

## Original Articles

## Resolving taxonomic ambiguities: Effects on rarity, projected richness, and indices in macroinvertebrate datasets

Christy S. Meredith<sup>a,\*</sup>, Anett S. Trebitz<sup>b</sup>, Joel C. Hoffman<sup>b</sup><sup>a</sup> National Research Council, U. S. Environmental Protection Agency, Office of Research and Development, Mid-Continent Ecology Division, 6201 Congdon Blvd, Duluth, MN 55804, USA<sup>b</sup> U. S. Environmental Protection Agency, Office of Research and Development, Mid-Continent Ecology Division, 6201 Congdon Blvd, Duluth, MN 55804, USA

## ARTICLE INFO

## Keywords:

Resolution  
Biodiversity  
Non-parametric  
Ambiguous  
Benthos  
Macro-invertebrate

## ABSTRACT

Biodiversity information is an important basis for ecological research and environmental assessment, and can be impacted by choices made in the manipulation and analysis of taxonomic data. Such choices include methods for resolving multiple redundant levels of taxonomic resolution, as typically arise with morphological identification of damaged or immature aquatic macro-invertebrates. In particular, the effects of these processing choices on number of rare taxa are poorly understood yet potentially significant to the estimation of projected taxa richness and related evaluations such as biodiversity conservation value and survey sufficiency. Using aquatic macro-invertebrate data collected for two nearshore areas of Lake Superior, we determined how multiple methods of resolving taxonomic redundancies influence two commonly-used estimates of projected richness, Chao1 and Chao2, which hinge on the ratio of taxa that are singletons to doubletons (i.e., just one versus two individuals found) or uniques versus duplicates (i.e., just one versus two occurrences). We also determined how choice of ambiguous taxa method, including some modified specifically to retain rare taxa and others taken from the literature, influenced effort to reach 95% of projected richness, site-level richness and abundance, and representative invertebrate IBI scores. We found that Chao1 was more sensitive to method choice than Chao2, because singleton and doubleton status was more frequently affected when taxa were deleted, merged, or re-assigned in the process of resolving taxonomic redundancies than was unique and duplicate status. Methods that eliminated redundant taxa at the site scale but not the study-area scale tended to overinflate study area and projected richness, and resulted in a significant loss of abundance. The method that aggregated or deleted redundant taxa depending on abundance resulted in a decrease in site and study area richness, abundance, and an underestimation of projected richness. Methods which re-assigned parents to common children retained a majority of richness and abundance information and a more realistic estimate of projected taxa richness; however, the identity of poorly-identified parents was imputed. All methods resulted in little effect to typical IBI scores. Overall, no one method is fully capable of removing spurious richness at the study-area scale while preserving all taxa occurrence, abundance and rarity patterns. Therefore, the most appropriate method for making comparisons among sites may be different than the most appropriate method for comparing among surveys or among study areas, or if a goal is to estimate projected taxa richness or retain rare taxa information.

## 1. Introduction

Measures of biological composition are fundamental to the characterization and assessment of aquatic ecosystems. Richness overall and within various taxonomic groups is frequently incorporated into indices of biotic integrity (IBIs – Novotny et al., 2005; Davies and Jackson 2006) and used to prioritize areas for conservation (Weigel et al., 2002; Abell et al., 2008). While it is recognized that sampling choices (e.g., equipment or spatial scale; Lammert and Allan, 1999; Blocksom and

Flotemersch, 2005; Silva et al., 2016) and taxonomic processing choices (e.g., level of subsampling and identification; Carter and Resh, 2001; King and Richardson, 2002) can influence apparent biological composition, another aspect of data management which is often overlooked can also have a large influence on the diversity and abundance of taxa—the resolving of ambiguous taxa. Redundant or “ambiguous” taxa occur when the same taxon is identified to different levels of resolution, such as when complete taxonomic keys only exist for certain life stages, damaged specimens are present, or even when limited resources result

\* Corresponding author.

E-mail address: [cmeredith@jrn.lter](mailto:cmeredith@jrn.lter) (C.S. Meredith).<https://doi.org/10.1016/j.ecolind.2018.10.047>

Received 19 January 2018; Received in revised form 19 October 2018; Accepted 22 October 2018

Available online 02 November 2018

1470-160X/ © 2018 Elsevier Ltd. All rights reserved.

in only a subset of a sample being fully identified to species. The presence of these ambiguous taxa may result in inflated estimates of richness or distorted patterns of rarity across sites (Cuffney et al., 2007). To resolve these redundancies, taxa can be merged, deleted, or re-assigned. However, alternative rule-sets for doing so can produce different patterns of organism identities, with the side effect of either reducing abundance or redistributing abundance information on which many IBI and other metrics depend.

Redundant taxa situations are frequently encountered for aquatic macro-invertebrate data. For instance, some specimens may be identified to the genus level (e.g., *Hexagenia*) while others are identified to the species level (e.g., *Hexagenia limbata*). If ambiguous taxa are not resolved, then the coarser-resolution taxon (i.e., the parent, *Hexagenia*) is counted as a unique taxon when it may actually be the same species as one of the finer-level taxon (i.e., the children, *H. limbata*). Typical methods for resolving ambiguous taxa include merging higher-resolution “children” with lower-resolution “parents,” deleting lower-resolution parents, re-assigning or dividing lower-resolution parents among children, or a combination of these (Cuffney et al., 2007). Removing parents and retaining children is generally considered suitable for preserving richness at the site-scale, but obviously discards abundance information. Conversely, methods that divide parents among children or assign parents to the most abundant children to remove spurious richness do not result in a loss of overall abundance, but may alter patterns of richness and the distribution of that abundance (Cuffney et al., 2007). The impact of resolving ambiguous taxa depends on whether redundancies are eliminated only at the site level or also at the dataset level. Accomplishing the latter may require assigning the identity of a poorly-resolved taxa at a site (i.e. the “parent”) to a more-resolved taxa from elsewhere in the study area (i.e. the “child”). This assignment may or may not be acceptable depending on the goals of the study.

Because they are of low abundance, rare taxa may be particularly affected by efforts to resolve ambiguous taxa because they are more likely to be deleted or merged with other taxonomic groups. Accordingly, choices for resolving ambiguous taxa can impact non-parametric estimates of projected species richness (aka, asymptotic richness) that rely on the number of rare taxa encountered with a given sampling effort to estimate the total species pool of an area. Projected richness estimates are commonly used in prioritizing biodiversity conservation efforts (Desmet and Cowling, 2004; Sambuichi and Haridasan, 2007), and to assess sampling effort sufficiency (Schreiber and Brauns, 2010; Ram et al., 2014). We recognize that these non-parametric estimates of projected taxa richness have been criticized because their precision and accuracy depend on both sampling effort (D’Alessandro and Fattorini, 2002) and distribution patterns, especially rarity and spatial aggregation (Reese et al., 2014; Walther and Moore, 2005; Chiarucci et al., 2003). However, they remain widely used in survey assessment, because they provide reasonable estimates of projected richness even for under-sampled biological communities, and are parsimonious and analytically tractable solutions to a problem that must otherwise be addressed with randomization and curve fitting (Chao et al., 2009).

Our goal is to investigate how alternative methods of resolving ambiguous taxa influence the number of rare taxa and projected taxa richness. Specifically, we consider methods that might avoid degrading rarity patterns by retaining rare taxa information when coarse-resolution taxa are re-assigned to finer-resolution taxa, as well as several methods recommended by Cuffney et al. (2007). Cuffney et al. (2007) laid substantial groundwork to understand the potential impact of choice of methods for resolving ambiguous taxa on estimates of richness, abundance, and IBIs. Ideally, the same method of resolving ambiguous taxa would be used for multiple data analysis goals. Thus, we also evaluate how the various methods compare with respect to their effect on taxa richness and abundance, IBIs, as well as projected taxa richness. Finally, given the complicated nature of resolving ambiguous

taxa in macroinvertebrate datasets, we make our computer programs to resolve ambiguous taxa available as an R package: (<https://github.com/christystarr/ambiguoustaxa>) archived at Zenodo (<https://doi.org/10.5281/zenodo.1472508>). We use macro-invertebrate data collected within two nearshore areas of Lake Superior in conducting our analyses.

## 2. Methods

### 2.1. Study design, sampling methods, taxonomic procedures

We performed our analysis separately for two separate study areas with potentially quite different benthic invertebrate composition: The St. Louis River Estuary (SLRE) and nearshore areas of Isle Royale. The SLRE, occurring at the confluence of the St. Louis River and Lake Superior, is characterized by mesotrophic waters, primarily soft substrates, and many potential aquatic invasive species (AIS) transport vectors given its adjacency to the urban areas of Duluth, MN and Superior, WI and the largest commercial seaport on the Great Lakes. Isle Royale is a remote island located in the northwest portion of Lake Superior, and is characterized by oligotrophic waters, a mixture of hard and soft substrates, and relatively few AIS transport vectors (i.e., mostly recreational boat traffic). Both study areas were sampled as part of a research effort aimed at developing protocols for multi-taxa AIS early detection monitoring.

