ELSEVIER

# Relationship between morphological taxonomy and molecular divergence within Crustacea: Proposal of a molecular threshold to help species delimitation

T. Lefébure [a,*], C.J. Douady [a], M. Gouy [b], J. Gibert [a]

[a] *Laboratoire d'Ecologie des Hydrosystèmes Fluviaux, UMR-CNRS 5023, Université Claude Bernard Lyon 1, F-69622 Villeurbanne Cedex, France*
[b] *Laboratoire de Biométrie et Biologie Evolutive, UMR-CNRS 5558, Université Claude Bernard Lyon 1, F-69622 Villeurbanne Cedex, France*

## Abstract

With today's technology for production of molecular sequences, DNA taxonomy and barcoding arose as a new tool for evolutionary biology and ecology. However, their validities still need to be empirically evaluated. Of most importance is the strength of the correlation between morphological taxonomy and molecular divergence and the possibility to define some molecular thresholds. Here, we report measurements of this correlation for two mitochondrial genes (COI and 16S rRNA) within the sub-phylum Crustacea. Perl scripts were developed to ensure objectivity, reproducibility, and exhaustiveness of our tests. Our analysis reveals a general correlation between molecular divergence and taxonomy. This correlation is particularly high for shallow taxonomic levels allowing us to propose a COI universal crustacean threshold to help species delimitation. At higher taxonomic levels this correlation decreases, particularly when comparing different families. Those results plead for DNA use in taxonomy and suggest an operational method to help crustacean species delimitation that is linked to the phylogenetic species definition. This pragmatic tool is expected to fine tune the present classification, and not, as some would have believed, to tear it apart.
© 2006 Elsevier Inc. All rights reserved.

*Keywords:* DNA taxonomy; DNA barcoding; COI; 16S; Crustacea; Cryptic species

## 1. Introduction

Species level is recognized as the major unit of biodiversity (Claridge et al., 1997). Nevertheless many species definitions exist (22 following Mayden, 1997) and there are no standardized operational criteria to delimit them (Sites and Marshall, 2004). Alpha taxonomy, and so biodiversity assessment, remains today mainly based on morphological characters. Morphology is a complex and non-neutral marker. Consequently, morphological taxonomy could lead to under- or over-estimation of biodiversity. With today's technology for production of molecular sequences, cryptic speciations have been reported from all phyla and seem to be a frequent bias associated with morphological taxonomy (i.e., in 2004 more than 200 publications reported cryptic diversity).

Some authors (e.g., Blaxter, 2004; Hebert et al., 2003; Tautz et al., 2003) suggested to use DNA in taxonomy to overcome those "impediments" and to face with the enormous quantity of living forms that remain to be classified whereas the number of taxonomists is declining. This approach has been strongly criticized (e.g., Lipscomb et al., 2003; Mallet and Willmott, 2003; Seberg et al., 2003; Will and Rubinoff, 2004). Nevertheless, almost all debates were developed on theoretical grounds while utility and consequences of the use of DNA in taxonomy have not been tested on broad datasets. Without entering in too many details, it is now rather clear that part of the conflicts between proposals and criticisms were linked to a problem

---

of definition. Under the terminology DNA taxonomy (or DNA barcoding) two quite independent tasks have been merged (e.g., DeSalle et al., 2005): (1) identify and assign specimens to taxonomic groups (e.g., species, families) that have been previously described (e.g., Savolainen et al., 2005), and (2) predict and classify new taxa using DNA. In this paper we will use the term DNA barcoding (Hebert et al., 2003; see www.barcodinglife.org/) when referring to the first goal (identification) and DNA taxonomy when referring to the second (prediction and classification). Theoretical aspects (DeSalle et al., 2005; Savolainen et al., 2005), methods (Blaxter et al., 2005; Steinke et al., 2005), and applied cases (e.g., Chase et al., 2005; Monaghan et al., 2005) of the DNA barcoding are today under quite an intense development. DNA barcoding goal is partly independent from the way the classification has been built. It "only" requires to define for each taxonomic group a set (possibly more than one if the taxon is para or polyphyletic and is therefore composed of more than one monophyletic unit) of molecular synapomorphies that could be used as taxonomic tags. In contrast, the second goal (classification) expects the molecular and current taxonomies to be congruent. In other words this second goal requires that a particular range of molecular divergence of a particular gene could be assigned to a particular taxonomic rank.

The relationship between taxonomy and molecular divergence has been already studied by Avise and collaborators. Avise and Johns (1999) demonstrated that comparable taxonomic ranks between animal phyla were not at all equivalent in molecular divergence. Johns and Avise (1998) comparing inter-specific divergence of Cytochrome *b* for vertebrates found poor equivalence of divergence across taxa. Nevertheless Avise and Walker (1999), always with Cytochrome *b*, within the vertebrates and for shallow taxonomic rank (species rank) found good relationship between taxonomic classification and molecular divergence, and concluded that "mtDNA and traditional taxonomic assignments tend to converge on what therefore may be real biotic units in nature". More recently, Hebert et al. (2004) proposed a standard sequence threshold of 10 times the mean intra-specific variations to delimit animal species. Nevertheless, methods and particularly sampling strategies of both approaches have been criticized (Hendry et al., 2000; Moritz and Cicero, 2004), and the degree of correlation between molecular divergence and taxonomy remains unclear.

The aims of the present study are to formally test the correlation between taxonomic ranks and molecular divergences and, if possible, to define molecular thresholds to help taxonomic decision. This test was developed for two mitochondrial genes, the COI [1] and 16S rRNA, on the highly diversified crustacean sub-phylum. Cryptic species are common in crustaceans (Burton and Lee, 1994;

Daniels et al., 2003; de Bruyn et al., 2004; Edmands, 2001; Jarman and Elliott, 2000; King and Hanner, 1998; Lee, 2000; Mathews et al., 2002; Müller, 2000; Penton et al., 2004; Rawson et al., 2003; Rocha-Olivares et al., 2001; Wares, 2001; Williams et al., 2001; Witt and Hebert, 2000). Crustaceans are also particularly abundant in extreme habitats, like subsurface water, where extreme conditions seem responsible of morphological convergences leading to important biodiversity under-estimation (Lefébure et al., in press; Proudlove and Wood, 2003). For these reasons crustaceans constitute a group for which DNA taxonomy or at least total evidence taxonomy (i.e., morphological plus molecular) could be highly valuable. Furthermore, most previous studies have focused on vertebrates (Avise and Walker, 1999; Hebert et al., 2004; Johns and Avise, 1998) which taxonomy has been intensively studied and for which "taxonomic impediments" are probably fewer than in other phyla. In this way vertebrates, and particularly vertebrates of the Northern hemisphere, may constitute a biased test (Harris and Froufe, 2005) and also perhaps a group for which DNA taxonomy would be less valuable. Unfortunately, invertebrate groups have also been far less sequenced, and only two genes (the mitochondrial COI and 16S) are today sufficiently sampled to be analyzed.

To ensure test objectivity, reproducibility, and exhaustiveness, we developed custom made Perl scripts (practical extraction and report language; http://www.perl.org/). All sequences available in public sequence databases were used, and only sequences meeting a priori defined criteria of length, position, similarity, and taxonomy were analyzed. Molecular divergences were compensated for multiple substitutions. Our results demonstrate that the current taxonomy is in global agreement with molecular differentiation, but that this relation is degrading at high taxonomic levels. We show that COI DNA variations are suitable to delimit what crustacean taxonomists have called species. A COI molecular threshold is therefore proposed to help the delimitation of new crustacean species.

## 2. Materials and methods

### 2.1. Sequence retrieval

All crustacean DNA sequences were extracted from GenBank on the 7th of June 2004 (representing a total of 11,885 sequences, to the exclusion of EST, STS, GSS, TPA, working draft, and patented sequences). A BLAST database was then built using formatdb 2.2.8 (NCBI BLAST package, ftp://ftp.ncbi.nlm.nih.gov/blast/). Complete COI and 16S sequences were extracted from five complete mitochondrial genome sequences of Cladocera (GenBank Accession No. NC_000844), Branchiopoda (GenBank Accession No. AB084514), Copepoda (GenBank Accession No. NC_003979), Decapoda (GenBank Accession No. AF150756), and Ostracoda (GenBank Accession No. AB114300) to represent an overall crustacean diversity. For

---

[1] *Abbreviations used:* S, intra species divergences; G, inter-species but intra-genera divergences; F, inter-genera but intra-family divergences; COI, cytochrome oxidase subunit I gene.

the COI and 16S, five discontiguous megaBLAST (Zhang et al., 2000) were performed using the program megablast 2.2.8 (NCBI BLAST package, ftp://ftp.ncbi.nlm.nih.gov/blast/) with as seed one of the five complete sequences and the following parameters: word size = 11, discontiguous word template length = 16, discontiguous template type = 0 (coding) for the COI and = 1 (non-coding) for the 16S, maximum number of reported sequences = 5000, and cutoff expectation value = 0.05. The sequences found by each seed were then joined and multiple occurrences removed (3037 sequences for the COI and 2181 for the 16S). In parallel, for each sequence the strand orientation (3′–>5′ or 5′–>3′) was collected from the BLAST output. Sequences shorter than 300 bp were removed and all sequences were put in the same orientation. Complete mitochondrial genome sequences were also excluded as they added little more information (25 sequences) but drastically slowed down alignment processes. Comparing to a traditional database querying method by keywords, this BLAST approach recovers approximately the same amount of sequences but ensures homology of the fragments.

