

# Assigning African elephant DNA to geographic region of origin: Applications to the ivory trade

Samuel K. Wasser\*<sup>†</sup>, Andrew M. Shedlock\*<sup>‡</sup>, Kenine Comstock\*<sup>§</sup>, Elaine A. Ostrander<sup>§</sup>, Benzeth Mutayoba<sup>¶</sup>, and Matthew Stephens<sup>||</sup>

\*Department of Biology, Center for Conservation Biology, University of Washington, Box 351800, Seattle, WA 98195; <sup>§</sup>Clinical Research and Human Biology Divisions Fred Hutchinson Cancer Research Center, 1100 Fairview Avenue North, D4-100, Seattle, WA 98109-1024; <sup>¶</sup>Department of Veterinary Physiology, Biochemistry, Pharmacology, and Toxicology, Sokoine University of Agriculture, P.O. Box 3017, Morogoro Tanzania; and <sup>||</sup>Department of Statistics, University of Washington, Seattle, WA 98195

Edited by Henry C. Harpending, University of Utah, Salt Lake City, UT, and approved August 24, 2004 (received for review May 7, 2004)

**Resurgence of illicit trade in African elephant ivory is placing the elephant at renewed risk. Regulation of this trade could be vastly improved by the ability to verify the geographic origin of tusks. We address this need by developing a combined genetic and statistical method to determine the origin of poached ivory. Our statistical approach exploits a smoothing method to estimate geographic-specific allele frequencies over the entire African elephants' range for 16 microsatellite loci, using 315 tissue and 84 scat samples from forest (*Loxodonta africana cyclotis*) and savanna (*Loxodonta africana africana*) elephants at 28 locations. These geographic-specific allele frequency estimates are used to infer the geographic origin of DNA samples, such as could be obtained from tusks of unknown origin. We demonstrate that our method alleviates several problems associated with standard assignment methods in this context, and the absolute accuracy of our method is high. Continent-wide, 50% of samples were located within 500 km, and 80% within 932 km of their actual place of origin. Accuracy varied by region (median accuracies: West Africa, 135 km; Central Savannah, 286 km; Central Forest, 411 km; South, 535 km; and East, 697 km). In some cases, allele frequencies vary considerably over small geographic regions, making much finer discriminations possible and suggesting that resolution could be further improved by collection of samples from locations not represented in our study.**

**T**rade in wildlife products such as ivory, bush meat, and whale meat is capable of rapidly decimating species on a continent-wide scale (1–3). Consequences of poaching often go undetected until the damage becomes too severe to correct, especially for wildlife in difficult to observe habitats such as rain forests or open ocean. The ability to determine origin of wildlife products could help address these problems, providing early warning of where greater law enforcement is needed, and helping guide international decisions regarding delisting species and resanctioning their legal trade. We describe a combination of genetic, sampling, and statistical methods for inferring the geographic origin of elephant DNA that can greatly assist such management problems.

High demand for ivory reduced the African elephant (*Loxodonta africana*) population from 1.3 million to 600,000 individuals between 1979 and 1987 (4). This circumstance prompted the Convention on International Trade in Endangered Species (CITES) to implement a ban on ivory trade in 1989. However, international pressure persisted to resanction trade even though the illicit ivory market continued to thrive (5–7). Three of the largest ivory seizures since the trade ban occurred since June 2002 (Singapore, Hong Kong, and Shanghai, China). Elephant poaching in Central Africa forests has been particularly severe in recent years, where poor visibility has hindered monitoring. Much of this ivory is presumed to be smuggled into international markets by West African countries (8). Monitoring the origin of ivory passing through the major ivory markets around the world would greatly assist efforts to contain such trade.

Our methods rely on (i) noninvasive techniques to acquire DNA from scat (9–10), enabling rapid development of a continent-wide geographic map of elephant allele frequencies; (ii) the ability to

acquire DNA from small amounts of ivory taken anywhere along the length of the tusk (11); and (iii) spatial smoothing methods to estimate allele frequencies at any location in a continuous region, using genotypes of reference samples from a number of sampling locations. Our smoothing approach allows us to improve on standard assignment methods (SAMs) (12), both in terms of assignment accuracy, and more importantly, allowing assignment of samples to locations where no reference samples are available; especially important in our application, because ivory seizures may have come from such locations.

Our method accurately discriminates among DNA samples originating from forest (*Loxodonta africana cyclotis*) and savanna elephants from the four major regions of Africa (West, Central, South, and East). In some cases, allele frequencies vary considerably over small geographic regions, making much finer distinctions possible, and suggesting that resolution could be further improved by collection of samples from locations not represented in our study. We also present the first compelling, to our knowledge, genetic evidence for continued hybridization between forest and savanna elephants, although such hybridization appears to be rare in our samples.

## Materials and Methods

**Samples.** Mitochondrial and microsatellite DNA can be isolated and amplified from small amounts of African elephant ivory taken nearly anywhere along a tusk (11). No attached tissue is required and old samples stored at ambient temperatures can be successfully analyzed. We collected genotypes at 16 microsatellite loci (13–15) by using DNA isolated from 350 tissue samples and 242 fecal samples collected from 28 locations in 16 African countries (Fig. 1). Of these samples, 399 had alleles amplify at seven or more loci, and were used in subsequent analyses. Approximately 85% of Africa's elephants reside in the 16 countries represented in our sample (16).

We used skin biopsy samples analyzed in Comstock *et al.* (14), which were collected by using the method of Karesh *et al.* (17) in full compliance with CITES regulations. Fecal samples were collected with no two samples being closer than 1 km apart to reduce chances of sampling multiple individuals from the same family group. Fecal samples were either preserved in 90% ethanol or 20% DMSO in TNE buffer, recorded for date and location and transported to the United States in compliance with U.S. Department of Agriculture/Animal and Plant Health Inspection Service regulations.

