

CONSIDERATIONS ON THE EVOLUTION OF QUALITATIVE MULTISTATE TRAITS

JOHN C. AVISE

Department of Zoology
University of Georgia
Athens, Georgia 30602

(Received 7-II-1979; revised 11-V-1979)

ABSTRACT

Simple models for the evolution of qualitative multistate traits are considered, in which the traits are permitted to evolve in time-dependent versus speciation-dependent fashion. Of particular interest are the means and variances of distances for these traits in evolutionary phylads characterized by different rates of speciation, when alternative characters are neutral with respect to fitness, and when the total number of observable characters is limited to small values. As attainable character states are increasingly restricted, mean distance (\bar{D}) in a phylad decreases, regardless of whether evolution is a function of time or of rate of speciation. The ratio of mean distances in species-rich and species-poor phylads of comparable evolutionary age (\bar{D}_R/\bar{D}_P) remains near one when differentiation is proportional to time, even when attainable character states are severely restricted. \bar{D}_R/\bar{D}_P also nears one as a result of restricting character states when differentiation is proportional to rate of speciation, but the effect is not severe unless the number of character states is very small and the probability of change per speciation very large. These and other results are discussed with reference to available data sets on qualitative multistate traits.

INTRODUCTION

Techniques of electrophoresis separate proteins by net electric charge and to some degree by size and shape. Proteins migrate to positions on a gel which are elucidated by protein-specific histochemical stains. Proteins (electromorphs) specified by a given gene locus may exhibit identical mobility, with the inference that the alleles encoding them are identical. This may often not be the case because of the redundancy of the genetic code and because many amino acid changes leave net protein charge unaffected. Alternatively, electromorphs of a given locus may be different in mobility, with the inference that the alleles encoding them differ by one or more nucleotides. (This may not be true in all cases due to the possibility of post-translational modification – Finnerty and Johnson, 1979). At any rate, there is no information inherent in absolute protein mobility which permits a decision about relative degree of protein similarity – electromorphs can only be scored as ‘same’ or ‘different.’ Thus electromorphs produced by a given gene probably belong to the broader class of qualitative multistate characters, whose states cannot be arrayed in some obvious order of degree of difference.

Since electrophoretic techniques are now routinely employed to estimate genetic distances among populations or species, it is important to ask what the possible consequences may be of restrictions in the number of observable electromorphs. From our own laboratory experiences, it often appears that electromorphs of a given protein exhibit a rather restricted range of gel mobilities, and that the total number of electromorphs is far less than the number of species in the phylad examined. It is certainly likely that the rate of electromorph evolution in such phylads is slow compared to the rate of species origin. However, the possibility should not be excluded *a priori* that electromorphs arise rapidly, but that for technical or other reasons (*i.e.*, convergence to a few states by natural selection) only a very finite number of electromorphs is observed. If the latter were true, attempts to reconstruct evolutionary relationships could be seriously compromised. Similar considerations apply to morphological or other qualitative multistate traits.

This report examines the effects of artificially limiting the number of qualitative character states attainable by members of an evolutionary phylad. Distances will be calculated for these characters between all pairs of living species in phylads of various sizes exhibiting varying rates of character state evolution. We will build upon previously introduced models (Avice and Ayala, 1975; Avice, 1978) which permit evolution to occur at either a) time-dependent or b) speciation-dependent rates. For simplicity, only neutral character states are considered. Such simplifying assumptions have proved to be of considerable heuristic value in a wide variety of population models.

THE MODELS

Imagine that evolutionary conversion between character states of a given trait occurs rapidly with respect to the age of a phylad, and that the number of attainable character states for the trait is small compared to the number of living species in the phylad. If the states are neutral and qualitative, interconversion among all attainable states will occur at random, and many distantly related pairs of species will share character states. The final mean distance and frequency distribution of distances among species in the phylad will depend on the number of attainable character states, the number of species in the phylad, and the evolutionary rate of conversion among character states (or the rate of 'saturation' of the available character space).

In the following discussion, a measure of distance analagous to that proposed by Nei (1972) for electrophoretic information is employed, and to complete the analogy for electrophoretic data, 'electromorphs' may be read for 'character states.' However, results should be generally applicable to qualitative multistate traits.

