

# Many species in one: DNA barcoding overestimates the number of species when nuclear mitochondrial pseudogenes are coamplified

Hojun Song<sup>\*†</sup>, Jennifer E. Buhay<sup>\*‡</sup>, Michael F. Whiting<sup>\*</sup>, and Keith A. Crandall<sup>\*</sup>

<sup>\*</sup>Department of Biology, Brigham Young University, Provo, UT 84602; and <sup>†</sup>Belle W. Baruch Institute for Marine Sciences, University of South Carolina, Columbia, SC 29208

Edited by W. Ford Doolittle, Dalhousie University, Halifax, NS, Canada, and approved July 14, 2008 (received for review March 28, 2008)

Nuclear mitochondrial pseudogenes (numts) are nonfunctional copies of mtDNA in the nucleus that have been found in major clades of eukaryotic organisms. They can be easily coamplified with orthologous mtDNA by using conserved universal primers; however, this is especially problematic for DNA barcoding, which attempts to characterize all living organisms by using a short fragment of the mitochondrial cytochrome *c* oxidase I (COI) gene. Here, we study the effect of numts on DNA barcoding based on phylogenetic and barcoding analyses of numt and mtDNA sequences in two divergent lineages of arthropods: grasshoppers and crayfish. Single individuals from both organisms have numts of the COI gene, many of which are highly divergent from orthologous mtDNA sequences, and DNA barcoding analysis incorrectly overestimates the number of unique species based on the standard metric of 3% sequence divergence. Removal of numts based on a careful examination of sequence characteristics, including indels, in-frame stop codons, and nucleotide composition, drastically reduces the incorrect inferences of the number of unique species, but even such rigorous quality control measures fail to identify certain numts. We also show that the distribution of numts is lineage-specific and the presence of numts cannot be known *a priori*. Whereas DNA barcoding strives for rapid and inexpensive generation of molecular species tags, we demonstrate that the presence of COI numts makes this goal difficult to achieve when numts are prevalent and can introduce serious ambiguity into DNA barcoding.

cytochrome *c* oxidase I | Decapoda | Orthoptera

The orthology of characters is one of the fundamental and implicit assumptions in the use of DNA sequence data to reconstruct phylogeny or to establish “barcodes” for species. If the orthology assumption is violated, that is, whether paralogous sequences are unknowingly treated as orthologs, incorrect inferences are inevitable (1). This is especially true for the DNA barcoding initiative, which relies on the premise that all organisms have a unique and identifiable molecular tag, namely, a short region of mitochondrial cytochrome *c* oxidase subunit 1 (COI) amplified by universal primers, and that one is comparing only orthologs among species when formulating barcodes (2). As such, DNA barcoding relies on the assumption that the COI fragments generated by PCR from genomic DNA represent orthologous copies of mitochondrial DNA (mtDNA). Increasing empirical evidence suggests; however, that this assumption does not always hold true and that there are a number of molecular evolutionary processes that can hinder correct amplification and identification of the orthologs (3), including (i) duplication of the gene of interest within the mitochondrial genome (4), (ii) heteroplasmy (5), (iii) bacterial infection biasing mtDNA variation (6), and (iv) nuclear integration of mtDNA (7, 8). If a portion of COI was duplicated in a given species, conventional PCR might amplify both the correct and duplicated COI fragments, thus introducing ambiguity into the barcoding whether the paralogous copy had diverged since duplication. Hetero-

plasmly is the presence of a mixture of more than one type of mitochondrial genome within a single individual, and the coamplification of divergent heteroplasmic copies of mtDNA would lead to an overestimation of the number of unique species under barcoding (3). Maternally inherited symbionts, such as *Wolbachia*, can cause linkage disequilibrium with mtDNA and, whether a population becomes infected with such symbionts, the mtDNA associated with the initial infection will spread throughout the population and result in the homogenization of mtDNA haplotypes (6). Among closely related species, these symbionts can break through the species barrier by hybridization followed by selective sweep, resulting in identical mtDNA sequences among different species, which would cause the underestimation of the number of unique species under barcoding (9). Whereas these three processes may be relatively uncommon and limited to a small number of organisms, a fourth process, the nuclear integration of mtDNA that gives rise to nuclear mitochondrial pseudogenes (numts), is a widespread phenomenon that has been reported in many eukaryotic clades (8, 10). The effect of numts on DNA barcoding, however, has not been systematically studied to date.