Sampling at Isle Royale was conducted using a stratified random survey design that distributed sites across three depth strata (0–2 m, 2–5 m, and 5–20 m) in each of five major embayments. Sampling at SLRE also followed a stratified random survey design that covered the urbanized portion of the system with two depth strata (0–2 m, > 2 m). A suite of equipment was used in both systems to span different substrate types, positions in the water column, and potential organism attractants (e.g., crevices, light). At Isle Royale, 252 sites were sampled in summer of 2012 using a combination of colonization plates, light traps, PONAR grabs, rock bags, sweep nets, and a vacuum pump. At SLRE, 94 sites were sampled in summer of 2012 and 2013 using benthic sleds, colonization plates, light traps, PONAR grabs, rock bags, and sweep nets.

All samples were elutriated or scraped, rinsed through a 500- $\mu$ m mesh net, and preserved in buffered formalin or ethanol. Laboratory processing consisted of picking invertebrates from the samples and identifying them via morphological features to the best taxonomic resolution possible given available taxonomic keys, specimen life-stage, and condition. Ten % of all samples had identification verified by a second staff taxonomist as part of quality assurance procedures; questionable specimens were sent to outside experts for verification. Given our interest in detecting potentially new and rare taxa and the known problems with incomplete counting and identification methods (e.g., King and Richardson, 2002), the majority of samples had all specimens counted and identified. However, for samples containing many chironomids or oligochaetes, identification and enumeration of these taxa was done via subsampling because the best-resolution morphological identification of chironomids and oligochaetes requires making slide mounts and examining them under a compound microscope. For subsampling, specimens were sorted into coarse taxonomic groups, then 20% from each group (or all specimens if group  $N \leq 10$ ) were mounted and identified. We use the subsampled counts directly to determine invertebrate composition (e.g., if 30 of 150 Oligochaeta were mounted and identified as 1 *Limnodrilus hoffmeisteri*, 2 *L. maumeensis*, and 27 “*Tubificinae* without setae”, the data as analyzed are 120 Oligochaeta parents and the actual counts for the three child taxa).

### 2.2. Estimating projected taxa richness

A commonly used function to estimate projected (asymptotic) richness is Chao1 (Chao et al., 2009), an abundance-based method which hinges on the ratio of singletons (one individual found) to

**Table 1**

Description of the various methods use to resolve ambiguous taxa. Suffixes “-S”, “-G”, and “-SG” indicating methods that respectively consider site level information only (S), study-level information only (G), or first site-level and then study-level information (SG). All methods are applied by iterating from lowest to highest taxonomic resolution.

Method	Description
RPMC	Remove Parents Merge Children: If parent outnumbers children in the study area, change the identity of its children to that of parent. If children outnumber parents in the study area, delete parents.
RPKC-S	Remove Parents Keep Children-S: At any site having children, delete entries representing parent. Leave parent taxa at sites where no children are identified.
RPKC-SG	Remove Parents Keep Children-SG: At any site having children, delete entries representing parent. Then assign remaining parents to their most numerous child in the study area.
DPAC-SG	Distribute Parents Among Children-SG: At any site having children, distribute the count of parent proportionally among children. Then assign any remaining parents to their most numerous (Chao1) or widespread (Chao2) child in the study area.
APTC-S	Assign Parents To Children-S: At any site having children, assign parent to the most numerous child at the site. Leave parent taxa at sites where no children are present.
APTC-G	Assign Parents To Children-G: Assign parent to most numerous (Chao1) or most widespread (Chao2) child in study area, without consideration for children at the same site.
APTC-SG	Assign Parents To Children-SG: At any site having children, assign parent to the most numerous child present at the site. Then assign remaining parents to the most numerous (Chao1) or widespread (Chao2) child in the study area.
APTC-SG1	As APTC-SG but if the most numerous/widespread child is a singleton or doubleton, or a unique/duplicate*, delete parent from all sites instead.
APTC-SG2	As APTC-SG but if the most numerous/widespread child is a singleton or doubleton, or a unique/duplicate*, delete parent from sites where it has no child, but reassign parent to child at sites where a child is present.

\* Implemented based on singleton/doubleton status for Chao1 estimator, but unique/duplicate status for Chao2 estimator.

doubletons (two individuals found). The related occurrence-based method, Chao2, hinges on the ratio of uniques (found at one site) to duplicates (found at two sites) (Chao et al., 2009). Because projected richness is assessed across entire data sets, redundancies (e.g., *Hexagenia* and *H. limbata*) must be removed at the dataset-scale even if they never co-occur in a site. Regardless of whether cross-dataset ambiguities are resolved by merging child-level taxa (potentially eliminating rare taxa) or by reassigning parents to children (making the children more common than they were), the ratio of rare to common taxa across a dataset can be altered (Gotelli and Colwell, 2011) which then affects projected taxa richness.

We calculated Chao1 and Chao2 and their associated confidence intervals (Estimate S software users guide; Colwell, 2013). Chao1 uses the ratio of taxa that are singletons to doubletons (one versus two individuals in study area) to predict the total number of taxa (i.e., projected taxa richness):

$$S_{est} = S_{obs} + f_1^2 / (2f_2)$$

where  $S_{est}$  represents the total number of taxa (detected + undetected),  $S_{obs}$  is the number of taxa detected,  $f_2$  is the number of doubletons, and  $f_1$  is the number of singletons (Chao et al., 2009). Chao2 uses the ratio of taxa that are uniques to duplicates (taxa found at one versus two sites):

$$S_{est} = S_{obs} + Q_1^2 / (2Q_2)$$

where  $Q_2$  is the number of duplicates and  $Q_1$  is the number of uniques (Chao et al., 2009). For Chao1 (or 2), the estimated number of undetected taxa increases as either the total number of singletons (uniques) or doubletons (duplicates) increases, or as the ratio of singletons: doubletons (uniques: duplicates) increases.

### 2.3. Methods for resolving ambiguous taxa

Every processing method to resolve ambiguous taxa removes redundancies at the site level. Some methods also resolve ambiguous taxa at the study area-level. For this analysis, we evaluated both site-level and study area-level methods, recognizing that study area-level methods are preferred for estimating projected taxa richness. All methods for resolving ambiguous taxa were carried out using automated scripts we developed in R (R Core Team, 2015).

The methods we developed, which we dubbed Assign Parent to Child (APTC) methods, are designed specifically to contrast effects of retaining or eliminating rare taxa. The APTC methods preserve real richness by neither deleting nor merging any child taxa, and retain

abundance information by reassigning parent counts to children rather than by deleting them. The APTC methods also generally preserve the status of rare taxa, because parents at the site are assigned to the most common child (most numerous-Chao1; most widespread-Chao2) rather than divided among all children (as in DPAC-SG). We tested both site-level (APTC-S) and study-area level variants of APTC (APTC-G, APTC-SG). With APTC-G, all parents are immediately assigned to the most common child in the study area, whereas with APTC-SG variants, parents are assigned to children at the site-level first and then any remaining parents are then assigned to the most common child at the study area. We tested two additional modifications of the APTC methods that avoid decreasing the rarity of very rare children by assigning parents to them. Respectively, these are APTC-SG1, which deletes rather than reassigns parents if the most common child in the study area is a singleton or doubleton (or unique or duplicate), and APTC-SG2, which deletes these parents at sites where the child does not occur but reassigns them to the child at sites where the parent and child co-occur. Note that for all the methods we implemented, we randomly chose a child to assign the parent to if multiple children were tied for most common.

We tested three additional methods recommended by Cuffney et al. (2007), namely Remove Parent Merge Child (RPMC), Remove Parent Keep Child (RPKC), and Divide Parents Among Children (DPAC) method (Table 1). Cuffney et al. (2007) suggested RPMC as the best method for the resolving ambiguities at the study area scale if imputing new child taxa is to be prevented. However, we anticipate that RPMC may have a large influence on projected richness because rare child records are lost by merging with parents. Cuffney et al. (2007) suggested RPKC as among the most effective methods for preserving real (as opposed to spurious) richness because parents are never retained. When applied at the dataset level (RPKC-SG rather than RPKC-S), this method may also influence projected richness, because identities are imputed to parents lacking site-level children, thereby changing child distribution patterns. Finally, Cuffney et al. (2007) recommended the DPAC-SG method as best for retaining abundance information. However, we anticipate the DPAC-SG method to result in a loss of taxa rarity when parents are divided among all children rather than assigned to the most common child at a site.

Given that the most common child at the study area could be described as either the most abundant or the most widespread (e.g., present at the most sites), we wanted to consider this in our analysis for all methods which assign parents to children in the study area. Prior to computing Chao1, which uses abundance information, we assigned remaining site-level children to the most abundant child in the study

area. For Chao2, which uses incidence information, we assigned remaining site-level children to the most widespread child (i.e., found at the most sites) in the study area. Also, for method APTC-SG1 and APTC-SG2, we defined rare children (i.e., those to whom parents would not get assigned) based on singleton or doubleton status for Chao1, and unique or duplicate status for Chao2. Because method DPAC could result in fractions of individuals, we considered taxa to be singletons if their abundance was exactly 1, and doubletons if abundance was  $> 1$  but  $\leq 2$ .