This paper was submitted directly (Track II) to the PNAS office.

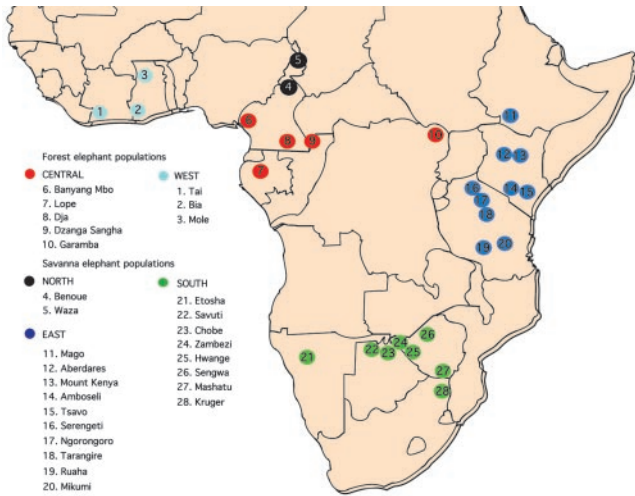
Freely available online through the PNAS open access option.

Abbreviations: CITES, Convention on International Trade in Endangered Species; LR, likelihood ratio; SAM, standard assignment method; CAM, continuous assignment method; ET, Etosha.

<sup>†</sup>To whom correspondence should be addressed. E-mail: wassers@u.washington.edu.

<sup>‡</sup>Present address: Department of Organismic and Evolutionary Biology, Harvard University, 26 Oxford Street, Cambridge, MA 02138.

© 2004 by The National Academy of Sciences of the USA



**Fig. 1.** Map of Africa showing the collection sites divided into five regions: West Africa (cyan), Central forest (red), and Central (black), South (green), and East (blue) savanna.

**DNA Extraction and Amplification.** DNA was extracted in duplicate from fecal samples by using a Qiagen Qiaquick Stool DNA purification kit (Qiagen, Valencia, CA) followed by a Gene Clean III (Bio101) nucleic acid isolation kit for each template isolation. Negative controls were included for every set of 10 extractions completed. We used standard three-step PCR following an adaptation of the multiple tubes approach (18) to identify/minimize allelic dropout (10, 19). All heterozygotes were scored twice and homozygous individuals were examined by several independent amplification products from two independent DNA isolations derived from the same fecal sample. This approach yielded consistent consensus scores for all multiple amplifications included in the final data set.

Amplification of target loci from <50 ng of genomic template was achieved by 35 cycles of standard three-step PCR in 25- $\mu$ l reactions employing 5'-6-carboxy fluorescein-labeled forward oligonucleotides (Applied Biosystems) and *Taq* polymerase antibody hotstart chemistry (Clontech) in an ABI 9700 Gene Amp PCR System. One microliter of PCR products was subjected to fragment analysis by using the GeneScan mode on an ABI capillary array genetic analyzer, model 3100. Allele sizes were scored by using the program GENOTYPER V. 3.7, with categories defined by the weighted average of histograms for each allele size bin (tolerance 0.5 bp).

Raw genotype data from acrylamide vs. Applied Biosystems capillary platforms would be expected to exhibit small allele size differences, consistent across samples within each locus. We reconciled data across platforms in two ways: (i) For each locus, allele size differences (acrylamide minus Applied Biosystems capillary platforms) were computed for control DNAs typed on both platforms, to establish a correction factor to be added to each Applied Biosystems call; and (ii) empirical population allele frequencies for samples from Benoue and Dzanga Sanga were compared across platforms, revealing a single discrepancy after correction as in (i) above: at locus FH40, allele size 229 was frequently observed in acrylamide genotypes, whereas an allele of size 227 was frequently observed in Applied Biosystems genotypes. Because no alleles of similar size occurred on either platform, we concluded these were the same allele.

**Statistical Procedures. Hybrid and assignment analyses.** We used a variation of a standard assignment procedure to assign population of origin to samples, and to identify putative forest-savannah hybrids (see also ref. 20). Let  $g$  denote the sample genotype, with alleles at locus  $l$  being  $(g_{11}, g_{12})$ , and let  $f_{jk}$  denote the frequency of

allele  $j$  at locus  $l$  in population  $k$ . (For the hybrid analysis, we pool our sampling locations into five populations: West and Central Forest, and Central, East, and South savanna, according to geographic location and habitat. For assignment analyses, we treat each sampling location as a separate population.) Assuming Hardy-Weinberg equilibrium and linkage equilibrium within each population, the probability of observing this sample, if the two parents came from populations  $k_1$  and  $k_2$ , can be written as:

$$L(k_1, k_2; g, f) = \Pi p_1(g_{11}, g_{12} | k_1, k_2, f), \quad [1]$$

where

$$\begin{aligned} p_l(i, j | k_1, k_2, f) &= (1 - \gamma) (p_l(i | k_1, f) p_l(j | k_2, f) \\ &\quad + p_l(j | k_1, f) p_l(i | k_2, f)) \\ &\text{if } i \neq j, \text{ and } p_l(i, i | k_1, k_2, f) \\ &= (1 - \gamma) p_l(i | k_1, f) p_l(i | k_2, f) \\ &\quad + \gamma(0.5) (p_l(i | k_1, f) + p_l(i | k_2, f)), \end{aligned}$$

with

$$p_l(i | k, f) = (1 - \delta) f_{jk} + \delta / m_l.$$

Here  $m_l$  is the number of observed alleles (across all populations) at locus  $l$ ,  $\delta$  is a genotyping error probability, and  $\gamma$  is the probability of only one allele amplifying. We fixed  $\delta = 0.05$ ,  $\gamma = 0.1$ . Although in this case, we obtained qualitatively similar results (not shown) when we ignored the possibility of genotyping error (i.e.,  $\delta = 0.0$ ,  $\gamma = 0.0$ ), in general, allowing for the possibility of genotyping error is an important and often overlooked aspect of these kinds of analyses.

