

1 **Best practices for analyzing imputed genotypes from low-pass sequencing in dogs**

2 Reuben M. Buckley¹, Alex C. Harris¹, Guo-Dong Wang^{2,3}, D. Thad Whitaker¹, Ya-Ping Zhang^{2,3},
3 and Elaine A. Ostrander^{1*}

4

5 ¹ Cancer Genetics and Comparative Genomics Branch, National Human Genome Research
6 Institute, National Institutes of Health, Bethesda, Maryland

7 ²State Key Laboratory of Genetic Resources and Evolution, Kunming Institute of Zoology,
8 Chinese Academy of Sciences, Kunming 650223, China

9 ³Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences,
10 Kunming 650223, China

11

12 Running title: Low-pass WGS and imputation in domestic dogs

13

14 *Correspondence: Elaine A. Ostrander, Ph.D., Cancer Genetics and Comparative Genomics
15 Branch, 50 South Drive, Building 50, Room 5351, Bethesda MD 20892. Phone: 3015945284;
16 Email: eostrand@mail.nih.gov

17

18

19 Keywords: Imputation, Low-pass Sequencing, Canine Genomics, Genome Sequencing

1 **Abstract**

2 Although DNA array-based approaches for genome wide association studies (GWAS) permit
3 the collection of thousands of low-cost genotypes, it is often at the expense of resolution and
4 completeness, as SNP chip technologies are ultimately limited by SNPs chosen during array
5 development. An alternative low-cost approach is low-pass whole genome sequencing (WGS)
6 followed by imputation. Rather than relying on high levels of genotype confidence at a set of
7 select loci, low-pass WGS and imputation relies on the combined information from millions of
8 randomly sampled low confidence genotypes. To investigate low-pass WGS and imputation in
9 the dog, we assessed accuracy and performance by downsampling 97 high-coverage (>15x)
10 WGS datasets from 51 different breeds to approximately 1x coverage, simulating low-pass
11 WGS. Using a reference panel of 676 dogs from 91 breeds, genotypes were imputed from the
12 downsampled data and compared to a truth set of genotypes generated from high coverage
13 WGS. Using our truth set, we optimized a variant quality filtering strategy that retained
14 approximately 80% of 14M imputed sites and lowered the imputation error rate from 3.0% to
15 1.5%. Seven million sites remained with a MAF > 5% and an average imputation quality score
16 of 0.95. Finally, we simulated the impact of imputation errors on outcomes for case-control
17 GWAS, where small effect sizes were most impacted and medium to large effect sizes were
18 minorly impacted. These analyses provide best practice guidelines for study design and data
19 post-processing of low-pass WGS imputed genotypes in dogs.

1 **Introduction**

2 The price per marker for a genotyping assay can have a large influence on the success of
3 genetic association studies. In dogs, DNA genotyping arrays, which provide hundreds of
4 thousands of genotypes at relatively low costs, are highly beneficial for mapping loci (Awano et
5 al. 2009; Hayward et al. 2016), characterizing genetic architecture (Boyko et al. 2010; Friedrich
6 et al. 2019), and defining breed and population structure (Shannon et al. 2015; Ali et al. 2020).
7 However, DNA genotyping arrays are limited by various known and unknown biases that occur
8 during marker selection and probe design that cannot be removed without redesigning a new
9 DNA array, which is an expensive and time-consuming process. An alternative similarly priced
10 approach is low-pass whole genome sequencing (WGS) and imputation (Martin et al. 2021).
11 Rather than assigning genotypes based on high confidence calls across a finite set of loci, low-
12 pass WGS combines information from millions of randomly sampled low confidence variant calls
13 to impute likely genotypes from a reference panel, comprised of a large collection of WGS
14 datasets representing potential haplotypes found within a population. Since low-pass WGS isn't
15 biased towards sampling specific loci, the only limiting factor is the reference panel used.
16 Therefore, the utility of previous datasets never diminishes.

17
18 Due to its flexibility and scalability, low-pass sequencing and imputation has been applied to
19 humans (Rubinacci et al. 2021; Wasik et al. 2021) and other mammalian species (Benjelloun et
20 al. 2019; Piras et al. 2020; Snelling et al. 2020; Nosková et al. 2021). Results in humans
21 demonstrate that low-pass WGS and imputation provide more accurate genotypes than those
22 imputed using array data, leading to increased power for genome wide association studies
23 (GWAS) and more accurate polygenic risk score calculation. Piras et al. (2020) used low-pass
24 WGS and imputation to identify candidate loci for canine idiopathic pulmonary fibrosis in West
25 Highland white terriers (CPSF7 and SDHAF2). While successful in this case, the existing 350

1 dog breeds present a unique problem for conducting GWAS studies, as the existing structure of
2 each breed, its history, and genome homogeneity are distinct (Ostrander et al. 2017). In the
3 absence of empirical evidence for developing optimal strategies for study design and data
4 processing, the probability of poor performance and misleading results is unknown. As many
5 dog breeds and populations have only been sequenced to low levels, the development of a
6 generalizable set of rules for low pass WGS imputation across breeds would convert much of
7 the existing data from low to high applicability, thus accelerating the dog as a genetic system for
8 studies of canine and human health.

