

# Genome-wide assessment of worldwide chicken SNP genetic diversity indicates significant absence of rare alleles in commercial breeds

William M. Muir<sup>a</sup>, Gane Ka-Shu Wong<sup>b,c</sup>, Yong Zhang<sup>c</sup>, Jun Wang<sup>c</sup>, Martien A.M. Groenen<sup>d</sup>, Richard P.M.A. Crooijmans<sup>d</sup>, Hendrik-Jan Megens<sup>d</sup>, Huanmin Zhang<sup>e</sup>, Ron Okimoto<sup>f</sup>, Addie Vereijken<sup>g</sup>, Annemieke Jungerius<sup>g</sup>, Gerard A.A. Albers<sup>g</sup>, Cindy Taylor Lawley<sup>h</sup>, Mary E. Delany<sup>i</sup>, Sean MacEachern<sup>e</sup>, and Hans H. Cheng<sup>e,1</sup>

<sup>a</sup>Department of Animal Sciences, Purdue University, West Lafayette, IN 47907; <sup>b</sup>University of Alberta, Department of Biological Sciences and Department of Medicine, Edmonton, Alberta T6G 2E9, Canada; <sup>c</sup>Beijing Institute of Genomics of Chinese Academy of Sciences, Beijing Genomics Institute, Beijing Proteomics Institute, Beijing 101300, China; <sup>d</sup>Animal Breeding and Genomics Centre, Wageningen University, 6709 PG Wageningen, The Netherlands; <sup>e</sup>United States Department of Agriculture, Agricultural Research Service, Avian Disease and Oncology Laboratory, East Lansing, MI 48823; <sup>f</sup>Cobb-Vantress, Inc., Siloam Springs, AR 72761; <sup>g</sup>Hendrix Genetics, 5831 CK Boxmeer, The Netherlands; <sup>h</sup>Illumina, Inc., San Diego, CA 92121; and <sup>i</sup>Department of Animal Science, University of California Davis, Davis, CA 95616.

Edited by James E. Womack, Texas A&M University, College Station, TX, and approved September 23, 2008 (received for review July 8, 2008)

**Breed utilization, genetic improvement, and industry consolidation are predicted to have major impacts on the genetic composition of commercial chickens. Consequently, the question arises as to whether sufficient genetic diversity remains within industry stocks to address future needs. With the chicken genome sequence and more than 2.8 million single-nucleotide polymorphisms (SNPs), it is now possible to address biodiversity using a previously unattainable metric: missing alleles. To achieve this assessment, 2551 informative SNPs were genotyped on 2580 individuals, including 1440 commercial birds. The proportion of alleles lacking in commercial populations was assessed by (1) estimating the global SNP allele frequency distribution from a hypothetical ancestral population as a reference, then determining the portion of the distribution lost, and then (2) determining the relationship between allele loss and the inbreeding coefficient. The results indicate that 50% or more of the genetic diversity in ancestral breeds is absent in commercial pure lines. The missing genetic diversity resulted from the limited number of incorporated breeds. As such, hypothetically combining stocks within a company could recover only preexisting within-breed variability, but not more rare ancestral alleles. We establish that SNP weights act as sentinels of biodiversity and provide an objective assessment of the strains that are most valuable for preserving genetic diversity. This is the first experimental analysis investigating the extant genetic diversity of virtually an entire agricultural commodity. The methods presented are the first to characterize biodiversity in terms of allelic diversity and to objectively link rate of allele loss with the inbreeding coefficient.**

alleles | biodiversity | poultry

**G**lobal production of chickens has experienced massive change and growth over the past 50 years. The commercial broiler and layer markets produce more than 40 billion birds annually to meet current worldwide consumer demands of more than 61 metric tons of meat and more than 55 million metric tons of eggs. In fact, poultry has become the leading meat consumed in the United States and most other countries and is the most dynamic animal commodity in the world; production has increased by 436% since 1970, more than 2.3 times and 7.5 times the corresponding growth in swine and beef, respectively (1). Selection for specific traits by poultry breeders was the key factor in the steep rise in productivity, accounting for up to 90% of the increase (2). For the industry to remain successful, sufficient genetic diversity must exist within companies, because (unlike in crop agriculture) introgression from noncommercial birds is rarely used.