### 2.2. Taxonomy

The taxonomy database of NCBI (version of the 02/07/2004) was used to associate family, genus, and species names with each sequence. The database was modified to give to the Porcellanidae the status of family (Martin and Davis, 2001). The sequences selected by BLAST were then split by family and, for statistical purposes, only families containing at least 50 sequences were analyzed (18 families for the COI and 14 for the 16S).

### 2.3. Family alignments and design of homologous blocks

#### 2.3.1. Initial family alignments

For each gene and family, individual sequences were aligned to the complete sequence from a close representative (same family when available). Pairwise alignments were performed with ClustalW 1.82 with default parameters (Thompson et al., 1994). Sequences with less than 50% similarity between the fragment and the complete sequence or with numerous stop codons or gaps for the COI were removed of the analysis. Those cases were rare and generally corresponded to COI pseudogenes (e.g., Williams and Knowlton, 2001). A complete family, the Hyalellidae, was also removed from the COI dataset as it contained numerous sequences belonging to undetermined species. Pairwise alignments were then successively aligned to each other using the profile option of ClustalW and complete sequences removed from the profile alignment. This produced for each family a first and global alignment between potentially non- or weakly overlapping sequences.

#### 2.3.2. Design of homologous blocks

Most sequences in the present dataset are partial COI and 16S fragments and cover regions that vary extensively. Gene regions with large taxonomic sampling were identified as follows. Each family alignment was profile aligned against the complete mitochondrial sequence of *Penaeus monodon* (GenBank Accession No. AF217843). Then, this sequence was used as reference to determine the first and last positions of our initial alignments. After plotting the positions (Supplementary material), sequence blocks were designed to obtain the best compromise between block length, number of sequences within blocks, and number of undetermined sites. Two blocks were defined in the COI gene between *Penaeus monodon*'s positions 100 and 580, and between positions 720 and 1260, and one block in the 16S rRNA gene at positions 760–1220. For statistical purposes, sequences shorter than 80% of the block length for the COI and 70% for the 16S (due to numerous deletions within some families), were removed. For the same reasons, family blocks containing less than 30 sequences were removed.

#### 2.3.3. Refined family alignments

The first alignment step was used to reference the first and last position of each fragment between possibly non-overlapping fragments, but produced alignments of poor quality. Therefore, within each resulting block and family

Table 1
Crustacean families analyzed

| Gene | Position | Family | Groups | Gen. | Sp. | Seq. |
|------|----------|--------|--------|------|-----|------|
| COI | 100–580 | Gammaridae | Amphipoda (O) | 8 | 11 | 57 |
| COI | 100–580 | Balanidae | Cirripedia (infCl) | 2 | 2 | 174 |
| COI | 100–580 | Parastacidae | Decapoda (O) | 4 | 17 | 40 |
| COI | 100–580 | Coronulidae | Cirripedia (infCl) | 1 | 2 | 80 |
| COI | 100–580 | Chthamalidae | Cirripedia (infCl) | 5 | 17 | 32 |
| COI | 100–580 | Harpacticidae | Copepoda (subCl) | 1 | 2 | 61 |
| COI | 100–580 | Asellidae | Isopoda (O) | 2 | 2 | 73 |
| COI | 100–580 | Paguridae | Decapoda (O) | 1 | 6 | 65 |
| COI | 100–580 | Daphniidae | Cladocera (subO) | 2 | 40 | 180 |
| COI | 100–580 | Idoteidae | Isopoda (O) | 3 | 7 | 50 |
| COI | 100–580 | Atyidae | Decapoda (O) | 2 | 2 | 40 |
| COI | 720–1260 | Aeglidae | Decapoda (O) | 1 | 55 | 167 |
| COI | 720–1260 | Alpheidae | Decapoda (O) | 3 | 64 | 270 |
| COI | 720–1260 | Parastacidae | Decapoda (O) | 2 | 4 | 41 |
| COI | 720–1260 | Penaeidae | Decapoda (O) | 12 | 27 | 104 |
| COI | 720–1260 | Portunidae | Decapoda (O) | 3 | 15 | 82 |
| COI | 720–1260 | Atyidae | Decapoda (O) | 2 | 3 | 48 |
| 16S | 760–1220 | Potamonautidae | Decapoda (O) | 1 | 14 | 29 |
| 16S | 760–1220 | Aeglidae | Decapoda (O) | 1 | 55 | 168 |
| 16S | 760–1220 | Bosminidae | Cladocera (subO) | 3 | 14 | 63 |
| 16S | 760–1220 | Alpheidae | Decapoda (O) | 3 | 33 | 70 |
| 16S | 760–1220 | Gammaridae | Amphipoda (O) | 2 | 8 | 57 |
| 16S | 760–1220 | Parastacidae | Decapoda (O) | 15 | 79 | 187 |
| 16S | 760–1220 | Penaeidae | Decapoda (O) | 12 | 39 | 121 |
| 16S | 760–1220 | Palaemonidae | Decapoda (O) | 2 | 21 | 61 |
| 16S | 760–1220 | Astacidae | Decapoda (O) | 3 | 5 | 57 |
| 16S | 760–1220 | Cambaridae | Decapoda (O) | 5 | 12 | 60 |
| 16S | 760–1220 | Ocypodidae | Decapoda (O) | 17 | 51 | 75 |
| 16S | 760–1220 | Varunidae | Decapoda (O) | 16 | 32 | 42 |
| 16S | 760–1220 | Porcellanidae | Decapoda (O) | 4 | 40 | 51 |
| 16S | 760–1220 | Cirolanidae | Isopoda (O) | 8 | 13 | 62 |

Gen., number of genera analyzed; sp., number of species; Seq., number of sequences; O, order; infCl, infraclass; and subCl, subclass.

(11 $COI_{100-580}$, 6 $COI_{720-1260}$, and 14 $16S_{760-1220}$, Table 1), we performed multiple alignments using ClustalW with default parameter values. There was no positional homology ambiguity for the COI, but some for the 16S. To quantify the level of homology ambiguity, we used the sensitivity approach implemented in Soap 1.1b2 (Löytynoja and Milinkovitch, 2001). Soap parameters were set as follows: gap penalties from 11 to 19, by steps of 2; extension penalties from 3 to 11, also by steps of 2; and 100% conservation.

### 2.4. Molecular divergences

Two sets of molecular divergences were calculated: patristic and pairwise distances. Patristic distances between two taxa is defined as the amount of divergence since they shared a common ancestor, i.e., the path-length distance between the two taxa along the tree. This approach requires to build a phylogeny of the studied taxa and subsequently to extract divergences from the reconstructed tree. To the opposite, pairwise distances only compare taxa by pair and give the observed number of differences that may be corrected or not following a model of molecular evolution. This last approach has the benefit of being fast, but it does not account for phylogenetic relationships between taxa and is thereby more sensible to bias such as multiple substitutions.

Within each family, Bayesian inferences were performed with MrBayes v3.0B4 (Huelsenbeck and Ronquist, 2001) using a GTR+G+I model of evolution. The tree-space was explored by using four chains over at least 1,000,000 generations sampled every 100. Burn-in value was fixed at 10% the total generation number after empirical determination of the convergence. When 1,000,000 generations were not sufficient, 1,500,000–2,000,000 generations were performed. Branch lengths of the most probable topologies were then estimated using a maximum likelihood criterion under a GTR+G+I model of evolution with PAUP* 4b10 (Swofford, 2002). From these trees, patristic distances between each taxa were extracted using APE library (Paradis et al., 2004) of R 2.0.0 (R Development Core Team, 2004). In parallel, four sets of pairwise distances, with relevant parameters estimated by maximum likelihood on the most probable topology of MrBayes, were computed with PAUP*: (1) GTR+G+I distances, (2) K2p+G distances, (3) K2p distances, and (4) uncorrected *p* distance.

