







SOFTWARE

prepR4pcm: An R Package for Preparing Data and Trees for Phylogenetic Comparative Methods

Shinichi Nakagawa^{1,2}  | Santiago Ortega¹  | Ayumi Mizuno¹  | Eduardo S.A. Santos¹  | Malgorzata Lagisz^{1,2}  | Bhavya Jain¹ | Jimuel Jr Celeste¹  | Sergio Poo Hernandez¹ 

¹COSSEE, Department of Biological Sciences, University of Alberta, Edmonton, AB, Canada; ²Evolution & Ecology Research Centre, School of Biological, Earth and Environmental Sciences, University of New South Wales, Sydney, NSW, Australia

Correspondence: Shinichi Nakagawa, snakagaw@ualberta.ca / s.nakagawa@unsw.edu.au

Running headline: Preparing data and trees for PCMs

Abstract

1. Phylogenetic comparative methods require species names in a trait dataset to match tip labels in a phylogenetic tree. Yet this apparently simple prerequisite is often one of the most fragile steps in a comparative workflow. Names may differ because of, for example, formatting, taxonomic revisions, synonyms, or spelling errors. If these differences are resolved informally, species can be lost from analyses, and the reasons for their loss can be difficult to reconstruct.
2. Here, we present `prepR4pcm`, an R package for preparing data and trees for phylogenetic comparative methods. The package reconciles species names through a staged procedure: exact matching, normalised matching, synonym lookup with local taxonomic databases, and optional fuzzy matching for likely spelling errors. Each decision is stored in a reconciliation object with the original name, matched name, match type, confidence score, and a short explanation. This object turns name matching from a hidden preprocessing step into an auditable part of the analysis.
3. `prepR4pcm` also supports the points where comparative workflows need human judgement. Users can inspect unresolved names, accept or reject suggested matches, add manual corrections, apply taxonomy crosswalks (which link names across taxonomic systems), compare reconciliation runs, and generate reports. The package then returns a matched data frame and pruned tree with the same species set, ready for phylogenetic generalised least squares, phylogenetic mixed models, phylogenetic meta-analysis, and related workflows. If users do not yet have a tree, `prepR4pcm` can retrieve trees from several sources, date trees when suitable information is available, and format tree-source citations.
4. We illustrate the workflow using bundled datasets with realistic name mismatches. `prepR4pcm` is available at <https://github.com/itchyshin/prepR4pcm> with documentation and vignettes covering data and tree reconciliation, tree retrieval, multi-tree workflows, and phylogenetic meta-analysis.

Keywords: *phylogenetic comparative analysis, species names, taxonomic harmonisation, data–tree reconciliation, provenance*

1 Introduction

Phylogenetic comparative methods (PCMs) allow researchers to study trait variation while accounting for shared evolutionary history (Felsenstein, 1985; Grafen, 1989). These methods range from phylogenetic generalised least squares (PGLS; Freckleton et al., 2002) to phylogenetic generalised linear mixed models or meta-analyses (e.g. Hadfield and Nakagawa, 2010; Nakagawa and Santos, 2012; Mizuno et al., 2025). Before fitting such models, researchers need two matched objects: a dataset with species names and a phylogenetic tree with tip labels representing the same species names. Yet this apparently simple prerequisite, matching species names, is often a stumbling block in a comparative workflow, because using scientific species names (binomial strings) as keys causes mismatches that can remove species from the analysis.

The name-matching problem often becomes visible only when the data and the tree are brought together. A species may be present in both sources, but fails to match because one source writes the name as “*Parus major*”, whereas another uses the tip label “*Parus_major*”. In other cases, the two sources may use different names for the same or corresponding taxon, for example, because of an older genus name, a synonym, or a different taxonomic authority. Also, species may be genuinely absent from the tree. These cases could have different biological and practical implications, yet they may have similar effects if simple filtering is applied to remove mismatched species, thereby reducing the number of species in the analysis. For PCMs, this loss matters not only because the retained species set has less information and power, but also because the data analysed and the phylogenetic structure used to analyse them, especially if the removed species are not randomly distributed.

Several R (R Core Team, 2025) packages address parts of this problem. Taxonomic tools, such as `taxize` (Chamberlain and Szöcs, 2013), `taxadb` (Norman et al., 2020), `Taxonstand` (Cayuela et al., 2012), and `bdc` (Ribeiro et al., 2022), help users validate, standardise, or clean species names; the underlying Global Names Architecture (Mozzherin et al., 2017) provides language-agnostic name-parsing and verification services that several of these tools draw on. Tree-retrieval tools, such as `rotl` (Michonneau et al., 2016), `rtrees` (Li, 2023), `clootl` (Miller et al., 2025), `fishtree` (Rabosky et al., 2018; Chang et al., 2019), and `datelife` (Sanchez Reyes et al., 2024), help users obtain phylogenies. PCM helper functions, such as `geiger::treedata()` (Harmon et al., 2008; Pennell et al., 2014) and `caper::comparative.data()` (Orme et al., 2025), help align data and trees for analysis. However, these tools are not designed to provide a single workflow that takes raw species names, from several data objects and produces aligned comparative objects together with an explicit, auditable record of name-matching decisions. The gap is therefore not the absence of useful tools, but the lack of a transparent conduit connecting between name reconciliation, tree preparation, and downstream PCM analysis (Table 1).

Here, we present `prepR4pcm`, an R package for preparing species-level data and phylogenetic trees for PCMs. The package connects three steps that are often handled separately: name reconciliation, tree preparation, and final data-tree alignment. First, it reconciles species names through a staged procedure that combines exact matching, name normalisation, synonym lookup, and optional fuzzy matching. Second, it records each decision in a reconciliation object, allowing users to inspect which names matched, changed, left unresolved, or corrected manually. Third, it applies the reconciliation to return a data frame and a pruned tree containing the same species set, ready for common comparative workflows, including PGLS, phylogenetic mixed models, and phylogenetic meta-analysis. When users do not yet have a suitable tree, `prepR4pcm`

73 can also retrieve trees from several sources, record name substitutions, heck the alignment
74 between the tree and species-level data, and format tree-source citations to credit the original
75 references.

