# Genetic and linguistic differentiation in the Americas

(population genetics/molecular anthropology/human evolution/Pacific Northwest)

R. H. WARD*, ALAN REDD*†, DIANA VALENCIA*, BARBARA FRAZIER*, AND SVANTE PÄÄBO‡

*Department of Human Genetics, University of Utah, Salt Lake City, UT 84112; ‡Zoologisches Institut der Universität München, Postfach 202136, D-80021 Munich, Germany; and †Department of Anthropology, Pennsylvania State University, University Park, PA 16802

ABSTRACT    The relationship between linguistic differentiation and evolutionary affinities was evaluated in three tribes of the Pacific Northwest. Two tribes (Nuu-Chah-Nulth and Bella Coola) speak Amerind languages, while the language of the third (Haida) belongs to a different linguistic phylum—Na-Dene. Construction of a molecular phylogeny gave no evidence of clustering by linguistic affiliation, suggesting a relatively recent ancestry of these linguistically divergent populations. When the evolutionary affinities of the tribes were evaluated in terms of mitochondrial sequence diversity, the Na-Dene-speaking Haida had a reduced amount of diversity compared to the two Amerind tribes and thus appear to be a biologically younger population. Further, since the sequence diversity between the two Amerind-speaking tribes is comparable to the diversity between the Amerind tribes and the Na-Dene Haida, the evolutionary divergence within the Amerind linguistic phylum may be as great as the evolutionary divergence between the Amerind and Na-Dene phyla. Hence, in the New World, rates of linguistic differentiation appear to be markedly faster than rates of biological differentiation, with little congruence between linguistic hierarchy and the pattern of evolutionary relationships.

A world-wide study of human populations has concluded that since major linguistic phyla exhibit a pattern of clustering that is congruent with genetic clusters defined by genetic markers, linguistic evolution has paralleled the genetic differentiation of our species (1, 2). In agreement with this conclusion, linguistic boundaries in Europe frequently coincide with gene frequency discontinuities (3). However, in the Americas, the relationship between genetic differentiation and linguistic affiliation is less clear cut. Three language phyla have been proposed for the Americas: Amerind, Na-Dene, and Inuit (4, 5). Within each phyla a number of subdivisions can be recognized: Amerind, with 51 stocks and 69 families, is the most linguistically diverse; and Inuit and Na-Dene, with only 2 stocks each, contain much less linguistic diversity. Although there is fairly good correspondence between linguistic affiliation and genetic relationship in Central American tribes (6), the Amerind-speaking tribes of South America generally fail to exhibit correspondence between their genetic relationship and linguistic affiliation at the family level (7, 8). In North America, where all three phyla exist, the situation is varied—some tribes exhibit good agreement between their genetic relationships and linguistic affiliations, but others do not (9). This study was designed to test the hypothesis that the relative magnitude of evolutionary divergence, within and between Pacific Northwest tribes, is congruent with their hierarchy of linguistic differentiation.

We chose the Pacific Northwest for several reasons. (i) This region exhibits the highest degree of linguistic diversity in North America (10) so that tribes within this geographic region speak languages that belong to different phyla and to different stocks. (ii) Since the tribal populations of this region have similar subsistence patterns and population sizes, factors influencing linguistic and genetic change will be similar in each group (10). (iii) Since the geographic distances between the tribes are small and roughly equal, the potentially misleading consequences of recent gene flow between tribes will not be further confounded by correlations between linguistic divisions and geographic separation.

We studied three Pacific Northwestern tribes whose linguistic affiliations span two phyla (Amerind and Na-Dene) and two Amerind families (Wakashan and Salishan) (4, 5). The Na-Dene language phylum is represented by the Haida, and the two Amerind language families (Wakashan and Salishan) are represented by the Nuu-Chah-Nulth and the Bella Coola, respectively. If linguistic and evolutionary divergence are congruent, the accepted linguistic classification predicts that the evolutionary divergence between the Haida and the two Amerind tribes will be substantially greater than the divergence between the two Amerind-speaking tribes. We defined the degree of evolutionary divergence by sequencing the control region of mitochondrial DNA,§ as the rapid rate of evolution of this molecule is suitable for the time scale during which major linguistic differentiation is expected to occur. In addition, the maternal mode of mitochondrial inheritance facilitates the construction of molecular phylogenies without the complicating effects of recombination.