COI sequences were also translated into amino-acid sequences using the invertebrate mitochondrial genetic code, and branch lengths of the Bayesian DNA trees were re-estimated under a JTT+G amino-acid model of evolution with Tree-Puzzle 5.2 (Schmidt et al., 2002). Patristic amino-acid distances between all taxa were then extracted following the above described procedure.

### 2.5. Sorting, distribution, and overlap of molecular divergences

Patristic distances were sorted in three categories: intra-specific distances (S), inter-species but intra-generic dis-

tances (G), and inter-generic but intra-familial distances (F). Next, distribution of these categories of distances was plotted by family using the boxplot representation of R. Boxplots (Tuckey, 1977) represent the overall shape of the dataset. It describes median (central bar), position of the upper and lower quartiles (called Q1 and Q3, central box), extremes of the data ("whiskers") and very extreme points of the distribution that can be considered as outliers (dots). Points are considered as outliers when they exceed $Q3 + 1.5$ IQR for the upper part of the distribution or $Q1 - 1.5$ IQR for the lower part, where IQR is the inter quartile range (i.e., Q3–Q1). Although distance distributions within families are not independent from each other, we performed Mann–Whitney tests between S, G, and F distributions to obtain a first statistical indication of the overlap between divergence distributions. Finally, we developed an assumption-free statistical approach to directly measure the overlap between our distributions and possibly locate the "best" threshold delimiting two distributions. This method consists in determining the percentage of samples of each distribution that are below (for the first distribution) or above (for the second distribution) a range of thresholds. This percentage is then considered to be the chance of success for each threshold to discriminate samples from a distribution. The best threshold—in fact the best compromise— is found where both success curves cross. The performance of this threshold is finally represented by the percentage of any samples that would have been correctly sorted using it. Two textbook cases are given in Fig. 1. Two fully overlapping or very close distributions will lead to a success between 50 and 60%. Overlapping but nonetheless differentiated distributions will produce a success between 60 and 80%. Then, weakly overlapping distributions will lead to a success between 80 and 90%. Finally, very different or entirely disjoint distributions will produce a success superior to 95%.

### 2.6. Artefacts and taxonomic bias

From a theoretical point of view, two main factors may bias our divergence assessment. First, a strong disequilibrium in the representation of some taxa could bend divergence distribution. To test for a possible artefact due to biased species representation (some species being represented by more than 100 sequences and others by a single one), we performed a second set of analyses where each taxa was given the same weight by computing mean divergence (i.e., mean S divergence per species, mean G divergence per couple of species, and mean F divergence per couple of genera). Second, the taxonomic classification may be incorrect or uncertain. Most common problems will result from (1) cryptic species, (2) taxa with multiple denominations or taxonomic ranks, and (3) paraphyletic or polyphyletic taxa. The significance of sibling species or multiple taxonomic designations bias was evaluated by reading published papers associated with sequences used in this analysis. Then all the recognized cryptic species and taxa
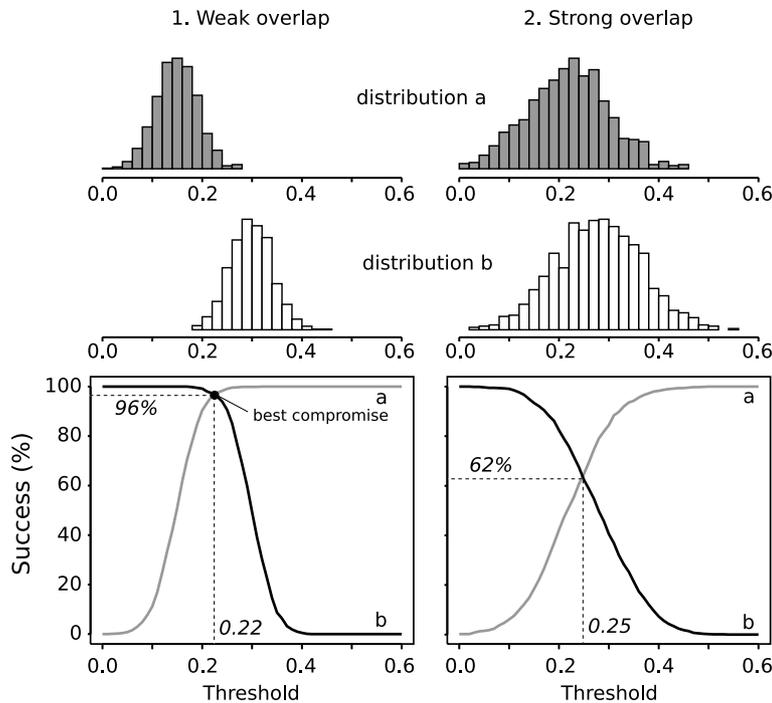
Fig. 1. Presentation of the proposed method to quantify the overlap between two distributions and to determine the best threshold value to discriminate them. The method is here applied for two textbook cases: a weak (left column) and a strong (right column) overlap. The method consists in determining the percentage of samples (*y*-axis) of each distribution that are below (for the distribution (a)) or above (for the distribution (b)) a range of threshold (*x*-axis). This percentage is then considered to be the chance of success for each threshold to discriminate (a) and (b) samples. The best threshold (compromise) is found where (a) and (b) success curves intersect. The performance of this threshold is finally obtained by the percentage of (a) and (b) samples that would have been correctly sorted using it. Here using a threshold of 0.22 for the weak overlap situation (left column), we would have successfully discriminated 96% of (a) and (b) samples, but a threshold of 0.25 for the strong overlap case would only correctly sort our samples 62 times out of 100.

with doubtful taxonomy were removed. The influence of presumably non-monophyletic taxa on the divergence distribution was finally tested by discarding taxa that appeared non-monophyletic in our Bayesian trees. Of course as the gene tree may not represent the species tree—because of unresolved topology, incomplete lineage sorting … (Funk and Omland, 2003)—this method will probably discard real monophyletic taxa. Nevertheless, this situation is supposed to be rare and furthermore should not bias our test.

## 2.7. Scripts and datasets

To ensure reproducibility all previously described steps were automatized through Perl and R scripts. Scripts and alignments are available upon request (lefebure@univ-lyon1.fr).

## 3. Results

### 3.1. Taxonomic range

The bioinformatic procedure developed led to the examination of two blocks of COI and one of 16S (Supplementary material). They are composed of 11 different families for the block $COI_{100–580}$, six families for $COI_{720–1260}$, and 14 families for the $16S_{760–1220}$, (Table 1). Most of these families

are decapods (20 out of 31). Nevertheless and especially for the block $COI_{100–580}$, a wide diversity of Crustacea are represented (e.g., Cirripedia, Copepoda, Amphipoda, Isopoda, and Cladocera). The sampling effort between families is variable, ranging from high diversity at the genus and species levels (e.g., Parastacidae within the $16S_{760–1220}$ represented by 15 genus and 79 species), to poor diversity at both levels (e.g., Coronulidae within the $COI_{100–580}$ represented by one genus and two species).