#### 2.4. Assessing outcome of methods for resolving ambiguous taxa

To illustrate the distribution of richness and abundance values obtained using each method for resolving ambiguous taxa, we plotted the mean proportion of original richness and abundance at each site  $\pm 1$  standard deviation (SD). The choice to assign site-level parents to the most numerous versus most widespread taxa in the study area affected the results for APTC-G (richness), and APTC-SG1 and APTC-SG2 (abundance), so for these methods, we present results using both approaches.

To estimate the effect of removing ambiguous taxa on projected taxa richness, we first counted the resulting number of singletons and doubletons, or uniques and duplicates. We then estimated projected taxa richness using Chao1 and Chao2, as well as the number of additional sites needed to reach 95% of total richness using the method provided by Chao et al. (2009). The amount of effort to detect the last five % of taxa is typically extremely high; therefore, we compared the effort needed to detect the majority of taxa (95%).

We also compared the ramifications of applying ambiguous taxa methods for IBI scores, the value of which could be affected by imputed taxonomic identity when parents are assigned to a child, loss of abundances when parents are deleted, and loss of taxonomic information when children are merged with parents. The metrics we evaluated included the number of EPT genera (i.e., belonging to the orders Ephemeroptera, Plecoptera, or Trichoptera), the number of Chironomidae genera, and the number of Odonata genera, as well as the percent of total counts at each site that was made up of each of these taxonomic groups. Since a typical means of evaluating ecological condition is to assign a score based on where the site falls relative to a range across sites, we compared the effect of ambiguous taxa resolution methods on IBI scores using the number of sites that changed their condition category. For metrics involving EPT taxa and Odonata, we considered the top third of the range to be in the good category and the bottom third of the range to be in the poor category, and the converse for the Chironomidae metrics. Because the RPKC-S method provides the most accurate estimate of number of taxa (ambiguities removed, no new taxa imputed), we compared all other taxonomic resolution methods to RPKC-S in assessing effects on richness-based IBI scores. We compared the abundance-based IBI scores to the original dataset.

#### 2.5. Assessing effect of subsampling on rare taxa

We use the subsampled counts directly to determine invertebrate composition (see above). However, subsampling occurred at 10 sites for SLRE and 16 sites for Isle Royale. A common procedure in handling subsampled taxonomic data is to extrapolate out to the whole sample by multiplying by the inverse of the subsample-fraction and dropping the coarse-sort counts (e.g., if 30 of 150 Oligochaeta were mounted and identified as 1 *Limnodrilus hoffmeisteri*, 2 *L. maumeensis*, and 27 “*Tubificinae* without setae”, they become 5 *L. hoffmeisteri*, 10 *L. maumeensis*, and 135 *Tubificinae* without setae, and 0 Oligochaeta). The two alternatives can produce different outcomes after resolving ambiguous taxa. The use-actual-count approach can result in parents that outnumber the subsampled children but still allow for children to be singletons and doubletons, whereas the extrapolation approach inherently removes the numerically dominant parents and eliminates the

possibility of children being singletons or doubletons (due to the multiplication). We examined the effects on rare taxa (singletons and doubletons, and uniques and duplicates) as well as projected taxa richness (Chao1, Chao2) if the direct versus subsample approach was used at subsampled sites. For this comparison, we used the dataset resulting from applying the APTC-SG method of resolving ambiguous taxa, which is the most conservative method of eliminating redundancies at the study area level.

### 3. Results

#### 3.1. Characteristics of original dataset

Among all samples and before resolving ambiguous taxa, there were 481 different taxa identified from Isle Royale samples and 344 identified from SLRE samples. Taxa spanned multiple orders and families of leeches (Hirudinea), oligochaete worms (primarily Tubificidae), mites (primarily Trombidiformes), insects (primarily Coleoptera, Diptera, Ephemeroptera, Hemiptera, Odonata, and Trichoptera), malacostracans (Amphipoda, Isopoda, Mysids), and mollusks (primarily Veneroida, Basommatophora, Heterostrophia, and Neotaenioglossa). Most specimens were identified to a fine level of resolution (Fig. 1). Only 10% of records for Isle Royale and 8% of records for SLRE were identified at a level coarser than order. For SLRE, these included taxa in the superorder Acariformes, and classes Bivalvia, Gastropoda, and Oligochaeta. For Isle Royale, these included taxa in the superorder Acariformes, and classes Gastropoda, Hirudinea, and Oligochaeta. For SLRE, approximately 32% of records represented an ambiguous parent, and 28% of these ambiguous parent records were in the family Chironomidae or a subfamily, genus, or tribe therein. For Isle Royale, approximately 42% of records represented an ambiguous parent, and 32% of these ambiguous parent records were in the family Chironomidae.

#### 3.2. Effect of ambiguous taxa methods on study area richness and abundance

A substantial number of records were merged, deleted, or re-assigned, depending on the ambiguous taxa method that we implemented (Table 2). Based on the difference in total taxa between the original dataset and the RPKC-SG method (which retained only non-ambiguous taxa groups in the study areas), 88 of 344 original taxa at SLRE were ambiguous taxa, and 137 of 481 original taxa at Isle Royale were ambiguous taxa. The DPAC method and all the study-area APTC method variants also retained all 256 of the non-ambiguous taxa for SLRE and all 344 of the non-ambiguous taxa for Isle Royale. The RPMC method removed all ambiguous taxa, but also removed multiple non-ambiguous

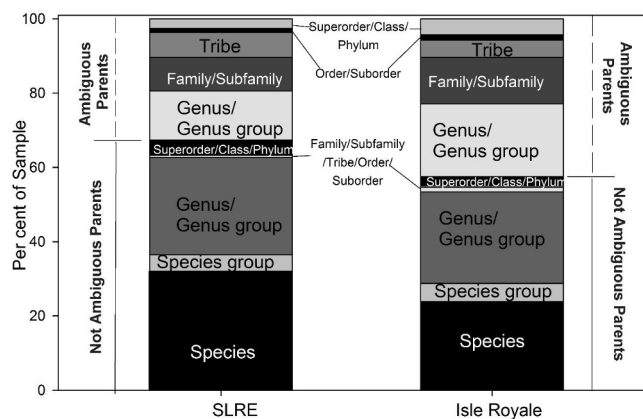


Fig. 1. Ambiguous parents spanned a range of taxonomic groups, with 32 per cent of records represented an ambiguous parent for SLRE and 42 per cent of records represented an ambiguous parent for Isle Royale.



**Table 2**

Description of changes made to the original data sets (percent records deleted, merged, or reassigned) and effects on study area richness and abundance. “NA” indicates not applicable. For APTC-SG1 and APTC-SG2, any values before the “/” represent the outcome if the method is applied to singletons and doubletons, while values after the “/” represent the outcome of the method is applied to uniques and duplicates (See Table 1). The estimate of non-spurious richness is denoted by an \*.

Method	% Records Deleted	% Records Merged	% Records Reassigned	# Unique Taxa	% of Original Abundance
<i>SLRE (nrecords = 2405, original abundance = 60209, 344 unique taxa)</i>					
RPKC-S	20.0	NA	NA	322	53.9
APTC-S	NA	NA	20.0	322	100.0
RPMC	18.6	28.01	NA	195	68.5
RPKC-SG	20.0	NA	12.6	256*	53.9
DPAC-SG	NA	NA	32.6	256	100.0
APTC-G	NA	NA	32.6	256	100.0
APTC-SG	NA	NA	32.6	256	100.0
APTC-SG1	1.7/5.2	NA	30.9/27.4	256	99.9/98.5
APTC-SG2	1.5/4.2	NA	31.4/28.5	256	99.9/99.1
<i>Isle Royale (nrecords = 6906; original abundance = 189116, 481 unique taxa)</i>					
RPKC-S	19.9	NA	NA	459	74.5
APTC-S	NA	NA	20.0	459	100.0
RPMC	20.0	9.9	NA	299	92.7
RPKC-SG	19.9	NA	22.6	344*	74.5
DPAC-SG	NA	NA	42.2	344	100.0
APTC-G	NA	NA	42.4	344	100.0
APTC-SG	NA	NA	42.4	344	100.0
APTC-SG1	2.0/8.8	NA	40.4/33.6	344	99.8/95.5
APTC-SG2	1.68	NA	40.8/34.5	344	99.8/95.9

taxa because children were merged into parents if the latter were numerous. With RPMC, 61 of the 256 unique non-ambiguous taxa were deleted for SLRE, and 45 of the unique non-ambiguous taxa were deleted for Isle Royale. Methods RPKC-S and APTC-S, which resolve redundancies only at the site scale, both substantially increase the number of unique taxa present across the study areas (322 at SLRE, and 459 at Isle Royale).

