# Spectral ordering and biochronology of European fossil mammals

Mikael Fortelius, Aristides Gionis, Jukka Jernvall, and Heikki Mannila

*Abstract.*—Spectral algorithms have been shown to work well in a wide range of situations that involve the task of ordering. When applied to the localities of a set of European Neogene land mammal taxa, spectral ordering relies almost entirely on the most common genera, depends on connectivity more than on length of taxon lists, and is robust to noise from rarer and less connected taxa. The spectral coefficients for localities are highly correlated with known geochronological ages. Although elementary compared with more sophisticated biochronological tools, spectral ordering allows a fast and standardized way to generate biochronological ordering of localities when other information than faunal lists is lacking. Compared with the conventional mammal Neogene (MN) units, spectral ordering of localities appears to lack distinct temporal boundaries in taxon content and render a much lower count of Lazarus events. If, as seems to be the case, biochronology depends mainly on the most common taxa and if evolutionary change is also most clearly reflected in them, then the main evolutionary patterns should be detectable at a modest level of sampling.

*Mikael Fortelius.\* Department of Geology and Institute of Biotechnology, Post Office Box 64, FIN-00014 University of Helsinki, Finland. E-mail: mikael.fortelius@helsinki.fi*

*Aristides Gionis and Heikki Mannila. HIIT Basic Research Unit, Department of Computer Science, Post Office Box 68, FIN-00014 University of Helsinki, Finland*

*Jukka Jernvall. Developmental Biology Program, Institute of Biotechnology, Post Office Box 56, FIN-00014 University of Helsinki, Helsinki, Finland and Department of Ecology and Systematics, Post Office Box 65, FIN-00014 University of Helsinki, Finland*

*\*Corresponding author*

## Introduction

Biochronology, in the original sense of Williams (1901) as ''the endurance of organic characters'' is the only stratigraphic method available when localities lack geological evidence of superposition or geochronologically datable materials, as is the case for much of the terrestrial European Neogene (Lindsay and Tedford 1989; Steininger et al. 1989, 1996; Steininger 1999). Traditionally, the solution has been the creation of more or less arbitrary biochronologic entities into which such fossil occurrences can be grouped. Over the last several decades, biochronologic ordering by numerical methods has become a practically feasible alternative to conventional biochronology, or even, it has been argued, to formal chronostratigraphic systems (Alroy 1998). These methods include, among others, the graph-theoretical unitary associations method (Guex and Davaud 1984; Savary and Guex 1991, 1999), disjunct distribution ordination and appearance event ordination (Alroy 1992, 1994, 1998, 2000; Azanza et al. 1997; Alroy et

al. 1998), parsimony analysis (Martinez 1995; Hooker 1996; Hooker and Weidmann 2000), and Bayesian methods (Halekoh and Vach 2004). So far, they have been used mostly in studies aimed to produce the best possible temporal ordering for a given set of fossil occurrences, and there has understandably been a desire and tendency to include as much information as possible in the procedure, including specialist knowledge such as local geochronologic tiepoints or systematic relationships of the taxa. To the extent that such information is required, this approach limits the applicability of the methods to data sets that are well known to the investigator.

A philosophically and methodologically distinct approach has been developed and progressively refined by John Alroy (Alroy 1992, 1994, 1998, 2000; Azanza et al. 1997; Alroy et al. 1998), who argues that similarity between faunal lists is a poor stratigraphic indicator and that biochronology should be based instead on taxon ranges estimated from the occurrence matrix. Rather than maximize

the fit between taxon lists and a similarity matrix, Alroy's approach is thus based on estimating taxon ranges and maximizing the fit of these hypothesized ranges to independent stratigraphic information, including known stratigraphic superposition of localities. Alroy (2000) gives the latest review of this approach.

Here our purpose is not to present an alternative to the sophisticated tools mentioned above. Rather, we wish to investigate the properties of the biochronological information embodied in taxon occurrence as such, this being the information on which conventional biochronologic systems in general and the MN (mammal Neogene) system in particular have been based. We use the relatively simple standard method of spectral clustering or ordering (Atkins et al. 1999), which has been shown to work well in a wide range of applications, including among others genomic sequencing, domain decomposition in finite element methods, clustering problems in data analysis, VLSI circuit design and simulation (see, e.g., Spielman and Teng 1996; Hagen and Kahng 1992). The data are similar to those of Alroy et al. (1998) and our results support the main findings of that study, but the method allows us to address the nature of the temporal signal in taxon occurrence data. We ask simple questions: How much of the conventional biochronology can spectral ordering capture, using only taxon-locality occurrence information? Does spectral ordering produce clusters corresponding to conventional biochronologic units of the MN system of European Neogene mammal chronology (Mein 1975, 1989; De Bruijn et al. 1992)? How is the biochronologic ordering and clustering of localities affected by the choice of taxa and localities? How does the spectral ordering compare to the MN system based on comparison with geochronological tie points? We also consider the implications of our findings for evolutionary studies using large data sets.