9
10 Here, we present an analysis of imputation accuracy of low-pass WGS in the context of canine
11 genomics and establish optimized approaches for study design and data processing. We
12 analyzed imputed genotypes from 97 test samples from 58 different breeds, many of which are
13 not included in the reference panel containing the haplotypes used for imputation. We assessed
14 the impact of minor allele frequency (MAF) on genotyping accuracy and determined whether it
15 was better to use MAFs generated from imputation or to use reference panel MAFs. Finally, we
16 investigate the impact of imputation errors on study design by determining the necessary
17 sample sizes and case-control ratios for a sufficiently powered case-control GWAS.

18

19 **Results**

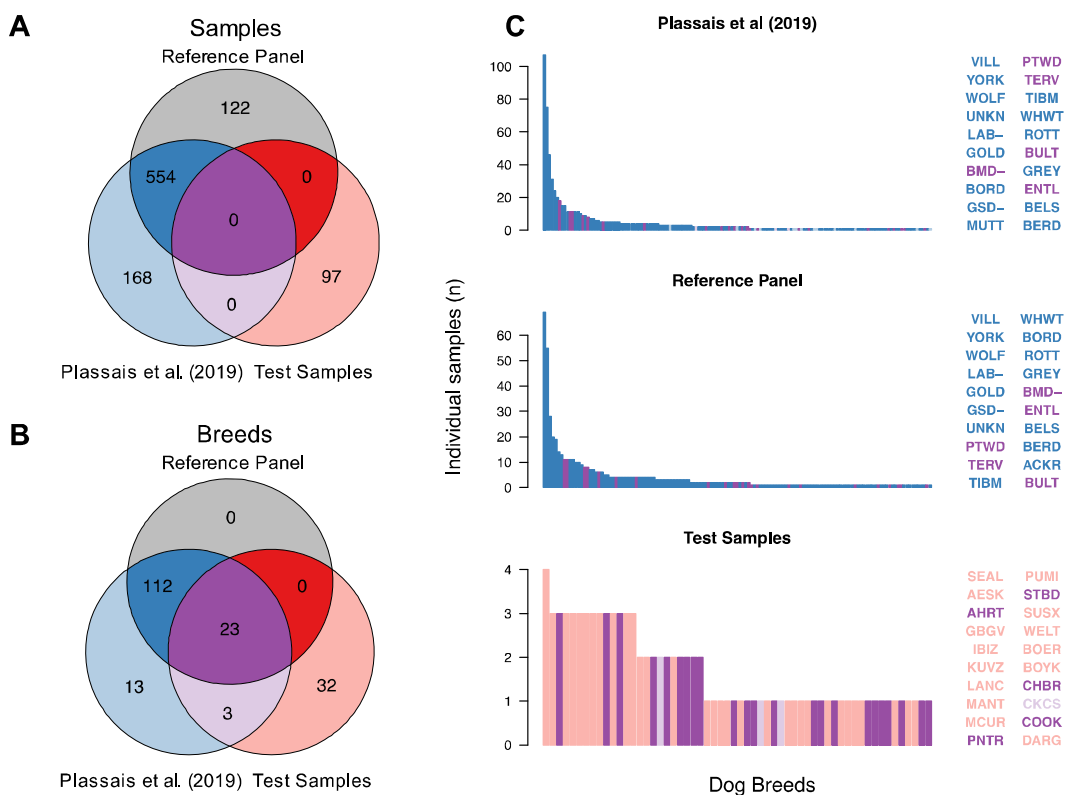
20 *Breed composition of the imputation reference panel and test datasets*

21 The test dataset used for assessing imputation accuracy consists of 97 samples that met
22 selection criteria (Methods). Of the 676 samples used in the reference panel, whose IDs were
23 provided by Gencove, Inc. (New York, NY), 554 were matched to a previously published dataset
24 **(Fig 1A) (Supplemental Table S1)** (Plassais et al. 2019). Since a large portion of samples were

1 shared between the Gencove reference panel and the published dataset, we opted to use the
2 publicly available VCF file as a stand-in for the reference panel VCF. The individual breeds
3 comprising the reference samples were compared to the breeds in the Plassais et al. (2019)
4 dataset and the breeds from our test dataset (**Fig 1B**). Only 13 breeds were identified as unique
5 to the published data set and they ranged from one to three members each (**Fig 1C**)
6 (**Supplemental Table S2**). The remaining breeds shared between the reference panel and the
7 published dataset typically contain similar numbers of individuals, with village dogs (VILL),
8 Yorkshire terriers (YORK), wolves (WOLF), Labradors (LAB-), golden retrievers (GOLD), and
9 unknown breeds (UNKN) being the six most popular breed designations in both datasets (**Fig**
10 **1C**). Within the test samples, 23 breeds were shared with the reference panel and four breeds
11 were shared with the Plassais et al. (2019) dataset only, while 32 breeds were unique to the test
12 samples (**Fig 1B**). In terms of member frequency per breed, the test samples had no more than
13 four Sealyham terriers, the most common breed within the test samples. In addition, the total
14 members per breed were relatively evenly distributed between breeds unique to the test
15 samples and breeds found in other datasets (**Fig 1C**).