The various methods of resolving ambiguous taxa also differed substantially in the total study-area abundance retained. When site-level parents were deleted because they were ambiguous (methods RPKC-S and RPKC-SG), only ~54% of original abundance remained for SLRE, and ~75% of original abundance remained for Isle Royale. The RPMC method retained ~68% of abundance for SLRE and ~93% of abundance for Isle Royale. The greater amount of abundance lost for SLRE compared to Isle Royale is because oligochaete parent taxa at SLRE were less numerous than children, and therefore deleted, where at Isle Royale oligochaete parent taxa were more numerous and children were merged with parents. Abundance information was largely retained for the study area-level methods DPAC-SG and the APTC group, even when parents were deleted rather than re-assigned to rare taxa (under APTC-SG1 and APT-SG2).

### 3.3. Effect of ambiguous taxa methods on site-level richness and abundance

To evaluate effects on site-level richness, we compared results from each method to the number of taxa retained with the RPKC-S/APTC-S methods, which provide the best estimate of non-ambiguous taxa at a site without imputing taxa identified from the study area. We found that the number of unique taxa at a site differed greatly depending on which ambiguous taxa methods were implemented (Fig. 2). Patterns were similar to those observed at the study area level. For both study areas, approximately 20% of the richness estimated at the site-level was from presence of ambiguous taxa (e.g., RPKC-S and APTC-S compared to the original dataset). Of the methods that removed redundancies at the study area scale, RPMC had a much larger effect on dataset-level taxa richness compared to other methods. On average, sites in SLRE had 18% fewer taxa than with the RPKC-S/APTC-S methods, and sites in Isle

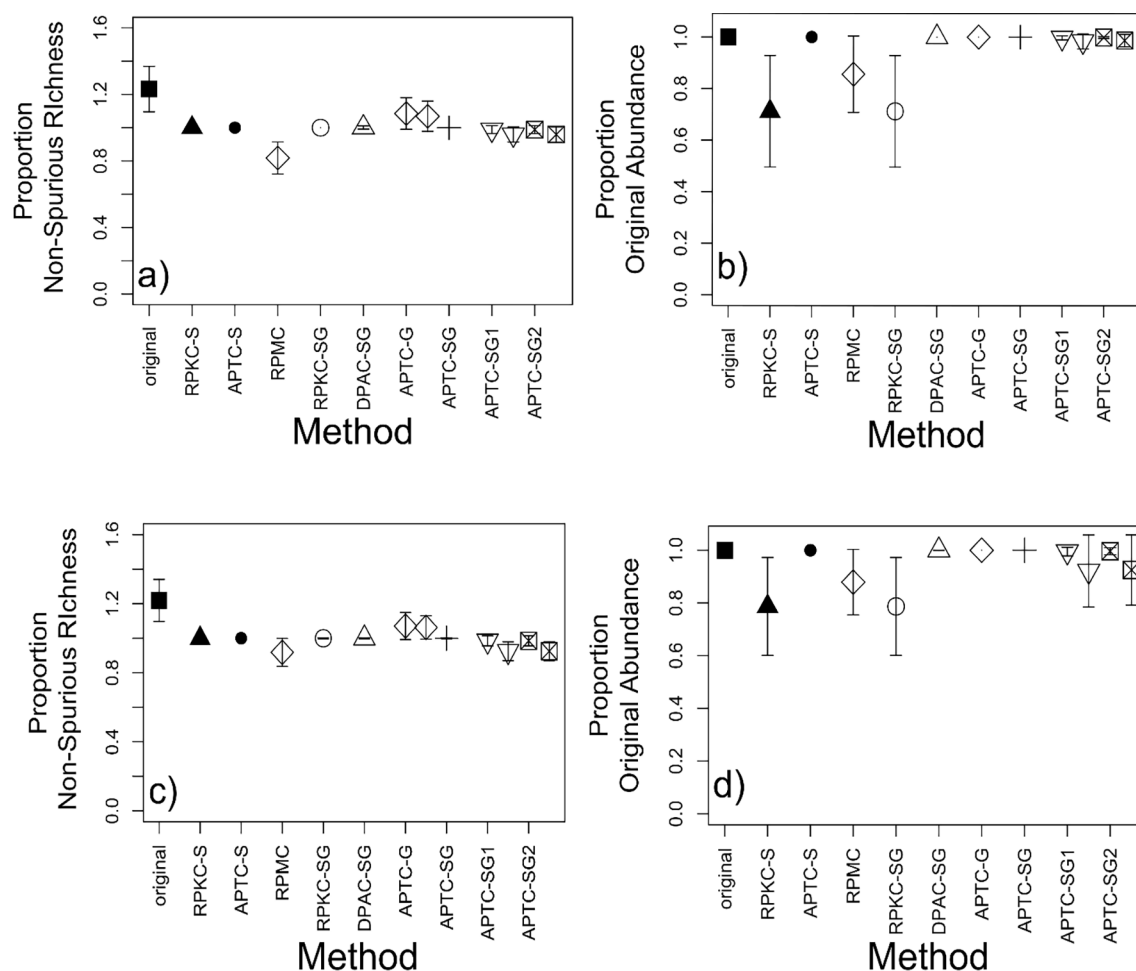
Royale had 9% fewer taxa than with the RPKC-S/APTC-S methods. The methods which assigned parents to the site before the dataset (DPAC-SG, RPKC-SG, APTC-SG) retained 100% of richness. For the method of immediately assigning parents to the most common child in the study area (APTC-G), estimated site-level richness changed slightly depending on whether ambiguous parents were assigned to the most numerous or the most widespread child. However, both methods overestimated richness, because parents were often assigned to a child not already at the site. Modifying the APTC-SG methods to retain rare taxa (APTC-SG1 & APTC-SG2) resulted in a slight decrease in richness (< 1%) at the site level when parents were not reassigned to singletons and doubletons, and a much larger decrease in richness when parents were not reassigned to uniques or duplicates (5% for SLRE and 8% for Isle Royale). This phenomenon occurred because uniques and duplicates were more prevalent than singletons or doubletons.

The amount of the original abundance retained at the site-level also differed greatly between ambiguous taxa methods (Fig. 2). Compared to the original dataset, the RPKC-S and RPKC-SG methods substantially reduced abundance at a site when all ambiguous parents at a site were removed (~28% for SLRE and ~21% for Isle Royale). All 100% of abundance was retained with the methods which re-assigned parents at the study area scale, but did not delete parents (DPAC-SG, APTC-G, APTC-SG). Some abundance was lost with the deletion of parents that would otherwise have been reassigned to a singleton or doubleton or unique or duplicate taxa (under APTC-SG1 and SG2). However, this decision had only minimal effects on abundance at the site scale (< 1% for singletons and doubletons; and 2% for uniques and duplicates at SLRE and 7% for uniques and duplicates at Isle Royale). Compared to the other study area-scale methods, the RPMC method resulted in a significant decrease in abundance (~14% for SLRE and ~12% for Isle Royale).

We also found substantial inter-site variation in the magnitude of richness and abundance effects. This variation is greatest for the RPKC methods, in which the standard deviation of abundance was 23% of mean abundance for Isle Royale and 25% of mean abundance for SLRE (Fig. 2). In general, the mean richness or abundance retained at individual sites (Fig. 2) was higher than the cumulative amount retained for the study area (Table 2), suggesting that many sites exhibited little change in taxa richness and abundance after implementing methods to resolve ambiguous taxa, while a smaller percentage of sites exhibited a large amount of change, presumably because taxonomic redundancies were concentrated in a small percentage of the sites.

### 3.4. Effect on rare taxa, projected taxa richness, and sampling effort to reach projected taxa richness

The choice of method to resolve ambiguous taxa substantially affected the resulting number of taxa with counts of one (a singleton) or two (a doubleton), which are the basis for estimating projected (asymptotic) richness and sampling effort based on Chao1 (Table 3). Methods RPKC-S and APTC-S, which resolve redundancies only at the site level, resulted in similar numbers of singletons and doubletons to the original dataset because parent taxa were left as unique taxa in the data set. The RPKC-S method produced a few additional singletons and doubletons because some taxa that were parents at one site but children at another site became singletons when the parents were removed. Of the study area-level methods, RPMC and DPAC-SG produced substantially lower numbers of singletons and doubletons than did APTC. For RPMC, this occurred because the least abundant of the parent–child pair were removed by deletion or merging. For the DPAC method, this occurred because dividing parents among all children present adds counts to rare child taxa as well as to common child taxa. The APTC-SG method, which assigns parents to the most abundant child at the site level before the study area level, produced fewer singletons and doubletons than the APTC-G method which considered only the study area, because the most abundant child at the site-scale is more likely to be a



**Fig. 2.** Mean proportion of original site-level richness and abundance retained for SLRE (A, B) and Isle Royale (C, D) varied greatly between methods (mean  $\pm$  1 SD), with site-level methods indicated by filled-in symbols, and study-area methods indicated by open symbols. For APTC-G, APTC-SG1, APTC-SG2, results differed depending on whether ambiguous parents were assigned to the a) most numerous children in the study area, or b) the most widespread children in the study area. The first point represents the result after if rule a was used, and the second point represents the result if rule b was used.

singleton or doubleton than the most abundant child at the study area scale. When parents were deleted rather than assigned to singleton or doubleton children (APTC-SG1), the singletons retained increased by 26% and 38% and doubletons retained by 25% and 40% for the SLRE and Isle Royale, respectively, compared to the method which ignored the effect on singletons and doubletons (e.g., APTC-SG). However, only re-assigning parents to a rare child if the child was found at the same site (APTC-SG2) resulted in a 12% and 13% reduction for singletons and a 17% and 15% reduction for doubletons, respectively.