## MATERIALS AND METHODS

**Population Samples.** Most of the traditional members of the Haida, whose language is classified as an "isolate" within the Na-Dene linguistic phylum, are located in two bands on the Queen Charlotte Islands, numbering some 1200 individuals, plus a small population in Alaska. The 41 individuals, who were originally sampled as part of a rheumatic disease survey, represent both Queen Charlotte communities and claimed genealogical descent from Haida individuals believed to have lived on the Queen Charlottes during the late 19th century. There are ≈600 traditional members of the Bella Coola, resident in small communities on the Bella Coola River on the British Columbia mainland. Their language, Bella Coola, is regarded as an isolate within the Salishan language family, which contains 16 languages. The 40 Bella Coola individuals in this study were also randomly selected from patients originally surveyed for rheumatic disease. The 63 Nuu-Chah-Nulth sequences, which have already been described in detail (11), form a representative sample for the 2400 traditional members of this Wakashan-speaking tribe. The Wakashan language family, which is sometimes grouped with Salishan as a separate Amerind subphylum, contains six dialects, or languages, of which three are contained within the Nuu-Chah-Nulth. Intratribal linguistic differentiation is

§The sequences reported in this paper have been deposited in the GenBank data base (accession nos. L20143–L20155).

less marked in the Haida and the Bella Coola, suggesting these two tribes have less population substructure than the Nuu-Chah-Nulth. While the Haida, Bella Coola, and Nuu-Chah-Nulth differ in their linguistic affiliation and in the finer details of their culture, all three tribes share the common motifs of the Pacific Northwest coast cultural area (12). These include a dependence on maritime resources, extensive coastal voyaging, and populations concentrated in discrete communities. In addition, all three tribes experienced roughly equal levels of population decline in the aftermath of European contact, with the Nuu-Chah-Nulth suffering slightly less than the other two groups (13). Hence, the contemporary size of these tribes can be regarded as reflecting their relative effective population sizes, prior to the disruption caused by the arrival of Europeans.

**DNA Sequencing and Data Analysis.** DNA was prepared from serum samples, as described (11). The Haida and Bella Coola serum samples had been stored for >25 years, attesting to the feasibility of using the PCR to recover molecular data from specimens that have been archived for many years. PCR was performed in a 25-$\mu$l reaction volume using the external primer pair L15926 and H16498 (11). Single-stranded DNA was generated by asymmetric PCR using the internal primers L15997 and H16401 (11), which were also used to prime the sequencing reaction for the first 360 nt of the control region, defined by positions 16,024–16,383 in the human reference sequence (14). The pairwise sequence differences presented here are derived from direct sequence comparisons, since correction for multiple hits gave results virtually identical to those obtained by direct comparison. Phylogenetic trees were constructed by application of the maximum likelihood and maximum parsimony algorithms in the PHYLIP package (15), with existence of clades evaluated by bootstrap analysis (16).

## RESULTS

**Intratribal Variation.** The DNA sequence of a 360-bp segment of the mitochondrial control region was determined for 41 individuals of the Haida and 40 individuals of the Bella Coola. The sequences revealed 10 mitochondrial lineages in the Haida sample and 11 lineages in the Bella Coola sample, 4 lineages being found in both tribes. Fig. 1 displays the 17 lineages, 4 of which were previously observed in the Nuu-Chah-Nulth (11), and 13 are, so far as we know, described for the first time in Amerindians (lineages 29–41). It should be noted that lineage 35 found in the Haida does not display the characteristic sequence motif of Amerindian populations and may be due to non-Amerindian admixture.

