

Structure-Related Differences between Cytochrome Oxidase I Proteins in a Stable Heteroplasmic Mitochondrial System

David O.F. Skibinski^{1,*}, Fabrizio Ghiselli², Angel P. Diz³, Liliana Milani², and Jonathan G.L. Mullins¹

¹Institute of Life Science, Swansea University Medical School, United Kingdom

²Department of Biological, Geological, and Environmental Sciences, University of Bologna, Italy

³Department of Biochemistry, Genetics and Immunology, University of Vigo, Spain

*Corresponding author: E-mail: d.o.f.skibinski@swansea.ac.uk.

Accepted: November 13, 2017

Abstract

Many bivalve species have two types of mitochondrial DNA passed independently through the female line (F genome) and male line (M genome). Here we study the cytochrome oxidase I protein in such bivalve species and provide evidence for differences between the F and M proteins in amino acid property values, particularly relating to hydrophobicity and helicity. The magnitude of these differences varies between different regions of the protein and the change from the ancestor is most marked in the M protein. The observed changes occur in parallel and in the same direction in the different species studied. Two possible causes are considered, first relaxation of purifying selection with drift and second positive selection. These may operate in different ways in different regions of the protein. Many different amino acid substitutions contribute in a small way to the observed variation, but substitutions involving alanine and serine have a quantitatively large effect. Some of these substitutions are potential targets for phosphorylation and some are close to residues of functional importance in the catalytic mechanism. We propose that the observed changes in the F and M proteins might contribute to functional differences between them relating to ATP production and mitochondrial membrane potential with implications for sperm function.

Key words: mtDNA, protein structure, amino acid properties, bivalves, doubly uniparental inheritance, heteroplasmy.

Introduction

Animal mitochondrial DNA (mtDNA) shows considerable diversity, including variation in size, structure, and gene content, which might have adaptive significance (Breton et al. 2014). Many species of bivalves have two types of mitochondria, one passed through the female line of descent (F type) and the other through the male line (M type) (Skibinski et al. 1994; Zouros et al. 1994). In this system (doubly uniparental inheritance [DUI]), females pass the F type to the progeny of both sexes, whereas males pass the M type to sons only. Females have F type mitochondria in both soma and germ line, whereas in males the germ line is homoplasmic for the M type with somatic tissues showing varying levels of heteroplasmy (Zouros 2013). There is ongoing research into the evolutionary and functional significance of the two genomes and their link with sex determination (Breton, Stewart, et al. 2011; Zouros 2013). Of potential relevance is that the F and M genomes show a variety of sequence differences. They can be highly diverged, exceeding 50% in some freshwater

mussels and often show major structural differences and rearrangements (Ghiselli et al. 2013; Zouros 2013; Breton et al. 2014). Moreover, M and F type specific mtORFans of putative viral origin and with a proposed role in DUI have been discovered (Breton, Ghiselli, et al. 2011; Ghiselli et al. 2013; Milani et al. 2013). Sequence differences between M and F genomes located in the intergenic regions and in the RNA-coding genes can affect mtDNA transcription and translation, potentially resulting in sex-specific mitochondrial expression, especially in conditions where the two mtDNAs are in homoplasmy, namely, in gonads and gametes (Milani et al. 2014).

Although major structural differences between the F and M genomes may be biologically important, functional differences could also be sought in normal sequence variation in mitochondrial proteins. Amino acid substitution is influenced by a variety of structural and functional constraints (Chelliah et al. 2004). These relate to the physicochemical properties of amino acids and are consistent with evidence that substitutions within chemical groups (conservative) are more

prevalent than those between groups (radical) (Hughes et al. 1990). The latter are more likely to have deleterious effects on structure and function and be removed by purifying selection (Zhang 2000). Information on protein structure and amino acid properties can augment DNA and protein sequence information in evolutionary studies in a variety of ways. Evolutionary trees built from atomic distances in superimposed known protein tertiary structures can be compared with those built from aligned sequences (Johnson et al. 1990; Balaji and Srinivasan 2007). Amino acid substitution tables can be built from different structural environments, for example, solvent accessible and inaccessible regions within proteins of known tertiary structure, and be used for constructing refined and better supported evolutionary trees (Overington et al. 1990; Thorne et al. 1996; Gong and Blundell 2008). In addition, if the tertiary structure is known for one of the sequences in an amino acid alignment it can be used to infer the structural environments for the other sequences (e.g., Mizuguchi et al. 1998; Melvin et al. 2008; Puslednik et al. 2012).

The chemical and physical properties of amino acids can also be used to improve evolutionary models for a variety of purposes (e.g., Koshi and Goldstein 1997; Xia and Li 1998; Liu and Wang 2006), for example, to define patches with functional regions on the surface of proteins (Pettit et al. 2007). Metrics, obtained by substituting amino acid physicochemical property values for the nominal amino acid symbols, can be summed up over the different amino acids within a region of a protein (McClellan and McCracken 2001; Atchley et al. 2005; McClellan 2013). This can aid the identification of causal components underlying the sequence variation in different regions of the protein by using approaches based on analysis of variance (ANOVA) to partition the variation (Atchley et al. 2005).

The aim of this study is to compare the structural features and differences in amino acid property values in the cytochrome oxidase I (COI) protein from the F and M mtDNA genomes. This protein is responsible for the transfer of electrons from cytochrome c to reduce molecular oxygen to water and is the final step in the electron transport chain. In the reaction, eight protons come from the mitochondrial matrix, four of which make water and four protons are released into the intermembrane space. This results in an electrochemical gradient across the membrane, which ATP synthase uses to make ATP (Chen 1988). COI is widely studied both in relation to the action of natural selection (Garvin et al. 2015) and in taxonomic studies (Hebert et al. 2003), and its mechanism of action is well understood (e.g., Tsukihara et al. 2003), though some details are not completely elucidated (Popovic 2013). We provide evidence that the F and M COI proteins show differences in the values of amino acid properties in relation to structural environment, which are concordant in four bivalve species with DUI. The differences include changes in hydrophobicity

and helicity in the intermembrane space and mitochondrial matrix regions of the M protein. We suggest that these changes may affect ATP production and the mitochondrial membrane potential with consequent effects on sperm function and implications for DUI and discuss the evolutionary forces that might be responsible.

Materials and Methods

Sequences Used in Analysis

To select representative sequences for analysis, 44 bivalve F and M COI amino acid sequences were aligned and a maximum likelihood (ML) tree constructed (supplementary fig. S1, Supplementary Material online). Four species each with F and M genomes were chosen within each of the three orders, Mytiloidea, Veneroidea, and Unionoidea. No part of the lineages separating the F and M COI sequences through their common ancestor within each of the chosen species overlaps with the corresponding lineages separating F and M genomes in the other species. In other words, the F and M genomes form monophyletic groups within each species (fig. 1A). This approach allows statistically independent comparisons of the F and M genomes in line with a proposal for the comparative method (Felsenstein 1985). Sequences from the *Mytilus edulis* group (*M. edulis*, *M. trossulus*, and *M. galloprovincialis*) were not chosen because of the complication of role reversals and rearrangements (Zouros 2013), *Mytilus californianus* was chosen instead. *Musculista senhousia* was added as a fourth species because it also satisfies the above criterion of nonoverlapping lineages, and the F and M COI sequences are sufficiently diverged to merit inclusion. The eight sequences used in analysis with NCBI accession numbers are *M. californianus* (Cal) (F: ACV65353 M: ACV65365), *M. senhousia* (Sen) (F: ACY00212 M: ACY00224), *Pyganodon grandis* (Gra) (F: ACQ91058 M: ACQ91071), and *Ruditapes philippinarum* (Phi) (F: BAB83795 M: BAB83782).

The eight bivalve sequences were aligned with a sequence from a known bovine structure (PDB ID: 1V54, Tsukihara et al. 2003) using T-coffee (Notredame et al. 2000). Terminal regions and a few internal sites with gaps were removed from the alignment. The terminal regions comprise on average 22 amino acids per sequence and are similar in length and amino acid sequence between genomes within species. The resulting alignment is 501 amino acid sites long, 13 sites less than the number of amino acids in 1V54 COI. To reconstruct predicted ancestral sequences of the F and M COI sequences within each species, the alignment without gaps was submitted to the FASTML server, which implements ML algorithms for this purpose (Ashkenazy et al. 2012).

Residue and Site-Specific Attributes

Residue-specific and site-specific attributes are used for the comparative analysis of the F and M COI proteins. In the