### 3.2. Intra-familial divergences

#### 3.2.1. COI DNA divergences

At the family level no alignment ambiguity was detected. Patristic divergences based on DNA sequences for $COI_{100–580}$ and $COI_{720–1260}$ range from 0 to 2.8 substitutions per site (Fig. 2). Within families, F distances are globally higher than G distances, which are globally higher than S distances (Table 2, all Mann–Whitney tests were highly significant, *p* value $<10^{-6}$). F and G distances frequently overlap (families Chthamalidae $COI_{100–580}$, Daphniidae $COI_{100–580}$, Alpheidae $COI_{720–1260}$, and Penaeidae $COI_{720–1260}$). Overlaps between S and G distances appeared less frequent with the notable exception of the Aeglidae ($COI_{720–1260}$, Fig. 2). However, the magnitude of this intra-family trend seems quite different between families: each family apparently has its own range of differentiation. Most S distances are below
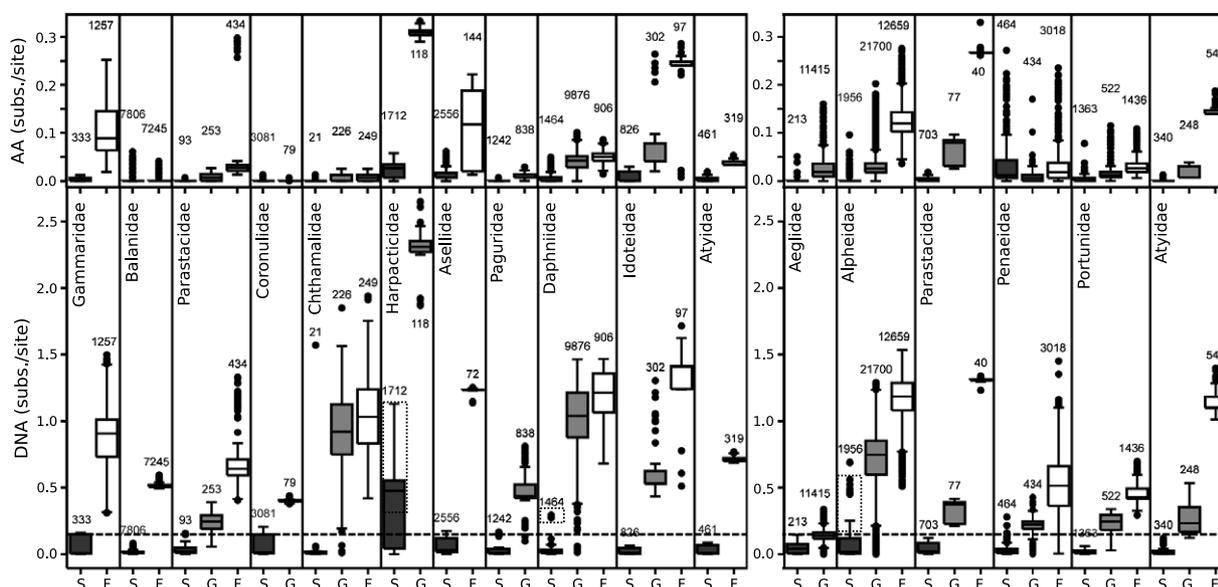
Fig. 2. Boxplot distribution of intra-species (S, in dark gray), inter-species but intra-genus (G, light gray), and inter-genera but intra-family (F, white) DNA (lower panel) or amino-acid (upper panel) divergences of 5′ (positions 100–580, left part) and 3′ (positions 720–1260, right part) ends of the COI for different families of crustaceans. Divergences are patristic distances measured on the most probable Bayesian trees with branch lengths re-estimated by maximum likelihood under a GTR+G+I model of evolution. The dashed line identifies a threshold of 0.15 subst./site. Hidden diversity described in the literature is delimited by dotted rectangles. Numbers above boxplots indicate the number of pairwise distances. The number of sequences, species, and genera per family is indicated in Table 1. Characteristics of boxplot representations are described in the Material and methods section.

Table 2
Patristic followed by pairwise k2p divergences within species (S), between species of the same genus (G), and between genus of the same family (F) of 5′ (positions 100–580) and 3′ (positions 720–1260, right part) ends of the COI, and 3′ ends of the 16S (positions 760–1220)

|                      |                | S | G | F |
|----------------------|----------------|--------------|--------------|--------------|
| $COI_{100-580}$      | Extreme lower  | 0–0          | 0.258–0.154  | 0.310–0.164  |
|                      | Median         | 0.015–0.013  | 1.016–0.251  | 0.520–0.203  |
|                      | Extreme upper  | 0.094–0.079  | 1.564–0.333  | 0.983–0.261  |
| $COI_{720-1260}$     | Extreme lower  | 0–0          | 0–0          | 0.085–0.134  |
|                      | Median         | 0.016–0.017  | 0.569–0.196  | 1.115–0.252  |
|                      | Extreme upper  | 0.079–0.064  | 1.290–0.320  | 1.534–0.361  |
| $16S_{760-1220}$     | Extreme lower  | 0–0          | 0–0          | 0–0.007      |
|                      | Median         | 0.026–0.021  | 0.069–0.037  | 0.621–0.222  |
|                      | Extreme upper  | 0.133–0.104  | 0.402–0.230  | 1.413–0.436  |

Patristic distances were measured on the most probable Bayesian trees with branch lengths re-estimated by maximum likelihood under a GTR+G+I model of evolution. Extreme lower and upper do not take into account outliers and are defined in the Material and methods.

0.15 substitutions per site with the exception of most Harpacticidae divergences ($COI_{100-580}$, Fig. 2). On the other hand, most G distances are higher than this threshold, with the exception of the already cited Aeglidae distances ($COI_{720-1260}$). Blocks 100–580 and 720–1260 of the COI exhibit similar patterns, and differences between blocks if they exist, are hidden by the important variations between families. Thus in subsequent analyses 5′ and 3′ blocks of the COI have been jointly analyzed.

### 3.2.2. COI AA divergences

Compared to DNA, amino-acid divergences are smaller (approximately 10 times less), with mostly no divergence within species, and G and F divergences often below 0.1 differences per site (Fig. 2). In the few instances where S, G, and F divergences are different from 0, the same pattern as DNA (i.e., enhancement of the differentiation with taxonomic rank) is observed. An exception to this rule is the Penaeidae of the block $COI_{720-1260}$, where S divergences are greater than G ones ($p$ value <0.1).

### 3.2.3. 16S divergences

Most of the 16S alignments contained less than 10% of sites that were unstable when tuning Clustalw parameters with Soap. Exceptions are Parastacidae (23%), Ocypodidae (27%), Porcellanidae (20%), and Cirolanidae (37%). Nevertheless, those sites were conserved to maintain generalization and comparison between families possible. 16S divergences range from 0 to 2.5 substitutions per site (Fig. 3). Like COI, divergence globally increases with taxonomic rank (i.e., S < G < F, Table 2, $p$ value <$10^{-5}$). However and unlike COI, the overlap between S and G divergence classes seems general. There are also overlaps between G and F divergences for some families (Parastacidae, Penaeidae, Cambaridae, Ocypodidae, Varunidae, and Porcellanidae, Fig. 3). Comparison between families led to the same observation as COI that divergence patterns (S < G < F) are not of the same magnitude between families.

### 3.3. Divergence overlaps and thresholds

### 3.3.1. S versus G divergences

Analysis of the overlap between S and G distributions for the COI at nucleotides level (Fig. 4), reveals that a
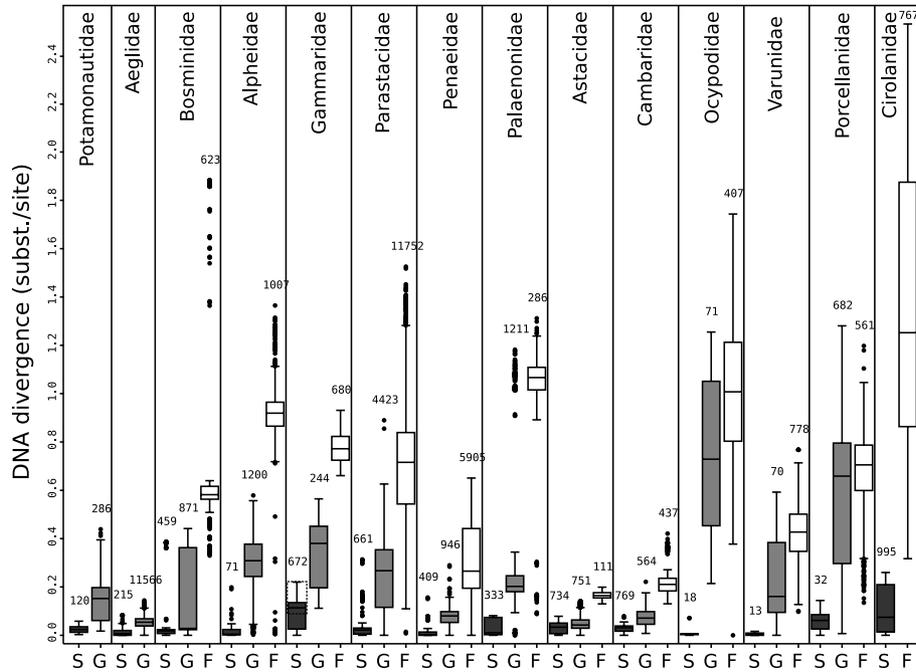
Fig. 3. Same as Fig. 2 applied to positions 760–1220 of the 16S gene and 14 crustacean families.
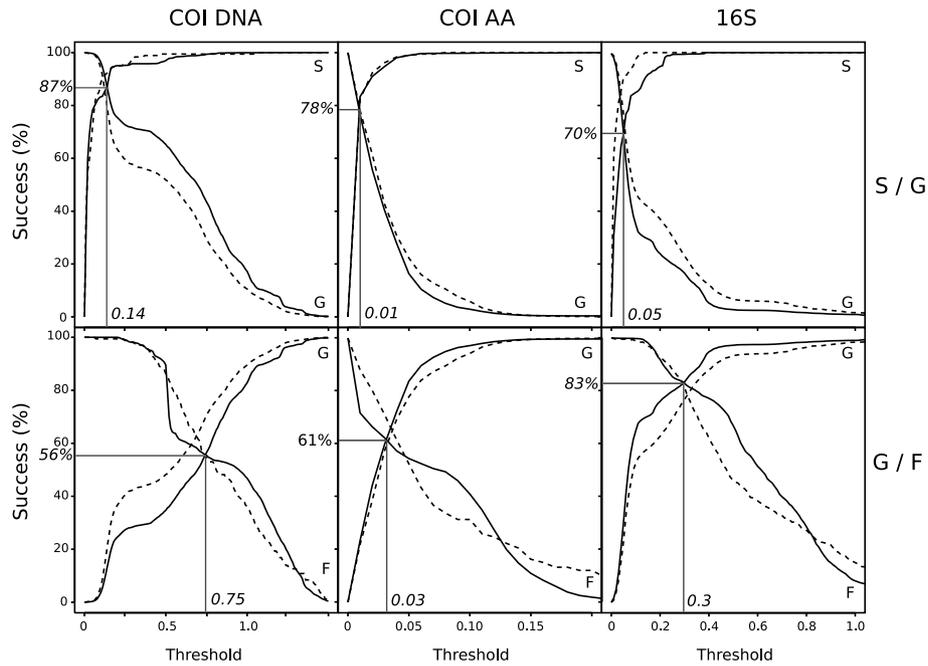