After accounting for the smaller sample size, the frequency distribution of the 11 lineages observed in the Bella Coola sample was similar to that observed in the Nuu-Chah-Nulth (11). Similarly, the diversity value, $h$ (17), of 90.4% estimated for the Bella Coola is only slightly less than the value of 95.4% estimated for the Nuu-Chah-Nulth. As before (11), the relative effective size of populations was compared by applying Ewens' sampling formula (18) to estimate $\theta$, where $\theta = 2N_f\mu$ and $N_f$ is the effective population size in terms of females and $\mu$ is the mutation rate. With 11 lineages observed in a sample of 40 individuals, $\theta$ for the Bella Coola is estimated to be 4.7, which is 25% of the value of $\theta = 19$ estimated for the Nuu-Chah-Nulth (11), in agreement with the relative population size of these two tribes. In contrast, the diversity value for the Haida is only 70.9%, while the estimated value of $\theta = 3.9$ is disproportionally small, considering the actual population size of this tribe. Finally, the deficit of lineages of intermediate frequency results in a frequency spectrum for the Haida that is inconsistent with the mutation-drift expectation given by the theory of Chakraborty and Weiss (19),

```
           1 1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 3 3 3 3        H   B
  6 6 8 0 0 4 6 6 6 6 8 9 0 3 4 6 6 7 7 8 9 0 0 3 3
  3 9 8 1 6 5 1 3 5 6 6 4 0 7 8 5 7 5 8 1 6 2 4 2 9
  -----------------------------------------------------
  T T C T G C C C C T T T C C T T C T C T C T G T C C T

 8  . . T . A . . . . . . . . . T . . . T . . . A . . . C    1   8
11  . . T . . . . . . . . . . . T . . . T . . . A . . . C   20   3
21  . . . . . . . . . . . . . . T . . . . . . . C . . C      2   5
22  . C . . . . . . . . . . . . T . . . . . . . C . .        -   2
29  . . T . . . . . . . . . . . T . . . T . . . A . . T C   10   -
30  . . T . . . . . . . . . . . T . . . T . T . A . . T C    1   -
31  . . T . . . . . T . . . . . T . . . T . . . A . . T C    1   -
32  . . T . . . . . . . C . . . T . . . T . . . A . . . C    1   -
33  . . T C . . . . . . . . . . T . . . T . . . A . . . C    1   -
34  C . . . . . . . . . . . . . T T . . . C . . C T . .      3   3
35  . . . . . . . . . . . . . . . . . . C . . . . . .        1   -
36  . . T . A . - T . C . . . . T . . . A . . . A . . . C    -   5
37  . . T . . T . . . . . . . . T . C . T . . . A . . . C    -   2
38  . . T . . . . . . . . . . . T . . C T . . . A . . . C    -   1
39  . . T . . . . . . . C . . . T . . . T . . . A . . T C    -   6
40  . . . . . . . . . . C . T . . . . . . . . . C . . C      -   3
41  . . T . . . . . C . C . . . . . . . . . . . . . .        -   2
```

FIG. 1. Definition of mitochondrial lineages found in the Haida and Bella Coola in terms of 25 variable positions in the control region, where position 63 corresponds to position 16,086 in the published human reference sequence (14). Dots indicate identity with the reference sequence, as defined by the upper sequence. The first four sequences, initially found in the Nuu-Chah-Nulth, are identified by their original ID numbers (11). The numbers in columns H and B indicate the number of individuals in the Haida and Bella Coola, respectively, determined to have a specific lineage. In the initial sample of 63 Nuu-Chah-Nulth (11), the first four lineages were observed in 2, 5, 3, and 3 individuals, respectively.

whereas the frequency spectrum of two Amerind tribes is consistent with the mutation-drift expectation.