We considered an individual a potential hybrid if the maximum of  $L(k_1, k_2)$  over  $k_1$  and  $k_2$  occurred for  $k_1$  in a savannah region and  $k_2$  in a forest region. For standard assignment analyses, we assigned each individual to the population  $k$  that maximized  $L(k, k)$ . In each case, we estimated  $f_{jk}$  as  $(n_{jlk} + 1) / (n_{+l+} + m_l)$ , where  $n_{jlk}$  is the number of times allele  $j$  is observed at locus  $l$  in population  $k$ , and  $n_{+l+} = \sum_{jk} n_{jlk}$ . When analyzing a particular sample, the alleles for that sample were ignored in computation of  $n_{jlk}$  (i.e., we used leave-one-out crossvalidation.)

**Spatial-smoothing-based estimates of allele frequencies.** Our smoothing-based method for estimating  $f_{jk}$  estimates allele frequencies at any location, by using all reference samples, with samples from nearby sampling locations being given more weight. Our method is based on writing

$$f_{jk}(\theta) = \exp(\theta_{jlk}) / \sum_{j'} \exp(\theta_{j'lk}), \quad [2]$$

where the  $\theta$  values corresponding to each locus-allele combination are assumed to be independent Gaussian processes. More specifically, the  $\theta_{jlk}$  have a joint normal distribution, with  $E(\theta_{jlk}) = \mu_{jl}$ , and  $Cov(\theta_{jlk}, \theta_{j'lk'}) = \sigma_{kk'}$ , with  $\theta_{jlk}$  and  $\theta_{j'l'k'}$  being independent if  $l \neq l'$  or  $j \neq j'$ . In this model,  $\mu$  controls the mean allele frequencies (across the study region) at each locus, and the population-specific allele frequencies are allowed to vary about this mean in a spatially correlated way. The value of  $\sigma_{kk'}$  controls the (expected) degree of similarity between allele frequencies in populations  $k$  and  $k'$ . We assume that  $\sigma_{kk'}$  depends only on the distance (in kilometers),  $d_{kk'}$ , between populations  $k$  and  $k'$ , with  $\sigma(d; \alpha) = (1/\alpha_0) \exp[-(\alpha_1 d)^{\alpha_2}]$ , where  $\alpha = (\alpha_0, \alpha_1, \alpha_2)$  are parameters to be estimated. See ref. 20 for other possible parameterizations. The value of  $\alpha_0$  controls the variability at  $d = 0$  (i.e., variability of regional frequencies from the mean);  $\alpha_0 = \infty$  gives no variation in allele frequencies. The parameters  $\alpha_1$  and  $\alpha_2$  control how the correlations decay with distance. Informally,

$\alpha_1$  controls the natural scale (e.g., kilometers, or tens of kilometers) on which the allele frequencies vary, and  $\alpha_2$  controls how quickly the correlations drop to 0. The case  $\alpha_2 = 0$  corresponds to no spatial correlation in frequencies.

The values of  $\alpha$  are estimated from the data in a Bayesian way (21). We use independent priors on components of  $\alpha$ :  $\alpha_0 \sim \Gamma(0.001, 0.001)$ ;  $\log_{10}(\alpha_1) \sim \text{Uniform}(0, 4)$ ;  $\alpha_2 \sim \text{Uniform}[0.1, 2]$ . We use independent  $N(0, 1/\beta)$  priors for components of  $\mu$ , and a  $\Gamma(0.001, 0.001)$  hyperprior on  $\beta$ .

We use Markov Chain Monte Carlo (see *Supporting Text*, which is published as supporting information on the PNAS web site) to sample from the posterior distribution of  $f$ , and estimate  $f_{ijk}$  by the mean of sampled values.

Previous work on modeling spatial variation in allele frequencies includes ref. 22. Our model differs from theirs in several ways, including (i) our parameterization (Eq. 1) treats every allele symmetrically, whereas they choose one allele against which every other allele's frequency is compared; and (ii) they use a regression to model  $\theta$ , whereas we use a Gaussian process.

**Continuous assignment method (CAM).** Our smoothing-based method can estimate allele frequencies at any location, including locations where no samples are available. We exploit this to develop our CAM, which allows samples to be assigned to any location, and not only locations with reference samples, thus overcoming a limitation of standard assignment tests (see also ref. 23). Let  $W$  denote the unknown position of origin of a sample to be located. We estimate  $W$  in a Bayesian framework, which accounts for uncertainty in estimated allele frequencies. We placed a uniform prior on  $W$  over all parts of the continent inhabited by forest elephants, or by savannah elephants, depending on whether the sample being analyzed was from a forest or savannah region. We extended the Markov Chain Monte Carlo algorithm to include  $W$  as an unknown quantity, and obtain a sample from the posterior distribution of  $W$  (*Supporting Text*). The median of the latitude and longitude of the sample give an estimate for  $W$ , and the spread gives an indication of the precision.

All assignment results were obtained by using leave-one-out crossvalidation, in which each sample in turn was treated as the sample whose location was unknown, whereas the other samples were assumed to have known location.

## Results

Previous genetic studies suggest that African elephants may be subdivided into at least two species: forest and savannah (14, 24, 25). However, observational studies have reported hybridization between the two groups on the edges of their range (26). If this were common it would impact the development of appropriate methods for tracking ivory samples. To estimate the extent of such hybridization, we computed for each sample a likelihood ratio (LR) to assess if the sample could be the offspring of one forest and one savannah elephant. Only two of our samples, both from Garamba, gave an LR favoring a hybrid origin, suggesting that in our sample hybridization is rare. However, one of these samples had overwhelming support for hybrid (central forest  $\times$  central savannah) origin (LR  $>14,000$  vs. pure central savannah; LR  $>4.5$  million vs. pure central forest), providing compelling genetic evidence for continued hybridization between the groups. Furthermore, except for these putative hybrid individuals, the forest or savannah origin of every sample was correctly determined. We therefore treated the two groups separately for subsequent analyses (e.g., for the smoothing method, allele frequencies at forest locations were estimated by using only samples from forest locations, and similarly for savannah).

