



<https://doi.org/10.11646/megataxa.6.2.1>

<http://zoobank.org/urn:lsid:zoobank.org:pub:EBAE6B9F-E1B1-4946-9540-04C579A6BB27>

## iTaxoTools 0.1: Kickstarting a specimen-based software toolkit for taxonomists

MIGUEL VENCES<sup>1\*</sup>, AURÉLIEN MIRALLES<sup>2</sup>, SOPHIE BROUILLET<sup>3</sup>, JACQUES DUCASSE<sup>4</sup>, ALEXANDER FEDOSOV<sup>5</sup>, VLADIMIR KHARCHEV<sup>6</sup>, IVAYLO KOSTADINOV<sup>7</sup>, SANGEETA KUMARI<sup>8</sup>, STEFANOS PATMANIDIS<sup>9</sup>, MARK D. SCHERZ<sup>10</sup>, NICOLAS PUILLANDRE<sup>11</sup> & SUSANNE S. RENNER<sup>12</sup>

<sup>1</sup>Department of Evolutionary Biology, Zoological Institute, Technische Universität Braunschweig, Mendelssohnstraße 4, 38106 Braunschweig, Germany

[m.vences@tu-braunschweig.de](mailto:m.vences@tu-braunschweig.de); <https://orcid.org/0000-0003-0747-0817>

<sup>2</sup>Institut de Systématique, Évolution, Biodiversité (ISYEB), Muséum national d'Histoire naturelle, CNRS, Sorbonne Université, EPHE, Université des Antilles 57 rue Cuvier, CP 50, 75005 Paris, France

[miralles.skink@gmail.com](mailto:miralles.skink@gmail.com); <https://orcid.org/0000-0002-2538-7710>

<sup>3</sup>Institut de Systématique, Évolution, Biodiversité (ISYEB), Muséum national d'Histoire naturelle, CNRS, Sorbonne Université, EPHE, Université des Antilles 57 rue Cuvier, CP 50, 75005 Paris, France

[sophie.brouillet@mnhn.fr](mailto:sophie.brouillet@mnhn.fr); <https://orcid.org/0000-0002-0845-4272>

<sup>4</sup>Independent researcher, 49 rue Eugène Carrière, 75018 Paris, France

[jacko21@aliceadsl.fr](mailto:jacko21@aliceadsl.fr); <https://orcid.org/0000-0002-6884-7844>

<sup>5</sup>A.N. Severtsov Institute of Ecology and Evolution, Russian Academy of Sciences, Leninsky prospect 33, 119071 Moscow, Russian Federation

[fedosovalexander@gmail.com](mailto:fedosovalexander@gmail.com); <https://orcid.org/0000-0002-8035-1403>

<sup>6</sup>Department of Evolutionary Biology, Zoological Institute, Technische Universität Braunschweig, Mendelssohnstraße 4, 38106 Braunschweig, Germany

[necrosoverign@gmail.com](mailto:necrosoverign@gmail.com); <https://orcid.org/0000-0002-1410-7381>

<sup>7</sup>GFBio - Gesellschaft für Biologische Daten e.V., c/o Research II, Campus Ring 1, 28759 Bremen, Germany

[ikostadi@gfbio.org](mailto:ikostadi@gfbio.org); <https://orcid.org/0000-0003-4476-6764>

<sup>8</sup>Department of Evolutionary Biology, Zoological Institute, Technische Universität Braunschweig, Mendelssohnstraße 4, 38106 Braunschweig, Germany

[kumari13@live.com](mailto:kumari13@live.com); <https://orcid.org/0000-0002-8165-8319>

<sup>9</sup>School of Electrical and Computer Engineering, National Technical University of Athens, Iroon Polytechniou St 9, 15780 Athens, Greece

[stefanpatman91@gmail.com](mailto:stefanpatman91@gmail.com); <https://orcid.org/0000-0003-3547-6140>

<sup>10</sup>Faculty of Mathematics and Natural Sciences, Institute for Biochemistry and Biology, University of Potsdam, Potsdam, Germany

[mark.scherz@gmail.com](mailto:mark.scherz@gmail.com); <https://orcid.org/0000-0002-4613-7761>

<sup>11</sup>Institut de Systématique, Évolution, Biodiversité (ISYEB), Muséum national d'Histoire naturelle, CNRS, Sorbonne Université, EPHE, Université des Antilles 57 rue Cuvier, CP 50, 75005 Paris, France

[nicolaspuillandre@gmail.com](mailto:nicolaspuillandre@gmail.com); <https://orcid.org/0000-0002-9797-0892>

<sup>12</sup>Department of Biology, Washington University, 1 Brookings Drive, Saint Louis, MO 63130, USA

[renner@lmu.de](mailto:renner@lmu.de); <https://orcid.org/0000-0003-3704-0703>

\* Corresponding author; [m.vences@tu-braunschweig.de](mailto:m.vences@tu-braunschweig.de)

### Abstract

While powerful and user-friendly software suites exist for phylogenetics, and an impressive cybertaxonomic infrastructure of online species databases has been set up in the past two decades, software targeted explicitly at facilitating alpha-taxonomic work, i.e., delimiting and diagnosing species, is still in its infancy. Here we present a project to develop a bioinformatic toolkit for taxonomy, based on open-source Python code, including tools focusing on species delimitation and diagnosis and centered around specimen

identifiers. At the core of iTaxoTools is user-friendliness, with numerous autocorrect options for data files and with intuitive graphical user interfaces. Assembled standalone executables for all tools or a suite of tools with a launcher window will be distributed for Windows, Linux, and Mac OS systems, and in the future also implemented on a web server. The initial version (iTaxoTools 0.1) distributed with this paper (<https://github.com/iTaxoTools/iTaxoTools-Executables>) contains graphical user interface (GUI) versions of six species delimitation programs (*ABGD*, *ASAP*, *DELINEATE*, *GMYC*, *PTP*, *tr2*) and a simple threshold-

clustering delimitation tool. There are also new Python implementations of existing algorithms, including tools to compute pairwise DNA distances, ultrametric time trees based on non-parametric rate smoothing, species-diagnostic nucleotide positions, and standard morphometric analyses. Other utilities convert among different formats of molecular sequences, geographical coordinates, and units; merge, split and prune sequence files, tables and species partition files; and perform simple statistical tests. As a future perspective, we envisage iTaxoTools to become part of a bioinformatic pipeline for next-generation taxonomy that accelerates the inventory of life while maintaining high-quality species hypotheses. The open source code and binaries of all tools are available from Github (<https://github.com/iTaxoTools>) and further information from the website (<http://itaxotools.org>).

**Key words:** integrative taxonomy, molecular diagnosis, species delimitation, ABGD, PTP, GMYC, TR2, DELINEATE, Limes

## Introduction

Bioinformatics has become the core of modern biology, especially in the context of high-throughput workflows that are becoming commonplace in many fields, in particular related to -omics approaches. The big data volumes obtained by these techniques require ever more efficient and sophisticated software, which is being developed and refined at a vigorous pace. In the field of systematics, powerful programs for phylogenetic analysis have been developed, and databases and data aggregators have been set up to deal with the massive globally-generated taxonomic dataset comprised of over one million species and many millions of specimen records. Furthermore, the increasing diversity and heterogeneity of the types of data used in taxonomy represents a poorly addressed challenge in terms of objective and rationalized data integration. Yet, only few bioinformatic tools so far have been tailored specifically to fit the work of taxonomists, who diagnose and name some 15,000–20,000 new species of organisms per year, a task that still is largely performed by single or small teams of (professional and amateur) researchers (Miralles *et al.* 2020). Most existing tools are aimed at valorizing or transferring raw taxonomic knowledge to the scientific community or to society, but few have been designed to facilitate the actual production of that knowledge in the form of basic taxonomic research. The most notable applications are for instance dedicated to the construction of identification keys (e.g. Dallwitz, 1974; Clark 2003; Delgado Calvo-Flores *et al.* 2006; Zhang *et al.* 2006; MacLeod 2008; Vignes Lebbe *et al.* 2015; Tofilski 2018), which in some groups help field identification. Only a handful of software packages (EDIT: [cybertaxonomy.eu](http://cybertaxonomy.eu), TaxonWorks: [taxonworks.org](http://taxonworks.org), Scratchpads: [scratchpads.org](http://scratchpads.org)) are aimed at facilitating descriptive work itself, but none of these are widely used;

furthermore, these programs do not include important aspects of the alpha-taxonomic workflow, such as species delimitation or molecular diagnosis (Miralles *et al.* 2020), which can be relevant for other fields, e.g. molecular ecology.

Although most taxonomic studies are still relying on morphology only (as shown in a recent review; Miralles *et al.* 2020), taxonomy increasingly integrates diverse lines of evidence (Padial *et al.* 2010), a procedure called integrative taxonomy by Dayrat (2005). Discovering, delimiting, diagnosing, and naming new species requires taxonomists to examine voucher specimens and associated catalogues, field books and pictures; take, tabulate and statistically analyze morphometric measurements; define, tabulate and document phenotypic character states; estimate geographical ranges based on specimen provenances; align and analyze DNA sequences; cross-reference all these types of data, and elaborate accurate specimen tables, species diagnoses and identification keys. Depending on the organism under study, more specialized approaches may be needed, such as comparing acoustic and visual signal repertoires of animals, or having to isolate and culture unicellular organisms. In addition, to fulfil standards of cybertaxonomy, data sets need to be archived in specialized repositories and new species names registered in online databases (Miralles *et al.* 2020). With rising best-practice standards, these many and varied tasks generally involve the use of different computer programs—and thus lead to an extra burden on taxonomists who may lack bioinformatic training. The present paper introduces a new multi-function tool intended for taxonomists, iTaxoTools, which is available both as single-module executables and as a launcher including all modules from <https://github.com/iTaxoTools/iTaxoTools-Executables>.

## The concept of iTaxoTools

We aim to develop a bioinformatic platform to facilitate the core work of taxonomists, that is, delimiting, diagnosing, and describing species. Our initiative produced an integrative taxonomy toolkit—iTaxoTools (Fig. 1; Table 1). The concept of iTaxoTools is based on taxonomic practice and aims at solving some of the technical difficulties that taxonomists are regularly confronted with, while minimizing the time spent in laborious tasks that can be automated. It rests on four pillars: (1) **fully open-source** code; (2) a **diversified** set of stand-alone programs (‘modules’) that in future versions will become increasingly interconnected; (3) a **specimen-centered** architecture, where at present tables (tab-delimited text files) with specimen identifier columns serve as main input format; and (4) a focus on **user-friendliness**, accessibility, and clear and transparent documentation.

All of the code developed by us is **fully open-source** and available from a dedicated GitHub repository (<https://github.com/iTaxoTools>). In the case of tools programmed by other researchers, we make this information transparent,

**TABLE 1.** Overview of the software tools and functionalities currently included in the 0.1 version of iTaxoTools. The majority of the tools can be run (i) command-line driven in Python. They are also distributed as (ii) a standalone executable (.exe) files, and (iii) as part of a full package with launcher window (Fig. 1) in a single executable or part of a folder. Furthermore, (iv) all tools will be implemented on a webserver. Note that functionalities for pre-existing species delimitation tools are explained in more detail in the original papers.