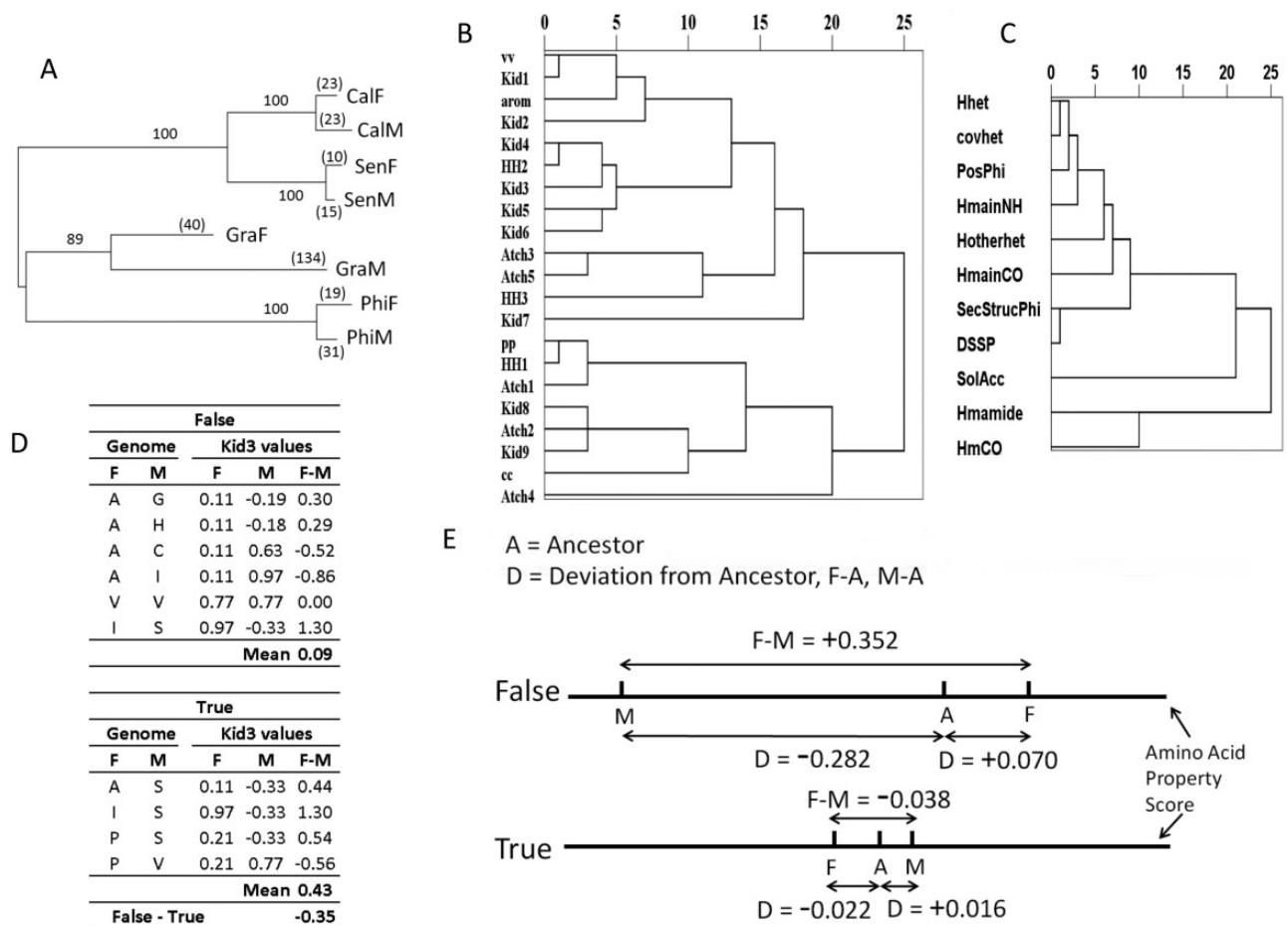


FIG. 1.—Amino acid properties, structural environments and evolution of bivalve COI proteins. (A) ML tree of F and M COI sequences in four bivalve sequences with ancestors constructed with MEGA (Kumar et al. 2016) using the general reversible mtDNA amino acid substitution model and gamma-distributed rate variation among sites (Adachi and Hasegawa 1996). Around 501 ungapped aligned sites were used. Branch lengths are in number of substitutions per site. Support values are based on 1,000 bootstrap samples. Species names abbreviated as in text. The numbers of amino acid differences of F and M proteins from their common ancestor within each species are shown on the terminal branches in brackets. (B) Dendrogram obtained by hierarchical clustering of amino acid properties based on values for 20 amino acids using average linkage between groups. Details of the properties are given in Materials and Methods (see Residue-Specific Attributes) with further information in supplementary table S1, Supplementary Material online. (C) Dendrogram obtained by hierarchical clustering of structural environments based on category assignments for 1V54 sequence. Details of the environments are given in Materials and Methods (see Site-Specific Attributes) with further information in supplementary table S1, Supplementary Material online. (D) Calculation of F–M and False–True for small made-up data set. Amino acid property values for Kid3 are used. (E) Example of breakdown of F–M into deviations from the ancestor in two categories (False and True). Numerical values are taken from table 4A for Hmamide group A. Note distances are not to scale.

alignment of the eight bivalve sequences used, there are $8 \times 501 = 4,008$ individual amino acid residues to which residue-specific attributes can be assigned, and 501 alignment sites to which site-specific attributes can be assigned.

Residue-Specific Attributes

The residue-specific attributes are chemical and physical amino acid property values for the 20 amino acids, derived from several studies, which capture information about the amino acid in a metric value (Sneath 1966; Grantham 1974;

Kidera et al. 1985; Haig and Hurst 1991; Xia and Li 1998; Atchley et al. 2005). We have given abbreviated names to these properties, which with brief description in brackets are cc (side chain composition), pp (polarity), vv (volume), arom (aromaticity), Atch1 (polarity), Atch2 (secondary structure), Atch3 (volume), Atch4 (amino acid composition), Atch5 (charge), Kid1 (bulk), Kid2, Kid3, Kid4 (hydrophobicity), Kid5, Kid6 (beta structure preference), Kid7 (alpha helix preference), Kid8, Kid9 (bend structure preference), HH1 (polarity), HH2 (hydropathy), and HH3 (isoelectric point). Further information on the properties is given in supplementary table

S1, Supplementary Material online. A high value on the metric scale would indicate a tendency toward the specified description. Thus for *w* (volume), tryptophan, a bulky amino acid has the highest value whereas glycine has the lowest value. Clustering of these 21 properties (fig. 1B) broadly reflects differences and similarities described in the cited studies. Many of the properties are derived by factor analysis and thus represent a combination of many individual physico-chemical properties. Correlation and coefficient of determination values among the properties are given in supplementary table S2, Supplementary Material online. Most of the values of coefficient of determination are less than 0.5, suggesting some but not necessarily high redundancy between the properties.

Site-Specific Attributes

To assign the site-specific attributes, the bivalve protein sequences are partitioned into different spatial regions called “structural environments,” which a priori are expected to reflect functional differences. Because the structures of the bivalve proteins are unknown, structural environment was estimated for COI using the known bovine structure 1V54. The program JOY (Mizuguchi et al. 1998) was used for this purpose, partitioning COI into categories for 11 structural environments. For illustration, consider an environment relating to main chain to main chain hydrogen bonding. Using a known or estimated protein structure, JOY would classify each amino residue as belonging to category “True” if it participated in such bonding and category “False” if it did not. We have given abbreviated names to these JOY environments, which with brief description in brackets are SecStrucPhi (secondary structure and phi angle), SolAcc (solvent accessibility), HmainCO (hydrogen bond from side chain to main chain CO group), HmainNH (hydrogen bond from side chain to main chain NH group), Hotherhet (hydrogen bond to nonstandard residue), Hhet (hydrogen bond to nonstandard residue), covhet (covalent bond to nonstandard residue), Hmamide (main chain to main chain hydrogen bond involving NH of specified residue), HmCO (main chain to main chain hydrogen bond involving CO of specified residue), DSSP (secondary structure using DSSP algorithm), and PosPhi (positive phi angle). The environment categories for SecStrucPhi and DSSP refer to participation in different types of secondary structure such as helix or coil. Further information on the environments is given in supplementary table S1, Supplementary Material online.

The COI protein was additionally partitioned into categories based on four environments, which potentially relate to selective constraint and purifying selection. The first of these is total amino acid diversity per site in the sample of eight sequences. The second and third relate to evolutionary conservation values derived using the programs ConSurf (Celniker et al. 2013) and FuncPatch (Huang and Golding 2015). As an exploratory approach, root mean square distance values for

atomic distances within superimposed COI structures including 1V54 were also calculated as the fourth constraint environment (see supplementary methods S1, Supplementary Material online), on the basis that higher constraint may be related to lower molecular distances in the superimposed structures. Contrasting categories based on these environments did not reveal the significant differences that are reported below for the JOY structural environments and are not considered further. A dendrogram showing clustering of the JOY structural environments, based on the category assignments for all sites over the 501 sites of 1V54 that align with the bivalve sequences, is given in figure 1C, with correlation values between environments in supplementary table S3, Supplementary Material online. Some clustering is expected, for example, of SecStrucPhi and DSSP, which reflect secondary structure. The application of JOY revealed that 7 of the 11 structural environments gave only small numbers of sites differing between F and M COI proteins in one of the contrasting categories, and preliminary analysis was not informative and so these structural environments were excluded from further analysis. The remaining four structural environments, SecStrucPhi, SolAcc, Hmamide, and HmCO, are the focus of further investigation. As illustration of the spatial location of category regions, these are shown marked on an image of COI for three of these environments, Hmamide and HmCO in figure 2A and SolAcc in figure 2B. Hmamide and HmCO are combined as they have some overlapping sites with significant but not high correlation (supplementary table S3, Supplementary Material online). These three environments feature most prominently in the later analysis. The Hmamide and HmCO True category residues are located in a hydrophobic region within the membrane. The False category residues are located toward the aqueous intermembrane space and mitochondrial matrix. These hydrophilic regions have a greater number of side-chain oxygen and nitrogen atoms available for hydrogen bonding with other molecules as compared with the True category (fig. 2C). The SolAcc True residues are on the external parts of the protein, the False residues are buried. Confirmation that JOY is effective in generating categories having differences potentially reflecting structure and function was obtained by comparing the amino acid distributions between categories for the pooled data set of eight bivalve sequences (supplementary table S4, Supplementary Material online). Chi-square contingency analysis reveals highly significant differences in amino acid distribution between categories for the four environments tested.

Calculation of F–M and False–True Differences

Much of the analysis is based on the comparison of mean property values between the contrasting categories of the different structural environments. To facilitate this approach, the difference in property values between F and M COI proteins (F–M) was calculated for each species for each

