

# Using phylogenetic information to predict species tolerances to toxic chemicals

GUILLAUME GUÉNARD,<sup>1,2,5</sup> PETER CARSTEN VON DER OHE,<sup>3</sup> DICK DE ZWART,<sup>4</sup> PIERRE LEGENDRE,<sup>2</sup> AND SOVAN LEK<sup>1</sup>

<sup>1</sup>Laboratoire Évolution et Diversité Biologique (EDB) UMR 5174 CNRS / Université Paul-Sabatier, 118 Rte. de Narbonne 4R3-b1, 31062 Toulouse Cedex 9, France

<sup>2</sup>Département des Sciences Biologiques, Université de Montréal, CP 6128 Succursale Centre-Ville, Montréal, Québec H3C 3J7 Canada

<sup>3</sup>UFZ, Department of Effect-Directed Analysis, Helmholtz Centre for Environmental Research—UFZ, Permoserstrasse 15, 04318 Leipzig, Germany

<sup>4</sup>Laboratory for Ecological Risk Assessment (LER), National Institute of Public Health and the Environment (RIVM), P.O. Box 1 NL-3720 BA Bilthoven, The Netherlands

**Abstract.** Tolerance to toxic substances is a characteristic of an organism that determines whether it is able to withstand the concentrations occurring in its environment. The measurement of tolerance is therefore of fundamental importance when assessing the impact of anthropogenic chemicals on ecosystems and ecological communities. Although an appreciable amount of information on species tolerance to chemicals has been collected through the last 50 years, substantial gaps remain in our knowledge of tolerance relative to the diversity of organisms inhabiting aquatic ecosystems and the great and increasing number of chemicals released in these ecosystems. Within that context, methods allowing one to reliably and accurately estimate a species' tolerance using other known characteristics would be valuable. In the present study we introduce an approach that uses phylogeny to estimate the tolerance of a species using that of a set of other species related to the focus species at different phylogenetic scales. We estimated phylogenies from molecular data (DNA sequences) or inferred them from taxonomy. Up to 83% of the among-species variation in tolerance (log-transformed median lethal concentration over 96 hours; LC<sub>50</sub>) was found to be phylogenetically structured and was therefore usable for making predictions. The ability of phylogenetic models to produce accurate estimates of species tolerances is apparently related to the availability of information within species groups and the variation in pesticide tolerance within these groups. Toxicity models integrating phylogeny therefore appear suitable to assist in risk assessment.

**Key words:** aquatic organisms; molecular characters; pesticides; phylogenetic eigenfunctions; phylogenetic model; phylogeny; risk assessment; tolerance.

## INTRODUCTION

Tolerance to toxic substances is a trait that determines the ability of organisms to withstand the level of pollutants occurring in their environment and is thus central to assessing the effects of toxicity on biodiversity (e.g., the calculation of species sensitivity distributions; von der Ohe and Liess 2004, Postuma et al. 2002). Tolerance is commonly approximated using bioassays, which are controlled experiments where individuals are exposed, for a given amount of time, to different concentrations of a substance or a mixture of substances and an effect is observed on a portion of the population (e.g., death of 50% of the population over 48 h of exposure, or inhibition of reproduction for 90% of the

population after 96 h of exposure). While being a useful trait for ecotoxicologists, estimating tolerance is costly (several thousands U.S. dollars needed per estimate), logistically challenging (lots of laboratory space and personnel must be mobilized), and sometimes impossible for all important species since specimens need to be raised in captivity or collected alive in nature. There are many substances known to be hazardous to organisms in the environment. The challenge faced by ecotoxicologists is to provide reliable estimates of tolerance for as many species–substance combinations as possible. This task is extremely difficult given the ever-increasing number of potentially hazardous compounds that are introduced each year and the often broad variety of organisms inhabiting the ecosystems affected by anthropogenic releases. It is therefore of interest to find alternatives to the exhaustive testing of species–substance combinations. Methods allowing the estimation of tolerance using other features of organisms—for instance, trait values (e.g., morphological, physiological, biochemical and/or ecological traits) and/or their phylogeny with respect to other species having known

Manuscript received 22 November 2010; revised 24 May 2011; accepted 26 May 2011. Corresponding Editor: M. E. Hellberg.

<sup>5</sup> Present address: Département des Sciences Biologiques, Université de Montréal, CP 6128 succursale Centre-Ville, Montréal, QC H3C 3J7 Canada.  
E-mail: guillaume.guenard@gmail.com

tolerance—may represent such alternatives. Here we consider methods for predicting tolerance using a statistical modeling approach based on phylogeny.

Tolerance is the result of multiple subordinate traits related to the uptake of pollutants by the organisms, their metabolism (e.g., transport, accumulation, sequestration, activation/inactivation), and excretion. Dependence on a wide array of such subordinate traits may generate character correlation and (positive or negative) phylogenetic autocorrelation. Modeling approaches can take advantage of these correlations to estimate tolerance, lessening the complexity associated with the numerous toxic substances and species co-occurring in the environment. *Character correlation* occurs when the phenotype value of a given trait is correlated with that of another trait as a consequence of, for instance, their common reliance on similar subordinate traits influenced by genetic (e.g., pleiotropy, linkage disequilibrium) or environmental (e.g., correlational selection; Lande and Arnold 1983) processes. The presence of character correlation implies that the value of a trait that is hard to measure can be, to some extent, estimated from that of a trait that is more easily obtained. That approach was used by Baird and Van den Brink (2007) to estimate tolerance (the median lethal concentration: the concentration that kills 50% of individuals of a population over a specified amount of time,  $LC_{50}$ ) using species' traits related to morphology, life history, physiology and feeding ecology. The second trait property, *phylogenetic autocorrelation*, implies that trait values show dependence with respect to species' positions in a phylogeny and may occur over multiple scales. Positive phylogenetic autocorrelation implies that closely related species share more similar trait values in comparison to more distant ones, as a consequence of evolution proceeding slowly by means of a series of small steps, over a long time period (Diniz-Filho et al. 1998, Blomberg et al. 2003, Buchwalter et al. 2008). Positive autocorrelation shows up as large-scale structures in phylogenetic trait signals. These large-scale structures are characterized by large differences between species pairs from different high-order taxonomic groups and small differences between species pairs from the same high-order taxonomic groups. However, closely related species can vary markedly in individual traits as a result of differentiation among parent species (e.g., interspecific competition; Svanbäck and Bolnick 2007). By contrast, negative phylogenetic autocorrelation implies that closely related species have more dissimilar trait values than more distant species. Negative autocorrelation appears as small-scale structures in phylogenetic trait signals. These small-scale structures are characterized by large differences occurring between closely related species pairs (e.g., from the same low-order taxonomic groups) and small differences occurring between loosely related species pairs. As for character correlation, the phylogenetic autocorrelation of a trait such as tolerance may be attributed to

subordinate traits, thereby reducing or enhancing its rate of change depending on the level of nonadditivity of the effects of those subordinate traits on higher-order traits. Hence, the effect of a change in a given subordinate trait may be dampened by that of other, more conserved, subordinate traits (leading to small differences among closely related species) whereas a change of a similar magnitude, but on a different subordinate trait, may have exacerbating effects (leading to substantial differences among closely related species). Phylogenetic autocorrelation was found to reliably describe the extinction threat to amphibians (Corey and Waite 2008) and the bioaccumulation of cadmium in insects (Buchwalter et al. 2008) and of trace elements in fish (Jeffree et al. 2010). However, and in spite of their anticipated relevance, predictive modeling approaches based on phylogenetic autocorrelation remain sparse.

The purpose of our present study is to develop a statistical modeling approach for making predictions of species' tolerances to toxic substances based on information available from other species and their common phylogeny, which can be obtained using different methods. We achieved this by providing assessments of (1) the fraction of variation in the tolerance of a set of species to toxic substances that can be modeled by phylogeny and of (2) the predictive power of tolerance models based on phylogeny. Considering the wide range of information and techniques now available to reconstruct the evolutionary relatedness of species, phylogenetic modeling of species tolerance may represent a critical step towards the improvement of toxicity assessment. The same approach could also be used to compute predictive models for any other species traits that exhibit phylogenetic autocorrelation.

## METHODS

### *Data sources and selection*

We used a database of concentrations associated with different toxicological endpoints and effects for various substances, aquatic species, and exposure times (de Zwart 2002). That database has been compiled from three sources: (1) AQUIRE (USEPA 1984) from U.S. Environmental Protection Agency, Mid-Century Ecology Division, (2) a compilation of pesticide toxicity made by the Centre for Substances and Risk Assessment (Netherlands National Institute of Public Health and the Environment; Crommentuijn et al. 1997, Tomlin 1997), and (3) another compilation of pesticide toxicity offered by the U.S. Environmental Protection Agency, Office of Pesticides Programs, Ecological Effects Branch. From that database we selected data of lethal concentration ( $LC_{50}$ ) after 96 hours while excluding all entries with inequality indications (i.e., greater than or smaller than). We selected that particular endpoint–effect combination in order to obtain the greatest number of substance–species combinations (7170 combinations over 8848 entries, with 1731 substances and 759 species involved). When multiple test values were

found for one substance, quality checks such as water solubility were employed to eliminate odd data entries (e.g., unit transformation errors). If values differed by more than a factor of 30 from the closest one in a group of at least two other references, we discarded the aberrant value in order to remove outliers from the data set. Of all the remaining values for a given substance, we took the geometric mean as the valid experimental value. The remaining selection procedure aimed at obtaining the largest set of species whose effect concentrations were available for as many substances as possible with no missing information. To obtain that species-by-compound table we first classified species and chemicals by decreasing order of number of effect concentrations available and investigated the topmost elements of the resulting lists.