<b>Tool</b>	<b>Purpose</b>	<b>Main functionalities</b>
dnacconvert	Converts among DNA sequence formats	<ul style="list-style-type: none"> <li>- Supports typical sequence formats (fasta, fastq, phylip, nexus)</li> <li>- Autocorrects typical errors in sequence files such as non-standard characters in sequence names.</li> <li>- Reads GenBank flat files and converts from and to tab-delimited files to manage sequences in spreadsheet editors.</li> <li>- Single-file conversion, batch conversion and conversion of copy-pasted files</li> <li>- Parses a large variety of formats of WGS84 geographical coordinates</li> <li>- Batch-conversion of coordinates in tables or copy-pasted lists which can contain coordinates in different formats (recognized by heuristic approach)</li> <li>- Main output in decimal degree format</li> </ul>
latlonconverter	Converts among different geographic coordinate formats	<ul style="list-style-type: none"> <li>- Can merge large files that usually cannot be opened in editors.</li> <li>- Works for any text file but includes additional features when processing fasta and fastq.</li> <li>- Allows for filtering sequences and sequence names with certain motifs and to include/exclude them in the merged file</li> </ul>
fastmerge	Merges DNA sequence files (fasta, fastq)	<ul style="list-style-type: none"> <li>- Can split large files of several GB in size that usually cannot be opened in editors into a series of equally sized smaller files</li> <li>- Designed for fasta and fastq, but works for any text file.</li> <li>- Allows for filtering sequences and sequence names with certain motifs and to include/exclude them in the split files</li> </ul>
fastsplit	Splits (large) DNA sequence files (fasta, fastq) into smaller files	<ul style="list-style-type: none"> <li>- Takes as input a tab-delimited file and a series of values of specimen identifiers</li> <li>- Removes all rows from the table where the column "specimen" (or other chosen column) agrees with any of the provided values</li> </ul>
specimentablepruner	Removes rows from tables based on a list of values for the row "specimen"	<ul style="list-style-type: none"> <li>- Takes as input two or more tab-delimited files, compares values in column "specimen" (or other chosen column) and merges into one table, combining values for same specimen number in the same row</li> <li>- Automatically checks for duplicate values of the same variable and specimen and issues warnings</li> </ul>
linebreaker	Changes among line break formats (Unix, Windows, Mac)	<ul style="list-style-type: none"> <li>- Takes as input any text file and changes all line breaks to the specified format</li> </ul>
simplestatscalculator	Performs a series of basic statistical analyses based on manually entered data	<ul style="list-style-type: none"> <li>- Values are typed or pasted into text boxes</li> <li>- Descriptive statistics (mean, median, standard deviation, and others)</li> <li>- Pairwise comparisons (t-test, U-test)</li> <li>- Comparisons of distributions (Chi-square, normality, Fisher's)</li> <li>- Corrections for multiple testing</li> </ul>

...Continued on the next page

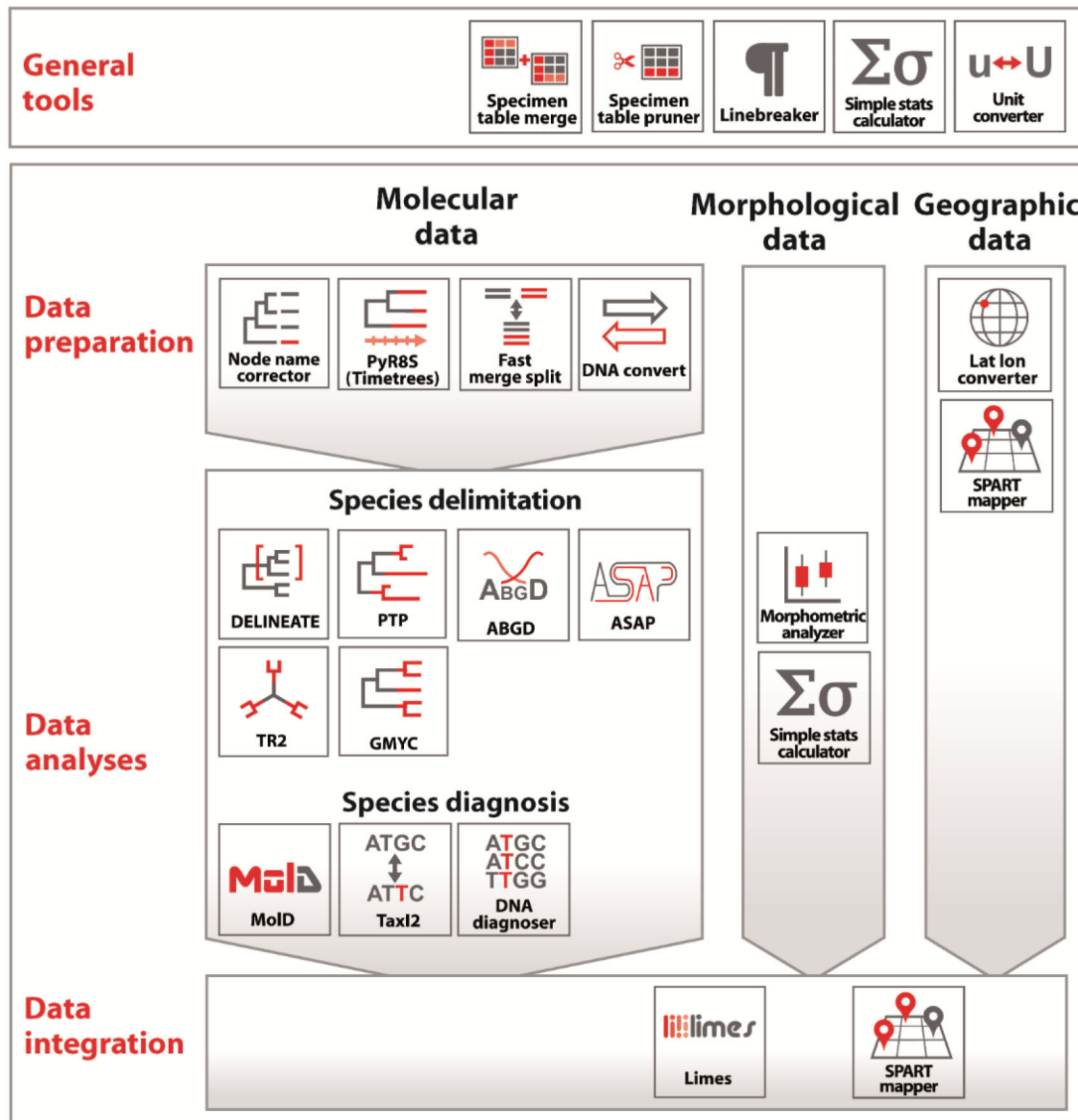
**TABLE 1.** (Continued)

<b>Tool</b>	<b>Purpose</b>	<b>Main functionalities</b>
unitconverter	Converts among different units	<ul style="list-style-type: none"> <li>- Values are typed into one field, all other fields show converted values in real time</li> <li>- Separate tabs for conversion of distance, volume, time, molarity, and others</li> </ul>
spartmapper	Computes a kml file from geographical coordinates and spart file	<ul style="list-style-type: none"> <li>- Takes as input a text file with decimal geographical coordinates and specimen identifiers, and a species partition (SPART) file</li> <li>- Outputs a kml file to show localities by species on Google Earth</li> </ul>
nodenamerecorrector	Removes special characters from terminal node names in Newick-formatted trees	<ul style="list-style-type: none"> <li>- Takes as input a Newick treefile, identifies node names, searches for characters in node names that are not standard alphanumerical, and replaces them with underscores</li> </ul>
pyr8s	Calculates ultrametric timetrees (chronograms) based on non-parametric rate-smoothing	<ul style="list-style-type: none"> <li>- Takes as input treefiles with branch length (phylograms)</li> <li>- Transforms into ultrametric using non-parametric rate smoothing, without the need to access original data (sequences)</li> <li>- User-friendly interface to set time constraints (calibrations) on nodes.</li> </ul>
TaxID	Calculates inter- and intraspecific distances and the barcoding gap based on pairwise-aligning DNA sequences	<ul style="list-style-type: none"> <li>- Takes as input aligned or unaligned sequence files in fasta or tab-delimited text format</li> <li>- For unaligned sequences, pairwise alignments are performed</li> <li>- Calculates pairwise genetic distances among all sequences</li> <li>- If tab file contains row with species names, inter- and intra-species distances are calculated and summarized, and the barcoding gap as well as some summary statistics of the barcoding gap calculated</li> <li>- Inter-species distances are calculated separately for species of the same genus vs. different genera</li> <li>- A histogram illustrating the barcoding gap is produced in editable PDF format</li> </ul>
morphometricanalyzer	Calculates a series of basic statistical comparisons of species based on morphometric data	<ul style="list-style-type: none"> <li>- Takes as input tab-delimited files with morphometric measurements (continuous variables)</li> <li>- Allows specifying if analyses should be done by species, by sex/stage, or by species and sex/stage</li> <li>- Calculates summary statistics, pairwise comparisons (t-tests, U-tests), ANOVAs, PCA and DA</li> <li>- Size-corrects values by calculating ratio against a reference measurement such as body size</li> <li>- Outputs boxplots and scatterplots of PCA and DA factors, by species and/or sex/stage in editable PDF format</li> <li>- Writes text output summarizing diagnostic characters (scientifically different measurements between species, with and without overlap of ranges)</li> </ul>
dnadiagnoser	Computes diagnostic sites for species from DNA sequences	<ul style="list-style-type: none"> <li>- Takes as input aligned or unaligned sequence files in fasta or tab-delimited text format</li> <li>- Unaligned sequences are pairwise aligned to reference sequence and differences recorded relative to position in reference</li> <li>- Summarizes variation within species and outputs diagnostic sites among species</li> <li>- Outputs unique diagnostic sites for the whole data sets, as well as diagnostic sites in pairwise comparisons among species</li> <li>- Output is given in the form of tables but also as text which can be used for species diagnoses in taxonomic papers</li> </ul>

...Continued on the next page

TABLE 1. (Continued)