Fig. 4. Analysis of the overlap between S (intra-species) and G (inter-species but intra-genus), and between G and F (inter-genera but intra-family), applied to the nucleotide, amino-acid variations of the COI and variations of the 16S gene. Solid lines refer to the raw data (highly represented taxa have more impact than weakly represented ones) whereas dashed lines present results for the data with corrected taxa sampling (all taxa have the same weight). Methods are described in the Material and methods section and in Fig. 1.

threshold of 0.14 substitution per site would differentiate intra from inter-species divergences in 87 times out of 100. This result indicates a weak overlap between S and G COI DNA divergences. To the opposite the best amino-acid COI S/G threshold is close to 0 (0.01 subst./site) and discriminates S and G divergence with a poor success (78%, Fig. 4). 16S variations show an even worst

performance: the best threshold (0.05 subst./site) only succeeds 70 times out of 100.

### 3.3.2. G versus F divergences

COI nucleotide and amino-acid variations appeared unable to differentiate G from F divergences (respectively 56 and 61% of success, Fig. 4). To the opposite, 16S distin-

guishes relatively well G and F divergences (83% of success) using a threshold of 0.3 subst./site.

### 3.4. Artefacts and taxonomic bias

#### 3.4.1. Taxonomic sampling bias

The impact of unbalanced taxonomic sampling was tested by computing mean divergences and so by giving to each taxa or couple of taxa the same weight (dashed lines in Fig. 4). This design did not modify greatly our previous observations to the exception of the 16S S/G threshold performance which raised to 85% (Fig. 4). These overall results indicate some homogeneity between taxa: the most represented taxa behave like the least represented ones. The case of the 16S indicates that some species with different sampling effort react quite differently (i.e., few well sampled taxa do not behave like the majority of the other taxa). Although it increases the success of the 16S S and G discrimination, that also emphasizes the poor generality of this 16S threshold.

#### 3.4.2. Non-monophyletic taxa

For the three thresholds leading to more than 80% of success (S/G for the COI and the 16S, and G/F divergences for the 16S), we removed all taxa that appeared non-monophyletic in our trees (36 species and 11 genera for the COI, and 44 species and 23 genera for the 16S). Discarding these taxa neither influenced greatly the thresholds nor their performances (first column in Fig. 5). This quite brutal procedure is likely to discard some false positive taxa but the overall stability of thresholds and success rates lets us believe that non-monophyletic taxa do not bias our test.

#### 3.4.3. Hidden diversity

Some authors of the datasets analyzed in this study suggested the existence of extreme divergence (Burton and Lee, 1994; Edmands, 2001) or cryptic speciation (Ganz and Burton, 1995; Müller, 2000; Penton et al., 2004; Rawson et al., 2003; Williams et al., 2001). The concerned families are for the $COI_{100-580}$: Coronulidae, Harpacticidae and Daphniidae, the Alpheidae for the $COI_{720-1260}$, and the Gammaridae for the 16S (involved divergences are highlighted by dotted rectangles in Figs. 2 and 3). We examined the influence of these recognized cases of hidden diversity by removing the concerned taxa from the S divergence computations (i.e., *Chelonibia testudinaria*, *Alpheus lottini*, *Tigriopus californicus* and *Daphnia obtusa* for the COI, and *Gammarus fossarum* for the 16S). Pruning such taxa clearly improved COI S/G threshold since the intersect of the corrected S (red dashes) and G (black plain) curves occurred around 92% instead of 87% (Fig. 5, second column). On the other
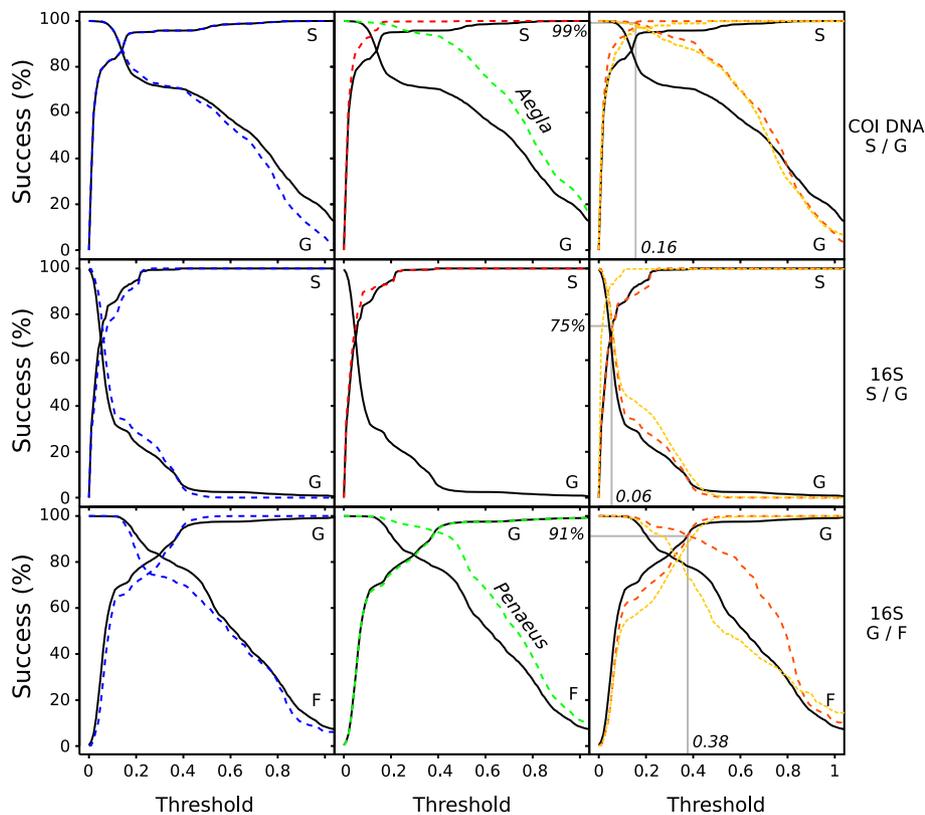


Fig. 5. Same as Fig. 4 for the cases of reduced overlap [i.e., between intra and inter-species for COI (COI DNA S/G, first row) and for 16S DNA variations (16S S/G, second row), and between intra and inter-genera 16S variations (16S G/F, third row)]. The impact of different potential biases is tested and compared to the original dataset (black plain curves): non-monophyletic taxa in the first column (blue dashes), cryptic species (red dashes) and peculiar taxa (green dashes) in the second column, and finally the sum of these potential biases with a balanced (yellow dashes) or unbalanced (orange dashes) taxonomic sampling (third column).

hand we observed almost no effect on the success rate of the 16S threshold for S and G distributions.

### 3.4.4. Genus Aegla

We previously noted that the family Aeglidae was the only family with an important overlap between COI S and G divergences, and also the one with particularly low G divergences (Fig. 2). This family is in fact only represented by the *Aegla* genus whose systematics have been recently studied (Pérez-Losada et al., 2004). Discarding this genus from the COI G divergences greatly improved the performance of the COI threshold (from 87 to 93% or 99% if cryptic species are also removed, Fig. 5), thus demonstrating that the overlap between S and G divergences was mainly due to hidden diversity and genus *Aegla*.

