhang GQ, Patuwo BE, Hu MY. 1998.** Forecasting with artificial neural networks: The state of

306  the art. *International Journal of Forecasting* **14:** 35-62.

307    **Zion B, Alchanatis V, Ostrovsky V, Barki A, Karplus I. 2007.** Real-time underwater sorting of

308        edible fish species. *Computers and Electronics in Agriculture* **56:** 34-45.

# Table 1(on next page)

Table 1


Features extracted

**Table 1.** Features extracted

| Type | Variable | Description |
|---|---|---|
| Geometrical | $A$ | Area |
| | $P$ | Perimeter |
| | $D$ | Diameter |
| | $C$ | Compatibility |
| | $Co$ | Compactness |
| | $S$ | Solidity |
| Texture | $u$ | Median |
| | $\delta^2$ | |
| | $E_{r,\theta}$ | |
| | $H_{r,\theta}$ | Variance |
| | $HG_{r,\theta}$ | Uniformity |
| | $I_{r|}$ | Entropy co-occurrence |
| | $\varphi_1$ | Homogeneity |
| | $I_1, I_2$ | Inertia |
| Morphological | | Hu1 |
| | | Ami1-Ami2 |

# Table 2(on next page)

table 2


FC (Fish collection); parameters used in neural network systems.

**Table 2.** FC (Fish collection); parameters used in neural network systems.

| Data set | Learning rate | Number of generations | Number of Hidden layers | Number of input layers | Number of output layers (# species) |
|---|---|---|---|---|---|
| FC-MZUSP | 0.2 | 95000 | 200 | 15 | 100 |
| FC-INPA | 0.15 | 100000 | 180 | 15 | 91 |
| FC-MNRJ | 0.25 | 78000 | 60 | 15 | 14 |
| FC-INVEMAR | 0.3 | 84000 | 250 | 15 | 189 |
| FC-CIUA | 0.12 | 90000 | 60 | 15 | 33 |
| FC-CRBMUV | 0.35 | 140000 | 300 | 15 | 172 |
| FC-MNCN | 0.2 | 110000 | 250 | 15 | 98 |
| FLAVIA | 0.1 | 50000 | 60 | 15 | 32 |
| BUTTERFLIES | 0. 5 | 50000 | 35 | 15 | 11 |

# Table 3(on next page)

Table 3


FC (Fish collection); results of ANN tests with species for 15 features

**Table 3.** FC (Fish collection); results of ANN tests with species tests for 15 features

| Data set | Species | Images | Average Percentage of images (Training / test) | | | |
|---|---|---|---|---|---|---|
| | | | 60/40 | 70/30 | 80/20 | 90/10 |
| FC-MZUSP | 100 | 1718 | 76.67 | 81.34 | 83.34 | 88.31 |
| FC-INPA | 91 | 1640 | 76.29 | 78.94 | 84.44 | 89.93 |
| FC-MNRJ | 14 | 422 | 82.62 | 87.18 | 90.56 | 91.65 |
| FC-INVEMAR | 189 | 1703 | 76.72 | 84.03 | 86.45 | 88.08 |
| FC-CIUA | 33 | 472 | 83.08 | 86.99 | 90.19 | 91.77 |
| FC-CRBMUV | 172 | 2392 | 77.36 | 85.21 | 87.29 | 88.85 |
| FC-MNCN | 98 | 959 | 72.34 | 86.21 | 88.15 | 89.11 |
| FLAVIA | 32 | 1800 | 68.79 | 88.48 | 91.61 | 92.87 |
| BUTTERFLIES | 11 | 92 | 73.62 | 80.43 | 88.83 | 93.25 |

# Figure 1

Figure 1

Samples of some species data set : 1)*Curimata mivartii 2)Leporinus striatus*3)*Ctecolucius hujeta*4)*Cinopotamus magdalenae*5)*Astyanax magdalenae*6)*Roeboides occidentalis*7)*Genycharax tarpon*8)*Cyphocharax magdalenae*9)*Hemibrycon decurrens* 10)*Brycon medemi*11)*Lebiasina multimaculata*12)*Hemibrycon dentatus*13)*Triporheus magdalenae*14)*Characidium phoxocephalum*15)*Leporinus muyscorum*16)*Hemibrycon boquiae*17)*Brycon henni*r18*Characidium caucanum*19)*Roeboides dayi*20)*Astyanax fasciatus*21)*Argopleura magdalenensis*22)*Apteronotus eschemeyeri*23)*Eigenmannia virescens.*