The first case of numts in Metazoa was reported in the grasshopper *Locusta migratoria* (11), in which a copy of a mitochondrial ribosomal RNA gene was found in the nuclear genome. Lopez *et al.* (12) found that nearly half of the mitochondrial genome (7.9 kb) was transferred to the nuclear genome in the domestic cat and coined the term “numts.” Since then, >82 eukaryotes have been reported to have numts (8). A BLAST search of mitochondrial sequences in the published nuclear genomes suggests that nearly 99% of the mitochondrial sequences were transferred to different parts of the nucleus in both human and mouse (10). Pamilo *et al.* (13) reported >2,000 possible numts in the honey bee genome and found a similarly large number of numt copies in the flour beetle genome. These findings collectively indicate that numts are extremely pervasive in nature and that there may be a large number of species with unrealized numts of the COI gene in the nucleus.

The possible existence of COI numts poses a serious challenge to DNA barcoding. The fact that the COI gene can be amplified

Author contributions: H.S., J.E.B., M.F.W., and K.A.C. designed research; H.S. and J.E.B. performed research; H.S. and J.E.B. analyzed data; and H.S., J.E.B., M.F.W., and K.A.C. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

Data deposition: The sequences reported in this paper have been deposited in the GenBank database (accession nos. EU589049–EU589148, EU583504–EU583573, EU583577–EU583678, EF207161–EF207162, and EF207165–EF207168).

Freely available online through the PNAS open access option.

<sup>†</sup>To whom correspondence may be addressed: Department of Biology, 692 Widtsoe Building, Brigham Young University, Provo, UT 84602. E-mail: hojun\_song@byu.edu.

This article contains supporting information online at [www.pnas.org/cgi/content/full/0803076105/DCSupplemental](http://www.pnas.org/cgi/content/full/0803076105/DCSupplemental).

© 2008 by The National Academy of Sciences of the USA

from diverse taxa by using a limited set of primers is heralded as one of the attractive features of this marker (14). It is true that relatively conserved regions within mtDNA allow the design of “universal” primers, which can amplify mitochondrial fragments from an unknown species (15). However, conserved primers can be a double-edged sword when numts are present because they can coamplify numts in addition to the target mtDNA (7, 8). If the nuclear integration of numts was an ancient and sufficient sequence divergence accumulated in the orthologous mtDNA, the conserved primers would be more likely to amplify numts in preference to mtDNA, which could possibly result in unambiguous, paralogous sequences (8). Despite this serious problem, numts have been dismissed as a minor concern for DNA barcoding (16) and the issue of numts has not been adequately addressed.

In this study, we investigate the effect of including numts in DNA barcoding in two divergent lineages of arthropods, insects, and crustaceans, which are known to have especially large numbers of numts (8, 17–19). We also examine the effect of numts at different levels of divergence: subfamily-level (grasshoppers) and species- and population-level (crayfish). Herein, we show that both grasshopper and crayfish species included in the study have numts of the COI gene and barcoding methods would incorrectly infer that single individuals belong to multiple, unique species. The prevalence of numts appears to be both species-specific and population-specific and the pattern of numt distribution is considerably different between lower-level and higher-level divergence among taxa. Finally, we demonstrate the importance of data exploration in DNA barcoding practice by examining sequence characteristics of numts.