Tool	Purpose	Main functionalities
PTP	Species delimitation based on Poisson tree processes	<ul style="list-style-type: none"> <li>- Uses as input a non-ultrametric tree with branch lengths (phylogram) in Newick or Nexus format</li> <li>- Models speciation on branching events in terms of number of mutations (inferred from branch lengths)</li> <li>- Bayesian and ML versions of PTP are implemented</li> </ul>
GMYC	Species delimitation based on the Generalized Mixed Yule Coalescent	<ul style="list-style-type: none"> <li>- Uses as input an ultrametric tree in Newick or Nexus format</li> <li>- Uses a likelihood approach to analyze the timing of branching events, seeking for significant switches between a Yule (interspecific) and a coalescent (intraspecific) branching structure.</li> </ul>
tr2	Species delimitation using Bayesian model comparison and rooted triplets	<ul style="list-style-type: none"> <li>- Takes as input a set of gene trees, and optionally a guide species tree</li> <li>- Calculates posterior probability scores for user-specified delimitation hypotheses.</li> <li>- Alternatively, finds the best delimitation under a guide tree specifying a hierarchical structure of species grouping.</li> </ul>
DELINEATE	Species delimitation by integrating an explicit model of speciation into the multipopulation coalescent	<ul style="list-style-type: none"> <li>- Takes as input a rooted ultrametric tree from a multispecies coalescent analysis, in Nexus or Newick format</li> <li>- Second input file is a table assigning specimens to species, or flagging their species identity as unknown</li> <li>- Outputs various alternative species partitions, ranked by probability</li> </ul>
ABGD	Species delimitation by automatic barcoding gap discovery	<ul style="list-style-type: none"> <li>- Takes as input a set of aligned sequences and calculates pairwise distances</li> <li>- Uses a coalescent model to identify the position of the most likely barcode gap, based on a maximal genetic intraspecific divergence defined a priori by the user.</li> <li>- Uses the DNA barcoding gap to propose species partitions.</li> </ul>
ASAP	Species delimitation from single-locus sequence data by the Assemble Species by Automatic Partitioning approach	<ul style="list-style-type: none"> <li>- Takes as input a set of aligned sequences and calculates pairwise distances</li> <li>- Proposes species partitions ranked by a new scoring system that uses no biological prior insight of intraspecific diversity.</li> </ul>
LIMES 2.0	Compare species partitions by different indexes and parsing/merge/export SPART files	<ul style="list-style-type: none"> <li>- Reads species partition (SPART) files, as well as species partition information in spreadsheet format</li> <li>- Computes <math>C_{\text{tax}}^{\text{mC}}</math>, <math>R_{\text{tax}}</math> and Match Ratio indexes</li> <li>- Can merge, extract and export SPART files</li> </ul>
MoID	Recovers DNA-based diagnoses for taxa from DNA sequence alignments	<ul style="list-style-type: none"> <li>- Recovers diagnostic combinations of nucleotides (DNCs) for pre-defined groups of DNA sequences, corresponding to taxa</li> <li>- Identifies pure diagnostic sites, minimal DNCs (mDNCs), and redundant DNCs (rDNCs), the latter fulfill pre-defined criteria of reliability</li> </ul>



**FIGURE 1.** Overview of the various tools implemented in iTaxoTools, and their scope. In the present version a focus is on molecular data analysis, but more functionalities to analyze and visualize morphological and geographic data will be implemented in the near future, while data integration remains the main focus for long-term implementation.

and the graphical user interface (GUI) we added specifies the original references and programmers. The current pre-release of compiled executables is available from <https://github.com/iTaxoTools/iTaxoTools-Executables>. See Table 2 for repositories of each individual tool. GitHub also offers an option for user feedback; users encountering bugs or having suggestions for additional functions of the tools can report them by creating an “Issue” in the respective GitHub repository which can then be dealt with by the iTaxoTools team.

The toolkit is **diversified**, including simple format converters of molecular or geographic data, text and spreadsheet merging and pruning, simple statistical analyses e.g. of morphometric data, but especially focuses on two main aspects: species delimitation and diagnosis, based on multiple kinds of data.

The distribution of the tools is likewise diversified, including command-line tools for those users familiar/comfortable with Python; standalone GUI executables

of each module for Windows, Linux, and Mac operating systems for those looking for a single functionality, e.g. a converter, to be called from a single and easily portable file—these tools will necessarily be “heavier” and slightly slower than command-line executables; and a single software package containing all libraries (currently developed for Windows and Linux), from which each module can be launched (Fig. 2). In the future, the latter software package will enable data transfer between different modules. As far as possible, the exchange files between these different tools are based on compatible standards, in order to minimize the number of format conversion steps (e.g. FASTA or NEXUS for sequences, NEWICK for tree topologies, and SPART for species partitions: Lipman & Pearson 1985; Maddison *et al.* 1997; Miralles *et al.* 2021). Additional efforts to make analytical and integrative processes more fluid through an ergonomically designed architecture represent nevertheless a priority task for future versions of iTaxoTools. The GUI software versions

**TABLE 2.** Repositories of the code of the tools included in the 0.1 version of iTaxoTools. The table also lists the main programmers involved in the development of each tool or its graphical user interface (GUI), and informs whether a tool was newly programmed for this project, adjusted from existing code (by adding a GUI plus sometimes additional functionalities), or included as original code and GUI without modification.

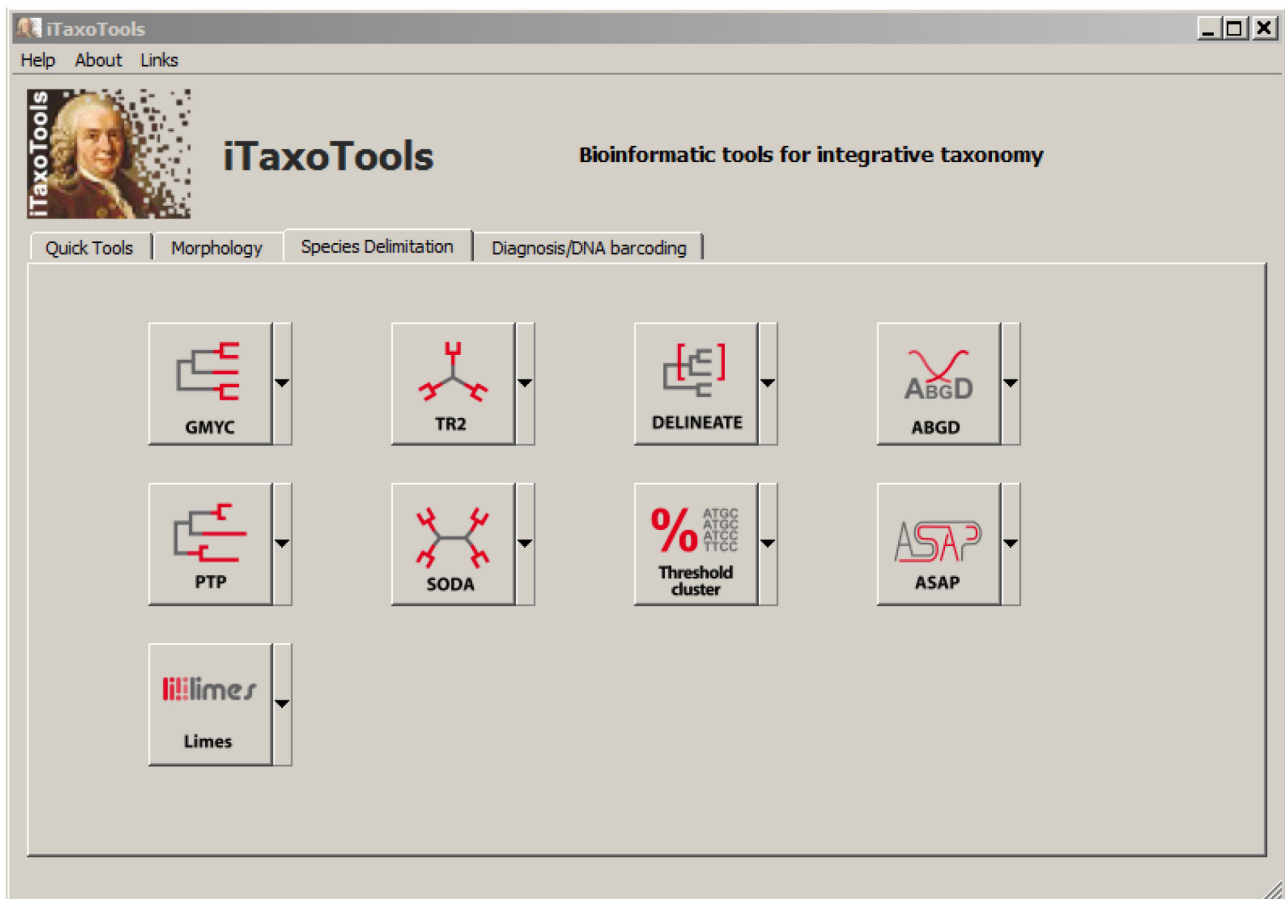
<b>Tool</b>	<b>New / Adjusted / Original</b>	<b>Github repository (original / modified)</b>	<b>Main programmers (original program) / GUI</b>
iTaxoTools executables (repository for the binaries in Windows, Linux and Mac format)	New	<a href="https://github.com/iTaxoTools/iTaxoTools-Executables">https://github.com/iTaxoTools/iTaxoTools-Executables</a>	NA
dnaconvert	New	<a href="https://github.com/iTaxoTools/DNAconvert">https://github.com/iTaxoTools/DNAconvert</a>	V. Kharchev
latlonconverter	New	<a href="https://github.com/iTaxoTools/latlon-converter">https://github.com/iTaxoTools/latlon-converter</a>	V. Kharchev
fastmerge	New	<a href="https://github.com/iTaxoTools/fastsplit-merge">https://github.com/iTaxoTools/fastsplit-merge</a>	V. Kharchev
fastsplit	New	<a href="https://github.com/iTaxoTools/fastsplit-merge">https://github.com/iTaxoTools/fastsplit-merge</a>	V. Kharchev
specimentablepruner	New	<a href="https://github.com/iTaxoTools/specimentablepruner">https://github.com/iTaxoTools/specimentablepruner</a>	V. Kharchev
specimentablemerger	New	<a href="https://github.com/iTaxoTools/specimentablemerger">https://github.com/iTaxoTools/specimentablemerger</a>	V. Kharchev
linebreaker	New	<a href="https://github.com/iTaxoTools/linebreaker">https://github.com/iTaxoTools/linebreaker</a>	S. Kumari
simplestatscalculator	New	<a href="https://github.com/iTaxoTools/simple_stat">https://github.com/iTaxoTools/simple_stat</a>	S. Kumari
unitconverter	New	<a href="https://github.com/iTaxoTools/UnitConverter">https://github.com/iTaxoTools/UnitConverter</a>	S. Kumari
spartmapper	New	<a href="https://github.com/iTaxoTools/Spartmapper">https://github.com/iTaxoTools/Spartmapper</a>	S. Kumari
nodenamecorrector	New	<a href="https://github.com/iTaxoTools/nodenamecorrector">https://github.com/iTaxoTools/nodenamecorrector</a>	V. Kharchev
pyr8s	New	<a href="https://github.com/iTaxoTools/pyr8s">https://github.com/iTaxoTools/pyr8s</a>	S. Patmanidis
TaxI2	New	<a href="https://github.com/iTaxoTools/TaxI2">https://github.com/iTaxoTools/TaxI2</a>	V. Kharchev
morphometricanalyzer	New	<a href="https://github.com/iTaxoTools/morphometricanalyzer">https://github.com/iTaxoTools/morphometricanalyzer</a>	V. Kharchev
dnadiagnoser	New	<a href="https://github.com/iTaxoTools/dnadiagnoser">https://github.com/iTaxoTools/dnadiagnoser</a>	V. Kharchev
PTP	Adjusted	<a href="https://github.com/zhangjiajie/PTP">https://github.com/zhangjiajie/PTP</a> <a href="https://github.com/iTaxoTools/PTP-pyqt5">https://github.com/iTaxoTools/PTP-pyqt5</a>	(J. Zhang) GUI: S. Kumari
GMYC	Adjusted	<a href="https://github.com/zhangjiajie/pGMYC">https://github.com/zhangjiajie/pGMYC</a> <a href="https://github.com/iTaxoTools/GMYC-pyqt5">https://github.com/iTaxoTools/GMYC-pyqt5</a>	(J. Zhang ) GUI: S. Kumari
tr2	Adjusted	<a href="https://github.com/xfujisawa/tr2-delimitation-git">https://github.com/xfujisawa/tr2-delimitation-git</a> <a href="https://github.com/iTaxoTools/pyqt5-tr2">https://github.com/iTaxoTools/pyqt5-tr2</a>	(T. Fujisawa) GUI: S. Kumari
DELINEATE	Adjusted	<a href="https://github.com/jeetsukumaran/delineate">https://github.com/jeetsukumaran/delineate</a> <a href="https://github.com/iTaxoTools/pyqt5-delineate">https://github.com/iTaxoTools/pyqt5-delineate</a>	(J. Sukumaran) GUI: S. Kumari
ABGD	Adjusted	<a href="https://bioinfo.mnhn.fr/abi/public/abgd/">https://bioinfo.mnhn.fr/abi/public/abgd/</a> <a href="https://github.com/iTaxoTools/ABGDpy">https://github.com/iTaxoTools/ABGDpy</a>	(S. Brouillet) GUI: S. Patmanidis
ASAP	Adjusted	<a href="https://bioinfo.mnhn.fr/abi/public/asap/">https://bioinfo.mnhn.fr/abi/public/asap/</a> <a href="https://github.com/iTaxoTools/ASAPy">https://github.com/iTaxoTools/ASAPy</a>	(S. Brouillet) GUI: S. Patmanidis
LIMES 2.0	Original	<a href="https://github.com/iTaxoTools/LIMES">https://github.com/iTaxoTools/LIMES</a>	J. Ducasse
MolD	Adjusted	<a href="https://github.com/SashaFedosov/MolD">https://github.com/SashaFedosov/MolD</a> <a href="https://github.com/iTaxoTools/MolD_pyqt5">https://github.com/iTaxoTools/MolD_pyqt5</a>	(A. Fedosov) GUI: S. Kumari