Table 1. Existing R tools for taxonomic name handling and phylogenetic tree retrieval relevant to phylogenetic comparative methods. The first three columns identify each tool and place it in context. Tool = R package or function; Role = the main task the tool was designed to perform; Coverage = the main names, taxa, tree sources, databases, or input objects covered by the tool. The remaining columns summarise features relevant to preparing data and trees for PCMs. Tree = accepts or works directly with a user-supplied `ape::phylo` tree; Prov. = records provenance or an audit trail for name-matching, tree retrieval, or alignment decisions; Fuz. = supports fuzzy matching for likely spelling errors; Syn. = supports synonym lookup or accepted-name resolution; Map = supports taxonomy mapping tables or crosswalks between naming systems; Fetch = retrieves a tree from an external source; Multi = can return, handle, or compare more than one tree, such as posterior samples or alternative phylogenies. ✓ = full support; ~ = partial or context-dependent support; – = absent.

Tool	Role	Coverage	Tree	Prov.	Fuz.	Syn.	Map	Fetch	Multi
<code>taxize</code>	Taxonomic name resolution	Multiple online taxonomic databases	–	–	~	✓	–	–	–
<code>rgnparser^a</code>	Name parsing and verification	Global Names services and linked name sources	–	~	✓	✓	~	–	–
<code>taxadb</code>	Local taxonomic name resolution	Local copies of multiple taxonomic databases	–	–	–	✓	–	–	–
<code>Taxonstand</code>	Plant name standardisation	Plant names, originally against The Plant List	–	~	~	✓	–	–	–
<code>bdc</code>	Biodiversity data cleaning	Occurrence records and associated taxonomic names	–	~	~	✓	–	–	–
<code>rotl</code>	Taxonomy and tree retrieval	Open Tree of Life taxonomy and synthesis tree	~	–	–	✓	–	✓	–
<code>rtrees</code>	Megatree assembly	Taxon-specific megatrees, including birds, mammals, fishes, plants, and others	~	~	–	–	–	✓	✓
<code>clootl</code>	Avian tree retrieval	Bird trees mapped to eBird/Clements taxonomy	~	–	–	–	–	✓	~
<code>fishtree</code>	Fish tree retrieval	Ray-finned fish phylogenies from the Fish Tree of Life	~	–	–	–	–	✓	✓
<code>datelife</code>	Dated-tree synthesis	Published dated phylogenies across taxa	~	~	–	~	–	✓	✓
<code>geiger::treedata</code>	Data–tree alignment	User-supplied comparative dataset and tree	✓	–	–	–	–	–	–
<code>caper::comparative.data</code>	PCM data-object creation	User-supplied comparative dataset and tree	✓	–	–	–	–	–	–
<code>prepR4pcm</code>	Audited PCM preparation pipeline	User datasets and trees, with optional tree retrieval	✓	✓	✓	✓	✓	✓	✓

^aR wrapper for Global Names Architecture services (Mozzherin et al., 2017); `taxize` also uses the Global Names Verifier internally. This table compares broad workflow features rather than the quality, completeness, or taxonomic authority of each tool.

76 2 Package description

77 2.1 Overview and design

78 `prepR4pcm` was designed around three principles: conservative matching, transparent records,
79 and reproducible outputs. These principles apply to both main parts of the package: reconciling
80 species names among data objects and trees, and preparing trees for downstream PCMs.

81 First, `prepR4pcm` is conservative. It gives priority to high-confidence matches and does not
82 replace them with weaker matches later in the workflow. Exact and normalised matches are
83 attempted before synonym lookup and optional fuzzy matching. Names are not changed without
84 being recorded, and uncertain matches are left for user review.

85 Second, the package is transparent. Each reconciliation returns a `reconciliation` object
86 that records what happened to each name. This object stores the original and matched names,
87 the match type, a score, a short explanation, the settings used in the run, and any manual
88 overrides. Users can print the object for a summary or extract the full mapping table for inspection
89 and reporting.

90 Third, the package is reproducible. The reconciliation object stores the settings used in the
91 run, together with the package, database, and tree-source information where relevant. This
92 design helps users repeat a workflow, compare alternative reconciliation runs, and report how
93 the final data and tree objects were produced.

94 The package uses `ape` (Paradis and Schliep, 2019) for phylogenetic tree handling, `cli` for
95 user messages, `rlang` for tidy evaluation, and `tibble` for data frames. Synonym matching uses
96 `taxadb` (Norman et al., 2020) when this suggested dependency is available.

97 2.2 Matching cascade

98 The core of `prepR4pcm` is a “four-stage matching cascade” (Figure 1). The cascade starts with
99 the safest matches (exact name matching), then moves to those that require more interpretation.
100 At each stage, only names that remain unmatched advance to the next. Thus, an exact match is
101 never replaced by a weaker match later in the pipeline.

102 **Stage 1: Exact matching.** The package first compares names exactly as written. These
103 matches receive a score of 1.0 (ranging from 0 to 1).

104 **Stage 2: Normalised matching.** The package then standardises names before matching them
105 again. This step deals with common formatting differences, such as spaces versus underscores,
106 inconsistent capitalisation, repeated spaces, authority strings, Open Tree of Life suffixes, and
107 infra-specific rank abbreviations. These matches also receive a score of 1.0, but they are labelled
108 as “normalised” because the original strings differed.

109 **Stage 3: Synonym matching.** Names that remain unmatched can be checked against local
110 taxonomic databases through `taxadb`. This step is useful when the dataset uses one accepted
111 name and the tree uses a synonym or an older name. Where possible, the package records
112 which source was used. Synonym matches receive a lower score than exact or normalised

113 matches because they depend on the taxonomic authority used. Synonyms also get a score of
114 1.0.

115 **Stage 4: Fuzzy matching.** Finally, the package can suggest near matches for likely spelling
116 errors. Fuzzy matching compares the genus and species epithet using a component-based
117 string distance. It is useful for finding typographical errors, but it can also suggest wrong matches.
118 For this reason, fuzzy matching is optional, and uncertain matches are flagged for review rather
119 than silently accepted. Scores below 1.0 are used for fuzzy matches.