### 3.4.5. Taxonomic status of the genus Penaeus

Within the family Penaeidae, the taxonomic status of the sub-genera of *Penaeus* (e.g., *Litopenaeus*, *Farfantepenaeus*, *Fenneropenaeus*, …) is especially unclear. These taxa are even listed as different genera in the NCBI taxonomy database. The most recent study remains evasive (Lavery et al., 2004), and thus pending on the final decision an important number of F divergences could finally become G ones. To test the influence of this taxonomic uncertainty, we removed this "genus" from the G and F 16S divergences (Fig. 5). Without *Penaeus* the 16S G/F overlap strongly decreased and the performance of the best threshold became quite acceptable (92% of success).

### 3.4.6. Sum of the potential biases

We finally jointly tested the impact of the potential biases (i.e., non-monophyly, hidden diversity, and specific taxa problems) with an unbalanced or balanced design (using mean divergence per taxa, Fig. 5). As described above, when excluding cryptic species, the genus *Aegla*, and non-monophyletic taxa, S and G COI DNA variations could be accurately segregated using a 0.16 subst./site threshold. This threshold and its performance appeared independent of the sampling design, thus demonstrating its robustness (Fig. 5). The joint removal of non-monophyletic and cryptic taxa did not modify the poor performance of the 16S S/G threshold. As for the original data (Fig. 4), the success rate of the best threshold strongly increased for the balanced design and still suggests that few taxa are behaving quite differently. Finally, the 16S G/F threshold without the uncertain *Penaeus* sequences and non-monophyletic taxa suggests its rather good performance (91% of success). Nevertheless this result disappeared with a balanced taxa sampling (80% of success), thus indicating it is an artefact generated by few highly represented and "well behaving" taxa.

### 3.5. Comparison between patristic and pairwise distances

While comparing COI pairwise distances computed with different models of evolution (no corrections: p-distances, a

model differentiating transitions and transversions: K2p, the same one but accounting for rate variations across sites: K2p+G, and a more complex considering the six reversible substitution types with rate variations across sites and invariant positions: GTR+G+I) to patristic distances (Fig. 6), we observed an improvement of the correlation with the complexity of the model. Within the range of this analysis, only the GTR+G+I pairwise distances seem to properly estimate patristic distances. However, the variance of the estimation also increases with the complexity of the model. We here recover the general link between the complexity of the model and the amount of necessary data (Posada and Buckley, 2004). Richly-parameterized models need large amount of data to become accurate, while simple models are more precise but tend to give results biased by homoplasy. In all instances, pairwise uncorrected p and K2p distances quickly reach saturation and seem very poor estimators of the molecular divergence. To test if the COI S/G threshold is dependent on how divergences are measured, we performed the same analysis but using the commonly used pairwise K2p divergences. Using patristic distances from a tree generated by the complex model of evolution or by K2p pairwise distances produced mostly the same results in term of threshold values and success rates (best threshold at 0.15 subst./site with a success of 98%). This absence of difference seems associated to the fact that S divergences are still well estimated using pairwise distances (Table 2). To the opposite G divergences are quickly
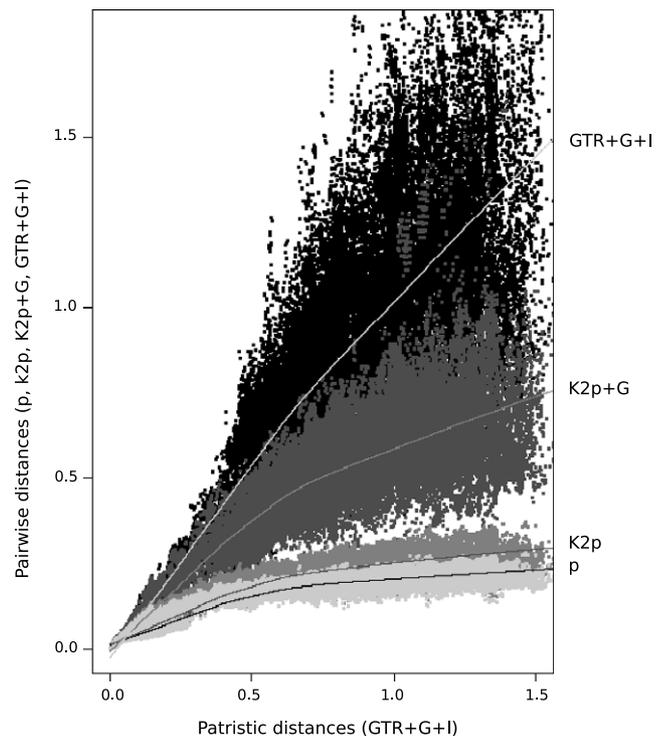


Fig. 6. Comparison between patristic COI divergences and different pairwise distances: p, K2p, K2p+G, and GTR+G+I. Patristic distances were measured on Bayesian trees with branch lengths re-estimated by maximum likelihood under a GTR+G+I model of evolution. Relevant parameters of pairwise distances were estimated by maximum likelihood.

under-estimated but not sufficient to increase overlaps between S and G divergences (Table 2).

## 4. Discussion

### 4.1. Limits of the analysis

This analysis is of course limited to crustaceans, and our sampling is far from being representative of the whole of Crustacea. We are also limited by the resolution of the NCBI taxonomic database and by some unresolved or uncertain taxonomic assignments (e.g., the status of the *Penaeus* sub-genus). Nevertheless it seems unlikely that the general patterns observed from Decapoda to Cirripedia are artefactual outcomes of our sampling or taxonomic references. This study also demonstrates that within a rigorous bioinformatic framework present genomic databases permit to test some important questions concerning DNA taxonomy and its potential applicability.

### 4.2. Which marker to measure molecular divergences?

#### 4.2.1. COI
Both ends of the COI gene appear to be appropriate molecular markers at several taxonomic scales but particularly at the species level. Yet these markers quickly get saturated around 0.3 subst./site (determined by the plateau of the uncorrected p distance, Fig. 6) they remain interesting markers for molecular divergence studies if saturation is compensated. To the opposite, COI aminoacid variations are small at the taxonomic scale we analyzed. Thus, they do not seem to be helpful for taxonomy below the family level.

#### 4.2.2. 16S
As a whole, the 16S evolves more slowly than the COI, and as a result the overlap between S and G distributions is important. On the other hand and for higher taxonomic ranks where the 16S could become a more efficient marker, we encountered some problem of homology during alignment in sites likely to correspond to loops. In such situations, the best option would be to remove such sites. But in doing so, we would loose generalization and any possible comparison between groups. Another way would be to separate stems and loops, and to analyze them separately. However, this would exceed the scope of this paper as no crustacean mitochondrial 16S rRNA secondary structure model is available in the literature. In consequence, the 16S does not appear as an efficient marker for molecular divergence assessments. Obviously, this does not refute its utility as a molecular barcode (see Steinke et al., 2005; Vences et al., 2005). 16S rDNA is probably easier to amplify than the COI, and is also probably a better source of synapomorphies in loop regions, but its potential for molecular divergence estimations seems more limited.

### 4.3. Relationship between taxonomy and molecular divergence

#### 4.3.1. Family and crustacean pattern
Within all families, our analysis shows for both COI and 16S a general increase of the molecular divergence with the taxonomic rank. This would suggest that morphological taxonomy is roughly in agreement with DNA evolution. Yet, this pattern is not perfect, and some divergence distributions at different taxonomic scales overlap. Whereas each family structure apparently respects a molecular hierarchy, the scale of divergence at each taxonomic level appears to vary extensively between families. As an example, in the $COI_{720–1260}$ block, molecular divergences follow the S < G < F ranking, but the Alpheidae and Portunidae families have completely different divergence scales (0.75 against 0.25 mean subst./site for G divergence and 1.2 against 0.4 mean subst./site for F divergences, Fig. 2).

#### 4.3.2. Crustacean S/G divergences
A detailed analysis of the overlap between S, G, and F COI divergences reveals that when cryptic species and the *Aegla* genus are removed, the COI DNA divergences become highly efficient to distinguish S from G divergences (99% of success, Fig. 5). The *Aegla* genus is clearly the only genus—out of 54 analyzed—to be composed of so weakly divergent species (Fig. 2). Thus, either this genus represents an extreme situation of quick morphological diversification and/or slow molecular evolution, or this genus has been over-split. This last hypothesis seems supported by a recent study of this genus (Pérez-Losada et al., 2004) where authors found eight paraphyletic taxa on 22 species sampled at more than one location although the authors had the opposite lecture of their results (i.e., that the genus was under-split, see Pérez-Losada et al., 2004, for a full account). The 16S rRNA gene clearly shows some pattern but is highly influenced by the taxonomic sampling. Furthermore, its best threshold is always below 0.06 subst. site, a value that is in practice unmanageable since experimental errors (e.g., amplification, sequencing) could greatly impact such a low value.