are designed to be stable over many different versions of the respective operating systems, e.g., from Windows 7 to Windows 10.

Alpha taxonomy is a primarily **specimen-centered** research field in which specimens—mostly single individual organisms or parts thereof, or cultured isolates composed of multiple individuals—are grouped into species. Consequently, iTaxoTools has implemented tab-delimited text as standard format for most tools, with one column indicating the specimen identifier. This will

in subsequent versions allow the user to save the output of different tools for each specimen, and combine these results for further analysis. The tab-delimited format also allows easy editing of the data tables in spreadsheet editors. This specimen-based architecture needed for alpha-taxonomic programs remains valid whether specimens are represented by physical vouchers, images, or in the future maybe by full genome sequences.

Simplicity and **user-friendliness** are at the core of the toolkit we are developing. Because most taxonomists



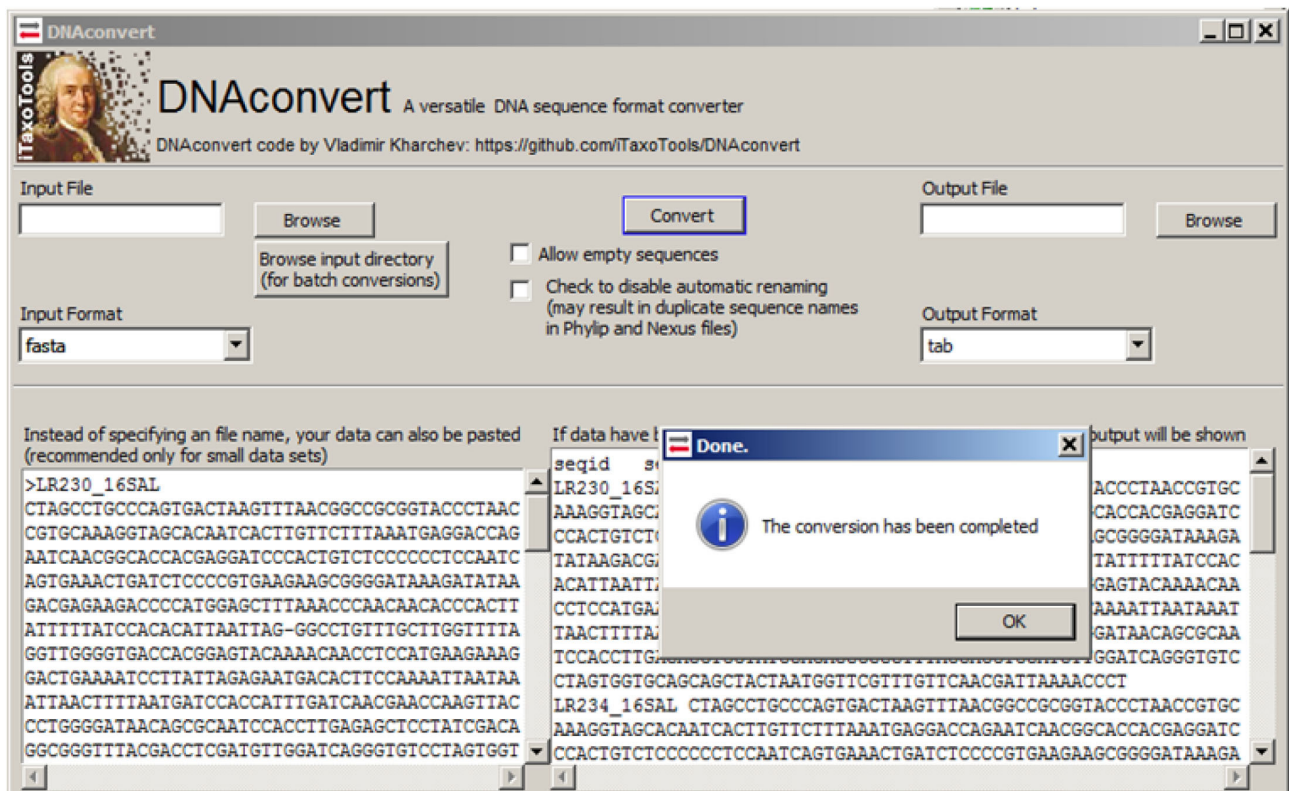
**FIGURE 2.** Main launcher window of iTaxoTools 0.1 with the option to start various species delimitation tools (additional tools can be started from the other tabs).

are not familiar with programming languages, such as Python, all our tools are accessible via GUI—analyses can therefore be carried out with a few intuitive mouse clicks, under default or custom settings, without the need to enter commands in a command line. We have added autocorrect routines to avoid the loss of time associated with the search for small misspellings or incorrect characters in input files that cause programs to fail. Furthermore, we provide a detailed manual with screenshots, along with a collection of example files for each tool, and a dedicated website with FAQs and other useful information (<http://itaxotools.org>); in the future, we will also prepare wikis and video tutorials for the most relevant tools. We chose Python as the main programming language for our package, because it combines good readability with simple-to-learn syntax, and we documented newly written code extensively, to allow its re-use by other programmers. This comes at the cost of speed that would have been achieved by using the C programming language, but our toolkit in this early phase is not designed to cope with huge genomic datasets or analyses with tens of thousands of specimens. Currently iTaxoTools is designed to provide support for the most common taxonomic research projects that discover and name a limited number of species only (Miralles *et al.* 2020), but will be extended to large-scale projects in the future. The executables are partly rather large files because they were compiled to include the entire Python libraries

they make use of. At this initial stage of the project this seemed to be the most reliable option to obtain tools with a stable performance on different platforms.

Considering that powerful programs exist for phylogenetic, phylogenomic, and DNA metabarcoding analyses, we did not attempt to include such functionalities in our toolkit. Similarly, we also did not focus on dedicated multiple sequence alignment programs or genome assemblers because (i) these bioinformatic tasks are more efficiently carried out by programs written in C language, (ii) GUI-driven programs and pipelines already exist for alignment and phylogeny (e.g., *PAUP*, *MEGA*, *BEAST*: Swofford 2003; Kumar *et al.* 2018; Bouckaert *et al.* 2019), genomics, and DNA metabarcoding (e.g., Anslan *et al.* 2017) and (iii) there is an active community both of commercial companies and academic research teams constantly extending these kinds of programs. We are, however, adding graphical user interfaces and new functionalities to other existing tools that are important for analyses in the context of systematics and that are not yet optimized with user-friendly GUIs. For instance, we have updated the code of *Partitionfinder* (Lanfear *et al.* 2016) from Python v. 2 to v. 3, aim to add a GUI also to the sequence alignment program *MAFFT* (Katoh & Standley 2013), and will explore options to include haploweb approaches (Flot *et al.* 2010). These developments will be added successively to iTaxoTools.





**FIGURE 3.** Screenshot of one of the newly programmed quick conversion tools, *dnaconvert*, which implements numerous autocorrect options to avoid sequence output files generating errors in downstream programs. *dnaconvert* also supports tab-delimited table input and its conversion to common sequence formats such as FASTA, NEXUS, or PHYLIP, to facilitate storage and management of sequences and sequence metadata in spreadsheet editors such as Microsoft Excel.