#### *Obtaining phylogenies*

Phylogenies can either be estimated using suitable characters, or obtained from the literature. A wide variety of phylogenetic inference methods now exists (e.g., maximum-parsimony, distance-based, maximum-likelihood, spectral, or Bayesian methods), whereas abundant, and rapidly increasing, information about molecular (DNA) characters is being made available on the Internet through organizations such as the U.S. National Center for Biotechnology Information (NCBI; *available online*).<sup>6</sup> Phylogenies can also be found within the molecular taxonomy literature or from the Tree of Life project (ToL; Maddison et al. 2007). Our methodology can be used with any of these sources of phylogenetic information as long as they are considered reliable and accurate.

We used two different approaches to obtain phylogenies in the present study. The first involved the estimation of a tree from DNA sequences using a maximum-likelihood approach (Felsenstein 1981, Felsenstein and Churchill 1996; analysis was performed using the software EMBOSS version 6.1.0–5 [Rice et al. 2000]). To do so, we obtained DNA sequences from NCBI's Nucleotide database which consisted, whenever available, of the entire mitochondrial genome as well as nuclear DNA sequences for 28S, 18S, and 5.8S ribosomal RNA transcripts and their internal transcribed spacers (ITS 1 and ITS 2). Then, we performed multiple sequence alignment on each gene separately using the computer program MUSCLE (version 3.7; Edgar 2004). Finally, we concatenated these aligned sequences into a super alignment of genes before estimating the tree. The resulting tree was used to assess the ability of phylogenetic autocorrelation to describe the tolerance of a set of species to multiple pesticides.

The second approach involved constructing a tree from information on taxonomic classification. For that purpose, we gathered information on a maximum of 19

taxonomic ranks from the ToL project for each species. Species with no available information for a given rank were assigned a generic taxon for that rank. We constructed the tree topology implied by the hierarchical structure of taxonomy and placed all taxa of a given rank at the same distance from the root.

Although the construction of a species tree from taxonomy may be the only solution available in the absence of suitable molecular information, readers must be warned that there are many situations in which these trees may not accurately represent the phylogeny. For instance, trees constructed from the taxonomy of species covering a wide range of high-order taxa may be congruent with molecular phylogenetics trees in term of their tree topology while their adequacy in representing branch lengths may remain questionable. The quality of a tree constructed from the taxonomy of species covering a narrower range of low-order taxa would be questionable both in terms of topology and branch lengths. As is the case for modeling methods in general, the modeling approach described herein assumes that the explanatory factor that is provided (i.e., the phylogeny), and on which it depends, is correct. In most practical situations, trees estimated from molecular phylogenetic methods should therefore be preferred over trees constructed using taxonomic classification.

#### *Constructing a phylogenetically explicit model*

We represented the structures of phylogenetic signals using eigenfunctions derived from a phylogenetic tree, a method also known as "phylogenetic eigenvectors regression" (PVR; Desdevises et al. 2003, Diniz-Filho et al. 1998; Diniz-Filho et al. (1998) used only the first few eigenfunctions obtained by principal-coordinate analysis (PCoA) to represent the phylogeny, whereas Desdevises et al. (2003) used all eigenfunctions, as in the method described in the present paper). These eigenfunctions were computed from the phylogenetic covariance matrix  $\mathbf{W}$  whose elements  $w_{i,j}$  correspond to the length of path leading from the root of the tree to the first common ancestor of species  $i$  and  $j$ . The eigenvalues and eigenvectors associated with  $\mathbf{W}$  after double centering were obtained by solving the equation

$$\mathbf{\Omega} = \mathbf{Q}\mathbf{W}\mathbf{Q} = \mathbf{U}\mathbf{D}_\lambda\mathbf{U}^\top$$

$$\mathbf{Q} = \mathbf{I}_n - \frac{1}{n}\mathbf{1}_n\mathbf{1}_n^\top \quad (1)$$

where  $\mathbf{U}$  is a matrix whose columns are eigenvectors, diagonal matrix  $\mathbf{D}_\lambda$  is a diagonal matrix of eigenvalues,  $\mathbf{Q}$  is a centering matrix calculated from an  $n \times n$  identity matrix  $\mathbf{I}_n$  and a vector of  $n$  1's  $\mathbf{1}_n$ ;  $n$  is the number of species and superscript  $\top$  denotes matrix transposition. As consequences of the symmetry of  $\mathbf{W}$  and its centering prior to eigenvalue decomposition,  $n - 1$  nonzero and mutually orthogonal unit vectors are obtained, defining an orthonormal basis against which trait variance can be decomposed with respect to phylogeny in a multiple-

<sup>6</sup> (<http://www.ncbi.nlm.nih.gov/>)

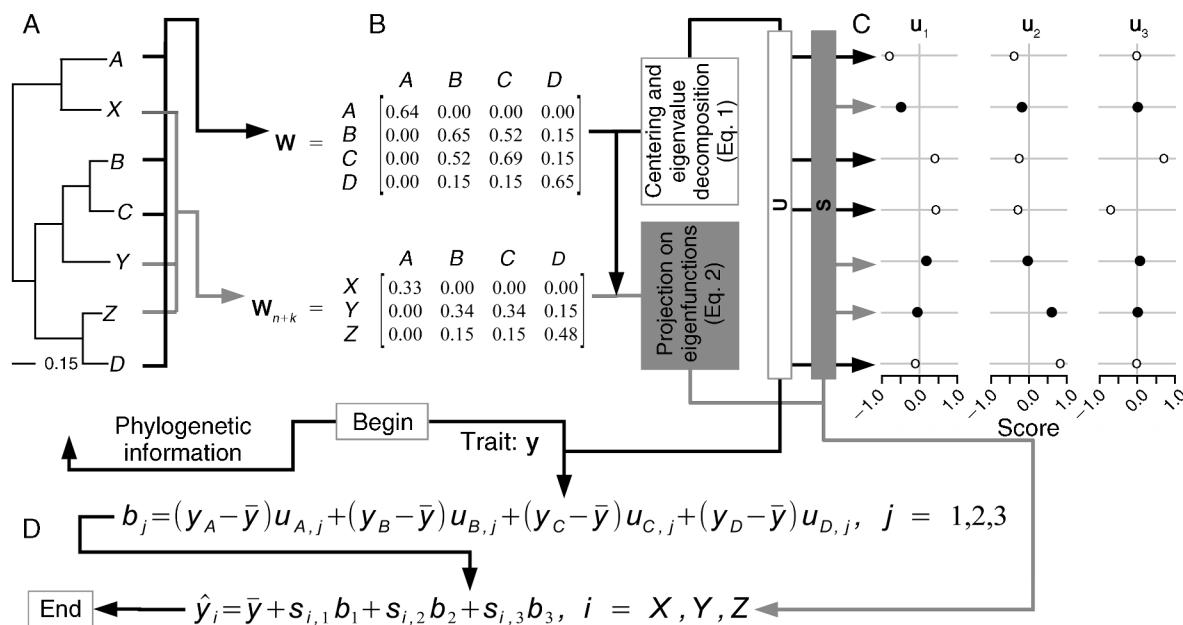


FIG. 1. Example illustrating the approach for modeling trait values using phylogeny. We start with trait values (vector  $y$ , mean trait value), which are known for species  $A$ – $D$  and are estimated for species  $X$ – $Z$  using phylogenetic information on all seven species. (A) The phylogenetic information is used to estimate a tree. (B) Phylogenetic covariance matrices (among species  $A$ – $D$ ,  $W$ ; and between species  $X$ – $Z$  and  $A$ – $D$ ,  $W_{n+k}$ ) are obtained from the tree. (C) These matrices are in turn used to obtain the species score matrices  $U$  (by eigenvalue decomposition after row and column centering on means where  $u_1$ ,  $u_2$ , and  $u_3$  are column vectors in the matrix; open circles) and  $S$  (by projection on the eigenfunctions defined for species  $A$ – $D$ ; solid circles). (D) Score matrix  $U$  is used to estimate the parameters  $b_j$  of any given eigenfunctions  $j$  in a linear model. The linear model is finally used to estimate trait values  $\hat{y}_i$  of species  $i$ , for which the trait is unknown, from its scores  $s_{i,j}$ .

scale fashion (Fig. 1). That approach is similar to a principal coordinate analysis based on a similarity matrix (Gower 1966).

In the models developed in the present study, the response variable  $y$  is a vector whose elements are experimental  $LC_{50}$  values for a given species and compound. This variable was regressed against a design matrix  $X$  involving two factors and their interaction. The first factor describes the fraction of  $LC_{50}$  variability that is associated strictly with mean toxicity of each compound for all the species involved in the model and was represented using Helmert orthogonal contrast variables. Each instance of a given compound was represented in the design matrix  $X$  by its scores on the contrast variables. The number of such contrast variables in the design matrix was  $m - 1$ , where  $m$  is the number of compounds considered. The second factor describes the  $LC_{50}$  variability that is associated strictly with the mean susceptibility of species for all the compounds involved in the model. It was represented in the design matrix using the species scores on each of the  $n - 1$  eigenvectors obtained from Eq. 1, the scores of any given species being repeated for each compound. The interaction term between the two factors describes the  $LC_{50}$  variability that is not accounted for simply by adding the mean toxicity of compounds with the mean susceptibility of species, thereby allowing the model to represent cases where different compounds affect the

species within the phylogeny in different ways. That interaction term was represented in the design matrix by the set of all possible  $(m - 1) \times (n - 1)$  element-wise multiplication of any Helmert contrast variable describing mean compound toxicity with any variable describing species susceptibility at a particular phylogenetic scale from their position in the phylogeny. The design matrix included a column of 1's to allow the estimation of the intercept of the model.