Definitions. Given a phylad X_i , let

t_i \equiv number of time units since origin (first speciation) of the phylad;

m_i \equiv number of time units between speciations (clads);

- $k_i \equiv t_i/m_i$ (k_i assumed to be an integer);
 $I_{ab} \equiv$ a measure of similarity between two extant species a and b ;
 $d_{ab} \equiv -\ln I_{ab} \equiv$ distance between two extant species a and b ;
 $N_i \equiv$ number of living species in the phylad;
 $S_i \equiv$ number of possible character states per trait available for assumption by members of the phylad;
 $P_{t_i} \equiv$ probability per unit time of a 'change' from one character state to another between two species; in other words, the percentage of traits in two species exhibiting a change of state during a time unit since their separation;
 $P_{c_i} \equiv$ probability per clad of a 'change' from one character state to another; in other words, the percentage of traits exhibiting a change of state during a speciation.

As we will employ these definitions, the 'changes' will usually, though need not, be reflected in a discernible difference between species. The probability that the 'change' increases distance between species is a function of the number of attainable character states (see models).

Thus (dropping subscripts for simplicity), the mean distance among living members of an evolutionary phylad (\bar{D}) is $\Sigma d/C$, where C is the number of pairwise species comparisons, given by

$$C = \frac{N(N-1)}{2} \quad (1)$$

and the variance of distance among living species is evaluated by

$$s_d^2 = \frac{1}{C} \Sigma (d - \bar{D})^2 \quad (2)$$

Since distance values among species belonging to an evolutionary phylad are correlated (Ohta and Kimura, 1971), variances calculated according to this formula cannot be appropriately evaluated by standard statistical tests. Rather, their variance simply represents a descriptive summary of dispersion.

In the following models, it is assumed that two species arise per speciation event, speciations occur at regular time intervals in all lineages, and no species go extinct.

Model 1.

Time divergence model. The distance between a pair of species is assumed to be proportional to the time elapsed since they last shared a common ancestor, discounted in the case of finite possible character states by the probability that 'changes' which have occurred have not resulted in different character states in a pair of species.

The probability that two species share a character state inherited directly from their ancestor without change, is simply $(1-P_i)^{im}$, where im is the age of their latest common ancestor. This may be called the similarity by descent, and is equivalent to the similarity between a pair of species when all

evolutionary change is divergent (infinite number of possible character states). For finite numbers of character states, this similarity must be inflated by the probability that if one or more changes have occurred since separation, the change has been to the identical character state of another species. For any pair of species this is given by $(1 - (1 - P_t)^{im})/S$ (this assumes that when a 'change' in character state occurs, any possible character can be achieved with equal probability). The overall similarity between two species (\bar{I}) is the sum of similarity by descent and similarity given by this latter term.

The distance between a pair of species ($d = -\ln I$), multiplied by the number of pairwise comparisons among species (given by terms of the expansion 2^{k+i-2} , $i=1, k$), yields the overall distance between pairs of species of given age within the phylad. The mean distance among living representatives of an evolutionary phylad is then the sum of overall distances between species pairs of given age, divided by the total number of pairwise species comparisons.

The average distance for members of a phylad is thus:

$$\bar{D} \text{ (time model)} = \frac{2}{N(N-1)} \sum_{i=1}^{i=k} \left(-\ln \left((1 - P_t)^{im} + \frac{1 - (1 - P_t)^{im}}{S} \right) \right) (2^{k+i-2}) \quad (3)$$

Typical numerical values of \bar{D} for phylads with assorted values of S , N , and P_t and with $t=8$ are presented in Fig. 1. Absolute mean distance in an evolutionary phylad consistently decreases as the number of attainable character states decreases. However, the ratio of mean distances in species-rich versus species-poor phylads (\bar{D}_R/\bar{D}_P) remains ≈ 1 , no matter how severely the character states are limited or how rapid the evolutionary conversion among them. Typical numerical values for the variances of distances in phylads with assorted values of S , N , P_t , and with $t=8$ are presented in Fig. 2. Variance of distance is consistently decreased as character states are restricted, and the smallest variances are observed in species-rich phylads.

Model 2.

Clad-divergence model. The distance between a pair of species is assumed to be proportional to the number of speciations in their evolutionary history discounted in the case of finite possible character states by the probability that 'changes' which have occurred have not resulted in different character states in a pair of species. An example of a procedure for calculating \bar{D} for a phylad with $S=100$, $N=16$, and $P_c=0.25$ is given in Table 1, and is exactly analogous to the procedure employed in the time-divergence model.

