

RESEARCH ARTICLE

Open Access



Developing a 670k genotyping array to tag ~2M SNPs across 24 horse breeds

Robert J. Schaefer¹, Mikkel Schubert², Ernest Bailey³, Danika L. Bannasch⁴, Eric Barrey⁵, Gila Kahila Bar-Gal⁶, Gottfried Brem⁷, Samantha A. Brooks⁸, Ottmar Distl⁹, Ruedi Fries¹⁰, Carrie J. Finno⁴, Vinzenz Gerber¹¹, Bianca Haase¹², Vidhya Jagannathan¹³, Ted Kalbfleisch¹⁴, Tosso Leeb¹³, Gabriella Lindgren¹⁵, Maria Susana Lopes¹⁶, Núria Mach⁵, Artur da Câmara Machado¹⁶, James N. MacLeod³, Annette McCoy¹⁷, Julia Metzger⁹, Cecilia Penedo¹⁸, Sagi Polani⁶, Stefan Rieder¹⁹, Imke Tammen¹², Jens Tetens^{20,21}, Georg Thaller²⁰, Andrea Verini-Supplizi²², Claire M. Wade¹², Barbara Wallner⁷, Ludovic Orlando^{2,23}, James R. Mickelson²⁴ and Molly E. McCue^{1*}

Abstract

Background: To date, genome-scale analyses in the domestic horse have been limited by suboptimal single nucleotide polymorphism (SNP) density and uneven genomic coverage of the current SNP genotyping arrays. The recent availability of whole genome sequences has created the opportunity to develop a next generation, high-density equine SNP array.

Results: Using whole genome sequence from 153 individuals representing 24 distinct breeds collated by the equine genomics community, we cataloged over 23 million de novo discovered genetic variants. Leveraging genotype data from individuals with both whole genome sequence, and genotypes from lower-density, legacy SNP arrays, a subset of ~5 million high-quality, high-density array candidate SNPs were selected based on breed representation and uniform spacing across the genome. Considering probe design recommendations from a commercial vendor (Affymetrix, now Thermo Fisher Scientific) a set of ~2 million SNPs were selected for a next-generation high-density SNP chip (MNEc2M). Genotype data were generated using the MNEc2M array from a cohort of 332 horses from 20 breeds and a lower-density array, consisting of ~670 thousand SNPs (MNEc670k), was designed for genotype imputation.

Conclusions: Here, we document the steps taken to design both the MNEc2M and MNEc670k arrays, report genomic and technical properties of these genotyping platforms, and demonstrate the imputation capabilities of these tools for the domestic horse.

Keywords: Equine genomics, Whole genome sequence, SNP-tagging, SNP chip, Variant recalibration, SNP discovery, SNP informativeness, SNP validation, Linkage disequilibrium

Background

Soon after the horse reference genome from Twilight, a female Thoroughbred, was sequenced using Sanger technology [1], a genotyping array (Illumina EquineSNP50 BeadChip) was developed to enable whole genome mapping using ~50k (54,602) single nucleotide polymorphism (SNP) markers from low-coverage Sanger sequence of 7 horses representing 7 different breeds (an Andalusian,

Arabian, Akhal-Teke, Icelandic, Standardbred, Thoroughbred and a Quarter Horse) [2]. Shortly thereafter, a slightly higher density array (Illumina Equine SNP70 BeadChip) with ~65k (65,157) informative SNP markers was developed. These widely used, now legacy, SNP arrays (with a combined 74,056 unique SNPs, ~74k hereafter [3]), have successfully enabled genetic studies examining domestication and selection [4, 5], disease and performance trait mapping [6–15], and population structure and dynamics [2, 16, 17]. However, extensive population structure and the low extent of linkage disequilibrium (LD) existing in many horse breeds severely limits conventional mapping

* Correspondence: mccu0173@umn.edu

¹Department of Veterinary Population Medicine, College of Veterinary Medicine, University of Minnesota, St. Paul, MN, USA

Full list of author information is available at the end of the article



approaches with the relatively low SNP density on the current SNP array [17]. Since the initial Sanger shotgun sequencing of Twilight, whole genome sequence (WGS) has been generated for hundreds of horses [15, 18–25], prompting the development of a new higher-density genotyping array for the horse.

Here we describe the steps taken, including careful and extensive variant filtering, to create both a 2 million (2M) SNP array and a 670 thousand (670k) SNP array from over 23 million variants discovered from whole genome sequence of 153 horses representing 24 breeds (See Additional file 1: Table S1). Genotypes from the 2 million SNP genotyping array (MNEc2M) in a cohort of 332 horses, from 20 actively-researched or economically-important breeds that represent known genetic diversity in the domestic horse [17], were used to select SNPs for inclusion on the commercially-available 670k SNP array (MNEc670k). The MNEc670k array was designed for accurate genotype imputation up to the higher density, 2M SNP set present on the MNEc2M array. We report summary statistics, broken down by breed, for both arrays as well as preliminary results on genotype imputation performance from the MNEc670k array to the SNP density on the MNEc2M array.

Results

Variant discovery from whole genome sequence

Whole genome, 100 base pair (bp), paired-end Illumina HiSeq reads were generated for 153 horses (including Twilight), representing 24 distinct breeds, at a depth between 1.7X and 64X, with a median depth of 13X (Additional file 1: Table S1). Read mapping was performed using the PALEOMIX genome mapping protocol to efficiently process samples in parallel and to assess individual sample quality. Each sample was mapped to the EquCab2.0 reference genome, which had been extended with an additional 7850 de novo assembled scaffolds (See Methods), to produce a total of nearly 48 billion unique reads aligned to the nuclear genome (See Additional file 1: Table S1). Variants were identified by extending the PALEOMIX framework to identify SNPs

using two variant callers (see Methods). To maximize efficiency of variant calling in individual breeds, and to minimize bias due to variable sequencing depth of coverage, individuals were broken up into 16 variant calling groups (see Additional file 1: Table S1, Variant Calling Group columns) by estimated depth of coverage and breed. Variants were called using permissive parameters in both the GATK UnifiedGenotyper [26] as well as SAMtools ‘mpileup’ utilities [27]. Approximately 23 million potential SNPs were in the intersection of SNP sets called by GATK and SAMtools. These ~23 million SNPs were kept for further analysis and validation (See Table 1; Additional file 6).

Precision of GATK QUAL scores for variants identified by both callers

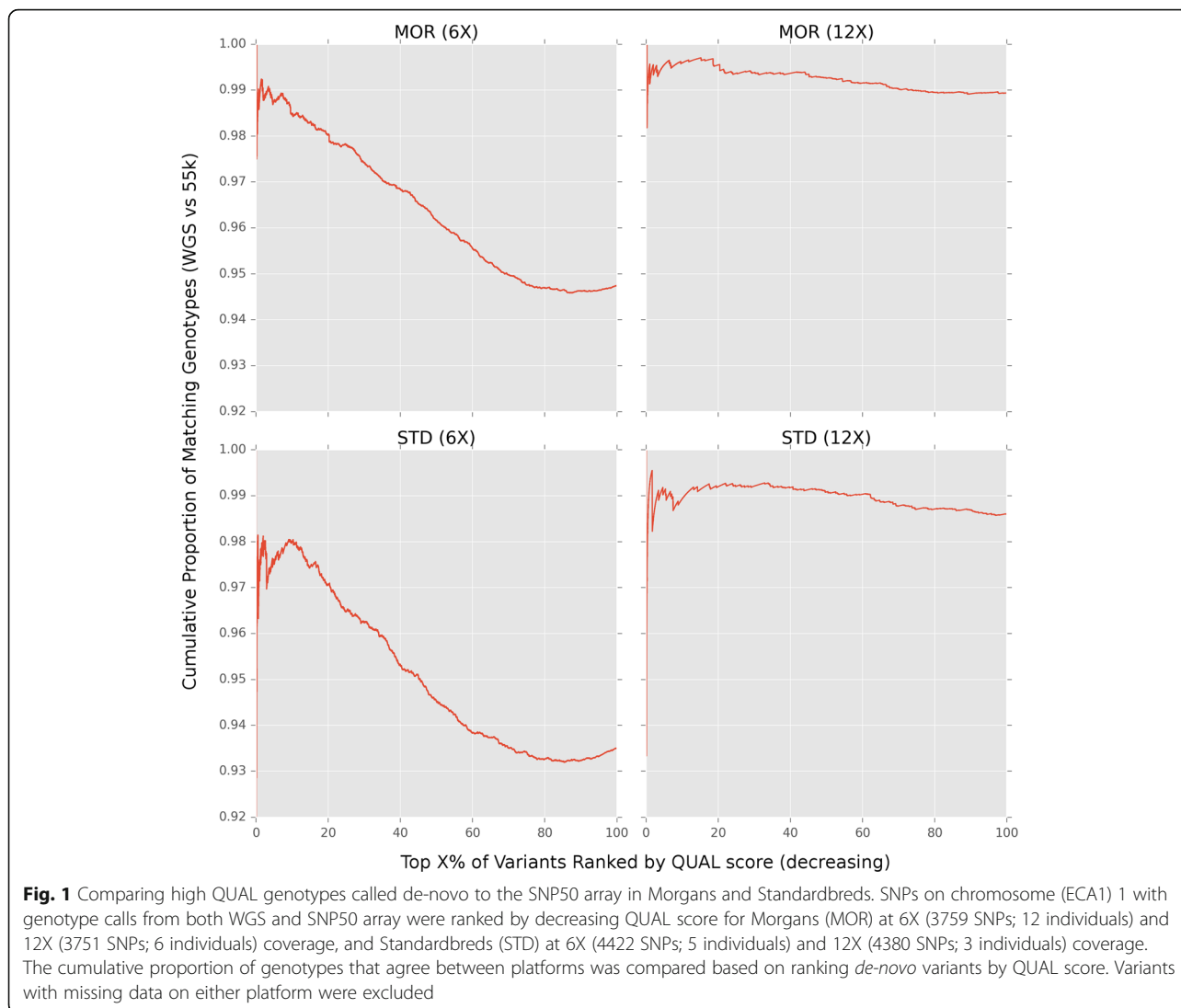
We evaluated the performance of GATK QUAL scores by comparing genotypes generated from WGS to genotypes generated from the legacy Equine 54 K SNP chip in 23 horses at two different sequencing depths. Variants detected on chromosome 1 (ECA1) by both WGS and 54k SNP chip were ranked by decreasing QUAL score. The proportion of genotype calls that agreed between 54K and WGS (i.e., precision) was compared based on ranked QUAL score (i.e., recall) [28]. For each ranked variant we evaluated the precision at that point, which is the cumulative proportion of genotypes that agree between the two genotyping methods (Fig. 1). We note that this approach evaluates the concordance between the two genotyping technologies and is thus unable to assess the accuracies of the underlying genotypes themselves. However, this approach does take into account the imbalance of false positives/negatives between genotypes called by WGS and by the 54K SNP chip [29].

Considering 100% of variants detected by both technologies, genotypes had an overall precision of ~99% for 12X-called variants and 94–95% for 6X-called variants (See Fig. 1, x-axis). This is consistent with results from a similar study that compared de novo variant genotypes to array based genotypes in the Franches-Montagnes horse breed [25]. Yet, many SNPs with very high QUAL

Table 1 SNP Sets used at various steps in array design

Set Name	Number of unique sites	Set Description	WGS Data	Array Data
23M	22,557,988	Possible/Discovered Variants	x	-
10M	11,435,936	Array Compatible	x	-
5M	5,443,950	Array Candidate	x	-
2M	2,001,826	Test Array	x	x
1.8M	1,846,988	Test Array Converted	x	x
670k	670,805	Commercial Array	x	x

Variants discovered from whole genome sequencing were filtered at various steps for quality control or using array design criteria. Six distinct sets of variants ranging from the initial ~23M high-quality, variants discovered from WGS to the 670k variants available in the commercial genotyping array are described throughout this manuscript



scores (e.g., SNPs within the top 10%; Fig. 1) had disagreeing genotype calls between WGS and 54K SNP chip. Additionally, when considering between 80 and 100% of variants ranked by QUAL score (Fig. 1), in both 6X MOR and STD comparisons, the proportion of matching genotypes increases, indicating that there are many variants with low QUAL scores and high precision. These results indicated that QUAL scores alone did not adequately rank variants, and that additional metrics were necessary to improve the reliability of SNPs ultimately chosen for the higher density genotyping arrays.