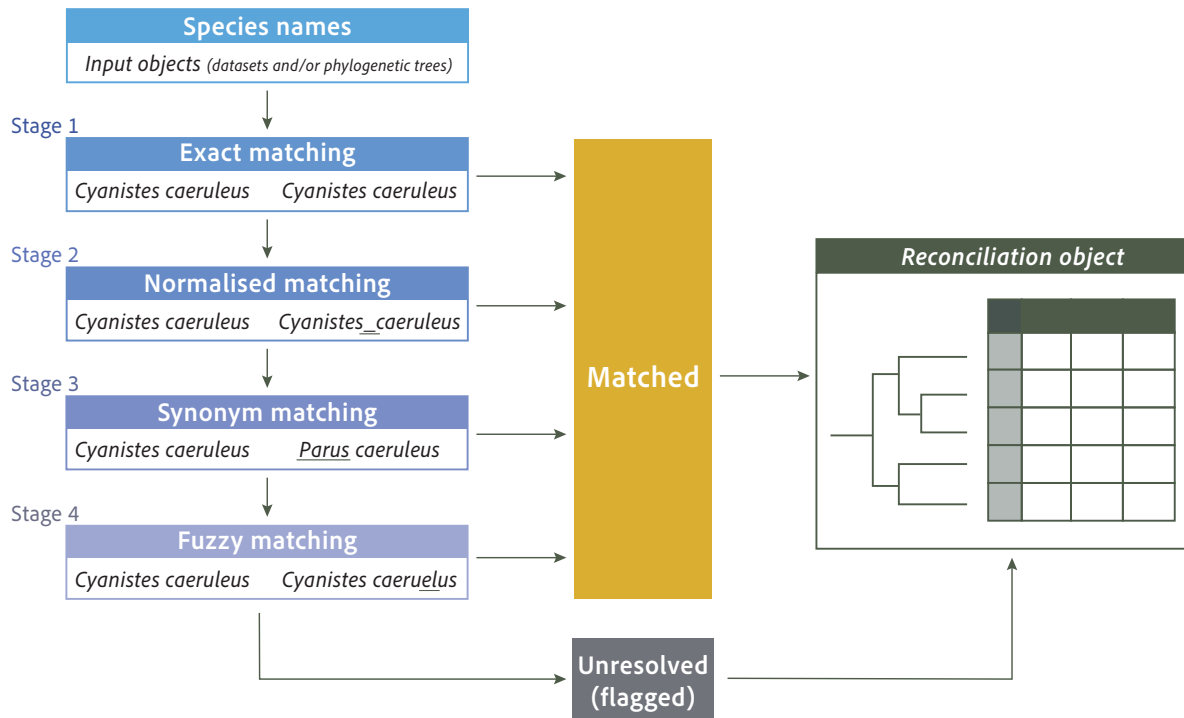


Figure 1. The four-stage matching cascade in prepR4pcm. Species names from input objects, such as datasets and phylogenetic trees, are passed through a conservative matching cascade. Exact matches are identified first, followed by matches after name normalisation, synonym-based matches, and optional fuzzy matches for likely spelling errors. The example species names illustrate the intended distinction among matching stages: identical spelling for exact matching (*Cyanistes caeruleus*), formatting differences for normalised matching (*Cyanistes_caeruleus*), synonymy or older nomenclature for synonym matching (*Parus caeruleus*), and a likely typographical error for fuzzy matching (*Cyanistes caeruelus*). Names matched at any stage, as well as names retained as unresolved or flagged cases, are stored in the reconciliation object so that the process can be inspected, reported, and reproduced.

120 2.3 The reconciliation object

121 The central output of `prepR4pcm` is a reconciliation object. This object records how names
122 were matched, changed, left unresolved, or corrected by the user. It is produced not only when a
123 dataset is reconciled with a tree, but also when users reconcile two datasets, two trees, several
124 datasets against one tree, or one dataset against several trees. Thus, the same object structure
125 is used across data–data, data–tree, tree–tree, multi-dataset, and multi-tree workflows.

126 A reconciliation object has four main components. The `$mapping` table records the names
127 from the input objects and the outcome for each name. The exact columns depend on the
128 workflow, but the table includes information such as original names, matched names, match
129 status, match type, score, source, and notes. The `$meta` list stores the settings used in the
130 run, together with timestamps and version information. The `$counts` table summarises how
131 many names were matched, changed, or left unresolved. The `$overrides` table records manual
132 changes made by the user.

133 The object separates reconciliation from application. Users can first inspect the name-
134 matching decisions, then decide whether to accept them, add corrections, or rerun the reconcilia-
135 tion with different settings. The object has `print()`, `summary()`, and `format()` methods. Printing
136 gives a short summary of match coverage. The function `reconcile_mapping()` returns the full
137 mapping table, which users can inspect, save, or include as supplementary material.

138 This object is the main difference between a simple alignment and an auditable reconciliation.
139 Standard PCM helper functions are useful for making data and trees ready for analysis, but they
140 are not designed to record all name-matching decisions. `prepR4pcm` keeps this record, so users
141 can see which names matched, which names changed, which names remained unresolved, and
142 which names were corrected by hand before the final data and tree objects are created.

143 2.4 Curation, augmentation, and multi-input workflows

144 Because the reconciliation object keeps the matching decisions separate from the final aligned
145 objects, users can curate the result before applying it. Automated matching should not be the
146 final word in every case. Taxonomy changes, spelling errors, and species splits and lumps
147 often require human judgement. `prepR4pcm` therefore provides tools for curation after the first
148 reconciliation.

149 Users can correct individual matches with `reconcile_override()` or apply many corrections
150 with `reconcile_override_batch()`. They can also use `reconcile_crosswalk()` to apply a
151 taxonomic crosswalk: a table that links names in one taxonomy to names in another. For
152 example, a crosswalk between BirdLife-based names and Jetz et al. / BirdTree tip labels can
153 show which BirdLife names correspond to tip labels in Jetz et al. bird phylogeny (Jetz et al.,
154 2012; Tobias et al., 2022).

155 Some studies need to deal with species that are present in the data but absent from the tree.
156 `reconcile_augment()` can add such species to a tree for exploratory analyses, for example, by
157 placing an unresolved species with a close relative. Internally, `prepR4pcm` can also delegate
158 grafting to the `V.PhyloMaker/U.PhyloMaker` suits when the required backbone and genus-family
159 information are available (Jin and Qian, 2019, 2022, 2023). This step should be treated as
160 a modelling assumption rather than an observed result. Users should therefore report which
161 species were added and repeat the analysis with and without augmented tips.