The average distance for living members of an evolutionary phylad when genetic change is a function of the number of speciation events separating species is thus:

$$\bar{D} \text{ (clad model)} = \frac{2}{N(N-1)} \sum_{i=1}^{i=k} \left(-\ln \left((1 - P_c)^{2i-1} + \frac{1 - (1 - P_c)^{2i-1}}{S} \right) \right) (2^{k+i-2}) \quad (4)$$

Table 1. Example of procedure for calculating \bar{D} when genetic distance is proportional to the number of clads separating species. The phylad has 20 possible character states, 16 living species, and a 25 per cent probability of a change of character per clad.

#clads separating species	probability of identity by descent		probability of identity by 'change' to a common state		overall similarity (I)	distance ($d = -\ln I$)	number pairwise species comparisons		overall distance	
	specific (2i-1, i=1, k)	general ((1-P _c) ²ⁱ⁻¹ ; i=1, k)	specific	general ($\frac{1-(1-P_c)^{2i-1}}{S}$; i=1, k)			specific	general (2 ^{k+i-2} ; i=1, k)		
1	0.750	(1-P _c) ¹	0.012	$\frac{1-(1-P_c)^1}{20}$	0.762	0.272	8	2 ³	2.18	
3	0.422	(1-P _c) ³	0.029	$\frac{1-(1-P_c)^3}{20}$	0.451	0.796	16	2 ⁴	12.74	
5	0.237	(1-P _c) ⁵	0.038	$\frac{1-(1-P_c)^5}{20}$	0.275	1.290	32	2 ⁵	41.28	
7	0.133	(1-P _c) ⁷	0.043	$\frac{1-(1-P_c)^7}{20}$	0.176	1.737	64	2 ⁶	111.16	
TOTALS							120	$\frac{N(N-1)}{2}$	167.36	
									$\bar{D} = \frac{167.36}{120} = 1.39$	

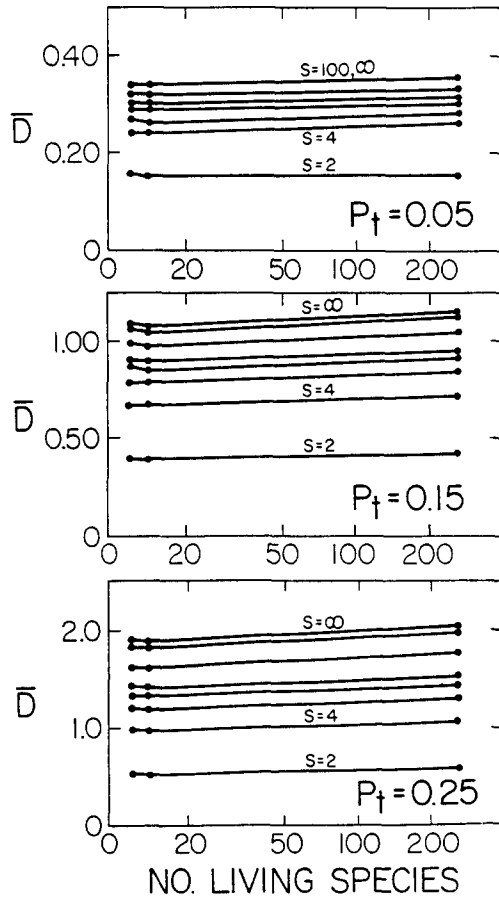


Fig. 1. Mean distances among living representatives of various evolutionary phylads of equal age, when distance is proportional to time, and with the number of assumable character states (S) ranging from 2 to infinity. The probability of 'change' per unit time from one character state to another between two species is P_t .

Typical numerical values of \bar{D} and s_d^2 for phylads with assorted values of S , N , and P_c are presented in Figs. 3 and 4. Mean distance increases as rate of speciation increases, but the rate of increase is slowed and limited as character states are restricted, the magnitude of the effect depending strongly on the rate of saturation of character states. Variances of distance typically exhibit a more complex pattern when character states are limited, first increasing and later decreasing as rate of speciation (reflected in the number of living species) increases. The magnitude of this effect also depends on how rapidly the limited character states can be assumed.