## Data and Methods

*Data Set.*—We used a data set of European late Cenozoic large land mammals derived from the NOW database (http://www.helsinki.fi/science/now) on 17 July 2003. We restricted the data set to the Eurasian conti-



FIGURE 1. Map of Europe showing the locations of the sites containing at least ten genera appearing in at least two sites and (data set G2S10). Numbers indicate eastern longitude (horizontal axis) and northern latitude (vertical axis), respectively.

nent and islands in the Mediterranean Sea, excluding localities with greater than 60° eastern longitude (Fig. 1). We also restricted the data set to the large mammal orders Primates, Creodonta, Carnivora, Perissodactyla, Artiodactyla, Proboscidea, Hyracoidea, and Tubulidentata.

We considered three different kinds of age: database age, real MN age, and geochronologic age. For each locality, we calculated a *database age* as the mean of the minimum and maximum ages given in the original downloaded file. By *MN age* we refer to the mean of the temporal boundaries of MN units according to the correlations given in Steininger et al. (1996). This is given only for the subset of localities assigned in the database to a single MN unit or an interval expressed in MN units. For the MN 9 type locality Can Llobateres we entered a regular MN 9 age in addition to the magnetostratigraphic age provided in the original NOW data set. We also compiled a new age variable by copying all geochronologic (radiometric or magnetostratigraphic) age data in the original data set to a separate variable. This new variable, referred to here as *geochronologic age*, was augmented by data taken from Appendix 2.1 of Steininger et al. (1996) (the main chronology used in the NOW data set) and from recent updates for a set of Greek localities (Sen et al. 2000; Koufos 2004). Because our purpose is not to produce an improved chronology but to investigate the na-

ture of biochronological ordering, we can make use of the slightly obsolete milestone compilation of Steininger et al. (1996) and thereby maximize the size of the data set while minimizing the error from mixing different correlation schemes or chronologies. The main limitation resulting from this choice is that internal comparisons (e.g., geochronologic age versus MN age) become meaningless, because the framework has been adjusted to be internally consistent. The data set is available as supplementary material at http://dx.doi.org/10.1666/04087.s1).

For the spectral ordering we selected further data subsets as follows. First we selected the genera that occurred at least $x$ times in the original data set; then we selected the sites in which at least $y$ such genera had been observed. The data set produced by this method is denoted $GxSy$, for different choices of $x$ and $y$. We used the combinations $(x,y) = (10,10)$, $(10,5)$, $(5,10)$, $(5,5)$, $(2,10)$, and $(2,2)$. Note that, e.g., in the G10S10 data set several genera occur fewer than ten times, as the selection on the number of genera is done first and then the sites are pruned.

*Spectral Ordering.*—Here we give only a short description of the spectral ordering method; more details can be found in Chung (1997) or Atkins et al. (1999); a practical "how to" is given in the appendix. Given $n$ observations $x_1,\ldots,x_n$ and a similarity measure $s(i,j)$ between all pairs of observations $x_i$ and $x_j$, the *Laplacian* matrix is defined to be the $n \times n$ matrix given by

$$L(i, j) = \begin{cases} -s(i, j), & i \neq j \\ \sum_k s(i, k), & i = j. \end{cases} \quad (1)$$

The spectral ordering method computes the eigenvector $v_{n-1} = (v_1,\ldots,v_n)$ that corresponds to the second smallest eigenvalue of the Laplacian matrix $L$. Then the observations $x_i$ are ordered according to the numerical values (increasing or decreasing) of their corresponding coordinates $v_i$ in the eigenvector $v_{n-1}$. The basis for this ordering relies on the fact that the eigenvector $v_{n-1}$ satisfies the conditions

$$\sum_i v_i = 0, \qquad \sum_i v_i^2 = 1, \quad (2)$$

and minimizes the quantity

$$\sum_i s(i, j)(v_i - v_j)^2. \quad (3)$$

The intuitive interpretation of the above conditions is that if two sites $i$ and $j$ are similar, then the difference in the values of the corresponding coordinates $v_i$ and $v_j$ tends to be small so that the contribution $s(i,j)(v_i - v_j)^2$ in the minimization objective is as small as possible. The constraints $\sum_i v_i = 0$, and $\sum_i v_i^2 = 1$ prevent the coordinates of $v_{n-1}$ from taking identical values, which is a trivial solution to the minimization problem. Following the above intuition, we use $v_i$ as a representative for the age of site $i$: if two sites are similar they are likely to have similar age. Note that if $(v_1,\ldots,v_n)$ is the solution to the above conditions, so is $(-v_1,\ldots,-v_n)$; thus the direction of time cannot be determined from the results of spectral ordering.