Identifying gold standard reference set for variant recalibration

In addition to QUAL type quality scores, GATK outputs additional metrics such as depth, quality of depth, Fisher strand position, map quality rank sum, and read position rank sum for each variant based on read statistics (See

Methods for details of these metrics). These values were used to train linear mixed models, using the GATK VariantRecalibrator [30], to assign a composite quality score (VQSLOD) to help detect type I and II errors.

Training these models required a “gold standard” reference set of known genotypes across multiple individuals. Lacking this resource in the domestic horse, we defined three high confidence, putative “gold standard” datasets with which to train the VariantRecalibrator: 1) SNPs on the legacy SNP50 chip; 2) WGS variants which were seen in four or more (4+) calling groups (Additional file 1: Table S1); and 3) WGS variants that were in the top 1 % of QUAL scores. Models were trained on features from “gold standard” variants present on chromosomes 2-32 and used to calculate VQSLOD scores for SNPs on chromosome 1. To assess the performance of VQSLOD scores generated by each training set, scores from each group were compared to each other in

addition to QUAL scores using horses genotyped on the 54K SNP Chip.

SNPs on chromosome 1 (ECA1) were ranked either by decreasing QUAL or by decreasing VQSLOD score generated from each of the three gold standard training groups. Figure 2 shows the proportion of matching genotypes between WGS and SNP50 platforms for each of the four groups. With the exception of 6X Morgans, SNPs with high VQSLOD scores agreed more often with genotypes called on the 54K SNP chip compared to SNPs with high QUAL scores alone. For example, the top 10% highest scoring VQSLOD scores had a higher concordance rate than QUAL scores in all cases except for MOR6X, where only two of the recalibrated scores marginally outperformed QUAL scores (See Fig. 2; see discussion). Additionally, VQSLOD variants did not drop below the overall discordance rate, showing an overall better ranking than QUAL scores alone (red curve, Fig. 2).

While an exact differentiation between the gold-standard training groups remains unclear (See Discussion), recalibrated VQSLOD scores were better able to differentiate high scoring, false positives than QUAL scores alone, making them a more informative metric for selecting the final SNPs for use on the high-density commercial array. VQSLOD scores were calculated for all 23 million variants, in all remaining horses, using WGS variants called in the 4+ breed calling groups (described above) as the training set (candidate 23 M SNP set). Scored variants were then filtered in several steps (described below) to select the subset of SNPs to be included on the new higher density arrays.

Preliminary 5M SNP selection for genotyping probe design

We analyzed the 23M candidate SNP set with the goal of generating a target set of ~5 million high-confidence SNPs that would be compatible with array design. Following

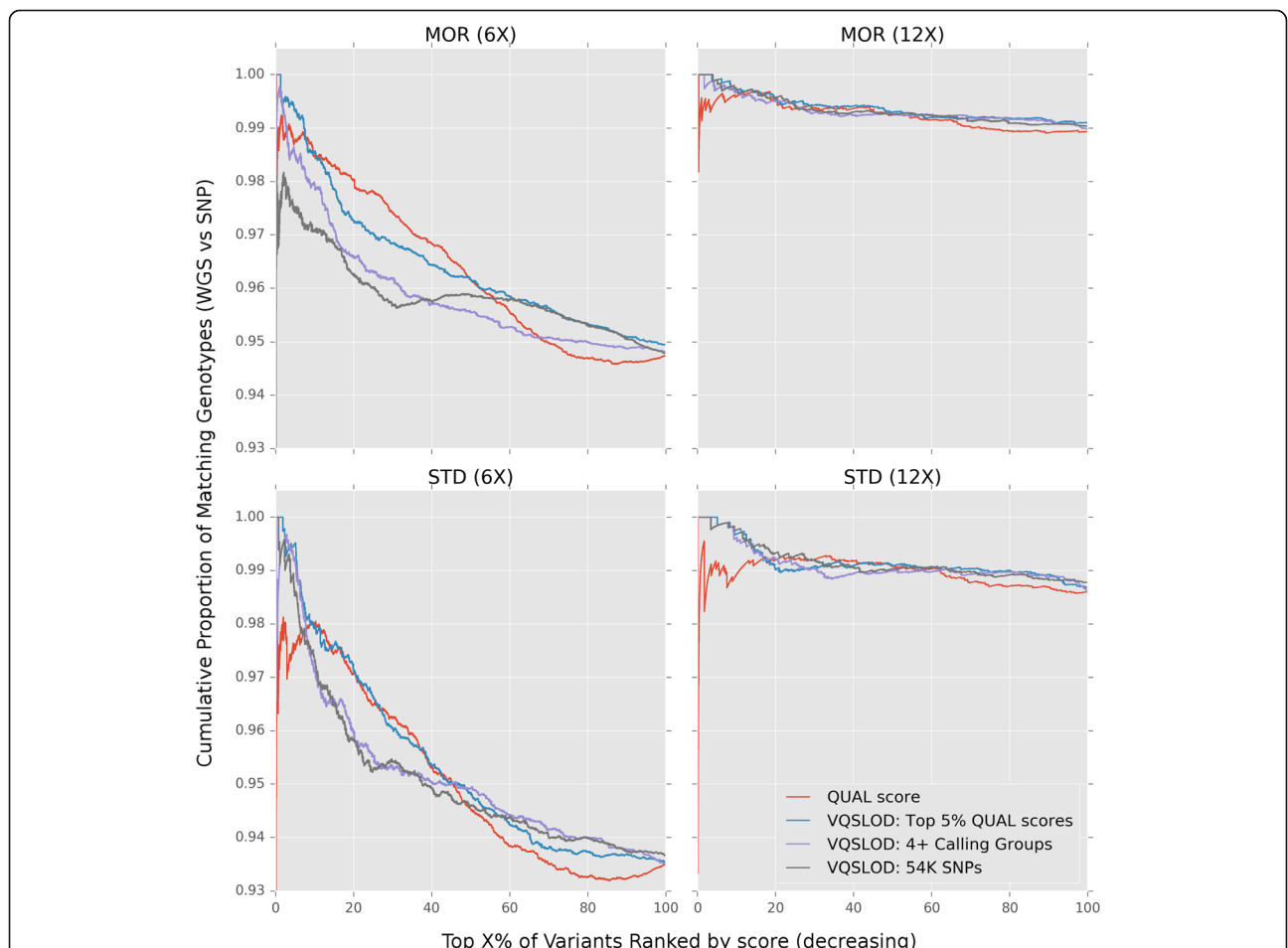


Fig. 2 Comparing QUAL ranked SNPs to VQSLOD ranked SNPs. VQSLOD scores were calculated from three different “gold standard” reference groups in Morgans (MOR) and Standardbreds (STD) using GATK VariantRecalibrator. Compared to QUAL scores (red line), high VQSLOD scored variants (top 10% variants by score) have a lower number of mismatched genotypes across the SNP50 BeadChip and variants discovered de novo

criteria provided by Affymetrix (now Thermo Fisher Scientific), 780 SNPs were filtered because they fell within repetitive regions; 9,342,733 SNPs were filtered because they were within 20 base pairs (bp) of another variant; and 4,858,273 SNPs were filtered for having a minor allele frequency below 1% in the 153 WGS horses. After filtering, approximately 10 million SNPs remained for further design selection (hereafter referred to as the 10 M SNP Set, Table 1). Of these ~10 million array-compatible candidates, 48,485 SNPs (65% of the total) overlapped with the sites available from the ~74k legacy SNP sets (although all legacy SNP sites were included in MNEc2M/670k array design). Additional file 9: Figure S1 shows the alternative allele frequency distribution of the 23M de novo and 10 million design candidate SNPs (10 M SNP set). The mean allele frequency is below 0.1, though there is a long tail of SNPs with a high alternative allele frequency. Interestingly, there are 445,421 SNPs (23M set) and 46,455 SNPs (10 M set) where the alternative allele frequency is 1.0. This is likely due to errors in the EquCab 2.0 reference sequence (Sanger reads from Twilight) as Twilight's whole genome Illumina sequence data (collectively at more than 25X), and whole genome sequences of all other horses, contained the alternate allele. These SNPs were excluded from array design.

To ensure backwards compatibility, we started with the ~74k SNPs that were present on the two legacy (54K/65K) arrays. To fill the remaining target of ~5 million probe design candidates, we added SNPs from the 10 M SNP set based on even genome distribution, informativeness among the variant calling groups (See Additional file 1: Table S1), and linkage disequilibrium (LD). LD was calculated for all pairs of SNPs within 10 kb of each other throughout the genome. SNPs that were in high LD ($r^2 > 0.90$) were filtered based on whether or not they were present in draft or pony variant calling groups to control for the fewer number of samples used in variant discovery. Previous SNP designs [2] show an underrepresentation of informative SNPs from these groups, which leads to poor mapping resolution in draft and pony breed groups (See Additional file 2: Table S2). SNPs discovered in both draft and ponies were prioritized over those discovered in one group and over those absent in both. If SNPs were discovered in an equal number of priority groups, VQSLOD scores were used to break ties.

After applying these SNP candidate criteria, 5,443,950 SNPs (5M set, see Table 1) were kept, of which 2,199,467 SNPs occurred in pony calling groups and 2,782,917 SNPs occurred in draft calling groups. A total of 1,695,347 SNPs occurred in both ponies and drafts. Flanking sequences (35 bp upstream and downstream) for these 5M filtered SNPs were submitted to Affymetrix for SNP probe design analysis.

SNP selection for the high density MNEc2M SNP array

SNP conversion recommendations for the 5 M SNPs were provided from Affymetrix in both forward and reverse strand directions. SNPs were assigned to one of four groups based on decreasing probability of successful probe design: 'recommended', 'neutral', 'not recommended', and 'not possible'. To achieve an even distribution of ~2 million SNPs, the equine reference genome (~2.7 Gb) was divided into approximately 54,000 50 kb windows and a target of 37 SNPs per window was established. SNPs within each window were chosen for inclusion in the MNEc2M SNP set using a greedy algorithm. Briefly, the ~74,000 SNPs on previous generation arrays, as well as SNPs within the equine major histocompatibility complex (MHC) region (ECA20:28.7-33.6 Mb), were given VIP status and were automatically included in both forward and reverse strand directions. SNPs were added to windows until the target of 37 SNPs was met. If a window had more than 37 candidates, SNPs were selected based on by Affymetrix recommendation group (See Methods for details). In total, given the above criteria, 2,001,826 high quality SNPs (2M SNP set) were chosen and submitted for probe design to comprise the MNEc2M genotyping array (Additional file 7).

SNP selection for the MNEc670k array

A cohort of 347 horses were genotyped on the MNEc2M high-density array using DNA isolated from blood ($n = 286$) or hair roots ($n = 61$) (see Methods). The ~2M SNP genotypes were split across three different physical arrays. Genotypes were called using Affymetrix Power Tools and sample quality control was assessed according to Affymetrix best practices (See Methods). For each sample, the three physical arrays were assessed separately for quality control. In total, 320 samples passed in all three arrays while 27 samples had one or more arrays that did not pass quality control. Of the 27 samples that failed, 7 samples had two passing arrays, 5 samples had one passing array, and 15 DNA samples failed to meet genotyping quality control metrics on all three physical arrays. Failed arrays were removed from the analysis. If a sample had at least one array, genotypes from those arrays were retained. In total, viable genotypes were produced for 332 horses.

Of the failed arrays, 25 had DNA isolated from hair roots and 2 had DNA isolated from blood. Hair root DNA had a lower average DNA concentration (Pico Green) when re-hydrated (2.7 ± 2.9 ng/ μ L) than did DNA samples isolated from blood (43.1 ± 55.5 ng/ μ L). To determine if sample origin (blood versus hair roots) or DNA concentration, or both, were associated with failure to pass genotyping quality control metrics, a logistic regression for sample success on DNA concentration and blood/hair status was performed. (Additional file 14: Figure S6). Both DNA source ($p \leq 4.05e-04$) and DNA concentration

($p \leq 2.77e-10$) significantly influenced the probability of samples producing quality genotypes. Both factors also had substantial coefficients in the model indicating a large magnitude of effect. In the logistic regression, the factor indicating hair root had a strong negative model coefficient of -2.70 while DNA concentration had a positive coefficient of 1.71 . (see Discussion).