The goal of this research was to determine the extent to which noncommercial and ancestral populations might contain potentially useful germplasm not found in commercial populations. Initially, in North America and Europe, chickens of numerous standard breeds

(e.g., Rhode Island Red, Single-Comb White Leghorn) were raised in small backyard flocks primarily for the production of eggs and meat as food, with others developed as game birds for sport and still others developed as fancy breeds for show. Beginning in the 1950s, modern poultry production emerged, with specialized industrial chicken breeds selected intensively for either meat-type (broiler) or egg-type (layer) chickens. All commercial white egg chicken lines are based in the White Leghorn breed, whereas brown egg chicken lines were initially selected from North American dual-purpose breeds (selected for both meat and egg qualities), such as Rhode Island Red and White Plymouth Rock, which originated from crosses between Asian and European breeds. Due to the negative genetic correlation between production (growth) and reproduction (egg number) (3), commercial poultry meat production uses crosses among specialized broiler lines. Lines selected primarily for growth traits are referred to as sire or male lines, because only males are used in the final commercial cross. The lines used for the female side of the cross are selected for both reproductive and growth traits and are referred to as dam or female lines. The male lines are derived from Cornish stock, originating from the British Cornish Indian Game breed, having a thick compact body type with a high proportion of breast muscle. The dam lines originate from many of the same dual-purpose breeds used for brown egg production (e.g., Barred Plymouth Rock, White Plymouth Rock, New Hampshire). Thus, the first tier of genetic diversity reduction was due to limited breed utilization.

The second tier of genetic diversity reduction is ongoing and due to breeding structure and within-line selection. The industry is structured such that the final commercial product is the result of intense within-line selection, followed by a pyramid expansion scheme. This scheme is designed so that the top or pure line level, a limited number of individuals are measured for critical production traits, because the collection of phenotypes is expensive and time-consuming. Genetic improvement based on these traits is performed within a line and is then multiplied by crossing with other

Author contributions: W.M.M., G.K.-S.W., M.G., and H.H.C. designed research; W.M.M., M.G., R.P.C., H.-J.M., A.V., A.J., C.T.L., and H.H.C. performed research; W.M.M., M.G., R.P.C., H.-J.M., H.Z., R.O., A.V., A.J., G.A.A., C.T.L., and H.H.C. contributed new reagents/analytic tools; W.M.M., G.K.-S.W., Y.Z., J.W., M.G., H.-J.M., M.E.D., S.M., and H.H.C. analyzed data; and W.M.M., G.K.-S.W., M.G., H.-J.M., H.Z., R.O., A.J., G.A.A., M.E.D., S.M., and H.H.C. wrote the paper.

The authors declare no conflict of interest.

This article is a PNAS Direct Submission.

<sup>1</sup>To whom correspondence should be addressed at: USDA, ARS, Avian Disease and Oncology Laboratory, 3606 E. Mount Hope Rd. East Lansing, MI 48823. E-mail: hans.cheng@ars.usda.gov.

This article contains supporting information online at [www.pnas.org/cgi/content/full/0806569105/DCSupplemental](http://www.pnas.org/cgi/content/full/0806569105/DCSupplemental).

© 2008 by The National Academy of Sciences of the USA

selected lines, for same or different traits, for three or four generations. At each generation, the number of offspring from a single bird can exceed 200. As a result, superior genetics of a single primary layer or broiler can be expanded more than a million-fold to produce end products of meat or eggs. Because these pure lines have dramatically different agronomic traits than noncommercial standard breeds, gene flow does not occur between commercial and noncommercial poultry, resulting in essentially closed breeding structures. Thus, inbreeding reduces genetic diversity within the pure lines, although poultry breeders work to avoid inbreeding to the greatest extent possible within closed populations.

Because inbreeding converts within-line genetic variability into between-line variability (4, 5), and because all commercial companies have many pure lines, regardless of the within-line inbreeding, multiple independent lines help preserve alleles within a company. But intense competition within the industry in recent decades has left only a few multinational companies remaining as genetic suppliers of the majority of commercial birds (6). Thus, this final tier limits preservation of alleles between lines.

Because of these multitiered diversity-reducing mechanisms, there is a realistic concern that genetic diversity for future needs may be compromised. Inadequate genetic diversity has had severe negative consequences in both plant and animal species. Oft-cited examples include the 1970 corn leaf blight outbreak due to the widespread use of the Texas male-sterile cytoplasm (7) and the high prevalence of bovine leukocyte adhesion deficiency (BLAD, an autosomal recessive hereditary disease) in Holstein cattle due to the carrier status of several prominent bulls used for artificial insemination (8).