### Functionalities implemented in iTaxoTools 0.1

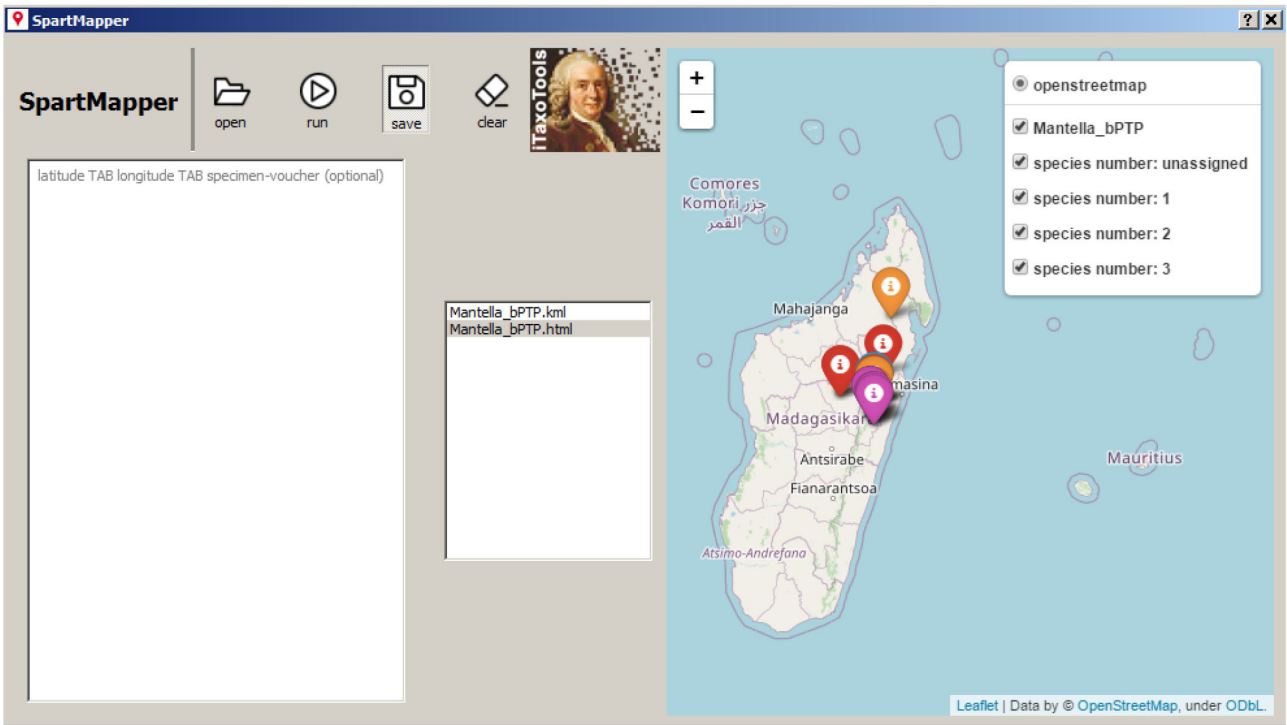
Our work on iTaxoTools is ongoing and will be intensified in the period 2021–2023 thanks to support by the DFG SPP 1991 TaxonOmics priority funding program. The current version, published along with the present article, already includes a series of functional tools that we predict will be useful in different steps of the alpha-taxonomic workflow, (1) Data Preparation (mainly Conversion), (2) Analysis, (3) Delimitation, and (4) Diagnosis.

#### Data Preparation

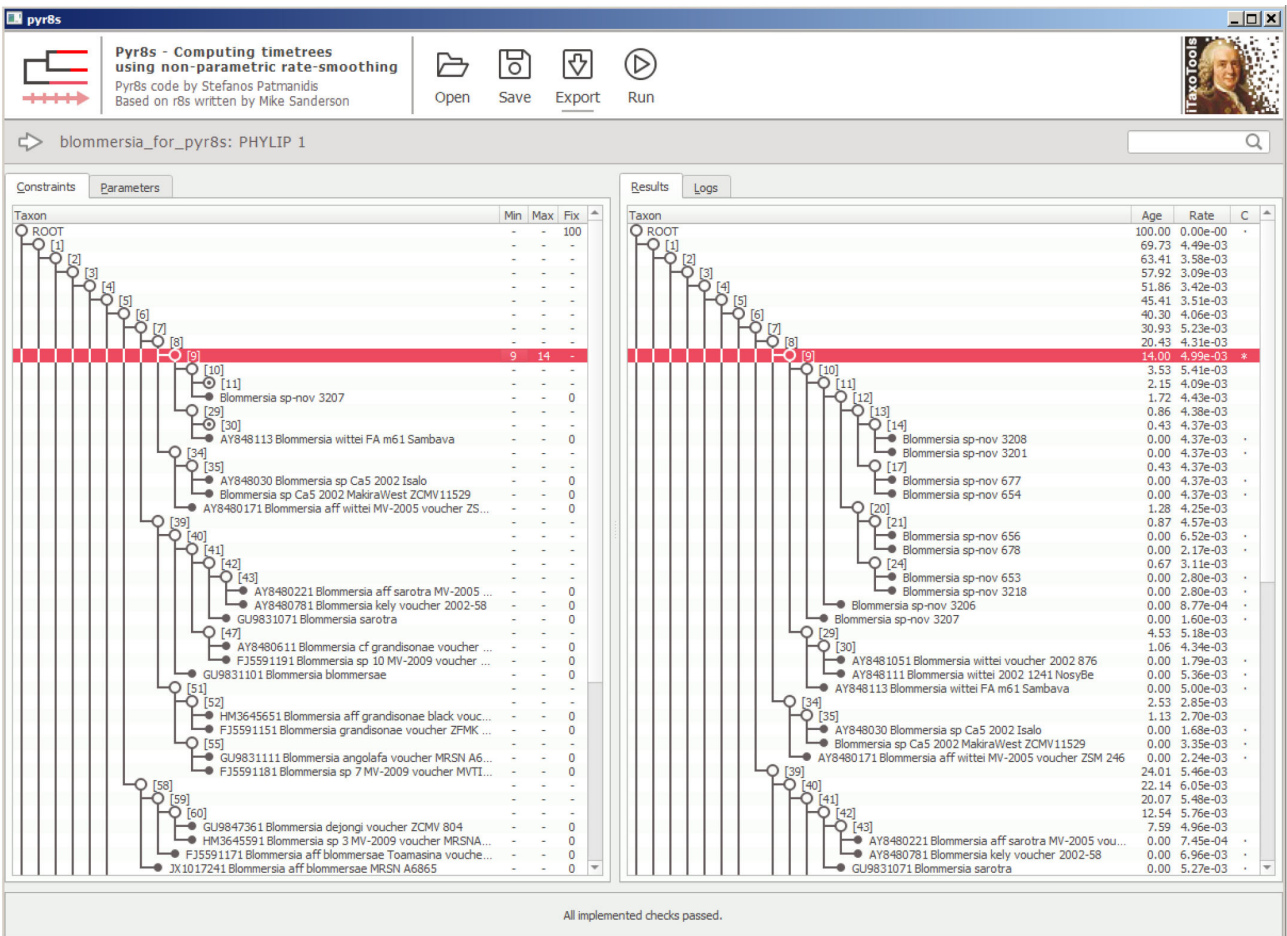
Several of our tools convert among data formats, with the major modules being *dnaconvert* for converting among common DNA sequence formats, *latlonconverter* for converting among geographic coordinate systems (elaborated upon below), and *pyr8s* for converting non-ultrametric trees to ultrametric. A collection of simpler tools includes *fastmerge* and *fastsplit* for splitting and merging large fasta and fastq files, including advanced filtering options by sequence name or sequence motif; *specimentablemerger* and *specimentablepruner* for splitting and merging tab-delimited text files by specimen identifiers; *linebreaker* for converting among Linux and Windows line-break styles (often necessary when processing input files from other bioinformatic tools); *nodenamerector* for replacing all non-standard ASCII

characters from Newick-format trees; and *unitconverter* for distance, time, volume, molarity, and other units.

*dnaconvert* is a versatile tool to convert DNA (and protein) sequence data among commonly used formats such as fasta, fastq, phylip, or nexus (Fig. 3). Compared to other sequence format converters, *dnaconvert* is particularly user-friendly in that it autocorrects numerous issues that usually create compatibility problems, e.g., by automatically replacing non-standard ASCII characters from sequence names or auto-renaming sequences in formats of limited sequence name length such as phylip. The main novelty is the support for tab-delimited files because in our experience, it is useful, for small to medium-sized taxonomy projects, to store and organize specimen-based DNA sequence information (DNA barcodes) in spreadsheet editors such as Microsoft *Excel* or its freeware equivalents *Libre Office / Open Office Calc*. From these spreadsheets, it is then easy to copy and paste the sequence, specimen-voucher, species and locality columns into *dnaconvert* and obtain a sequence file for analysis, e.g. in fasta format, with all respective information concatenated in the sequence name. The program also supports a format in which these metadata are bracketed as required for uploading the sequence data along with metadata to the NCBI Genbank repository (i.e., via Submission Portal or BankIt). Lastly, the program also converts Genbank flatfiles into a tabular format, allowing the user to immediately have all relevant



**FIGURE 4.** Screenshot of spartmapper, a tool that plots distribution records from geographical coordinates on a map and categorizes the records based on a species partition provided as SPART file (Miralles *et al.* 2021). The program allows live view and produces a kml file to visualize the records in Google Earth or Google Maps.



metadata associated with the Genbank record in separate columns in a spreadsheet.

*latlonconverter* allows batch conversion of geographic coordinates from a large number of different formats into standard decimal format as required by most geographical information system (GIS) programs. By performing a series of autocorrections of possible typos and then using a heuristic approach, *latlonconverter* is able to recognize and transform many idiosyncratic formats of geographical coordinates as they are commonly found in specimen databases containing geographical information taken by different researchers. With *spartmapper*, geographical coordinates combined with a species partition file (SPART; Miralles *et al.* 2021) can be previewed on a map, and then transformed into a kml file that plots all localities on Google Earth and visualizes the geographical distribution of the respective species hypotheses (Fig. 4).

*pyr8s* is one of our flagship modules (Fig. 5). For many evolutionary analyses, but also for species delimitation (e.g. *GMYC*), ultrametric phylogenies are required where non-ultrametric trees are available. This conversion is rather complex and can be time-intensive. While numerous programs exist to calculate time trees (e.g., *MCMCtree*, *BEAST*, *MEGAX*; Yang & Rannala 2006; Bouckaert *et al.* 2019; Kumar *et al.* 2018), they usually require DNA sequence information in addition to a previously inferred phylogenetic tree. For iTaxoTools, we opted to recuperate a vintage approach, non-parametric rate smoothing (NPRS), initially developed by Sanderson (1997) and later implemented as part of the program *r8s* (Sanderson 2003). This method only requires a phylogenetic tree as input, with the option to add one or more time calibration points. NPRS has previously been implemented in the R package *ape* (Paradis *et al.* 2004), but was removed from the latter and from the newest releases of *r8s* due to licensing issues. Specifically, the original version of *r8s* relied on a modified implementation of Powell's conjugate direction method which was incompatible with open-source licensing (Powell 1964; Gill *et al.* 1981; Press *et al.* 1992). In the GUI-driven tool *pyr8s*, the NPRS algorithm has been newly coded, making use of the open-source libraries DendroPy (Sukumaran & Holder 2010) and SciPy (Virtanen *et al.* 2020), thus resolving the previous licensing issues. This new version provides a GUI for user-friendly setting of time constraints, includes a Python interface for lower-level analysis and maintains support for *r8s*-formatted input files.

### Analysis

We include several data analysis modules: *TaxI2* for calculation of pairwise distances among individuals, and *morphometricanalyzer* for basic morphometric analyses (elaborated upon below). For convenience, we also include *simplestatscalculator*, a utility for quick, basic statistical analyses of manually entered or pasted data.