After quality control steps for samples, genotypes were called for all 2,001,826 SNPs. Genotypes were assessed for clustering quality using the *Metrics.R* script provided by Affymetrix (see Methods). In total, 92.2% of SNPs on the MNEc2M array passed quality control producing a set of 1,846,988 high-quality SNPs (1.8 M; see Table 1) genotyped on the 332 horses remaining in the analysis (see Methods).

SNPs exclusive to a single WGS variant calling group were validated (polymorphic), on average, at a rate of 80%. Yet, validation rates were much higher in SNPs discovered in multiple calling groups (>96%) (See Additional file 3: Table S3). Genotypes produced from the 332 horses with passing genotypes on the MNEc2M array (“SNP” genotypes) were combined with genotypes discovered from the 153 whole genome sequence horses (“WGS” genotypes) to create a dataset containing 1.8 M genotypes for 485 horses (“WGS + SNP” genotypes). These variants were analyzed to select a subset of SNPs for inclusion on the MNEc670k array, with the intent that this array would be designed for genotype imputation to the full 2M SNP set.

To maximize the information content of the MNEc670k array SNPs, multi-marker r^2 statistics were calculated on the 1.8M high-quality candidate SNPs to identify ‘tagging SNPs’ that allowed for efficient reconstruction of genomic haplotypes in the 485 horses both within and across breeds. Tagging SNPs that reconstructed haplotypes across all 485 horses (inter-breed tag SNPs) were identified using the software FastTagger [31]. In total, 355,903 tagging SNPs were needed to reconstruct haplotypes present across all the horses in the cohort with a minor allele frequency (MAF) > 0.01 and tag-SNP $r^2 > 0.99$. These SNPs were included on the MNEc670k imputation array (see Additional file 4: Table S4; Inclusion criteria: Inter).

Haplotype tagging SNPs were also examined at the breed-specific level. Horses were split into 15 *tagging breed groups* based on the minimal sample size necessary to perform SNP tagging (see Additional file 5: Table S5). Tagging SNPs were identified separately in each tagging breed group using FastTagger, then a subset of population specific tag SNPs were identified using the software Multi-Pop-TagSelect [32]. Combined, a total of 1,754,075 SNPs were needed to reconstruct fine-level, breed-specific haplotypes in all of these breed groups (MAF > 0.10 and tag SNP $r^2 > 0.90$; See Additional file 8). Breeds varied in the number of tag SNPs needed to reconstruct haplotypes.

Table 2 shows the number of tagging SNPs required to reconstruct haplotypes in each of the tagging breed groups. The Ponies, Draft and Quarter Horses (tagging breed groups) required the most tagging SNPs, each requiring over 350,000 SNPs to reconstruct breed-specific haplotypes, while Thoroughbred, Icelandic and Lusitano tagging breed groups each required less than 150,000 tag SNPs, to reconstruct haplotypes.

Tagging SNPs that were informative in 5 or more tagging breed groups were included on the MNEc670k array ($n = 206,822$; see Additional file 4: Table S4; inclusion criteria: Intra). An additional 13,993 SNPs that tagged haplotypes in four of the breeds requiring a larger number of tag SNPs (Quarter Horse, Pony, Morgan, Standardbred) were also included in the array (inclusion criteria: Diverse). Additionally, 7394 SNPs were included due to their location within the equine MHC region (inclusion criteria: MHC), 16,398 SNPs were included to increase SNP density in 12,104 (24.3%) 50 kb genomic windows to at least 8 SNPs (inclusion criteria: Density), and 70,295 SNPs were retained for backwards compatibility with legacy arrays (inclusion group: VIP). Collectively, 670,805 SNPs were included on the MNEc670k commercial array.

Imputation accuracy from the MNEc670k SNP set to the MNEc2M SNP

Genotype imputation accuracy from the MNEc670k array to the MNEc2M array was quantified using the 485 horse reference population. For each of the tagging breed groups a subset of 1/3 random individuals were masked down to the MNEc670k SNP set. Genotypes

Table 2 Number of breed specific tagging SNPs

Breed Group	Number Of Tag SNPs
Thoroughbred	144,175
Lusitano	148,097
Icelandic	148,206
F. Montagne	199,244
Arabian	199,264
Belgian	217,882
Marremanno	223,568
Standardbred	245,149
Trotter	256,790
Warmblood	304,510
Morgan	335,677
Land Race	338,040
QuarterHorse	366,702
Draft	370,701
Pony	387,279

Number of tagging SNPs required to reconstruct haplotypes in each breed (MAF > 0.10 and 0.90 r^2 , See Methods)

were then imputed to the MNEc2M SNP using the reference population after removing the individuals being imputed (See Methods). Imputation accuracy was measured as total genotype concordance across imputed individuals (correctly inferred genotypes/total number of imputed genotypes; see Methods for details).

Genotype imputation accuracy from the MNEc670k SNP set to the MNEc2M SNP set ranged between 96.6 and 99.4% in the 15 breeds tested (see Table 3). Tagging breed groups with over 99% mean imputation accuracy included Arabians, Belgians, Lusitano, Maremmano, Pony and Thoroughbred. No tagging groups were below 98% with the exception of the Draft group. This breed group contained both continental European draft breeds as well as British Isles draft breeds which have been previously shown to have distinct sub-population structure [17]. Random sampling for imputation validation in the Draft group included the only two Percheron samples (M1542 and M1545), which underperformed (average 92.6% accuracy) compared to the other Draft samples in the imputed group (average 98.5% accuracy), performing similarly to the other heavy horse breeds (e.g. Belgian). We suspect that an increased representation of Percheron samples in the reference population would increase imputation performance for individuals within that tagging breed group (See Discussion).

The effect of allele frequency was assessed by measuring the Pearson correlation between the imputed minor allele

dose and the true minor allele dose for SNPs binned by minor allele count (Additional file 10: Figure S2). Comparing dosage in terms of allele count, here, allows for direct comparison across populations that have varying allele frequencies based on the number of individuals. As observed with imputation in other species [33], SNPs with a low number of observed minor alleles in the population have a lower overall imputation accuracy, though imputation accuracy quickly improves with a higher allele count. With an alternate allele count of 8 or more (~2% minor allele frequency in this population), imputation accuracy was above 90% in most tagging breed groups (Additional file 10: Figure S2).

SNP properties of the MNEc2M and MNEc670k arrays

SNPs in gene coding regions

SNP positions for both arrays were compared to 26,991 predicted and annotated gene models from EquCab2. Of the 2,001,826 SNPs on the MNEc2M array, 591,521(29.5%) SNPs were within 17,128(63.5%) gene models. Likewise, of the 670,805 SNPs on the MNEc670k array, 192,681(28.7%) SNPs were within 14,758(54.7%) gene boundaries. In comparison, of the legacy 74,056 SNPs, 20,950(28%) were within 8249(30%) annotated gene models.

SNP informativeness and inter-SNP distance

Inter-SNP distance was calculated for the MNEc2M, as well as the MNEc670k, SNP sets at different levels of MAF (see Fig. 3 and Table 4) to assess the distribution of SNPs across the genome. On average, 1250 and 3756 bp, respectively, separated variants on the two arrays. Informativeness, defined as the number of SNPs with at least one heterozygote, was calculated for the same MAF cut-offs (Table 4). Inter-SNP distance, as well as informativeness, were broken down by breed for both the MNEc2M SNP set (See Additional file 11: Figure S3) and MNEc670k SNP set (See Additional file 12: Figure S4).

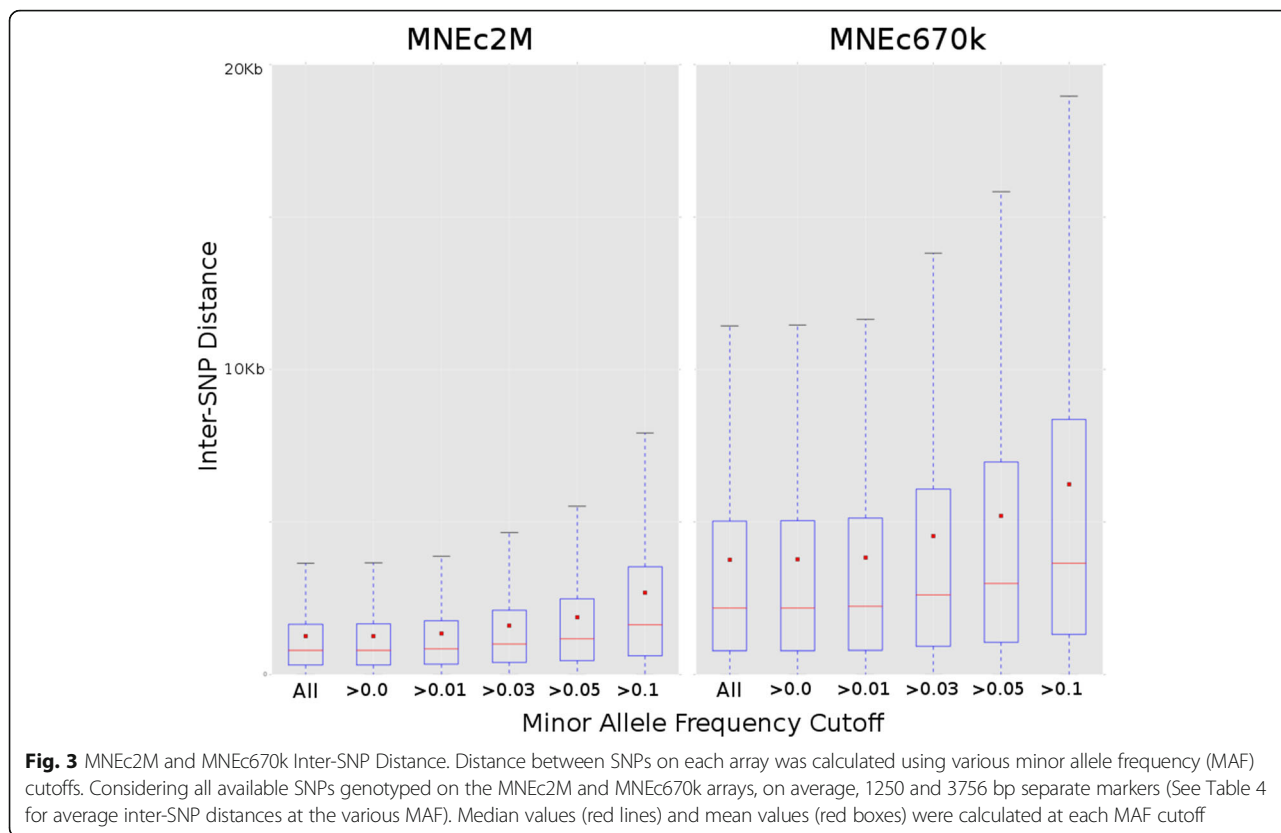
Alternate allele frequency

Frequency of alternate (non-reference) SNP alleles were calculated on both arrays using the full 485 sample dataset (WGS + SNP), genotypes derived from whole genome sequence only (WGS Only), and from samples genotyped on the 2M test array described above (SNP Only; see Additional file 5: Table S5). Kernel density estimations (KDE) of the alternate allele frequency of SNPs on the MNEc2M array showed a mean frequency between 0.20 and 0.28 (See Fig. 4 and Table 5). In general, regardless of genotyping source, alternative allele frequencies shared similar distributions. Frequency distributions exhibit long tails indicating substantial numbers of samples with non-reference alleles. Alternate allele (ALT) distributions also exhibit a lower median

Table 3 Imputation accuracy of the MNEc670k SNP genotyping array

Tagging Group	670 K to 2 M	Num Imputed Samples
Land Race	0.981 +/- 0.001	3
Arabian	0.993 +/- 0.0009	13
Belgian	0.992 +/- 0.0005	7
Draft	0.9658 +/- 0.0140	6
F. Montange	0.988 +/- 0.0017	10
Icelandic	0.989 +/- 0.0012	6
Lusitano	0.992 +/- 0.0004	7
Maremmano	0.994 +/- 0.0005	9
Morgan	0.988 +/- 0.0036	20
Pony	0.991 +/- 0.0025	19
Quarter Horse	0.983 +/- 0.0033	25
Standardbred	0.989 +/- 0.0045	13
Thoroughbred	0.991 +/- 0.0068	9
Warmblood	0.985 +/- 0.0116	9
Trotter	0.9857 +/- 0.0046	9