To achieve our objectives, we used the recent chicken genome sequence (9), the identification of more than 2.8 million single-nucleotide polymorphisms (SNPs) (10), and the ability to perform high-throughput genotyping to evaluate the existing genetic diversity in commercial pure lines. Using analytical methods that account for inbreeding and SNP ascertainment bias, we found that commercial poultry breeds have considerably less allelic diversity compared with noncommercial breeds, due primarily to the first tier of narrowing genetic diversity, that is, the limited number of chicken breeds that went into the formation of modern commercial lines. A possible strategy for preserving and accessing more genetic diversity is discussed.

## Results

**SNP Verification and Genotyping Performance.** All but 14 of the 2580 DNA samples collected from commercial pure lines, experimental chickens, and standard breeds were genotyped successfully (0.54% sample failure rate). Of the 3072 SNPs spaced evenly throughout the chicken genome and examined [see [supporting information \(SI\) Table S1](#)], 2733 provided results, for a success rate (89.0%) that is within the expected 5%–10% loss range because of multiplex amplification issues. The reproducibility rate was 99.996% based on plate and other controls. A comparison of the allele calls with the control DNAs (those used in the actual SNP discovery process) indicated that 2428 of the 2706 SNPs (89.7%) were in full agreement. A minor allele frequency (MAF) of  $\geq 2\%$  was observed for 2416 of the 2733 working SNPs (88.4%); 182 SNPs were monomorphic (6.7%), leaving 2551 SNPs segregating in this collection. No significant difference in allele frequency distributions were observed between the tolerant coding nonsynonymous SNP (cnSNPs) and all of the remaining SNPs.

**Reconstructing Allele Frequencies for the Hypothetical Ancestral Population (HAP).** Results of the unweighted pair-group method using arithmetic averages (UPGMA) clustering of samples are given in [Table S2](#). The effect of number of clusters on resulting allele frequency distribution is shown in [Fig. S1](#). Level N ([Table S2](#)) was our *a priori* clustering distance based on known relationships of broiler lines. For one level below N and two levels above N,

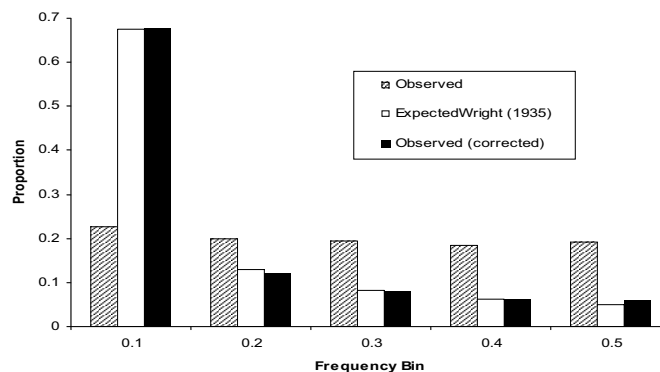


Fig. 1. Observed, corrected, and expected allele frequency distributions.

clustering level had little effect on allele frequency distribution; however, at the highest level (Q), the distribution was severely skewed to the right. For a SNP discovery depth of 2, an approximate uniform distribution would be expected (11), thus, level Q clearly is incorrect. For all levels below Q examined, the distribution was approximately uniform, but with a slight skew toward more alleles in the lowest frequency bin.

### Distribution of Allele Frequencies and Ascertainment Bias Correction.

Allele frequency distributions for the observed and after ascertainment bias correction are shown in [Fig. 1](#). Because the ancestral state of the alleles was not known, the distribution was folded based on MAF. When corrected for ascertainment, a folded U-shaped distribution resulted, which, when fit to Wright's distribution (12):  $\phi(q) = 4Nvq^{4Nv-1}$ , was nearly exact for a parameter estimate of  $4Nv = 0.184$ . A number of tests use estimates of Wright's distribution and variations thereof to infer divergence from the neutral model as a method of detecting positive selection (13–15). This is the first time that an estimate of this parameter was done with a data set of sufficient size to approximate the distribution in economically important chickens. These data tend to support the neutral model even for animals that have been highly selected. Two possible explanations for this are that the proportion of the genome actually under selection may be quite small, or that the SNP loci used were in linkage equilibrium with quantitative trait loci under selection. The possibility that SNP were in linkage equilibrium with quantitative trait loci is supported by results showing linkage disequilibrium in these populations can extend to  $<0.1$  cM, whereas the SNPs in our study were spaced  $\approx 1$  cM apart (16, 17).