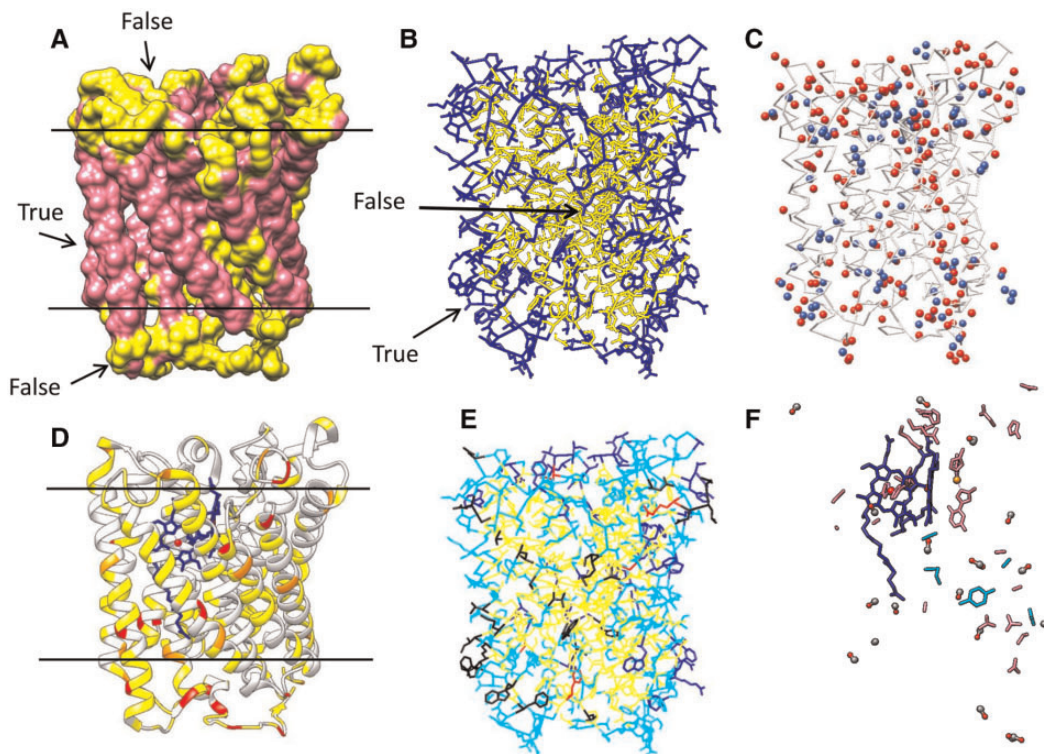


Fig. 2.—COI subunit (bovine 1V54 and bivalve molluscs) with annotations made in Chimera (Pettersen et al. 2004). The 1V54 COI protein is 514 amino acids long, but residues 1–7, 16, 139, 177–178, 287, 514 have been trimmed from the alignment. The orientation is the same in all images with the intermembrane space at the top and the matrix at the bottom. The transmembrane region is delimited by horizontal lines in *A* and *D*. The heme *a* and *a*3 centers with Fe and Cu are present in some images but partly obscured because of the orientation. (*A*) Environments Hmamide and HmCO regions. Yellow, category False for either environment; Pink, category True for both environments. (*B*) Environment SolAcc regions. Blue, category True, exterior residues; Yellow, category False, interior buried residues. (*C*) Location of side chain oxygen (in red) and nitrogen (in blue) atoms showing preponderance in hydrophilic regions at top and bottom. (*D*) Sites with radical differences between F and M marked on 1V54 in color. Numbers of differences are summed over species and over nine properties of Taylor (1986) (small, tiny, negative, positive, polar, charged, hydrophobic, aromatic, and aliphatic). Yellow, 1 or more differences; Orange, >4 differences; Red, >5 differences. (*E*) Environment SolAcc regions with subunit interface residues indicated. Cyan, category True, exterior residues; Yellow, category False, interior buried residues; Blue, SolAcc category True sites, which are also at mtDNA–mtDNA subunit interfaces; Black, SolAcc category True sites, which are also at mtDNA–nuclear subunit interfaces; Red, SolAcc category False sites, which are at any type of subunit interface. (*F*) Alanine to Serine (A in F genome, S in M genome) or Serine to Alanine (SA) substitutions made visible in relation to functional important sites including the D and K proton entry pathways, the heme *a* and heme *a*3 sites, putative proton and water exit pathways, and electron transfer pathway derived from the NCBI entry for 1V54 taking account of Tsukihara et al. (2003). The eight bivalve sequences have been modelled and superimposed on the 1V54 known structure as template using Modeller (Webb and Sali 2014; see supplementary methods S2, Supplementary Material online, for further details), and residue side-chains are revealed. Blue, heme with Fe (red) and Cu (brown); Dumbbells, AS or SA substitutions with oxygen of serine (red) and carbon of serine and alanine (gray); Pink or Cyan, functionally important residues; Cyan, functionally important residues within 5 Å of any AS or SA substitution.

alignment site. Use of $F-M$ corrects for variation between sites analogous to the approach in a paired t -test. A simple example of the calculation of $F-M$ is given in figure 1*D* for a small made-up data set for aligned F and M sequences. The figure also shows the computation of the quantity $False-True$, the difference between categories, which is used in further calculations as described below. A diagrammatic representation of $F-M$ for two categories for an environment is illustrated in figure 1*E*. In the example shown, $F-M$ is positive for category False ($F > M$) and negative for category True ($F < M$). The values of $F-M$ are also broken down into the deviations from the ancestors derived from FASTML. In different subsequent analyses, the averaging of $F-M$ over sites is done in two

ways, first over the 281 sites that differ between the F and M proteins, and second over all the 2,004 sites in the alignment for which the F and M proteins can be compared ($4 \text{ species} \times 501 \text{ sites} = 2,004$). The first method gives a mean value for only those sites that are variable, which may be more pertinent if many of the sites within a category region are invariant, the second for the entire physical region of a category.

Results

Sequence Differences and Sources of Variation

The alignment of the eight sequences used for analysis is 501 sites long. Around 291 sites show variation, that is having at

least two different amino acids in the alignment of eight sequences. The numbers of sites in which COI sequences differ between F and M genomes within species for 1, 2, 3, or all 4 species are respectively 135, 56, 10, and 1. Multiplying these vectors and summing gives 281 sites differing between genomes within species out of $291 \times 4 = 1,164$, a proportion of 0.24. Because each amino acid residue can be substituted by property values, a three-way ANOVA can be used to gauge the relative overall effect of the factors, genome, species, and environment. Mean partial eta squared values, the percentage of the total variation in the dependent variable attributable to a specified independent variable, over all properties and environments are given in table 1. The error is large because it reflects amino acid differences between sites. All factors have small effect as judged by the value of partial eta squared (%), and that for genome is smaller than that for species and environment. More optimistically, the genome effect is about 6% of the species value, and the interaction of genome and species is also higher, suggesting that the genome effect may in part be species dependent.

F–M for Entire Protein Alignment

The difference in property values between the F and M COI proteins was first examined over the entire protein alignment. Confidence intervals (CIs) of F–M were computed for the 281 site and 2,004 site data sets. For all 21 properties, the CIs overlap zero for both data sets (supplementary table S5, Supplementary Material online), and the means over properties are -0.181 (95% CI $-0.549, 0.188$) and -0.025 (95% CI $-0.078, 0.027$), respectively. The second data set has lower absolute mean values because it includes many sites where $F-M=0$ and has narrower CIs because the sample size is greater. These results suggest no difference between the F and M COI proteins in property values over the whole COI protein. However, this averaging may hide differences in the value of F–M between the categories of structural environments.

F–M Compared between Categories for Each Property

Further analysis thus focused on comparing the value of F–M between different regions of the protein, that is, between the True and False categories for the different structural environments (see fig. 1D for illustration of computation of F–M, and fig. 1E for further clarification). For each property for each environment, the value of F–M was compared between categories averaging over all the alignment sites within categories and using either ANOVA or the nonparametric Kruskal–Wallis test. The resulting *P*-values for each amino acid property were then combined over environments using Fisher’s combining probabilities test (table 2 with further details in supplementary table S6, Supplementary Material online). The properties that have significant *P*-values in table 2 are retained for further analysis. Some of these are

Table 1

Summary of Partial Eta Squared Values for Structural Environments

Source of Variation	Partial Eta Squared (%)
Species	0.18
Genome	0.01
Environment	1.37
Species × Genome	0.03
Species × Environment	0.21
Genome × Environment	0.02
Species × Genome × Environment	0.04

NOTE.—Values given are averages of all environments and amino acid properties derived from a three-way ANOVA on data set of 501 sites × 8 sequences.

Table 2

Amino Acid Properties Ranked According to Significance in Comparison of F–M between Structural Environment Categories

Amino Acid Property	Fisher <i>P</i> -Value
Kid9	0.000
Kid8	0.000
cc	0.000
HH2	0.000
Atch2	0.000
Kid7	0.000
Kid5	0.000
HH1	0.000
Kid4	0.000
pp	0.000
Kid2	0.004
Atch1	0.005
Kid3	0.024
HH3	0.370
Atch4	0.466
vv	0.494
Atch3	0.737
Kid6	0.875
Atch5	0.914
Kid1	0.914
arom	0.941

NOTE.—Summary of *P*-values for different amino acid properties are from Fisher’s combining probabilities test over structural environments. Color fill in amino acid property column marks clustered groups of individual properties in figure 1B. Pink fill, *P*-values ≤ 0.05 . Further information is given in supplementary table S6, Supplementary Material online.

combined into groups that reflect the closeness of their clustering (fig. 1B). The groups are named Group A (comprising Kid8, Kid9, Atch2, and cc), Group B (Kid3, Kid4, Kid5, and HH2), and Group C (pp, Atch1, and HH1). In subsequent analyses, mean values for Groups A, B, and C are obtained by averaging over properties within groups. Properties with nonsignificant *P*-values in table 2 are not considered further.

Concordance of Category Differences across Species

Concordant and parallel changes between the F and M proteins across species would provide evidence of a general

rather than species-specific evolutionary effect. To test for this, the mean values of F–M for the retained properties were calculated for the categories for each environment, separately for each of the four species. These F–M values were then compared between categories across the four species using a paired *t*-test. An example to illustrate the concordant differences between False and True across species is given in [table 3](#) for the environment HmCO and the Group A and Group B properties. For the Group A property values, the mean F–M is higher for category False than True. A similar pattern occurs for Group B with True higher than False, though one of the species (Sen) is different for Kid4 and Kid5. The corresponding values for all properties and environment categories are given in supplementary table S7, Supplementary Material online, and the *P*-values from all the paired *t*-tests in supplementary table S8, Supplementary Material online. The structural environment SecStrucPhi does not show as good concordance or as many low *P*-values as SolAcc, Hmamide, and HmCO and is excluded from further consideration. The values of False–True are plotted for SolAcc, Hmamide, and HmCO in [figure 3A](#) from which the concordance between the four species is clearly evident.

Mean F–M Values and Deviation from Ancestor

Summary mean values of F–M for the False and True categories are given in [table 4A](#) with a further statistical interpretation in [table 5](#). Of the 15 mean values for False in [table 4A](#), 9 are significantly different from zero (bold and underlined in [table 4A](#)). None of the True values are significant. For eight of these significant values, all four species are concordant with the same sign for the difference between F–M for False and for True (indicated by a “4” in the concordance rows). For Hmamide and HmCO that show a similar pattern of variation, for all 10 paired False and True mean values, the absolute value of False is >True (row b, [table 5](#)). SolAcc does not show this effect so clearly (row a). In 12 of the 15 paired False and True mean values, the sign of the mean is different (row c). This suggests that a change in evolution that makes $F > M$ in one category is accompanied by an opposite change making $M > F$ in the other category consistent with the absence of a difference between F and M over the protein as a whole (see above and supplementary table S5, Supplementary Material online).

The value of F–M has also been broken down into deviation from the ancestor (shown as $D = F - A$ and $D = M - A$ in [fig. 1E](#)). In this illustration, the ancestor is positioned between F and M though it is also possible in evolution for the ancestor to have the most extreme value. The values of D are given in [table 4A](#) with statistical interpretation in [table 5](#). The absolute value of M–A is more often greater than F–A particularly for category False (rows d and f, [table 5](#)). The signs of F–A and M–A are more often different within categories

Table 3

Species Concordance for F–M for HmCO Groups A and B

		Amino Acid Property			
Group A					
Species	Category	cc	Kid8	Kid9	Atch2
Cal	False	0.231	0.873	0.591	0.461
	True	–0.096	0.132	0.140	0.040
Sen	False	0.160	0.129	0.201	0.289
	True	–0.099	–0.192	–0.164	–0.349
Phi	False	0.301	0.336	0.333	0.295
	True	–0.050	–0.022	0.047	0.168
Gra	False	0.118	0.338	0.377	0.325
	True	–0.230	–0.313	–0.161	–0.261
Paired <i>t</i> -test <i>P</i> -value		0.001	0.016	0.005	0.031
Group B					
Species	Category	Kid3	Kid4	Kid5	HH2
Cal	False	–0.121	–0.503	–0.341	–1.973
	True	0.004	–0.111	–0.094	–0.200
Sen	False	–0.187	–0.001	–0.156	–0.657
	True	–0.040	–0.081	–0.190	–0.188
Phi	False	–0.109	0.025	–0.296	0.200
	True	0.106	0.045	0.160	0.472
Gra	False	–0.155	–0.297	–0.222	–0.731
	True	–0.028	0.013	0.042	0.309
Paired <i>t</i> -test <i>P</i> -value		0.005	0.250	0.104	0.078

NOTE.—Shading indicates the higher of each paired value of False and True for a property.