16
17 To evaluate breed representation at a higher level, we determined the clade membership of all
18 known breeds across the test dataset and the reference panel dataset (Methods) (**Table 1**)
19 (**Supplemental Table S3**). Altogether, the Pinscher and Hungarian clades were unrepresented
20 in the reference panel. This is, perhaps, because each of these clades contains only two breeds
21 (Parker et al. 2017). Other clades of concern due to underrepresentation include the American
22 terrier, Asian toy, small spitz, and toy spitz clades, which have less than three representative
23 samples within the known breeds of the reference panel. Finally, there were 31 dogs from 14
24 breeds within our test samples with no previously assigned clade (**Supplemental Table S4**).
25 Since most of these breeds are either European or have recent known European ancestry, they

- 1 can likely be assigned to previously identified clades. Together, this data suggests the Plassais
- 2 et al. (2019) dataset likely represents many of the same haplotypes found in the reference panel
- 3 used for imputation, and that the test samples represent a mixture of shared and unique breeds,
- 4 appropriate for determining the impact of imputation on low-pass sequencing.



5

6 **Fig 1: Test samples belong to a wide variety of breeds with most breeds likely not found**

7 **within the imputation reference panel. A)** Sample membership within each dataset.

8 Reference panel IDs could not always be linked to a publicly available dataset. **B)** Breed

9 membership among each dataset. Reference panel dogs whose IDs could not be linked to a

10 publicly available sample have no breed label. **C)** Breed frequency across each dataset. Using

11 the colors from the Venn diagram in B, bar colors represent the population a specific breed can

12 be found in. Labels to the left of each bar chart identify the 20 most common breeds. Breeds in

13 bar charts are sorted by most to least common.

14

15

1 **Table 1:** Clade representation of reference and test datasets
2

Clades	Parker et al. (2017) ¹		Plassais et al. (2019) ²		Test Samples	
	Samples	Breeds	Samples	Breeds	Samples	Breeds
Alpine	26	3	20	4	5	3
American Terrier	16	3	1	1	4	2
American Toy	20	2	5	2	0	0
Asian Spitz	83	9	25	9	1	1
Asian Toy	44	5	2	2	2	2
Continental Herder	44	5	25	4	2	2
Drover	34	4	15	4	0	0
European Mastiff	139	16	24	11	7	5
Hungarian	9	2	0	0	3	1
Mediterranean	98	14	12	6	7	3
New World	45	7	24	7	3	2
Nordic Spitz	64	5	13	10	7	4
Pinscher	12	2	0	0	1	1
Pointer Setter	88	12	18	10	5	3
Poodle	72	8	20	6	3	2
Retriever	66	7	48	7	3	2
Scent Hound	71	8	14	6	1	1
Schnauzer	20	2	5	2	0	0
Small Spitz	14	2	1	1	1	1
Spaniel	44	5	14	5	2	1
Terrier	140	18	100	15	7	6
Toy Spitz	32	4	2	2	2	2
UK Rural	145	16	42	12	0	0
Unplaced* ³	20	2	7	4	0	0
Unknown breed* ⁴	0	-	13	-	0	-
Village Dogs* ⁵	0	-	69	-	0	-
Mix Breed* ⁶	0	-	6	-	0	-
Wild Canids* ⁷	9	2	29	2	0	0
No Clade Info* ⁸	0	0	0	0	31	14

- 3
4 ¹ Analysis that initially defined breed clade membership
5 ² Dogs in the reference panel that are also included in Plassais et al. (2019)
6 ³ Dog breeds that formed their own branch in previous phylogenetic analysis and were not a
7 member of a clade
8 ⁴ Dogs with no corresponding breed information
9 ⁵ Non-breed dogs sampled from 14 distinct geographic regions
10 ⁶ Dogs with mixed breed ancestry
11 ⁷ Group consists of grey wolves and golden jackals and breed label is used to differentiate these
12 two different species
13 ⁸ Breeds not included in previous phylogenetic analyses and therefore not assigned clade
14 membership

1 * Groups of samples that either do not form a monophyletic clade or have not been included in
2 previous phylogenetic analyses

3

4 *Downsampling and imputation*

5 The 97 test samples were each downsampled to a coverage level of 1x and underwent
6 imputation using loimpute as part of the Gencove, Inc. platform (Wasik et al. 2021). The
7 average read coverage of WGS variants and imputed variants was 17.5x and 1.06x,
8 respectively (**Supplemental Table S5**). A single sample, Pointer06, had a mean coverage of
9 1.67x, an outlier compared to next highest coverage dog which was 1.30x, suggesting
10 Pointer06's original coverage levels were incorrectly estimated. However, since variation in
11 coverage level is a potential outcome of low-pass sequencing, Pointer06 was retained for
12 further analyses. Imputation returned 53,649,170 million (M) variant sites, consisting of
13 35,875,925 SNV and 17,773,245 indel sites. Most sites were homozygous for the reference
14 allele across all samples and were therefore removed from the analysis, leaving a total of
15 14,845,499 SNVs and 7,946,973 indels. Alternatively, genotype calling of high coverage WGS
16 data for the test samples resulted in 18,476,517 SNVs and 12,831,692 indels.