*Similarity Function between Sites.*—The notion of similarity between occurrence lists has raised lots of discussion (see Alroy 2000). In the application of spectral ordering to the taxonomic matrix we use as the similarity function between the occurrence lists $x_i$ and $x_j$ the number of taxa that occur in both $x_i$ and $x_j$, normalized by the square roots of total number of occurrences of all taxa occurring in $x_i$ and $x_j$. That is, if $c(x_i, x_j)$ is the number of taxa that are occur in both sites and $|t(x_i)|$ is total number of occurrences of the taxa that occur in $x_i$, then the similarity $s(x_i,x_j)$ between $x_i$ and $x_j$ is defined to be

$$s(x_j, x_j) = \frac{c(x_i, x_j)}{|t(x_i)|^{1/2}|t(x_j)|^{1/2}} \quad (4)$$

Using other similarity functions (such as the dot product cosine) between the occurrence vectors gave very similar results.

*Properties of Spectral Ordering.*—From the results of Atkins et al. (1999) it is straightforward to show that if there is an ordering of the sites that does not cause any Lazarus events (temporal gaps in the occurrence of a taxon), then the spectral ordering is such. Also, if a taxon occurs at sites $i$ and $j$, it creates a constraint between their ordering, as the similarity of the sites will be non-zero. Such constraints have an effect also through chains of co-occurrences: if two sites $i$ and $j$ are linked

indirectly by a series of taxa (taxon 1 occurs both at sites $i$ and $k_1$, taxon 2 at sites $k_1$ and $k_2$,. . . , taxon $h-1$ at sites $k_{h-2}$ and $k_{h-1}$, and taxon $h$ at sites $k_h$ and $j$), then there is a connection between the spectral coefficients of $i$ and $j$. That is, orderings in which the two sites are close are preferred; and the strength of the constraint caused by the chain is (at worst) inversely proportional to the square of the length of the chain.

The spectral ordering method has a close similarity to the pioneering approach of Alroy (1992). However, there are two key differences. First, Alroy's algorithm starts by computing an ordering for the species, and then it computes an ordering for the sites (the score of each site is the average of the values of the species values in that site). Second, Alroy's algorithm uses an iterative process to compute a fixed vector of the *stochastic* species-species similarity matrix, whereas the spectral algorithm uses the Laplacian of the sites-sites similarity matrix. The stochastic species-species matrix is a very powerful approach, but the iterative process cannot be interpreted as directly as an optimization problem like the one described by the conditions defined by equations (2) and (3).

*Likelihood Models for Orders.*—Given an ordering of the sites, we considered the occurrence probability of a taxon to be 0 before its first occurrence in the order and after its last occurrence. Between these, the occurrence probability was defined as $p = o/(l-f+1)$, where $l$ is the number of the site with the last occurrence, $f$ the number of the site with the first occurrence, and $o$ the total number of occurrences. The log-likelihood of the taxon is then

$$\log p^{l-f+1} = (l-f+1) \log p. \qquad (5)$$

*Lazarus Event Counts.*—For an ordering of the sites, the count of Lazarus events for a taxon is the number of sites between the first and last occurrences of the taxon in which the taxon does not occur. Note that we count each absence in a contiguous sequence of absences of the taxon individually.