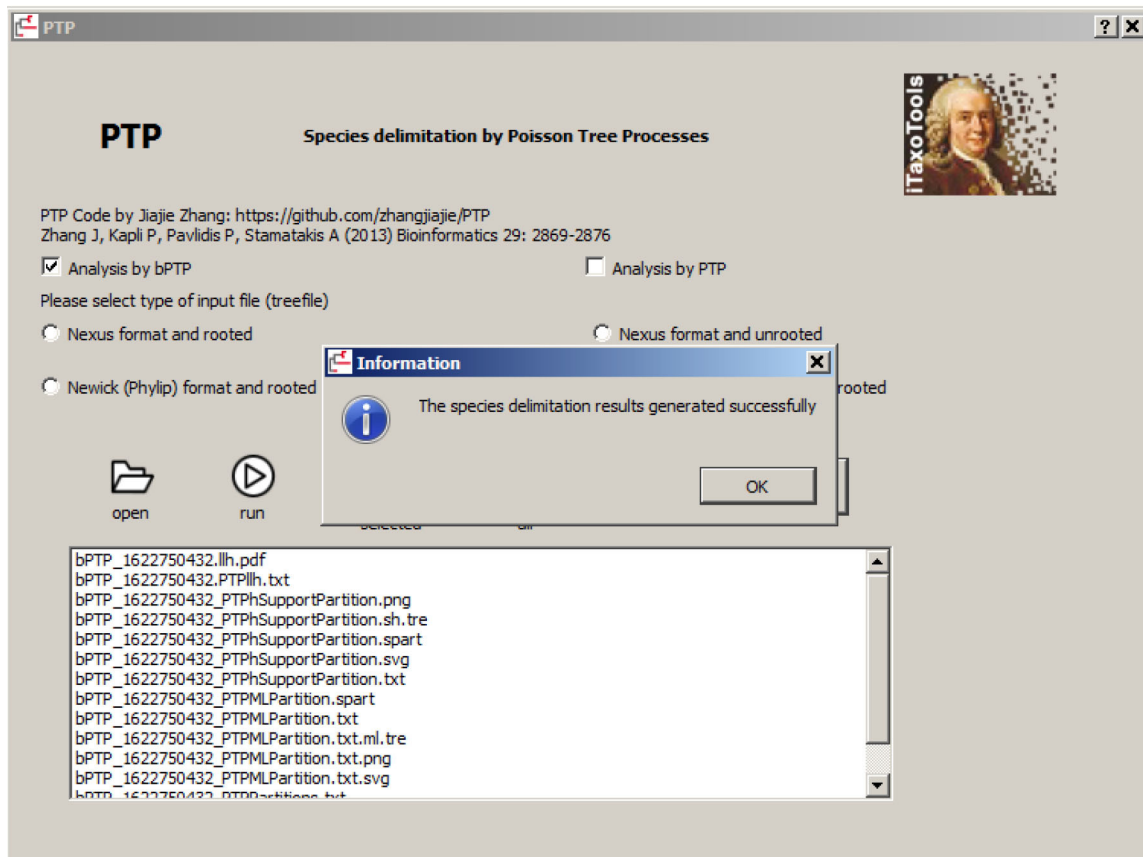
*TaxI2* is a tool for pairwise sequence comparison. To analyze DNA barcoding data sets, Steinke *et al.* (2005) proposed the program *TaxI*, which performs pairwise alignments between sequences and calculates pairwise

distances based on these alignments. Compared to a multiple sequence alignment (MSA) the authors argued that these distance calculations might be more accurate in the case of highly divergent markers including multiple insertions and deletions, such as stretches of mitochondrial ribosomal RNA genes. The pure-Python tool *TaxI2* performs similar calculations, with numerous added functionalities such as support for pre-MSA aligned data sets. The tool has two main analysis modes: First, following the original *TaxI* approach, it compares a set of sequences against a reference database, via pairwise alignments, identifies for each query the closest reference sequence, and calculates various genetic distances among the two. Second, it can also perform all-against-all comparisons of a set of sequences. In this latter approach, sequences can be added in tab-delimited table format along with species name, and from these data the program calculates within-species, between-species, and between-genus distances. Various metrics and graphs defining the barcode gap in a given data set are also included in the output. The program furthermore performs a simple threshold-based clustering of DNA sequences into OTUs, following the approach previously implemented in *TaxonDNA* (<http://taxondna.sourceforge.net/>; Meier *et al.* 2006), and outputs the resulting species partition as SPART file (Miralles *et al.* 2021)

*morphometricanalyzer* is our tool for exploratory analysis of morphometric datasets. Integrative taxonomists do not only use molecular data. In many cases, a limited number of one-dimensional morphometric measurements such as body length and width (or leaf length and width in plants) are taken and compared among groups of individuals. For simple statistical analyses, we have included the tool *morphometricanalyzer* which performs a series of exploratory routine comparisons from morphometric data. It takes as input tab-delimited text files with species hypotheses and a series of other optional categories, and then performs automatically a series of statistical comparisons between species (and between other categories), such as calculations of means, medians, standard deviation, minimum and maximum values; pairwise Mann-Whitney U-tests and Student's t-tests between all pairs of species; a simple Principal Component analysis; and calculation of ratios among original values as a means to size-correct them, followed by statistical comparison of these size-corrected values. Finally, the program also outputs pre-formulated taxonomic diagnoses, with full-text sentences specifying by which morphometric value or ratio a species/population differs from other species/populations with statistical significance, or without value overlap. It would also be possible to explore non-morphological (e.g. bioacoustic) data with this tool, although it is primarily developed for morphometrics.

### Delimitation

A special emphasis in the first development phase of iTaxoTools is species delimitation, a burgeoning field in systematics. The available species delimitation algorithms



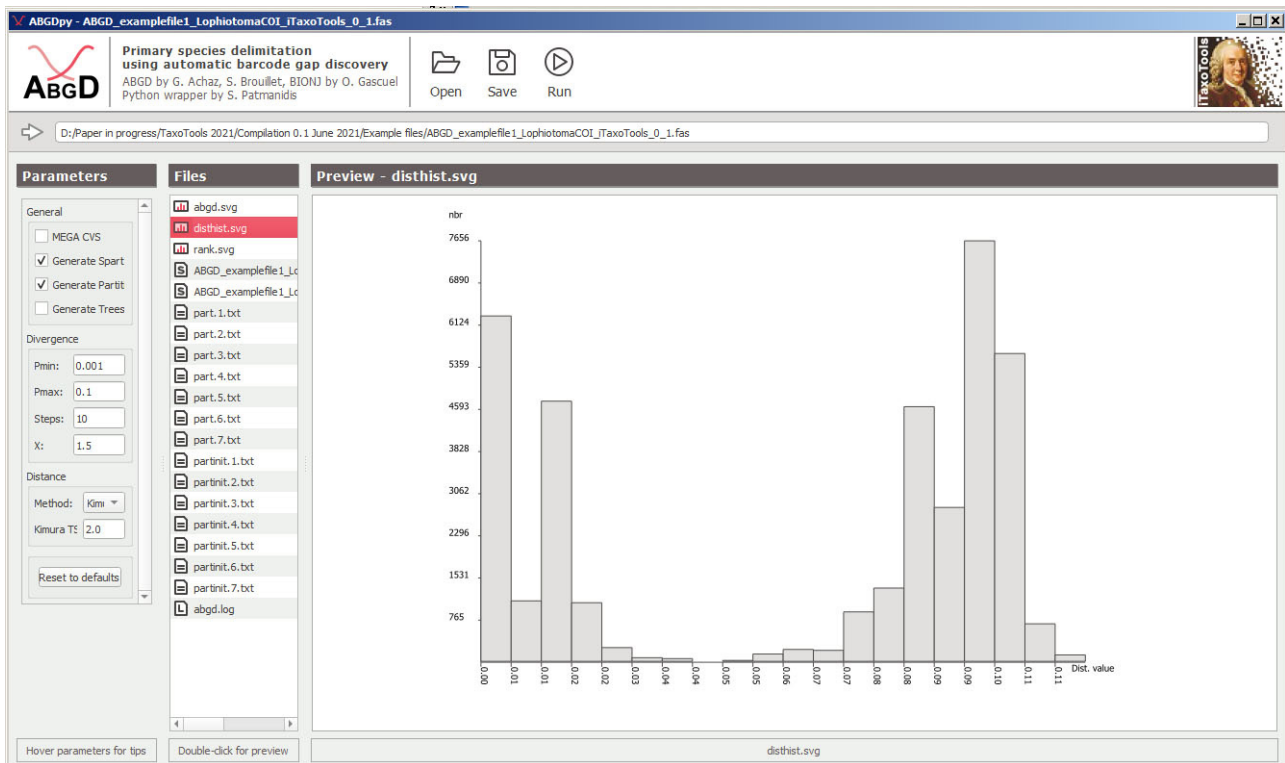
**FIGURE 6.** Screenshot of the GUI-based version of PTP, a program that delimits species from non-ultrametric trees. The original Python code of PTP was written by Zhang *et al.* (2013); iTaxoTools adds the GUI, as well as functionality to export species partition in the SPART format (Miralles *et al.* 2021).

are intended to explore the most likely species partitions within a given dataset and under specific assumptions. They mostly use DNA sequences and implement various species boundary detection criteria, such as the position of the barcode gap, the shape of the coalescent tree, or monophyly. They will therefore likely provide disparate results, and sometimes tend to overestimate the number of species in a data set (e.g., Miralles *et al.* 2013); indeed, they may delimit populations rather than species (Sukumaran & Knowles 2017). Yet, such automated delimitation may play a role in objectivizing the formulation of initial species hypotheses that can then be subsequently tested in an integrative taxonomy pipeline. In the first version of iTaxoTools, we have focused on tools already available in Python programming language. For these tools, we added user-friendly GUIs and slightly extended the functionality, for example, by enabling them to output species partition information in the standardized SPART format proposed by Miralles *et al.* (2021). The current version of iTaxoTools includes GUI-enhanced versions of *PTP* (Zhang *et al.* 2013) (Fig. 6) and *GMYC* (Pons *et al.* 2006; Fujisawa & Barraclough 2013; Python version J. Zhang) that delimit species from single-locus trees; *tr2* (Fujisawa *et al.* 2016) and *DELINEATE* (Sukumaran *et al.* 2020) that use coalescence-based approaches on multiple gene trees; and *ABGD* (Puillandre *et al.* 2012) (Fig. 7) and *ASAP* (Puillandre *et al.* 2021) that are alignment-based and

rely on calculations of genetic distances. iTaxoTools also includes *LIMES 2.0*, a program to handle and compare species partitions obtained by these various approaches (Ducasse *et al.* 2020, Miralles *et al.* 2021).