162 The package also supports studies with more than one dataset or more than one tree.
163 `reconcile_data()` reconciles two data frames, `reconcile_merge()` joins reconciled datasets,
164 `reconcile_multi()` reconciles several datasets against one tree, and `reconcile_to_trees()`
165 reconciles one dataset against several trees. These functions are useful when users want to
166 compare alternative trees or combine trait datasets from different sources.

167 2.5 Tree retrieval and dating

168 Many users begin with a trait dataset without a suitable phylogenetic tree. `pr_get_tree()` helps
169 in this situation by retrieving a candidate phylogeny from public sources. At present, the package
170 supports five backends: `rotl`, `rtrees`, `clootl`, `fishtree`, and `datelife`. These backends differ
171 in taxonomic coverage, tree sources, and whether they return a single tree or a set of trees.

172 Before a backend is called, `prepR4pcm` standardises the input species names with `pr_normalize_names()`.
173 This step helps avoid failures caused by simple formatting differences. When taxonomic name
174 resolution changes an input name, the package records the replacement, linking the user's
175 original name to the resolved name. The package also checks that the reported matched and
176 unmatched names remain tied to the user's input list.

177 Multi-tree support depends on the backend. For example, some sources can return posterior
178 samples of trees, whereas others return a single tree. The function `pr_date_tree()` can add
179 time calibration to a topology using `datelife` where suitable information is available. The
180 function `pr_cite_tree()` formats citations for tree sources in plain text, Markdown, or BibTeX.

181 2.6 Reporting, comparison, and export

182 A reconciliation is most useful when it can be inspected, compared, and reported. `reconcile_summary()`
183 prints a detailed summary to the console, and `reconcile_report()` writes a self-contained
184 HTML report. `reconcile_suggest()` proposes near matches for unresolved names, `reconcile_diff()`
185 compares two reconciliation runs, and `reconcile_splits_lumps()` flags cases where one tax-
186 onomy treats a name as one species, but another taxonomy treats it as several species, or *vice*
187 *versa*.

188 Once users are satisfied with the reconciliation, `reconcile_apply()` returns the two ob-
189 jects required for most PCMs: a data frame and a pruned tree with the same set of species.
190 `reconcile_export()` writes the aligned data, tree, and mapping table to disk. For tree com-
191 parisons, `pr_tree_compare()` summarises overlap in tip sets and differences in topology. For
192 phylogenetic meta-analysis, `pr_phylo_cor()` builds a phylogenetic correlation matrix that can
193 be used in packages such as `metafor` (Viechtbauer, 2010). Together, these functions move
194 users from raw species names to auditable, aligned comparative objects (Table 2).

Table 2. The 31 available functions in `prepR4pcm` package, organised by functionalities. Functions prefixed `reconcile_` enable the matching and alignment pipeline; functions prefixed `pr_` enable tree retrieval, dating, comparison, and supporting utilities.

Function	Family	Description
<code>reconcile_tree</code>	Reconciliation	Reconcile a data frame against a phylogenetic tree
<code>reconcile_data</code>	Reconciliation	Reconcile two data frames
<code>reconcile_trees</code>	Reconciliation	Reconcile two phylogenetic trees
<code>reconcile_to_trees</code>	Reconciliation	Reconcile one dataset against a list of trees
<code>reconcile_multi</code>	Reconciliation	Reconcile multiple datasets against one tree
<code>reconcile_apply</code>	Alignment	Produce matched data frame and pruned tree
<code>reconcile_merge</code>	Alignment	Join two reconciled datasets
<code>reconcile_override</code>	Curation	Manually correct an individual match
<code>reconcile_override_batch</code>	Curation	Apply batch corrections from CSV or data frame
<code>reconcile_crosswalk</code>	Curation	Convert a taxonomy crosswalk to overrides
<code>reconcile_augment</code>	Augmentation	Graft missing species onto a phylogenetic tree for exploratory analyses
<code>reconcile_mapping</code>	Inspection	Extract the full mapping tibble
<code>reconcile_summary</code>	Inspection	Generate a detailed summary report
<code>reconcile_plot</code>	Inspection	Plot match composition (bar or pie chart)
<code>reconcile_suggest</code>	Review	Suggest fuzzy candidates for unresolved names
<code>reconcile_review</code>	Review	Interactively accept/reject flagged matches
<code>reconcile_diff</code>	Review	Compare two reconciliation runs
<code>reconcile_splits_lumps</code>	Review	Detect taxonomic splits and lumps
<code>reconcile_report</code>	Reporting	Produce self-contained HTML report
<code>reconcile_export</code>	Reporting	Write aligned data, tree, and mapping to disk
<code>pr_get_tree</code>	Tree retrieval	Fetch a tree from external sources / backends: <code>rot1/rtrees/cloot1/fishtree/datelife</code> (manual selection or auto)
<code>pr_get_tree_status</code>	Tree retrieval	Probe which backends are installed and reachable
<code>pr_date_tree</code>	Dating	Add DateLife calibration to an existing topology
<code>pr_cite_tree</code>	Citations	Format citations for the sources of phylogenies used (select formats: text/Markdown/BibTeX)
<code>pr_tree_compare</code>	Comparison	Pairwise Jaccard, Robinson–Foulds, and bipartition-matched branch-length tree agreement
<code>pr_phylo_cor</code>	Correlation	Build the phylogenetic correlation matrix from a tree for <code>metafor/MCMCg1mm</code>
<code>pr_tree_cache_dir</code>	Cache	Get or set the cache directory
<code>pr_tree_cache_status</code>	Cache	List cache entries
<code>pr_tree_cache_clear</code>	Cache	Wipe the cache
<code>pr_normalize_names</code>	Utility	Standardise scientific-name formatting
<code>pr_extract_tips</code>	Utility	Extract tip labels from a phylogenetic tree

3 Worked examples

3.1 Example 1: Reconciling a dataset with an existing tree

We first demonstrate the core workflow using a subset of AVONET bird morphological data (Tobias et al., 2022) and a subset of the Jetz et al. bird phylogeny (Jetz et al., 2012), both bundled with the package. This example is deliberately simple, but it illustrates a common problem in comparative workflows. AVONET uses spaces in species names, whereas the tree uses underscores. Exact matching therefore fails, even when the biological species are the same.