#### 4.3.3. Crustacean G/F divergences
Unlike the S/G divergences, the overlap between G and F divergences is important. Measured by the COI, this overlap is near complete (only 56% of success, Fig. 4). At this level of divergence, the COI appears to be fully saturated (Fig. 6). Despite compensation of saturation by a complex model of evolution, COI divergences are probably inaccurately evaluated and may remain under-estimated at this scale. Thereby the overlap may be over-estimated by saturation. The 16S, on the other hand, is more successful (91 and 82% percent of success with the unbalanced and balanced taxa sampling respectively, Fig. 5). However, the numerical value of the threshold is too drastically impacted by the sampling (from 0.30 to 0.38 subst./site) to be seen as a general trend.

### 4.4. Using molecular divergence in taxonomy and biodiversity assessment?

#### 4.4.1. A COI threshold to help species delimitation

Together with Hebert et al. (2003), we conclude that nucleic acid sequences of COI are of much interest for taxonomy. The weak overlap between S and G COI divergences indicates that in crustaceans there is a morphological and molecular entity, called species, that could be delimited using a rule of 0.16 subst./site in the COI gene. We suggest that this threshold could be used to help the delimitation of new or uncertain species. The proposed criterion is that two monophyletic groups divergent by more than 0.16 subst./site in the COI gene, as measured by patristic distances, have a strong probability to belong to different species. As seen previously, a fairly good approximation could be made using rather crude pairwise divergence measurements (Table 2). Nevertheless, K2p pairwise divergences can be misled by various biases such as multiple substitutions. On the other hand, the more accurate patristic distances involve the burden of reconstructing a phylogenetic tree. Yet this tree being done, the patristic approach will also inform about monophyly of the clades and potentially their support. This two elements are of most importance for DNA taxonomy. It may also be advocated, that the difficulty to compute phylogenetic trees under realistic models are today greatly reduced by the development of fast software such as PhyML (Guindon and Gascuel, 2003) or MrBayes (Huelsenbeck and Ronquist, 2001). Furthermore, the development of user friendly programming languages such as Perl and R make the use of patristic distances compatible with any DNA taxonomy projects.

#### 4.4.2. DNA taxonomy

As we previously stated, this study is not intended as a test for DNA barcoding validity. Therefore the debate about whether the COI can produce molecular synapomorphies to identify taxa is not to be settled here. On the other hand this study produces some answers about incomplete lineage sorting of ancestral polymorphisms, gene introgression, large changes in the substitution rate between lineages, and paralogy effects. Effects that should all disrupt DNA taxonomy (Mallet and Willmott, 2003; Moritz and Cicero, 2004; Will and Rubinoff, 2004). Indeed if these phenomena were frequent enough so that DNA taxonomy is pointless, we should not observe a universal species molecular threshold as we did. Furthermore, we could also argue that species delimitation based on sole maternally inherited molecule is biased. This is indeed likely to be true (e.g., in case of asymmetric hybridization or sex-biased gene flow), however our results indicate that these artefacts are also quite rare or at least have little effects in crustaceans. Species delimitation based on COI is by essence imperfect and only an approach combining at least two genes, a mitochondrial and a nuclear one, would make the proposed criterion more robust. This implies further investigations that are yet limited by the reduced amount of crustacean nuclear sequences available in databases. It is also seen as a major issue, when using COI for both delimitation or identification, that this molecule will not be able to distinguish recent species (e.g., Mathews et al., 2002). However, this constitutes a challenge for any taxonomy and not only for DNA based ones.

#### 4.4.3. Link with the species definition

Our proposed criterion to help delimiting species is not linked to any particular species definition. In this way the word "species" could have been replaced by OTU (operational taxonomic units). However, the observed link between present taxonomy (likely to have been influenced by different species concepts) and molecular divergences suggests that this pragmatic approach can be used as an objective tool to help taxonomy in difficult situations (e.g., absence of specialist, morphological convergences). Because our criterion directly relies upon monophyletic units, it is strongly linked to the phylogenetic species definition (PSD) with the difference that a quantitative criterion, that is a molecular threshold, is added. Thus, instead of considering the smallest diagnosable monophyletic unit as a species, we could only consider monophyletic units that contain taxa diverging by less than 0.16 subst./site of the COI. Consequently, any species defined using our criterion is likely to be considered as a species under the PSD. Agapow et al. (2004) argue that the increasing use of the PSD should lead to an increase in species numbers. This taxonomic inflation could have important repercussions on conservation and macro-ecology, particularly when the taxonomic changes are biased toward certain groups (Isaac et al., 2004). Our analysis revealed that the joint use of the PSD and a molecular threshold should not lead to a dramatic increase in species number. In this way the introduction of molecular thresholds could help for the acceptance of the PSD as an operational species definition.

### Acknowledgments

### Appendix A. Supplementary data

Supplementary data associated with this article can be found, in the online version, at doi:10.1016/j.ympev.2006.03.014.

### References

Agapow, P.M., Bininda-Emonds, O.R., Crandall, K.A., Gittleman, J.L., Mace, G.M., Marshall, J.C., Purvis, A., 2004. The impact of species concept on biodiversity studies. Q. Rev. Biol. 79, 161–179.

Avise, J.C., Johns, G.C., 1999. Proposal for a standardized temporal scheme of biological classification for extant species. Proc. Natl. Acad. Sci. USA 96, 7358–7363.

Avise, J.C., Walker, D., 1999. Species realities and numbers in sexual vertebrates: perspectives from an asexually transmitted genome. Proc. Natl. Acad. Sci. USA 96, 992–995.

Blaxter, M.L., 2004. The promise of a DNA taxonomy. Philos. Trans. R. Soc. B Biol. Sci. 359, 669–679.

Blaxter, M., Mann, J., Chapman, T., Thomas, F., Whitton, C., Floyd, R., Eyualem-Abebe, 2005. Defining operational taxonomic units using DNA barcode data. Philos. Trans. R. Soc. B 360, 1935–1943.

Burton, R.S., Lee, B.N., 1994. Nuclear and mitochondrial gene genealogies and allozyme polymorphism across a major phylogeographic break in the copepod *Tigriopus californicus*. Proc. Natl. Acad. Sci. USA 91, 5197–5201.

Chase, M.W., Salamin, N., Wilkinson, M., Dunwell, J.M., Prasad Kesanakurthi, R., Nadia Haidar, N., Savolainen, V., 2005. Land plants and DNA barcodes: short-term and long-term goals. Philos. Trans. R. Soc. B 360, 1889–1895.

Claridge, M.F., Dawah, H.A., Wilson, M.R., 1997. Species: The Units of Biodiversity. Chapman & Hall, London.

Daniels, S.R., Gouws, G., Stewart, B.A., Coke, M., 2003. Molecular and morphometric data demonstrate the presence of cryptic lineages among freshwater crabs (Decapoda: Potamonautidae: *Potamonautes*) from the Drakensberg mountains, South Africa. Biol. J. Linnean Soc. 78, 129–147.

de Bruyn, M., Wilson, J.A., Mather, P.B., 2004. Huxley's line demarcates extensive genetic divergence between eastern and western forms of the giant freshwater prawn, *Macrobrachium rosenbergii*. Mol. Phylogenet. Evol. 30, 251–257.

DeSalle, R., Egan, M.G., Siddall, M., 2005. The unholy trinity: taxonomy, species delimitation and DNA barcoding. Philos. Trans. R. Soc. B 360, 1905–1916.

Edmands, S., 2001. Phylogeography of the intertidal copepod *Tigriopus californicus* reveals substantially reduced population differentiation at northern latitudes. Mol. Ecol. 10, 1743–1750.

Funk, D.J., Omland, K.E., 2003. Species-level paraphyly and polyphyly: frequency, causes and consequences, with insights from animal mitochondrial DNA. Annu. Rev. Ecol. Evol. Syst. 34, 397–423.

Ganz, H.H., Burton, R.S., 1995. Genetic differentiation and reproductive incompatibility among baja california populations of the copepod *Tigriopus californicus*. Mar. Biol. 123, 821–827.

Guindon, S., Gascuel, O., 2003. A simple, fast, and accurate algorithm to estimate large phylogenies by maximum likelihood. Syst. Biol. 52, 696–704.

Harris, J.D., Froufe, E., 2005. Taxonomic inflation: species concept or historical geopolitical bias? Trends Ecol. Evol. 20, 6–9.

Hebert, P.D.N., Cywinska, A., Ball, S.L., deWaard, J.R., 2003. Biological identifications through DNA barcodes. Proc. R. Soc. Lond. B Biol. Sci. 270, 313–321.