**Inbreeding.** Estimation of the allele frequencies for the sampled loci were based on the HAP. For a correct estimate, sufficient diversity must be sampled to be representative of at least two alternative independent lineages. The accuracy of the estimate (bias) is not dependent on the number of lineages sampled, because the expected average allele frequency over independent lineages is the allele frequency in the HAP; however, the precision (variance) of the estimate improves with the number of independent lineages sampled. The estimate is dependent on correctly separating those samples into representative strata (lineages).

Estimation of inbreeding is dependent on several factors, including the formula used, due to potential ascertainment bias. This issue of estimation was resolved by using three approaches for finding  $F_{IT}$ : (1) per individual based on the reduction in total heterozygosity across loci, then averaged across individuals; (2) per locus based on the reduction in heterozygosity, then averaged across loci; and (3) per locus based on the reduction in variance. The regression of  $F_{IT}$  estimated from methods 2 and 3 ([Fig. S2](#)) resulted in good agreement, with an  $R^2$  of 98% and a slope of  $1.04 \pm .04$ , which is not significantly different from 1, as expected. In addition, the

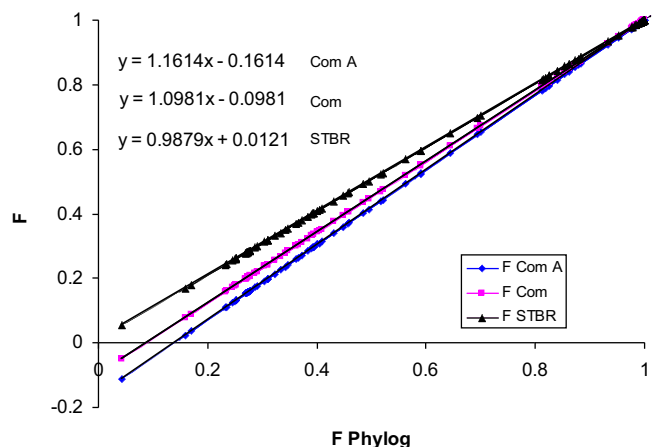


Fig. 2. Estimation of  $F$  with different subpopulations contributing to the HAP.

regression of  $F_{IT}$  estimated from method 1 on method 2 (Fig. S3) resulted in a slope of  $1.02 \pm .05$ , which also is not significantly different from 1. Because all three methods are in good agreement, method 1 was used.

The impact of alternative clustering and known insufficient sampling on estimates of inbreeding was examined by reconstituting the HAP with various subsamples and using an alternative clustering method. The first subsample was based only on the standard breed populations (STBR), the second subsample was based only on industry pure lines (COM), and the third subsample was based on all sampled pure lines within only a single commercial broiler company (COM A). These inbreeding estimates were regressed on lineages defined by UPGMA clustered at level N; the results are shown in Fig. 2. These regressions show that, as expected, bias resulted when sampling was not representative of the HAP. This bias increased as the samples deviated more greatly from a representative sample of the true HAP. The use of just the STBR lines as lineages resulted in nearly identical estimates as those from the UPGMA-reconstituted HAP, with a difference of  $<2\%$ . When only commercial lines (broiler and layer) were used to define the HAP, then the resulting bias was almost 10%. Finally, when a single company attempted to estimate inbreeding using

these methods by combining all lines within the company, then the bias exceeded 16%.

In all cases, the bias was downward; that is, the amount of inbreeding would have been underestimated. Thus, by extrapolation, we conclude that if our samples are not representative of the true HAP, then our estimates will be biased downward as well, resulting in a conservative estimate of the true level of inbreeding. Our analysis could overestimate the level of inbreeding if our samples have greater genetic variability compared with the HAP; but because genetic sampling (inbreeding) always reduces genetic variability, it is not likely that our samples represent greater variability compared with the HAP.