Core workflow: data–tree reconciliation

```
library(prepareR4pcm)
library(caper)

# Load the example trait data and phylogenetic tree
data(avonet_subset)
data(tree_jet2)

# Step 1: Reconcile species names between the data and tree
# fuzzy = TRUE allows fuzzy suggestions, but resolve = "flag"
# keeps uncertain matches for review rather than accepting them silently
rec <- reconcile_tree(
  x = avonet_subset,
  tree = tree_jet2,
  x_species = "Species1",
  fuzzy = TRUE,
  resolve = "flag")

# Step 2: Inspect how names were matched
# This table shows whether names were matched exactly, after
# normalisation, by synonym lookup, by fuzzy matching, or left unresolved
mapping <- reconcile_mapping(rec)
table(mapping$match_type)

# Step 3: Apply the reconciliation
# This returns an aligned data frame and a pruned tree with the
# same species set
aligned <- reconcile_apply(
  rec,
  data = avonet_subset,
  tree = tree_jet2,
  species_col = "Species1",
  drop_unresolved = TRUE)

# Step 4: Use the aligned objects in a comparative analysis
cd <- comparative.data(
  phy = aligned$tree,
  data = aligned$data,
  names.col = "Species1",
  vcv = TRUE)
```

203

```
model <- pglis(log(Mass) ~ log(Wing.Length),  
data = cd)
```

204

205 The output of `table(mapping$match_type)` shows that all retained species were matched
206 through normalisation:

```
normalised  unresolved_x  
657         262
```

207

208 In this example, `reconcile_tree()` compares 919 data names with 657 tree tips. It matches
209 657 data names to tree tips, all through normalised matching, because the main difference is
210 formatting. The remaining 262 species names from the AVONET data object are unresolved
211 because they are present in the AVONET subset but not in the tree subset. Although fuzzy
212 matching is enabled in the code, no retained species are matched by fuzzy matching in this
213 example. After `reconcile_apply()` is run, both the data frame and the tree contain the same
214 657 species.

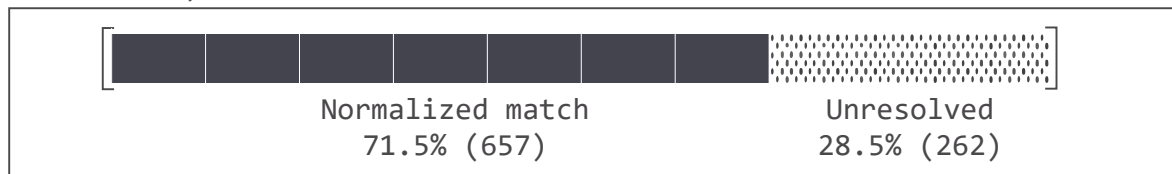
215 3.2 Inspecting the audit trail

216 This example illustrates why an audit is useful. A simple final count would tell us that 657
217 species were retained. The reconciliation object tells us more: it shows that the retained species
218 matched after a formatting change, that no exact matches occurred, and that 262 data rows
219 were removed because they had no matching tree tip in this example.

220 The same audit trail becomes more important in less tidy examples. The `reconciliation`
221 object can also record synonym matches, fuzzy suggestions, manual overrides, unresolved
222 names, taxonomic sources, database versions, and user-supplied corrections. Users can
223 inspect this information with `reconcile_mapping()`, summarise it with `reconcile_summary()`,
224 or generate an HTML report with `reconcile_report()`.

225 Figure 2 shows two outputs from this workflow: the match composition and a schematic pre-
226 view of the HTML report. Together, these outputs show both the numerical result of reconciliation
227 and the record of decisions that led to the final data–tree pair.

a. Match composition



b. Reconciliation report preview

Reconciliation Report

Generated by prepR4pcm

Type	data vs tree
Source x	avonet_subset
Source y	phylo (657 tips)
Authority	col
Reconciliation date	year-month-day hh:mm:ss

Match summary

Match type	Count	% of x
Exact	0	0.0%
Normalized	657	71.5%
Synonym	0	0.0%
Fuzzy	0	0.0%
Manual	0	0.0%
Flagged	0	0.0%
Unresolved (x only)	262	28.5%
Unresolved (y only)	0	
Total matched	657	71.5%

Normalized matches (657)

Formatting differences resolved automatically.

Name (x)	Name (y)	Notes
Acanthiza apicalis	Acanthiza_apicalis	'Acanthiza apicalis' normalised to 'Acanthiza apicalis'
Acanthiza chrysorrhoa	Acanthiza_chrysorrhoa	'Acanthiza chrysorrhoa' normalised to 'Acanthiza chrysorrhoa'
		⋮

Figure 2. Example outputs from a prepR4pcm reconciliation workflow. (a) Match composition for the worked example using an AVONET data subset containing 919 species names and a subset of the Jetz et al. (2012) phylogeny containing 657 tip labels. In this example, 657 data names were matched to the 657 tree tips through normalised matching, mainly because AVONET uses spaces in binomials, whereas the tree uses underscores. The remaining 262 species names in the AVONET data subset were unresolved data-only names, meaning that they were present in the data subset but absent from the tree subset. (b) Schematic preview of the HTML reconciliation report generated by prepR4pcm. The report summarises run metadata, source objects, taxonomic authority, match-type counts, unresolved names, and examples of names resolved through normalised matching.

228 3.3 Example 2: Retrieving trees before reconciliation

229 The first example assumes that the user already has a phylogenetic tree. In many projects,
230 however, users begin with a trait dataset but without a suitable tree. In this case, `prepR4pcm` can
231 retrieve a candidate tree before the reconciliation step. The code below retrieves a posterior
232 sample of fish trees and formats the relevant citations.

```
# Retrieve 50 trees for three fish species from the fishtree backend
trees <- pr_get_tree(
  c("Salmo salar", "Esox lucius",
    "Oncorhynchus mykiss"),
  source = "fishtree",
  n_tree = 50)

# Check that the result is a posterior sample of trees
class(trees$tree) # "multiPhylo"
length(trees$tree) # 50

# Format citations for the tree source
cat(pr_cite_tree(trees, format = "markdown"))
```

233
234 The retrieved tree object can then be used in the same way as a user-supplied tree.
235 For example, users can pass the retrieved tree to `reconcile_tree()`, inspect the resulting
236 reconciliation object, and then use `reconcile_apply()` to create aligned data and tree ob-
237 jects.