To assess the potential for genetic data to infer sample origins, we used our smoothing assignment procedures to attempt to infer, from the genotype data, the sampling location of each sample (Tables 1 and 3). Forest elephants were accurately assigned to their actual forest of origin (Table 1), suggesting considerable population

**Table 1. Classification matrix for samples from forest locations by using the smoothing method**

	Estimated location								Accuracy, %
	TI	BI	MO	BM	DJ	DS	LO	GA	
TI	2	0	0	0	0	0	0	0	100
BI	0	9	0	0	0	0	0	0	100
MO	2	0	9	1	0	0	0	0	75
BM	0	0	1	10	0	0	0	0	91
DJ	2	0	0	0	6	1	0	0	67
DS	4	0	1	2	1	34	8	4	63
LO	0	0	0	0	1	1	13	0	87
GA	2	0	0	0	0	2	0	16	80
Average accuracy									83

structure, despite the fact that all West samples, and 34% of the Central forest samples, were from DNA isolated from scat, which amplified less reliably than tissue samples (0.22 vs. 0.06 missing loci on average; see Table 5). All seven Central forest samples whose location was wrongly inferred to be somewhere in West Africa had genotype data at  $<10$  loci, which may partially account for their less accurate assignment. Fortunately, DNA appears to be well preserved in ivory (10), providing higher amplification success than from scat.

Savanna elephants were more difficult to assign to their specific sampling location. However, when samples were wrongly assigned the estimated location was typically near to the actual location. Further, some locations [e.g., Etosha (ET), Mashatu (MA), and to a lesser extent Kruger (KR)] had a large proportion of samples correctly assigned, suggesting that they are genetically somewhat distinct from other locations we sampled. In the case of ET and MA, this finding is presumably because they are separated from other sampling locations by distance and habitat. In the case of KR, this discovery may be due to a huge influx of elephants from Mozambique in the 1960s (27).

To demonstrate the advantages of our smoothing method, we compare it with results (Tables 2 and 4) from a SAM (see *Materials and Methods*). (Although several alternative SAMs exist, and they would give quantitatively different results, any method that estimates allele frequencies in each population separately, using only samples from that population, will suffer the problems that we identify here due to small reference sample sizes.)

Comparing Tables 1–4, our smoothing-based method gives higher average assignment accuracy than the SAM (83% vs. 68% for forest locations and 35% vs. 30% for savanna). However, besides producing lower assignment accuracy, the SAM has two other important problems that are alleviated by smoothing. First, it has a strong systematic bias against assigning samples to locations where reference sample sizes are small. For example,

**Table 2. Classification matrix for samples from forest locations by using the SAM**

	Estimated location								Accuracy, %
	TI	BI	MO	BM	DJ	DS	LO	GA	
TI	0	0	0	0	0	2	0	0	0
BI	0	9	0	0	0	0	0	0	100
MO	0	1	9	1	0	1	0	0	75
BM	0	0	1	7	0	3	0	0	64
DJ	0	0	0	0	7	2	0	0	78
DS	0	0	0	1	1	47	3	2	87
LO	0	0	0	0	0	5	10	0	67
GA	0	0	0	0	1	4	0	15	75
Average accuracy									68

**Table 3. Classification matrix for samples from savannah locations by using the smoothing method**

	Estimated location																			Accuracy, %	
	BE	WA	MG	AB	MK	AM	TZ	SE	NG	TA	MI	RU	SW	HW	ZZ	CH	SA	MA	KR		ET
BE	9	4	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	69
WA	5	15	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	71
MG	0	0	2	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	1	50
AB	0	0	3	1	2	1	4	1	1	0	1	0	3	0	0	0	1	0	0	1	5
MK	0	1	0	0	1	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	33
AM	0	0	2	0	1	3	1	1	1	0	0	1	0	0	0	1	0	1	0	1	21
TZ	1	0	4	0	4	1	3	1	0	0	1	0	1	0	0	0	0	0	0	0	18
SE	0	1	3	2	3	0	3	0	1	1	0	0	1	0	0	0	0	2	1	0	0
NG	0	0	1	1	0	0	0	1	2	0	2	4	0	1	0	0	1	0	0	2	13
TA	0	0	1	0	1	1	0	2	1	1	2	2	1	0	0	0	0	0	0	0	8
MI	0	0	0	0	1	0	1	0	0	1	4	1	2	0	0	0	0	1	1	0	33
RU	0	0	0	0	0	0	0	0	1	0	0	1	0	0	0	0	0	0	0	0	50
SW	0	0	1	0	0	0	0	0	1	0	1	6	2	0	1	0	2	0	0	2	12
HW	0	0	1	0	0	0	0	0	0	0	1	1	0	1	1	1	2	2	1	1	0
ZZ	0	0	0	0	0	0	0	0	0	0	0	0	1	8	1	2	0	1	0	0	61
CH	0	0	0	0	0	0	0	0	0	0	0	0	0	1	2	1	6	0	1	1	8
SA	0	1	0	0	0	0	0	0	0	0	0	0	2	1	3	1	1	0	2	4	6
MA	0	1	0	0	0	0	0	0	0	0	0	1	0	0	0	0	19	0	0	0	90
KR	0	0	0	0	0	0	0	0	0	0	1	0	0	2	1	0	0	1	11	0	68
ET	0	0	0	0	0	0	0	0	0	0	0	0	1	0	1	1	0	2	11	0	68
Average accuracy																					35