The effect of the different methods on the number of uniques and duplicates was similar to the effect on the number of singletons and doubletons (Table 3). The RPKS-S and APTC-S methods produced the same or slightly higher numbers than the original data set. Again, the RPMC method most drastically reduced projected taxa richness, because uniques and duplicates are eliminated both by deleting less-numerous parents and merging less-numerous children into a more numerous parent. However, DPAC gave essentially the same result as the APTC study area methods with regard to uniques and duplicates. This difference regarding the effect on singletons and doubletons occurs because dividing parents among all children at a site (as DPAC does) does not alter the occurrence frequency of rare children across sites, although it does increase counts enough to affect singletons and doubletons. Deleting parents rather than re-assigning them if the child was a unique or duplicate (e.g., APTC-SG1 vs APTC-SG) slightly increased the number of uniques and duplicates retained (by 21% and 14% for SLRE and 19% and 4% for Isle Royale). However, only re-assigning

parents to a child if the child was found at the same site (APTC-SG2) did not further change the number of uniques and duplicates retained (as it did the number of singletons and doubletons), because the number of sites where the rare child was found remained the same regardless of whether a parent was deleted or re-assigned.

The aforementioned differences among methods in the number of rare taxa retained could substantially affect the projected taxa richness based on both Chao1 and Chao2 (Table 3). The two methods that resolved ambiguous taxa at the site-level only (APTC-S and RPKC-S) produced the highest projected richness, but much of this richness is spurious because it is an artifact of taxonomic redundancies that persist across the dataset. In contrast, the Assign Parent to Child (APTC) methods, which resolved redundancies both within and across sites, yielded substantial reductions in projected richness relative to the original dataset. Chao1 and Chao2 produced similar projected taxa richness among APTC method variants (Table 3) because neither the total number of non-ambiguous taxa nor the ratio of singletons to doubletons or uniques to duplicates varied greatly among them. The method of Remove Parent Keep Children at the study area scale (RPKC-G) had similar projected richness to the APTC methods, given similar total number of taxa as well as number of singletons and doubletons. As was expected from the effect on rare taxa, the RPMC method had a projected richness that was lower ( $\sim 23\%$  and  $\sim 9\%$  for Chao1 at SLRE and Isle Royale;  $\sim 23\%$  and  $\sim 13\%$  for Chao2 at SLRE and Isle Royale) than the APTC-G/S and RPKC-SG methods, due to some cases of rare children being merged with a parent. For Chao1, implementing the DPAC-SG

**Table 3**

Numbers of rare taxa (Chao1: singletons and doubletons, Chao2: uniques and duplicates), estimated total projected taxa richness at 100% of effort, and projected effort (additional samples needed) to reach 95% of richness using the Chao1 and Chao2 estimators, after implementing each method to resolve ambiguous taxa. The additional projected effort increases as the ratio of singletons to doubletons increases.

Method	Chao1					Chao2				
	No. Singletons	No. Doubletons	Projected Taxa Richness	Projected Effort (95%)	Singleton/Doubleton Ratio	No. Uniques	No. Duplicates	Projected Taxa Richness	Projected Effort (95%)	Uniques/Duplicates Ratio
<i>St Louis Estuary</i>										
original	66	37	404	183	1.78	120	46	498	316	2.61
RPKS-S	67	40	377	179	1.68	124	48	480	321	2.58
APTC-S	59	38	367	160	1.55	124	48	480	321	2.58
RPMC	33	17	226	188	1.94	64	23	283	331	2.78
DPAC-SG	24	20	270	97	1.20	79	29	362	319	2.72
RPKS-G	47	25	300	189	1.88	81	29	367	329	2.79
APTC-G	47	25	300	190	1.88	80	30	361	313	2.67
APTC-SG	42	24	292	170	1.75	81	29	367	329	2.79
APTC-SG1	53	30	302	187	1.77	98	33	399	367	2.97
APTC-SG2	47	28	295	171	1.68	98	33	399	367	2.97
APTC- sub	40	18	300	207	2.22	80	30	361	313	2.67
<i>Isle Royale</i>										
original	91	42	579	590	2.17	149	61	662	780	2.44
RPKS-S	93	43	559	605	2.16	150	62	639	784	2.42
APTC-S	78	37	540	537	2.11	151	62	642	791	2.44
RPMC	55	25	360	592	2.20	87	39	395	702	2.23
DPAC-SG	31	33	358	384	0.94	102	47	454	688	2.17
RPKS-G	63	34	397	507	1.85	101	46	454	693	2.20
APTC-G	63	31	407	546	2.03	101	46	454	693	2.20
APTC-SG	54	27	397	505	2.00	101	46	454	693	2.20
APTC-SG1	75	38	431	668	1.97	121	48	495	832	2.52
APTC-SG2	61	31	410	576	1.97	121	48	495	832	2.52
APTC- sub	52	28	391	467	1.86	101	46	454	693	2.20

method lowered the estimate of projected richness (7.5% lower for SLRE and 7% lower for Isle Royale) because rare taxa were reduced when parents were apportioned among them at the site-level. This reduction did not occur for Chao2, because adding parental counts to children at a site did not affect the child's occurrence frequency across sites.

### 3.5. Effect on effort to reach 95% of richness

The amount of additional effort needed to reach 95% of projected richness was sensitive to the number and ratio of singletons and doubletons (and uniques and duplicates). For SLRE, the estimated number of samples to reach 95% of projected taxa richness using Chao1 ranged from 3 for the DPAC-SG method to 96 for the APTC-G method (Table 3 and Fig. 3). For Isle Royale, the range was from 92 samples with the DPAC-SG method to 416 samples with the APTC-SG1 method. For Chao2, the difference in number of samples to reach 95% richness ranged from 219 for the APTC-G method to 273 for the APTC-SG1/APTC-SG2 methods, and from 436 for the DPAC-SG method to 580 for the APTC-SG1/APTC-SG2 methods (Table 3). The larger range among methods for Chao1 than Chao2 was related to the large reduction in singletons, doubletons, and their ratio produced by method DPAC-SG.

These effects also differed by study area. For Chao1, not assigning a parent to a child that is singleton or doubleton (APTC-SG1 or -SG2) had a larger effect on the estimated number of samples to reach 95% of projected taxa at Isle Royale compared to SLRE. With the exception of the methods that did not re-assign parents to rare taxa (APTC-SG1 & APTC-SG2), differences in projected taxa richness across methods were comparatively less for Chao2. However, the choice to retain rare taxa with the APTC-SG1 and APTC-SG2 methods resulted in a substantially greater ratio of uniques than duplicates, and thus a greater amount of effort need to reach 95% of richness compared to the other methods.