# Figure 2

Figure 2

Samples of our data set: 1)*Phyllostachys edulis*2)*Aesculus chinensis*3)*Berberis anhweiensis*4)*Cercis chinensis*5)*Indigofera tinctoria*6)*Acer Dalmatum*7)*Phoebe zhennan*8)*Kalopanax septemlobus*9)*Cinnamomum japonicum*10)*Koelreuteria paniculata*11)*Ilex macrocarpa*12)*Pittosporum tobira*13)*Chimonanthus praecox*14)*Cinnamomum camphora*15)*Viburnum awabuki*16)*Osmanthus fragrans*17)*Cedrus deodara*18)*Ginkgo biloba*19)*Lagerstroemia indica*20)*Nerium oleander*21)*Podocarpus macrophyllus*22)*Prunus yedoensis*23)*Ligustrum lucidum*24)*Tonna sinensis*25)*Prunus persica*26)*Manglietia fordiana*27)*Acer buergerianum*28)*Mahonia bealei*29)*Magnolia grandiflora*30)*Populus Canadensis*31)*Liriodendron chinense*32)*Citrus reticulate.*

# Figure 3

Figure 3

Samples of our data set: 1) *Agraulis vanillae* 2) *Anthocharis midea* 3) *Ascia monuste* 4) *Danaus gilippus* 5) *Danaus plexippus* 6) *Dryas iulia* 7) *Enodia portlandia* 8) *Glutophrissa Drusilla* 9) *Heliconius charithonia* 10) *Pieres rapae* 11) *Pontia protodice*.
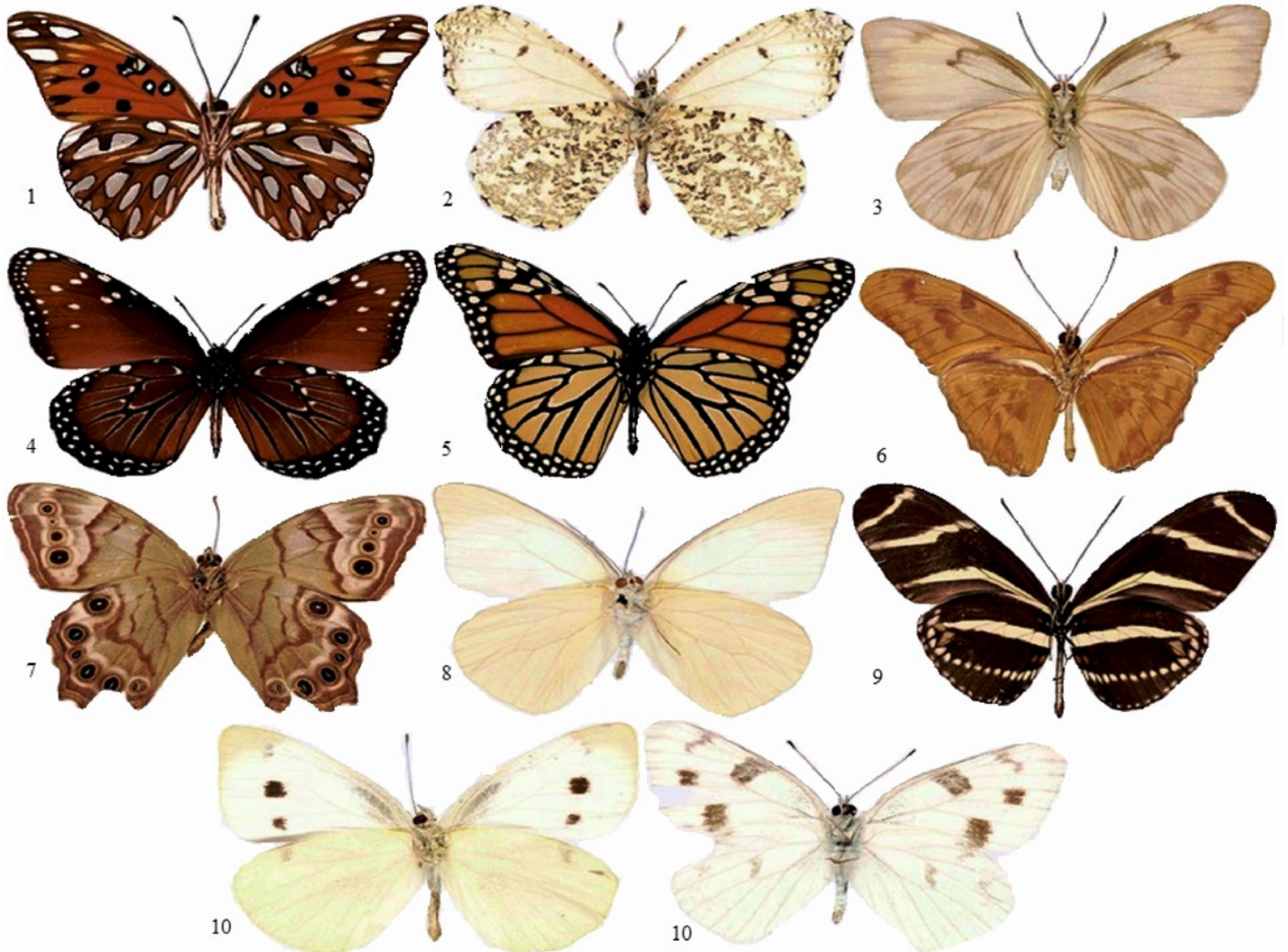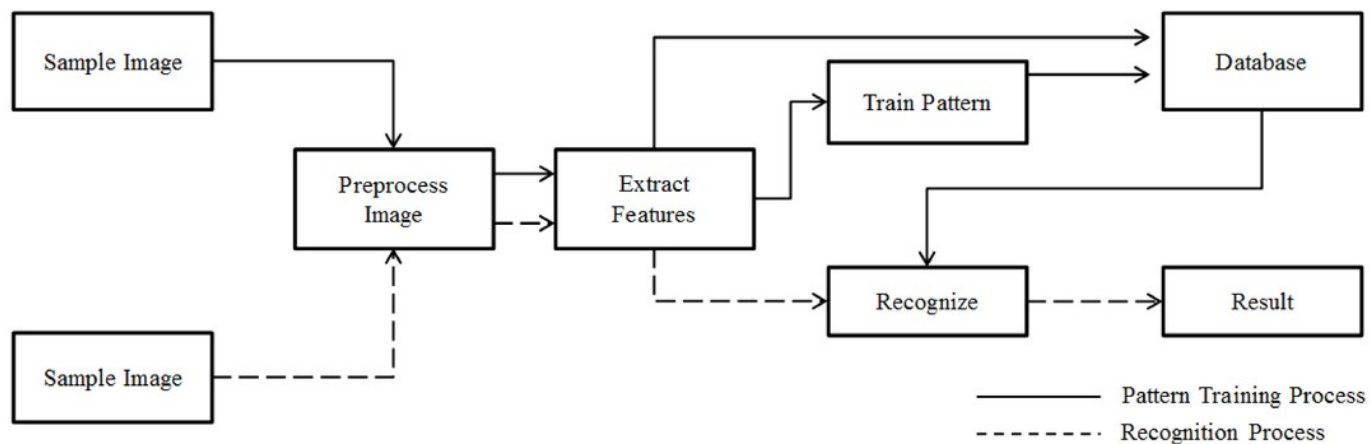
# Figure 4