in contrast to the smoothing method, the SAM assigned none of the 399 samples to the 4 sampling locations (TI, MG, MK, and RU) with <5 reference samples. Our explanation for this is that small reference samples give inaccurate allele frequency estimates at these locations, which result in artificially small likelihoods for observed genotypes. Because each sample is assigned to the location giving the highest likelihood for its genotype, the SAM tends not to assign samples to these sampling locations. To further test this premise, we compared assignment results for samples from the DS location ( $n = 54$ ), by using (i) the full reference sample at DS, and (ii) a reduced reference sample of only three individuals. When the reduced reference sample is used, the SAM assigns 0/51 of the remaining samples to DS, in striking contrast to the 47/54 that were assigned there when

using the full reference sample (Table 1). In comparison, reducing the size of the reference sample at DS has a much less dramatic effect on the smoothing method (corresponding numbers are 9/51 and 34/54).

The SAM's tendency to assign samples to locations with larger reference samples actually increases the total number of reference samples correctly assigned because, by definition, locations with larger reference samples are overrepresented in the samples being assigned. In our application, what matters in practice is not accuracy on our reference samples, but accuracy for future samples (e.g., from ivory seizures). Most likely, such samples will tend to come from locations with fewer reference samples because it may be harder to find reference samples in locations where poaching has dramatically reduced the elephant population. For this reason, the

**Table 4. Classification matrix for samples from savannah locations by using the SAM**

	Estimated location																			Accuracy, %	
	BE	WA	MG	AB	MK	AM	TZ	SE	NG	TA	MI	RU	SW	HW	ZZ	CH	SA	MA	KR		ET
BE	5	8	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	38
WA	3	18	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	85
MG	0	0	0	0	0	0	1	2	0	0	0	0	0	0	0	0	0	0	0	1	0
AB	0	0	0	5	0	0	7	1	1	0	0	0	1	0	0	0	1	1	0	0	29
MK	0	1	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0
AM	0	0	0	3	0	4	0	1	3	0	0	0	0	0	0	0	1	0	1	1	28
TZ	0	0	0	8	0	1	4	2	0	0	0	1	0	0	0	0	0	0	0	0	25
SE	0	1	0	2	0	0	5	4	2	0	0	0	1	0	0	0	0	2	1	0	22
NG	0	0	0	1	0	1	1	2	3	0	0	0	3	1	0	0	1	0	1	1	20
TA	0	0	0	1	1	2	1	2	2	1	1	0	2	0	0	0	0	0	0	0	8
MI	0	0	0	1	0	0	2	1	1	1	2	0	2	0	0	0	1	0	1	0	16
RU	0	0	0	0	0	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0
SW	0	0	0	0	0	0	2	0	1	0	1	0	8	0	0	0	2	0	0	2	50
HW	0	0	0	0	0	0	0	0	0	0	0	0	2	1	2	0	2	1	2	1	9
ZZ	0	0	0	1	0	0	1	0	0	0	0	0	0	1	6	1	2	0	1	0	46
CH	0	0	0	1	0	0	0	0	0	0	0	1	0	1	0	7	0	1	1	0	0
SA	0	1	0	0	0	1	1	0	0	0	0	3	1	1	4	0	0	1	2	0	0
MA	0	1	0	0	0	0	0	0	0	0	0	1	0	0	0	0	19	0	0	0	90
KR	0	0	0	0	0	0	0	0	1	1	0	0	0	0	0	1	0	3	10	0	62
ET	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	2	0	2	11	0	68
Average accuracy																					30

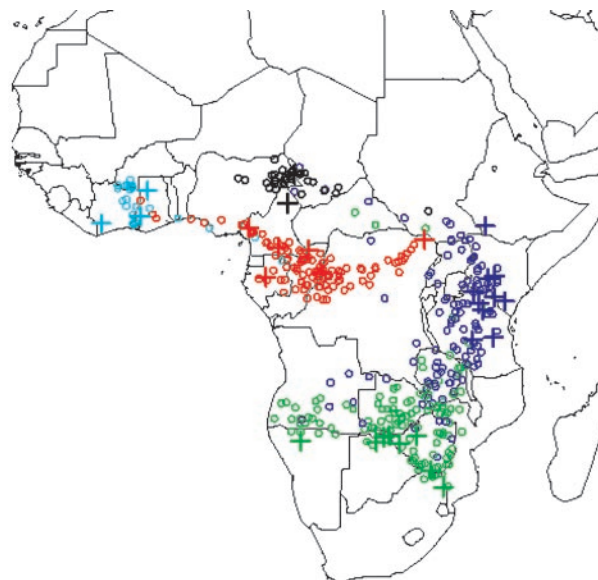
total number of reference samples correctly assigned is not a good measure of how the methods will perform on future samples. If we assume that future samples are equally likely to come from any of our sampling locations (which both assignment schemes assume implicitly by assigning samples to the location with the highest likelihood) then the average accuracy across sampling locations, quoted above and in Tables 1–4, estimates the expected accuracy of the methods on future samples.

A second problem with the SAM is that the artificially small likelihoods arising from inaccurate allele frequency estimates can cause gross overconfidence in incorrect conclusions. One measure of the confidence in an assignment to one location vs. another is the LR for the sample genotype in the two locations. If the LR for the assigned population vs. the true population is large ( $>1,000$  for instance), then we might confidently, but wrongly, conclude that the true population is not the source of the sample. This finding is important in our application because it may be helpful to exclude a location as a potential source for a tusk even if the actual source remains uncertain. Of the samples that come from sampling locations with  $<10$  reference samples, 9/29 have an LR  $>1,000$ , favoring the wrong location for the SAM. In contrast, 0/29 have an LR  $>1,000$  by using the smoothing method. Similarly, for the experiment with a reduced reference sample in DS, 26/51 samples had LRs  $>1,000$ , favoring the wrong location by the SAM, vs. 2/51 by the smoothing method.