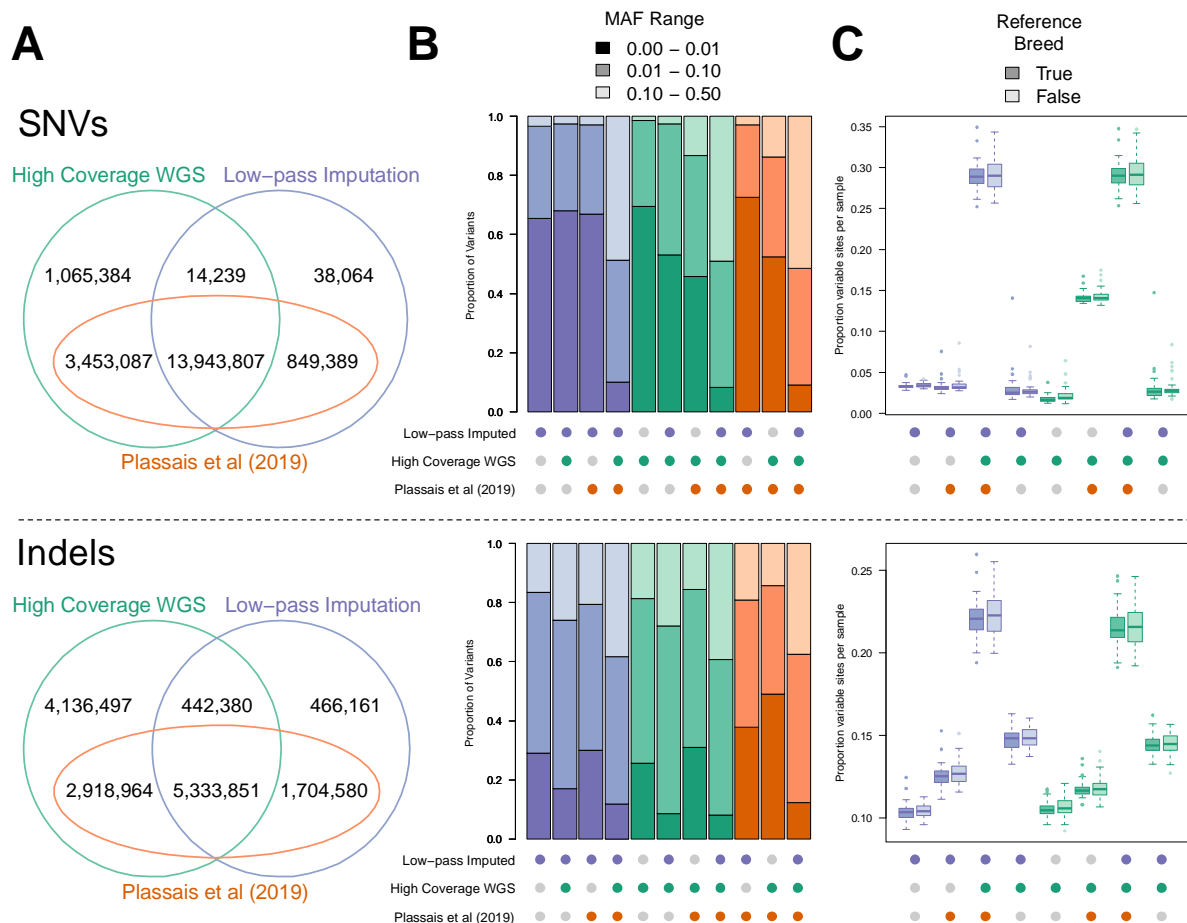
17

18 To analyze the accuracy of imputation, we identified variants shared between the following
19 groups: low-pass imputation, high coverage WGS, and the Plassais et al. (2019) dataset. The
20 majority of SNVs, 13,943,807, were shared between all three variant groups, which represented
21 93.9% of all imputed SNVs. Similarly, the majority of indels, 5,333,851, were also shared
22 between the three variant groups. However, shared indels represented only 67.1% of all
23 imputed indels, a smaller fraction than observed with SNVs. Importantly, 99.6% of all imputed
24 SNV sites and 88.6% of indel sites were present within the Plassais et al. (2019) dataset,
25 indicating its utility as a stand in for the reference panel (**Fig 2A**).

1
2 Although there were high levels of agreement between the low-pass imputation and the high
3 coverage WGS variant groups, many sites were specific to only one variant group. Since these
4 sites were removed from our downstream analysis, we measured their MAF distributions to
5 determine the extent of their potential impact on imputation accuracy measurements. Variant
6 group specific sites usually had MAFs < 0.01 , whereas variants shared across all groups usually
7 had higher MAFs (**Fig 2B**). This indicated that high coverage WGS specific sites were likely
8 individual specific sites, and absent from the reference panel used for imputation. Although a
9 large number of these sites are also found in the Plassais et al. (2019) dataset, they may belong
10 to breeds present in the test samples that were not present in the imputation reference panel.
11 Low-pass imputation specific sites are due to imputation errors, as these sites are imputed as
12 variable, but are actually genotyped as homozygous for the reference allele across all high
13 coverage WGS samples. Since these variants tend to present as rare alleles (MAF < 0.01) in
14 the imputed data, they are easily filtered out by MAF cutoffs typically used in association
15 analyses and, thus, have only a minor impact on analysis.

16
17 To determine whether breed composition of the reference panel has an impact on variant
18 imputation within out test samples, we counted the number of variant sites per individual for
19 each combination of the variant groups and compared the results between reference panel
20 breeds and non-reference panel breeds. Results showed that each carried a similar number of
21 variant sites for both SNVs and indels, indicating that the breed composition of the current
22 reference panel has little impact on variant imputation (**Fig 2C**). Ultimately, the total number of
23 variant sites per samples varied according to MAF distributions, as shown in **Fig 2B**.