In order to avoid over-fitting, a column subset of the design matrix was selected when constructing the models. We obtained that subset by first including the factor representing the compounds (i.e., the intercept and Helmert contrasts) and then performing a forward-stepwise selection, using  $F$  tests, of the variables representing phylogeny and the interactions between compounds and phylogeny. Family-wise (corrected)  $P$  values of the inference tests performed for the stepwise addition of variables were obtained using the sequential Bonferroni procedure (Holm 1979). Finally, proportions of variation associated with the compounds, the phylogeny, and the compound–phylogeny interactions were estimated as their respective adjusted coefficients of determination. That approach is meant to provide a column subset  $X_S$  of the design matrix  $X$  that best fitted the response variable while avoiding over-fitting. It does not, however, allow one to make predictions of  $LC_{50}$  values for additional species.

*Making predictions*

The approach to make predictions for additional species involves four steps. Firstly, the positions of the new species in the phylogenetic tree have to be taken from a previous analysis or estimated. In the case that the position has to be estimated, the new species must be added to the established phylogenetic tree (i.e., the one used to calculate **W**) without modifying the topology and branch lengths of the subset tree for the original species. Warning should be made here that redoing/repeating phylogenetic analysis with one or more additional species often results in the alteration of the original subset tree. Under these circumstances, the orthonormal basis must be recalculated and any model based on it rebuilt. We avoided this issue by including in the phylogenetic analysis, from the beginning, the species for which predictions were to be made; it was then possible to select the subset tree of the *n* species with known response variable to estimate the phylogenetic model, and then use the positions of the remaining *q* species to make predictions. Secondly, a  $q \times n$  matrix **W**<sub>*n+k*</sub> whose elements *w*<sub>*n+k,j*</sub> are the lengths of the paths leading from the root of the tree to the first common ancestor of a new species *k* and a species *j* within the model, is calculated. Thirdly, the projection scores **S**<sub>*n+k*</sub> of the new species on the *n* - 1 eigenfunctions underlying the eigenvectors in **U** are obtained following Gower's approach for adding new observations in an existing principal coordinate analysis, by rearranging Eq. 1 and performing a partial substitution of matrix **W** by **W**<sub>*n+k*</sub> (Gower 1969; also Fig. 1C: solid circles):

$$S_{n+k} = \left\{ W_{n+k} - \frac{1}{n} (\mathbf{1}_q \mathbf{1}_n^T W + W_{n+k} \mathbf{1}_n \mathbf{1}_n^T) + \frac{1}{n^2} \mathbf{1}_q \mathbf{1}_n^T W \mathbf{1}_n \mathbf{1}_n^T \right\} U D_k^{-1} \tag{2}$$

Finally, the last step involves using the scores of the new species as explanatory variables to calculate predictions. Note that using scores obtained from species found to be outside the originally established phylogeny to make predictions involves extrapolation beyond the known range of phylogenetic variation of traits and should thus be avoided. Besides those involving phylogenetic eigenfunctions, other approaches (based, for instance, on generalized least-squares regression or autoregression) have been proposed to test for phylogenetic signals (e.g., Blomberg et al. 2003, Zheng et al. 2009) and to estimate trait values (e.g., Martins and Hansen 1997, Garland and Ives 2000, Rohlf 2001, Bokma 2008).

*Constructing phylogenetic models through cross-validation*

The previously described framework provides the possibility of using cross-validation as an alternative to forward stepwise multiple regression to obtain a phylogenetically explicit predictive model. Cross-validation allows a straightforward assessment of the ability of

the approach to make predictions for new species while avoiding the issue of over-fitting the model. Such an approach involves (1) removing one species at a time from an original data set, (2) calculating linear model coefficients (**b**) using the remaining species, (3) predicting the value of the response from the removed species, and (4) reiterating the first three steps for every species. In that case, linear coefficients (**b**) and standardized linear coefficients (**β**) of the relationship between the response variable *y* (LC<sub>50</sub> in the present study) and the eigenvectors describing phylogeny (**U**) are calculated as:

$$b = U^T [y_i - \bar{y}]$$

$$\beta = \frac{1}{\sqrt{[y_i - \bar{y}]^T [y_i - \bar{y}]}} U^T [y_i - \bar{y}] \tag{3}$$

where *y<sub>i</sub>* is the trait value for a given species *i*,  $\bar{y}$  is the mean trait value, and *y* is one of the response variables, and the predicted values of the response variable (*y*<sub>pred</sub>) are obtained from

$$y_{pred} = \bar{y} + S_{n+k} b. \tag{4}$$

Since the observed values of the response variable are not involved in the calculation of their respective predictions, that approach has the advantage of conserving the independence of the observed and predicted values under the null hypothesis that the response is unrelated to phylogeny. Although that approach allows the use of every single eigenfunction in the models, it does not, however, guarantee that all of them are relevant for making predictions. A simple method to obtain more generalizable models is to truncate the vector of linear coefficients **b** by assigning 0 to its elements that are associated with square standardized linear coefficients (**β**<sup>2</sup>) that are below a threshold chosen to minimize the mean squared error of the model (the mean of the squared differences between predicted and observed values), thereby filtering out irrelevant eigenfunctions. The cross-validation procedure was illustrated by selecting LC<sub>50</sub> values (96 h) for the pesticide carbaryl on all available species in the database and constructing a tree representing their phylogeny from information on their taxonomy.

*Comparing observed with predicted tolerance*

The comparison of observed and predicted tolerance values was performed at two levels. Firstly, a global comparison of these values was made through the examination of the confidence intervals of the slope and intercept of a linear regression line with observed values on the ordinates and predicted values on the abscissa using log<sub>10</sub>-transformed LC<sub>50</sub> values on a molecular basis. Secondly, a comparison was performed at the observation level by calculating the deviation factor *d* of a species *i* as

$$d_i = \begin{cases} 10^{(y_{pred\ i} - y_{obs\ i})} - 1 & \text{if } y_{pred\ i} \geq y_{obs\ i} \\ 1 - 10^{(y_{obs\ i} - y_{pred\ i})} - 1 & \text{if } y_{pred\ i} < y_{obs\ i} \end{cases} \tag{5}$$

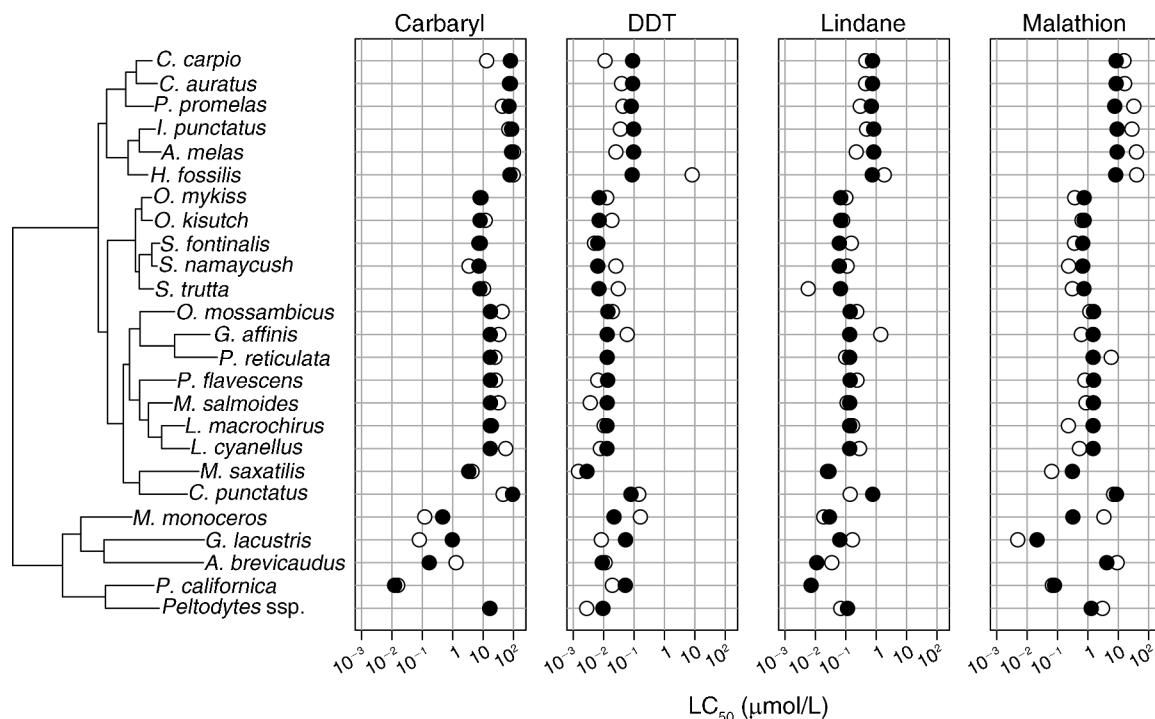


FIG. 2. The model for the 25 aquatic species for the pesticides carbaryl, malathion, DDT, and lindane. Open circles and solid circles represent the observed and fitted values, respectively. For the full species names, see *Results: Phylogenetic analysis*.  $LC_{50}$  is the lethal concentration that kills 50% of individuals of a population over a specific amount of time.

where  $y_{obs}$  are the observed values and  $y_{pred}$  are those predicted by the model, both on a  $\log_{10}$  scale. The deviation factor is the number of times tolerance is overestimated (positive values) or underestimated (negative values) by the model. For example, a value close to 0 means that the tolerance observed for a species is in close agreement with that predicted by the phylogenetic model. Similarly, a value of +10 means that the tolerance observed for a species is 10 times lower than that predicted by the model while a value of -2 means that the tolerance observed for a species is twice as high as that predicted by the model.