Breed specific imputation accuracy (mean +/- s.e.m.) of genotypes from MNEc670k to MNEc2M SNP sets. In each tagging breed group, 1/3 of samples genotypes were masked to lower density SNP sets and removed from the reference population of 485 horses. Imputation was performed using Beagle 4.0 and concordance was determined with VCFtools



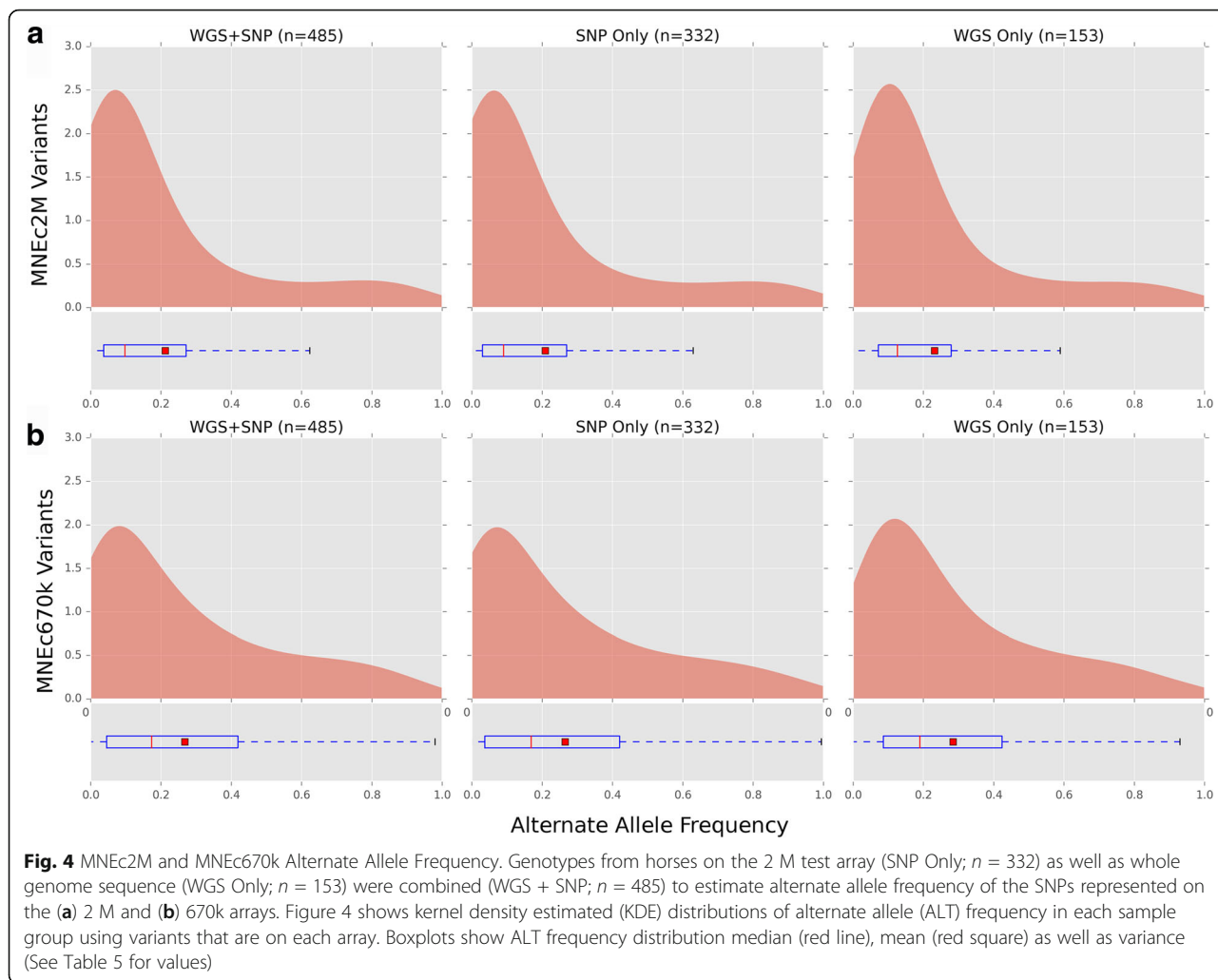
frequency (red line) than mean frequency (red bar), a common property of right skewed distributions [34]. Alternate allele frequency distribution was also broken down by breed for both the MNEc2M array (Fig. 5) and the MNEc670k array (Additional file 13: Figure S5) using genotypes derived from the WGS + SNP sample set. Allele frequency is balanced across breeds, though there

were minute differences. For example, the median MAF for Thoroughbreds was 3-12% lower than all other breeds, though this is not unexpected given the reference genome is a Thoroughbred. Despite minor differences, all breeds had long tails and similar distributions of allele frequencies indicating a balanced and representative SNP selection for GWAS.

Table 4 MNEc2M and MNEc670k Inter-SNP distance at various minor allele frequency cutoffs

Chip	MAF	Mean InterSNP Distance	Median InterSNP distance	Number of SNPs at MAF
MNEc2M	All SNPs	1250	785	1,986,984
	MAF > 0	1255	787	1,978,913
	MAF > 0.01	1334	835	1,862,844
	MAF > 0.03	1590	991	1,562,205
	MAF > 0.05	1876	1162	1,324,205
	MAF > 0.10	2676	1623	928,235
MNEc670k	All SNPs	3756	2172	661,349
	MAF > 0	3768	2178	659,278
	MAF > 0.01	3837	2226	647,481
	MAF > 0.03	4534	2606	547,858
	MAF > 0.05	5199	2980	477,719
	MAF > 0.10	6240	3651	398,055

Inter-SNP distance was calculated between SNPs informative at minor allele cutoffs greater than 0, 0.01, 0.03, 0.05 and 0.10. The number of SNPs included at this MAF cutoff is included. Distance and informativeness was re-calculated on both MNEc2M and MNEc670k arrays which were further broken down by tagging breed group (See Additional file 5: Table S5)



Linkage disequilibrium decay by breed

Linkage disequilibrium was measured using genotype r^2 between all 1.8M SNPs within 1 Mb of one another within and across breeds. LD across breeds (i.e., the WGS + SNP sample sets) (see Fig. 6, SNP and WGS curves) decayed faster than LD within any given breed (remaining curves). Within-breed calculations demonstrated that Quarter Horses and the Pony breeds had the lowest LD between SNPs at long distances, decaying to

below 0.10 at 1 Mb, while Thoroughbreds had the highest LD at all distances considered.

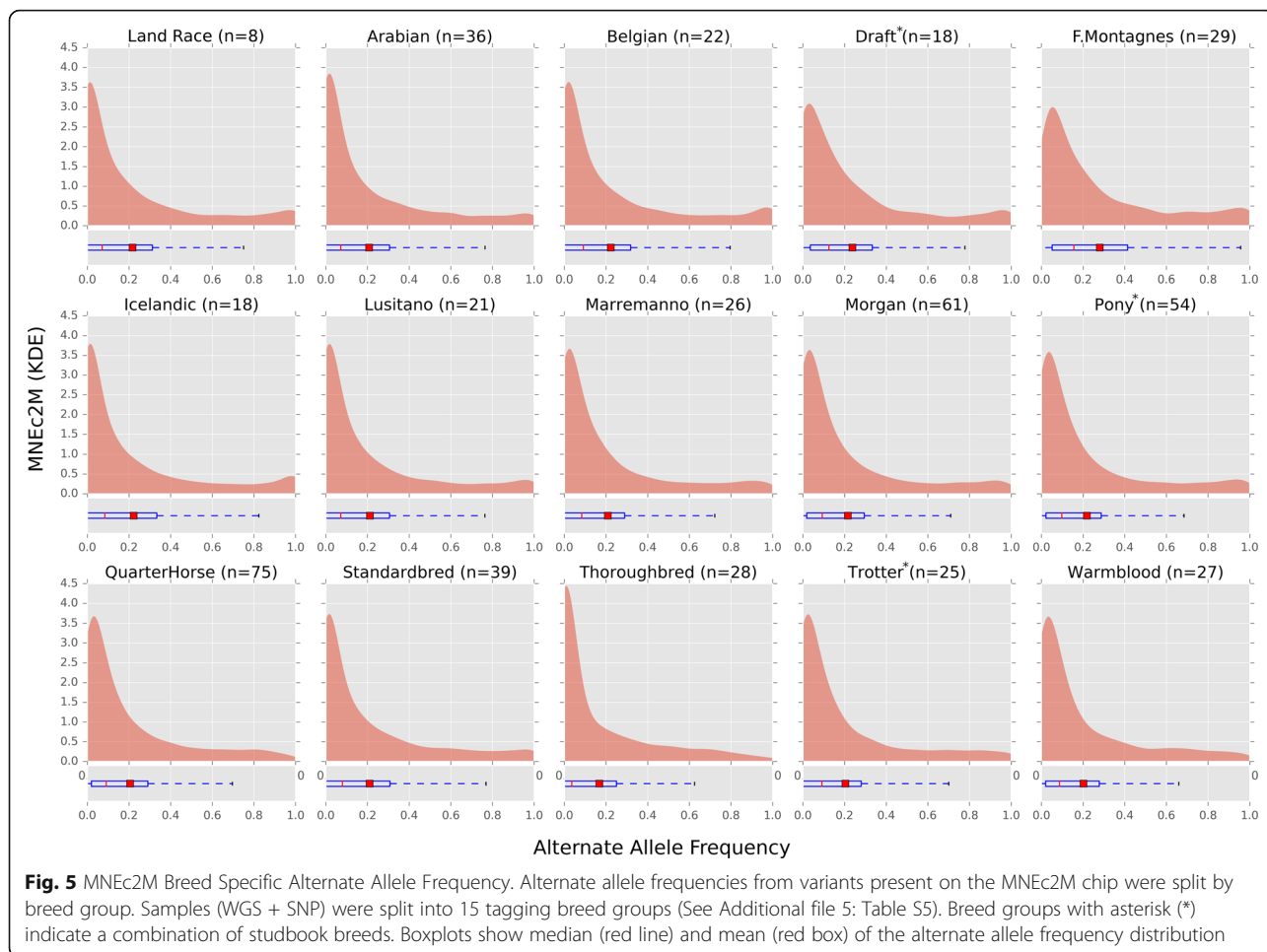
Discussion

Our goal was to provide a high-quality, standardized SNP array designed for imputation to overcome limitations in SNP density that under-power many genome mapping projects in the horse. To do this, we utilized whole genome sequencing data from 156 horses representing actively-

Table 5 MNEc2M and MNEc670k variant mean and median alternate allele frequency

SNP Chip	Sample Split	Mean ALT Allele Frequency	Median ALT Allele Frequency
MNEc 2M Variants	WGS + SNP	0.2115920561	0.0975609756
	SNP Only	0.208498773	0.0900621118
	WGS Only	0.2313047405	0.1258278146
MNEc 670k Variants	WGS + SNP	0.267280485	0.1729559748
	SNP Only	0.2647029832	0.1682389937
	WGS Only	0.2836436744	0.1895424837

Average (mean and median) values for MNEc2M and MNEc670k arrays broken down by genotype information available from WGS, CHIP or WGS + CHIP