For comparison, other methods of defining strata were used, including principal component analysis (PCA), as described by Price *et al.* (18), who showed that PCA can be used to correct for stratification and neighbor-joining clustering (see Figs. S4–S6). Clusters also were confirmed by bootstrapping an UPGMA tree using 5000 replicates (see Fig. S7). Bootstrapping values of 90% were set at as the cutoff. Comparing clusters based on bootstrapping cutoff values and those determined from *a priori* knowledge shows that they are very similar and suggests that their separation is strongly supported by the data. In addition, regression of inbreeding estimates based on UPGMA and PCA were within 3%. As such, all methods gave very similar results, indicating that the allele frequency estimates in the HAP are somewhat robust to the clustering method and, by extension, the estimated level of inbreeding within subpopulations. Estimates of inbreeding for each line are given in Table S3.

**Proportion of Missing Alleles.** The proportions of alleles missing ( $\Omega$ ) estimated using SNP weights (SNP.WTs; see Table S4) are given in Table S5. These results demonstrate at least 20% inbreeding for each line on average, with a corresponding 60% reduction in allelic diversity. The relationship between  $\Omega$  and inbreeding is shown in Fig. 3. The regression shows a linear relationship between the proportion of missing alleles and  $F$ , with alleles missing at a proportional rate of 50% per unit increase in  $F$ , but with an intercept of 50% loss. This was surprising, because when  $F = 0$ , then clearly  $\Omega = 0$ , indicating the existence of an extreme nonlinear relationship between  $F = 0$  and 0.2. This relationship was examined using simulations with an initial distribution of allele frequencies based on Wright's equation (12) with  $4Nv = 0.184$ , as estimated in the Results section on distributions. Results from these simulations

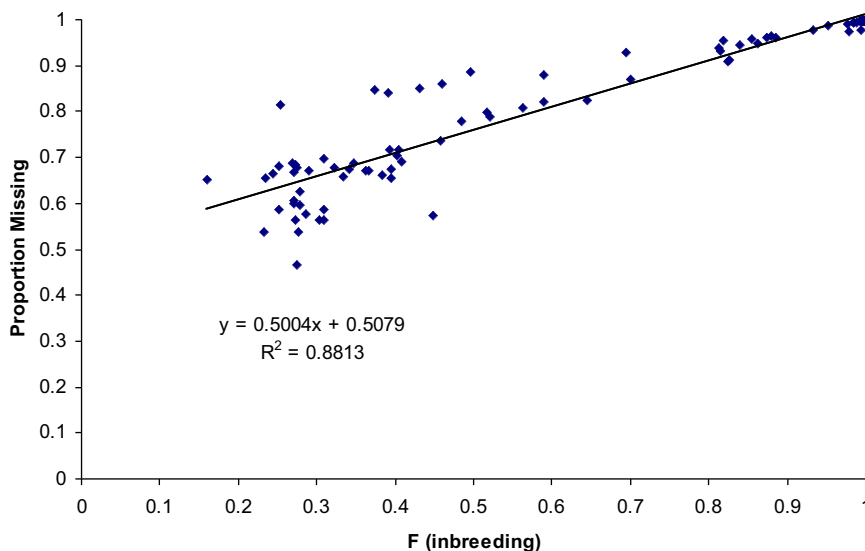


Fig. 3. Empirical relationship between inbreeding and missing alleles ( $\Omega$ ).



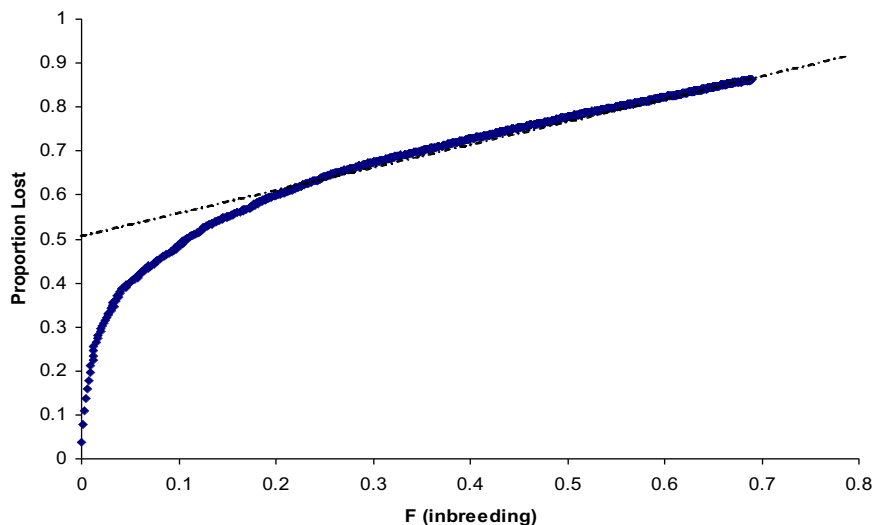


Fig. 4. Simulation results showing lost alleles ( $\Omega$ ) with inbreeding. The dotted line gives the slope and intercept of the linear portion of the curve.