Figure 4

System architecture

# Figure 5

Figure 5

Image processing 1) jpg image, 2) Image background is removed, 3) grayscale image, 4) smoothing filter, 5) median filter, 6) binarized image, 7) contour image 8) skeletonized image.
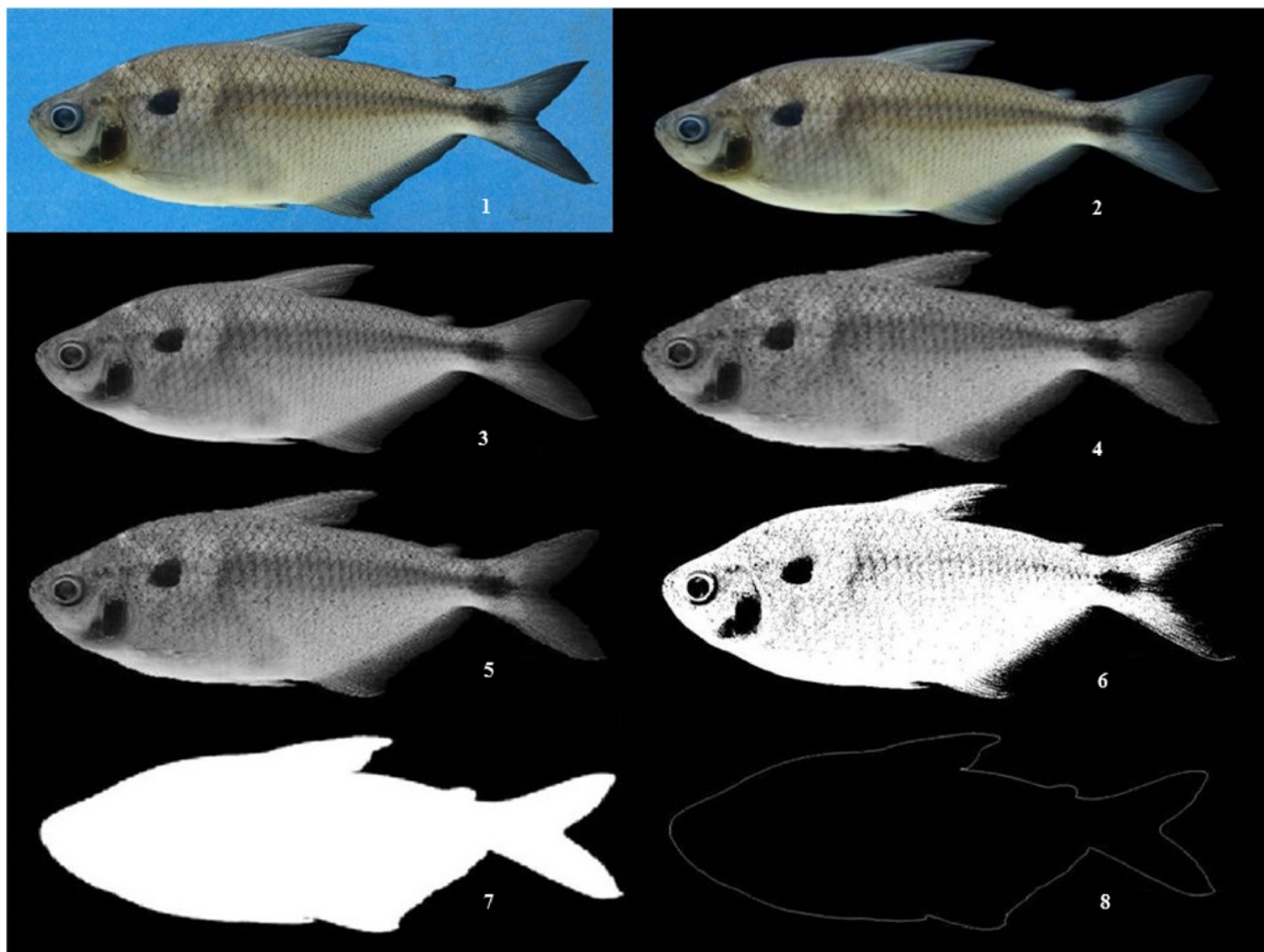
# Figure 6