The distributions of pairwise sequence differences within the three tribes also suggest a different evolutionary history for the Haida compared to the two Amerind tribes (Fig. 2). The Nuu-Chah-Nulth and Bella Coola have virtually identical distributions of intrapopulation sequence differences, with similar mean pairwise sequence differences (5.3 and 5.1 substitutions, respectively), modal values of six or seven substitutions, and <6% of the pairwise differences involving lineages that differ by only a single nucleotide. In contrast, in the Haida sample >34% of the pairwise comparisons involve lineages differing by a single nucleotide and <5% of the comparisons involve lineages that differ by six or seven substitutions. Also, 29% of the pairwise comparisons within the Haida involve identical lineages. As a consequence, the mean pairwise sequence difference in the Haida (2.5 substitutions) is only half the value observed in the two Amerind tribes (Table 1).

**Intertribal Variation.** With 36 of the 41 lineages being unique to a specific tribe, only a small proportion of mitochondrial lineages are found in more than one population. However, the three lineages observed in all three tribes include the most frequent lineage observed in the Haida, and the two most frequent lineages observed in the Bella Coola sample (Fig. 1). Even so, only a small proportion of pairwise lineage comparisons between tribes involves identical lineages—1.9% when comparing the lineages of the two Amerind tribes and 3.8% and 5.3% when the Haida lineages are compared to those of the Nuu-Chah-Nulth and Bella Coola, respectively. The mean pairwise sequence differences among tribes are all very similar (Table 1), with the distances between the two Amerind tribes being slightly greater than either is from the Haida.

**Tree Analysis.** Fig. 3 presents a phylogenetic tree for the 41 lineages, constructed under the maximum likelihood criteria (15). A bootstrapped parsimony analysis indicated that 31 lineages in this tree fall into four clusters that are supported by >49% of 750 replicate samples. These four clusters are the same as those originally observed in the Nuu-Chah-Nulth (11) but have now been joined by lineages from the other two tribes. In addition, 10 lineages that cannot be assigned to
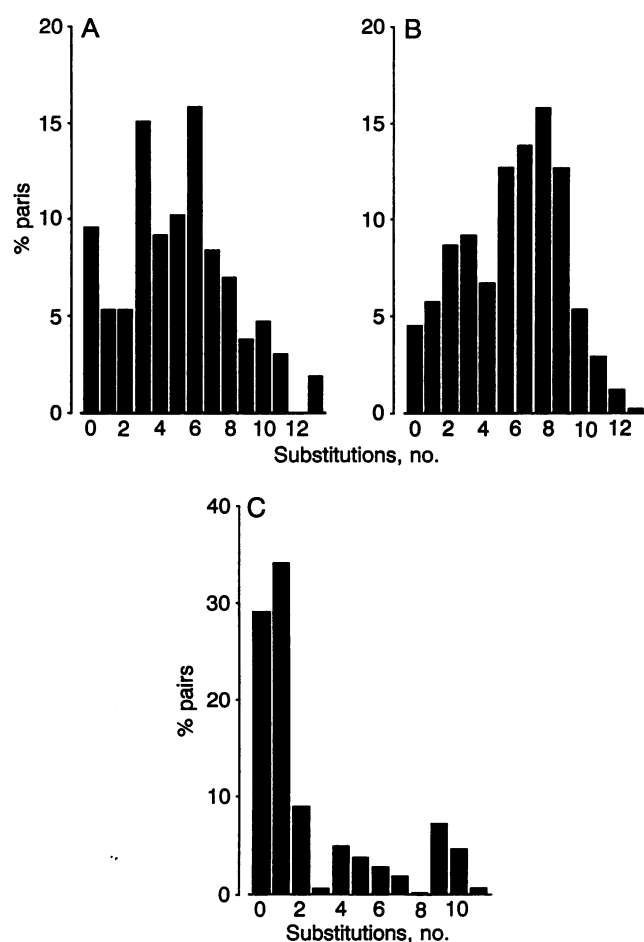
FIG. 2.   Distribution of pairwise sequence difference within the three tribal populations. The number of nucleotide differences between sequences is shown on the $x$ axis, and the percentage of pairs exhibiting a given number of nucleotide differences is indicated on the $y$ axis. (*A*) Bella Coola ($n = 40$). (*B*) Nuu-Chah-Nulth ($n = 63$). (*C*) Haida ($n = 41$).

clusters by bootstrap analysis also exist in the tree. When the tribal affiliations of the mitochondrial lineages are evaluated, it is noted that whereas lineages unique to the two Amerind tribes are distributed throughout the tree, apart from lineage 35, the unique Haida lineages are constrained to a single cluster and differ from one another by only one or two substitutions. It is also noted that the lineages common to all three tribes (lineages 8, 11, and 21) are positioned at internal nodes of the tree. In general, the molecular phylogeny indicates that the lineages found in the Na-Dene-speaking Haida are less dispersed and exhibit less evolutionary divergence, compared to the lineages found in the two Amerind-speaking tribes.

