# Use of Genetic Polymorphisms to Assess the Genetic Structure and Breed Composition of Crossbred Animals

WEERASINGHE MUDIYANSELAGE SHALANEE PRIYANGA WEERASINGHE

Bachelor of Veterinary Medicine and Animal Science, University of Peradeniya, 2003

Master of Philosophy in Animal Science, University of Peradeniya, 2010

A thesis submitted for the degree of Doctor of Philosophy of the

University of New England

June 2015



School of Environmental and Rural Science

# Abstract

This thesis explores the accuracy of methods to estimate the breed composition of crossbred animals which have unknown pedigree. Herein I present the use of SNP technologies to estimate the breed composition of small-holder crossbred dairy cattle in developing countries for the first time. Before this could be done there was a need to determine: what are the accuracies of different methods for estimating breed composition? The genetic structure of animals, the design of reference populations, the number of SNP markers and the model selected has possible consequences for estimation of breed composition. Once the effect of the above factors on the accuracy of estimation of breed composition is identified, it is possible to estimate with confidence the breed composition of crossbred animals that have no recorded pedigree. The overall aim of this thesis was to investigate the use of high-density SNP data to understand the livestock breed's population structure and estimate the breed composition of crossbred animals.

The first block of research was to determine the genetic structure of Australian sheep breeds and estimate the accuracy and bias of SNP-based estimates of breed composition in F1 sheep, using Sheep-CRC data as a model. These results confirmed that principal component analysis (PCA) provided useful visualisation of the genetic diversity of Australian sheep breeds and crosses, and identified a variety of errors in the data. I therefore, recommend the use of PCA for initial examination of population structure of the data and identification of data errors before performing a detailed genomic analysis. Subsequent results indicated that for a simple cross the breed composition can be estimated with high accuracy and low bias when using high density SNP data. This was true even when breeds unrelated to the cross being analysed, were included in the analyses.

The second block of research examined the accuracy and bias of estimates of breed composition of alternative F1 crosses and three-way crosses in sheep. It examined the effect of population stratification of the crossbred population on the accuracy and bias of estimates. The results showed breed composition of composite crossbred sheep can also be estimated with high accuracy and low bias. The estimates of breed composition of crossbreds from purebred parents showed higher accuracy and less bias than the crossbreds from composite parent breeds. However, the accuracy was slightly increased and bias reduced when the ancestors of the composite breed were included in the admixture model. My results also showed that the design of the reference population has a considerable impact on the estimates of breed composition of crossbreds. The estimates of breed composition of crossbreds among genetically distant parent breeds were somewhat more accurate and less biased than those among genetically closer breeds. Overall, for crosses among breeds that are clearly genetically distinct, errors of estimation of breed proportion should have S.D. in the range of 0.02 to 0.03. My results can be extrapolated to predict that the S.D. could be substantially lower than I observed, when crosses are between highly diverse breeds, such as between *Bos indicus* and *Bos taurus* cattle. Finally, when all the crossbred data were combined into a single population and estimates of each crossbred group were obtained in a single analysis, estimates still had high accuracy and low bias.

The third and fourth blocks of research applied the best models from the first and second blocks of my research to estimate the breed composition of different populations of crossbred dairy cattle used by small-holder dairy farmers in East Africa. These crossbred populations have little information on the admixture levels, due to the lack of pedigree recording across the many generations of crossbreeding in the region. First, in both research blocks, the genetic structure of indigenous and crossbred cattle populations was examined. Then the dairy proportion and individual breed composition were estimated using the ADMIXTURE program. Chapter 3 focuses on the genetic structure and admixture levels of

iii

small-holder cattle in Kenya and Uganda. One of the most significant findings to emerge from this study was that the correlation of the estimated dairy proportion across the various types of analysis is greater than 0.99, indicating that these estimates are highly robust to model assumptions. The average dairy proportion of crossbred cows is 0.69 with S.D. 0.21, whereas the bulls sampled from the same area as crossbred cows show a comparatively higher dairy proportion (0.79 with S.D. 0.18). Another important finding was that the farmers' prediction of the dairy proportion of individual animals explains only 16% of the real variation of the dairy proportion. It demonstrates that many farmers have a poor understanding of the degree of crossbreeding of their animals. Also, the average dairy proportion is lower than expected and the variation much higher than expected in these populations of crossbred animals. The estimates of individual dairy breed proportion (Holstein vs Ayrshire vs Jersey, etc.) reflected what was expected in terms of the recorded history of different breed use in different locations. Although, the estimates obtained were dependent on the model used, and without an independent validation of the estimates, they were considered unreliable for use in analyses of breed performance.