*Comparison against Database Age and MN Age.*—The spectral ordering method produces a continuous estimate of the age of the site,

TABLE 1. Spectral ordering compared with orderings produced by the MN system and by geochronology for different subsets of the data. Correspondence is generally quite good, with correlation coefficients ranging from 0.94 to 0.99 for all comparisons. Abbreviations: gl, the lower limit for the number of occurrences of a genus; sl, the lower limit for the number of genera per site; gn, the number of genera occurring in at least gl sites in the original data set; sn, the number of sites containing at least sl genera occurring in at least gl sites; cDB, the correlation between the database age of sites and the coefficient of spectral ordering; NGC, the number of sites in this selection for which geochronologic age is available; cGC, the correlation between the geochronologic age of sites and the coefficient of spectral ordering; NMN, the number of sites in this selection for which there is an MN age available; cMN, correlation between the spectral ordering coefficient and the MN class of the site.

| gl | sl | gn | sn | cDB | NGC | cGC | NMN | cMN |
|----|----|-----|-----|------|-----|------|-----|------|
| 10 | 10 | 139 | 124 | 0.96 | 21 | 0.98 | 119 | 0.97 |
| 10 | 5 | 139 | 259 | 0.95 | 35 | 0.97 | 230 | 0.96 |
| 5 | 10 | 198 | 136 | 0.97 | 22 | 0.99 | 125 | 0.97 |
| 5 | 5 | 201 | 273 | 0.96 | 35 | 0.98 | 240 | 0.96 |
| 2 | 10 | 281 | 147 | 0.97 | 22 | 0.99 | 132 | 0.97 |
| 2 | 2 | 285 | 512 | 0.94 | 46 | 0.97 | 444 | 0.94 |

whereas in the MN system several sites are assigned to the same class. Because most database ages are based on time units (mostly MN units), database ages also tend to be assigned to classes. To compare the likelihoods of the spectral ordering and the MN system we produced 100 random orderings respecting the MN ordering and computed the likelihoods and Lazarus counts for those. The same method was used for the database ages.

*Outlier Removal.*—The spectral ordering method produces degenerate results if the data contain isolated components, i.e., sets of sites and genera such that the genera occur only on those sites and on those sites only these genera occur. There has to be some connectivity in the site-genus matrix for the method to work. If the spectral coefficient of a site differed by more than two standard deviations from 0, the site was removed from the data set.

## Results

Table 1 shows the basic statistics of the data sets and compares the orderings produced by the spectral method against other orderings. The spectral ordering in general corresponds quite well to both database age and MN class:

FIGURE 2. Spectral ordering compared with conventional age estimates. A, Database age of the sites as a function of the spectral ordering coefficient, for sites containing at least ten genera appearing in at least two sites and (data set G2S10). Sites that lack an MN assignment are indicated by a crossed symbol. B, The spectral ordering coefficient and the geochronologic age of the site for the sites in the G2S10 data set for which geochronologic age information is available.

the correlation is between 0.94 and 0.97. Figure 2A shows in detail the relationship of the database age and the coefficient of the spectral ordering. The spectral coefficient clearly separates three temporal groups, corresponding to the early and middle Miocene, the late Miocene, and the Plio-Pleistocene. Within these temporal groups there is also clear trend from older to younger. However, the spectral ordering does not reproduce lower-order strati-

graphic units, such as individual MN units (e.g., MN 10 versus MN 11). Comparison of spectral ordering against the real MN ages and the geochronologic ages of the sites for which these data are available is also shown in Table 1 and in Figure 2B. We observe that the correlation between the MN age or geochronologic age and the spectral coefficient is very high, between 0.97 and 0.99 for all the subsets.

The likelihoods of the different orderings are given in Table 2, as are the number of Lazarus counts for the different orderings. For both metrics, the spectral ordering outperforms the other orderings by a clear margin. There are no clear spatial trends for the spectral ordering coefficients, although sites from eastern Europe (longitude > 20°E) are overrepresented among the outliers.

Tables 1 and 2 show that the results do not vary much when the selection criteria for genera are changed. The major effects in the data seem to be largely unaffected by removal of rarer genera. For example, removing from the G10S10 data set the 58 genera with the smallest occurrence counts leaves 81 genera. The spectral ordering for this truncated data set has correlation 0.999 with the ordering for the G10S10 data set. Thus the genera with few occurrences had little effect on the results. Using both genera and species combined gave a very similar result, whereas species alone gave lower correlation coefficients (data not shown).

We compared the spectral ordering against the MN system also by clustering the spectral coefficients into 16 temporally ordered clusters. The clustering of the coefficients into 16 classes was obtained by using the standard *k*-means algorithm (Duda et al. 2001). Then we computed the presence/absence matrix for these clusters and the original genera. For the best clustering the Lazarus count for the clustered data is 104, whereas for the MN classification the Lazarus count is 157.