## DISCUSSION

**Sample Size.** When attempting to reconstruct population history from the distribution of DNA sequence variability

Table 1.   Mean pairwise sequence differences (substitutions) in three Pacific Northwest tribes

| | | Between | |
| Tribe | Within | Haida | Bella Coola |
| --- | --- | --- | --- |
| Haida | 2.5 | | |
| Bella Coola | 5.1 | 4.1 | |
| Nuu-Chah-Nulth | 5.3 | 4.7 | 5.5 |



FIG. 3.   Maximum likelihood phylogeny for the 41 lineages found in the three Pacific Northwest tribes. This phylogeny was estimated using the PHYLIP package (15), assuming a 1:30 transversion to transition ratio and using !Kung sequence no. 2 (11) as an outgroup to root the tree. Nodes that fail to give statistically reliable estimates of branching order are indicated by open boxes. Lineages are numbered in accordance with Fig. 1. When analyzed by maximum parsimony, the data gave equivalent results, with all major features of the tree in agreement. The symbols at the tips of the branches indicate in which tribe a lineage was observed: squares, Nuu-Chah-Nulth; circles, Bella Coola; triangles, Haida.

within and among populations, it is critical that the sampling strategy be designed to yield data that is both comparable and informative. This requires that the population samples be large enough to give an adequate estimate of the molecular diversity. To gauge the relationship between sample size and the number of mitochondrial lineages detected, we followed Chakraborty *et al.* (20) in computing the number of lineages expected to be observed in a sample of size $n$ as: $K - \Sigma_{i=1}^{K}$ $(1 - x_i)^n$, where $K$ is the number of lineages in the tribal population and $x_i$ is the frequency of the $i$th lineage. To bracket the range of lineage diversity likely to be found in tribal populations, we used the lineage distribution observed in the Haida to represent tribes with limited mitochondrial sequence diversity, while tribes with more extensive diversity were approximated by our complete Nuu-Chah-Nulth data base, which contains 36 lineages distributed among 120 individuals (21). In the absence of known population frequen-

cies, we used the observed sample frequencies, thereby generating conservative lower bounds for the proportion of lineages expected to be observed in a sample of given size.

As shown in Fig. 4, relatively modest sample sizes will identify most of the lineages present in tribes containing limited mitochondrial sequence diversity, with 63% of lineages being detected with a sample size of 25. However, for tribes with extensive diversity, a sample size of 25 only detects 40% of the lineages and sample sizes of 70 are required to detect two-thirds of the lineages. A sample size of 40 represents a reasonable compromise for most tribes, since this is expected to identify from 50 to 75% of lineages, depending on the actual degree of mitochondrial diversity within the tribe (Fig. 4). Consequently, the diminished sequence diversity observed in the Haida is unlikely to be due to an inadequate sample size.

Although the results in Fig. 4 are most relevant for Pacific Northwest tribes, preliminary results from other geographic regions in the New World suggest that a sample size of 40 will identify the majority of lineages that exist within contemporary tribal populations (data not shown). Hence, attempts to unravel the evolutionary history of aboriginal populations throughout the world (22) will require similar-sized samples to be confident that the distribution of sequence diversity reflects evolutionary relationships, rather than inadequacies of the sampling design.