Hebert, P.D.N., Stoeckle, M.Y., Zemlak, T.S., Francis, C.M., 2004. Identification of birds through DNA barcodes. PLoS Biol. 2, e312.

Hendry, A.P., Vamosi, S.M., Latham, S.J., Heilbuth, J.C., Day, T., 2000. Questioning species realities. Conserv. Genet. 1, 67–76.

Huelsenbeck, J.P., Ronquist, F., 2001. MRBAYES: Bayesian inference of phylogenetic trees. Bioinformatics 17, 754–755.

Isaac, N.J.B., Mallet, J., Mace, G.M., 2004. Taxonomic inflation: its influence on macroecology and conservation. Trends Ecol. Evol., 464–469.

Jarman, S.N., Elliott, N.G., 2000. DNA evidence for morphological and cryptic Cenozoic speciations in the Anaspididae, 'living fossils' from the Triassic. J. Evol. Biol. 13, 624–633.

Johns, G.C., Avise, J.C., 1998. A comparative summary of genetic distances in the vertebrates from the mitochondrial cytochrome *b* gene. Mol. Biol. Evol. 15, 1481–1490.

King, J.L., Hanner, R., 1998. Cryptic species in a "living fossil" lineage: taxonomic and phylogenetic relationships within the genus *Lepidurus* (Crustacea: Notostraca) in North America. Mol. Phylogenet. Evol. 10, 23–36.

Lavery, S., Chan, T.Y., Tam, Y.K., Chu, K.H., 2004. Phylogenetic relationships and evolutionary history of the shrimp genus *Penaeus s.l.* derived from mitochondrial DNA. Mol. Phylogenet. Evol. 31, 39–49.

Lee, C.E., 2000. Global phylogeography of a cryptic copepod species complex and reproductive isolation between genetically proximate "populations". Evolution 54, 2014–2027.

Lefébure, T., Douady, C.J., Gouy, M., Trontelj, P., Briolay, J., Gibert, J., in press. Phylogeography of a subterranean amphipod reveals cryptic diversity and dynamic evolution in extreme environments. Mol. Ecol.

Lipscomb, D., Platnick, N., Wheeler, Q., 2003. The intellectual content of taxonomy: a comment on DNA taxonomy. Trends Ecol. Evol. 18, 65–66.

Löytynoja, A., Milinkovitch, M.C., 2001. SOAP, cleaning multiple alignments from unstable blocks. Bioinformatics 17, 573–574.

Mallet, J., Willmott, K., 2003. Taxonomy: renaissance or tower of babel? Trends Ecol. Evol. 18, 57–59.

Martin, J.W., Davis, G.E., 2001. An updated classification of the recent Crustacea. Natural History Museum of Los Angeles County, Los Angeles.

Mathews, L.M., Schubart, C.D., Neigel, J.E., Felder, D.L., 2002. Genetic, ecological, and behavioural divergence between two sibling snapping shrimp species (Crustacea: Decapoda: *Alpheus*). Mol. Ecol. 11, 1427–1437.

Mayden, R.L., 1997. A hierarchy of species concepts: the denouement in the saga of the species problem. In: Claridge, M.F., Dawah, H.A., Wilson, M.R. (Eds.), Species: The Units of Biodiversity. Chapman &Hall, London, pp. 381–424.

Monaghan, M.T., Balke, M., Gregory, T.R., Vogler, A.P., 2005. DNA-based species delineation in tropical beetles using mitochondrial and nuclear markers. Philos. Trans. R. Soc. B 360, 1925–1933.

Moritz, C., Cicero, C., 2004. DNA barcoding: promise and pitfalls. PLoS Biology 2, e354.

Müller, J., 2000. Mitochondrial DNA variation and the evolutionary history of cryptic *Gammarus fossarum* types. Mol. Phylogenet. Evol. 15, 260–268.

Paradis, E., Claude, J., Strimmer, K., 2004. Ape: analyses of phylogenetics and evolution in R language. Bioinformatics 20, 289–290.

Penton, E.H., Hebert, P.D., Crease, T.J., 2004. Mitochondrial DNA variation in North American populations of *Daphnia obtusa*: continentalism or cryptic endemism? Mol. Ecol. 13, 97–107.

Pérez-Losada, M., Bond-Buckup, G., Jara, C.J., Crandall, K.A., 2004. Molecular systematics and biogeography of the southern South American freshwater "crabs" *Aegla* (Decapoda: Anomura: Aeglidae) using multiple heuristic tree search approaches. Syst. Biol., 767–780.

Posada, D., Buckley, T.R., 2004. Model selection and model averaging in phylogenetics: advantages of Akaike information criterion and Bayesian approaches over likelihood ratio tests. Syst. Biol. 53, 793–808.

Proudlove, G., Wood, P.J., 2003. The blind leading the blind: cryptic subterranean species and DNA taxonomy. Trends Ecol. Evol. 18, 272–273.

R Development Core Team, 2004. R: a language and environment for statistical computing.

Rawson, P.D., Macnamee, R., Frick, M.G., Williams, K.L., 2003. Phylogeography of the coronulid barnacle, *Chelonibia testudinaria*, from loggerhead sea turtles, *Caretta caretta*. Mol. Ecol. 12, 2697–2706.

Rocha-Olivares, A., Fleeger, J.W., Foltz, D.W., 2001. Decoupling of molecular and morphological evolution in deep lineages of a meiobenthic harpacticoid copepod. Mol. Biol. Evol. 18, 1088–1102.

Savolainen, V., Cowan, R.S., Vogler, A.P., Roderick, G.K., Lane, R., 2005. Towards writing the encyclopaedia of life: an introduction to DNA barcoding. Philos. Trans. R. Soc. B 360, 1805–1811.

Schmidt, H.A., Strimmer, K., Vingron, M., von Haeseler, A., 2002. TREE-PUZZLE: maximum likelihood phylogenetic analysis using quartets and parallel computing. Bioinformatics 18, 502–504.

Seberg, O., Petersen, G., Humphries, C.J., Knapp, S., Stevenson, D.W., Scharff, N., Andersen, N.M., 2003. Shortcuts in systematics? a commentary on DNA-based taxonomy. Trends Ecol. Evol. 18, 63–65.

Sites, J.W., Marshall, J.C., 2004. Operational criteria for delimiting species. Annu. Rev. Ecol. Evol. Syst. 35, 199–227.

Steinke, D., Vences, M., Salzburger, W., Meyer, A., 2005. TaxI: a software tool for DNA barcoding using distance methods. Philos. Trans. R. Soc. B 360, 1975–1980.

Swofford, D.L., 2002. PAUP*. Phylogenetic Analysis Using Parsimony (*and Other Methods), Version 4. Sinauer Associates, Sunderland.

Tautz, D., Arctander, P., Minelli, A., Thomas, R.H., Vogler, A.P., 2003. A plea for DNA taxonomy. Trends Ecol. Evol. 18, 70–74.

Thompson, J.D., Higgins, D.G., Gibson, T.J., 1994. CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, positions-specific gap penalties and weight matrix choice. Nucleic Acids Res. 22, 4673–4680.

Tuckey, J.W., 1977. Exploratory Data Analysis. Addison–Wesley, Boston.

Vences, M., Thomas, M., van der Meijden, A., Chiari, Y., Vieites, D.R., 2005. Comparative performance of the 16S rRNA gene in DNA barcoding of amphibians. Front. Zool. 2, 5.

Wares, J.P., 2001. Patterns of speciation inferred from mitochondrial DNA in North American *Chthamalus* (Cirripedia: Balanomorpha: Chthamaloidea). Mol. Phylogenet. Evol. 18, 104–116.

Will, K.W., Rubinoff, D., 2004. Myth of the molecule: DNA barcodes for species cannot replace morphology for identification and classification. Cladistics 20, 47–55.

Williams, S., Knowlton, N., 2001. Mitochondrial pseudogenes are pervasive and often insidious in the snapping shrimp genus *Alpheus*. Mol. Biol. Evol. 18, 1484–1493.

Williams, S.T., Knowlton, N., Weigt, L.A., Jara, J.A., 2001. Evidence for three major clades within the snapping shrimp genus *Alpheus* inferred from nuclear and mitochondrial gene sequence data. Mol. Phylogenet. Evol. 20, 375–389.

Witt, J.D.S., Hebert, P.D.N., 2000. Cryptic species diversity and evolution in the amphipod genus *Hyalella* within central glaciated North America: a molecular phylogenetic approach. Can. J. Fish. Aquat. Sci. 57, 687–698.

Zhang, Z., Schwartz, S., Wagner, L., Miller, W., 2000. A greedy algorithm for aligning DNA sequences. J. Comput. Biol. 7, 203–214.