### Discussion

The main principles on which the MN system is said to be based are (1) ordering relative to fixed reference points, either abstract reference horizons or designated actual type localities (Thaler 1966; Fahlbusch 1976), and (2)

TABLE 2. Likelihoods and Lazarus counts for spectral ordering compared with orderings from database age and MN age. For both metrics, spectral ordering gives the best result. Abbreviations: gl, the lower limit for the number of occurrences of a genus; sl, the lower limit for the number of genera per site; Ls, likelihood of the spectral ordering; LMN, likelihood of the MN ordering; Lage, likelihood of the ordering based on the database age estimates; Lazs, number of Lazarus events in the spectral ordering; LazMN, number of Lazarus events in the MN ordering; Lazage, number of Lazarus events in the database age ordering.

| gl | sl | Ls | LMN | Lage | Lazs | LazMN | Lazage |
|----|----|------|------|------|------|-------|--------|
| 10 | 10 | −4881 | −5153 | −4998 | 3792 | 4174 | 3974 |
| 10 | 5 | −9038 | −9573 | −9416 | 9728 | 10,906 | 10,563 |
| 5 | 10 | −6008 | −6455 | −6275 | 5220 | 5901 | 5622 |
| 5 | 5 | −10,723 | −11,340 | −11,132 | 13,003 | 14,638 | 14,147 |
| 2 | 10 | −6904 | −7429 | −7234 | 6398 | 7314 | 6969 |
| 2 | 2 | −16,660 | −17,610 | −17,323 | 30,568 | 34,886 | 33,621 |

ordering based on stage-of-evolution of assemblages or individual lineages (Mein 1975, 1989; De Bruijn et al. 1992). The true stratigraphic nature of the MN ''units'' has been subject to considerable discussion and opinions have diverged widely (Fahlbusch 1991; Steininger 1999), but it is fair to say that it is usually thought to contain more information than is present in a genus-level presence/absence matrix of large mammal taxa and localities. Our comparison with the MN system must be regarded as highly conservative in that we use only presence/absence information at the genus level and exclude the stratigraphically important small mammals entirely.

The temporal sequence of localities produced by spectral ordering of large mammal occurrences is similar in general terms to that produced for an earlier but similar data set by disjunct distribution ordination (Alroy et al. 1998), suggesting that the degree of temporal resolution extracted from the taxonomic information by both methods is similar. Spectral ordering produces a temporal sequence that is similar to those obtained from the overall best age estimate (database age) or from the MN system only. All three perform about equally well with respect to geochronologic age, but in terms of likelihood the spectral ordering outperforms the others by a clear margin.

Our results add to a growing body of evidence that the temporal signal that can be extracted from taxonomic occurrence does not produce clusters corresponding to the conventional biochronologic units. Azanza et al. (1997) felt that ordination by nonmetric multidimensional scaling brought out the European land mammal ages of the latest Miocene to early Pleistocene, but not the MN units, and their plot of disjunct distribution ordination fails to show any distinct clusters for this interval. Disjunct distribution ordination applied to an earlier NOW large mammal data set by Alroy et al. (1998) also failed to show distinct temporal clusters, although the correspondence with MN assignations was generally good (r = 0.933 for the whole set of 654 localities). It thus appears increasingly probable that the boundaries of conventional biochronologic units, usually defined by the first appearance of a limited number of key taxa, do not typically represent major discontinuities in the whole taxonomic occurrence matrix. Although our results have no direct bearing on the *a priori* position that MN units are purely taxonomic constructs and as such cannot have real temporal boundaries (De Bruijn et al. 1992), they do imply that the *taxonomic* boundaries of conventional low-order biochronological units tend to be arbitrary choices rather than reflect distinct horizons of faunal change.

A comparison of subsets selected on the basis of the occurrence counts and connectivity of taxa (Table 1) reveals that all have high correlation with database age and its major component, the MN system. The highest values (0.95–0.97) were obtained for the three sets with the highest connectivity (S10), but even for the set in which only singletons had been removed (G2S2) it was 0.94. Moreover, the correlation between spectral age and geochronologic age was remarkably high, 0.97–0.99 for all sets. It therefore seems that the spectral ordering uses primarily the connectivity of

the most common taxa, and that adding less connected taxa adds relatively little noise. This supports the practice of Alroy (1992) that removing singletons is a sufficient minimum requirement for numerical ordering of occurrence matrices. As described above, removing genera with few occurrences and randomly removing single data points also had little effect on the spectral order produced.