For any number of character states, the maximum attainable mean distance (\bar{D} max) can also be determined for both the time-divergence and clad-divergence models. \bar{D} max is observed when the possible character states

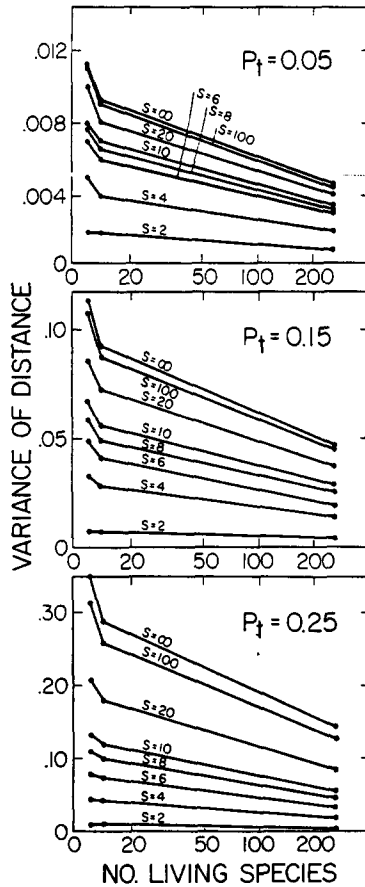


Fig. 2. Variances of distance among living representatives of various evolutionary phylads of equal age, when distance is proportional to time, and with the number of assumable character states (S) ranging from 2 to infinity. The probability of 'change' per unit time from one character state to another between two species is P_t .

are most equitably distributed among species and is given by

$$\bar{D} \max = -\ln \frac{1}{S_i}. \quad (5)$$

This corresponds to the point of maximum entropy or diversity of the Shannon-Wiener information measure for a given S_i . When genetic distance is a function of time, maximum character state diversity (and hence $\bar{D} \max$) occurs when the rate of conversion among character states is sufficiently rapid (instantaneous) with respect to age of the phylad. When genetic distance is a function of speciation rate, $\bar{D} \max$ occurs in a phylad with sufficiently large numbers of species. A schedule of values for $\bar{D} \max$ with the S values used in this study is given in Table 2. These values are rapidly approached in the phylads of Figs. 1 and 3 as P_t and P_c become large.