### Diagnosis

The diagnosis of new species—rather than its lengthy description—represents one of the most important parts of the alpha-taxonomic process, and in all Nomenclatural Codes, diagnoses can be based on molecular, as well as morphological characters (Renner 2016). Several software tools have been proposed to extract diagnostic nucleotide positions of clades and species, either phylogeny-based (*CAOS*; Sarkar *et al.* 2008) or primarily alignment-based (*Mold*, *Fastachar*, *DeSignate*; Fedosov *et al.* 2019; Merckelbach & Borges 2020; Hütter *et al.* 2020). In order to facilitate the use of such DNA characters in differential diagnoses of new species, we implemented a crucial new tool for DNA taxonomy named *dnadiagnoser*. Compared to other tools, *dnadiagnoser* has various functionalities to improve the use of DNA characters in species diagnosis. It takes as input tab-delimited text files in which one column specifies the unit for analysis (typically the species), and provides as output pre-formulated text sentences which specify (i) in a pairwise fashion, all the diagnostic sites of one species against all other species, and (ii) the unique



**FIGURE 7.** Screenshot of the GUI-based version of ABGD, a program that delimits species by detecting the barcoding gap from pairwise single-locus sequence distances (Puillandre *et al.* 2012). For this tool, the original ABGD code written in C was wrapped with a Python GUI and compiled as standalone executable. The different output files produced by ABGD (text and graphs) can be selected and pre-viewed within the GUI.

diagnostic sites (if any) that differentiate a species against all other species. These text sentences can then directly be used in species diagnoses. As a further innovation, *dnadiagnoser* interprets one of the sequences in the input alignment as a reference sequence and outputs the diagnostic sites relative to this sequence. To facilitate such comparisons, the program also includes a series of standard reference sequences (such as the full *Homo sapiens* COI or *cox1* gene) and allows as input unaligned sequences, which are then pairwise aligned against the reference sequence to identify diagnostic positions and label them according to their position in the reference sequence, a procedure that works reliably in sets of sequences with no or only few insertions or deletions such as COI. The original Python code also allows the user to define their own set of reference sequences by simply adding to an existing list, and this option will also be included in future versions of the GUI-driven binaries. In addition, we have programmed a GUI for *Mold* (Fedosov *et al.* 2019), a program that is tailored for recovering DNA-based diagnoses in large DNA dataset, and is capable of identifying diagnostic combinations of nucleotides (DNCs) in addition to single (pure) diagnostic sites. The crucial and unique functionality of *Mold* allows assembling DNA diagnoses that fulfil pre-defined criteria of reliability, which is achieved by repeatedly scoring diagnostic nucleotide combinations against datasets of in-silico mutated sequences.

## Future extensions

Our goal with this paper is to make the tools we have developed available to the community as soon as possible so they may be critically evaluated and improved. The next developments will be in three fields: (i) **Geography:** iTaxoTools will not compete with geographical information systems (GIS), but there are a number of recurrent and rather simple geographical analyses in alpha-taxonomy that can be facilitated by bioinformatic tools, in particular calculation of linear distances among sites and of the surface (minimum convex polygon) of a distribution range of a species, based on a set of georeferenced locality points, and most importantly, a simple graphical editor that outputs publication-ready distribution maps, with customizable colors and symbols for different species, from a set of georeferenced locality records. For more sophisticated analyses, connecting iTaxoTools (via data formats such as SPART) with dedicated toolboxes for analysis of spatial biodiversity data such as *SDMToolbox* (Brown *et al.* 2017) could allow *e.g.* for comparative niche modelling of alternative species partitions. (ii) **Extraction of diagnostic traits from specimen data:** Besides molecular diagnosis with *Mold* and *dnadiagnoser*, we plan to develop a tool that automatically outputs diagnoses based on (specimen-based) categorical data sets of morphological characters. (iii) **Connection to other programs:** We also plan to explore options to connect iTaxoTools to the *DELTA* (DEscription Language

for TAXonomy) software package (Coleman *et al.* 2010). *DELTA* is a format for coding descriptive taxonomic information that, however, is primarily species-based (not specimen-based as iTaxoTools). A series of programs have been developed on this basis, spearheaded by M. Dallwitz at CSIRO (Canberra, Australia) (Dallwitz 1974, 1980). The new Free *DELTA* platform launched in 2000 (<http://freedelta.sourceforge.net/>) includes options for editing and maintenance of data sets in *DELTA* formats, as well as utilities for data conversion, interactive identification of taxa, automated generation of diagnostic keys, and descriptions. Especially, information on species-specific molecular and morphological characters identified in iTaxoTools could be seamlessly coded in *DELTA*, making use of *pydelta* (<http://freedelta.sourceforge.net/pydelta/>).

The biggest gap in taxonomy software so far is the integrative aspect in the sense of Dayrat's (2005) concept of integrative taxonomy (see also Padial *et al.* 2010). That is, the many available species delimitation programs all output a species hypothesis based on one analytical approach—usually based on only a molecular data set, with a few exceptions such as *iBPP*, which can integrate morphometric and molecular data (Solís-Lemus *et al.* 2015). Approaches that combine information from different lines of evidence into species delimitation are exceedingly scarce. One example, *DELINEATE* (Sukumaran *et al.* 2020), allows the user to fix a series of species hypotheses (i.e., firmly assign a series of specimens to species) while letting the other specimens “float” freely in the analysis and assign them to either one of the previously defined species, or to a new species. Such an option of “prior” species delimitation should be universally available to users as a manual option (e.g., if evidence comes from field or experimental data on hybridization, genomic information, or other data that are yet difficult to code or implement in species delimitation software), similar to what is implemented in *DELINEATE*. But ideally, automated proposals of firm *a priori* evidence for two specimens to either belong to two species, or to the same species, could also be elaborated by the software—for instance, using evidence such as sympatric geographical occurrence without gene flow, full concordance between genetic and morphological characters, or unusually high (for the taxon at hand) genetic distances. We plan to develop concepts for such analysis priors (i.e. a tentative formalization of the deductive reasoning used by integrative taxonomists when they are manually integrating sometimes discordant lines of evidence), and start implementing them in an iTaxoTools webserver pipeline (on our website, <http://itaxotools.org/>), in the next years.

Importantly, our project is open for other developers to join, and for the taxonomic community as a whole to provide suggestions. Fully aware of the wide disparity of practices, “philosophies”, and data-types used within the discipline, we especially welcome proposals of additional tools that could help to streamline and accelerate the whole process of delimiting and naming species (whether it concerns the initial step of data acquisition, their treatment, their analyses, or their final submission

to a dedicated repository). Only practicing taxonomists know which parts of the alpha-taxonomic workflow for their group of taxa is particularly time-consuming, and where time and effort is lost with repetitive or error-prone manual tasks that could be equally well performed by a computer program automatically—and thereby formulate requirements for such dedicated programs.

## Perspectives for iTaxoTools

The different taxonomic tools made available here are performing analyses offline on a local computer (and in the future will also be available on a webserver), but without linking to external resources. True next-generation taxonomy will require linking specimen-based taxonomy software with online resources and databases, and scaling the analyses to data of many thousands of specimens. On the one hand, this involves archiving newly acquired data in dedicated repositories (Miralles *et al.* 2020). But on the other hand, it means aggregating DNA sequences, morphological characters, images, and increasingly -omics data (e.g., Lendemer *et al.* 2020) for each specimen identifier (as spearheaded by the Distributed System of Scientific Collections initiative; <https://www.dissco.eu/>), and then entering these large-scale cyberspecimen data into species delimitation, diagnosis, and naming pipelines. The process could be coupled with machine-learning programs to automatically extract diagnostic traits e.g. from images, with data aggregators such as GBIF (<https://www.gbif.org/>) and online tools such as Map of Life (<https://mol.org/>), or Timetree of Life (<http://timetree.org/>) to obtain geographical and temporal context, and distribution models for alternative species hypotheses. These bioinformatic opportunities may gain power under a view of species as probabilistic hypotheses that may allow defining probability thresholds of integrative taxonomic analysis above which lineages can be confidently named as species by semi-automated pipelines. While the current version of iTaxoTools is far from this vision, it may represent a seed for developing the necessary environment, and a sandbox to test software tools with the potential to significantly accelerate the inventory of life.

## Acknowledgments

We are grateful to Mike Sanderson for information on the initial code of *r8s*, and Rashid Kratou for advice. This study was supported by the Deutsche Forschungsgemeinschaft (grant VE247/16-1–HO 3492/6-1 and RE 603/29-1) in the framework of the ‘TaxonOmics’ priority program.