All calculations and statistical analyses were performed using the R language and environment (version 2.10.1; R Development Core Team 2010). Database queries were done using package RMySQL (version 0.7-4; James and DebRoy 2009) and phylogenetic analyses with package ape (version 2.4-1; Paradis et al. 2004).

## RESULTS

### Data mining

The best data set that we found involved the pesticides malathion (CAS: 121-75-5), DDT (CAS: 50-29-3), lindane (CAS: 58-89-9), and carbaryl (CAS: 63-25-2), and 25 aquatic animal species, including 20 bony fish, 3 crustacean, and 2 insect species (see Appendix: Table A1 for details). We built a first model (hereafter referred as model 1) using these data. In order to study the response of models to a varying number of substances with respect

to species, we assembled three additional data sets by adding chemical substances, which resulted in a subsequent reduction in the number of species. We obtained the first additional data set by adding parathion (CAS: 56-38-2; 18 species: 13 fishes, 3 crustaceans, and 2 insects), the second data set by further adding dieldrin (CAS: 60-57-1; 14 species: 10 fishes, 3 crustaceans, and 1 insect), and the third data set by further adding rotenone (CAS: 83-79-4) and toxaphene (CAS: 8001-35-2; 11 species: 9 fishes, 1 crustacean, and 1 insect). The three additional models built from these data sets are hereafter referred as models 2, 3, and 4, respectively.

DNA sequences were available for 23 of the 25 species and the most widespread were those for the cytochrome oxidase subunit 1 (COX1; 19 species) and the mitochondrial large (21 species) and small (18 species) ribosomal RNA subunits (see Appendix: Table A2 for details). On average, 20.56 of the 44 sequences were available at a specific level (range, 5–44). We completed that set of sequences by borrowing sequences from other species within the same genus (2.08 sequences on average) or family (3.16 sequences on average), while an average of 18.2 sequences remained missing. The resulting super alignment included from 3246 to 24 433 base pairs (median, 16 503 base pairs).

### Phylogenetic analysis

The tree obtained from DNA sequences placed most species within their known taxonomic group (Fig. 2).

TABLE 1. Statistical test results associated with the phylogenetic models describing among-pesticides and among-species variation of  $LC_{50}$  (96 h), and their associated coefficient of multiple determination ( $R^2$ ) and adjusted coefficients of multiple determination ( $R^2_{adj}$ ).

Factor	<i>F</i>	df	<i>P</i>		$R^2$	$R^2_{adj}$	
			Test-wise	Family-wise		Individual factor	All factors
Model 1: 4 pesticides, 25 species							
Pesticide	129.407	3	<0.0001		0.602	0.590	0.847
Phylogeny	23.815	5	<0.0001	<0.0001	0.185	0.141	0.000
Interaction	16.558	3	<0.0001	<0.0001	0.077	0.048	0.000
Residual		88			0.136		
Model 2: 5 species, 18 pesticides							
Pesticide	33.489	4	<0.0001		0.478	0.453	0.683
Phylogeny	18.429	3	<0.0001	<0.0001	0.197	0.169	0.000
Interaction	10.169	1	0.002	0.03	0.036	0.025	0.000
Residual		81			0.289		
Model 3: 6 pesticides, 14 species							
Pesticide	22.222	5	<0.0001		0.520	0.489	0.612
Phylogeny	12.070	1	0.0009	0.006	0.056	0.045	0.000
Interaction	14.470	1	0.0003	0.003	0.068	0.056	0.000
Residual		76			0.356		
Model 4: 8 pesticides, 11 species							
Pesticide	26.296	7	<0.0001		0.562	0.524	0.734
Phylogeny	18.460	1	<0.0001	0.002	0.056	0.045	0.734
Interaction	16.314	3	<0.0001	<0.0001	0.150	0.119	0.734
Residual		76			0.232		
Model 1T: 4 pesticides, 25 species†							
Pesticide	102.061	3	<0.0001		0.602	0.590	0.805
Phylogeny	24.485	3	<0.0001	<0.0001	0.144	0.117	0.805
Interaction	13.013	3	<0.0001	<0.0001	0.077	0.048	0.805
Residual		90			0.177		

† The additional model 1T corresponds to model 1 but was constructed on a phylogeny obtained from taxonomy rather than from molecular characters.

The tree was rooted at the separation between arthropods and (bony) fish. The first separation on the arthropod subtree occurred between crustaceans and insects, and the second for crustaceans at the subordinal level between eucarids (the speckled shrimp, *Metapenaeus monoceros*) and peracarids (represented by orders isopoda, *Asellus brevicaudus*, and Amphipoda, *Gammarus lacustris*). On the fish subtree, the first separation occurred between ostariophysians and the remaining two teleost suborders, i.e., Protacanthopterygii and Acanthopterygii. Within the ostariophysians the separation first occurred at the ordinal level between cypriniforms and siluriforms, each represented by a single family (Cyprinidae and Ictaluridae, respectively). The second separation on the fish subtree occurred between protacanthopterygians, which is represented by the the family Salmonidae (order Salmoniformes), and acanthopterygians. On the salmonids subtree, the first separation occurred between genus *Oncorhynchus* (rainbow trout and coho salmon) and genera *Salmo* (brown trout) and *Salvelinus* (brook trout and lake trout), with the second separation occurring between the latter genera. Discrepancies of the constructed phylogeny with respect to taxonomy occurred on the acanthopterygians subtree. First, the striped bass (*Morone saxatilis*, family Moronidae) and spotted snakehead (*Channa punctatus*,

family Channidae) separated from other species of order perciformes rather than the species from order Cyprinodontiformes (both family Poeciliidae: the mosquito-fish, *Gambusia affinis* and the guppy, *Poecilia reticulata*), as expected by taxonomy. Cyprinodontiformes species remained clustered within perciforms up to the subordinal level where they separate from the Mozambique tilapia (*Oreochromis mossambicus*, family Cichlidae). Following taxonomy, the striped bass and spotted snakehead were expected separate from other perciforms at the subordinal level. These apparent discrepancies may outline the limit of the current DNA data set for reconstructing the phylogeny of these species; we explored their possible impact on the modeling approach herein described by recalculating model 1 using a tree obtained from taxonomy (hereafter referred to as “model 1T”).

#### Phylogenetic models of tolerance

The models describing  $LC_{50}$  variability among pesticides (where “ $LC_{50}$ ” means lethal concentration for 50% of the population over a specified period) and as a function of species’ phylogenetic structure explained from 61% (model 3) to 85% (model 1, Fig. 2) of the observed variation in tolerance to pesticides (Table 1). By comparison, a model using the mean  $LC_{50}$  of all

organisms for each of the pesticides (i.e., factor pesticide) only explained from 45% (model 2) to 49% (model 1) of that variability. The addition of phylogenetic information thus represents improvements ranging from 12% (model 3) to 26% (model 1) of the total  $LC_{50}$  variability, with phylogeny explaining from 24% (model 3) to 63% (model 1) of  $LC_{50}$  variability within pesticides. Model 1T slightly differed from model 1, but led to similar conclusions.

We found 67 species whose  $LC_{50}$  values (96 h) for pesticide carbaryl were available to illustrate the cross-validation procedure. These species included 35 fish, 18 crustaceans, 8 insects, 4 mollusks, 1 amphibian, and 1 annelid. We estimated the  $\beta^2$  threshold for the truncation of the vector of linear coefficients graphically as 0.00052 from a plot of cross-validated mean standard error obtained by repeating the calculations for thresholds ranging from 0 (all eigenfunctions retained) to 0.015 (the expected  $\beta^2$  if all 67 eigenfunctions were equally relevant) in steps of 0.00001. The resulting cross-validated models explain 83% of the observed variation of the  $\log_{10} LC_{50}$  for carbaryl among these species (adjusted  $R^2$ ; Fig. 3). The regression slope (1.06; 95% confidence limits 0.94 and 1.18) and intercept ( $-0.01$ ; 95% CL  $-0.15$ , 0.13) of the relationship between predicted and observed value was consistent with those of a 1:1 relationship and are not suggestive of a substantial prediction bias by the approach. Its ability to predict  $LC_{50}$  accurately within taxonomic group differed among high-order taxonomic groups, with the model representing from only 13% ( $P > 0.05$ ) of  $\log_{10} LC_{50}$  variability among mollusk species up to 81% ( $P = 0.001$ ) of that among insects species (fish, 40% and  $P < 0.0001$ ; crustaceans, 68% and  $P < 0.0001$ ). The median deviation factor was 0.09 (range  $-46$  to 10), and ranged from  $-1.84$  (mollusks) to 0.57 (insects) to 0.11 (crustaceans) to 0.05 (fish) among the four groups with more than one representative species (Fig. 4). Overall, predictions for 64 out of 67 species (95.5%) had a deviation factor between  $-10$  and  $+10$  whereas a deviation factor between  $-1$  and  $+1$  was obtained for 41 (61.2%) species. The fish was the group whose tolerance was the most accurately represented by the models (median absolute deviation factor: 0.70), followed by crustaceans (0.94), insects (1.10), and mollusks (2.37).