238 The package code, documentation, and worked vignettes are available in the anonymised
239 repository at <https://anonymous.4open.science/r/prepR4pcm-2A33/>. The vignettes demon-
240 strate data–tree reconciliation, tree retrieval, multi-tree workflows, tree comparison, bird and
241 mammal trait workflows, and phylogenetic meta-analysis.

242 4 Bundled data

243 Importantly, `prepR4pcm` includes small example datasets for vignettes, worked examples, and
244 tests. These datasets are not intended to replace the full published datasets. Rather, they
245 provide realistic examples in which species names do not always match perfectly.

246 The bird data examples include an AVONET data subset (Tobias et al., 2022), a NestTrait
247 subset (Chia et al., 2023), a plumage lightness subset (Delhey et al., 2019), a Jetz et al.
248 tree subset (Jetz et al., 2012), a subset of the complete avian phylogeny mapped to the
249 eBird/Clements Checklist v2025 (McTavish et al., 2025; Clements et al., 2025), and a BirdLife to
250 BirdTree taxonomy crosswalk distributed with AVONET (Tobias et al., 2022; Jetz et al., 2012).
251 These objects support the worked example above and the bird-trait vignette.

252 The mammal examples include subsets from the Amniote life-history database (Myhrvold
253 et al., 2015), PanTHERIA (Jones et al., 2009), and TetrapodTraits (Moura et al., 2024), along
254 with a mammal tree used for demonstration (Upham et al., 2019). Users who rely upon these
255 example datasets in their published work should cite the original data providers and, where
256 needed, replace the example objects with the full original datasets.

257 The examples are useful because they are not perfectly clean. Some species names differ
258 in formatting, some datasets follow different taxonomies, and some species are present in one
259 source but not another. They therefore show the kinds of problems that users often face in real
260 PCM workflows.

261 5 Discussion

262 Preparing data for PCMs often begins with an important yet underreported step: matching
263 species names between a dataset and a phylogenetic tree. Researchers usually report the
264 final number of species analysed, but they may not report how many species were lost during
265 matching, which names failed to match, or how ambiguous names were resolved. These missing
266 details make it difficult to reproduce the analysis or to assess whether the final species set could
267 affect the results.

268 `prepR4pcm` makes this step visible. The package reconciles names, records the decisions,
269 and returns aligned objects for downstream analysis. The key outputs are not only a pruned
270 tree and a matched data frame but also a `reconciliation` object that can be inspected, saved,
271 shared, and reported. This object records exact matches, normalised matches, synonym
272 matches, fuzzy suggestions, unresolved names, and manual corrections. It also allows users
273 to report how many names were matched at each stage, which names remained unresolved,
274 which decisions were corrected manually, and which taxonomic authority or database version
275 was used. The full mapping table can be archived as supplementary material so that readers
276 can trace how the final species set was obtained.

277 The same idea applies to tree retrieval. When users retrieve trees using `prepR4pcm`, the
278 package records the tree source and the name substitutions made along the way. This record
279 helps users write a clearer methods section and provides a more complete audit trail for the
280 analysis. In this sense, `prepR4pcm` treats name matching, tree preparation, and final data
281 alignment as parts of one workflow rather than as separate preprocessing steps.

282 Several limitations remain. Fuzzy matching can produce incorrect matches, so users should
283 review fuzzy suggestions before accepting them. Tree augmentation also requires care. Here,
284 tree augmentation means adding species that are present in the dataset but absent from the
285 phylogenetic tree. Such additions may be based on genus-level placement, a close relative, a
286 taxonomic crosswalk, or information from another more detailed published tree. These choices
287 are assumptions, not observed phylogenetic relationships. Users should therefore report which
288 species were added, how they were added, and what source or rule was used. They should also
289 inspect the resulting tree visually, because plotting the augmented tree is often a useful sanity
290 check. Where possible, analyses should be repeated with and without augmented species, and
291 sometimes with alternative augmentation strategies.

292 Synonym matching also depends on the coverage and version of the taxonomic database
293 used. Users should always report the authority, database version, and any manual changes
294 made during reconciliation. Future versions will expose the Global Names Architecture ([Moz-
295 zherin et al., 2017](#)) as alternative backends, including `rgnparser` for the parsing step in
296 `pr_normalize_names()` and the Global Names Verifier for synonym resolution against a larger
297 set of sources, alongside the existing `taxadb`-based path.

298 By making reconciliation and related decisions explicit, `prepR4pcm` helps users move from

299 raw species names to analysis-ready comparative objects. More importantly, it helps users show
300 how they got there. This audit trail should make PCM workflows easier to verify, repeat, and
301 report.

302 **Author contributions**

303 **Author Statement**

304 Conceptualization: SN, BJ, SPH. Data curation: SN, AM, SO. Formal analysis: SN, AM, SO.
305 Funding Acquisition: SN, ML. Investigation: All. Methodology: SN. Project administration: SN,
306 ESAS, SPH. Supervision: SN, ESAS, ML. Validation: All. Visualization: AM. Writing - Original
307 Draft: SN. Writing - Review & Editing: All.

308 **Acknowledgements**

309 We thank the developers of the upstream R packages that `prepR4pcm` wraps or builds on.

310 **Statements**

311 **Conflict of interest:** The authors declare no conflict of interest.

312 **Peer review:** The peer review history for this article is available at [URL].

313 **Data availability:** All data and code are available at [https://github.com/itchyshin/](https://github.com/itchyshin/prepR4pcm)
314 [prepR4pcm](https://github.com/itchyshin/prepR4pcm). A Zenodo DOI will be added at acceptance.

315 **Use of generative AI**

316 Generative AI tools, including Claude Code and Codex, were extensively used during software
317 development to assist with code drafting, refactoring, documentation, and test development.
318 These tools were not used as independent sources of scientific evidence. AI-assisted code was
319 checked by the authors through unit tests, worked examples, vignettes, manual inspection, and
320 comparisons with expected outputs. The authors take full responsibility for the accuracy of the
321 package, analyses, and manuscript.