(rows e and g). For Hmamide and HmCO the absolute value of F–A (or M–A) is more often greater for False than True (rows i and l, [table 5](#)). SolAcc does not show this effect so clearly (rows h and k). There is more often a difference in sign for F–A between False and True (row j) and also for M–A (row m) consistent with the absence of a difference between F and M over the protein as a whole (supplementary table S5, Supplementary Material online). Finally, the [table 5](#) data can be broken down into 15 2×2 tables for each property and environment combination. The cell with the highest absolute value is recorded and summed over tables (row n). The False and (M–A) cell has highest value in 11 of the 15 tables.

The (F–A) and (M–A) values in [table 4A](#) are the means over four species. A multiway ANOVA with sources of variation Property (5 classes), Environment (3), Category (2), Genome (2), and Species (4) on the absolute values was carried out. Genome and Category are both highly significant ($P = 0.000$) as is Genome \times Category ($P = 0.007$) consistent with the analysis from [table 5](#). No significant result is obtained in the corresponding ANOVA of the real values, again consistent with observations of different signs ([table 5](#)). Values of deviation from the ancestor for F and M proteins and for False and True categories are illustrated in [figure 3B](#) and confirm that the observed trends are consistent across species.

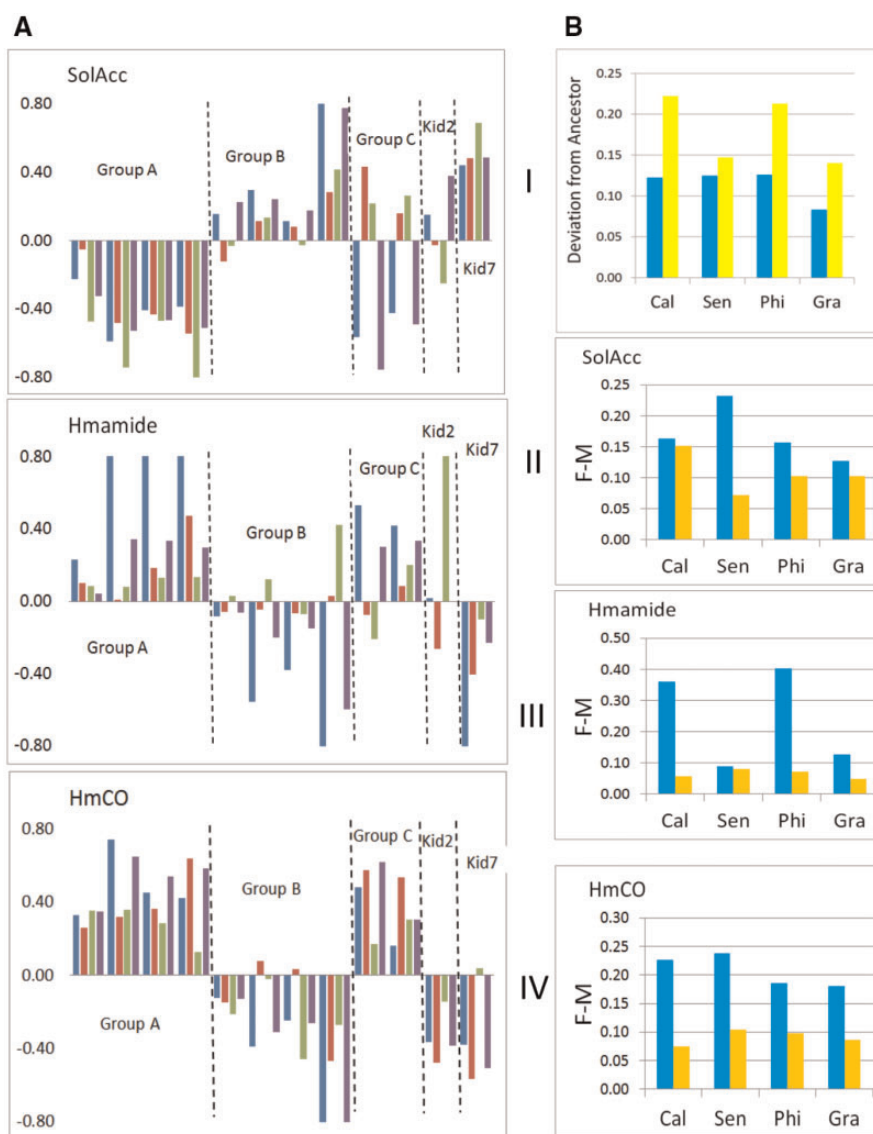


FIG. 3.—Amino acid property differences between categories and species. (A) Histograms for False–True for three structural environments for the data set of 281 variable sites. The amino acid property values are from the left in order cc, Kid8, Kid9, Atch2, Kid3, Kid4, Kid5, HH2, pp, HH1, Kid2, Kid7, with species in the order Cal (blue), Sen (red), Phi (green), and Gra (purple). Plotted values are derived from those in supplementary table S7, Supplementary Material online. For clarity, some high and low values are truncated on the graphs. (B) Histograms derived from data underlying table 4A but shown separately for each species. I: Absolute deviation from ancestor for F genome (F–A, left blue fill) and M genome (M–A, right yellow fill) averaged over properties and environments. II–IV: Absolute value of deviation from ancestor for False (left blue fill) and True (right orange fill) for environments SolAcc, Hmamide, and HmCO averaged over properties.

Physicochemical Interpretation of Mean Category Values

For a physicochemical interpretation of the values of F–M, it is necessary to establish the meaning of high and low values on the amino acid property scales. This has been done using the interpretation of the scales made by the authors who devised them (supplementary table S1, Supplementary Material online), supplemented by two other classification systems (from Taylor 1986 and Malkov et al. 2008) and also taking into consideration the amino acid chemical groups. The amino acids are sorted according to the property values from high

to low values and lined up with these classification systems (supplementary table S9, Supplementary Material online). This allows a summary headline classification to be assigned to high or low property values. Whichever of the F or M COI proteins has the higher property value as judged by the F–M mean value in table 4A is assigned as best match to the headline classification as shown in table 4B. The entries in table 4B give a consistent picture showing greater helicity and hydrophobicity (or nonpolarity) for the M protein in external residues or those in the matrix or intermembrane space

Table 4

Summary of Statistics and Analysis of F–M Values for Amino Acid Properties for Three Structural Environments

Environment	Group A		Group B		Group C		Kid2		Kid7			
	False	True	False	True	False	True	False	True	False	True		
(A)	Mean amino acid properties (F–M) for three structural environments											
SolAcc	Mean	<u>-0.288</u>	0.180	0.055	-0.181	0.015	0.161	0.010	-0.052	<u>0.342</u>	-0.181	
	Ancestor	D=F–A	-0.065	0.030	0.039	-0.061	-0.078	0.113	0.042	-0.005	0.091	-0.015
		D=M–A	0.223	-0.151	-0.017	0.120	-0.093	-0.049	0.032	0.047	-0.251	0.167
	Concordance	False–True	4		4		2		2		4	
Hmamide	Mean	<u>0.352</u>	-0.038	-0.268	-0.054	0.250	0.050	0.099	-0.044	<u>-0.439</u>	0.083	
	Ancestor	D=F–A	0.070	-0.022	-0.125	-0.016	0.227	0.032	0.011	0.006	-0.095	0.060
		D=M–A	-0.282	0.016	0.143	0.039	-0.023	-0.018	-0.088	0.050	0.345	-0.023
	Concordance	False–True	4		3		2		2		4	
HmCO	Mean	<u>0.335</u>	-0.088	<u>-0.345</u>	0.014	<u>0.361</u>	-0.033	<u>-0.254</u>	0.088	<u>-0.285</u>	0.069	
	Ancestor	D=F–A	0.109	-0.033	-0.108	-0.005	0.134	0.025	-0.041	0.030	-0.091	0.037
		D=M–A	-0.226	0.055	0.238	-0.019	-0.228	0.058	0.213	-0.059	0.194	-0.032
	Concordance	False–True	4		4		4		4		3	
(B)	Protein best matching the headline classification											
Headline classification	Helical		Hydrophobic		Nonpolar		Hydrophobic		Helical			
SolAcc	F	M	= (F)	M	= (M)	M	= (F)	= (M)	F	M		
Hmamide	M	= (F)	M	= (M)	M	= (M)	= (F)	= (M)	M	= (F)		
HmCO	M	= (F)	M	= (F)	M	= (F)	M	= (F)	M	= (F)		
(C)	Ancestral values for properties											
SolAcc	<u>-0.372</u>	-0.189	<u>0.936</u>	0.814	<u>0.156</u>	0.323	<u>0.246</u>	0.158	<u>0.217</u>	0.009		
Hmamide	0.210	<u>-0.355</u>	0.144	<u>1.048</u>	1.230	<u>0.018</u>	0.029	<u>0.233</u>	-0.166	<u>0.116</u>		
HmCO	-0.129	<u>-0.314</u>	0.682	<u>0.966</u>	0.436	<u>0.161</u>	0.041	<u>0.267</u>	-0.027	<u>0.127</u>		

NOTE.—(A) Mean; values of F–M for False and True categories for three environments (SolAcc, Hmamide, and HmCO) for amino acid properties (groups A, B, C and Kid2 and Kid7) averaged over sites and four species for the data set of 281 variable sites. The values for groups A, B, and C are arithmetic mean values of the individual properties forming each group. False and True: structural environment categories. Bold underline: 95% CIs for F–M do not overlap the test value = 0. Broadly similar patterns of significance and concordance are obtained for the 2,004 site data set. Ancestor: D = F–A and D = M–A, deviations from ancestor. Concordance: the number of species having the same sign for the difference False–True. Sample sizes of sites for categories for 281 and 2,004 site data sets are SolAcc (False = 84, 912; True = 197, 1,092), Hmamide (False = 59, 364; True = 222, 1,640), HmCO (False = 98, 536; True = 183, 1,468). (B) Protein best matching the headline classification, depending on relative values on property scales. Decisions for False in bold, = (F), = (M): decision given but difference between F and M proteins small. In comparison with supplementary table S9, Supplementary Material online, where high values of Group A indicate non-helicity, the Group A scale is reversed so that the headline classifications of Group A and Kid7 match. Similarly, the Group C scale is reversed to match Group B and Kid2. (C) Ancestral mean values for properties averaged across the four species for sites differing between F and M proteins. Bold underline: False or True have greater helicity (for Group A and Kid7) or greater hydrophobicity (for Group B, Group C, and Kid2). Note that as in supplementary table S9, Supplementary Material online, a more negative value of Group A indicates greater helicity and a more negative value for Group C indicates greater polarity.