#### DISCUSSION

The approach herein described exemplifies how phylogeny could be used to predict tolerance to pesticides and other chemical substances. In spite of the relatively modest number of representative species available, the results of our present study suggest that the phylogenetic structuring of tolerance, quantified in terms of  $LC_{50}$ , accounted for almost one fourth to almost two thirds of the residual variation within sets of 4–8 pesticides. When cross-validated against a single pesticide, carbaryl, the phylogenetic prediction ap-

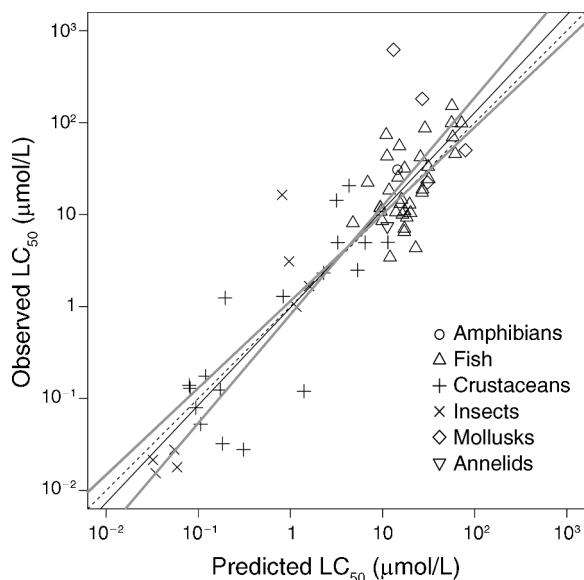


FIG. 3. Relationship between predicted and observed  $\log_{10}(LC_{50})$  for carbaryl. The regression line is solid black; confidence limits of the slope are solid gray; the 1:1 line is the dashed line.

proach provided good estimates of observed  $LC_{50}$  values taken from published laboratory studies, using a reasonable amount of empirical information. Given the ever-increasing availability of molecular information, more particularly in the form of DNA sequences, these results highlight an opportunity to stretch our current usage of the existing tolerance data through phylogenetic-based estimation for species of unknown tolerances. The phylogenetic modeling framework developed in the present study seems, at least under certain circumstances, robust to discrepancies in its prediction basis (i.e., the phylogenetic tree), as illustrated by similarity of the results obtained by model 1 and model 1T, which was based on taxonomy. The robustness of phylogenetic models towards misspecified phylogenies has also been recently demonstrated for the phylogenetic generalized least-squares regression, another method to construct phylogenetic models (Stone 2011).

The approach we used was, in part, borrowed from that of the phylogenetic comparative method, whose purpose is to study the relationships between traits by means of comparisons across species, while correcting for their respective phylogenetic autocorrelation. In our present study the fraction of trait variation that is organized with respect to phylogeny is exploited for making predictions. We ought to mention here that autocorrelation implies the violation of the assumption of independence of observations and may thus affect the outcome of statistical tests. It has been recognized that phylogenetic autocorrelation may render invalid the statistical tests of correlation between species traits (Feldsenstein 1985). This represents a serious shortcom-

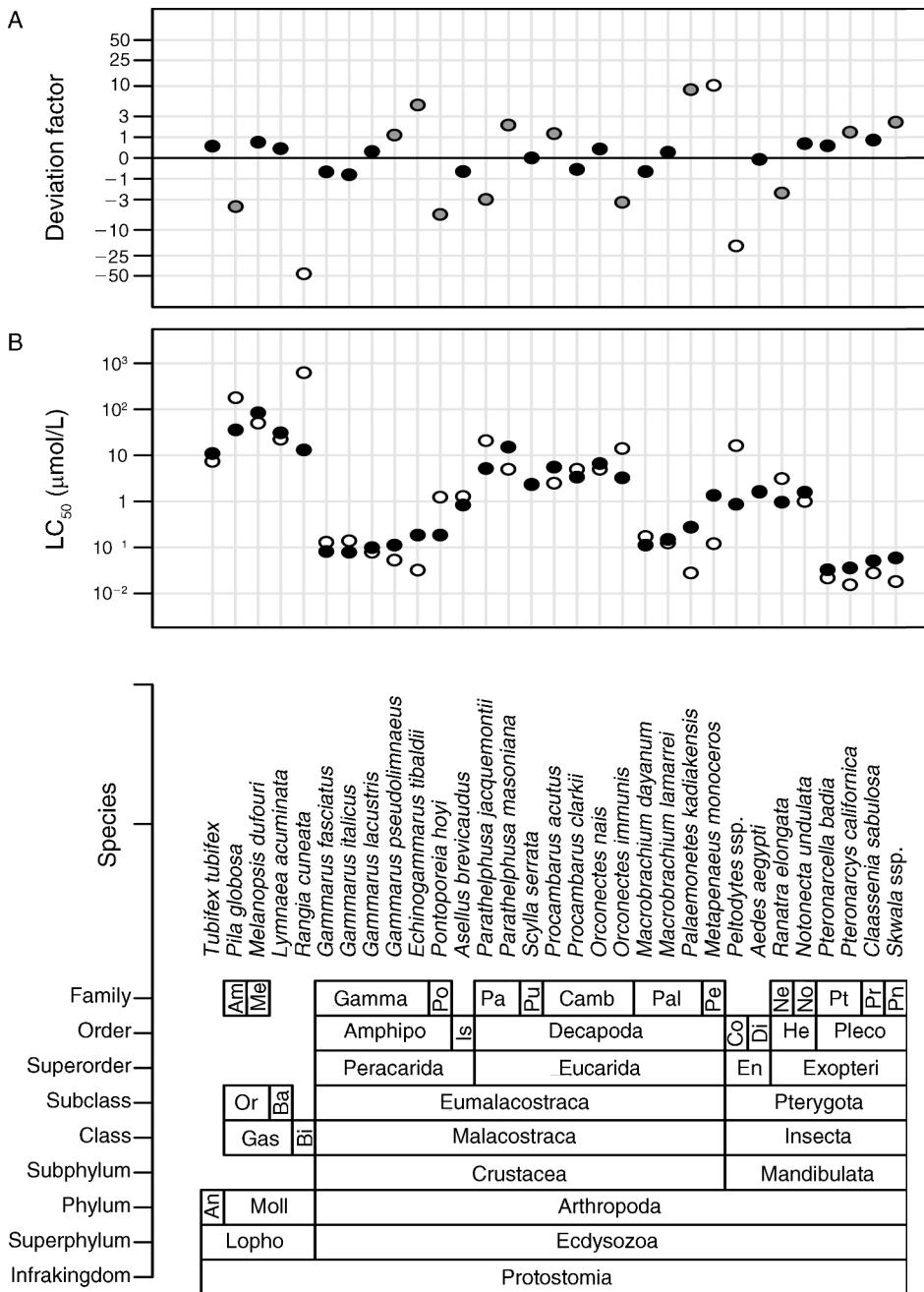


FIG. 4. (A) Deviation factor, i.e., the number of times tolerance is over- or underestimated by the phylogenetic model (overestimation, positive values; underestimation, negative values). Symbol key: open circle, absolute value > 10; shaded circle, 1 < absolute value < 10; solid circle, absolute value < 1. (B) The LC<sub>50</sub> values with respect to their taxonomy. Symbol key: open circle, observed values; solid circle, predicted values. Species abbreviations key: for *superphylum*, Lopho = Lophotrochozoa; for *phylum*, An = Annelida, Moll = Mollusca; for *class*, Gas = Gastropoda, Bi = Bivalvia, Am = Amphibia; for *subclass*, Or = Orthogastropoda, Ba = Basommatophora; for *superorder*, En = Endopterygota, Exopteri = Exopterygota, Proacantho = Protacanthopterygii; for *order*, Amphipo = Amphipoda, Is = Isopoda, Co = Coleoptera, Di = Diptera, He = Hemiptera, Pleco = Plecoptera, Cy = Cyprinodontiformes, Salmonifo = Salmoniformes, Silu = Siluriformes; and for *family*, Am = Ampullariidae, Me = Melanopsidae, Gamma = Gammaridae, Po = Pontoporeiidae, Camb = Cambaridae, Pal = Palaemonidae, Pe = Penaeidae, Ne = Nepidae, No = Notonectidae, Pt = Pteronarcyidae, Pr = Perlidae, Pn = Perlodidae, Os = Osphronemidae, Ci = Cichlidae, Pc = Percidae, Te = Terapontidae, Cent = Centrarchidae, Mo = Moronidae, Ch = Channidae, Cl = Clariidae, He = Heteropneustidae, Ic = Ictaluridae).