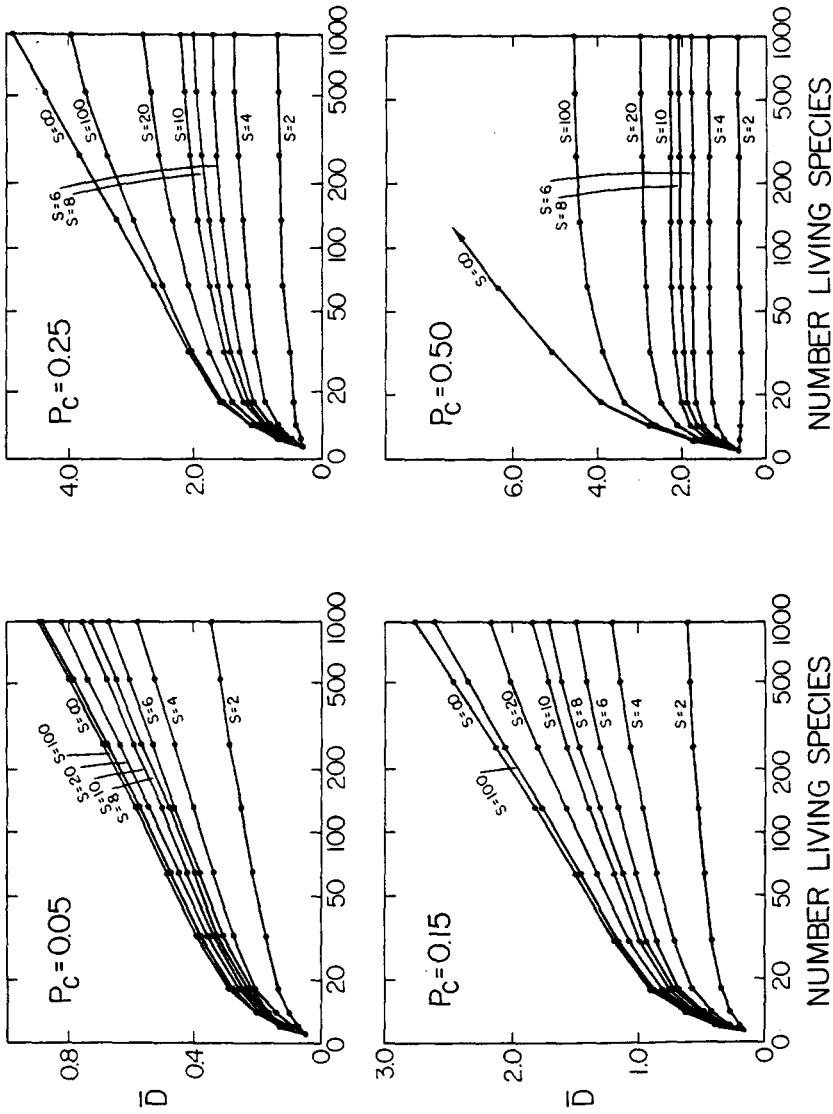


Fig. 3. Mean distances among living representatives of various evolutionary phylads, when distance is proportional to rate of speciation, and with the number of assumable character states (S) ranging from 2 to infinity. The probability of 'change' per speciation from one character state to another is P_c .

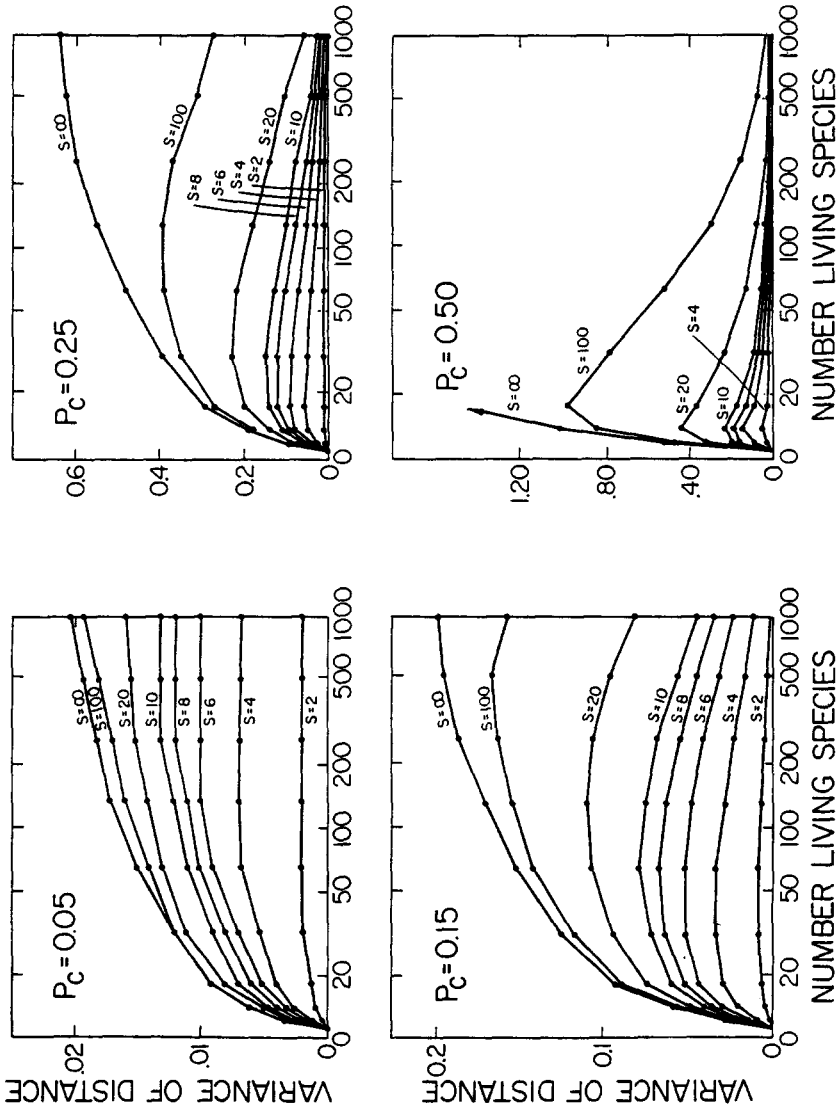


Fig. 4. Variances of distance among living representatives of various evolutionary phylads, when distance is proportional to rate of speciation, and with the number of assumable character states (S) ranging from 2 to infinity. The probability of 'change' per speciation from one character state to another is P_c .

Table 2. Maximum attainable mean distances (\bar{D} max) for evolutionary phylads with S possible character states. \bar{D} max occurs when character state diversity is maximized (\bar{D} max = $-\ln \frac{1}{S}$).