(Hmamide and HmCO False, and SolAcc True). By contrast the F protein shows greater helicity for buried residues (SolAcc False). To get an overall picture of hydrophobicity and helicity in the categories prior to divergence of F and M proteins, the mean property values over all residues were calculated for the ancestral sequences at the sites, which differ between F and M (table 4C). The buried residues (SolAcc False, and Hmamide and HmCO True) have greater helicity and hydrophobicity as expected, compared with the alternative categories representing external residues or those in the matrix or intermembrane space.

A summary of the more marked changes in the F and M proteins from table 4 are mapped onto a diagram of COI (fig. 4). Of general significance is that the changes in the M protein, which tend to be larger (see also fig. 3B), are in a direction to minimize the difference between the categories. Thus, for Hmamide and HmCO, the M protein has a higher hydrophobicity change in the False category (fig. 4B), which is

the more hydrophilic part of the protein. The F protein shows a smaller change in the opposite direction. All the five directional changes in the M protein and the four directional changes in the F protein in figure 4 show this same tendency, which is that M changes toward the mean of the two categories, whereas F changes to increase the difference between the two categories.

Influence of Individual Amino Acids on Variation in F–M

To gauge the influence of individual amino acids on variation in F–M, the mean category values have been partitioned into separate contributions from each amino acid. The F–M values associated with individual amino acids in the F and M sequences are summed for each amino acid over all sites in the F sequence and separately over all sites in the M sequence. The absolute values of the F–M are also summed over all amino acids. The individual amino acid contributions are then

Table 5

Analysis of F–M and Deviations from Ancestor

(A)	Mean F–M Compared between False and True Categories			Yes	No
a	Absolute value	SolAcc	False > True	2	3
b		Hmamide/HmCO	False > True	10	0
c	Sign		False ≠ True	12	3
(B)	M–A Compared With F–A Within the False and Within the True Categories				
d	Absolute value		False (M–A) > (F–A)	12	3
e	Sign		False (M–A) ≠ (F–A)	13	2
f	Absolute value		True (M–A) > (F–A)	10	5
g	Sign		True (M–A) ≠ (F–A)	12	3
(C)	The Value of F–A (or M–A) Compared Between False and True				
h	Absolute value	SolAcc	F–A False > True	3	2
i		Hmamide/HmCO	F–A False > True	10	0
j	Sign		F–A False ≠ True	10	5
k	Absolute value	SolAcc	M–A False > True	3	2
l		Hmamide/HmCO	M–A False > True	10	0
m	Sign		M–A False ≠ True	11	4
(D)	Highest Absolute Value Within 2 × 2 Tables (False, True) Versus ((F–A), (M–A))				
n	False and (F–A)	False and (M–A)	True and (F–A)	True and (M–A)	
	1	11	1	2	

NOTE.—Analysis of table 4A values. Yes and No: number of values from table 4 in accordance with the specified condition (e.g., False > True for category). (F–A) and (M–A): deviations of the F and M values from their ancestor within species, see figure 1D. Row labels a–n: see text. (A) Analysis of the 15 paired False and True values in the three rows labelled “Mean” in table 4A. (B) Analysis of the 15 paired (M–A) and (F–A) values for False and True separately. Analysis based on the absolute values and, in a separate test, the sign. (C) Analysis of the 15 paired False and True values separately for (F–A) and (M–A). Analysis based on the absolute values and, in a separate test, the sign. (D) Comparison of the four values in the 15 2 × 2 tables (False, True) versus ((F–A), (M–A)), for example, the four values –0.065, 0.030, 0.223, –0.151 for SolAcc Group A is one such table.

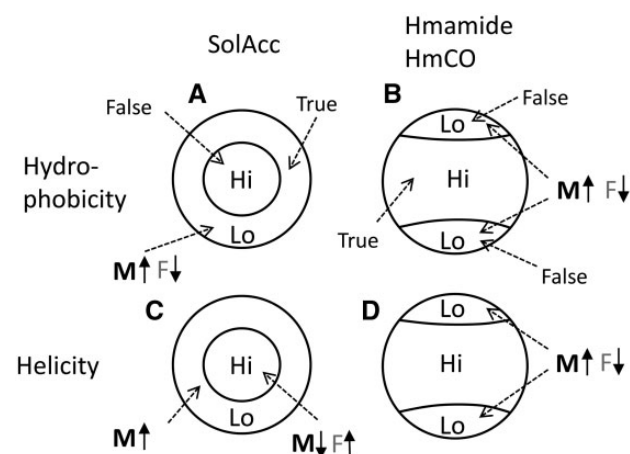


Fig. 4.—Diagrammatic representation of COI with quantitatively largest amino acid property value changes derived from table 4. (A) and (C): inner circle delimits category False for SolAcc. (B) and (D): top and bottom regions cut off by arcs delimit category False for Hmamide and HmCO. Hi and Lo: high or low hydrophobicity and helicity in the ancestors of F and M. The representation is based on combining property values for Group A and Kid7 as helicity and Group B, Group C, and Kid2 as hydrophobicity. Up and down solid arrows: value has increased or decreased in that category region compared with ancestor, with bold and gray letters indicating the larger and smaller absolute changes, respectively.

computed by dividing the sum for each amino acid by the overall sum of absolute values. A demonstration of the method is given in supplementary table S10A, Supplementary Material online, using the made-up data set

of figure 1D. In the interpretation, focus is on those values which are large and have the same sign as the overall F–M value thus contributing to the excess of F over M or vice-versa. A summary of the results of this analysis is given in figure 5 where amino acids with high contributions are indicated, with the underlying data analysis in supplementary table S10B, Supplementary Material online. A high contribution in absolute terms will be favored when the amino acid has a high value on the amino acid property scale and when it has a high frequency, and will also be a function of the specific pairings of amino acids at sites in the alignment.

Amino acids showing the most marked contributions are the aliphatic amino acids (G, A, V, L, I) and the aliphatic hydroxyl amino acid serine (S). Contingency tests for association were carried out to compare the amino acid frequency distributions of F and M given in figure 5 (further details in supplementary table S10B, Supplementary Material online). Some significant results were obtained for individual amino acids: SolAcc False (alanine A, $P=0.031$), Hmamide True (A, $P=0.028$), HmCO True (A, $P=0.023$, serine S, $P=0.042$), and for the entire contingency table for True ($P=0.029$). The analysis of figure 5 was extended by calculating the contributions for precise F with M amino acid pairings in the alignment, of which there are 96 (supplementary table S10C, Supplementary Material online). Contributions of more than 10% are flagged in supplementary table S10C, Supplementary Material online, and these are AS (i.e., A in F protein and S in M protein), for SolAcc False Groups A, B, and Kid7 and

(2008) used alteration of physical and chemical properties or change in a contact site with another subunit as evidence of a likely change in cytochrome oxidase activity. Six of nine observed amino acid mutations were concluded to have potentially changed activity. By this approach, some of the many radical amino acid differences observed here between the F and M COI proteins may be of functional significance. A classification of the residues in 1V54 was made using information derived from its NCBI entry taking account of the study of Tsukihara et al. (2003). This identifies residues functionally important in the catalytic mechanism, as well as mtDNA-encoded subunit interface residues (i.e., COI with COII and COIII) and COI with nuclear DNA-encoded subunit interfaces.

Subunit interface residues are mapped onto 1V54 in figure 2E. Seventy-one out of 76 interface sites occur near the surface of the protein, and of these 30 showed radical amino acid differences. The three mitochondrial coded subunits play an important catalytic function in the complex. They are regulated by their own redox state (Allen 2015) as well as being regulated by the nuclear encoded subunits (Ludwig et al. 2001). Studies of K_a/K_s ratios may throw light on whether selection and constraint differ between different types of interfaces perhaps to optimize interactions between residues coded by the mtDNA and nuclear genomes, but a consistent picture has not yet emerged (Schmidt et al. 2001; Aledo et al. 2014). In this study, no consistent differences are found comparing F and M COI proteins in the value of K_a/K_s between mtDNA–mtDNA and mtDNA–nuclear DNA interfaces. Following the approach used for comparing between structural environment categories (table 4A), no significant differences were observed in F–M for property values for hydrophobicity or helicity between the two types of interface sites, or between interface and noninterface sites.

The functionally important residues are mapped onto 1V54 in figure 2F. For those F and M COI sites aligning with the residues functionally important in the catalytic mechanism of 1V54 (31 identified here), the K_a/K_s ratio between F and M is 0.004 indicating high conservation as expected of such sites. Twenty-three of these 31 are buried within the COI protein (SolAcc False), and of these only one showed an amino acid difference between F and M and then in only one species. The bivalve sequences as targets were also modelled on the 1V54 structure as template using Modeller (Sali and Blundell 1993; Webb and Sali 2014) and the models and 1V54 structure superimposed in Chimera (Pettersen et al. 2004; see supplementary methods S2, Supplementary Material online, for further information). This allows the calculation of the distances of specific residues in the bivalve models from functional residues in 1V54. Following the approach used for table 4A, no significant differences were observed in F–M for property values for hydrophobicity and helicity between bivalve model sites closer to and further away than either 2 or 4 Å from the 1V54 functionally important sites shown in figure 2F. Further focus was thus on the substitutions between alanine (A) and

serine (S), which are found to play a relatively large quantitative role in the amino acid property differences between the F and M proteins (fig. 5). The 13 AS and 7 SA substitutions together with the functional important 1V54 residues alone are made visible in figure 2F. The substitutions are scattered throughout the protein but many are close to functional sites (e.g., within 5 Å). There is a preponderance of such substitutions within the more hydrophobic regions of the protein (SolAcc False, Hmamide/HmCO True). In silico analysis using the program NetPhos (Blom et al. 1999) suggested that 13 of the 20 substitutions had serine residues, which had potential to be phosphorylated with associated potential kinase binding sites.