**Age of Tribal Groups.** The data from the Bella Coola confirm that Amerind-speaking tribes can contain a substantial amount of mitochondrial sequence diversity, with average sequence difference approaching 60% of the value in Subsaharan Africans (11). Further, since the values of $\theta$ estimated for these two tribes reflect the ratio of their standing population sizes, it appears that these two Pacific Northwest tribes have had a similar evolutionary history. Preliminary results from three other Amerind-speaking tribes in the Southwestern and Southeastern United States (data not shown) confirm that the high level of molecular diversity observed within the Nuu-Chah-Nulth (11) is probably typical of most Amerind-speaking tribes of North America.

In contrast, relative to their population size the Haida exhibit much more limited levels of sequence diversity. Unlike the two Amerind tribes, the majority of pairwise sequence comparisons involve closely related lineages and



FIG. 4.    Relationship between sample size and number of lineages observed in the sample. The upper curve represents the proportion of lineages expected to be observed when sampling from a tribal population with limited mitochondrial sequence diversity (e.g., Haida), and the lower curve gives the expected proportion when sampling from a tribe with extensive mitochondrial diversity (e.g., Nuu-Chah-Nulth). The shaded area represents the range in the proportion of lineages expected to be observed for a sample size of 40 individuals.

the mean sequence divergence within the Haida is only 50% of that observed in the Nuu-Chah-Nulth or the Bella Coola (Fig. 2). These observations, along with the disproportionally small estimate for $\theta$, suggest that the Na-Dene-speaking Haida had a population history that was different from the two Amerind-speaking groups. There are three possible explanations for the different distribution of lineage diversity in the Haida: (*i*) They had experienced a more severe bottleneck in the recent past; (*ii*) they have extensive population substructure; (*iii*) they have a more recent origin. As noted above, the postcontact population decimation of the Haida was no worse than that suffered by the Bella Coola and not much worse than that suffered by the Nuu-Chah-Nulth (13). Further, there is no evidence to suggest that during the past several thousand years the ancestral population of the contemporary Haida was substantially smaller, or more isolated, than the ancestral populations of other Pacific Northwest tribes (23, 24). Furthermore, the lack of marked linguistic differentiation within the Haida suggests that their smaller effective size and deficit of moderately frequent lineages is not due to population substructure. Hence, the more likely explanation for the restricted distribution of mitochondrial sequence diversity in the Haida is that this group originated more recently than the two Amerind-speaking groups. In this connection, we note that the Dogrib, a northern Athapaskan-speaking Na-Dene group, also exhibit limited mitochondrial diversity (25), suggesting a shallow time depth for Na-Dene groups as a whole.

By assuming comparable population sizes and similar levels of population subdivision and migration over time, the average divergence time between mitochondrial lineages will tend to reflect the evolutionary age of the population. Previously, we had used empirical comparisons between human and chimpanzee sequences to obtain a tentative estimate of the evolutionary rate for the first 360 nt of the control region as 1% pairwise sequence divergence per 30,500 years (11). However, a more recent estimate based on applying the coalescent technique to a set of Amerindian sequences gave a substantially faster rate of $1.4 \times 10^{-5}$ per nucleotide per year, corresponding to 1% pairwise sequence divergence per 8950 years (26). By using this statistically derived estimate, the distribution of intratribal pairwise sequence differences implies that the average ancestry of mitochondrial DNA molecules in the contemporary Haida may extend back little more than 6200 years. Conversely, the average ancestry of mitochondrial lineages in the two Amerind-speaking tribes would extend back twice as far, to 13,000 years. While these estimated divergence times refer to the ancestry of the mitochondrial lineages rather than the absolute ages of the tribal populations, they suggest that the relative ancestry of the Na-Dene-speaking Haida is much more recent than that of the two Amerind-speaking tribes.

Tree analyses can give additional information about the relative age of molecular lineages and, consequently, about the age of the populations in which the lineages are found. If the ancestral populations of the tribes diverged in the distant past, then in the absence of gene flow, the mitochondrial lineages sampled from the contemporary tribes will tend to cluster within the molecular phylogeny. Conversely, if the population divergences leading to the contemporary tribal groups are relatively recent, lineages from different tribes will show no tendency to cluster. Furthermore, the phylogenetic position of lineages that occur in more than one tribe can be used to gauge whether the lineage sharing is due to recent gene flow or is a consequence of ancestral lineages being retained in several populations. Ancestral lineages will tend to be nodal in position, whereas shared lineages that occur at the tips of the phylogeny are more likely to have become dispersed due to recent gene flow.