24



1

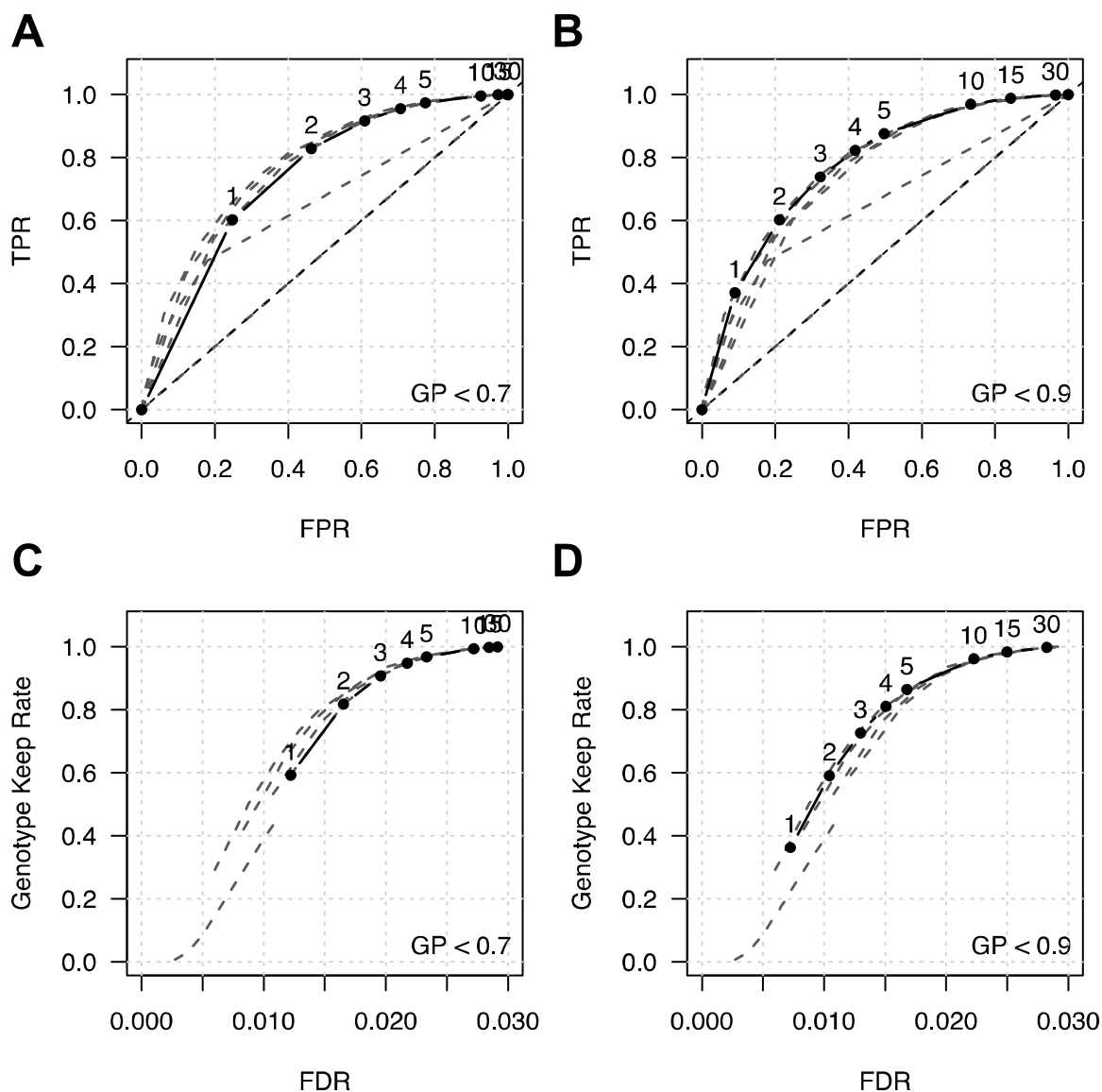
2 **Fig 2: Genomic variant positions and their corresponding alleles are consistent across**
3 **datasets.** **A)** Venn diagrams for SNVs and indels showing variants unique and shared across
4 datasets. Datasets include the high coverage WGS variant sites and low-pass imputed variant
5 sites found across the 97 test samples and variants discovered in Plassais et al. (2019).
6 Variants were identified as shared across datasets if the variant position, reference allele, and
7 alternate allele were identical. **B)** MAF distribution of each variant group from A. Variant groups
8 are indicated by colored circles beneath the bar chart. Groups contain variants which are the
9 intersect between the colored circles and the set difference between the grey circles. The colors
10 of each bar indicate the dataset used to calculate the MAF distribution and the shading level
11 indicates the relevant MAF range. **C)** Sites per sample in each variant group, where variant
12 groups are presented as in B. Sites per sample are measured as the proportion of total sites
13 within the relevant variant group that contain a non-reference allele for a particular sample.
14 Samples have also been divided into two groups based on whether the respective breed also
15 belongs to the Plassais et al. (2019) dataset and is therefore likely used in the imputation
16 reference panel.
17