Discussion

The amino acid properties considered here differ significantly in value between the categories of structural environments representing different regions of the protein (table 2). For environments Hmamide and HmCO, the categories contrast the more hydrophilic regions associated with the matrix and intermembrane space (False) with the region within the membrane (True) (figs. 2A and 4). For SolAcc, the contrast is between external regions (True) and buried regions (False) (figs. 2B and 4). The calculated difference between F and M proteins (F–M) is frequently significantly different from zero in individual categories (table 4A), and the value of F–M itself differs between categories consistently across species (fig. 3A). When the property changes are broken down into deviations of F and M from the ancestor (M–A and F–A, fig. 1E), greater change in the absolute value is observed for the M protein and the category False (tables 4A and 5, fig. 3B). In addition, although some amino acids make a relatively large contribution to the variation (e.g., alanine and serine) many other amino acids have small individual effects but which sum to a large value overall (fig. 5 and supplementary table S10, Supplementary Material online).

Category differences are to a great extent concordant in the four species studied (fig. 3A, table 3, supplementary table S7, Supplementary Material online). The absolute values of False and True also differ concordantly (fig. 3BII–IV). The absolute value of deviation from the ancestor is also greater for the M than F COI protein in all four species (fig. 3BI). The concordant changes give confidence that there may be shared causes occurring in parallel in the four species. One possible shared parallel cause of faster evolution of the COI M protein could be a higher mutation rate in the M genome resulting from oxidative damage in sperm or a higher number of mitotic divisions in the male germ line (Zouros 2013). This is consistent with some evolutionary studies of F and M genome sequences in *Mytilus* (Quesada et al. 1998), though Stewart et al. (1996) found no evidence of faster synonymous substitution in the M genome.

The changes in hydrophobicity and helicity, particularly in the M protein (tables 4 and 5, fig. 4) can be considered from a functional viewpoint. In anthropoid primates, Schmidt et al. (2005) found relatively many charged to uncharged amino acid substitutions in COI at the cytochrome c binding site on the intermembrane side of the protein. It was proposed that these would increase hydrophobicity and affect electron transfer from cytochrome c to COI. It was further proposed that this might be an adaptive change to reduce and hinder OXPHOS activity and consequent free radical damage in the brain, which uses a relatively large amount of oxygen in the anthropoids. In a comparison of COI between marine and freshwater copepods, McClellan (2013) found 11 radical amino acid changes consistent with a decrease in hydrophobicity around the proton input channel and possibly a less energetically expensive route for proton transfer during adaptation to freshwater. By analogy to these studies, the changes in hydrophobicity and helicity in the M protein (fig. 4) might have functional consequences as a result of effects on the movement of protons, electrons, or water molecules or alteration of the mitochondrial membrane potential. This could in turn affect ATP production by ATP synthase. Sperm uses ATP for movement, the acrosome reaction, and a variety of metabolic functions (Visconti 2012). Thus, reduction in ATP output could have negative consequences for sperm carrying the M genome. In humans, there is evidence that mtDNA mutations reducing ATP production may cause reduced sperm motility (Ruiz-Pesini et al. 2000). Conversely higher ATP production might lead to increased swimming speed or endurance, which in turn may lead to higher chance of successful fertilization. Changes in helicity might also have impact on protein stability because of requirements associated with hydrogen bonding (Pace et al. 2014). Proteins can only tolerate small changes in conformation operating under conditions of marginal stability, which may be of advantage in evolution (Taverna and Goldstein 2002). However, because of the effects of stability loss in a structure that is already marginally stable, a large proportion of missense mutations in proteins are likely to affect function (Tokuriki and Tawfik 2009). Thus, the decrease in helicity in buried residues (fig. 4C, SolAcc False) may have consequences for M COI protein stability and thence function.

Relaxation of selective constraint and purifying selection with increased genetic drift and fixation of nearly neutral mutations (Ohta 1992) as the cause of the changes in amino acid property values in the M protein would be consistent with the observed directional changes in property values. The False and True categories (see fig. 2A and B) have ancestral physicochemical differences reflected in their overall property values (table 4C). The changes in the M protein tend to be such as to decrease these differences with a difference in sign of M–A for both the False and True categories (tables 4A and 5). For example, the regions of the M protein that are less hydrophobic in the ancestor evolve to become more

hydrophobic, the less helical regions tend to become more helical, and the more helical regions tend to become less helical (fig. 4). Given that the maintenance of these differences is important for normal functioning, their erosion by fixation of nearly neutral mutations would be expected to affect function adversely, altering mitochondrial membrane potential and ATP production with consequent biological effects, for example, a reduction in sperm performance. Other studies also indicate faster evolution of the M genome as a result of lower functional constraint and less purifying selection (Stewart et al. 1996; Zouros 2013). A contributory factor may be a lower effective population size for the M than F genome, for example, due to a narrower bottleneck of mtDNA copy number in sperm (Stewart et al. 1996) compared with eggs. There is no evidence however that the M genome is nonfunctional. In sperm of *R. philippinarum*, one of the species used here, there is experimental evidence for the existence of membrane potential consistent with OXPHOS activity (Milani and Ghiselli 2015). The M genome is abundant in somatic tissues of *R. philippinarum* (Ghiselli et al. 2011) and transcriptionally active in male tissues (Milani et al. 2014). The M genome is also expressed in male but less so in female tissues of *Mytilus* species (Dalziel and Stewart 2002; Obata et al. 2011). It has been suggested that selective constraint may be lower in the M genome because it functions in fewer tissues (Stewart et al. 1996; Zouros 2013), but the relative influences of selective pressures in different tissues are unknown.

Given the importance of mitochondrial function in gonads and sperm of these broadcast spawning species (Ghiselli et al. 2013; Zouros 2013), parallel positive selection, for example, toward improved sperm performance through enhanced ATP production, can also be considered as a cause of the changes in amino acid property values. The associated high membrane potential could also be used as a signal for preferential partitioning of sperm-derived mitochondria into the primordial germ cells in males of DUI species (Milani 2015). However, in this study the K_a/K_s ratio is less than 1 for comparisons between F and M proteins and a site-specific analysis with HyPhy provides no evidence of positive selection. This is consistent with evidence of high conservation of the COI protein in general across eukaryotes (Pierron et al. 2012) and in bivalves specifically (Plazzi et al. 2016), though meta-analyses have revealed examples of positive selection in almost all mtDNA-encoded proteins (Garvin et al. 2015; James et al. 2016). Consideration could also be given to the possibility that the positive selection acts at linked regions in the mtDNA molecule (genetic draft or hitchhiking), which may be highly prevalent in mtDNA (Bazin et al. 2006). Such an explanation would require that positive selection in another mtDNA-encoded protein would cause the same parallel amino acid property changes between F and M COI proteins in all four of the studied species. In the F COI protein, the changes in property values away from the ancestral values tend to increase the difference between False and True

categories. For example, for SolAcc the category False comprises buried residues in a highly helical region and the F protein shows an increase of helicity in this region (fig. 4C). The changes are small in magnitude, but are difficult to explain by relaxation of selection and drift, and are more in line with an adaptive change enhancing the function of the F protein. The difference in sign of F–M in the False and True categories (tables 4A and 5) is consistent with the absence of differences between the F and M protein overall (supplementary table S5, Supplementary Material online). This could reflect compensatory evolution (Ivankov et al. 2014) with stabilizing selection maintaining the overall F–M property value within a tolerable range and deleterious mutations in one category being balanced by advantageous mutations in the other category.

Individual amino acid substitutions may also have a context-dependent functional influence on the catalytic mechanism. Because the radical differences between F and M are scattered throughout the protein (fig. 2D), many fulfil a criterion for potential influence, being close to functionally important sites (Schmidt et al. 2005; Melvin et al. 2008). Many of these substitutions involve alanine and serine, which play a relatively large quantitative role in the amino acid property differences between the F and M COI proteins (fig. 5). Serine residues occur quite frequently in functional centers because the OH side chain group can form hydrogen bonds with other residues or substrates. In a study of different isolates of the nematode *C. elegans*, Dingley et al. (2014) discovered a single A to S substitution in COI, which had a significant effect on cytochrome oxidase activity with other phenotypic consequences. The serine was close to the binding site for a kinase known to affect mitochondrial membrane potential and ATP synthase. In this study, the *in silico* analysis also suggested potential for serine phosphorylation. Phosphorylation of enzymes involved in OXPHOS may play an important role in its regulation (Acin-Perez et al. 2009). Tyrosine phosphorylation by kinases is involved in sperm capacitation and the acrosome reaction (Naz and Rajesh 2004). Thus, this frequently occurring substitution may be a good potential candidate in the search for functional effects.

Chapman et al. (2008) evaluated physicochemical properties in a specific extension of the M genome COII protein in unionoidean bivalves, and although purifying selection was dominant, the properties helical contact area and partial specific amino acid volume showed evidence of positive selection. Similarly, relaxation of selection and drift and positive selection might together be combined in an integrated explanation of the change in the M COI protein. Residues in the external regions of proteins are generally less conserved than buried sites in the core or important functional residues (Toth-Petroczy and Tawfik 2011), and the distribution of selection coefficients for mtDNA shows a peak at near neutrality but with many values greater or less than zero (Tamuri et al. 2012). Thus, given a lower effective population size in males (Stewart et al. 1996), the many substitutions in the external

regions of the protein associated with the category differences (table 4A and fig. 4) might have smaller selection coefficients, some decreasing and some increasing amino acid properties, and be nearly neutral. Their fixation by drift would degrade the variation in hydrophobicity between the different regions of the M protein and adversely affect function. A consequent strong selective pressure to maintain membrane potential and ATP production could result in compensating adaptive substitutions with larger selection coefficients in the slower evolving protein core, such as those involving alanine and serine, to enhance the catalytic mechanism. This could involve an additive effect on fitness of the spatially separated sites or involve epistatic interactions, which are generally assumed to be widespread in protein evolution (Starr and Thornton 2016). Examples of studies providing approaches to mechanistic understanding of such interactions are the analysis of compensatory interactions between mutations around a heme pocket in the protein CY51 affecting resistance to azoles (Mullins et al. 2011) or the analysis of the restoration of channel activity involving allosteric interaction between extra-membrane and transmembrane domains in viruses (To et al. 2017).