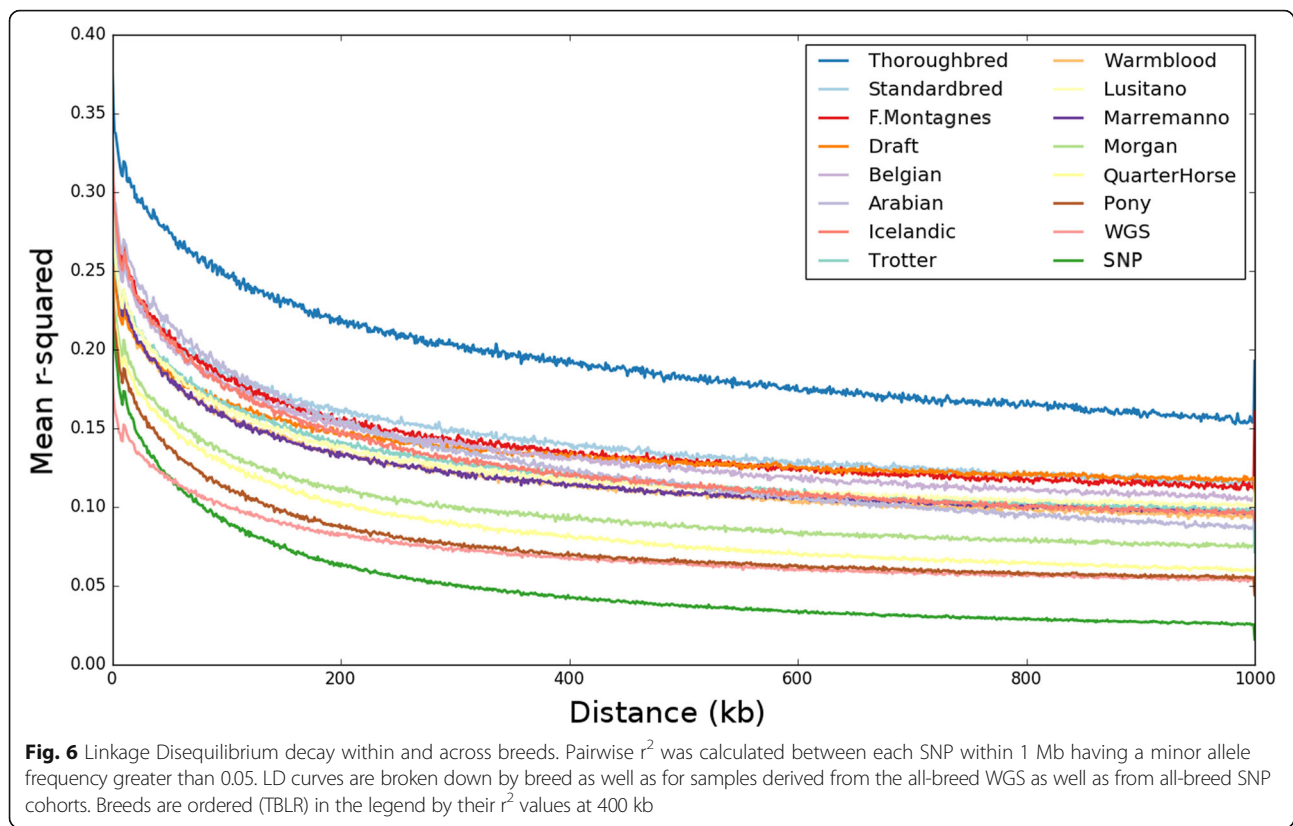
researched and economically-important breeds that were collaboratively collected by 17 laboratories within the equine genetics community. These data, in turn, enabled a large scale, whole genome sequence mapping pipeline with *de-novo* variant calling resulting in a starting set of over 23 million variants. Variant filtering was based on genome coverage, breed representation, linkage disequilibrium, and feasibility of array probe design. This resulted in the successful design of a fully backwards-compatible, high-density, 2 million SNP genotyping array (MNEc2M). Using a test cohort of 485 horses, we identified “tag” SNPs that reconstructed haplotypes both across diverse breeds and within breeds, then selected a subset of ~670k SNPs for design of the commercially available MNEc670k array. The use of tagging SNPs in this commercial array ensures its utility as an imputation tool up to a SNP density of at least 2 million.

Several successful GWAS in other agricultural animal species have been reported using this same Affymetrix® Axiom® HD genotyping array technology [35–37]. High-density ~ 670k genotyping arrays were used in domestic cattle to identify structural variation [38, 39]. Other

studies have leveraged this technology to discover and refine loci associated with production traits [40–42]. Based on reports from other domestic animals, coupled with preliminary analyses performed here, we anticipate similar performance and increased power for map-based studies to soon follow in the horse.

Assembly correction of the equine genome

Initial filtering for array compatible SNPs immediately reduced the pool from ~23 million SNPs that were discovered from de novo variant calling to 10 million SNPs compatible with array design. These purely technical constraints, such as removing all variants that were in close proximity to one another, substantially reduced the amount of variants considered here. While we thoroughly characterized the SNPs that were ultimately included on the genotyping array, the entire 23M SNP set has been submitted to dbSNP and the European Variation Archive. Additional analyses are being performed to further enhance knowledge of equine genetic diversity at the WGS level. This includes investigation of the >400,000 sites where all WGS alleles differed from the reference allele



determined by Sanger sequencing of Twilight. If the discrepancies in these data are verified, this information will be useful for further annotation and error correcting of the current equine genome reference assembly as well as for future genome assemblies. Here, we focused on variants that were compatible with array design, however, we anticipate that improved genome annotation and finer scale haplotype maps will result from this larger SNP set.

Determining highly precise, gold standard SNP sets

Quality control and evaluation of conversion rates on the previous arrays have shown that many of the SNPs assayed on the legacy arrays are non-informative or not polymorphic in many breeds (See Additional file 2: Table S2). These SNPs could truly occur at low frequencies within those populations, could have been singleton mutations in the individuals included in the legacy SNP discovery panel, or may have been false positive sites where no true genomic variation occurs. Compared to the SNP discovery panel used in the legacy arrays, the cohort of 153 horses used here provided a valuable opportunity to identify a precise set of SNPs that would be informative across many breeds.

While variant discovery from WGS had substantially higher overall depth of coverage than the SNP50 discovery effort, many breeds still had a modest number of individuals or relatively low average depth of coverage

(Additional file 1: Table S1). Still, genotype comparisons at variant sites from WGS to legacy arrays showed an overall high congruency rate (94% at 6X and 99% at 12X; Fig. 1), however, many SNPs with high QUAL score did not agree between WGS and legacy arrays. Since we were unable to determine the source of error in genotype mismatches between the WGS discovered variants and the legacy SNPs, especially those with high QUAL scores, we used additional information such as how many variant calling groups a SNP was discovered in as well as sequencing read level information to determine the criteria for including WGS discovered SNPs on the MNEc2M and MNEc670k arrays.

Distinguishing marginally higher quality variants, even at overall validation rates of over 94% becomes important when selecting a small fraction of variants to be included in a commercial array. Selecting a final subset (670k) containing only 6% of the initial ~10 million array-compatible variants posed an opportunity to control for overall false discovery rate during SNP selection. Variant recalibration allowed successful identification of variants with high QUAL scores that were less discordant between WGS and SNP50 genotypes. In the absence of high-confidence training SNP sets (i.e., dense HapMap data), we defined several different “gold standard” SNP sets for variant recalibration. While variant recalibration performed better than QUAL scores alone, there was not a single training set

that clearly out-performed the others (Fig. 2). However, using recalibrated VQSLOD scores coupled with a focus on identifying high precision SNPs, did result in a high conversion rate of probes (92.3%) on the 2M test array. As more WGS data are generated in more horses allowing for further comparisons between SNPs genotyped using different technologies, and SNPs discovered here can be further validated, the ideal SNP training set for variant recalibration will become more formalized.

Genotyping success in blood versus hair root DNA

DNA samples genotyped on the MNEc2M array were primarily derived from blood, but also came from hair roots. Both DNA sources had failed samples, however, a substantially higher fraction of hair root samples produced poor genotyping rates. Hair root DNA also tended to have lower DNA concentrations when samples were re-hydrated. To maximize information, we chose to genotype all submitted samples, even though samples from hair roots did not meet minimal DNA quantity guidelines specified by Affymetrix, though samples were dropped if they did not meet genotyping quality control thresholds.

To determine if sample origin (blood versus hair roots), DNA quantity, or both, were associated with failure to pass genotyping quality control metrics, a logistic regression for sample success on DNA concentration (Pico Green) and blood/hair status was performed (Additional file 14: Figure S6). It is clear that higher DNA concentrations increased the probability of genotyping success. However, while blood samples were more likely to produce passing samples regardless of DNA concentration, at adequate concentrations, hair root samples were also highly likely to produce passing genotypes.

We also noted variation in genotyping quality of specific SNPs based on tissue origin. During array design, SNPs were not disqualified that potentially performed better in blood versus hair root, e.g. 'PolyHighResolution' in one and not the other (Additional file 4: Table S4). While it is difficult to determine if certain probes perform better using DNA isolated from one tissue versus the other, in the samples tested here, DNA isolated from blood was preferred over hair root when DNA concentration is low or questionable.

Setup for precision imputation

Although the legacy equine arrays were not designed for imputation, several previous studies have demonstrated their utility in imputing genotypes, which, if performed on a larger scale across many breeds, could greatly improve the chances of success in horse genome mapping studies. A study in Standardbreds, Quarter Horses and

Thoroughbreds showed high fidelity imputation from legacy arrays (54K and 65k) to a higher marker density (74k) using a reference population of 248 horses [3]. Another study found that genotype imputation in Thoroughbreds was feasible from a very low density (1-3 K markers) to a legacy SNP set (70 K), although it was impacted by the minor allele frequency of the SNPs being imputed, as well as potentially complex LD structure [43].

We masked variants genotyped on the MNEc2M array down to the MNEc670k SNP set in a test cohort of 485 horses and found a high overall imputation accuracy (96-99%) up to the 2M density, across several different breeds of horses (See Table 3). It is important to note that the 96% average imputation accuracy in the Draft horses was mainly due to underperformance (92.6%) in the Percheron horses. This is not a surprising result as Percherons were underrepresented in the Draft breed group (2 out of 18 individuals) thus removing both horses from the imputation reference population for validation yielded poor results. The inclusion or exclusion of samples in the reference population can significantly impact imputation accuracy. Here, it was critical to maximize the available data during the SNP discovery process, especially for those breed groups that had limited sample representation.

While our imputation scenarios focused on imputing to the MNEc2M SNP set, imputation to higher SNP densities (>2M) is feasible. A recent study using whole genome sequencing from 44 Franches-Montagnes and Warmbloods imputed SNPs from the legacy SNP50 genotypes to nearly 13 M variants with 95% accuracy. WGS breed representation in this study varied between 1 and 29 individuals. In the future, it will be critical to expand the number of WGS samples in the reference population described here in order to attain proper breed representation for reliable imputation to higher SNP densities.

As sample size continues to plague the characterization of complex equine phenotypes, despite falling genotyping costs, the development of this affordable, backwards compatible, 670k imputation array will allow the inclusion of already genotyped individuals in higher-powered analysis, providing an economical trade-off for studies that must choose between more markers or more samples. We expect that imputation will complement future GWAS studies and mapping studies designed based on haplotype structure. In the future, improvements in imputation software as well as development of equine-specific imputation protocols that leverage breed information (e.g. recombination maps) will further increase accuracy of genotype imputation. Furthermore, as additional WGS are integrated into reference populations, imputation performance will continue to improve.

Conclusions

Here we report, through a community effort, the leveraging of WGS data from 153 individuals with a 13X median coverage to achieve new variant discovery in the domestic horse. The use of WGS enabled us to generate orders of magnitude more data than previous technologies allowed. Empirical SNP properties were used to produce composite scores for each SNP based on machine learning approaches, allowing better detection of false positive variant calls that could not be achieved using generic QUAL scores alone. Thus, from a starting set of over 23 million SNPs, we identified a set of ~5 million SNPs which were considered for array design. With probe design recommendations from Affymetrix, we further filtered this list to 2M SNPs (MNEc2M), which adequately represented the breeds used in variant discovery, were evenly spaced across the genome, and had the highest chance of conversion in the array implementation.

Using a test cohort of 332 horses, we used the MNEc2M SNP set to identify haplotype tagging SNPs both within and across breeds. Choices made in 2M array design, particularly with regard to variant filtering, resulted in greater than 92% of the SNPs on the 2M array (MNEc2M) returning high-quality genotypes. Filtering further, we designed the MNEc670k genotyping array to contain SNPs which allowed for accurate imputation. Together these genotyping platforms represent the next generation of genomic array technologies for the domestic horse.