## References

- Anslan, S., Bahram, M., Hiiesalu, I. & Tedersoo, L. (2017) PipeCraft: Flexible open-source toolkit for bioinformatics analysis of custom high-throughput amplicon sequencing data. *Molecular Ecology Resources*, 17, e234–e240. <https://doi.org/10.1111/1755-0998.12692>
- Bouckaert, R., Vaughan, T.G., Barido-Sottani, J., Duchêne, S., Fourment, M., Gavryushkina, A., Heled, J., Jones, G., Kühnert, D., De Maio, N., Matschiner, M., Mendes, F.K., Müller, N.F., Ogilvie, H.A., du Plessis, L., Popinga, A., Rambaut, A., Rasmussen, D., Siveroni, I., Suchard, M.A., Wu, C.H., Xie, D., Zhang, C., Stadler, T. & Drummond, A.J. (2019) BEAST 2.5: An advanced software platform for Bayesian evolutionary analysis. *PLoS Computational Biology*, 15, e1006650. <https://doi.org/10.1371/journal.pcbi.1006650>
- Brown, J.L., Bennett, J.R. & French, C.M. (2017) SDMtoolbox 2.0: the next generation Python-based GIS toolkit for landscape genetic, biogeographic and species distribution model analyses. *PeerJ*, 5, e4095. <https://doi.org/10.7717/peerj.4095>
- Calvo-Flores, M.D., Contreras, W.F., Galindo, E.G. & Pérez-Pérez, R. (2006) XKey: A tool for the generation of identification keys. *Expert Systems with Applications*, 30, 337–351. <https://doi.org/10.1016/j.eswa.2005.07.034>
- Clark, J.Y. (2003) Artificial neural networks for species identification by taxonomists. *Biosystems*, 72, 131–147. [https://doi.org/10.1016/S0303-2647\(03\)00139-4](https://doi.org/10.1016/S0303-2647(03)00139-4)
- Coleman, C.O., Lowry, J.K. & Macfarlane, T. (2010) DELTA for Beginners: An introduction into the taxonomy software package DELTA. *ZooKeys*, 45, 1–75. <https://doi.org/10.3897/zookeys.45.263>
- Dallwitz, M.J. (1974) A flexible computer program for generating identification keys. *Systematic Zoology*, 23, 50–57. <https://doi.org/10.1093/sysbio/23.1.50>
- Dallwitz, M.J. (1980) A general system for coding taxonomic descriptions. *Taxon*, 29, 41–46. <https://doi.org/10.2307/1219595>
- Dayrat, B. (2005) Toward integrative taxonomy. *Biological Journal of the Linnean Society*, 85, 407–415. <https://doi.org/10.1111/j.1095-8312.2005.00503.x>
- Ducasse, J., Ung, V., Lecointre, G. & Miralles, A. (2020) LIMES: a tool for comparing species partition. *Bioinformatics*, 36, 2282–2283. <https://doi.org/10.1093/bioinformatics/btz911>
- Fedosov, A., Achaz, G. & Puillandre, N. (2019) Revisiting use of DNA characters in taxonomy with MOLD - a tree independent algorithm to retrieve diagnostic nucleotide characters from monolocus datasets. *bioRxiv*, 838151. <https://doi.org/10.1101/838151>
- Flot, J.F., Couloux, A. & Tillier, S. (2010) Haplowebs as a graphical tool for delimiting species: a revival of Doyle's "field for recombination" approach and its application to the coral genus *Pocillopora* in Clipperton. *BMC Evolutionary Biology*, 10, 372. <https://doi.org/10.1186/1471-2148-10-372>
- Fujisawa, T., Aswad, A. & Barraclough, T.G. (2016) A rapid and scalable method for multilocus species delimitation using Bayesian model comparison and rooted triplets. *Systematic Biology*, 65, 759–771. <https://doi.org/10.1093/sysbio/syw028>
- Fujisawa, T. & Barraclough, T.G. (2013) Delimiting species using single-locus data and the Generalized Mixed Yule Coalescent approach: a revised method and evaluation on simulated data sets. *Systematic Biology*, 62, 707–724. <https://doi.org/10.1093/sysbio/syt033>
- Gill, P.E., Murray, W. & Wright, M.H. (1981) *Practical Optimization*. Academic Press, New York, 401 pp.
- Hütter, T., Ganser, M.H., Kocher, M., Halkic, M., Agatha, S. & Augsten, N. (2020) DeSignate: detecting signature characters in gene sequence alignments for taxon diagnoses. *BMC Bioinformatics*, 21, 151. <https://doi.org/10.1186/s12859-020-3498-6>
- Katoh, K. & Standley, D.M. (2013) MAFFT Multiple Sequence Alignment Software Version 7: improvements in performance and usability. *Molecular Biology and Evolution*, 30, 772–780. <https://doi.org/10.1093/molbev/mst010>
- Kumar, S., Stecher, G., Li, M., Knyaz, C. & Tamura, K. (2018) MEGA X: Molecular Evolutionary Genetics Analysis across computing platforms. *Molecular Biology and Evolution*, 35, 1547–1549. <https://doi.org/10.1093/molbev/msy096>
- Lanfear, R., Frandsen, P.B., Wright, A.M., Senfeld, T. & Calcott, B. (2016) PartitionFinder 2: new methods for selecting partitioned models of evolution for molecular and morphological phylogenetic analyses. *Molecular Biology and Evolution*, 34, 772–773. <https://doi.org/10.1093/molbev/msw260>
- Lendemer, J., Thiers, B., Monfils, A.K., Zaspel, J., Ellwood, E.R., Bentley, A., LeVan, K., Bates, J., Jennings, D., Contreras, D., Lagomarsino, L., Mabee, P., Ford, L.S., Guralnick, R., Gropp, R.E., Revelez, M., Cobb, N., Seltmann, K. & Aime, M.C. (2020) The extended specimen network: a strategy to enhance US biodiversity collections, promote research and education. *BioScience*, 70, 23–30. <https://doi.org/10.1093/biosci/biz165>
- Lipman, D.J. & Pearson, W.R. (1985) Rapid and sensitive protein similarity searches. *Science*, 227, 1435–1441. <https://doi.org/10.1126/science.2983426>
- MacLeod, N. (Ed.) (2008) *Automated Taxon Identification in Systematics: Theory, Approaches and Applications*. CRC Press, Boca Raton FL, USA, 350 pp.
- Maddison, D.R., Swofford, D.L. & Maddison, W.P. (1997) Nexus: an extensible file format for systematic information. *Systematic Biology*, 46, 590–621. <https://doi.org/10.1093/sysbio/46.4.590>
- Merckelbach, L.M. & Borges, L.M.S. (2020) Make every species count: fastachar software for rapid determination of molecular diagnostic characters to describe species. *Molecular Ecology Resources*, 20, 1761–1768. <https://doi.org/10.1111/1755-0998.13222>
- Meier, R., Kwong, S., Vaidya, G. & Ng, P.K.L. (2006) DNA Barcoding and taxonomy in Diptera: a tale of high intraspecific variability and low identification success. *Systematic Biology*, 55, 715–728. <https://doi.org/10.1080/10635150600969864>
- Miralles, A., Bruy, T., Wolcott, K., Scherz, M.D., Begerow, D., Beszteri, B., Bonkowski, B., Felden, J., Gemeinholzer, B., Glaw, F., Glöckner, F.O., Hawlitschek, O., Kostadinov, I.,

- Nattkemper, T.W., Printzen, C., Renz, J., Rybalka, N., Stadler, M., Weibulat, T., Wilke, T., Renner, S.S. & Vences, M. (2020) Repositories for taxonomic data: Where we are and what is missing. *Systematic Biology*, 69, 1231–1253. <https://doi.org/10.1093/sysbio/syaa026>
- Miralles, A. & Vences, M. (2013) New metrics for comparison of taxonomies reveal striking discrepancies among species delimitation methods in *Madascincus* lizards. *PLoS ONE*, 8, e68242. <https://doi.org/10.1371/journal.pone.0068242>
- Miralles, A., Ducasse, J., Brouillet, S., Flouri, T., Fujisawa, T., Kapli, P., Knowles, L.L., Kumari, S., Stamatakis, A., Sukumaran, J., Lutteropp, S., Vences, M. & Puillandre, N. (2021) SPART, a versatile and standardized data exchange format for species partition information. *BioRxiv*. <https://doi.org/10.1101/2021.03.22.435428>
- Padial, J.M., Miralles, A., De la Riva, I. & Vences, M. (2010) The integrative future of taxonomy. *Frontiers in Zoology*, 7, e16. <https://doi.org/10.1186/1742-9994-7-16>
- Paradis, E., Claude, J. & Strimmer, K. (2004) APE: Analyses of Phylogenetics and Evolution in R language, *Bioinformatics*, 20, 289–290. <https://doi.org/10.1093/bioinformatics/btg412>
- Pons, J., Barraclough, T.G., Gomez-Zurita, J., Cardoso, A., Duran, D.P., Hazell, S., Kamoun, S., Sumlin, W.D. & Vogler, A.P. (2006) Sequence-based species delimitation for the DNA taxonomy of undescribed insects. *Systematic Biology*, 55, 595–609. <https://doi.org/10.1080/10635150600852011>
- Powell, M.J.D. (1964) An efficient method for finding the minimum of a function of several variables without calculating derivatives. *The Computer Journal*, 7, 155–162. <https://doi.org/10.1093/comjnl/7.2.155>
- Press, W.H., Flannery, B.P., Teukolsky, S.A. & Vetterling, W.T. (1992) *Numerical Recipes in C*. Cambridge University Press, New York. 2nd ed, 1018 pp.
- Puillandre, N., Brouillet, S. & Achaz, G. (2021) ASAP: assemble species by automatic partitioning. *Molecular Ecology Resources*, 21(2), 609–620. <https://doi.org/10.1111/1755-0998.13281>
- Puillandre, N., Lambert, A., Brouillet, S. & Achaz, G. (2012) ABGD, Automatic Barcode Gap Discovery for primary species delimitation. *Molecular Ecology*, 21, 1864–1877. <https://doi.org/10.1111/j.1365-294x.2011.05239.x>
- Renner, S.S. (2016) A return to Linnaeus's focus on diagnosis, not description: The use of DNA characters in the formal naming of species. *Systematic Biology*, 65, 1085–1095. <https://doi.org/10.1093/sysbio/syw032>
- Sanderson, M.J. (1997) A non-parametric approach to estimating divergence times in the absence of rate constancy. *Molecular Biology and Evolution*, 14, 1218–1231. <https://doi.org/10.1093/oxfordjournals.molbev.a025731>
- Sanderson, M.J. (2003) r8s: inferring absolute rates of molecular evolution and divergence times in the absence of a molecular clock. *Bioinformatics*, 19, 301–302. <https://doi.org/10.1093/bioinformatics/19.2.301>
- Sarkar, I.N., Planet, P.J. & Desalle, R. (2008) CAOS software for use in character-based DNA barcoding. *Molecular Ecology Resources*, 8, 1256–1259. <https://doi.org/10.1111/j.1755-0998.2008.02235.x>
- Solís-Lemus, C., Knowles, L.L. & Ané, C. (2015) Bayesian species delimitation combining multiple genes and traits in a unified framework. *Evolution*, 69, 492–507. <https://doi.org/10.1111/evo.12582>
- Steinke, D., Salzburger, W., Vences, M. & Meyer, A. (2005) TaxI - A software tool for DNA barcoding using distance methods. *Philosophical Transactions of the Royal Society London, Series B*, 360, 1975–1980. <https://doi.org/10.1098/rstb.2005.1729>
- Sukumaran, J. & Knowles, L.L. (2017) Multispecies coalescent delimits structure, not species. *Proceedings of the National Academy of the U.S.A.*, 114, 1607–1612. <https://doi.org/10.1073/pnas.1607921114>
- Sukumaran, J., Holder, T.M. & Knowles, L.L. (2020) Incorporating the speciation process into species delimitation. <https://github.com/jetsukumaran/delineate>.
- Sukumaran, J. & Holder, M.T. (2010) DendroPy: A Python library for phylogenetic computing. *Bioinformatics*, 26, 1569–1571. <https://doi.org/10.1093/bioinformatics/btq228>
- Swofford, D.L. (2003) *PAUP\*. Phylogenetic Analysis Using Parsimony (\*and Other Methods)*. Version 4. Sinauer Associates, Sunderland, Massachusetts.
- Tofilski, A. (2018) DKey software for editing and browsing dichotomous keys. *ZooKeys*, 735, 131–140. <https://doi.org/10.3897/zookeys.735.21412>
- Vignes Lebbe, R., Chesselet, P. & Diep Thi, M.H. (2015) Xper3: new tools for collaborating, training and transmitting knowledge on botanical phenotypes. In: Rakotoarisoa, N.R., Blackmore, S., Riéra, B. (Eds) *Botanists of the 21st Century*. UNESCO, Paris, 11 pp.
- Virtanen, P., Gommers, R., Oliphant, T.E., Haberland, M., Reddy, T., Cournapeau, D., Burovski, E., Peterson, P., Weckesser, W., Bright, J., van der Walt, S.J., Brett, M., Wilson, J., Millman, K.J., Mayorov, N., Nelson, A.R.J., Jones, E., Kern, R., Larson, E., Carey, C.J., Polat, Y., Feng, Y., Moore, E.W., VanderPlas, J., Laxalde, D., Perktold, J., Cimrman, R., Henriksen, I., Quintero, E.A., Harris, C.R., Archibald, A.M., Ribeiro, A.H., Pedregosa, F., van Mulbregt, P. & SciPy 1.0 Contributors (2020) SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nature Methods*, 17, 261–272. <https://doi.org/10.1038/s41592-019-0686-2>
- Yang, Z. & Rannala, B. (2006) Bayesian estimation of species divergence times under a molecular clock using multiple fossil calibrations with soft bounds. *Molecular Biology and Evolution*, 23, 212–226. <https://doi.org/10.1093/molbev/msj024>
- Zhang, J., Kapli, P., Pavlidis, P. & Stamatakis, A. (2013) A general species delimitation method with applications to phylogenetic placements. *Bioinformatics*, 29, 2869–2876. <https://doi.org/10.1093/bioinformatics/btt499>
- Zhang, X.-B., Chen, X.-X. & Cheng, J.-A. (2006) Lucid Phoenix: A tool for building and deploying interactive, multimedia keys through internet. *Entomotaxonomia*, 28, 231–234.