It is obvious from our results (Tables 1, 2) that the connectivity (locality occurrence) of the genera included has much more influence on the spectral order than has the length of the faunal lists. For example, halving the minimum number of genera per faunal list from ten to five increases the Lazarus count from 3792 to 5220 (just over one-third), whereas halving the minimum number of occurrences of each genus from ten to five increases the count to 9728 (more than twice). The same pattern can also be seen when inspecting the fraction of Lazarus events out of the total number of cells in the data matrix. This is intuitively easy to accept, because connectivity can be thought of as the glue that holds the occurrence matrix together. It is also noteworthy that the same pattern holds true for the MN system and database age: both show almost exactly the same relationship although at a slightly higher level of Lazarus counts (Table 2). This agrees well with the finding of Alroy (1992) that the presence of taxa across many lists is more critical for biochronology in general than the length of individual taxon lists.

Spectral ordering minimizes a function related to the number of Lazarus events in the sequence, so it is not surprising that the Lazarus event count should be lower for the spectral order than for the conventional sequences. Because the MN ages represent classes whereas the spectral coefficients represent a continuum, a meaningful comparison requires binning of the coefficients into classes. When this is done, a remarkably great difference is found: 104 Lazarus events for the best binning of spectral coefficients into 15 units against 157 for the MN system, well over one-third more. In principle this could mean either that the spectral ordering is significantly better than the MN system, or that the true matrix contains significantly more gaps than the one produced by spectral ordering. In the latter case, the gaps could be either real, reflecting major geographic range shifts at the scale of MN units, or artificial, reflecting only uneven sampling. Although it is impossible without additional information to reject the hypothesis that spectral clustering produces an unrealistically continuous stratigraphic occurrence pattern for the taxa, it is difficult to suggest a reason for such a result. A possibility that cannot be addressed here is that the inclusion of small mammals would force the localities into a sequence more like that of the MN. Given the high correlation between spectral coefficients and geochronologic ages (Table 1), the conservative conclusion seems to be that spectral ordering is at least as good as its conventional alternatives.

Our results strongly support the view that numerical ordination methods could be a helpful tool to generate a relative temporal sequence when the only alternative is assignment to conventional biochronologic units. The choice of method will depend on the nature of the data and the preferences of the user, although it is obvious that good correlation with independent geochronologic age estimates is a minimum requirement.

It is clear that any limitations that apply generally to numerical ordination techniques must fundamentally apply to all ordination procedures using the same information. In this respect, the temporal accuracy and precision of spectral ordering are a good indication of what can be realistically expected of any ordering that relies exclusively on biochronology. The dominant influence of the most common taxa (those with the highest locality-occurrence counts) for the spectral ordering is methodologically fortunate, because it is these taxa that are also the least sensitive to sampling effects. The taxa that are most common in each stratigraphic unit are also the ones that most strongly show evolutionary trends in response to environmental change (Jernvall and Fortelius 2002; Vermeij and Herbert 2004). Together, these findings suggest that the main patterns of evolutionary change should be detectable with relatively modest sampling levels.

## Acknowledgments

## Literature Cited

Alroy, J. 1992. Conjunction among taxonomic distributions and the Miocene mammalian biochronology of the Great Plains. Paleobiology 18:326–343.

———. 1994. Appearance event ordination: a new biochronologic method. Paleobiology 20:191–207.

———. 1998. Diachrony of mammalian appearance events: implications for biochronology. Geology 26:23–26.

———. 2000. New methods for quantifying macroevolutionary patterns and processes. Paleobiology 26:707–733.

Alroy, J. J., R. L. Bernor, M. Fortelius, and L. Werdelin. 1998. The MN system: regional or continental? Mitteilungen der Bayerischen Staatssammlung von Paläontologie und historischen Geologie 38:243–258.

Atkins, J. E., E. G. Boman, and B. Hendrickson. 1999. A spectral algorithm for seriation and the consecutive ones problem. SIAM Journal on Computing 28:297–310.

Azanza, B., M. T. Alberdi, E. Cerdeño, and J. L. Prado. 1997. Biochronology from latest Miocene to middle Pleistocene in the western Mediterranean area: a multivariate approach. Pp. 567–574 in J.-P. Aguilar, S. Legendre, and J. Michaux, eds. Actes du Congrès BiochroM '97. Mémoires et Travaux de l'Institut de Montpellier, Montpellier.

Chung, F. R. K. 1997. Spectral graph theory. CBMS Regional Conference Series in Mathematics no. 92. American Mathematical Society, Providence, R.I.

De Bruijn, H., R. Daams, G. Daxner-Höck, V. Fahlbusch, L. Ginsburg, P. Mein, and J. Morales. 1992. Report of the RCMNS working group on fossil mammals, Reisensburg 1990. Newsletters in Stratigraphy 26(2/3):65–118.