Figure 6
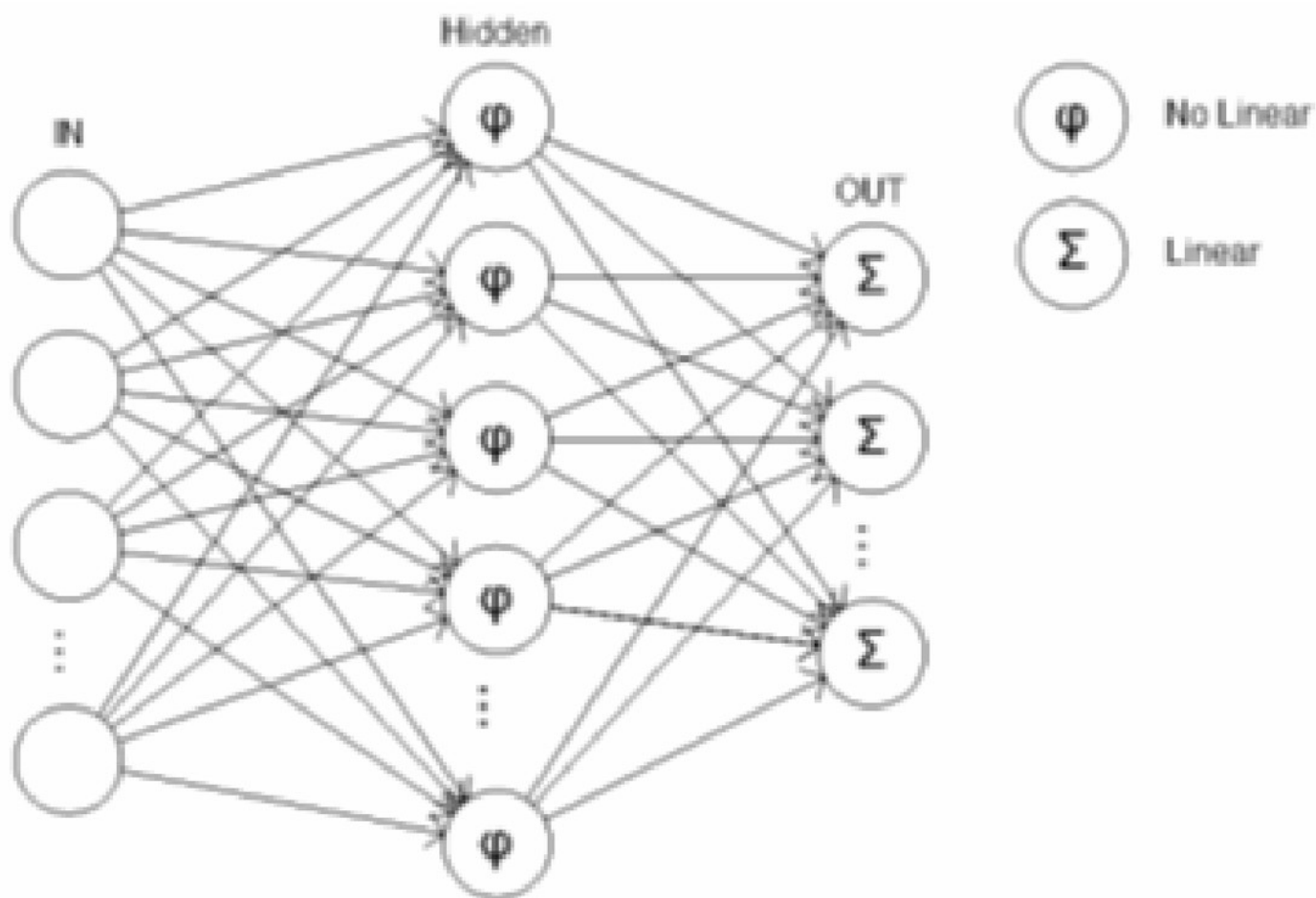
General architecture of a multilayer perceptron.