Number of possible character states in a phylad (S)	\bar{D} max
2	0.69
4	1.39
6	1.79
8	2.08
10	2.30
20	2.99
100	4.60
∞	∞

RESULTS

Limiting the number of assumable neutral character states severely depresses absolute mean distances among living members of an evolutionary phylad when the rate of conversion between character states is sufficiently great, no matter whether evolutionary change is a function of time or of rate of speciation. However, the effect on the ratios of mean distances in species-rich versus species-poor phylads of comparable evolutionary age (\bar{D}_R/\bar{D}_P) depends strongly on whether change is time-independent or clad-dependent.

When all evolutionary change is divergent (∞ possible character states), $\bar{D}_R/\bar{D}_P \approx 1$ according to the time-divergence model, and $\bar{D}_R/\bar{D}_P \gg 1$ according to the speciation-dependent model (these results correspond to those previously presented by Avise and Ayala (1975), using a somewhat different approach). In the time-divergence model, \bar{D}_R/\bar{D}_P remains very near one, even when S is severely restricted and P_t is large (Fig. 1). However, in the clad-divergence model, limiting S has the effect of depressing \bar{D}_R/\bar{D}_P below expectations based on infinite possible character states. Limiting numbers of available character states diminishes the distinctness of the theoretical predictions of the time-divergence and clad-divergence models, but the effect is not severe unless S is very small and P_c unreasonably large. Examples of the approach to confluence of clad-divergence and time-divergence predictions are illustrated in Fig. 5.

Limiting character states also has the effect of decreasing the variance of distances in evolutionary phylads. In fact, when evolutionary divergence is a function of rate of speciation, species-rich phylads may exhibit a smaller variance than species-poor phylads, exactly the opposite of predictions when $S = \infty$ (Fig. 4; see also Avise, 1977). Thus $s_{d_R}^2/s_{d_P}^2$ can be < 1 even for reasonably sized phylads ($N = 25-250$, for example) with reasonable levels of change per speciation (*i.e.* $P_c = 0.15$), as long as S is small. Since $s_{d_R}^2/s_{d_P}^2$ is less than one when divergence is a function of time, the qualitative distinctness of time-divergence versus clad-divergence predictions about variance ratios in speciose and depauperate phylads may be seriously compromised as a result of restrictions on numbers of available character states.

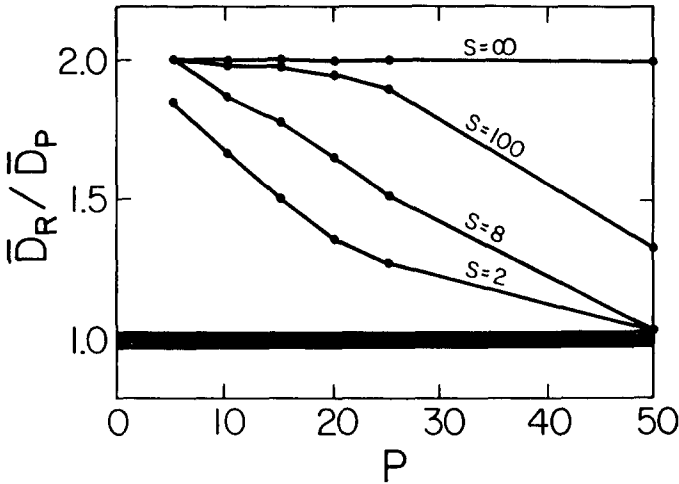


Fig. 5. Expected ratios of distances in a speciose phylad of 128 species versus a depauperate phylad of 16 species calculated according to the 'clad-divergence' model. The solid bar indicates the comparable expected ratios according to the 'time-divergence' model when the phylads are of equal evolutionary age. P is the probability per speciation (or per unit time) of 'change' in character state.

DISCUSSION

For a given trait, the number of character states in an evolutionary phylad may appear limited for a variety of reasons. First, the rate of origin of new character states could be slow compared to age of the phylad. Second, new character states might arise more rapidly but be actively sorted by natural selection to a relatively few states conferring the highest fitnesses upon their bearers. Third, techniques employed to monitor particular traits may be insensitive to some character state differences which do in fact exist. For example, the 'electromorphs' assayed on electrophoretic gels may often be allelically heterogeneous (King and Ohta, 1975; Singh *et al.*, 1976, Coyne, 1976). Of course these three factors may be operating simultaneously to limit observable character states in a given situation.