(see Fig. 4) clearly show the nonlinearity with low levels of inbreeding. If only data with  $F > 0.2$  were used, then the same relationship was found as with the observed data, that is, an intercept of 0.5 and a slope of 0.5, as indicated by the dotted line. The simulation results lend validity to using the SNP-WT method as a metric to quantify missing allelic diversity in populations. Equally important is the conclusion that major effects on allelic diversity are incurred by relatively minor amounts of inbreeding, followed by a loss that is linear with inbreeding.

These results make sense when viewed as a departure from equilibrium. The HAP is assumed to be in a dynamic equilibrium, in which the rate of loss of alleles due to inbreeding is balanced by new mutations. This means that the vast majority of allelic diversity in the HAP is rare. Any change in mutation rate or effective population size will alter that equilibrium. If we assume that mutation rates are approximately constant per generation but effective population size fluctuates, and, in particular, if the effective population size is reduced, as is the case for most domestic species, then rare alleles are lost preferentially. This result is verified in Fig. S8, which shows for a representative line that rare alleles are eliminated first. The effect of missing rare alleles as it relates to addressing future commercial poultry needs is unknown; however, rare alleles have been relevant for some production traits in other livestock species (19–22), and reduction of genetic diversity is not favorable for identifying genetic resistance factors to new or emerging infectious diseases.

Assuming that phenotypic variation in traits are due to single base changes in or around functional genes (quantitative trait nucleotide [QTN]), and that the ability to respond to future challenges is reflected by single base changes that have not realized their full evolutionary potential for a given environment, then the rate of frequency reduction of neutral SNPs should be reflective of the rate of frequency reduction of QTN alleles, provided that the QTN alleles are neutral. Thus, the evolutionary potential of a commercial population can be inferred by the absence of random SNPs, provided that their effects on fitness are neutral. Therefore, information on missing random SNPs can be used as an indicator of missing neutral functional alleles in the genome. But alleles that are not neutral may behave much differently, depending on the strength and direction of selection. Besides validating the allelic frequency reduction results, this method has the added functional attribute of being applicable to any poultry population using these SNPs, because the weights are now known.

**Recovery of Genetic Diversity.** To explore whether genetic diversity could be reconstituted within existing commercial pure lines, *in silico* groups were generated by combining all lines within a company, across a breed category, and as a single all-encompassing commercial group. Expected heterozygosity ( $H_s$ ) of the combined lines was calculated based on average allele frequencies across those lines. This value was then used to compute a biased estimate of the inbreeding coefficient ( $F_{st}$ ); that is, targeted genotyping of SNP discovered from a small sample may overestimate the expected heterozygosity ( $H_s$ ), resulting in an underestimate of ( $F_{st}$ ) (23). Therefore, downward-biased estimates of  $F_{st}$  are given in Table S6 for combinations of tested lines within a company and across the industry. This result suggests that combining all lines across all companies would result in a population with a conservative estimate of  $\approx 10\%$  inbreeding coefficient. This equates to missing at least 50% of the alleles present in the HAP (Fig. 4). Combining all commercial lines *in silico* is the same as creating a HAP based only on the commercial lines. In this case, the commercial-HAP (C-HAP) is a subset of the predomestication HAP and represents the inbreeding that occurred in the first tier of narrowing genetic diversity, that is, the few breeds that contributed to modern poultry breeding programs. These results suggest that among domesticated lines, the larger reservoir of allelic genetic diversity will be found outside breeds contributing to commercial poultry, that is, STBR.