# Figure 7

Figure 7

Relationship between the success rate and the number of neurons for each neural network.
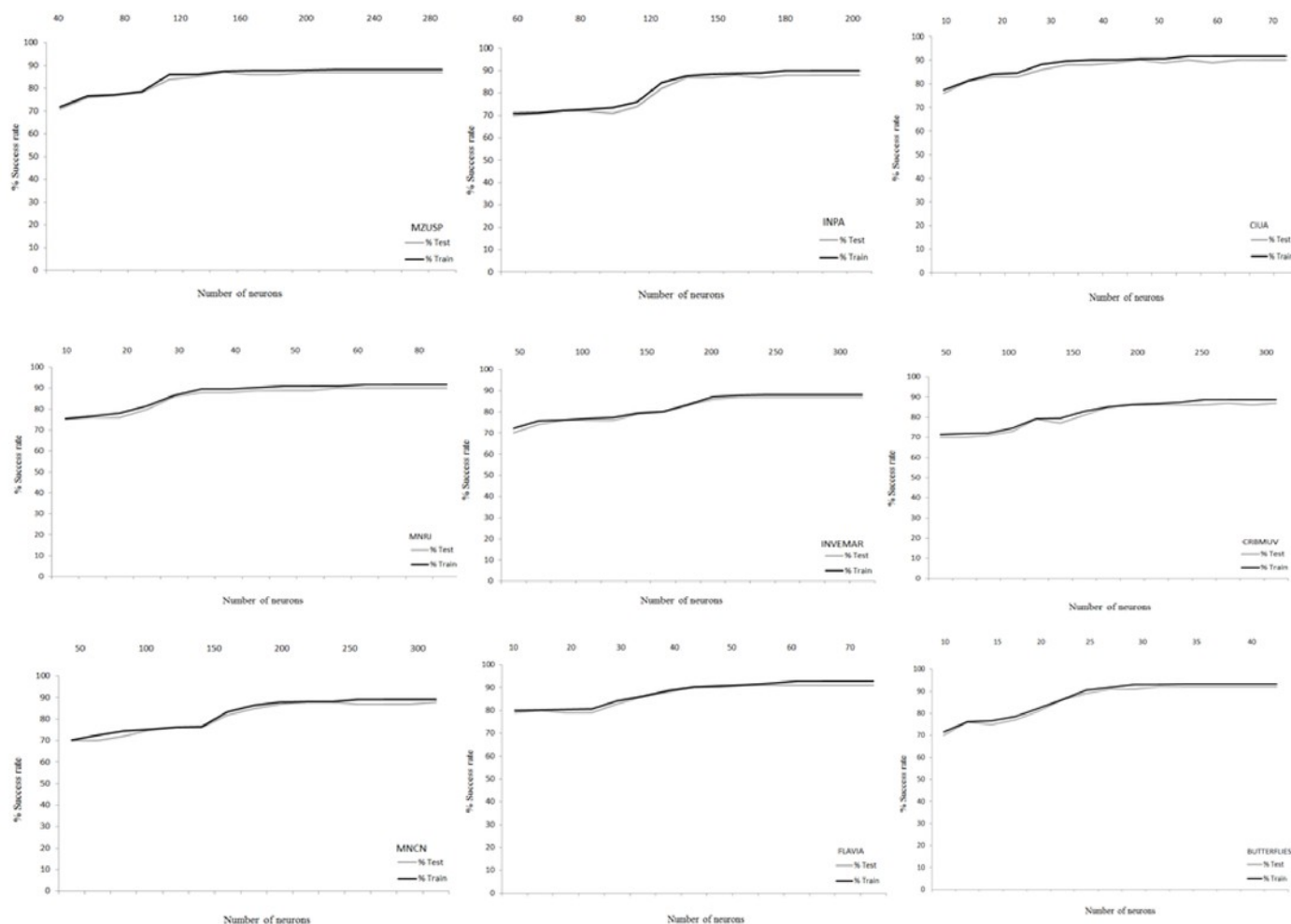
# Figure 8

An example of species confusion in the genus Astyanax 1)*Astyanax magdalenae*, 2)*Astyanax caucanus*, 3)*Astyanax fasciatus*, and4) *Astyanax microlepis.*