322 **ORCID**

323 Shinichi Nakagawa [0000-0002-7765-5182](https://orcid.org/0000-0002-7765-5182); Malgorzata Lagisz [0000-0002-3993-6127](https://orcid.org/0000-0002-3993-6127); Santiago
324 Ortega [0000-0002-3518-276X](https://orcid.org/0000-0002-3518-276X); Bhayvay Jain [0009-0000-1021-4326](https://orcid.org/0009-0000-1021-4326) Jimuel Jr Celeste [0000-
325 0003-0224-4124](https://orcid.org/0000-0003-0224-4124); Ayumi Mizuno [0000-0003-0822-5637](https://orcid.org/0000-0003-0822-5637); Eduardo S. A. Santos [0000-0002-0434-
326 3655](https://orcid.org/0000-0002-0434-3655) Sergio Poo Hernandez [0000-0003-0155-645X](https://orcid.org/0000-0003-0155-645X)

References

- 327
- 328 Luis Cayuela, Íñigo Granzow-de la Cerda, Fabio S. Albuquerque, and Duncan J. Golicher.
329 Taxonstand: An R package for species names standardisation in vegetation databases.
330 *Methods in Ecology and Evolution*, 3:1078–1083, 2012. doi: 10.1111/j.2041-210X.2012.00232.x.
- 331 Scott A. Chamberlain and Eduard Szöcs. taxize: taxonomic search and retrieval in R.
332 *F1000Research*, 2:191, 2013. doi: 10.12688/f1000research.2-191.v2.
- 333 Jonathan Chang, Daniel L Rabosky, Stephen A Smith, and Michael E Alfaro. An r package and
334 online resource for macroevolutionary studies using the ray-finned fish tree of life. *Methods in*
335 *Ecology and Evolution*, 10(7):1118–1124, 2019.
- 336 Stephanie Yuan Chia, Yi-Ting Fang, Yi-Ting Su, Pei-Yu Tsai, Chia Hsieh, Shu-Han Tsao, Jia-
337 Yang Juang, Chih-Ming Hung, and Mao-Ning Tuanmu. A global database of bird nest traits.
338 *Scientific Data*, 10:923, 2023. doi: 10.1038/s41597-023-02837-1.
- 339 J. F. Clements, P. C. Rasmussen, T. S. Schulenberg, M. J. Iliff, J. A. Gerbracht, D. Lep-
340 age, A. Spencer, S. M. Billerman, B. L. Sullivan, M. Smith, and C. L. Wood. The
341 ebird/clements checklist of birds of the world: v2025. [https://www.birds.cornell.edu/
342 clementschecklist/download/](https://www.birds.cornell.edu/clementschecklist/download/), 2025. Downloaded [insert date].
- 343 Kaspar Delhey, James Dale, Mihai Valcu, and Bart Kempenaers. Reconciling ecogeographical
344 rules: rainfall and temperature predict global colour variation in the largest bird radiation.
345 *Ecology Letters*, 22(4):726–736, 2019. doi: 10.1111/ele.13233.
- 346 Joseph Felsenstein. Phylogenies and the comparative method. *American Naturalist*, 125:1–15,
347 1985. doi: 10.1086/284325.
- 348 Robert P. Freckleton, Paul H. Harvey, and Mark Pagel. Phylogenetic analysis and comparative
349 data: a test and review of evidence. *American Naturalist*, 160:712–726, 2002. doi: 10.1086/
350 343873.
- 351 Alan Grafen. The phylogenetic regression. *Philosophical Transactions of the Royal Society of*
352 *London. Series B*, 326:119–157, 1989. doi: 10.1098/rstb.1989.0106.
- 353 Jarrod D. Hadfield and Shinichi Nakagawa. General quantitative genetic methods for compara-
354 tive biology: phylogenies, taxonomies and multi-trait models for continuous and categorical
355 characters. *Journal of Evolutionary Biology*, 23:494–508, 2010. doi: 10.1111/j.1420-9101.2009.
356 01915.x.
- 357 Luke J. Harmon, Jason T. Weir, Chad D. Brock, Richard E. Glor, and Wendell Challenger.
358 GEIGER: investigating evolutionary radiations. *Bioinformatics*, 24:129–131, 2008. doi: 10.
359 1093/bioinformatics/btm538.
- 360 Walter Jetz, G. H. Thomas, J. B. Joy, K. Hartmann, and A. O. Mooers. The global diversity of
361 birds in space and time. *Nature*, 491:444–448, 2012. doi: 10.1038/nature11631.
- 362 Yi Jin and Hong Qian. V.PhyloMaker: an R package that can generate very large phylogenies
363 for vascular plants. *Ecography*, 42(8):1353–1359, 2019. doi: 10.1111/ecog.04434.