Methods

Whole-genome sequencing and mapping pipeline

Paired-end Illumina Hi-Seq 100 base pair whole genome sequences were generated for 153 horses representing 24 distinct breeds (Additional file 1: Table S1). Raw reads for each individual were mapped using the PALEOMIX pipeline and aligned to an extended EquCab version 2.0 genome (see below). Specifically, reads for each individual were processed separately by sequencing lane and filtered for quality control using the program Adapter Removal [44], which removed PCR adapters, trimmed low quality base-pairs, and collapsed overlapping paired-end reads into a single high-quality read, and the resulting reads were filtered by length. Passing reads were mapped to each reference file using BWA [45]. Paired-end reads for which the mate was filtered were mapped in single-end mode. PCR duplicates were detected and removed, the resulting bam-files were merged, and reads were realigned around detected indels.

Extended EquCab 2.0 reference genome

An extended version of the EquCab2 reference genome was used which included the 31 autosomal chromosomes, ECAX, and equine chromosome unknown 1 (ChrUn1)

[1], together with an additional 7850 contigs designated as ChrUn2 generated from unmapped Twilight genomic DNA reads which were de novo assembled using the Velvet assembler [46]. These additional contigs were required to meet any of the following criteria: 1) longer than 1000 base pairs with no BLAST alignment (bit score > 99) to the human, canine, or bovine genomes; 2) longer than 1000 base pair and either a single BLAST alignment to a human, canine, or bovine chromosome or multiple alignments that mapped to a single chromosome; 3) between 500 and 999 base pairs and a BLAST alignment to a single human, canine, or bovine chromosome with a bit score > 499; 4) between 500 and 999 bp with alignment to a single chromosome where the total coverage of the aligned region included more than 80% of the coding length of an existing human, canine, or bovine annotated protein-coding gene.

PALEOMIX *vcf* pipeline implementation and python source code

Programs within PALEOMIX are abstracted as nodes within the program to allow files to be generated within a temporary directory and only to be moved to the final directory upon successful completion. In addition to nodes for the read alignment mapping programs, additional nodes were created to run variant calling programs implemented by GATK [26] and SAMtools [27]. Additional PALEOMIX nodes were created to perform variant recalibration as well as to assess precision versus recall for Morgan and Standardbred breed groups (Figs. 1 and 2). Source code for the extended PALEOMIX nodes are available at <https://github.com/schae234/pypipeline>.

Variant calling and validation PALEOMIX nodes

The PALEOMIX computational framework was further extended to process alignments and produce variant calls. Individuals were split into 16 different calling groups based on both breed and sequencing depth in order to minimize biases due to coverage and population stratification. Variant call files (.vcf) were produced for each group using both GATK Unified “Genotyper with” “–stand_call_conf” set to 30, “–stand_emit_conf” set to 10, and “–dcov” set to 200. Variants called by GATK were retained if they were independently called by SAMtools “mpileup” with default calling parameters. This allowed possible false negatives to be assessed downstream using other quality metrics. Of the 26,884,885 SNPs called by SAMtools and 31,506,364 SNPs called by GATK, 22,557,988 variants were called independently by both callers and retained for further analysis.

Variant recalibration (VQSLOD scores)

Both GATK UnifiedGenotyper as well as SAMtools assign generic quality scores (QUAL) to each discovered

variant, which is the posterior probability that a true variant exists given the pileup of reads at a given locus using base pair quality and expected allelic distribution of samples. Using GATK, additional metrics were generated for each variant including coverage (DP), quality of depth (QD), Fisher strand bias (FS), mapping quality rank sum (MQRankSum) and read position rank sum (ReadPosRankSum). Definitions from the GATK manual for each metric are below.

Coverage (DP) – Total, unfiltered depth of coverage.

Quality by Depth (QD) – Variant confidence (from QUAL field) divided by depth of non-reference samples.

Fisher Strand Bias (FS) – Measure of strand bias, i.e. the variation seen on only the forward or reverse strand.

Mapping Quality Rank Sum (MQRankSum) – The rank sum test for mapping qualities.

Read Position Rank Sum (ReadPosRankSum) – The rank sum test for the distance of the variant from the end of reads.

These parameters were used to train linear mixed models using GATK VariantRecalibrator which produces a composite variant quality metric called a VQSLOD score which can be summarized as using the formula $VQSLOD \sim DP + QD + FS + MQRankSum + ReadPosRankSum$. We defined three “gold standard” datasets: 1) SNPs previously genotyped on the 54K SNP chip, 2) variants called in four or more calling groups, 3) and variants within the 99th percentile of QUAL scores. Models were trained on chromosomes 2-31 and validated on chromosome 1 in both Morgan and Standardbred breed groups as a subset of each breed group had whole genome sequence data at target depths of 6X (12 Morgans and 6 Standardbreds) and 12X (6 Morgans and 4 Standardbreds) target coverage (actual coverage MOR6X: 5.23-6.96X, mean = 6.11X; STD6X: 4.61-5.89X, mean = 5.29X; MOR12X: 10.42-15.14X, mean = 12.96X; STD12X: 11.21-11.99X, mean = 11.63X). These “gold standard” training SNPs were used to assess the precision of different quality metrics of de novo variant calling.

Preliminary 5 M high-density SNP selection

A preliminary list of approximately 5 million SNPs was prepared for probe design (conversion) recommendation by Affymetrix. VQSLOD scores were generated for all ~23 million initial bi-allelic candidate SNPs. SNPs were filtered out if they 1) fell within repetitive regions designated by RepeatMasker 3.3.0 [47] (780 SNPs); 2) fell within 20 bp of another variant (9,342,733 SNPs); and 3) had fewer than 2 observed instances of minor allele (4,858,273). Of the remaining 11,435,936 SNPs, further filtering was applied to obtain equal coverage throughout the genome. Pairwise LD was calculated between SNPs within 10 kb windows. If SNPs within a window had an r^2 value of over 0.90, priority was given to SNPs

that were in called in Draft or Pony groups. If a SNP was in an equal number of priority groups, VQSLOD scores were used to break ties. After all filtering criteria were applied 5,443,950 SNPs remained.

High-density 2M SNP selection

From the ~5 million candidate SNPs, Affymetrix provided four classes of recommendations based on the probability the SNP would be convertible through probe design (‘recommended’, neutral, ‘non-recommended’, ‘not possible’). A recommended SNP has: a probability of conversion >0.60, no interfering polymorphisms within 24 bases, and unique flanking sequence. Non-recommended SNPs have: either duplicate flanking sequence, a probability of conversion <0.40, interfering polymorphisms within 21 bases, or more than 2 interfering polymorphisms within 24 bases. A ‘not possible’ designation is given to probes which probes cannot be created, and neutral recommendations cover all other cases. Furthermore, SNPs with alleles of A/T or C/G required two probes to differentiate between allele states and were tagged as allele-specific SNPs. Recommendations were generated for both forward and reverse strands based on the above criteria. Groups of SNPs were ranked within 50 kb windows and probe design criteria were chosen using a greedy algorithm until 37 SNPs were chosen per window using the following criteria: 1) VIP SNPs previously designed on 54/65k chip (regardless of recommendation or strand), 2) SNPs for known Mendelian traits such as coat color (regardless of recommendation or strand), 3) SNPs designated as ‘recommended’ by Affymetrix and designable with one probe, 4) SNPs with a ‘neutral’ recommendation from Affymetrix designable with one probe, 5) any SNP within the equine MHC region, 6) SNPs requiring multiple probes. If a SNP was equally designable in forward and reverse strand, the forward strand was chosen for interpretability. Using these criteria, 2,001,826 SNPs were chosen for array design.

2M SNP test array and sample quality control

Using probes designed according to the above criteria, 347 horses from 20 breeds were genotyped on the 2M array using DNA isolated from blood ($n = 286$) using the Genra PureGene Blood kit. DNA from hair roots ($n = 61$) was isolated using the Genra Puregene Blood Kit with a modified version of the Genra Puregene Mouse Tail purification protocol. The amount of Proteinase K was increased to 20 μ L, isopropanol to 650 μ L, and ethanol to 500 μ L. The DNA hydration solution was reduced to 20-30 μ L depending on size of DNA pellet. Quality control metrics were calculated on arrays grouped by tissue type (blood vs hair) as well as array batches, and arrays were dropped for various reasons at multiple QC steps. “DishQC” evaluates call rates between

A/T probes and C/G probes of non-polymorphic sites to assess background probe contrast and was calculated using Affymetrix Power Tools. Samples with DQC scores below 0.60 were removed from further evaluation. Passing arrays were then genotyped using approximately 20,000 high-confidence probes provided by Affymetrix in their R1 release package (See <https://www.thermofisher.com/order/catalog/product/550583#/legacy=affymetrix.com>). If the number of arrays per tissue/array group was above 96 samples, generic priors were used for probe-set genotyping; otherwise, SNP specific priors were used (provided in the R1 package released by Affymetrix). Sample arrays with 20 K call rates below 0.97 were removed from further analysis.

Remaining samples ($n = 332$) were genotyped on all probes using a similar approach. Generic or SNP specific priors (provided in R1 package) were used for groups containing more or less than 96 samples, respectively, and samples with call rates below 0.97 were removed from the analysis. Successfully genotypes samples were then grouped by genotyping plates to check for batch effects. Plates with pass rates below 95% in blood or 93% in hair were dropped.

Up to 4 probes were used to genotype each SNP across samples passing quality control criteria. Probe performance was calculated using the *'Metrics.R'* script provided by Affymetrix which assessed several quality criteria. Twelve criteria were used to assign each SNP into 6 categories representing probe conversion types. High-quality probes fell in one of three categories in descending order of quality: 'PolyHighResolution' SNPs had good cluster resolution and at least two examples of the minor allele, 'MonoHighResolution' SNPs had good SNP clustering but less than two samples with the minor allele, and 'NoMinorHom' had good cluster resolution but no samples with homozygous, minor alleles. Poor quality SNPs were qualified as having off target variation (OTV), a call rate below threshold, or a combination of poor performing cluster properties. SNPs with multiple probes were assigned a best probe based on the above high-quality conversion types with the added constraint that SNPs with no minor homozygote did not result in extreme Hardy-Weinberg values ($p \leq 10e-5$; Chi-Squared Test). 1,846,988 SNPs passed the above quality control metrics.

670k commercial array SNP selection

Genotype information derived from the horses on the 2M test array ($n = 332$) were combined with SNPs in horses called by whole genome sequencing ($n = 153$) resulting in a total of 485 horses genotyped at 1,846,988 SNPs. Tag SNPs were calculated using FastTagger [31] using two different tagging scenarios. Tagging SNPs

informative across populations (inter-population) were identified by running FastTagger on the entire dataset. Parameters provided to FastTagger identified inter-population tag-SNPs at a resolution of down to 0.01 minor allele frequency and up to 0.99 r^2 .

To identify SNPs tagging population specific haplotypes (intra-population), samples were split into 15 tagging breed groups based on available sample size: Land Race, Arabian, Belgian, Draft, Franches-Montagne, Icelandic, Lusitano, Maremmano, Morgan, Pony, Quarter Horse, Standardbred, Thoroughbred, Trotter, and Warmblood (See Additional file 5: Table S5). FastTagger was run with parameters to differentiate tag-SNPs within each population at a resolution of 0.10 minor allele frequency and 0.90 r^2 . Representative SNPs from each population specific set of tag SNPs were collapsed using the program 'Multi-Pop-TagSelect' which assesses overlap between sets of tag SNPs and identifies the subset of SNPs which near-minimally spans the set of populations [32]. Intra-breed specific SNPs were kept as long as they were tagged haplotypes in five or more breed groups.

Quarter Horses, Ponies, Morgans, and Standardbreds needed a much higher number of intra-breed specific SNPs to tag haplotypes, indicating they were the most diverse breeds. Additional tag SNPs ($n = 13,993$) were included in the final commercial array if they tagged haplotypes in three or more of these diverse breed groups.

Genotype imputation between the MNEc670k and MNEc2M arrays

A reference population of the 485 horses with genotypes at ~2M SNPs (generated either by the MNEc2M array or whole-genome sequence) was used as a reference population for imputation. Horses from each different breed were used to test the imputation accuracy of the array. A random set of 1/3 of individuals were removed from each breed group and masked down to the MNEc670k SNP set as well as the SNP65 SNP set. Genotypes were then imputed using the remaining (non-masked) individuals as a reference population. Imputation was performed using Beagle 4.0 [48] using default parameters. Genotype concordance was calculated using VCFtools (0.1.15) using the '-diff-indv-discordance' option [49]. Briefly, concordance, as calculated by VCFtools, is the proportion of non-missing SNPs in both datasets that have the same allele calls. Missing genotypes are not considered in the concordance calculation and phase is not taken into account (e.g. A/T and T/A are concordant). VCFtools reports the proportion of matching genotypes per individual to the number of sites that are shared across input SNP sets.