Our smoothing assignment method alleviates some of the problems with the SAM by improving allele frequency estimates at locations with small reference samples (and, in fact, at all locations in our study; see *Supporting Text*). However, it shares perhaps the biggest drawback of the SAM for our application: the implicit assumption that each sample comes from one of our 28 sampling locations. This assumption may not hold for tusks of unknown origin. Therefore, perhaps the biggest advantage of our smoothing method for estimating allele frequencies is that it leads to our CAM that does not make this assumption, but allows that the sample may have arisen from any location in the elephant's range. Fig. 2 shows the estimated location of each sample, color-coded according to actual broad geographic region of origin: West Africa, Central forest, and Central, South, and East savannah, using the CAM. Table 5 summarizes the distances between estimated location and actual sampling location of each individual, within each region. Overall, 50% of samples were located to within 499 km of their actual origin, and 80% to within 932 km. Accuracy was greatest for samples from Central and Western regions (Table 5, see top row per region). This result might have been expected for the forest populations because they exhibit considerably greater genetic diversity than their savannah counterparts (14), and their habitat seems likely to produce greater barriers to gene flow (28). The high accuracy for central savanna elephants is presumably due to isolation from their East and South counterparts by long distances of forest and desert habitat. In contrast, the larger discrepancies for some individuals from East and Southern Africa reflect greater genetic similarity between these regions, possibly due to relatively few barriers to gene flow in these savanna habitats (29).

In practice, it will be important to know how much confidence to place in estimated sample origins. To help with this hypothesis, the CAM gives a set of plausible locations for any sample. For example, Fig. 3A shows 100 plausible locations for samples from Bia, ET, and Mikumi. For each sample, these 100 plausible locations are drawn from the set of all possible locations, weighted according to their probability. The tight clustering of the points corresponding to the Bia sample indicate high confidence that the sample comes from near Bia, the intermediate clustering for the ET sample indicate slightly less certainty, whereas the higher dispersion for the Mikumi sample indicates that we could not be confident of its precise origin.

Finally, we examined how the CAM might perform for samples from locations not included in our reference database. To perform this assessment, we applied it to each sample in turn, but ignored



**Fig. 2.** Estimated locations of elephant tissue and fecal samples from across Africa when assignments are allowed to vary anywhere within the elephants' range. All tissue and scat samples ( $n = 399$ ) successfully amplified at seven or more loci. Sampling locations are indicated by a cross and are color-coded according to actual broad geographic region of origin: West Africa, Central forest, and Central, South, and East savannah (color-coded as in the Fig. 1 legend). Assigned location of each individual sample is shown by a circle and is color-coded according to its actual region of origin. The closer each circle is to crosses of the same color, the more accurate is that individual's assignment.

all samples that were from the same location as the sample being assigned when estimating allele frequencies. For example, when attempting to estimate the origin of a sample from ET, we first excluded all other ET samples from the data set used to generate allele frequencies. As expected, this results in a decrease in average assignment accuracy (Table 5, bottom row per region). For some locations, such as ET (and Garamba, data not shown) the effects of removing the samples from that location are particularly strong because the data available at other locations do not allow accurate allele frequency estimates to be obtained for these populations. For other locations, such as Mikumi, the effects of removing the samples at that location are small because samples at nearby locations provide adequate allele frequency estimates. This finding is illustrated in Fig. 3B, which shows the effect of excluding all other samples from a location on the confidence in assignments for the samples previously shown in Fig. 3A. The confidence in the assignment of the sample from Mikumi is barely affected, in contrast to the samples from Bia and ET. The results for the sample from Bia in Fig. 3 illustrate the complex structure of forest populations: when Bia samples are included in the assignment, the sample is precisely assigned to Bia, indicating that Bia is somewhat genetically distinct from nearby Mole; when Bia is removed the precision of the assignment decreases, although the sample is still confidently assigned to West Africa. Collectively, this result indicates that the sample is more similar to the other samples from Mole than to samples from Central forest regions.

## Discussion

Some of the most pressing needs in elephant conservation include timely identification of current poaching "hot spots" where greater law enforcement is needed; monitoring impacts of international trade decisions on elephant poaching throughout the African continent; determining whether declared government stockpiles are being illegally traded and replenished; determining whether sanctioned one-time sales include nonsanc-

**Table 5. Number of kilometers within which 20%, 50%, and 80% of samples could be correctly located by the continuous assignment method**

Area	Method	20%	50%*	80%	Tissue/scat	
					Samples	Loci*
Overall	Included	212	499	932	315/84	5.1/12.5
	Excluded	404	731	1,263		
West	Included	39	135	331	0/23	A/12.8
	Excluded	392	559	1,891		
Central forest	Included	203	411	720	72/37	5.3/11.0
	Excluded	432	725	1,117		
Central savanna	Included	98	286	405	26/6	4.8/13.8
	Excluded	247	338	450		
East	Included	339	697	1,293	95/18	5.3/14.8
	Excluded	431	844	1,450		
South	Included	227	535	933	120/0	15/NA
	Excluded	459	836	1,379		

\*When all neighboring samples from that subpopulation were included versus excluded from the calculation of geographic-specific allele frequencies. Number of tissue and scat samples and mean number of loci amplified per respective sample are also indicated. NA, not applicable.

tioned tusks originating from other locations; and determining whether stockpiles of illegal ivory across Africa are being consolidated and exported. The accuracy with which genetic methods can determine the origin of DNA isolated from small amounts of ivory, tissue, or scat could greatly contribute to each of the above elephant conservation and management needs.