Duda, R., P. Hart, and D. Stork. 2001. Pattern classification. Wiley, New York.

Fahlbusch, V. 1976. Report on the International Symposium on Mammal Stratigraphy of the European Tertiary. Newsletters in Stratigraphy 5:160–167.

———. 1991. The meaning of MN-zonation: considerations for a subdivision of the European continental Tertiary using mammals. Newsletters in Stratigraphy 24:159–173.

Guex, J., and E. Davaud. 1984. Unitary associations method: the use of graph theory and computer algorithm. Computer and Geoscience 10:69–96.

Hagen, L., and A. B. Kahng. 1992. New spectral methods for ratio cut partitioning and clustering. IEEE Transactions on Computed Aided Design 11:1074–1085, 1992.

Halekoh, U., and W. Vach. 2004. A Bayesian approach to seriation problems in archaeology. Computational Statistics and Data Analysis 45:651–673.

Hooker, J. J. 1996. Mammalian biostratigraphy across the Paleocene-Eocene boundary in the Paris, London and Belgian basins. In R. W. O. B. Knox, R. M. Corfield, and R. E. Dunay, eds. Correlation of the Early Paleogene in northwestern Europe. Geological Society of London Special Publication 101:205–218.

Hooker, J. J., and M. Weidmann. 2000. The Eocene mammal faunas of Mormont, Switzerland. Schweizerische Paläontologische, Abhandlungen 120:1–143.

Jernvall, J., and M. Fortelius. 2002. Common mammals drive the evolutionary increase of hypsodonty in the Neogene. Nature 417:538–540.

Koufos, G. 2004. Late Miocene mammal events and biostratigraphy in the Eastern Mediterranean. Pp. 343–372 in J. W. F. Reumer and W. Wessels, eds. Distribution and migration of Tertiary mammals in Eurasia: a volume in honour of Hans de Bruijn. Deinsea, Utrecht, The Netherlands.

Lindsay, E., and R. Tedford. 1989. Development and application of land mammal ages in North America and Europe, a comparison. Pp. 601–624 in Lindsay et al. 1989.

Lindsay, E. H., V. Fahlbusch, and P. Mein, eds. 1989. European Neogene mammal chronology. Plenum, New York.

Martinez, J. N. 1995. Biochronologie et méthode de parcimonie. Bulletin de la Société Géologique de France 166:517–526.

Mein, P. 1975. Résultats du groupe de travail des vertébrés: Biozonation du Neogène méditerranéen á partir des Mammifères. Pp. 78–81 in J. Senes, ed. Report on activity of the RCMNS Working Group (1971–1975). Bratislava.

———. 1989. Updating of MN zones. Pp. 73–90 in Lindsay et al. 1989.

Savary, J., and J. Guex. 1991. BioGraph, un nouveau programme de construction des corrélations biochronologiques basées sur les associations unitaires. Bulletin du Laboratoire Géologique de la Université de Lausanne 313:317–340.

———. 1999. Discrete biochronological scales and unitary associations: description of the BioGraph computer programme. Mémoires de Géologie (Lausanne) 34:1–282.

Sen, S., G. Koufos, D. Kostopoulos, and L. De Bonis. 2000. Magnetostratigraphy of late Miocene continental deposits of the Lower Axios valley, Macedonia, Greece. Geological Society of Greece Special Publication 9:197–206.

Spielman, D. A., and S.-H. Teng. 1996. Spectral partitioning works: planar graphs and finite element meshes. Proceedings of the 37th Annual Symposium on Foundations of Computer Science, pp. 96–105. ACM Press, New York.

Steininger, F. F. 1999. Chronostratigraphy, geochronology and biochronology of the Miocene ''European Land Mammal Mega-Zones'' (ELMMZ) and the Miocene ''Mammal-Zones (MN-Zones).'' Pp. 9–24 in G. E. Rössner and K. Heissig, eds. The Miocene land mammals of Europe. Dr. Friedrich Pfeil, Munich.

Steininger, F. F., R. L. Bernor, and V. Fahlbusch. 1989. European Neogene marine/continental chronologic correlations. Pp. 15–46 in Lindsay et al. 1989.

Steininger, F. F., W. A. Berggren, D. V. Kent, R. L. Bernor, S. Sen, and J. Agustí. 1996. Circum-Mediterranean Neogene (Miocene-Pliocene) marine-continental chronologic correlations of European mammal units. Pp. 7–46 in R. L. Bernor, V. Fahlbusch, and H.-W. Mittmann, eds. The evolution of western Eurasian Neogene mammal faunas. Columbia University Press, New York.