Additional files

Additional file 1: Table S1. Whole genome sequencing (WGS) samples. Whole genome sequencing samples with horse identification and read depth. 153 Individuals representing 24 breeds (Twilight is included as 4 entries). Table includes horse identifier, breed, contributing laboratory and coverage statistics for each individual in both the nuclear and mitochondrial genomes. Variant calling groups indicate which individuals were grouped together during variant discovery. (XLSX 21 kb)

Additional file 2: Table S2. SNP50 breed informativeness. Impact of breed on SNP50 informativeness. Table shows informativeness ($MAF \geq 0.05$) of the legacy, SNP50 SNP set for light, draft, and pony breeds described by McCue et al. [1, 2]. Mean number of informative SNPs per breed was ~42,000 (~77%). 19,427 SNPs were informative in all 14 breeds. (XLSX 8 kb)

Additional file 3: Table S3. SNP validation by discovery by variant calling group. SNPs discovered from whole genome sequence (See calling groups in Additional file 1: Table S1) were validated by assessing the minor allele frequency in individuals genotyped on the MNEc2M test array. (XLSX 8 kb)

Additional file 4: Table S4. MNEc670k SNP information. Contains information on MNEc670k SNPs. Includes Affymetrix probe ID, genetic coordinates, MNEc Identifier, tissue origin indicating high quality genotyping in either hair root or blood, Affymetrix assigned genotyping cluster resolution, and inclusion criteria: Intra – SNP tags intra-breed haplotype, Inter – SNP tags inter-breed haplotype, MHC – SNP within Equine Major Histocompatibility Complex region, Diverse – SNP tags haplotype in 3 or more of diverse breeds that need many tagging SNPs (Quarter Horse, Pony, Morgan, Standardbred; See Additional file 4: Table S4), Density – SNP included to target at least 8 SNPs per 50 kb genomic window (targeting uniform genomic coverage). (XLSX 32510 kb)

Additional file 5: Table S5. Sample breed and tagging breed groups. Samples ($n = 485$) from either whole genome sequence (WGS) or from the 2M test array (SNP) were assigned to one of 15 tagging breed groups based on their reported breed. (XLSX 21 kb)

Additional file 6: Position of variants discovered from WGS data. File contains chromosome, bp position, distance from previous discovered SNP and alternate allele frequency for the 23 million SNP set. (TSV 828110 kb)

Additional file 7: MNEc2M SNP information. Contains information for SNPs included on the MNEc2M array. Columns include MNEc identifiers, EquCab 2.0 coordinates, and fields indicating if the SNP was discovered in Draft and Pony groups (which were under-represented on the legacy arrays, See Additional file 2: Table S2). (CSV 82434 kb)

Additional file 8: Breed specific tagging SNPs. File containing informative tagging SNPs broken down by breed. (CSV 172641 kb)

Additional file 9: Figure S1. Alternate allele frequency for 23 M and 10 M SNP sets. Histograms show the alternative allele frequency for the ~23 million SNPs discovered by WGS and the ~10 million SNPs that were compatible with array design (see Table 1). A high number of SNPs were observed at 100% alternate allele frequency in a sample cohort that included deep sequencing (see Additional file 1: Table S1) of Twilight (the horse used for the reference genome) indicating likely Sanger sequence errors. (PNG 87 kb)

Additional file 10: Figure S2. The effects of minor allele on imputation accuracy. SNPs were binned by minor allele count of the marker in the reference panel for each tagging breed group. SNPs were masked down to the 670k set then imputed up to ~2M. For SNPs in each minor allele count bin, a Pearson correlation was calculated between the imputed minor allele dosage and the true minor allele dosage. The x-axis in each panel are on a log scale to show the relationship for low minor allele count SNPs. (PNG 264 kb)

Additional file 11: Figure S3. MNEc2M Inter-SNP distance and informativeness by breed. Distance between SNPs on the MNEc2M array were calculated by breed group (See Additional file 5: Table S5) at various minor allele frequency cutoffs (0, 0.01, 0.03, 0.05, 0.10). Breed groups with asterisk (*) indicate a combination of studbook breeds. (PNG 1152 kb)

Additional file 12: Figure S4. MNEc670k Inter-SNP distance and informativeness by breed. Distance between SNPs on the MNEc670k array were calculated by breed group (Additional file 5: Table S5) at various minor allele frequency cutoffs (0, 0.01, 0.03, 0.05, and 0.10). Breed groups with asterisk (*) indicate a combination of studbook breeds. (PNG 1117 kb)

Additional file 13: Figure S5. MNEc670k Breed Specific Alternate Allele Frequency. Alternate allele frequencies from variants present on the MNEc670k chip were split by breed group. Samples (WGS + SNP) were split into 15 tagging breed groups (See Additional file 5: Table S5). Breed groups with asterisk (*) indicate a combination of studbook breeds. (PNG 814 kb)

Additional file 14: Figure S6. Logistic regression of DNA Sample success. A logistic regression was fit, predicting the probability of a sample passing quality control using DNA concentration (Pico Green) and sample source as independent variables. (PNG 60 kb)

Abbreviations

ALT: Alternate allele; bp: Base pair; ChrUn1: Equus caballus unknown chromosome; ECA: *Equus caballus* chromosome; GATK: Genome analysis tool kit; Gb: giga-bases; KDE: Kernel density estimation; LD: Linkage disequilibrium; MAF: Minor allele frequency; MHC: Major histocompatibility complex; MNEc2M: the 2 million SNP array developed here; MNEc670k: the 670 thousand SNP array developed here; MOR: the Morgan horse breed; QUAL: Generic quality score output by GATK and SAMtools; SNP: Single Nucleotide Polymorphism; STD: the Standardbred horse breed; VIP: Very important probe; VQSLOD: Variant quality score log-odds; WGS: Whole genome sequencing

Acknowledgements

These genomic tools are the result of a collaborative effort that involved a number of laboratories within the equine genetics community, many of which provided support or insight at various steps for this community resource. We would like to acknowledge Chad Dow from Affymetrix for the support he and his team provided while designing the MNEc2M array.

Funding

Support for the generation of whole genome sequence came from the following sources:

- USDA NIFA project 2012-67,015-19,432 and Minnesota Agricultural Experiment Station Multistate project MIN-62-090.
- The National Animal Genome Project (NRSP8) through the equine genome coordinator: USDA-NRSP8 (2013-2018) horse-technical-committee coordinator funds.
- The Danish Council for Independent Research, Natural Sciences (Grant 4002-00152B); the Danish National Research Foundation (Grant DNRF94); Initiative d'Excellence Chaires d'attractivité, Université de Toulouse (OURASI), and; the European Research Council (ERC-CoG-2015-681,605).
- The Bavarian Ministry State Ministry for Food and Agriculture, and Forestry (A/13/39).
- The Laboratory of Molecular Evolution, The Koret School of Veterinary Medicine, The Hebrew University of Jerusalem, Israel) for contributing pure-bred Arabian whole-genomes on behalf of The Israel Science Foundation (ISF) grant #1365/10.
- The Swedish Research Council Formas (221-2013-1661) and the Swedish Research Council VR (621-2012-4666).

Funding sources played no role in the design of this study or the collection, analysis, and the interpretation of data and in writing the manuscript.

Availability of data and materials

Whole genome sequences are available in the following NCBI BioProjects: PRJEB14779, PRJNA273402, and PRJEB10098. Additional sequences have been restricted in availability due to pre-existing material transfer agreements and can be requested by contacting the contributing investigator in Additional file 1: Table S1. Genotypes for horses on the MNEc2M array will be released upon publication. Genome positions for all 23 million discovered SNPs have been submitted to dbSNP as well as the European Variation Archive.

Authors' contributions

Experimental concept and design: MEM, JRM, LO, RJS, MS; Sample collection and data contribution: EB, DLB, EB, GKB, GB, SAB, OD, RF, CJF, VG, BH, VJ, TK, TL, GL, MSL, NM, ACM, JNM, AM, JM, CP, SP, SR, IT, JT, GT, AVS, CMW, BW, LO, JRM, MEM; Data analysis and interpretation: RJS, MS, MEM; Computational support: RJS, MS; Manuscript writing and figures: RJS, JRM, MEM; Manuscript review: All authors read and approved the final manuscript.

Ethics approval and consent to participate

DNA samples were previously collected with approval from the Animal Care and Use Committees at the respective institutions. All animal work was performed in accordance and with approval from international and national governing bodies at the institutions where the samples were collected (University of Minnesota Institutional Animal Care and Use Committee (IACUC); University of California, Davis Institutional Animal Care and Use Committee (protocol #17491); University of Kentucky Institutional Animal Care and Use Committee (IACUC); Ethics Committee for Animal Experiments in Uppsala, Sweden (Number C121/14); Institutional animal care and use committee at Cornell University (protocol 2008-0121); University of California, Davis IACUC 19205; Hebrew University's approval number AG-23476-07; Institutional Animal Care and Use Committee (IACUC), the Lower Saxony state veterinary office- registration number 11A 160/7221.3-2.1-015/11, 8.84-02.05.20.12.066; University of Sydney Animal Ethics Committee: AEC APPROVAL NUMBER: N00/9-2009/3/5109; permit no. BE75/16, veterinary service of the Canton of Bern; Institutional ethics committee of the University of Veterinary Medicine Vienna Good Scientific Practice guidelines and national legislation; Italian Ministry of Agricultural, Food and Forestry Policies (Mipaaf); Ethical Committee of the Canton of Bern (BE33/07, BE58/10 and BE10/13) No commercial animals were used in this study. Written informed client consent describing the purpose and duration of the study, procedures, potential risks and benefits and containing study contact information were obtained from private owners.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interests.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Author details

¹Department of Veterinary Population Medicine, College of Veterinary Medicine, University of Minnesota, St. Paul, MN, USA. ²Centre for GeoGenetics, Natural History Museum of Denmark, University of Copenhagen, Copenhagen, Denmark. ³Maxwell H. Gluck Equine Research Center, Department of Veterinary Science, University of Kentucky, Lexington, KY, USA. ⁴School of Veterinary Medicine, University of California-Davis, Davis, CA 95616, USA. ⁵Unité de Génétique Animale et Biologie Intégrative-UMR1313, INRA, Université Paris-Saclay, AgroParisTech, 78350 Jouy-en-Josas, France. ⁶The Robert H. Smith Faculty of Agriculture, Food and Environment, The Robert H. Smith Faculty of Agriculture, Food and Environment, The Koret School of Veterinary Medicine, The Hebrew University, 76100 Rehovot, Israel. ⁷Institute of Animal Breeding and Genetics, Department of Biomedical Sciences, University of Veterinary Medicine Vienna, Vienna, Austria. ⁸Department of Animal Science, University of Florida, Gainesville, FL, USA. ⁹Institute for Animal Breeding and Genetics, University of Veterinary Medicine, Hannover, Germany. ¹⁰Lehrstuhl für Tierzucht der Technischen Universität München, Liesel-Beckmann-Strasse 1, 85354 Freising, Germany. ¹¹Swiss Institute of Equine Medicine, Department of Clinical Veterinary Medicine, Vetsuisse Faculty, University of Bern, and Agroscope, Länggassstrasse 124, 3001 Bern, Switzerland. ¹²School of Life and Environmental Sciences, Faculty of Veterinary Science, University of Sydney, Regimental Drive, B19-301 RMC Gunn, Sydney, NSW 2006, Australia. ¹³Institute of Genetics, University of Bern, 3001 Bern, Switzerland. ¹⁴Department of Biochemistry and Molecular Biology, School of Medicine, University of Louisville, Louisville, KY 40202, USA. ¹⁵Department of Animal Breeding and Genetics, Swedish University of Agricultural Sciences, Uppsala, Sweden. ¹⁶Biotechnology Centre of Azores, University of Azores, Angra do Heroísmo, Portugal. ¹⁷Department of Veterinary Clinical Medicine, College of Veterinary Medicine, University of Illinois at Urbana-Champaign, Champaign,

IL 61802, USA. ¹⁸Veterinary Genetics Laboratory, University of California Davis, Davis, CA, USA. ¹⁹Agroscope, Swiss National Stud Farm, 1580 Avenches, Switzerland. ²⁰Institute of Animal Breeding and Husbandry, Christian-Albrechts-University Kiel, Hermann-Rodewald-Strasse 6, 24098 Kiel, Germany. ²¹Department of Animal Sciences, Functional Breeding Group, Georg-August University Göttingen, Burckhardtweg 2, 37077 Göttingen, Germany. ²²Department of Veterinary Medicine - Sport Horse Research Centre, University of Perugia, Perugia, Italy. ²³Laboratoire d'Anthropobiologie Moléculaire et d'Imagerie de Synthèse, CNRS UMR 5288, Université de Toulouse, Université Paul Sabatier, 31000 Toulouse, France. ²⁴Department of Veterinary and Biomedical Sciences, College of Veterinary Medicine, University of Minnesota, St. Paul, MN, USA.