In our sample, three lineages are found in all three tribes, and another two are found in two tribes each (Fig. 1). Since the three ubiquitous lineages not only occur at nodal positions within the phylogeny (Fig. 3) but also are among the most numerous lineages in the region, they are most likely ancestral. This is also true for lineage 22, which is shared by the two Amerind tribes. In contrast, lineage 34, which is shared by the Haida and the Bella Coola, is the sole lineage, suggesting recent gene flow since it occurs at a tip in the phylogeny. Consequently, there is little evidence to suggest that recent gene flow has obscured the evolutionary relationships between these tribes. In the absence of extensive gene flow, the "interdigitation" of tribally unique lineages on the molecular phylogeny suggests a relatively recent origin for all three tribes. However, the more restricted distribution of the Haida-specific lineages (especially marked if lineages 34 and 35 are deemed to have entered the Haida by admixture) suggests that the Haida have an even more recent origin than the two Amerind tribes.

**Lack of Congruence Between Genes and Language.** Since language replacement has probably been rare during human evolution (27), it can be assumed that population differentiation has generally preceded linguistic differentiation. Hence, ancient population divergences are likely to result in well-defined genetic clusters that coincide with linguistic boundaries, as has been found for major ethnic groups (1, 2). However, given the antiquity and geographic separation of the major ethnic groups, this congruence could still occur even if the rate of linguistic divergence was substantially different from the rate of genetic divergence. To determine whether genetic and linguistic divergences truly proceed in concert, we have chosen to study Native American populations, which not only display a rich linguistic flora but also have a relatively short evolutionary history that is likely to approximate the time span required for significant linguistic differentiation.

Our data imply that the Na-Dene-speaking Haida are substantially younger than the two Amerind tribes. This is compatible with the view that the ancestors of Amerind speakers preceded the ancestors of Na-Dene speakers to the Americas (25, 28). However, if the pattern of genetic differentiation is consistent with the linguistic hierarchy, then the sequence divergence between the Haida and the two Amerind tribes is expected to be greater than the divergence between the Amerind tribes. The Haida lineages would also form a distinct clade in the molecular phylogeny. Clearly, neither prediction is fulfilled. The intertribal sequence differences are roughly equal among all three tribes, and the distribution of mitochondrial lineages in the phylogeny implies that the ancestors of the Na-Dene originated from an Asiatic population that had earlier contributed a substantial fraction (but not all) of the mitochondrial ancestry of Amerind speakers.

Thus, the hierarchy of genetic relationships fails to match the linguistic hierarchy, implying that the time span required for languages to differentiate into phyla (Na-Dene versus Amerind) and families (Wakashan versus Salishan) may be substantially shorter than that required for molecular differentiation. We conclude that the tempo and mode of genetic and linguistic divergence will frequently be discordant when compared at a high level of resolution (i.e., over a short time span). We speculate that this is because linguistic diversity is generated in a fundamentally different way from genetic diversity. If language change is viewed as a cultural phenomenon that is driven by social and historic events, it is likely to exhibit bursts of change as well as periods of stasis. In contrast, molecular evolution proceeds at a more even pace.

Whether or not the fundamental tempo of linguistic evolution is truly discordant from the pace of genetic evolution can only be addressed by carrying out similarly detailed regional comparisons, in the Americas and elsewhere.