## Results and Discussion

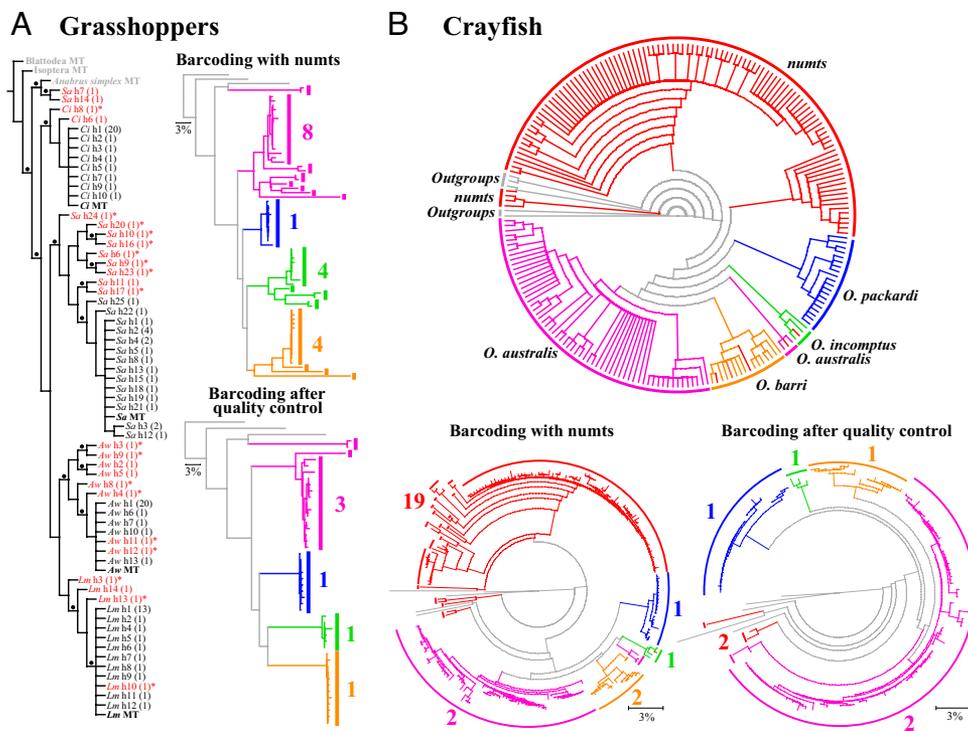
**Coamplification of Numts with Orthologous mtDNA.** Our results strongly suggest that a large number of paralogous haplotypes of various divergences are coamplified with the orthologous mtDNA sequences when conserved primers are used in both grasshoppers and crayfish, which can be identified by the presence of indels, point mutations, and in-frame stop codons [supporting information (SI) Table S1]. The majority of the coamplified paralogs can be easily considered nonfunctional numt haplotypes because of the presence of in-frame stop codons, which is especially evident in our crayfish data in which 97.3% of paralogs have stop codons. A large number of these numts have unusually high numbers of point mutations (mean = 65.52,  $n = 110$ ), suggesting that nuclear integration of mtDNA would result in random accumulation of nucleotide changes. In the grasshopper data, however, there are many paralogs that cannot readily be categorized as numts because they lack in-frame stop codons and differ from the orthologous mtDNA sequences by one or two nucleotides. If the same haplotype that appears to be functional other than the ortholog is repeatedly found from a single individual, one can suspect heteroplasmy (5). Indeed, heteroplasmy seems to explain the presence of certain paralogous haplotypes in *Schistocerca americana*. However, there are many paralogs represented by single haplotypes with small nucleotide differences in all four grasshopper species. Although, it is unlikely that these are *Taq* polymerase errors during PCR because of the high fidelity polymerase we used (0.015% error rate or 0.0732 bp per reaction); we cannot rule out the possibility of PCR error amplified by additional cloning (5, 18). Also, heteroplasmy might in fact be a plausible explanation for these haplotypes because we limited our study to only 30 clones per grasshopper species, thus not exploring the full extent of heteroplasmic diversity.