The fourth research block examined the estimation of breed composition of crossbred dairy cattle in Tanzania and Ethiopia, which were believed to have different small-holder dairy cattle populations compared to Kenya and Uganda. Here, I explored in more detail than was possible with the data in the third research chapter, the genetic structure of East African indigenous cattle. The results here and in the third research chapter are consistent with what are believed to be the major historical migrations of *Bos indicus* and *Bos taurus* cattle on the African continent. I found that several indigenous breeds were contaminated with the recent introduction of European breed germplasm. I propose methods for using SNP-based breed composition estimates to aid genetic conservation programs. Furthermore, I discovered that the molecular estimates of breed composition of the Tanzania synthetic Mpwapwa breed confirmed its reported origins and also found that it constituted a relatively well-defined breed. I observed Holstein

and Friesian breeds to have the most significant impact on small-holder crossbred cattle in Ethiopia and Tanzania. That the average and variance of total dairy proportion differed between the two countries and differed markedly between locations within the countries, has important implications for the design of programs to improve small-holder (SHD) management and genetics in different areas.

In conclusion, these results demonstrate the ability to identify the genetic structure and estimate the breed composition of crossbred cattle which do not have pedigree records. In this regards the design of the reference population and the number of SNP markers have a considerable effect on the accuracy of estimations. This knowledge has the potential to be used in livestock genetic improvement in the developing world. The further research could be carried out to design the cheap SNP test that could be used to genotypes the crossbred animals owned by small-holder farmers in these countries. Additionally, genotypes of more *Bos indicus* and African *Bos taurus* cattle breeds relevant to the sample collected regions could be included in the SNP panel to prevent the bias of estimations.

# Declaration

I certify that the substance of this thesis has not been submitted previously for any degree and is not currently being submitted for any degree or qualification.

I also certify that any help received and all sources used in preparing this thesis have been acknowledged therein.

W.M.S.P.Weerasinghe

12th June 2015

# Acknowledgments

First, I would like to thank my principal supervisor, Professor John Gibson, for all the encouragement and support he has given me. I thank him for his enthusiastic attitude and his knowledge which has enabled me to learn both the theoretical and applied aspects of scientific research, something I will value for my entire life. Every discussion and consultation with Professor Gibson left me both enlightened and compelled to think more deeply about my research. Second, I would like to express my heartfelt gratitude to associate supervisors, Associate Professor Cedric Gondro for his encouragement, support and never-ending patience. Thanks for showing me an inspirational journey in computer programming. I would also like to extend my gratitude to Dr Gilbert Jeyaruban who gave me strength and courage throughout my candidature especially when I was in difficult situations.

I would like to express my sincere gratitude to the Australian Government for granting me an Endeavour Postgraduate Research Scholarship which has provided me with the financial assistance to complete my PhD with so much ease and comfort. I also thank the University of New England, Australia, for providing all the infrastructure and facilities during my candidature period.

I would also like to thank my fellow scientists and friends who have helped me with my thesis.

My hard-working parents have sacrificed their lives for my sister, brother and me to provide us better education, and with unconditional love and care. I love them so much, and I would not have made it this far without them. I especially thank my dad for all the courage and care he has given me and my children. I will forever be grateful. I know you were devastated when we came to Australia, but you kept it from us. I miss you forever, I know wherever you are, you will be proud of me. I have completed your dream.

This thesis would not have been possible without the loving and selfless understanding and patience of my husband, Shantha. My lovely angels Ashini and Piyumi, thank you so much for your patience and unfailing love.

**Table of Contents**

# List of Tables

# List of Figures

xvii

xviii