1 *Filtering strategies to optimize accuracy*

2 Filtering strategies were optimized by analyzing the relationships between imputation accuracy,
3 genotyping confidence as determined by max genotype probability (GP), and low confidence
4 genotype thresholds (Methods) (**Supplemental Fig S1**). Receiver operator characteristic (ROC)
5 curves were calculated for seven GP thresholds and eight low confidence genotype thresholds
6 (**Supplemental Fig S2**). Each curve followed a similar trajectory. However, at a specific GP
7 threshold, low confidence genotype thresholds lead to different outcomes for true positive rates
8 (TPRs) and false positive rates (FPRs). For example, at $GP > 0.7$, a low confidence genotype
9 threshold of two is required for a $TPR > 0.8$ and a $FPR > 0.4$ (**Fig 3A**). By comparison, to meet
10 those same criteria at a threshold of $GP > 0.9$, a low confidence genotype threshold of four is
11 required (**Fig 3B**). These results indicate that at higher GP thresholds, a greater level of
12 robustness is achieved when selecting a low confidence genotype threshold, as small changes
13 in threshold values do not lead to large changes in the number of variants removed.

14
15 False discovery rate (FDR) and keep rate were also analyzed for the thresholds mentioned
16 above (Methods) (**Supplemental Fig S3**). This was done to assess how much data was lost by
17 filtering and to assess the number of imputation errors that remain within the dataset after
18 filtering. Similar to results shown in ROC curves, a higher GP threshold provided a higher level
19 of robustness for selecting a low confidence genotype threshold. For example, at $GP < 0.7$, a
20 low confidence genotype threshold of two leads to a keep rate > 0.8 and an approximate FDR of
21 0.015 (**Fig 3C**). To achieve a similar keep rate and FDR at $GP < 0.9$, a low confidence genotype
22 threshold of four is required (**Fig 3D**). These results show that higher GP thresholds are more
23 suitable for filtering sites based on the number of low confidence genotypes. This strategy
24 allows for fine tuning of filtered results without a loss in accuracy. For the remainder of our

- 1 analysis, we used a confidence filter of $GP < 0.9$ and a low confidence genotype threshold of
- 2 four, which roughly corresponds to a genotyping error rate of 5%.



- 3
- 4 **Fig 3: Performance of filtering strategies for reducing imputation errors. A)** ROC curve,
- 5 where genotypes with $GP < 0.7$ are identified as low confidence. Numbers above each point
- 6 represent low confidence rate thresholds for removing sites. Sites with a total number of low
- 7 confidence genotypes greater than or equal to the threshold are removed. Grey dashed lines
- 8 represent ROC curves for other confidence threshold values. **B)** ROC curve for confidence
- 9 threshold set at $GP < 0.9$. **C)** The proportion of variants remaining after filtering genotypes at
- 10 $GP < 0.7$ and the corresponding FDR. As in C, the numbers above each point represent low
- 11 confidence rate threshold values and grey dashed lines represent curves for other confidence
- 12 thresholds. **D)** Proportion of variants remaining and their corresponding FDR after filtering at
- 13 $GP < 0.9$.

1 *Minor allele frequency impacts imputation accuracy*

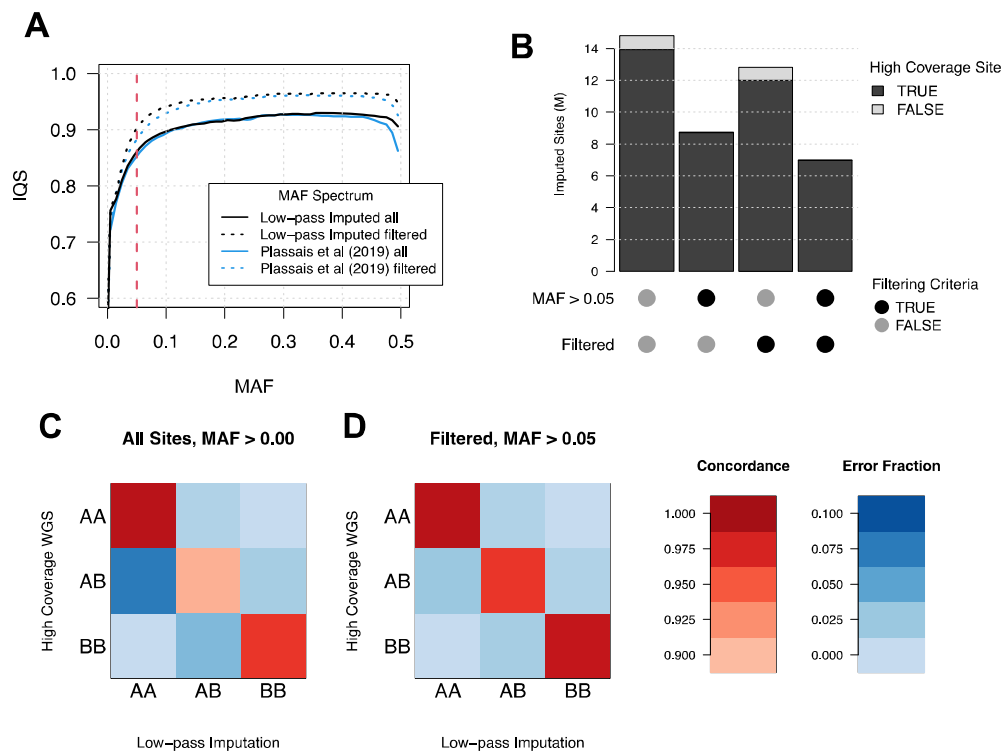
2 Imputation often performs poorly for rare alleles where statistical support is lacking. To
3 determine the impact of low MAFs on low-pass imputation in dogs, we analyzed imputation
4 accuracies using imputation quality score (IQS), a statistic that controls for allele frequencies by
5 taking chance agreement into account (Lin et al. 2010). In addition, population MAF estimates
6 were sourced from two different datasets, the imputed genotypes and the Plassais et al. (2019)
7 dataset. These were chosen as each is available for similar low-pass imputation analyses.