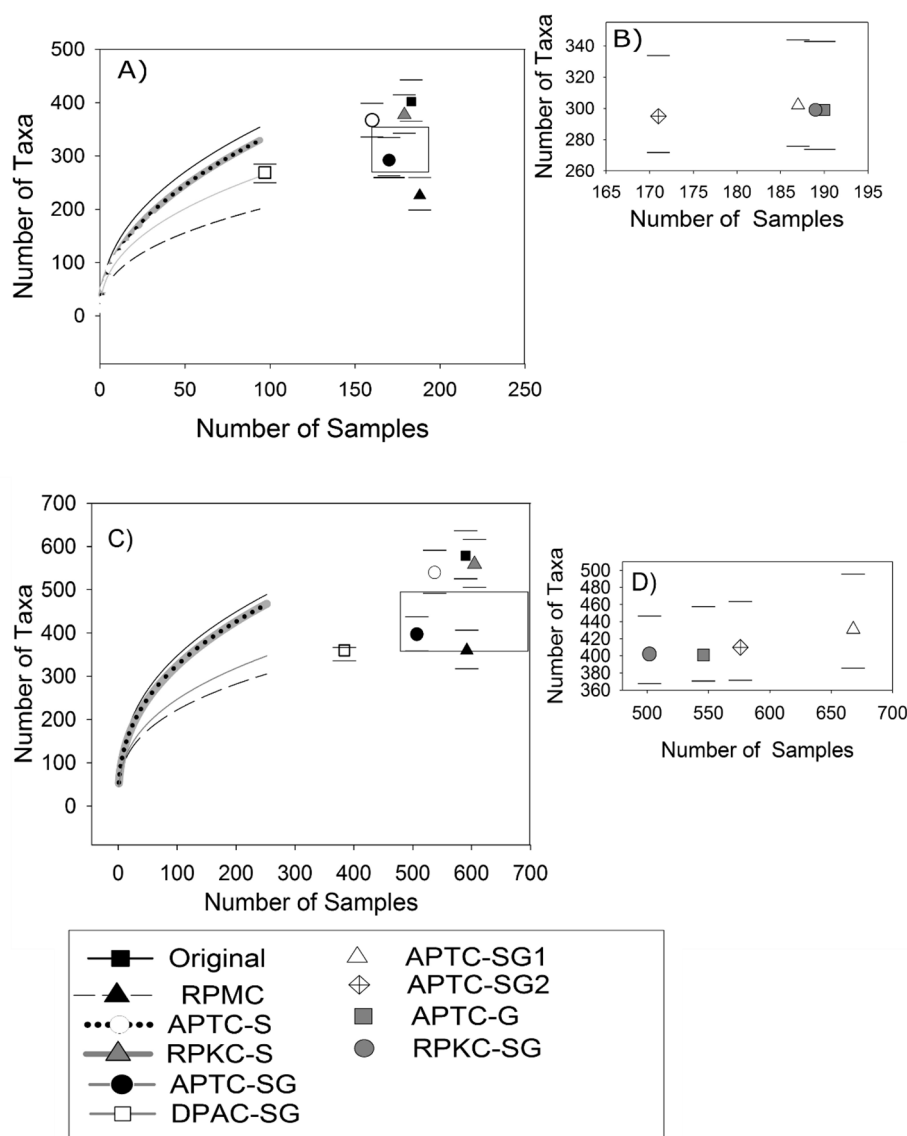
### 3.6. Effect on IBI scores

Ambiguous taxa methods can influence prioritization of sites by changing where a site falls within a range of values for an IBI metric (Table 4). Abundance-based IBI metric scores can be affected by any methods that delete relevant parent-taxa rather than reassigning their counts to a child (e.g., RPMC, RPKC-S, RPKC-SG, and APTC-SG1 and APTC-SG2), although the extent of the effect depends on how many parent taxa are in that situation. Richness-based IBI scores can be affected by any method which deletes rather than re-assigns unique taxa at a site (RPMC, APTC-SG1, APTC-SG2), or which re-assigns taxa to another taxa in the study area without first considering whether children are present at the site, thereby creating a new site-level taxa at some sites (APTC-G).

To summarize across richness-based and abundance-based IBI metrics, we found no effects on the distribution of sites across score categories for the methods APTC-S, DPAC-SG, and APTC-SG. Methods RPKC-S and RPKC-SG exhibited changes only for abundance-based IBI metrics and method APTC-G exhibited changes only for richness-based IBI metrics. Methods APTC-SG1, APTC-SG2, and RPMC all had potential for effects on both abundance-based and richness-based IBI categories but these effects were most frequent and largest in magnitude for the RPMC method, with the three APTC metrics being affected less often. Because the APTC-SG1 and APTC-SG2 methods delete parent taxa only when the potential children are singletons and doubletons or uniques and duplicates, these methods less often produced category-changes for abundance-based IBI metrics (i.e., relative to the original abundance in the dataset) than did RPMC, RPKC-S, and RPKC-SG (Table 3).

### 3.7. Effect of alternate treatment of subsampled data

When we compared the outcome of applying method APTC-SG using estimated whole sample counts (i.e., counts multiplied by the



**Fig. 3.** Methods showed differences in projected taxa-richness using the Chao1 estimator, and required sampling effort as illustrated by taxa accumulation curves for SLRE (A, B) and Isle Royale (C, D). The methods RPMC, DPAC-SG, APTC-S, RPKC-S, and APTC-SG are shown in main graphs A and C. Dataset-level methods with highly-overlapping projected-taxa-richness (RPKC-SG, APTC-G, APTC-SG, APTC-SG1, APTC-SG2) are found within the rectangle and are further depicted in graphs B and D.

inverse of the subsample fraction) versus actual counts at subsampled sites, we found an overall reduction in singletons and doubletons. For SLRE, there were two fewer singletons and six fewer doubletons; for Isle Royale, there were nine fewer singletons and three fewer doubletons. Because the APTC-SG method involves assigning parents to the most numerous child, the rare taxa affected were those that were not the most numerous children of ambiguous parents.

The ratio of singletons to doubletons increased at SLRE (2.22 vs. 1.75 without multiplying out) whereas it decreased at Isle Royale (1.86 vs. 2.00 without multiplying out). As a consequence, Chao1 also increased at SLRE (from 292 to 300), but decreased at Isle Royale (from 407 to 391). These results illustrate the importance of both the ratio of singletons: doubletons as well as the number of singletons and doubletons, and suggest that the effect of subsample extrapolation on singletons and doubletons and the projected richness computed from them depend on taxa composition of the study area. These effects could be larger if subsampling was conducted at a high number of sites than in our study.

We found no change in the number of uniques or duplicates at either SLRE or Isle Royale between data treatments. Subsample extrapolation does not alter the pattern of child taxa occurrence across sites.

#### 4. Discussion

##### 4.1. Choosing a method to resolve ambiguous taxa

Quantifying regional or site-specific biodiversity – using documented sample-processing methods and comparable datasets – is fundamental to ecosystem characterization, ecological assessment, exotic species detection, and conservation planning (e.g., Gotelli and Colwell, 2011; Cao et al., 1998; Hoffman et al., 2011; Pressey et al., 2007). An important aspect of attaining comparable datasets is to remove taxonomically ambiguous taxa. All methods of resolving these taxonomic ambiguities proceed by some combination of deleting, merging, or re-assigning taxa, which alters the total number of taxa and their pattern of incidence and abundance across sites and datasets. Our study makes



**Table 4**  
Number of sites which changed IBI category from good → average/poor, and from good/average → poor after ambiguous taxa methods were applied. The top row gives the number of sites originally in the first IBI category of the two being compared. For the richness-based metric, the comparison is relative to the RPKC-S method (which removed spurious richness), whereas for the abundance-based metrics, the comparison is to the original dataset. For methods with varying outcomes (\*), if a “/” is present for APTC-G, the first value is the result if the parent is assigned to the most numerous child, and the second value is if the parent is assigned to the most widespread child. For APTC-SG1 and APTC-SG2, the first value is if the parent is assigned to the most numerous child, or deleted according to the singleton/doubleton rule (see Table 1) while the second value is if the parent assigned to the most widespread child, or deleted according to the unique/duplicate rule (see Table 1).

SLRE	Richness-based						Abundance-based											
	No. EPT Genera			No. Chironomid Genera			No. Odonata Genera			Percent EPT			Percent Chironomidae			Percent Odonata		
	good → avg/ poor	avg/ good → poor	avg/ good → poor	good → avg/ poor	avg/ good → poor	avg/ good → poor	good → avg/ poor	avg/ good → poor	avg/ good → poor	good → avg/ poor	avg/ good → poor	avg/ good → poor	good → avg/ poor	avg/ good → poor	avg/ good → poor	good → avg/ poor	avg/ good → poor	avg/ good → poor
Starting #	5	81	53	5	6	78	5	73	61	6	92	2	1	1	1	1	1	1
RPKC-S	–	–	–	–	–	–	1	5	6	0	1	0	0	0	0	0	0	0
APTC-S	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
RPMC	2	2	0	0	1	1	4	11	4	1	0	0	1	0	0	0	0	0
RPKC-SG	0	0	0	0	0	0	1	5	6	0	1	0	0	1	1	0	0	0
DPAC-SG	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
APTC-G*	0	0	5/3	3/2	0/3	0	0	0	0	0	0	0	0	0	0	0	0	0
APTC-SG	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
APTC-SG1*	0/1	0/3	0	0	0	0	0/1	0/3	0/2	0/1	0	0/1	0/1	0	0	0	0	0
APTC-SG2*	0/1	0/3	0	0	0	0	0/1	0/3	0/2	0/1	0	0/1	0/1	0	0	0	0	0
Isle Royale	good → avg/ poor	avg/ good → poor	good → avg/ poor	avg/ good → poor	good → avg/ poor	avg/ good → poor	good → avg/ poor	avg/ good → poor	good → avg/ poor	avg/ good → poor	avg/ good → poor	good → avg/ poor	avg/ good → poor	avg/ good → poor	avg/ good → poor	good → avg/ poor	avg/ good → poor	avg/ good → poor
Starting #	4	211	143	9	1	247	6	234	222	5	246	3	1	1	1	1	1	1
RPKC-S	–	–	–	–	–	–	3	1	10	1	1	0	0	0	0	0	0	0
APTC-S	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
RPMC	1	6	0	6	0	2	0	0	5	0	0	0	0	0	0	0	0	0
RPKC-SG	0	0	0	0	0	0	3	1	10	1	1	1	1	1	1	1	1	1
DPAC-SG	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
APTC-G*	0	0	15/13	0/4	0	0	0	0	0	0	0	0	0	0	0	0	0	0
APTC-SG	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
APTC-SG1*	1/1	0/15	0	0	0	0	0	0	0/4	0	0	0	0	0	0	0	0	0
APTC-SG2*	1/1	0/15	0	0	0	0	0	1/0	0/3	0	0	0	0	0	0	0	0	0

an important contribution to our understanding of the effects of ambiguous taxa resolution methods with respect to projected taxa richness and associated ecological applications (i.e., survey completeness assessment, IBI metrics) by developing methods designed to retain rare taxa information and evaluating existing methods with known effects on species richness and abundance (i.e., the suite of methods developed by Cuffney et al. (2007).