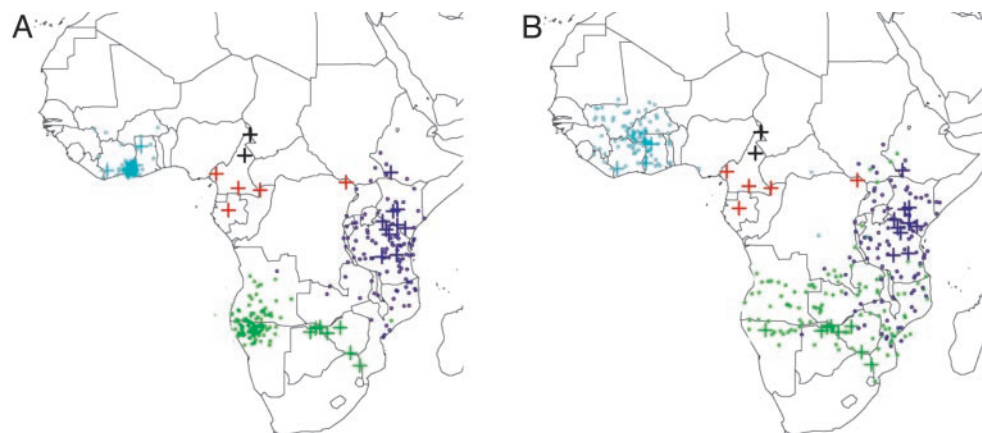
Simply discriminating between trade in forest and savanna elephant ivory should have considerable management significance in light of genetic evidence supporting the reclassification of forest and savanna elephants as separate species (14, 24–26). However, the genetic differentiation apparent among forest blocks (Table 2) suggests that we may soon be able to discern individual forests in which poaching is most heavily concentrated. These data may also assist authorities in tracking the bush meat trade, which increasingly includes elephant meat (8). Subanalyses of smaller tusks may be particularly useful here; unlike the ivory trade, bush meat poaching tends to occur independent of tusk size. Determining the proportion of African vs. Asian elephant ivory being sold in markets throughout the world is another important, and feasible use of these techniques (14).

A tremendous amount of information could be obtained simply by genetically monitoring the geographic composition of seizures as well as that of ivory moving through the major world ivory markets over time (e.g., Africa: Abidjan, Lagos, Dakar, Cairo, and Addis

Ababa; Asia: Bangkok, Guangzhou, Tokyo, and Osaka; E. Martin, personal communication; see also ref. 5). Such information could greatly assist CITES by complementing efforts of the Monitoring Illegal Killing of Elephants (MIKE) (30) and Elephant Trade Inventory System (ETIS) (31) groups and the at a substantially reduced cost. The estimated cost of MIKE alone is \$13.8 million over 6 years (32). These costs will only increase as populations continue to decline and elephants become increasingly difficult to detect. By contrast, DNA analyses are relatively inexpensive (currently approximately \$100 per tusk, plus labor), likely to decrease with time, and reference sites do not require resampling over time.

Our methods make a number of modeling assumptions, some of which may hold only approximately in practice. However, the results in Figs. 2 and 3 and Tables 1–5 are based on the empirical performance of the method on samples of known origin, and are therefore valid independent of modeling assumptions, and indeed of other factors such as genotyping error and nonamplification of alleles.

The fact that accuracy is considerably greater for samples from locations in our reference database (Table 5) demonstrates the crucial importance of obtaining samples from regions not currently represented in our data. Whereas the considerable structure apparent in our results, particularly in forest regions (Table 2), holds great promise for genetic methods to determine sample location



**Fig. 3.** Representation of confidence of assignments when all neighboring samples from that subpopulation (represented by a single cross) were included (A) vs. excluded (B) from the calculation of geographic-specific allele frequencies for ET (green), Bia (turquoise), and Mikumi (blue). The 100 color-coded circles are random draws from the set of all possible locations weighted according to their probability. The concentration of these 100 circles in any given area gives a guide to the probability that the sample arose from that area under each condition.

rather precisely, denser sampling will be required for this information to be fully exploited. Use of noninvasive means of acquiring DNA from feces should greatly facilitate these efforts, and the spatial smoothing used by our CAM means that relatively small numbers of samples from each location should suffice. This idea is particularly important when using samples such as scat where amplification success is uneven.

In conclusion, considerable breakthroughs in wildlife conservation and management should result from the methods we describe to track the geographic origin of poached ivory. Nearly all of the applications suggested above can be addressed with our existing data, although a few may require more reference samples from strategic areas to realize their full potential. The ability to acquire DNA from scat, along with tools to facilitate its collection over large remote areas (9), can also make such methods achievable in a timely manner for other species throughout the world.

The methods described here are implemented by using the software package SCAT (Smoothed and Continuous Assignment Tests), which can be accessed at [www.stat.washington.edu/stephens/software.html](http://www.stat.washington.edu/stephens/software.html).