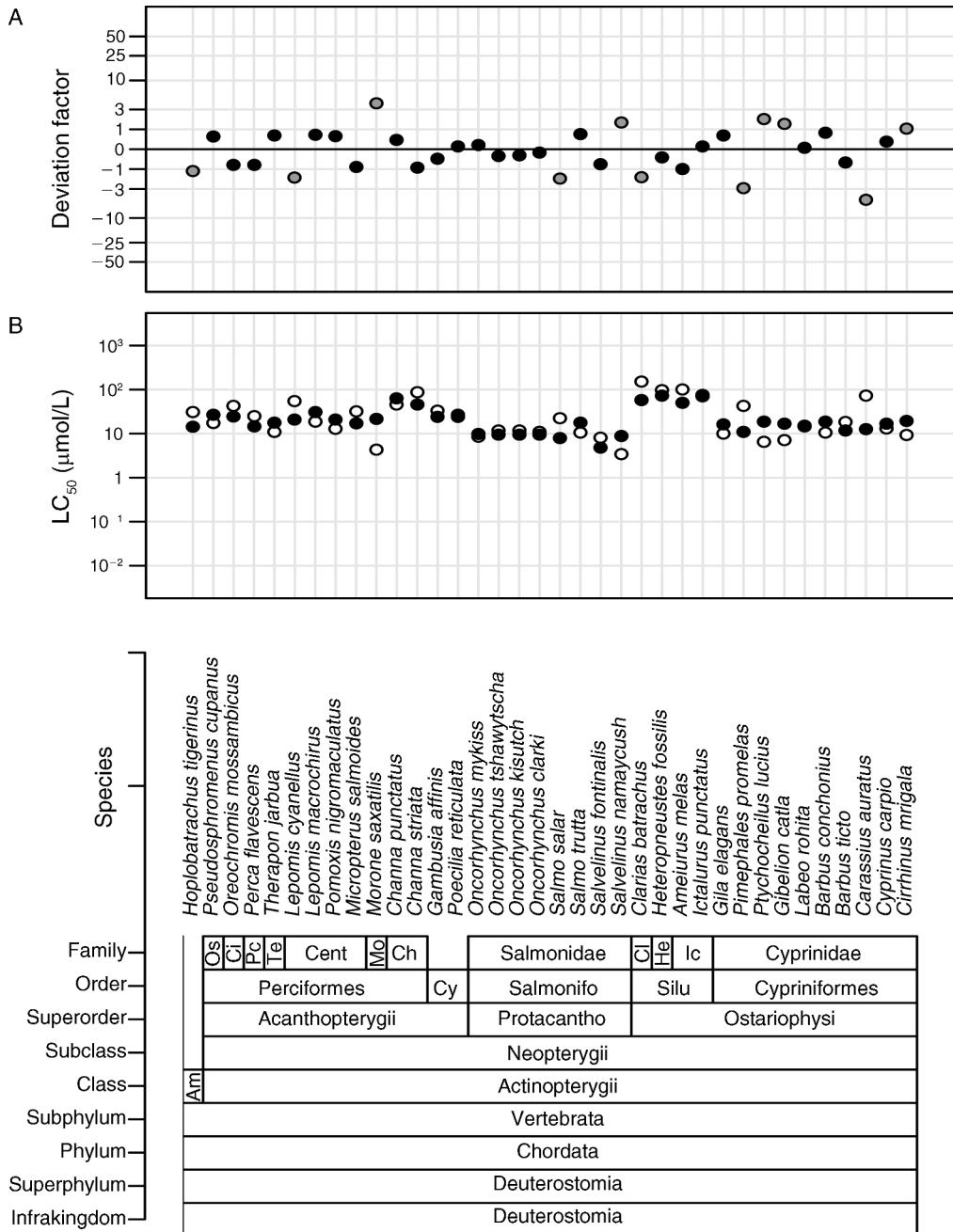


FIG. 4. Continued.

ing that sometimes fails to be addressed when using character correlation for predicting tolerance from other species' traits (e.g., Baird and Van den Brink 2007). As a potential solution, a model may use a phylogeny in conjunction with auxiliary traits related to tolerance to pesticides, possibly enhancing the capacity of the former. When constructing models involving auxiliary traits, however, one has to keep in mind that any trait used as an explanatory variable may itself be phylogenetically autocorrelated. For example, body size has

been shown to be related to a wide range of physiological and ecological attributes (Peters 1983) and may affect tolerance as well. However, the magnitude of body size is heavily structured by phylogeny at large scale and body size may also vary markedly, but within similar orders of magnitude, at smaller phylogenetic scales (e.g., within a family). Since both the tolerance and body size may be driven by the same phylogenetic structures, the parameters (intercept and slope) of a regression involving these traits are likely to be biased

and not representative of the general relationship between them. A general solution when integrating auxiliary traits in models is to use phylogenetic eigenfunctions, which correspond to the eigenfunctions selected for the phylogenetic model, as co-variables when estimating the relationship between the response (i.e., tolerance) and the auxiliary trait (e.g., body mass). Using phylogeny in that manner allows one to partial-out the phylogenetic components of the variation of these species traits before using them as explanatory variables. For instance, body size sometimes varies greatly during the ontogeny of organisms such as aquatic animals, and therefore is irrespective of their phylogeny. Hence, if tolerance to a given pollutant is related with body size, and one builds a model using many individuals of different sizes to represent each species, an important portion of the variation observed for tolerance cannot be represented using a phylogeny, but will be suitably accounted for by body size. In such a situation, a model involving body size as an auxiliary trait would explain a greater portion of the variation in tolerance than one involving phylogeny alone.

The ability of a phylogenetic model to make reliable predictions for a given taxonomic group may not only depend on the number of representative species involved in its construction, but also on the structure of the trait variation along the tree used to represent the phylogeny. For instance, 81% of the variability in the tolerance to carbaryl among insect species was explained by the phylogeny. The relatively good accuracy of the model for predicting the tolerance of insect species to carbaryl is driven by the small tolerance of the four plecopteran species (mean  $LC_{50} = 0.020 \mu\text{mol/L}$ ) compared to that of hemipterans (mean  $LC_{50} = 1.75 \mu\text{mol/L}$ ). Moreover, the large-scale component of the phylogenetic tolerance signal did accurately predict the tolerance of the only amphibian species (the Indian bullfrog; *Hoplobatrachus tigerinus*) to carbaryl from that of the other vertebrate species (fish). On the other hand, the poor performance of the model at predicting the tolerance among mollusks is seemingly the consequence of its inability to predict the greater tolerance of the Atlantic rangia (*Rangia cuneata*), the only bivalve species available, with respect to the other three gastropod species. These examples illustrate the two main requirements of the phylogenetic approach to accurately model the value of a trait such as tolerance: the accuracy of the method is dependent both on the degree of the phylogenetic autocorrelation of the trait (i.e., how much of the trait value is inherited from ancestral species) and the adequacy of the sampling (i.e., the number of members within taxonomic groups among which large differences in trait value are observed or expected). For instance, a phylogenetic model is expected to be inaccurate at evidencing very sensitive or tolerant species pertaining to a highly variable and poorly sampled genus. To this end, it is noteworthy that phylogenetic models cannot predict instances where outstandingly resistant populations

arise by natural selection, such as resistance to pesticides (Ferro 1993, Nandula 2010) or cases of populations living in heavily polluted environments and showing high tolerance to local pollutants (e.g., Nacci et al. 2010). In these cases, a phylogenetic model can nevertheless be useful as a baseline to qualify organisms as resistant or sensitive when their observed tolerances are higher or lower than predicted by the model, respectively. Also noteworthy is the fact that the approach described in our present study carries the assumption, which is common among statistical modeling methods, that the set of species under study forms a representative sample of a larger group of species for which we want to estimate tolerance (i.e., the statistical population). In some groups, the tolerance of ubiquitous species occurring close to—and/or bearing economical value for—humans, may be better studied than that of rare species. Hence, a model involving a sample of exceptionally tolerant (or sensitive) species will consistently overestimate (or underestimate) the tolerance of species for whom tolerance data are not available.

Besides its direct application for predicting a single toxicological effect and endpoint, the approach described in our present study remains applicable in a multiple-effects or multiple-endpoints model framework. Here we will suggest two approaches by which it can be achieved, although others may be applicable. The first possibility would be to use multivariate regression of a species  $\times$  effect or species  $\times$  endpoints response matrix instead of a single response vector as used in the present study. Such a relatively simple approach allows one to obtain several models describing the different effects and/or endpoints at once. The second, more elaborate possibility would be to combine the information on many different endpoints for a given species and calculate metrics describing their relationships to one another (e.g., the log ratio between concentration for observing effects  $x$ ,  $y$ ,  $z$  and  $LC_{50}$ , over the same amount of time), or with respect to a common tolerance baseline, and then applying multivariate regression to the resulting species  $\times$  metrics response matrix. For both approaches, it would be possible to further the analysis of the results obtained by subjecting their resulting multivariate fitted and residual values to principal-components analysis. The combination of these two methods, multivariate regression and principal-components analysis, is known in community ecology as redundancy analysis (Rao 1964, Legendre and Legendre 1998).

Although the phylogenetic-eigenfunctions approach considered in our present study relies on known chemicals for which toxicity was assessed empirically from bioassays, its flexible nature also allows it to be transposed to other frameworks based, for instance, on toxic modes of action (TMOA) or quantitative structure–activity relationships (QSAR; Russom et al. 1997, Schultz et al. 2003, von der Ohe et al. 2005, de Roode et al. 2006, Ajmani et al. 2009). TMOA refers to the

metabolic function that is the most adversely disturbed by a given chemical and most readily leads to the observed effect on the whole organism. Hence, different TMOA can be used as levels of a linear-model factor, with individual chemicals acting through the same mode nested within its respective level (i.e., its respective mode of action). Models thus obtained could provide insight on how the sensitivity towards particular TMOA is structured into phylogeny and which groups are the most or the least susceptible, etc. QSAR models seek to predict the biological activity of a compound from descriptors of its chemical structure. Since biological activity may vary among organisms as a consequence of their particular biochemical traits, it is conceivable that including phylogenetic eigenfunctions as a new set of parameters in QSAR models may improve their ability to predict the impact of new compounds on organisms from a given set of taxonomic groups. If such an approach proves successful, it would provide environmental protection agencies with more dependable tools to more readily screen across the growing list of emerging industrial compounds. Furthermore, organism-specific QSAR may benefit the chemical industry by providing insights on the theoretical innocuousness of compounds under development on the organisms that would specifically be exposed to it.

#### ACKNOWLEDGMENTS

We are thankful to Cândida Shinn, Steven C. Walker, and two anonymous reviewers for their helpful comments on earlier versions of the manuscript. This study was supported by the Marie Curie Research Training Network—Keybioeffects (MRTN-CT-2006-035695) and the Integrated Project MODELKEY (contract 511237-GOCE) of the 6th framework program of the European Commission. G. Guénard also received support from the “Fond Québécois de Recherche sur la Nature et la Technologie (FQRNT)” and P. von der Ohe received financial support through a “Deutsche Forschungsgesellschaft” (DFG—Bonn, Germany) fellowship (PAK 406/1).