If the proportion of numts is high compared with the orthologous mtDNA fragments in a given PCR product, it is possible to generate unambiguous paralogous sequences (8). This is exacerbated when the conserved primers preferentially amplify

numts because of relatively ancestral sequence similarity of numts or divergence within the primer regions of the orthologous mtDNA. In this case, typical indicators of different PCR products, such as multiple bands on gels and double peaks, background noise, and ambiguity in sequence chromatograms, will not be present; hence, paralogous sequences can be mistaken as orthologous mtDNA. In fact, this exact phenomenon was observed in 18 crayfish individuals of *Orconectes barri* and *Orconectes australis* from which numts were amplified and cleanly sequenced without cloning.

**Phylogenetic Analyses and Distribution of Numts.** For the grasshopper data, the parsimony and the Bayesian analyses both recovered the monophyly of the orthologous mtDNA and haplotypes for three of four species (Fig. 1A); although, the topology was different in the placement of *S. americana* and *Calliptamus italicus* clades between the two methods. In each species, the largest clade was the polytomous clade consisting of the mtDNA ortholog and several similar haplotypes. The remaining haplotypes formed highly structured clades within each species, and this pattern was especially evident in *Acrida willemsei* and *S. americana*. For crayfish data, the Bayesian analysis recovered a topology mostly congruent with the parsimony analysis (Fig. 1B) and both analyses found a large clade of numt haplotypes (84 in parsimony and 82 in Bayesian) and a small clade of numts (18 in both analyses), distinctly divergent from the clearly defined clades of the orthologous mtDNA of four species. These numt clades consisted of the haplotypes from *O. australis*, *O. barri*, and *Orconectes packardi*, which were not necessarily grouped either by the species or the populations. Only four numt haplotypes were placed among orthologs (one in *Orconectes incomptus*, one in *O. australis*, and two in *O. barri*). A clade consisting of three numt haplotypes of *O. barri* was robustly placed near the root of trees in both analyses. Among the orthologous mtDNA clades, three of four species formed monophyletic clades, with the exception of *O. australis* that had one large clade sister to *O. barri* and a small clade basal to the *australis* + *barri* clade. Based on the phylogenetic analyses, number of indels, point mutations, in-frame stop codons, and sequence divergence, it is possible to conclude: among grasshoppers, *Locusta migratoria* has three numts, *A. willemsei* has six, *C. italicus* has two, and *S. americana* has at least 11 numts; among crayfish, *O. australis* has 60, *O. barri* has 46, *O. incomptus* has one, and *O. packardi* has four numts. On average, 32.54% and 41.88% of haplotypes generated from grasshoppers and crayfish, respectively, were numts (Table S1). It is important to note that this is a conservative estimate of the number of numts per species because it is limited by the number of individuals and clones we generated.

Both the grasshopper and crayfish data suggest that there can be multiple types of numts present within single individuals that vary considerably in nucleotide composition, suggesting multiple independent transfer events from the mitochondrial genome to the nucleus (8, 17). Moreover, our data suggest that these independent nuclear integration events can give rise to a family of numts that can diverge at different substitution rates. For example, we found that the relationships among the numt haplotypes of *S. americana* are highly structured and a similar pattern was observed in other grasshopper and crayfish species. Not only can independent transfer occur multiple times, but it can occur at very different phylogenetic levels. Whereas many numts of *S. americana* are closely related to the orthologous mtDNA, two numt haplotypes form a strong clade with the mtDNA ortholog of *Anabrus simplex*, which belongs to a different suborder within Orthoptera. Similarly, three numt haplotypes of *O. barri* were placed near outgroups belonging to different crayfish genera. These findings collectively suggest that there could have been an ancient nuclear integration event and that enough time has passed for these numts to have accumu-