Apart from phylogenetic and molecular evolution approaches, a variety of experimental techniques are available for investigating the functional and biological effects of genetic variation in cytochrome oxidase, and some of these may be applicable in the analysis of F and M genome proteins in future. Mutations in COI are known to cause genetic conditions in humans (Zhen et al. 2015) or affect cytochrome oxidase activity in mouse cells (Acin-Perez et al. 2003). Performance of mtDNA genomes can be compared on the same or different nuclear backgrounds in cell cybrids (Kenyon and Moraes 1997) or in whole organisms obtained by backcrossing (Dingley et al. 2014). Such techniques are difficult to present in bivalves that are more difficult to breed in the laboratory. However measurements of swimming speed on sperm carrying different mtDNA genomes have been made in *M. edulis* (Everett et al. 2004; Jha et al. 2008) and this could be pursued in the species studied here, with the hypotheses of relaxed versus positive selection having different predictions on performance of F and M carrying sperm. Many spectrophotometric and cyto-histochemical methods are available for measuring cytochrome oxidase activity in isolated mitochondria, cells or tissue sections (Lanza and Nair 2009) and which can be extended with the use of cytochrome oxidase inhibitors (Pacelli et al. 2011). Techniques are also available for measuring ATP production directly or making accurate measurements of respiration in mitochondria or cells (Lanza and Nair 2009; TeSlaa and Teitell 2014). Such approaches seem more feasible for bivalves. Experimental techniques could be supplemented with molecular modelling to identify novel residue interactions or structural changes (Mullins et al. 2011; To et al. 2017) or with molecular dynamics simulations which have been used to study cytochrome oxidase function

(Arnarez et al. 2013), and identify water exit pathways in bovine cytochrome oxidase (Sugitani and Stuchebrukhov 2009).

In future, mtDNA will be sequenced in many other species with DUI, and it could be possible to extend the use of metrics in the place of amino acid codes as in this study to test more refined hypotheses on regional variation in the protein combined with the diverse experimental and simulation techniques outlined above. If the F and M COI proteins are structurally and functionally different, it will be important to investigate how the different proteins maintain the same primary function (OXPHOS) while adapting to different cellular environments (e.g., spermatozoon and egg). In this study, the changes in amino acid property values in the F and M COI proteins are frequently different in sign and direction within categories (rows e and g, table 5 and fig. 4). This could suggest an interaction due to physical proximity of the two proteins and their genomes, so that evolution in one direction in the F protein favors compensatory adaptive evolution in the opposite direction in the M protein to conserve some as yet unidentified property or function. This proximity could occur most readily in males with DUI, perhaps in the fertilized egg, or in heteroplasmic mitochondria. The natural coexistence of two diverged mtDNAs in the same organism in DUI also opens new areas of investigation into mito-nuclear interactions and co-evolution, for example, of proteins involved in fertilization and sex determination and identified through genomics and proteomics (Ghiselli et al. 2012; Diz et al 2013). This could be through interactions with different alleles of the same nuclear genes in different tissues, in accordance with growing evidence of tissue- and environment-specific modulation of intraindividual mito-nuclear interactions in animals (Wolff et al. 2014).

Supplementary Material

Supplementary data are available at *Genome Biology and Evolution* online.

Acknowledgments

We thank Andrej Sali and Ben Webb for advice and discussion, Andrej Sali for suggesting the use of JOY, and Kenji Mizuguchi for advice on the implementation of JOY. We thank Dennis Lavrov, the Associate Editor of GBE and three anonymous reviewers for their suggestions and helpful comments. The work by F.G. was supported by the Italian Ministry of Education, University and Research, MIUR—FIR Programme no. RBFR13T97A; L.M. was supported by MIUR—SIR Programme no. RBS14G0P5; A.P.D. was supported by the Spanish “Ministerio de Economía y Competitividad” (code AGL2014-52062-R), Fondos Feder and Xunta de Galicia (“Grupos de Referencia Competitiva” ED431C 2016-037).

Literature Cited

- Acin-Perez R, et al. 2003. An intragenic suppressor in the cytochrome c oxidase I gene of mouse mitochondrial DNA. *Hum Mol Genet.* 12(3):329–339.
- Acin-Perez R, et al. 2009. Cyclic AMP produced inside mitochondria regulates oxidative phosphorylation. *Cell Metab.* 9(3):265–276.
- Adachi J, Hasegawa M. 1996. Model of amino acid substitution in proteins encoded by mitochondrial DNA. *J Mol Evol.* 42(4):459–468.
- Aledo JC, Valverde H, Ruíz-Camacho M, Morilla I, López FD. 2014. Protein-protein interfaces from cytochrome c oxidase I evolve faster than nonbinding surfaces, yet negative selection is the driving force. *Genome Biol Evol.* 6(11):3064–3076.
- Allen JF. 2015. Why chloroplasts and mitochondria retain their own genomes and genetic systems: colocation for redox regulation of gene expression. *Proc Natl Acad Sci U S A.* 112(33):10231–10238.
- Arnarez C, Marrink SJ, Periole X. 2013. Identification of cardiolipin binding sites on cytochrome c oxidase at the entrance of proton channels. *Sci Rep.* 3(1):1343.
- Ashkenazy H, et al. 2012. FastML: a web server for probabilistic reconstruction of ancestral sequences. *Nucleic Acids Res.* 40(W1):W580–W584.
- Atchley WR, Zhao JP, Fernandes AD, Druke T. 2005. Solving the protein sequence metric problem. *Proc Natl Acad Sci U S A.* 102(18):6395–6400.
- Balaji S, Srinivasan N. 2007. Comparison of sequence-based and structure-based phylogenetic trees of homologous proteins: inferences on protein evolution. *J Biosci.* 32(1):83–96.
- Bazin E, Glemin S, Galtier N. 2006. Population size does not influence mitochondrial genetic diversity in animals. *Science* 312(5773):570–572.
- Blom N, Gammeltoft S, Brunak S. 1999. Sequence and structure-based prediction of eukaryotic protein phosphorylation sites. *J Mol Biol.* 294(5):1351–1362.
- Breton S, Ghiselli F, et al. 2011. Evidence for a fourteenth mtDNA-encoded protein in the female-transmitted mtDNA of marine mussels (*Bivalvia: Mytilidae*). *PLoS One* 6:e19365.
- Breton S, Stewart DT, et al. 2011. Novel protein genes in animal mtDNA: a new sex determination system in freshwater mussels (*Bivalvia: Unionoidea*)? *Mol Biol Evol.* 28:1645–1659.
- Breton S, et al. 2014. A resourceful genome: updating the functional repertoire and evolutionary role of animal mitochondrial DNAs. *Trends Genet.* 30(12):555–564.
- Celniker G, et al. 2013. ConSurf: using evolutionary data to raise testable hypotheses about protein function. *Isr J Chem.* 53(3–4):199–206.
- Chapman EG, et al. 2008. Extreme primary and secondary protein structure variability in the chimeric male-transmitted cytochrome c oxidase subunit II protein in freshwater mussels: evidence for an elevated amino acid substitution rate in the face of domain-specific purifying selection. *BMC Evol Biol.* 8:165.
- Chelliah V, Chen L, Blundell TL, Lovell SC. 2004. Distinguishing structural and functional restraints in evolution in order to identify interaction sites. *J Mol Biol.* 342(5):1487–1504.
- Chen LB. 1988. Mitochondrial-membrane potential in living cells. *Annu Rev Cell Biol.* 4:155–181.
- Dalziel AC, Stewart DT. 2002. Tissue-specific expression of male-transmitted mitochondrial DNA and its implications for rates of molecular evolution in *Mytilus* mussels (*Bivalvia: Mytilidae*). *Genome* 45(2):348–355.
- Dingley SD, et al. 2014. Mitochondrial DNA variant in COX1 subunit significantly alters energy metabolism of geographically divergent wild isolates in *Caenorhabditis elegans*. *J Mol Biol.* 426(11):2199–2216.
- Diz AP, et al. 2013. Proteomic analysis of eggs from *Mytilus edulis* females differing in mitochondrial DNA transmission mode. *Mol Cell Proteomics* 12(11):3068–3080.