For whatever reasons, it is true that many phylads exhibit fewer character states than extant species. For example, we have recently assayed 60 species of minnows (Cyprinidae) inhabiting eastern North America, and observed using electrophoretic techniques an average of only 9.5 common electromorphs (character states, S) per locus in the phylad. Among 19 sunfish (Centrarchidae) species similarly assayed, an average of 7.7 common electromorphs per locus were observed (these kinds of observations on electromorphs are not unique or unusual - among 11 species of kangaroo rats (*Dipodomys*) $\bar{S}=5.0$ per locus (Johnson and Selander, 1971) and among 15 species of *Peromyscus*, $\bar{S}=6.0$ per locus (Avise *et al.*, 1974). Certainly similar observations would apply to a variety of qualitative multistate characters in these and other phylads.

The approach employed in this report permits an examination of the

effects of finite attainable neutral character states on distances among living members of evolutionary phylads. This approach may in some cases aid in determining the factors responsible for a given observed distribution of distances. A specific example of such an application will be provided for previously observed genetic distances in two evolutionary phylads of fishes of roughly similar geologic age and different rates of speciation – the species-rich *Notropis* (100 species) and the species-poor *Lepomis* (11 species) (summaries in Avise and Smith, 1977 and Avise, 1977).

Among 10 species of *Lepomis*, mean genetic distance at 14 loci is $\bar{D} = 0.626 \pm 0.028$. Among 47 species of *Notropis* at a similar number of loci, $\bar{D} = 0.619 \pm 0.007$; hence $\bar{D}_R/\bar{D}_P = 0.989$. The results were interpreted as consistent with predictions of the time-divergence model for infinite possible electromorphs, and inconsistent with predictions of the clad-divergence model. Under what circumstances of restricted and neutral character states, if any, could this conclusion be incorrect? As just noted, observed \bar{S} in the two phylads is approximately equal to 8.

The appropriate points to search for on the graphs are those with $\bar{D} \approx 0.60$ over the range of $N = 10-120$, and $S = 8$. This occurs for the time-divergence model, with $P_t \approx 0.10$. Within the approximate range of absolute \bar{D} and $\bar{S} = 8$ under the clad-divergence graphs, the ratio \bar{D}_R/\bar{D}_P in this case should equal 2.78 (see Fig. 3, $P_c = 0.05$), little different from theoretical predictions of 2.85 with infinite possible electromorphs, but greatly different from the observed ratio of 0.99. To look at it another way, in order for the expected ratio under the clad-divergence model to seriously approach the observed ratio with $\bar{S} = 8$, P_c would have to equal nearly 0.50, but in this case the absolute values of \bar{D} would be much higher than those observed. Thus, a joint examination of absolute \bar{D} 's, and of ratios of \bar{D} 's in the two phylads, eliminates the possibility of rapid interconversion among neutral character states as a serious factor confounding the conclusion that time is a more important predictor of genetic differentiation than is rate of speciation for these phylads.

Could the observed distances in sunfish and minnows still be compatible with a clad divergence model if greatly different rate constants (P_c 's) are operative in the two phylads? The answer is 'yes', since in the absence of definitive evidence on the exact branching pattern of these phylads of fishes, any final value of \bar{D} could, at least in theory, be interpreted (partitioned) as having arisen only during speciation episodes. With $\bar{S} = 8$, a \bar{D} of approximately 0.60 is observed for phylads of about 8–16 species (such as *Lepomis*) when $P_c \approx 0.15$ (Fig. 3). Interestingly, this value of P_c closely approximates our best estimate of the amount of genetic differentiation accompanying speciation in the genus *Lepomis*, obtained using an entirely different empirical approach. Two well-differentiated but not yet reproductively isolated subspecies of bluegill, *Lepomis macrochirus macrochirus* and *L. m. purpureus*, exhibit a genetic identity of $\bar{I} = 0.843$, and the very closely related but 'good' species *Lepomis marginatus* and *L. megalotis* show $\bar{I} = 0.852$ (Avise and Smith, 1977). Since these subspecies or species are separated by essentially one speciation, $P_c = 1 - 0.85 = 0.15$.