Received: 20 February 2017 Accepted: 13 July 2017

Published online: 27 July 2017

References

- Wade CM, Giulotto E, Sigurdsson S, Zoli M, Gnerre S, Inslund F, et al. Genome sequence, comparative analysis, and population genetics of the domestic horse. *Sci*. 2009;326:865–7.
- McCue ME, Bannasch DL, Petersen JL, Gurr J, Bailey E, Binns MM, et al. A high density SNP array for the domestic horse and extant Perissodactyla: utility for association mapping, genetic diversity, and phylogeny studies. Georges M, editor. *PLoS Genet*. Public Library of Science; 2012;8:e1002451.
- McCoy AM, McCue ME. Validation of imputation between equine genotyping arrays. *Anim Genet*. 2014;45:153.
- Schubert M, Jónsson H, Chang D, Der Sarkissian C, Ermini L, Ginolhac A, et al. Prehistoric genomes reveal the genetic foundation and cost of horse domestication. *Proc Natl Acad Sci*. National Academy of Sciences; 2014;111:201416991.
- Jónsson H, Schubert M, Seguin-Orlando A, Ginolhac A, Petersen L, Fumagalli M, et al. Speciation with gene flow in equids despite extensive chromosomal plasticity. *Proc Natl Acad Sci U S A*. National Academy of Sciences; 2014;111:18655–60.
- McCoy AM, Beeson SK, Splan RK, Lykkjen S, Ralston SL, Mickelson JR, et al. Identification and validation of risk loci for osteochondrosis in standardbreds. *BMC Genomics*. BioMed Central. 2016;17:41.
- McQueen CM, Dindot S V, Foster MJ, Cohen ND. Genetic susceptibility to Rhodococcus equi. *J Vet Intern Med*. Wiley-Blackwell; 2015;29:1648–59.
- Hauswirth R, Haase B, Blatter M, Brooks SA, Burger D, Drögemüller C, et al. Mutations in MITF and PAX3 cause "splashed white" and other white spotting phenotypes in horses. Barsh GS, editor. *PLoS Genet*. Public Library of Science; 2012;8:e1002653.
- Hill EW, McGivney BA, Gu J, Whiston R, Machugh DE. A genome-wide SNP-association study confirms a sequence variant (g.66493737C>T) in the equine myostatin (MSTN) gene as the most powerful predictor of optimum racing distance for thoroughbred racehorses. *BMC Genomics*. BioMed Central; 2010;11:552.
- Lykkjen S, Dolvik NI, McCue ME, Rendahl AK, Mickelson JR, Roed KH. Genome-wide association analysis of osteochondrosis of the tibiotarsal joint in Norwegian Standardbred trotters. *Anim Genet*. 2010;41(Suppl 2):111–20.
- Raudsepp T, McCue ME, Das PJ, Dobson L, Vishnoi M, Fritz KL, et al. Genome-wide association study implicates testis-sperm specific FKBP6 as a susceptibility locus for impaired Acrosome reaction in stallions. Barsh GS, editor. *PLoS Genet*. Public Library of Science; 2012;8:e1003139.
- Lykkjen S, Dolvik NI, McCue ME, Rendahl AK, Mickelson JR, Røed KH. Equine developmental orthopaedic diseases—a genome-wide association study of first phalanx plantar osteochondral fragments in Standardbred trotters. *Anim Genet*. 2013;44:766–9.
- Signer-Hasler H, Flury C, Haase B, Burger D, Simianer H, Leeb T, et al. A genome-wide association study reveals loci influencing height and other conformation traits in horses. *PLoS One*. Public Library of Science; 2012;7:e37282.
- Corbin LJ, Blott SC, Swinburne JE, Sibbons C, Fox-Clipham LY, Helweggen M, et al. A genome-wide association study of osteochondritis dissecans in the thoroughbred. *Mamm Genome*. Springer-Verlag; 2012; 23:294–303.
- Finno CJ, Stevens C, Young A, Affolter V, Joshi NA. SERPINB11 Frameshift Variant Associated with Novel Hoof Specific Phenotype in Connemara Ponies. *PLoS Genetics*. Public Library of Science; 2015;1–17.
- Kader A, Li Y, Dong K, Irwin DM, Zhao Q, He X, et al. Population variation reveals independent selection toward small body size in Chinese Debao pony. *Genome Biol. Evol.* Oxford University Press; 2016;8:42–50.

17. Petersen JL, Mickelson JR, Cothran EG, Andersson LS, Axelsson J, Bailey E, et al. Genetic diversity in the modern horse illustrated from genome-wide SNP data. *PLoS One*. Public Library of Science; 2013;8:e54997.
18. Schubert M, Ermini L, Der Sarkissian C, Jónsson H, Ginolhac A, Schaefer R, et al. Characterization of ancient and modern genomes by SNP detection and phylogenomic and metagenomic analysis using PALEOMIX. *Nat Protoc*. 2014;9:1056–1082.
19. Orlando L, Ginolhac A, Zhang G, Froese D, Albrechtsen A, Stiller M, et al. Recalibrating Equus evolution using the genome sequence of an early middle Pleistocene horse. *Nature*. 2013; 499:74–8.
20. Sarkissian C, Der, Ermini L, Schubert M, Yang MA, Librado P, Fumagalli M, et al. Evolutionary Genomics and Conservation of the Endangered Przewalski's Horse. *Current Biology*. 2016;25:2577–83.
21. Librado P, Der Sarkissian C, Ermini L, Schubert M, Jónsson H, Albrechtsen A, et al. Tracking the origins of Yakutian horses and the genetic basis for their fast adaptation to subarctic environments. *Proc Natl Acad Sci*. 2015;112: 201513696.
22. Metzger J, Tonda R, Beltran S, Agueda L, Gut M, Distl O. Next generation sequencing gives an insight into the characteristics of highly selected breeds versus non-breed horses in the course of domestication. *BMC Genomics*. 2014;15:562.
23. Doan R, Cohen ND, Sawyer J, Ghaffari N, Johnson CD, Dindot SV. Whole-genome sequencing and genetic variant analysis of a quarter horse mare. *BMC Genomics*. 2012;13:78.
24. Andersson LS, Larhammar M, Memic F, Wootz H, Schwochow D, Rubin C-J, et al. Mutations in DMRT3 affect locomotion in horses and spinal circuit function in mice. *Nature*. 2012;488:642–6.
25. Frischknecht M, Neuditschko M, Jagannathan V, Drögemüller C, Tetens J, Thaller G, et al. Imputation of sequence level genotypes in the Franches-Montagnes horse breed. *Genet Sel Evol*. 2014;46:63.
26. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A, et al. The genome analysis toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res*. 2010;20:1297–303.
27. Li H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics*. 2011;27:2987–93.
28. Raghavan V, Bollmann P, Jung GS. A critical investigation of recall and precision as measures of retrieval system performance. *ACM Trans Inf Syst*. 1989;7:205–29.
29. Saito T, Rehmsmeier M. The precision-recall plot is more informative than the ROC plot when evaluating binary classifiers on imbalanced datasets. *PLoS One*. 2015;10:1–21.
30. DePristo MA, Banks E, Poplin R, Garimella K V, Maguire JR, Hartl C, et al. A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat Genet*. *Nature Research*; 2011;43:491–8.
31. Liu G, Wang Y, Wong L. FastTagger: an efficient algorithm for genome-wide tag SNP selection using multi-marker linkage disequilibrium. *BMC Bioinformatics*. 2010;11:66.
32. Howie BN, Carlson CS, Rieder MJ, Nickerson DA. Efficient selection of tagging single-nucleotide polymorphisms in multiple populations. *Hum Genet*. 2006;120:58–68.
33. Browning BL, Browning SR. Genotype imputation with millions of reference samples. *Am J Hum Genet*. *The American Society of Human Genetics*; 2016; 98:116–26.
34. von Hippel PT. Mean, Median, and Skew: Correcting a Textbook Rule. *J. Stat. Educ*. *American Statistical Association*. 2005;13.
35. Kranis A, Gheyas AA, Boschiero C, Turner F, Yu L, Smith S, et al. Development of a high density 600K SNP genotyping array for chicken. *BMC Genomics*. 2013;14:59.
36. Groenen MAM. Development of a high-density Axiom® porcine genotyping array to meet research and commercial needs. *Plant Anim. Genome XXIII Conf*. San Diego, CA: Plant & Animal Genome XXIII Conference; 2015.
37. Rincon G, Weber KL, Eenennaam A L Van, Golden BL, Medrano JF. Hot topic: performance of bovine high-density genotyping platforms in Holsteins and jerseys. *J Dairy Sci*. Elsevier; 2011;94:6116–21.
38. Salomón-Torres R, González-Vizcarra VM, Medina-Basulto GE, Montaña-Gómez MF, Mahadevan P, Yaurima-Basaldúa VH, et al. Genome-wide identification of copy number variations in Holstein cattle from Baja California, Mexico, using high-density SNP genotyping arrays. *Genet Mol Res*. 2015;14:11848–59.
39. Salomon-Torres R, Villa-Angulo R, Villa-Angulo C. Analysis of copy number variations in Mexican Holstein cattle using axion genome-wide Bos 1 array. *Genomics Data*; 2016;7:97–100.
40. Romé H, Varenne A, Hérault F, Chapuis H, Alleno C, Dehais P, et al. GWAS analyses reveal QTL in egg layers that differ in response to diet differences. *Genet Sel Evol*. *BioMed Central*; 2015;4:783.
41. Lu D, Akanno EC, Crowley JJ, Schenkel F, Li H, De Pauw M, et al. Accuracy of genomic predictions for feed efficiency traits of beef cattle using 50K and imputed HD genotypes. *J Anim Sci*. 2016;94:1342–53.
42. Li G, Li D, Yang N, Qu L, Hou Z, Zheng J, et al. A genome-wide association study identifies novel single nucleotide polymorphisms associated with dermal shank pigmentation in chickens. *Poult Sci*. 2014;93:2983–7.
43. Corbin LJ, Kranis A, Blott SC, Swinburne JE, Vaudin M, Bishop SC, et al. The utility of low-density genotyping for imputation in the thoroughbred horse. *Genet Sel Evol*. *BioMed Central*; 2014;4:69.
44. Lindgreen S. AdapterRemoval: easy cleaning of next-generation sequencing reads. *BMC Res Notes*. 2012; 5:337.
45. Li H, Durbin R. Fast and accurate short read alignment with burrows-wheeler transform. *Bioinformatics*. 2009;25:1754–60.
46. Zerbino DR, Birney E. Velvet : Algorithms for de novo short read assembly using de Bruijn graphs. *Genome Research*. 2008;821–9.
47. Smit, AFA, Hubley, R & Green, P. RepeatMasker Open-4.0.2013-2015 <http://www.repeatmasker.org>.
48. Browning SR, Browning BL. Rapid and accurate haplotype phasing and missing-data inference for whole-genome association studies by use of localized haplotype clustering. *Am J Hum Genet*. 2007;81:1084–97.
49. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant call format and VCFtools. *Bioinformatics*. 2011;27:2156–8.

Submit your next manuscript to BioMed Central and we will help you at every step:

- We accept pre-submission inquiries
- Our selector tool helps you to find the most relevant journal
- We provide round the clock customer support
- Convenient online submission
- Thorough peer review
- Inclusion in PubMed and all major indexing services
- Maximum visibility for your research

Submit your manuscript at
www.biomedcentral.com/submit