### Discussion

The results of the two analyses, which used different approaches, indicate that commercial pure lines of chicken, both broiler (meat) and layer (egg) lines, are missing significant genetic diversity found in noncommercial chickens. We explored possible strategies for companies to restore genetic diversity within lines by crossing multiple pure lines. Crossing combines the diversity preserved among lines, thereby restoring some or all of the within-line variability depending on the number of lines maintained and crossed. In addition, it is possible that industry consolidation will continue, meaning that gene flow could occur across companies in the future. However, as shown in Table S6, such *in silico* crosses, if done across the entire poultry industry, could reduce the inbreeding coefficient to 10%, but this reduction does not translate into a large recovery of missing alleles. The minimum missing alleles were determined by interpolating the estimated inbreeding coefficients from Table S6 onto the allele loss from Fig. 4. This method provides a conservative estimate of missing alleles, because the *in silico* estimates of inbreeding are biased downward by the ascertainment bias. As

shown in Fig. 4, an inbreeding coefficient as low as 0.10 results in an allele loss of almost 50% from a population experiencing inbreeding. Thus, even in the unlikely and hypothetical situation in which all commercial birds were combined into a single population, a limited increase in allelic diversity would result; that is, a large proportion of genetic diversity would not be present in early lines used for the formation of commercial breeds. An independent assessment of 65 diverse chicken populations showed that commercial birds form their own clusters with very low admixtures with other clusters (24). These findings indicate that the poultry industry, across both the egg and meat pure line stocks, has a narrowed genetic reservoir and possibly a reduced capacity to respond to future industry needs.

Interestingly, the question arises as to whether modern agricultural practices further contribute to this diversity reduction. Although some pure lines are highly inbred, others show very moderate levels. Crosses among lines within a company would result in an inbreeding level between 14% and 15%, as opposed to crossing of all lines across all breeds, which would result in a 10% level of inbreeding. Thus, on average, modern agriculture has contributed less than 5% to the level of inbreeding despite intense levels of selection, closed populations, and industry consolidation. It is worthwhile to note that these findings do not preclude future genetic progress, especially given the results of long-term selection studies in maize, which show continued phenotypic response after 100 generations of intense selection (25). Therefore, new mutations may provide needed genetic variability and contribute to a lack of a perceived “selection wall” for growth and reproduction traits (26). But our findings do raise concerns about traits attributed only to rare alleles, such as resistance to certain infectious diseases, which may be missing in commercial poultry. Under these conditions, there may be no easy way for the industry to access the relevant genetic diversity other than by introgression (slow) or direct genetic manipulation (controversial). Certainly, as a source for rare alleles, our findings reemphasize the need for support and planning for ongoing, new, or novel efforts to maintain genetic diversity using noncommercial and native poultry populations. Future food production challenges are unpredictable and likely will include new diseases or more virulent recurring diseases, environmental changes, changes to animal welfare and consumer preferences, as well as expansion of poultry-related nutritional demands from a global society, necessitating alternatives. Therefore, a healthy genetic reservoir in food-producing animals remains as crucial as ever. Indeed, noncommercial flocks, including those found in many underdeveloped and developing countries, potentially represent the reservoir opportunity for alleles “missing” from commercial pure line stocks.

## Materials and Methods

**Chickens.** To survey the extant biodiversity of commercial poultry, an extensive collection of DNA from commercial pure lines was assembled. Four major breeding companies (three broiler breeders and one layer breeder), which together account for ≈90% of meat-type and ≈40% of egg-type chickens supplied commercially worldwide, each provided material from 40 selected birds in each of 9 pure lines. Furthermore, to establish a baseline for diversity, additional DNA was collected from a Red Jungle Fowl line (the progenitor of domestic chickens (27)), standard breeds, and experimental lines derived from commercial and standard breeds, which yielded a total of 2580 unique individuals; 1440 commercial birds (representing male and female broilers, white and brown egg layer pure lines), 1136 experimental and standard breed chickens, and 4 controls (UCD 001 #256 Red Jungle Fowl, the sequenced bird; Chinese Silkie, commercial broiler, and experimental White Leghorn, the actual birds used in the SNP discovery process). [Table S7](#) shows how the lines were grouped and coded.

**SNP Selection and Genotyping.** To obtain SNPs evenly spaced throughout the chicken genome, the genome sequence (WASHUC1) was divided into 3072 bins, taking into account the recombination rate per chromosome. For each bin, three SNPs from the 2.8 million SNP data set identified previously by Wong *et al.* (10) were selected (see [Table S1](#)). Preference was given to high-confidence SNPs in genes, especially those judged to be tolerant cSNPs, which accounted for 1124 assays. All SNPs were evaluated for assay suitability, and a single suitable SNP was

selected from each bin. In addition, 34 SNPs in genes of interest were evaluated. The DNAs were genotyped at Illumina.