- Everett EM, Williams PJ, Gibson G, Stewart DT. 2004. Mitochondrial DNA polymorphisms and sperm motility in *Mytilus edulis* (Bivalvia: Mytilidae). *J Exp Zool Part A* 301(11):906–910.
- Felsenstein J. 1985. Phylogenies and the comparative method. *Am Nat.* 125(1):1–15.
- Garvin MR, Bielawski JP, Sazanov LA, Gharrett AJ. 2015. Review and meta-analysis of natural selection in mitochondrial complex I in metazoans. *J Zool Syst Evol Res.* 53(1):1–17.
- Ghiselli F, Milani L, Passamonti M. 2011. Strict sex-specific mtDNA segregation in the germ line of the DUI species *Venerupis philippinarum* (Bivalvia: Veneridae). *Mol Biol Evol.* 28(2):949–961.
- Ghiselli F, et al. 2012. De novo assembly of the manila clam *Ruditapes philippinarum* transcriptome provides new insights into expression bias, mitochondrial doubly uniparental inheritance and sex determination. *Mol Biol Evol.* 29(2):771–786.
- Ghiselli F, et al. 2013. Structure, transcription, and variability of metazoan mitochondrial genome: perspectives from an unusual mitochondrial inheritance system. *Genome Biol Evol.* 5(8):1535–1554.
- Gong S, Blundell TL. 2008. Discarding functional residues from the substitution table improves predictions of active sites within three-dimensional structures. *PLoS Comput Biol.* 4(10):e1000179.
- Grantham R. 1974. Amino-acid difference formula to help explain protein evolution. *Science* 185(4154):862–864.
- Haag-Liautard C, et al. 2008. Direct estimation of the mitochondrial DNA mutation rate in *Drosophila melanogaster*. *PLoS Biol.* 6(8):1706–1714.
- Haig D, Hurst LD. 1991. A quantitative measure of error minimization in the genetic-code. *J Mol Evol.* 33(5):412–417.
- Hebert PDN, Cywinska A, Ball SL, DeWaard JR. 2003. Biological identifications through DNA barcodes. *Proc Biol Sci.* 270(1512):313–321.
- Huang YF, Golding GB. 2015. FuncPatch: a web server for the fast Bayesian inference of conserved functional patches in protein 3D structures. *Bioinformatics* 31(4):523–531.
- Hughes AL, Ota T, Nei M. 1990. Positive Darwinian selection promotes charge profile diversity in the antigen-binding cleft of class I major-histocompatibility-complex molecules. *Mol Biol Evol.* 7:515–524.
- Ivankov DN, Finkelstein AV, Kondrashov FA. 2014. A structural perspective of compensatory evolution. *Curr Opin Struct Biol.* 26:104–112.
- James JE, Piganeau G, Eyre-Walker A. 2016. The rate of adaptive evolution in animal mitochondria. *Mol Ecol.* 25(1):67–78.
- Jha M, Côté J, Hoeh WR, Blier PU, Stewart DT. 2008. Sperm motility in *Mytilus edulis* in relation to mitochondrial DNA polymorphisms: implications for the evolution of doubly uniparental inheritance in bivalves. *Evolution* 62(1):99–106.
- Johnson MS, Sali A, Blundell TL. 1990. Phylogenetic-relationships from 3-dimensional protein structures. *Method Enzymol.* 183:670–690.
- Kenyon L, Moraes CT. 1997. Expanding the functional human mitochondrial DNA database by the establishment of primate xenomitochondrial cybrids. *Proc Natl Acad Sci U S A.* 94(17):9131–9135.
- Kidera A, Konishi Y, Oka M, Ooi T, Scheraga HA. 1985. Statistical-analysis of the physical-properties of the 20 naturally-occurring amino-acids. *J Protein Chem.* 4(1):23–55.
- Konrad A, et al. 2017. Mitochondrial mutation rate, spectrum and heteroplasmy in *Caenorhabditis elegans* spontaneous mutation accumulation lines of differing population size. *Mol Biol Evol.* 34(6):1319–1334.
- Koshi JM, Goldstein RA. 1997. Mutation matrices and physical-chemical properties: correlations and implications. *Proteins* 27(3):336–344.
- Kumar S, Stecher G, Tamura K. 2016. MEGA7: Molecular Evolutionary Genetics Analysis version 7.0 for bigger datasets. *Mol Biol Evol.* 33(7):1870–1874.
- Lanza IR, Nair KS. 2009. Functional assessment of isolated mitochondria in vitro. *Method Enzymol.* 457:349–372.
- Liu N, Wang TM. 2006. Protein-based phylogenetic analysis by using hydrophathy profile of amino acids. *FEBS Lett.* 580(22):5321–5327.
- Ludwig B, et al. 2001. Cytochrome c oxidase and the regulation of oxidative phosphorylation. *Chembiochem* 2(6):392–403.
- Malkov SN, Zivković MV, Beljanski MV, Hall MB, Zarić SD. 2008. A re-examination of the propensities of amino acids towards a particular secondary structure: classification of amino acids based on their chemical structure. *J Mol Model.* 14(8):769–775.
- McClellan DA. 2013. Directional Darwinian selection in proteins. *Bioinformatics* 14(Suppl 13):S6.
- McClellan DA, McCracken KG. 2001. Estimating the influence of selection on the variable amino acid sites of the cytochrome b protein functional domains. *Mol Biol Evol.* 18(6):917–925.
- Melvin RG, Katewa SD, Ballard JWO. 2008. A candidate complex approach to study functional mitochondrial DNA changes: sequence variation and quaternary structure modeling of *Drosophila simulans* cytochrome c oxidase. *J Mol Evol.* 66(3):232–242.
- Milani L. 2015. Mitochondrial membrane potential: a trait involved in organelle inheritance? *Biol Lett.* 11(10):20150732.
- Milani L, Ghiselli F. 2015. Mitochondrial activity in gametes and transmission of viable mtDNA. *Biol Direct* 10:22.
- Milani L, Ghiselli F, Guerra D, Breton S, Passamonti M. 2013. A comparative analysis of mitochondrial ORFans: new clues on their origin and role in species with doubly uniparental inheritance of mitochondria. *Genome Biol Evol.* 5(7):1408–1434.
- Milani L, Ghiselli F, Iannello M, Passamonti M. 2014. Evidence for somatic transcription of male-transmitted mitochondrial genome in the DUI species *Ruditapes philippinarum* (Bivalvia: Veneridae). *Curr Genet.* 60(3):163–173.
- Mizuguchi K, Deane CM, Blundell TL, Johnson MS, Overington JP. 1998. JOY: protein sequence-structure representation and analysis. *Bioinformatics* 14(7):617–623.
- Mullins JGL, et al. 2011. Molecular modelling of the emergence of azole resistance in *Mycosphaerella graminicola*. *PLoS One* 6(6):e20973.
- Naz RK, Rajesh PB. 2004. Role of tyrosine phosphorylation in sperm capacitation/acrosome reaction. *Reprod Biol Endocrinol.* 2:75.
- Notredame C, Higgins DG, Heringa J. 2000. T-Coffee: a novel method for fast and accurate multiple sequence alignment. *J Mol Biol.* 302(1):205–217.
- Obata M, Sano N, Komaru A. 2011. Different transcriptional ratios of male and female transmitted mitochondrial DNA and tissue-specific expression patterns in the blue mussel, *Mytilus galloprovincialis*. *Dev Growth Differ.* 53(7):878–886.
- Ohta T. 1992. The nearly neutral theory of molecular evolution. *Annu Rev Ecol Syst.* 23(1):263–286.
- Overington J, et al. 1990. Applications of environment specific amino-acid substitution tables to identification of key residues in protein tertiary structure. *Curr Sci.* 59:867–874.
- Pace CN, et al. 2014. Contribution of hydrogen bonds to protein stability. *Protein Sci.* 23(5):652–661.
- Pacelli C, et al. 2011. Tight control of mitochondrial membrane potential by cytochrome c oxidase. *Mitochondrion* 11(2):334–341.
- Pettersen EF, et al. 2004. UCSF chimera—a visualization system for exploratory research and analysis. *J Comput Chem.* 25(13):1605–1612.
- Pettit FK, Bare E, Tsai A, Bowie JU. 2007. HotPatch: a statistical approach to finding biologically relevant features on protein surfaces. *J Mol Biol.* 369(3):863–879.
- Pierron D, et al. 2012. Cytochrome c oxidase: evolution of control via nuclear subunit addition. *Biochim Biophys Acta* 1817(4):590–597.
- Plazzi F, Puccio G, Passamonti M. 2016. Comparative large-scale mitogenomics evidences clade-specific evolutionary trends in mitochondrial DNAs of Bivalvia. *Genome Biol Evol.* 8(8):2544–2564.
- Pond SLK, Frost SDW. 2005. Not so different after all: a comparison of methods for detecting amino acid sites under selection. *Mol Biol Evol.* 22(5):1208–1222.

- Popovic DM. 2013. Current advances in research of cytochrome c oxidase. *Amino Acids* 45:1073–1087.
- Puslednik L, Yeates DK, Faith DP, Ballard JWO. 2012. Protein-protein interactions of the cytochrome c oxidase DNA barcoding region. *Syst Entomol.* 37(1):229–236.
- Quesada H, Warren M, Skibinski DOF. 1998. Nonneutral evolution and differential mutation rate of gender-associated mitochondrial DNA lineages in the marine mussel *Mytilus*. *Genetics* 149:1511–1526.
- Ruiz-Pesini E, et al. 2000. Human mtDNA haplogroups associated with high or reduced spermatozoa motility. *Am J Hum Genet.* 67(3):682–696.
- Sali A, Blundell TL. 1993. Comparative protein modeling by satisfaction of spatial restraints. *J Mol Biol.* 234:779–815.
- Schmidt TR, Wu W, Goodman M, Grossman LI. 2001. Evolution of nuclear- and mitochondrial-encoded subunit interaction in cytochrome c oxidase. *Mol Biol Evol.* 18(4):563–569.
- Schmidt TR, et al. 2005. Rapid electrostatic evolution at the binding site for cytochrome c on cytochrome c oxidase in anthropoid primates. *Proc Natl Acad Sci U S A.* 102(18):6379–6384.
- Skibinski DO, Gallagher C, Beynon CM. 1994. Sex-limited mitochondrial DNA transmission in the marine mussel *Mytilus edulis*. *Genetics* 138(3):801–809.
- Sneath PH. 1966. Relations between chemical structure and biological activity in peptides. *J Theor Biol.* 12(2):157–195.
- Starr TN, Thornton JW. 2016. Epistasis in protein evolution. *Protein Sci.* 25(7):1204–1218.
- Stewart DT, Kenchington ER, Singh RK, Zouros E. 1996. Degree of selective constraint as an explanation of the different rates of evolution of gender-specific mitochondrial DNA lineages in the mussel *Mytilus*. *Genetics* 143(3):1349–1357.
- Sugitani R, Stuchebrukhov AA. 2009. Molecular dynamics simulation of water in cytochrome c oxidase reveals two water exit pathways and the mechanism of transport. *Biochim Biophys Acta* 1787(9):1140–1150.
- Tamuri AU, dos Reis M, Goldstein RA. 2012. Estimating the distribution of selection coefficients from phylogenetic data using sitewise mutation-selection models. *Genetics* 190(3):1101–1115.
- Taverna DM, Goldstein RA. 2002. Why are proteins marginally stable? *Proteins* 46(1):105–109.
- Taylor WR. 1986. The classification of amino-acid conservation. *J Theor Biol.* 119(2):205–218.
- TeSlaa T, Teitell MA. 2014. Techniques to monitor glycolysis. *Method Enzymol.* 542:91–114.
- Thorne JL, Goldman N, Jones DT. 1996. Combining protein evolution and secondary structure. *Mol Biol Evol.* 13(5):666–673.
- To J, et al. 2017. Channel-inactivating mutations and their revertant mutants in the envelope protein of infectious bronchitis virus. *J Virol.* 91:e02158-16.
- Tokuriki N, Tawfik DS. 2009. Stability effects of mutations and protein evolvability. *Curr Opin Struct Biol.* 19(5):596–604.
- Toth-Petroczy A, Tawfik DS. 2011. Slow protein evolutionary rates are dictated by surface-core association. *Proc Natl Acad Sci U S A.* 108(27):11151–11156.
- Tsukihara T, et al. 2003. The low-spin heme of cytochrome c oxidase as the driving element of the proton-pumping process. *Proc Natl Acad Sci U S A.* 100(26):15304–15309.
- Visconti PE. 2012. Sperm bioenergetics in a nutshell. *Biol Reprod.* 87(3):72.
- Webb B, Sali A. 2014. Comparative protein structure modeling using MODELLER. *Curr Protoc Bioinformatics* 47:5.6.1–32.
- Wolff JN, Ladoukakis ED, Enriquez JA, Dowling DK. 2014. Mitonuclear interactions: evolutionary consequences over multiple biological scales. *Philos Trans R Soc Lond B Biol Sci.* 369(1646):20130443.
- Xia XH, Li WH. 1998. What amino acid properties affect protein evolution? *J Mol Evol.* 47(5):557–564.
- Zhang JZ. 2000. Rates of conservative and radical nonsynonymous nucleotide substitutions in mammalian nuclear genes. *J Mol Evol.* 50(1):56–68.
- Zhen XM, et al. 2015. Increased incidence of mitochondrial cytochrome c oxidase 1 gene mutations in patients with primary ovarian insufficiency. *PLoS One* 10(7):e0132610.
- Zouros E. 2013. Biparental inheritance through uniparental transmission: the doubly uniparental inheritance (DUI) of mitochondrial DNA. *Evol Biol.* 40(1):1–31.
- Zouros E, Oberhauser Ball A, Saavedra C, Freeman KR. 1994. An unusual type of mitochondrial DNA inheritance in the blue mussel *Mytilus*. *Proc Natl Acad Sci U S A.* 91(16):7463–7467.

Associate editor: Dennis Lavrov