To estimate projected taxa richness, which is calculated based on taxa data from a collection of sites in a study area, taxonomic ambiguities should be resolved at the study area scale. Methods which resolve ambiguities only at the site scale, such as RPKC-S and APTC-S, increase study-area richness because of the persistence of parents lacking site-level children. Methods which consider among-site information in reassigning parents (i.e. APTC-G and the APTC-SG variants) resolve redundancies data-set wide and therefore are more suitable for estimating projected taxa richness. However, these methods are not suitable for comparing distributions of individual taxa, because the occurrence and relative abundance patterns among children are distorted when the identity of the most abundant study-area child is imputed to parents without site-level children. The RPKC-SG method also resolves taxonomic ambiguities at the study-wide scale, but shares the drawbacks of the APTC-G and APTC-SG methods and has the additional drawback of eliminating entirely from the dataset the counts of parents that do have site-level children.

Two methods recommended by Cuffney et al. (2007) seem not suitable for projecting taxa richness. Cuffney et al (2007) recommended DPAC-SG to preserve relative abundance among taxa, but we found this method eliminates singleton and doubleton taxa and thus causes the value of Chao1 to decrease. The DPAC-SG method can be used with Chao2, however, since dividing parent counts among children does not alter their unique or duplicate status, and the most widespread study-area child (to which the remaining parents are assigned) is unlikely to be a unique or duplicate. Cuffney et al. (2007) recommended the RPMC method to balance preserving richness and abundance, but we found that RPMC significantly reduced projected richness because multiple rare children were merged into common parent taxa. An important rationale for using RPMC is to not pass on data for which there is uncertainty as to whether a taxon was actually present at a site (i.e., imputed identities for parents lacking site-level children). However, the RPMC method also deleted some low-abundance child taxa entirely, which is an obvious loss of abundance information. While no other single method removes taxonomic redundancy study-area wide without imputing any identities to parents, we recommend using multiple different methods in place of RPMC.

#### 4.2. Additional method considerations for bio-assessment

When deciding whether methods for resolving ambiguous taxa are appropriate for bio-assessment, the resolution at which organisms are identified in relation to the resolution needed to calculate a particular metric is paramount. Abundance-based IBI scores are affected by the deletion of counts of any taxa that contributes to the metric, as evidenced by score-category changes for percent Ephemeroptera, Plecoptera, and Trichoptera (EPT), Chironomidae, or Odonata taxa under the RPMC, RPKC-S, and RPKC-SG methods. In contrast, abundance-based IBI scores are unaffected by reassigning the counts of parents to children (as in DPAC and the APTC method suite) as long as the taxonomic resolution of the parents is at least as good as that required for the metric. Abundance-based IBI scores would be adversely affected if parent identifications were so coarse that the counts were included in the IBI only when reassigned to children, but this was not an issue with our dataset.

Richness-based IBI scores are conservatively estimated by counting the fully identified children at each site, or if no children are present, assuming the parents represent a single, child taxa. Of the methods we tested, those that consistently produced richness-based IBI scores

different from this conservative choice were RPMC (by merging multiple children into a single more numerous parent) and APTC-G (by imputing to parents an identity different than the children already present). Method APTC-SG1 and APTC-SG2 can also affect richness-based IBI scores because the parent contribution to richness can be lost if the only dataset level child is a singleton or doubleton (or unique or duplicate) not found at the same site (in which case it would be spurious, but deleted (SG1); or re-assigned (SG2)). If parent identifications are too coarse to be included in an IBI, then any method that imputes identities to parents lacking site-level children would affect richness-based IBI scores; but, as noted above, this was not an issue with our dataset. Imputing parents to children could also affect bio-assessment metrics that rely on finely-resolved taxonomic information (e.g., species and genus-level tolerance values), such as with the Hilsenhoff Biotic index or an Observed/Expected metric (Hilsenhoff, 1987; Clarke et al., 2003).

#### 4.3. Possible effect of study area characteristics and survey design on outcome

The taxonomic composition of the study area can interact with the taxonomic resolution method to influence projected taxa richness estimates. For example, we found that the decrease in singletons and doubletons was greater for SLRE than for Isle Royale when the Remove Parent Merge Child (RPMC) method was implemented. As a result, the Chao1 value was more similar to the other data-set level methods for SLRE than Isle Royale. Also, Chao1 values were more similar among the APTC variants at SLRE than at Isle Royale, because in the SLRE few singletons or doubletons had ambiguous parents at the same site (in which case they would lose singleton or doubleton status with the APTC-SG2 versus APTC-SG1 methods). As such, we recommend that when choosing a method of resolving taxonomic ambiguities, specific characteristics of the dataset be considered, including the level of rarity and the numerical ratio of parents to children, which may vary inherently among study areas but which also depends on sample collection methods and details of the laboratory enumeration (e.g., taxonomic expertise and resolution, use of subsampling).

The finding that the estimated number of samples to detect 95% of projected richness varied substantially between methods of resolving ambiguous taxa has implications for survey design and evaluation. For example, under the binational 2012 Great Lakes Water Quality Agreement, US and Canadian federal agencies are tasked with implementing monitoring programs for exotic aquatic species, and projected taxa richness (Chao1) is being used to estimate the sampling effort needed to detect taxa of a given rarity (which equates to a given level of establishment for exotic species), as well as develop more efficient monitoring programs (Hoffman et al., 2011; Ram et al., 2014). Given the sensitivity of estimated sampling effort to detect 95% of projected taxa, appropriate caution is warranted when determining how many additional samples are needed as part of a monitoring program. In addition to taxonomic processing choices, other factors may influence estimates of projected taxa richness. For instance, the extent and timing of sampling can influence the number and ratio of singletons or doubletons (or uniques and duplicates) documented within a study area (Hoffman et al., 2016); as could the sampling mesh size used, with larger mesh simplifying the enumeration job (Pinna et al., 2013) but also reducing the number of rare taxa collected. Investigators should carefully consider both the data methods and sampling strategy when interpreting estimates of future effort required to detect a desired level of taxa richness.

#### 4.4. Impact of method modifications to retain rare taxa

In highly diverse study areas, it may be necessary to develop nuanced methods of resolving ambiguous taxa (e.g., a case-by-case determination of why finer resolution identification was not possible),

or to let ambiguous parents stand as unique taxa. In our analysis, the tested modifications to method APTC-SG intended to preserve rare taxa (i.e., to not allow singletons and doubletons or uniques and duplicates to be eliminated, APTC-SG1 and APTC-SG2) had relatively little effect on projected-taxa-richness. However, these modifications could have a larger impact in places with extremely high species richness and rarity. For instance, after 30-years of surveying ants in the tropics, Longino et al. (2002) found that many taxa still remained rare (e.g., singletons and doubletons) and the taxa accumulation curve had still not reached an asymptote. In this case, choosing whether or not to assign parents to rare children could result in a much larger effect on projected richness than we observed. In addition, the reason why parents are “ambiguous” in the first place may affect the method choice. In our study, ambiguous parents were typically damaged specimens or unmounted chironomids and oligochaetes, which were assigned to a child already present. In other study areas, ambiguous parents may be unidentifiable because they have features not present in current taxonomic keys; if they are in fact new species, then they are by definition different from the previously identified children.

#### 4.5. Successful implementation requires consistent and detailed supporting information

We urge careful documentation of methods chosen and derivative datasets generated to ensure repeatability and comparability among studies. We encourage investigators to provide these details in their publications, rather than to make general statements which could refer to one of several methods. Comparisons of taxonomic composition between studies or locations are meaningful only if they use the same method of resolving taxonomic redundancies. For repeatability, resolving ambiguous taxa is best accomplished with a computer program supplied with a taxonomic-hierarchy data file. Depending on the complexity of the original taxonomic data, this hierarchy file may need to be expanded beyond the familiar phylum-class-order-family-genus-species levels for the automated processing to work. For example, some of the identifications returned by the taxonomists were at levels between the major elements of the taxonomic hierarchy (e.g., superorder “Acariformes”, genus group “Bezzia/Palpomyia”, species group “Cladotanytarsus vanderwulpi group”), which required us to construct correspondingly resolved hierarchies.

#### 4.6. Broad guidelines for choosing a method to resolve ambiguous taxa

We summarize our results by offering some broad guidelines for choosing among methods to resolve ambiguous taxa (Table 5). This choice should involve careful consideration of the benefits and drawbacks with respect to the particular study goals, the taxonomic resolution at which the data were collected, and the importance of rare taxa. No single method both avoids imputing taxonomic identities to parents at the site level and fully resolves ambiguous richness while preserving real (i.e., taxonomically confirmed) richness at the dataset level; accordingly, tradeoffs must be made or multiple methods must be employed to meet various analysis goals. For our research, which uses projected taxa richness as a basis for evaluating taxa detection probabilities under alternative survey designs, we will choose a method which assigns ambiguous parents to the most common or abundant child at a site, and then imputes remaining parents to the most common or abundant child at the study area. Increased attention to this aspect of taxonomic data analysis will inform broader understanding of the applicability and comparability of biodiversity endpoints important for ecosystem characterization, conservation planning, and rare species detection.