#### LITERATURE CITED

- Ajmani, S., K. Jadhav, and S. Kulkarni. 2009. Group-based QSAR (G-QSAR): Mitigating interpretation challenges in QSAR. *QSAR and Combinatorial Science* 28:36–51.
- Baird, D. J., and P. J. Van den Brink. 2007. Using biological traits to predict species sensitivity to toxic substances. *Ecotoxicology and Environmental Safety* 67:296–301.
- Blomberg, S., T. Garland, and A. Ives. 2003. Testing for phylogenetic signal in comparative data: behavioral traits are more labile. *Evolution* 57:717–745.
- Bokma, F. 2008. Detection of “punctuated equilibrium” by Bayesian estimation of speciation and extinction rates, ancestral character states, and rates of anagenetic and cladogenetic evolution on a molecular phylogeny. *Evolution* 62:2718–2726.
- Buchwalter, D. B., D. J. Cain, C. A. Martin, L. Xie, S. N. Luoma, and T. Garland, Jr. 2008. Aquatic insect ecophysiological traits reveal phylogenetically based differences in dissolved cadmium susceptibility. *Proceedings of the National Academy of Science USA* 105:8321–8326.
- Corey, S. J., and T. A. Waite. 2008. Phylogenetic autocorrelation of extinction threat in globally imperilled amphibians. *Diversity and Distributions* 14:614–629.
- Crommentuijn, T., M. D. Polder, and E. J. van de Plassche. 1997. Maximum permissible concentrations and negligible concentrations for metals, taking background concentrations into account. Report number 601501001. National Institute of Public Health and the Environment (RIVM), Bilthoven, The Netherlands. (<http://www.rivm.nl/bibliotheek/rapporten/601501001.pdf>)
- de Rooede, D., C. Hoekzema, S. de Vries-Buitenweg, B. van de Waart, and J. van der Hoeven. 2006. QSARs in ecotoxicological risk assessment. *Regulatory Toxicology and Pharmacology* 45:24–35.
- Desdevises, Y., P. Legendre, L. Azouzi, and S. Morand. 2003. Quantifying phylogenetically structured environmental variation. *Evolution* 57:2647–2652.
- de Zwart, D. 2002. Observed regularities in SSDs for aquatic species. Pages 133–152 in L. Posthuma, G. W. Suter, and T. Traas, editors. *Species sensitivity distributions in ecotoxicology*. Lewis Publishers, Boca Raton, Florida, USA.
- Diniz-Filho, J. A. F. 2001. Phylogenetic autocorrelation under distinct evolutionary processes. *Evolution* 55:1104–1109.
- Diniz-Filho, J. A. F., C. E. R. de Sant’Ana, and L. M. Bini. 1998. An eigenvector method for estimating phylogenetic inertia. *Evolution* 52:1247–1262.
- Edgar, R. C. 2004. MUSCLE: multiple sequence alignment with high throughput. *Nucleic Acids Research* 32:1792–1797.
- Felsenstein, J. 1981. Evolutionary trees from DNA sequences: a maximum likelihood approach. *Journal of Molecular Evolution* 17:368–376.
- Felsenstein, J. 1985. Phylogenies and the comparative method. *American Naturalist* 125:1–15.
- Felsenstein, J., and G. A. Churchill. 1996. A hidden Markov model approach to variation among sites in rate of evolution. *Molecular Biology and Evolution* 13:93–104.
- Ferro, D. N. 1993. Potential for resistance to *Bacillus thuringiensis*: Colorado potato beetle (Coleoptera: Chrysomelidae)—a model system. *American Entomologist* 39:38–44.
- Garland, T., Jr., and A. R. Ives. 2000. Using the past to predict the present: confidence intervals for regression equations in phylogenetic comparative methods. *American Naturalist* 155:346–364.
- Gower, J. C. 1966. Some distance properties of latent root and vector methods used in multivariate analysis. *Biometrika* 53:325–338.
- Gower, J. C. 1969. Adding a point to vector diagrams in multivariate analysis. *Biometrika* 55:582–585.
- Holm, S. 1979. A simple sequentially rejective multiple test procedure. *Scandinavian Journal of Statistics* 6:65–70.
- James, D. A., and S. DebRoy. 2009. RMySQL: R interface to the MySQL database (R package version 0.7-4). (<http://CRAN.R-project.org/package=RMySQL>)
- Jeffree, R. A., F. Oberhansli, and J.-L. Teyssie. 2010. Phylogenetic consistencies among condrichthyan and teleost fishes in their bioaccumulation of multiple trace elements from seawater. *Science of the Total Environment* 408:3200–3210.
- Lande, R., and S. J. Arnold. 1983. The measurement of selection on correlated characters. *Evolution* 37:1210–1226.
- Legendre, P., and L. Legendre. 1998. *Numerical ecology*. Second English edition. Elsevier Science, Amsterdam, The Netherlands.
- Maddison, D. R., K. S. Schulz, and W. P. Maddison. 2007. The tree of life web project. *Zootaxa* 1668:19–40.
- Martins, E. P., and T. F. Hansen. 1997. Phylogenies and the comparative method: a general approach to incorporating phylogenetic information into the analysis of interspecific data. *American Naturalist* 149:646–667.
- Nacci, D. E., D. Champlin, and S. Jayaraman. 2010. Adaptation of the estuarine fish *Fundulus heteroclitus* (Atlantic killifish) to polychlorinated biphenyls (PCBs). *Estuaries and Coasts* 33:853–864.

- Naudula, V. K. 2010. Glyphosate resistance in crops and weeds: history, development, and management. Wiley Publishing, Hoboken, New Jersey, USA.
- Paradis, E., J. Claude, and K. Strimmer. 2004. APE: analyses of phylogenetics and evolution in R language. *Bioinformatics* 20:289–290.
- Peters, R. H. 1983. The ecological implications of body size. Cambridge University Press, New York, New York, USA.
- Posthuma, L., G. W. Suter, and T. Traas. 2002. Species sensitivity distributions in ecotoxicology. Lewis Publishers, Boca Raton, Florida, USA.
- R Development Core Team. 2010. R: a language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria. (<http://www.R-project.org>)
- Rao, C. R. 1964. The use and interpretation of principal component analysis in applied research. *Sankhyā, Series A* 26:329–358.
- Rice, P., I. Longden, and A. Bleasby. 2000. EMBOSS: the European molecular biology open software suite. *Trends in Genetics* 16:276–277.
- Russom, C. L., S. P. Bradbury, S. J. Broderius, D. E. Hammermeister, and R. A. Drummond. 1997. Predicting modes of toxic action from chemical structure: acute toxicity in the fathead minnow (*Pimephales promelas*). *Environmental Toxicology and Chemistry* 16:948–967.
- Schultz, T. W., M. T. D. Cronin, and T. I. Netzeva. 2003. The present status of QSAR in toxicology. *Journal of Molecular Structure THEOCHEM* 622:23–38.
- Stone, E. A. 2011. Why the phylogenetic regression appears robust to tree misspecification. *Systematic Biology* 60:245–260.
- Svanbäck, R., and D. I. Bolnick. 2007. Intraspecific competition drives increased resource use diversity within a natural population. *Proceedings of the Royal Society of London B* 274:839–844.
- Tomlin, C. D. S., editor. 1997. Pesticide manual. British Crop Protection Council, Farnham, UK.
- USEPA [United States Environmental Protection Agency]. 1984. AQUIRE, aquatic information retrieval toxicity database. Project description, guidelines, and procedures. By R. C. Russo and A. Pilli. EPA 600/8-84-021. U. S. Environmental Protection Agency, Environmental Research Laboratory, Duluth, Minnesota, USA. (<http://www.epa.gov/nscep/index.html>)
- von der Ohe, P. C., R. Kuhne, R. U. Ebert, R. Altenburger, M. Liess, and G. Schuurmann. 2005. Structural alerts – a new classification model to discriminate excess toxicity from narcotic effect levels of organic compounds in the acute daphnid assay. *Chemical Research in Toxicology* 18:536–555.
- von der Ohe, P. C., and M. Liess. 2004. Relative sensitivity distribution of aquatic invertebrates to organic and metal compounds. *Environmental Toxicology and Chemistry* 23:150–156.
- Zheng, L., A. R. Ives, T. Garland, B. R. Larget, Y. Yu, and K. Cao. 2009. New multivariate tests for phylogenetic signal and trait correlations applied to ecophysiological phenotypes of nine *Manglietia* species. *Functional Ecology* 23:1059–1069.

#### APPENDIX

Two tables containing the log<sub>10</sub>-transformed LC<sub>50</sub> values used in the study and the information about the DNA sequences used to estimate phylogeny (*Ecological Archives* A021-142-A1).

- 1 **Guillaume Guénard, Peter Carsten von der Ohe, Dick de Zwart, Pierre Legendre, and Sovan**
- 2 **Lek. Using phylogenetic information to predict species tolerances to toxic chemicals. *Ecological***
- 3 ***Applications.***
- 4 Appendix A. Two tables containing the  $\log_{10}$ -transformed  $LC_{50}$  values of used in the study and the
- 5 information about the the DNA sequences used to estimate phylogeny.

6 Table A1. Values of  $\log_{10}\{LC_{50} (\mu\text{mol}\cdot\text{L}^{-1})\}$  obtained from the database for eight pesticides and  
7 available to build the models. When more than a single value was available, the number of averaged  
8 values is indicated in parenthesis and the standard deviation is shown.