Thaler, L. 1966. Les rongeurs fossiles du Bas-Languedoc dans leur rapports avec l'histoire des faunes et la stratigraphie du Tertiaire d'Europe. Mémoires de la Muséum National d'Histoire Naturelle C 17:1–295.

Vermeij, G. J., and G. S. Herbert. 2004. Measuring relative abundance in fossil and living assemblages. Paleobiology 30:1–4.

Williams, H. S. 1901. Discrimination of time value in geology. Journal of Geology 9:570–585.

*Appendix*

**Spectral Ordering "How To"**

In this appendix we provide a practical guide on how the spectral-ordering algorithm can be implemented using a mathematical/statistical toolbox like MATLAB. The method is explained as a simple three-step process.

**1. Setting up the data.** We assume that the data are given in a matrix format. For example, if the data contain information about $n$ localities and $m$ taxa, the data matrix would have the following form.

|       | Loc-1 | Loc-2 |         | Loc-n |         |
|-------|-------|-------|---------|-------|---------|
|       | 1     | 1     | $\cdots$ | 0     | Taxon-1 |
| $D =$ | 1     | 0     | $\cdots$ | 1     | Taxon-2 |
|       | $\cdots$ |     |         |       | $\cdots$ |
|       | 0     | 1     | $\cdots$ | 0     | Taxon-m |

The meaning in the above example is that Taxon-1 appears in Localities 1 and 2, Taxon-2 appears in Localities 1 and $n$, and so on. For concreteness we also write the data matrix in a column format $D = [l_1, l_2, \ldots, l_n]$, where $l_i$ is the taxonomic list of the $i^{th}$ locality, or in other words, the $i^{th}$ column vector of the matrix $D$.

**2. Computing the Laplacian matrix.** We first estimate the locality-locality similarity matrix $S$. The matrix $S$ is an $n$x$n$ *symmetric* matrix whose $(i,j)$ entry is computed by the formula

$$s(x_i, x_j) = \frac{c(x_i, x_j)}{|t(x_i)|^{1/2}|t(x_j)|^{1/2}}.$$

Here, the operator $c(\bullet,\bullet)$ indicates the *dot-product* operation between vectors, and $|t(l_i)|$ is the total number of occurrences of the taxa that occur in $l_i$. In MATLAB:

```
> T = diag(1./sqrt(sum(D)));
> S = T*D'*D*T;
```

Next, we compute the $n$x$n$ *diagonal* matrix $A$, whose $(i,i)$ entry is the sum of the $i^{th}$ row of matrix $S$ (and all off-diagonal entries of $A$ are set equal to 0). Now the Laplacian matrix $L(D)$ is the difference of the similarity matrix $S$ from the diagonal matrix $A$.

$$L(D) = A - S$$

The above computation in MATLAB would be

```
> A = diag(sum(S'));
> L = A-S;
```

**3. Computing the spectral coefficients.** First we find the eigenvector-eigenvalue pairs of the Laplacian matrix $L(D)$. Assume that the eigenvector-eigenvalue pairs are denoted by $(\lambda_1, \mathbf{v}_1)$, $(\lambda_2, \mathbf{v}_2)$, $\ldots$, $(\lambda_{n-1}, v_{n-1})$, $(\lambda_n, \mathbf{v}_n)$, sorted from the largest eigenvalue to the smallest, *i.e.*, $\lambda_1 = \lambda_2 = \ldots = \lambda_{n-1} = \lambda_n$. From the structure of the Laplacian matrix $L(D)$, we know that it is always $\lambda_n = 0$ and that $\mathbf{v}_n$ is a constant vector. The spectral coefficients are the coordinates of the eigenvector $\mathbf{v}_{n-1}$, that is, the eigenvector that corresponds to the *second smallest* eigenvalue of $L(D)$ (in the literature the vector $\mathbf{v}_{n-1}$ is also known as the *Fiedler vector*). For our data set, the $i^{th}$ coordinate of $\mathbf{v}_{n-1}$ is taken to be the spectral coefficient of the $i^{th}$ locality.

In MATLAB, the eigenvalue computation can be done using the command eig.

```
> [Vec,Val] = eig(L);
```

And we still need to sort the eigenvalues and take the eigenvector that correspond to the second smallest eigenvalue. This can be done using the command sort.

```
> [s,I] = sort(diag(val));
> Fiedler = vec(:,I(2));
```