**Fig. 1.** Phylogenetic and barcoding analyses based on orthologous mtDNA COI and paralogous numt haplotypes from grasshoppers and crayfish. (A) Grasshoppers: the cladogram on the left is a strict consensus of 41 MPTs ( $L = 1002$ ;  $CI = 0.54$ ;  $RI = 0.85$ ). Dots above branch indicate the nodes with the bootstrap value of  $>75$  and posterior probability of  $>95\%$ . Orthologous mtDNA is indicated in bold and putative numts are indicated as red terminals. Number in parenthesis represents the number of identical copies for a particular haplotypes (h) and asterisk indicates ones with in-frame stop codons. When DNA barcoding analysis (NJ analysis based on K2P distances) is performed on the complete dataset, the number of unique species inferred based on 3% sequence divergence (colored numbers next to the vertical bars) is overestimated (barcoding with numts). After the removal of the haplotypes with indels and in-frame stop codons (barcoding after quality control), the number of unique species inferred under DNA barcoding is drastically reduced. Purple, *Schistocerca americana* (Sa); blue, *Calliptamus italicus* (Ci); green, *Acrida willemsi* (Aw); orange, *Locusta migratoria* (Lm); and gray, outgroups. (B) Crayfish: the circular cladogram on top is the strict consensus of 94 MPTs ( $L = 1064$ ;  $CI = 0.39$ ;  $RI = 0.91$ ). Terminals are colored to indicate species. Purple, *Orconectes australis*; orange, *O. barri*; green, *O. incomptus*; blue, *O. packardii*; and gray, outgroups. All numt haplotypes are indicated as red terminals. Similarly, DNA barcoding overestimates the number of unique species when numts are included, but the removal of numts reduces the inferred number of species. Notice that even after rigorous quality control, the inferred number of unique species is actually higher than the actual number of species, suggesting that some numts are difficult to identify.

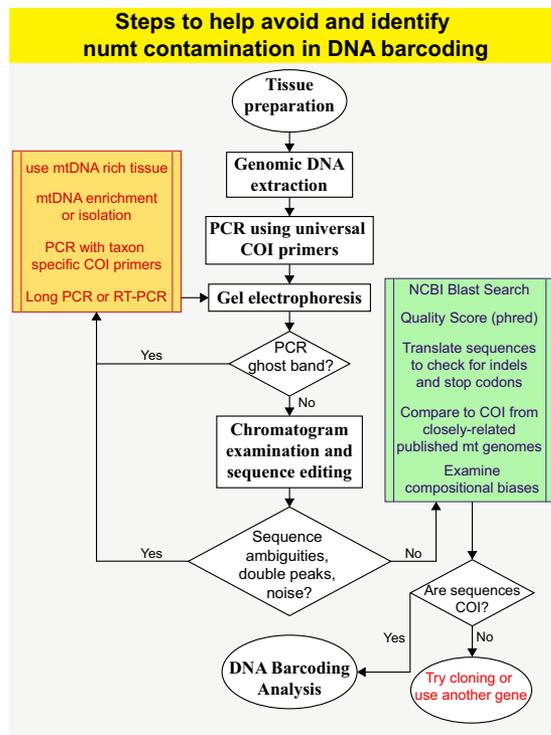
lated substantial sequence divergence from the orthologous mtDNA.

Our two datasets differ at the level of divergence among the ingroup species along with the distribution of numts deduced from the phylogenetic analyses. Except for two numt haplotypes of *S. americana* that grouped within an outgroup, all grasshopper numt haplotypes strongly grouped with their orthologous mtDNA. However, this pattern is not observed in the four closely related crayfish species. Only a small portion of numt haplotypes grouped with their orthologs, whereas the majority formed clades among themselves with no apparent population or species-specific groupings. In other words, numt haplotypes sequenced from different crayfish individuals from different populations and different species form monophyletic groups. This finding implies that the closely related crayfish species share similar types of numts that must predate the speciation events. Both patterns have been reported from other studies looking at various mitochondrial genes in diverse metazoan lineages at different phylogenetic levels (18, 20). We conclude that the distribution pattern of numts within a given group of organisms cannot be predicted *a priori*, but depends on the timing and the frequency of nuclear integration, which can clearly predate and postdate speciation events.