8
9 Both sources of MAFs expressed similar IQSs, where the largest differences are due to whether
10 imputed data was filtered. At MAFs > 0.1, the IQS of the unfiltered genotypes plateaued at
11 approximately 0.91, whereas filtered genotypes plateaued at approximately 0.95 (**Fig 4A**). At
12 MAFs < 0.1, the most accurate method for analyzing imputed genotypes is to use MAFs
13 generated from imputed genotypes. At a MAF of 0.05, only the filtered genotypes that were
14 partitioned according to imputed MAFs had an IQS > 0.9, providing the most accurate results for
15 low MAF genotypes.

16
17 Next, we analyzed the impact of MAF and filtering on imputation accuracy for different
18 genotypes. Imputation accuracy was poorest for heterozygous genotypes, especially at low
19 minor allele frequencies, indicating that heterozygous genotypes were least likely to be correctly
20 imputed. Conversely, homozygous reference imputation was most accurate at lower MAFs,
21 which is likely due to increased chance agreement between the high number of reference
22 genotypes (**Supplemental Fig. S4**). In addition, non-reference concordance, mean r^2 , and IQSs
23 were measured across chromosome 38 so that results can be compared to other analyses that
24 use a different accuracy measurement (**Supplemental Fig. S5**). All three measures show

1 higher levels of imputation accuracy at higher MAFs and a clear improvement for accuracy
2 measurements for filtered genotypes. Together these results demonstrate that removing sites
3 with a MAF < 5% retains a highly accurate set of genotypes, whereas filtering on GP values
4 increases overall imputation accuracy. In addition, we show imputation errors are more likely for
5 heterozygous sites, and MAF estimates derived from imputed genotypes are suitable for filtering
6 by MAF.

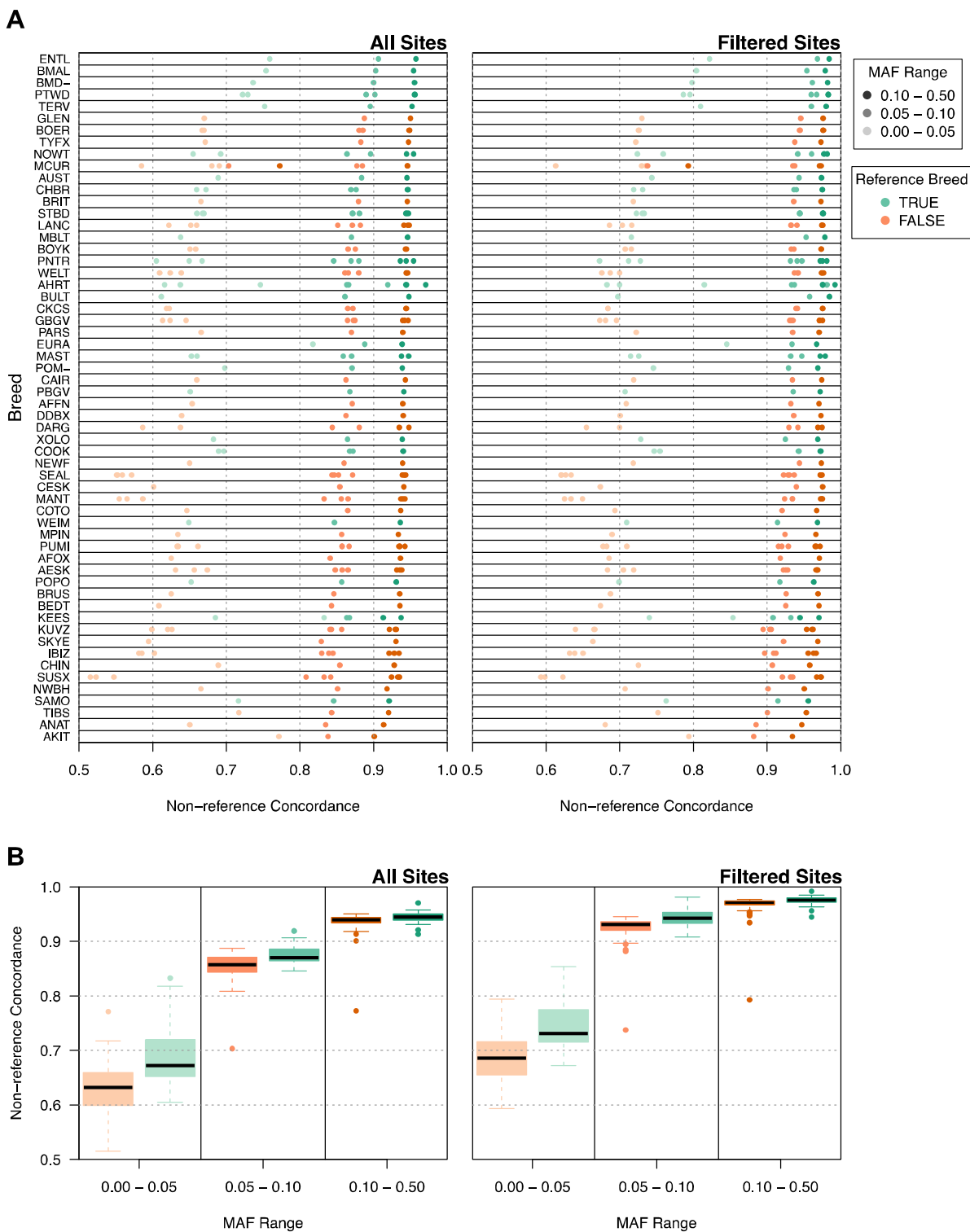
7
8 An important consideration when filtering variants is the total number of sites remaining.
9 Removing variants with a MAF < 0.05 and filtering out variants according to GP values results in
10 a total of 7M remaining sites. Most sites are removed due to the MAF cutoff rather than GP
11 filtering. In addition, almost all sites remaining after filtering correspond to a site found in the
12 high coverage WGS dataset, highlighting the quality and accuracy of the final dataset (**Fig 4B**).
13 An additional impact from filtering is the reduced variability between genotype specific
14 imputation errors and concordance rates. For example, prior to filtering the overall concordance
15 rate for high coverage WGS heterozygous genotypes with low-pass imputed genotypes was
16 0.908 with 0.072 of those imputed as homozygous reference, and 0.020 imputed as
17 homozygous alternate (**Fig 4C**). Conversely, after filtering, 0.964 of these genotypes were
18 imputed correctly, with approximately 0.023 genotypes imputed as homozygous reference and
19 0.013 genotypes as homozygous alternate (**Fig 4D**). Therefore, after filtering according to GPs
20 and MAFs, overall concordance rates are increased relative to unfiltered genotypes and
21 imputation errors are more evenly spread across the other genotypes.