For the minnows with $\bar{S}=8-10$, a \bar{D} of approximately 0.60 is observed for phylads of about 125-250 species (such as *Notropis* or all North American minnows of the subfamily Leuciscinae) when $P_c \approx 0.05-0.07$ (Fig. 3). Again, surprisingly, this value of P_c closely approximates our current best estimates of the amount of genetic differentiation accompanying speciation among minnows, obtained using a different approach. Despite being placed in different genera, *Hesperoleucus symmetricus* and *Lavinia exilicauda* are very closely related, and have probably speciated very recently (Avise *et al.*, 1975). Between these species, $\bar{I}=0.948$; hence $P_c \approx 0.05$. Furthermore, 9 other pairs of closely related cyprinid species have now been observed with $\bar{D} < 0.10$ (Avise, 1977).

We suspect, however, that the apparent agreement in estimates of P_c by the two approaches is largely fortuitous, since it does not take into account unknown but almost certain differences in these real fish phylads from those artificial phylads simulated in the models; namely, that speciations do not occur at regular time intervals and that no species go extinct. We prefer to emphasize the more conservative conclusion that rapid interconversion among neutral electromorphs is not a serious factor confounding estimated distances in these phylads, and that if rate constants are similar in the two phylads, time is a more important predictor of genetic divergence than is rate of speciation.

The question remains why the observed number of electromorphs per locus appears small in a variety of evolutionary phylads. Limitations in electrophoretic resolution could restrict observable electromorphs to a low number, but if evolutionary conversion among these electromorphs were rapid with respect to age of the phylad, \bar{D} should closely approximate \bar{D} max for any phylad with given \bar{S} . As shown above, this situation seems unlikely. (For minnows and sunfish with $\bar{S} \approx 8$, \bar{D} max = 2.08, not even closely approximated by observed values of \bar{D}). It is more difficult to evaluate the possibility of convergence due to selection for certain common electromorphs, since the exact outcome would be strongly dependent upon the particular mode of selection specified. It is now widely recognized that electrophoretic data, even when analyzed phenetically, yield information that in broad perspective rather closely approximates the probable evolutionary relationships of a group of species as determined through independent information. The slow rate of change among character states relative to the age of the assayed phylads must account for the majority of the similarity between closely related extant species.

ACKNOWLEDGEMENTS

Work was made possible by grant support from the American Philosophical Society. I especially want to thank Dr. Mike Clegg for comments on the manuscript.

REFERENCES

- Avise, J. C. (1978). Variances and frequency distributions of genetic distance in evolutionary phylads. *Heredity*, 40, p. 255–267.
- Avise, J. C. (1977). Is evolution gradual or rectangular? Evidence from living fishes. – *Proc. Nat. Acad. Sci.*, 74, p. 5083–5087.
- Avise, J. C. & F. J. Ayala (1975). Genetic change and rates of cladogenesis. – *Genetics* 81, p. 757–773.
- Avise, J. C., J. J. Smith, & F. J. Ayala (1975). Adaptive differentiation with little genic change between two native California minnows. – *Evolution* 29, p. 411–426.
- Avise, J. C. & M. H. Smith (1977). Gene frequency comparisons between sunfish (family Centrarchidae) populations at various stages of evolutionary divergence. – *Syst. Zool.*, 26, p. 319–335.
- Avise, J. C., M. H. Smith, & R. K. Selander (1974). Biochemical polymorphism and systematics in the genus *Peromyscus*. VI. The boylii species group. – *J. of Mammal.* 55, p. 751–763.
- Coyne, J. A. (1976). Lack of genic similarity between two sibling species of *Drosophila* as revealed by varied techniques. – *Genetics* 84, p. 593–607.
- Finnerty, V. & G. Johnson (1979). Post translational modification as a potential explanation of high levels of enzyme polymorphism: xanthine dehydrogenase and aldehyde oxidase in *Drosophila melanogaster*. – *Genetics*, 91, p. 695–722.
- Johnson, W. E. & R. K. Selander (1971). Protein variation and systematics in kangaroo rats (genus *Dipodomys*). – *Syst. Zool.* 20, p. 377–405.
- King, J. L. & T. Ohta (1975). Polyallelic mutational equilibria. – *Genetics* 79, p. 681–691.
- Nei, M. (1972). Genetic distance between populations. – *Amer. Nat.* 106, p. 283–292.
- Ohta, T. & M. Kimura (1971). On the constancy of the evolutionary rate of cistrons. – *J. Mol. Evol.* 1, p. 18–25.
- Singh, A. S., R. C. Lewontin, & A. A. Felton (1976). Genetic heterogeneity within electrophoretic 'alleles' of xanthine dehydrogenase in *Drosophila pseudoobscura*. – *Genetics* 84, p. 609–629.