1. Dublin, H., Milliken, T. & Barnes, R. F. W. (1995) *Four Years After the CITES Ban: Illegal Killings of Elephants, Ivory Trade and Stockpiles* (International Union for Conservation of Nature and Natural Resources, Gland, Switzerland).
2. Robinson, J. G. & Bennett, E. L. (1999) *Hunting for Sustainability in Tropical Forests* (Columbia Univ. Press, New York).
3. Baker, C. S., Cipriano, F. & Paulmbi, S. R. (1996) *Mol. Ecol.* **5**, 671–685.
4. Douglas-Hamilton, I. (1987) *Pachyderm* **8**, 1–10.
5. Said, M. Y. (1995) *African Elephant Database* (International Union for Conservation of Nature and Natural Resources, Gland, Switzerland).
6. Martin, E. & Stiles, D. (2000) *The Ivory Markets of Africa* (Save the Elephants, Nairobi).
7. Douglas-Hamilton, I. (1988) *African Elephant Database Project: Phase Two, December* (●●●, Nairobi).
8. Courable, M., Hurst, F. & Milliken, T. (2003) *More Ivory Than Elephants: Domestic Ivory Markets in Three West African Countries* (Traffic International, Cambridge, U.K.).
9. Wasser, S. K., Houston, C. S., Koehler, G. M., Cadd, G. G. & Fain, S. R. (1997) *Mol. Ecol.* **6**, 1091–1097.
10. Wasser, S. K., Davenport, B., Ramage, E. R., Hunt, K. E., Parker, M., Clarke, C. & Stenhouse, G. (2004) *Can. J. Zool.* **82**, 475–492.
11. Comstock, K. E., Ostrander, E. A. & Wasser, S. K. (2003) *Cons. Biol.* **17**, 1–4.
12. Paetkau, D., Calvert, W., Sterling, I. & Strobeck, C. (1995) *Mol. Ecol.* **4**, 347–354.
13. Comstock, K. E., Wasser, S. K. & Ostrander, E. A. (2000) *Mol. Ecol.* **9**, 1004–1006.
14. Comstock, K. E., Georgiadis, N., Pecon-Slattery, J., Roca, A. L., Ostrander, E. A., O'Brien, S. J. & Wasser, S. K. (2002) *Mol. Ecol.* **11**, 2489–2498.
15. Nyakaana, S. & Arcander, P. (1998) *Mol. Ecol.* **7**, 1436–1437.
16. Barnes, R. F. W., Craig, G. C., Dublin, H. T., Overton, G., Simons, W. & Thouless, C. R. (1998) *African Elephant Database* (International Union for Conservation of Nature and Natural Resources/Species Survival Commission African Elephant Specialist Group, International Union for Conservation of Nature and Natural Resources, Gland, Switzerland and Cambridge, U.K.).
17. Karesh, W., Smith, F. & Frazier-Taylor, H. (1989) *Cons. Biol.* **1**, 261–262.
18. Taberlet, P., Griffin, S., Goossens, B., Questiau, S., Manceau, V., Escaravage, N., Waits, L. P. & Bouvet, J. (1996) *Nucleic Acids Res.* **24**, 3189–3194.
19. Fernando, P., Vidy, T. N., Rajapakse, C., Dangolla, A. & Melnick, D. J. (2003) *J. Hered.* **94**, 115–123.
20. Anderson, E. C. & Thompson, E. A. (2000) *Genetics* **16**, 1217–1229.
21. Diggle, P. J., Tawn, J. A. & Moyeed, R. A. (1998) *J. R. Stat. Soc. A* **47**, 299.
22. Vounatsou, P., Smith, T. & Gelfand, A. E. (2000) *Biostatistics* **1**, 177–189.
23. Cornuet, J. M., Piry, S., Luikart, G., Estoup, A. & Solignac, M. (1999) *Genetics* **153**, 1989–2000.
24. Roca, A. L., Georgiadis, N., Pecon-Slattery, J. & O'Brien, S. J. (2001) *Science* **293**, 1473–1475.
25. Barriel, V., Thuet, E. & Tassy, P. (1999) *Evolution (Lawrence, Kans.)* **322**, 447–454.
26. Grubb, P., Groves, C. P., Dudley, J. P. & Shoshani, P. (2000) *Elephants* **2**, 1–4.
27. Brack, L. (1997) *Management Plan for the Kruger National Park: Policy Proposals Regarding Issues Relating to Biodiversity Maintenance, Maintenance of Wilderness Qualities, and Provision of Human Benefits*, (South African National Parks, ●●●), Vol. III.
28. Eggert, A., Rasner, C. A. & Woodruff, D. S. (2002) *Proc. R. Soc. London Ser. B* **269**, 1993–2006.
29. Georgiadis, N., Bischof, L., Templeton, A., Patton, J., Karesh, W. & Western, D. (1994) *J. Hered.* **85**, 100–104.
30. Traffic Report (1998) *Monitoring of Illegal Killing of Elephants (MIKE)* (Traffic International, Cambridge, U.K.).
31. Traffic Report (2000) *Elephant Trade Inventory System (ETIS)* (Traffic International, Cambridge, U.K.).
32. ●●●. (●●●) *CITES Document 11.31.2. Monitoring of Illegal Trade and Illegal Killing* (●●●).

We thank two anonymous reviewers for helpful comments; Mark Handcock for discussions on geostatistical approaches to spatial smoothing; L. Andre, S. Blake, L. Eggert, N. Georgeatis, D. Ishengoma, R. Ruggiero, and D. Woodruff for providing tissue and fecal samples for this study; and A. Turkalo, J. M. Fay, R. Weladji, W. Karesh, M. Lindeque, W. Versvelt, K. Hillman Smith, F. Smith, M. Tchamba, S. Gartlan, P. Aarhaug, A. M. Austmyr, Bakari, Jibrila, J. Pelletier, L. White, M. Habibou, M. W. Beskreo, D. Pierre, C. Tutin, M. Fernandez, R. Barnes, B. Powell, G. Doungoubé, M. Storey, M. Phillips, B. Mwasaga, and A. Mackanga-Missanzou for assistance in tissue collections. We also thank the governments of Botswana, Cameroon, Central African Republic, Congo (Brazzaville), Democratic Republic of Congo, Ethiopia, Gabon, Ghana, Ivory Coast, Kenya, Namibia, South Africa, Tanzania, and Zimbabwe for permission to collect samples. This work was supported by U.S. Fish and Wildlife Service Grants 1448–98210-98-G145 and 98210-G794 (to S.K.W.), the International Elephant Foundation and the Woodland Park Zoo, National Institutes of Health Training Grant T32 HG0035 (to K.C.), the Department of Molecular Biotechnology at the University of Washington, and the Fred Hutchinson Cancer Research Center. Tissue sampling was supported by the National Geographic Society, the European Union (through the Wildlife Conservation Society), the National Science Foundation, and the U. S. Fish and Wildlife Service (to N.G.).