The distribution of numts in both datasets suggests that their prevalence may be lineage-specific. For example, we sequenced numts from individuals collected from 11 of 56 cave crayfish populations (Tables S2 and S3). Among the four species, numts

were sequenced for 7 of 34 localities in *O. australis* (southernmost region of range, primarily caves in Alabama), 2 of 7 in *O. barri* (southernmost region of range, caves in Tennessee), 1 of 3 in *O. incomptus*, and 1 of 12 in *O. packardii*. This observed pattern does not necessarily mean that the remaining 45 populations are free of numts but it does mean that there is a nonrandom population-specific variation in the level of numt prevalence. Nuclear integration of mtDNA happens at the level of the individual (8) and a large population size can effectively dilute the amount of numts in a given population. In this case, the proportion of numts is much smaller than that of the orthologous mtDNA in a given individual, rendering the numt coamplification less likely, even with a possibility of the presence of plesiomorphic numts in the nucleus. However, whether a population experiences genetic drift because of an extreme bottleneck, numts can be fixed in a few founders, resulting in a disproportionately high level of numts. Another intrinsic factor, nuclear genome size, might also play a role that results in uneven distribution of numts. Whereas all grasshopper species have numts, *A. willemsi* and *S. americana* have especially high numbers that are divergent from each other. Bensasson *et al.* (21) suggested that a positive correlation between the number of numts and nuclear genome size might exist and it is possible that these two species might have larger nuclear genomes than the others. Alternatively, inherent species-specific differences in the frequency of DNA transfer from mitochondria to the nucleus and in the rate of loss of numts in the nucleus have also been





**Fig. 2.** Suggested steps to help avoid and identify numts in DNA barcoding analysis. Whereas these steps will help reduce the chance of sequencing numts instead of the target COI, they are not guaranteed to remove all numts. Each resulting sequence must be examined as part of quality control protocols. If numts are rampant, then the isolation of COI sequences becomes difficult and it may be best to use other genes. When interpreting the results from DNA barcoding analysis, it is important to survey congruence with other molecular markers, morphology, ecology, and behavior.

should employ when using mtDNA for barcoding studies (Fig. 2). However, our suggested quality control measures against numts are neither simple nor rapid, which is at odds with the goal of DNA barcoding initiatives. DNA barcoding is a tool to aid rapid biological identification, which should be used in conjunction with other information including morphology, behavior, and ecology, and the use of other information will help reduce incorrect molecular inferences.

**Concluding Remarks.** The possible coamplification of numts is clearly a major impediment to DNA barcoding. To be fair, this is a problem that all PCR-based studies face, including phylogeography and phylogenetic studies using mtDNA. This is why both fields have largely rejected sole reliance on a single marker and emphasized congruence among multiple markers. The problem is exacerbated because the variation in the prevalence of numts appears to be a widespread phenomenon. Richly and Leister (10) surveyed numts in sequenced eukaryotic genomes and found that the number of numts ranges from none in the mosquito *Anopheles gambiae* to >500 in human. Although there are little or no reported cases of numts in groups such as flies (10), chicken (26), and fishes (10), a large number of eukaryotic clades including plants (29), birds (30), nonavian reptiles (31), mammals (12, 20), and arthropods (8, 11, 13, 17–19) were shown to have numts. In our study, the variation is not only clade-specific, but also species-specific and population-specific. From a barcoding perspective, the presence of numts can be disastrous. Because the DNA barcoding initiative attempts to barcode all life forms, including both organisms with known numts and other organisms

that potentially have numts, this issue cannot simply be ignored. Otherwise, the number of single individuals that are inferred to be multiple species because of numt contamination may become the legacy of the DNA barcode movement.