Species	Carbaryl	DDT	Lindane	Malathion	Parathion	Dieldrin	Rotenone	Toxaphene
<i>Asellus brevicaudatus</i> <sup>3</sup>	0.110 ± 0.033 (2)	-1.954	-1.464	0.958	0.589 ± 0.275 (2)	-1.882		
<i>Carassius auratus</i> <sup>4</sup>	1.868 ± 0.051 (2)	-1.413 ± 0.097 (3)	-0.346	1.204 ± 0.306 (2)	0.798	-2.326	0.100	-1.471
<i>Channa punctatus</i> <sup>1</sup>	1.660 ± 0.11 (6)	-0.849 ± 0.376 (2)	-0.858 ± 0.194 (2)	0.845 ± 0.126 (7)			-1.356*	
<i>Cyprinus Carpio</i>	1.117 ± 0.104 (6)	-1.954	-0.336 ± 0.173 (2)	1.183 ± 0.219 (3)	0.465	0.197	-1.317	-2.049
<i>Gambusia affinis</i>	1.521 ± 0.678 (2)	-1.223 ± 0.128 (4)	0.150 ± 0.221 (4)	-0.218	0.443 ± 0.402 (2)	-1.089	-1.365	-1.514 ± 0.201 (4)
<i>Gammarus lacustris</i>	-1.100	-2.079 ± 0.477 (2)	-0.782	-2.309	-1.639 ± 0.282 (2)	0.173 ± 0.091 (2)	0.819	-1.202
<i>Heteropneustes fossilis</i> <sup>2</sup>	1.994 ± 0.006 (2)	0.914	0.266	1.605 ± 0.126 (5)	1.951			
<i>Ameiurus melas</i> <sup>1</sup>	1.997	-1.600 ± 0.275 (2)	-0.657	1.592			-0.006*	-1.853*
<i>Ictalurus punctatus</i> <sup>4</sup>	1.842 ± 0.053 (2)	-1.449 ± 0.121 (6)	-0.315 ± 0.505 (2)	1.450 ± 0.016 (2)	0.959	-1.928	-1.665 ± 0.516 (2)	-2.099 ± 0.382 (3)
<i>Lepomis cyanellus</i> <sup>4</sup>	1.746	-2.109	-0.545	-0.276	0.504	-1.788	-0.447	-1.503
<i>Lepomis macrochirus</i> <sup>4</sup>	1.275 ± 0.142 (6)	-1.986 ± 0.119 (12)	-0.775 ± 0.09 (5)	-0.644 ± 0.108 (3)	-0.420 ± 0.664 (2)	-1.886 ± 0.131 (4)	-0.913 ± 0.324 (3)	-1.981 ± 0.118 (6)
<i>Metapenaeus monoceros</i> <sup>3</sup>	-0.922 ± 0.015 (3)	-0.792	-1.723	0.530	0.573	-1.007		
<i>Micropterus salmoides</i> <sup>4</sup>	1.502	-2.444 ± 0.107 (3)	-0.958	-0.064	0.328	-2.037	-0.444	-2.316
<i>Morone saxatilis</i> <sup>3</sup>	0.637 ± 0.060 (2)	-2.831	-1.600	-1.192 ± 0.108 (4)	-1.214	-1.286		-1.929* ± 0.044 (2)

Species	Carbaryl	DDT	Lindane	Malathion	Parathion	Dieldrin	Rotenone	Toxaphene
<i>Oncorhynchus kisutch</i> <sup>1</sup>	1.072 ± 0.262 (2)	-1.728 ± 0.226 (2)	-1.102	-0.192 ± 0.096 (2)			-0.804*	-1.714*
<i>Oncorhynchus mykiss</i> <sup>4</sup>	0.933 ± 0.115 (6)	-1.895 ± 0.138 (11)	-0.995 ± 0.037 (2)	-0.431 ± 0.065 (6)	0.546 ± 0.136 (2)	-2.329 ± 0.189 (5)	-1.835 ± 0.231 (4)	-1.736 ± 0.327 (3)
<i>Peltodytes sp.</i> <sup>2</sup>	1.215	-2.556	-1.163	0.481	-1.619			
<i>Perca flavescens</i> <sup>1</sup>	1.404	-2.189 ± 0.588 (2)	-0.631	-0.099			-1.119*	-1.538*
<i>Pimephales promelas</i> <sup>4</sup>	1.634 ± 0.113 (3)	-1.363 ± 0.106 (2)	-0.524	1.511 ± 0.063 (3)	0.796 ± 0.111 (2)	-1.689 ± 0.312 (2)	-1.407 ± 0.332 (5)	-1.635 ± 0.125 (4)
<i>Poecilia reticulata</i> <sup>1</sup>	1.386 ± 0.067 (3)	-1.880 ± 0.199 (2)	-1.004 ± 0.256 (2)	0.766 ± 0.206 (2)		-1.918* ± 0.091 (3)		
<i>Pteronarcys californica</i> <sup>4</sup>	-1.813 ± 0.190 (2)	-1.711	-2.137 ± 0.327 (2)	-1.170 ± 0.349 (2)	-1.483 ± 0.262 (3)	-2.882	-0.107 ± 0.091 (2)	-2.255
<i>Salmo trutta</i> <sup>1</sup>	1.019 ± 0.477 (2)	-1.518	-2.233	-0.515				-2.125*
<i>Salvelinus fontinalis</i> <sup>2</sup>	0.909 ± 0.222 (3)	-2.300	-0.817	-0.440	0.715		-1.507* ± 0.557 (2)	
<i>Salvelinus namaycush</i> <sup>2</sup>	0.535	-1.597	-0.958	-0.638	0.819			
<i>Oreochromis mossambicus</i> <sup>1</sup>	1.626	-1.711	-0.640 ± 0.068 (2)	0.062 ± 0.118 (2)		-1.619* ± 0.038 (2)	-0.693*	

9 <sup>1</sup>: species used only for model #1 (and 1T)

10 <sup>2</sup>: species used for models #1 and #2

11 <sup>3</sup>: species used for models #1, #2, and #3

12 <sup>4</sup>: species used for all 4 models

13 \*: unused

14 Table A2. Sources for molecular characters for the 25 aquatic animal species: number of sequences  
 15 found for the species, or substituted from a related species of the same genus or family, total number of  
 16 sequences found and missing, and Genbank accession codes (parenthesis G and F: sequence borrowed  
 17 within genera or families, respectively).

Species name	Taxonomic resolution			Found / Missing	Accession codes
	Species	Genera	Families		
<i>Ameiurus melas</i>	7	0	0	7 / 37	AY184263, AY705821, DQ421854, DQ421876, EU524419
<i>Asellus brevicaudus</i>	0	3	33	36 / 8	AF255701 (G), AF259529 (F), DQ305105 (G), FJ749279 (G), NC_008412 (F)
<i>Carassius auratus</i>	41	0	0	41 / 3	AF047349, EF100727, NC_006580
<i>Channa punctatus</i>	11	2	0	13 / 31	AB196280, AY763724 (G), AY763770 (G), EU216546, EU342184, EU417796, EU836885
<i>Cyprinus carpio</i>	44	0	0	44 / 0	AF133089, NC_001606
<i>Gambusia affinis</i>	39	0	0	39 / 5	AP004422
<i>Gammarus lacustris</i>	5	1	0	6 / 38	AF228046 (G), AY529073, AY926671, AY926784, EF582869
<i>Heteropneustes fossilis</i>	8	0	0	8 / 36	AF520826, AJ876377, DQ119383, FN677932, GQ461897
<i>Ictalurus punctatus</i>	40	0	0	40 / 4	AF021880, NC_003489
<i>Lepomis cyanellus</i>	5	1	0	6 / 38	AB271768 (G), AY115973, AY517733, AY742522, AY742616, EU524705
<i>Lepomis macrochirus</i>	14	0	0	14 / 30	AB167815, AB167816, AY517740, AY742530, AY742623, AY828968, EU524732
<i>Metapenaeus monoceros</i>	0	1	42	43 / 1	AF124597 (F), AY264904 (G), EU920969 (F), NC_002184 (F)
<i>Micropterus salmoides</i>	39	1	0	40 / 4	EU502753 (G), NC_008106
<i>Morone saxatilis</i>	11	0	3	14 / 30	AF147741, AF240746, AY072684, AY138963, AY538941, DQ028057, EU524145, L60529, X74147 (F), X74148 (F), X74149 (F)
<i>Oncorhynchus</i>	40	0	0	40 / 4	AF030250, NC_009263

Species name	Taxonomic resolution			Found / Missing	Accession codes
<i>kisutch</i>					
<i>Oncorhynchus mykiss</i>	43	0	0	43 / 1	AF308735, NC_001717, OMU34341
<i>Oreochromis mossambicus</i>	42	0	0	42 / 2	AF497908, AY597335, DQ397880
<i>Peltodytes ssp.</i>	0	6	0	6 / 38	AJ318668 (G), AY071790 (G), AY071816 (G), AY745649 (G), AY745665 (G), EU797379 (G)
<i>Perca flavescens</i>	8	0	1	9 / 35	AF045357, AY225721, AY520099, AY538950, AY726669, EU524238, NC_008111 (F), Y14728
<i>Pimephales promelas</i>	7	1	0	8 / 36	AF126355, AY102292, AY102302, AY216557 (G), AY430235, AY855349, EU525095, GQ275159
<i>Poecilia reticulata</i>	19	1	0	20 / 24	DQ983928, EF017485, EF017585, EU751921 (G), GQ855720, GU179192
<i>Pteronarcys californica</i>	6	34	0	40 / 4	AY521812, AY521880, EF623110, EF623261, EF623427, EU099983, NC_006133 (G)
<i>Salmo trutta salvelinus fontinalis</i>	40	0	0	40 / 4	DQ009482, NC_010007
<i>Salvelinus fontinalis</i>	39	0	0	39 / 5	NC_000860
<i>Salvelinus namaycush</i>	6	1	0	7 / 37	AF174610, AF297989, DQ451375, EU522418, FJ620124, NC_000861 (G), U61182
Average:	20.56	2.08	3.16	25.80 / 18.20	