**Reconstructing Allele Frequencies for the HAP.** To determine inbreeding, loss of heterozygosity, and proportion of alleles missing, it was first necessary to reconstruct a HAP (28) as a reference. Neutral drift theory posits that if an ancestral population is divided into a number of independent subpopulations, then the average allele frequency across a random sample of such subpopulations will remain unchanged, and it provides an unbiased estimate of allele frequencies in the original base population (12). Because our samples (lines) are not independent, it was important to recombine these samples in a manner consistent with the subdivision structure in which they arose. Because the lineages and relationships between lineages for our samples were relatively unknown, it was necessary to reconstruct this information from the data. This reconstruction was performed using cluster analysis. For clustering of lines (samples) into lineages, genetic distances between samples  $i$  and  $i'$  were computed as  $D_{ii'} = (\sum_j (p_{ij} - p_{i'j})^2)^{1/2}$ , where  $p_{ij}$  is the allele frequency at the  $j$ th locus in the  $i$ th population. The distances were then clustered using the UPGMA clustering method of SAS. Next, it was necessary to determine how many clusters were in the samples. Although a number of methods are available to achieve this goal, all have limitations and restricting assumptions, and none is universally accepted as the best for all situations. Thus, we used an empirical clustering criteria based on prior knowledge of the poultry industry. For example, it is known that all broiler male lines across the industry have a White Cornish ancestry in common, whereas all white egg-type lines have a Single-Comb White Leghorn ancestry in common; therefore, we set a genetic distance between clusters such that all broiler male lines were in the same cluster as our criterion, then used this cluster distance to differentiate all clusters. The effect of clustering criteria on results was examined by comparing outcomes that would have been obtained had the number of clusters been more or less than that set as our criterion. PCA (16), another clustering method, also was conducted for comparison.

Allele frequencies were averaged over samples first within clusters, then across clusters. These averages provided our best estimate of allelic frequencies in the HAP. These estimates are a biased representation of the allele frequency distribution due to ascertainment bias, however.

**Ascertainment Bias Correction.** Ascertainment bias is relevant for SNPs because of the way in which SNPs are discovered. In poultry, SNPs were discovered by comparing sequences from essentially only two chromosomes, one from three birds that were sampled sequenced at 1/4X coverage (10) and the other being the Red Jungle Fowl 6.6X whole genome sequence (9). Although a high stringency was applied to the reads, these SNPs are putative as they could be due to sequencing errors. Verification is needed to confirm that the loci are polymorphic, which was one of the goals of this research. Because the results of SNP genotyping are based on polymorphisms observed from a limited number of individuals (two for poultry), these represent conditional probabilities, that is, the probability of observing an SNP in a randomly genotyped individual  $j$  given that it was observed in the previously sequenced individuals  $s$  and  $s'$ . As such, the observed frequencies will tend to overestimate the frequencies of common alleles and underestimate those of rare alleles. Ascertainment bias correction is necessary to obtain the true probability distribution of SNP frequencies in the HAP. The correction was applied to these data using the methods provided by Nielsen and colleagues (11, 29). In essence, this procedure estimates the actual SNP frequency had all of the birds been sequenced rather than genotyped and had all resulting SNP observed been scored.

**Proportion of Missing Allele Calculations.** The corrected allele frequency distribution of SNPs in the HAP presents a standard for comparing the effect of inbreeding on that distribution. Let  $L_i$  be the observed number of loci with frequency  $i/2N$  and let  $N$  be the total number of individuals sampled. The corrected relative frequencies ( $C_i$ ) of samples in those bins were found using the Fortran program AS.BIAS given in [SI Materials](#) based on the formula of Nielsen and colleagues (11, 29) for a SNP discovery depth of 2. The depth was 2 because SNP were discovered based on only two birds at a time (10). The proportion of the corrected frequency distribution represented by each SNP was  $SNP\_WT_i = (C_i/L_i)$ . These SNP.WTs sum to 1 over loci in the HAP; as such, these SNP represent sentinels of biodiversity. If absent in a subpopulation, they represent that proportion of the original distribution missing in that subpopulation. The proportion of missing alleles in any subpopulation is found by first scoring each allele, based on MAF, as present ( $z_i = 1$ ) or absent ( $z_i = 0$ ) in the subpopulation, then weighted using the SNP-WTs from the HAP:

$$\Omega = \sum_{i=1}^{L_i} SNP\_WT_i(z_i)$$

The weighted mean estimates the proportion of alleles missing ( $\Omega$ ) in the subpopulation relative to the HAP.