## Materials and Methods

**Taxon Sampling.** To study the evolution and distribution of numts at higher-level divergence, we included four grasshopper species belonging to four different subfamilies of Acrididae. To establish the orthology of mtDNA, we used the taxa whose partial or complete mitochondrial genomes have been sequenced: *Acrida willemsei* (Acridinae, EU589053), *Calliptamus italicus* (Calliptaminae, EU589054), *Locusta migratoria* (Oedipodinae, EU589051), and *Schistocerca americana* (Cyrtacanthacridinae, EU589055). We generated numts from single individuals per species and used the same individuals that the complete mitochondrial genomes were sequenced from with an exception of *L. migratoria*. As outgroups, we used the COI regions of a Mormon cricket *Anabrus simplex* (EU589052), a cockroach *Gromphadorhina portentosa* (EU589049), and a termite *Mastotermes darwinensis* (EU589050). To study the numts at population- and species-level divergence, we included a total of 119 individuals of four closely related species belonging to the cave crayfish genus *Orconectes*, collected from 56 localities along the Cumberland Plateau of the Southern Appalachians: *O. australis*, *O. barri*, *O. incomptus*, and *O. packardii*. As outgroups, we included *O. limosus* (AF517105), *Procambarus simulans* (EU583575), and three species of the genus *Cambarus* (*C. gentryi* [DQ411785], *C. tenebrosus* [EU583576], and *C. bartonii* [EU583574]). COI from the complete mtDNA genome of *Cherax destructor* (NC.011243) was used for reference. GenBank accession numbers for the haplotypes are EU589057–EU589148 (grasshoppers) and EU583504–EU583573, EU583577–EU583678, EF207161–EF207162, and EF207165–EF207168 (crayfish). Details about numt amplification can be found in *SI Methods*.

**Characterization of Numts.** To ensure the quality and identity, each haplotype was blasted by using MegaBLAST option against the nucleotide collection (nr/nt) as implemented in the National Center for Biotechnology Information website (<http://www.ncbi.nlm.nih.gov/blast/Blast.cgi>). Only the haplotypes that had high similarity to the COI sequence were used for the further analyses. For example, the blast search revealed that eight cloned sequences (3 from *L. migratoria* and 5 from *A. willemsei*) were of either bacterial or unknown origins and these sequences were treated as cloning error and removed from further analyses. We characterized the haplotypes in Sequencher 4.6 for length, nucleotide composition, and number of in-frame stop codons. Putative indels and point mutations were estimated by comparing the haplotype sequences against the mtDNA orthologs. The number of unique haplotypes was determined and the sequence divergence from the orthologous mtDNA sequence for each species was calculated under Kimura 2-parameter (K2P) model in MEGA 3.1 (32), as routinely used in barcoding studies, despite this model being under fit relative to the data (see below).

**Data Analysis.** To study divergence pattern of numts, we performed phylogenetic analyses in both parsimony and Bayesian frameworks. For both grasshopper and crayfish datasets, the unique haplotypes and the mtDNA orthologs were aligned in MUSCLE (33) by using default parameters. For grasshopper data, we created a matrix of 69 terminals (7 mtDNA orthologs and 62 unique haplotypes) and 475 aligned nucleotides. For crayfish data, we created a matrix of 215 terminals (5 outgroups and 210 unique haplotypes of four species) and 663 aligned nucleotides. Within the parsimony framework, the aligned sequence data were analyzed with gaps treated as missing, by using search algorithms implemented in TNT ([www.zmuc.dk/public/phylogeny](http://www.zmuc.dk/public/phylogeny)). To assess support, we calculated standard bootstrap values based on 1,000 replicates (100 random-addition TBR replicates each) and Bremer support values, both in TNT. Within the Bayesian framework, we analyzed the datasets by using the program MrBayes 3.1 (34) after selecting best-fit models of nucleotide evolution under the AIC criteria by using MrModeltest 2.2 (program distributed by J.A.A. Nylander, Evolutionary Biology Centre, Uppsala University). The analyses consisted of running four simultaneous chains for 20 million generations for grasshopper data (GTR+G) and six simultaneous chains for 30 million generations for crayfish data (HKY+I+G), both sampling every 1,000 generations. Four independent identical Bayesian runs were performed to ensure convergence on similar results and the nodal support was assessed by using the posterior probability generated from a consensus tree of the sampled trees past burn-in determined by using Tracer 1.4 (<http://beast.bio.ed.ac.uk>).

To study how the presence of numts might influence the inferences from DNA barcoding, we performed a neighbor-joining (NJ) analysis under K2P