1. Cavalli-Sforza, L. L., Piazza, A., Menozzi, P. & Mountain, J. (1988) *Proc. Natl. Acad. Sci. USA* **85**, 6002–6006.
2. Cavalli-Sforza, L. L., Minch, E. & Mountain, J. (1992) *Proc. Natl. Acad. Sci. USA* **89**, 5620–5624.
3. Barbujani, G. & Sokal, R. R. (1990) *Proc. Natl. Acad. Sci. USA* **87**, 1816–1819.
4. Greenberg, J. H. (1987) *Language in the Americas* (Stanford Univ. Press, Stanford, CA).
5. Ruhlen, M. (1987) *A Guide to the World's Languages* (Stanford Univ. Press, Stanford, CA), Vol. 1.
6. Barrantes, R., Smouse, P. E., Mohrenweiser, H. W., Gershowitz, H., Azofeifa, J., Arias, T. D. & Neel, J. V. (1989) *Am. J. Hum. Genet.* **46**, 63–84.
7. Chakraborty, R. (1976) *Nature (London)* **264**, 350–352.
8. Black, F. L., Salzano, F. M., Berman, L. L., Gabbay, Y., Weimar, T. A., Franco, M. H. L. P. & Pandey, J. P. (1983) *Am. J. Phys. Anthropol.* **60**, 627–635.
9. Spuhler, J. N. (1979) in *The First Americans: Origins, Affinities and Adaptations*, eds. Laughlin, W. S. & Harper, A. L. (Fischer, New York), pp. 135–183.
10. Nichols, J. (1990) *Language* **66**, 475–521.
11. Ward, R. H., Frazier, B. L., Dew-Jaeger, K. & Pääbo, S. (1991) *Proc. Natl. Acad. Sci. USA* **88**, 8720–8724.
12. Driver, H. E. (1969) *Indians of North America* (Univ. of Chicago, Chicago).
13. Duff, W. (1964) *The Indian History of British Columbia* (Provincial Museum of British Columbia, Victoria).
14. Anderson, S., Bankier, A. T., Barrell, B. G., de Bruijn, M. H. L., Coulson, A. R., Drouin, J., Eperon, I. C., Nierlich, D. P., Roe, B. A., Sanger, F., Schreier, P. H., Smith, A. J. H., Staden, R. & Young, I. G. (1981) *Nature (London)* **290**, 457–465.
15. Felsenstein, J. (1991) *PHYLIP: Phylogenetic Inference Package* (Univ. of Washington, Seattle), Version 3.4.
16. Felsenstein, J. (1985) *Evolution* **39**, 783–791.
17. Nei, M. (1987) *Molecular Evolutionary Genetics* (Columbia Univ. Press, New York), pp. 178–179.
18. Ewens, W. J. (1972) *Theor. Popul. Biol.* **3**, 87–112.
19. Chakraborty, R. & Weiss, K. M. (1991) *Am. J. Phys. Anthropol.* **86**, 497–506.
20. Chakraborty, R., Smouse, P. E. & Neel, J. V. (1988) *Am. J. Hum. Genet.* **43**, 709–725.
21. Valencia, D. (1992) M.S. thesis (Univ. of Utah, Salt Lake City).
22. Cavalli-Sforza, L. L., Wilson, A. C., Cantor, C. R., Cook-Deegan, R. M. & King, M. C. (1991) *Genomics* **11**, 490–491.
23. Dumond, D. (1983) in *Ancient Native Americans*, ed. Jennings, J. (Freeman, San Francisco), pp. 43–93.
24. Fladmark, K. R. (1975) *A Palaeoecological Model for Northwest Coast Prehistory* (National Museum of Man, Ottawa).
25. Torroni, A., Schurr, T. G., Yang, C., Szathmary, E. J. E., Williams, R. C., Schanfield, M. S., Troup, G. A., Knowler, W. C., Lawrence, D. N., Weiss, K. M. & Wallace, D. C. (1992) *Genetics* **130**, 153–162.
26. Lundstrom, R., Tavaré, S. & Ward, R. H. (1992) *Proc. Natl. Acad. Sci. USA* **89**, 5961–5965.
27. Renfrew, C. (1988) *Archaeology and Language* (Cambridge Univ. Press, New York).
28. Greenberg, J. H., Turner, C. G., II, & Zegura, L. Z. (1986) *Curr. Anthropol.* **27**, 477–497.