#### Acknowledgements

We thank Gerald Shepard, Adam Frankiewicz, and Brent Gilbertson

**Table 5**  
Summary of which methods for resolving ambiguous taxa are appropriate for which study goals or reasons why not. In this table, the word “rare” is used specifically as shorthand for taxa that are singletons and doubletons, or uniques and duplicates.

Goal:	RPKC-S	APTC-S	RPMC	RPKC-SG	DPAC-SG	APTC-G	APTC-SG	APTC-SG1	APTC-SG2
Eliminate redundancies study-area wide	No, cross-site redundancy remains	No, cross-site redundancy remains	Yes	Yes	Yes	Yes	Yes	Yes	Yes
Document site occupancy using only fully-identified taxa	Yes	Yes	No, some children lost with merging	No, can impute children	No, can impute children	No, can impute children	No, can impute children	No, can impute children	No, can impute children
Preserve counts, e.g., for abundance-based IBI metric	No, deletes site-level parents	Yes*	No, some parents deleted	No, some parents deleted	Yes*	Yes*	Yes*	Mostly, deletes parent of rare*	Mostly, may delete parent of rare*
Preserve true site richness, e.g., for richness-based IBI metrics	Yes	Yes	No, some children lost with merging	Yes	Yes*	No, may add site-level children	Yes*	Mostly, deletes parent of rare*	Mostly, May delete parent of rare
Project richness with Chao1 using singleton/doubleton info	No, cross-site redundancy remains	No, cross-site redundancy remains	No, some children lost with merging	Mostly, parent imputed to least-rare child	No, parent often imputed to rare	Mostly, parent imputed to least-rare child	Mostly, parent imputed to least-rare child	Yes	Mostly, parent sometimes added to rare
Project richness with Chao2 using unique/duplicate info	No, cross-site redundancy remains	No, cross-site redundancy remains	No, some children lost with merging	Mostly, parent imputed to least-rare child	Mostly, parent imputed to least-rare child	Mostly, parent imputed to least-rare child	Mostly, parent imputed to least-rare child	Yes	Yes

\* Note: Suitability for IBI metrics assumes that both parents and children are IBI targets; see text for discussion.

for the taxonomy work, and the many individuals who contributed to the Isle Royale and St. Louis River Estuary sampling campaigns. The views expressed in this article are those of the authors and do not necessarily reflect the views or policies of the U.S. Environmental Protection Agency.

## References

- Abell, R., Thieme, M.L., Revenga, C., Bryer, M., Kottelat, M., Bogutskaya, N., Coad, B., Mandrak, N., Balderas, S.C., Bussing, W., Stiassny, M.L.J., Skelton, P., Allen, G.R., Unmack, P., Naseka, A., Ng, R., Sindorf, N., Robertson, J., Armijo, E., Higgins, J.V., Heibel, T.J., Wikramanayake, E., Olson, D., López, H.L., Reis, R.E., Lundberg, J.G., Sabaj Pérez, M.H., Petry, P., 2008. Freshwater ecoregions of the world: a new map of biogeographic units for freshwater biodiversity conservation. *Bioscience* 58, 403–414.
- Blocksom, K.S., Flotemersch, J.E., 2005. Comparison of macroinvertebrate sampling methods for nonwadeable streams. *Environ. Monit. Assess.* 102, 243–262.
- Cao, Y., Williams, D.D., Williams, N.E., 1998. How important are rare species in aquatic community ecology and bioassessment? *Oceanography* 43, 1403–1409.
- Carter, J.L., Resh, V.H., 2001. After site selection and before data analysis: sampling, sorting, and laboratory procedures used in stream benthic macroinvertebrate monitoring programs by USA state agencies. *J. N. Am. Benthol. Soc.* 20, 658–682.
- Chao, A., Colwell, R.K., Lin, Chih-Wei, Gotelli, N.J., 2009. Sufficient sampling for asymptotic minimum species richness estimators. *Ecology* 90, 1125–1133.
- Chiarucci, A., Enright, N.J., Perry, G.L.W., Miller, B.P., Lamont, B.B., 2003. Performance of nonparametric species richness estimators in a high diversity plant community. 9: 283–295.
- Clarke, R.T., Wright, J.F., Furse, M.T., 2003. RIVPACS models for predicting the expected macroinvertebrate fauna and assessing the ecological quality of rivers. *Ecol. Model.* 160, 219–233.
- Colwell, R., 2013. EstimateS: Statistical estimation of species richness and shared species from samples. Version 8.2. User's Guide and application. 2009.
- Cuffney, T., Bilger, M., Haigler, A., 2007. Ambiguous taxa: effects on the characterization and interpretation of invertebrate assemblages. *J. N. Am. Benthol. Soc.* 26, 286–307.
- D'Allessandro, L., Fattorini, L., 2002. Resampling estimators of species richness from presence-absence data: why they don't work. *Metron-Int. J. Stat. LX*, 3–17.
- Davies, S., Jackson, S., 2006. The biological condition gradient: a descriptive model for interpreting change in aquatic ecosystems. *Ecol. Appl.* 16, 1251–1266.
- Desmet, P., Cowling, R., 2004. Using the species-area relationship to set baseline targets for conservation. *Ecol. Soc.* 9, 11.
- Gotelli, N.J., Colwell, R.K., 2011. Estimating species richness. *Biol. Divers.* 12, 39–54.
- Hilsenhoff, W.L., 1987. An improved biotic index of organic stream pollution. *Great Lakes Entomol.* 20, 31–40.
- Hoffman, J.C., Kelly, J.R., Trebitz, A.S., Peterson, G.S., West, C.W., 2011. Effort and potential efficiencies for aquatic non-native species early detection. *Can. J. Fish. Aquat. Sci.* 68, 2064–2079.
- Hoffman, J.C., Schloesser, J., Trebitz, A.S., Peterson, G., Gutsch, M., Quinlan, H., Kelly, J.R., 2016. Sampling design for early detection of aquatic invasive species in Great Lakes ports. *Fisheries* 41, 26–37.
- King, R.S., Richardson, C.J., 2002. Evaluating subsampling approaches and macroinvertebrate taxonomic resolution for wetland bioassessment. *J. N. Am. Benthol. Soc.* 21, 150–171.
- Lammert, M., Allan, A., 1999. *J. Environ. Manage.* 23, 257–270.
- Longino, J.T., Coddington, J., Colwell, R.K., 2002. The ant fauna of a tropical rain forest: estimating species richness three different ways. *Ecology* 83, 689–702.
- Novotny, V., Bartošová, A., O'Reilly, N., Ehlinger, T., 2005. Unlocking the relationship of biotic integrity of impaired waters to anthropogenic stresses. *Water Res.* 39, 184–198.
- Pinna, M., Marini, G., Rosati, I., Neto, J.M., Patrício, J., Marques, J.C., Basset, A., 2013. The usefulness of large body-size macroinvertebrates in the rapid ecological assessment of Mediterranean lagoons. *Ecol. Indic.* 29, 48–61.
- Pressey, R.L., Cabeza, M., Watts, M.E., Cowling, R.M., Wilson, K.A., 2007. Conservation planning in a changing world. *Trends Ecol. Evol.* 22, 583–592.
- R Core Team, T., 2015. R: A Language and Environment for Statistical Computing. R Foundation for statistical computing, Vienna, Austria.
- Ram, J.L., Banno, F., Gala, R.R., Gizicki, J.P., Kashian, D.R., 2014. Estimating sampling effort for early detection of nonindigenous benthic species in the Toledo Harbor Region of Lake Erie. *Manage. Biol. Invasions* 5, 209–216.
- Reese, G.C., Wilson, K.R., Flather, C.H., 2014. Performance of species richness estimators across assemblage types and survey parameters. *Glob. Ecol. Biogeogr.* 23, 585–594.
- Sambuichi, R.H., Haridasan, M., 2007. Recovery of species richness and conservation of native Atlantic forest trees in the cacao plantations of southern Bahia Brazil. 16: 3681–3701.
- Silva, D.R.O., Ligeiro, R., Hughes, R.M., Callisto, M., 2016. The role of physical habitat and sampling effort on estimates of benthic macroinvertebrate taxonomic richness at basin and site scales. *Environ. Monit. Assess.* 188, 340.
- Schreiber, J., Brauns, M., 2010. How much is enough? Adequate sample size for littoral macroinvertebrates in lowland lakes. *Hydrobiologia* 649, 365–373.
- Walther, B.A., Moore, J.L., 2005. The concepts of bias, precision, and accuracy, and their use in testing the performance of species richness estimators, with a literature review of estimator performance. *Ecography* 28, 815–829.
- Weigel, B.M., Henne, L.J., Martinez Rivera, L.M., 2002. Macroinvertebrate-based index of biotic integrity for protection of streams in west-central Mexico. *J. N. Am. Benthol. Soc.* 21, 686–700.