-
- 364 Yi Jin and Hong Qian. V.PhyloMaker2: An updated and enlarged R package that can generate
365 very large phylogenies for vascular plants. *Plant Diversity*, 44(4):335–339, 2022. doi: 10.1016/
366 j.pld.2022.05.005.
- 367 Yi Jin and Hong Qian. U.PhyloMaker: An R package that can generate large phylogenetic trees
368 for plants and animals. *Plant Diversity*, 45(3):347–352, 2023. doi: 10.1016/j.pld.2022.12.007.
- 369 Kate E. Jones, Jon Bielby, Marcel Cardillo, Susanne A. Fritz, Justin O’Dell, C. David L.
370 Orme, Kamran Safi, Wes Sechrest, Elizabeth H. Boakes, Chris Carbone, Christina Con-
371 nolly, Michael J. Cutts, Janine K. Foster, Richard Grenyer, Michael Habib, Christopher A.
372 Plaster, Samantha A. Price, Elizabeth A. Rigby, Janna Rist, Amber Teacher, Olaf R. P.
373 Bininda-Emonds, John L. Gittleman, Georgina M. Mace, and Andy Purvis. PanTHERIA: a
374 species-level database of life history, ecology, and geography of extant and recently extinct
375 mammals. *Ecology*, 90(9):2648, 2009. doi: 10.1890/08-1494.1.
- 376 Daijiang Li. rtrees: an r package to assemble phylogenetic trees from megatrees. *Ecogra-
377 phy*, 2023(7):e06643, 2023. doi: 10.1111/ecog.06643. URL [https://github.com/daijiang/
378 rtrees](https://github.com/daijiang/rtrees).
- 379 Emily Jane McTavish, Jeff A. Gerbracht, Mark T. Holder, Marshall J. Iloff, Denis Lepage, Pam
380 Rasmussen, Benjamin Redelings, Luna L. Sanchez-Reyes, and Eliot T. Miller. A complete
381 and dynamic tree of birds. *Proceedings of the National Academy of Sciences*, 122(18):
382 e2409658122, 2025. doi: 10.1073/pnas.2409658122.
- 383 François Michonneau, Joseph W. Brown, and David J. Winter. rotl: an R package to interact
384 with the Open Tree of Life data. *Methods in Ecology and Evolution*, 7:1476–1481, 2016. doi:
385 10.1111/2041-210X.12593.
- 386 Eliot Miller, Emily Jane McTavish, and Luna L. Sanchez-Reyes. cloodl: Fetch and explore the
387 cornell lab of ornithology open tree of life avian phylogeny, 2025. URL [https://github.com/
388 eliotmiller/cloodl](https://github.com/eliotmiller/cloodl).
- 389 Ayumi Mizuno, Szymon M. Drobniak, Coralie Williams, Malgorzata Lagisz, Shinichi Nakagawa,
390 Xiang-Yi Li Richter, and Carolin Kosiol. Promoting the use of phylogenetic multinomial
391 generalised mixed-effects model to understand the evolution of discrete traits. *Journal of
392 Evolutionary Biology*, 38:1699–1715, 2025. doi: 10.1093/jeb/voaf116.
- 393 Mario R. Moura, Karoline Ceron, Jhonny J. M. Guedes, Rosana Chen-Zhao, Yanina V. Sica,
394 Julie Hart, Wendy Dorman, Julia M. Portmann, Pamela González-del Pliego, Ajay Ranipeta,
395 Alessandro Catenazzi, Fernanda P. Werneck, Luís Felipe Toledo, Nathan S. Upham, João
396 F. R. Tonini, Timothy J. Colston, Robert Guralnick, Rauri C. K. Bowie, R. Alexander Pyron, and
397 Walter Jetz. TetrapodTraits: A comprehensive trait database with imputed values for global
398 tetrapods. *PLOS Biology*, 22(7):e3002658, 2024. doi: 10.1371/journal.pbio.3002658.
- 399 Dmitry Y. Mozzherin, Alexander A. Myltsev, and David J. Patterson. “gnparser”: a powerful
400 parser for scientific names based on parsing expression grammar. *BMC Bioinformatics*, 18:
401 279, 2017. doi: 10.1186/s12859-017-1663-3.

402 Nathan P. Myhrvold, Elita Baldrige, Benjamin Chan, Dhileep Sivam, Daniel L. Freeman, and
403 S. K. Morgan Ernest. An amniote life-history database to perform comparative analyses with
404 birds, mammals, and reptiles. *Ecology*, 96(11):3109, 2015. doi: 10.1890/15-0846R.1.

405 Shinichi Nakagawa and Eduardo SA Santos. Methodological issues and advances in biological
406 meta-analysis. *Evolutionary Ecology*, 26(5):1253–1274, 2012.

407 Kari E. A. Norman, Scott Chamberlain, and Carl Boettiger. taxadb: a high-performance local
408 taxonomic database interface. *Methods in Ecology and Evolution*, 11:1153–1159, 2020. doi:
409 10.1111/2041-210X.13440.

410 David Orme, Rob Freckleton, Gavin Thomas, Thomas Petzoldt, Susanne Fritz, Nick Isaac, and
411 Will Pearse. *caper: Comparative Analyses of Phylogenetics and Evolution in R*, 2025. R
412 package version 1.0.4.

413 Emmanuel Paradis and Klaus Schliep. ape 5.0: an environment for modern phylogenetics and
414 evolutionary analyses in R. *Bioinformatics*, 35:526–528, 2019. doi: 10.1093/bioinformatics/
415 btty633.

416 Matthew W. Pennell, Jonathan M. Eastman, Graham J. Slater, Joseph W. Brown, Josef C. Uyeda,
417 Richard G. FitzJohn, Michael E. Alfaro, and Luke J. Harmon. geiger v2.0: an expanded suite
418 of methods for fitting macroevolutionary models to phylogenetic trees. *Bioinformatics*, 30:
419 2114–2116, 2014. doi: 10.1093/bioinformatics/btu181.

420 R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for
421 Statistical Computing, Vienna, Austria, 2025. URL <https://www.R-project.org/>.

422 Daniel L. Rabosky, Jonathan Chang, Pascal O. Title, Peter F. Cowman, Lauren Sallan, Matt
423 Friedman, Kristin Kaschner, Cristina Garilao, Thomas J. Near, Marta Coll, and Michael E.
424 Alfaro. An inverse latitudinal gradient in speciation rate for marine fishes. *Nature*, 559:392–395,
425 2018. doi: 10.1038/s41586-018-0273-1.

426 Bruno R. Ribeiro, Santiago José Elías Velazco, Karlo Guidoni-Martins, Geiziane Tessarolo,
427 Lucas Jardim, Steven P. Bachman, and Rafael Loyola. bdc: a toolkit for standardizing,
428 integrating and cleaning biodiversity data. *Methods in Ecology and Evolution*, 13:1421–1428,
429 2022. doi: 10.1111/2041-210X.13868.

430 Luna L. Sanchez Reyes, Brian C. O’Meara, Joseph W. Brown, and Emily Jane McTavish.
431 DateLife: Leveraging databases and analytical tools to reveal the dated tree of life. *Systematic*
432 *Biology*, 73(2):470–485, 2024. doi: 10.1093/sysbio/syae015.

433 Joseph A. Tobias, Catherine Sheard, Alex L. Pigot, et al. AVONET: morphological, ecological
434 and geographical data for all birds. *Ecology Letters*, 25:581–597, 2022. doi: 10.1111/ele.13898.

435 Nathan S. Upham, Jacob A. Esselstyn, and Walter Jetz. Inferring the mammal tree: Species-level
436 sets of phylogenies for questions in ecology, evolution, and conservation. *PLOS Biology*, 17
437 (12):e3000494, 2019. doi: 10.1371/journal.pbio.3000494.

438 Wolfgang Viechtbauer. Conducting meta-analyses in r with the metafor package. *Journal of*
439 *statistical software*, 36:1–48, 2010.