1
 2 **Fig 4: Imputation accuracy according to minor allele frequency and genotype. A)**
 3 Imputation accuracy according to imputed and Plassais et al. (2019) MAFs for all sites and
 4 quality filtered sites. Imputation accuracy is measured as mean imputation quality score (IQS),
 5 an imputation accuracy statistic that accounts for the probability an allele is correctly imputed by
 6 chance. The red dotted line indicates a MAF of 0.05. **B)** The number of sites remaining after
 7 filtering for MAF > 0.05 and for low confidence genotypes < 5% as indicated by the “Filtered”
 8 label. Bar colors represent imputed sites that were either found or missing from the high
 9 coverage WGS dataset. **C)** Concordance and error rates for all genotypes, expressed as a
 10 fraction of the total number of high coverage WGS genotypes. **D)** Concordance and error rates
 11 for genotypes in sites with < 5% low confidence genotypes and MAFs > 0.05. Rates are
 12 expressed as a fraction of the number of high-coverage WGS genotypes that meet the
 13 corresponding filtering criteria.

14
 15 To determine whether the breed composition of the reference panel used for imputation impacts
 16 imputation accuracy, we measured non-reference concordance according to breed and MAF.
 17 Overall, breeds whose members showed the highest levels of imputation accuracy prior to
 18 filtering, such as the Entlebucher sennenhund, Belgian malinois, Bernese mountain dog,
 19 Portuguese water dog, and Belgian tervuren, also had members within the imputation reference
 20 panel (**Fig 5A**) (**Supplemental Table S6**). Importantly, four of these five breeds were among

1 the top 20 most highly represented breeds within the reference panel, with each containing at
2 least nine members each within the reference panel (**Fig 1C**). Alternatively, breeds with the
3 lowest levels of imputation accuracy usually did not contain members with the reference panel
4 (**Fig 5A**). Moreover, breeds with low imputation accuracy that did have members in the
5 reference panel, such as the Samoyed and Keeshond, had poorer representation than high
6 imputation accuracy breeds, with both breeds containing only three members each within the
7 reference panel (**Supplemental Table S2**). Furthermore, imputation accuracy rates of reference
8 panel breed members were consistently higher than non-reference panel breed members
9 across all MAF ranges (**Fig 5B**). These results indicate the importance of breed representation
10 for improving imputation accuracy.



1

2 **Fig 5: Imputation accuracy of dog breeds.** A) Individual dog breed imputation accuracy. Dog
 3 breeds are displayed on the Y axis with imputation accuracy on the X axis as non-reference
 4 concordance. Accuracy rates are displayed for all sites (left) and sites that remain after quality
 5 filtering (right). The shading of each data point indicates imputation accuracy of SNVs within a

1 specific MAF range. Green data points show breeds identified in the reference panel, while
2 orange points show breeds not found in the reference panel. Breeds are ranked according to
3 their median imputation accuracy for all sites. Imputation accuracies are displayed for each
4 member of the breed. **B)** Imputation accuracy of reference and non-reference breeds according
5 to MAF.

6

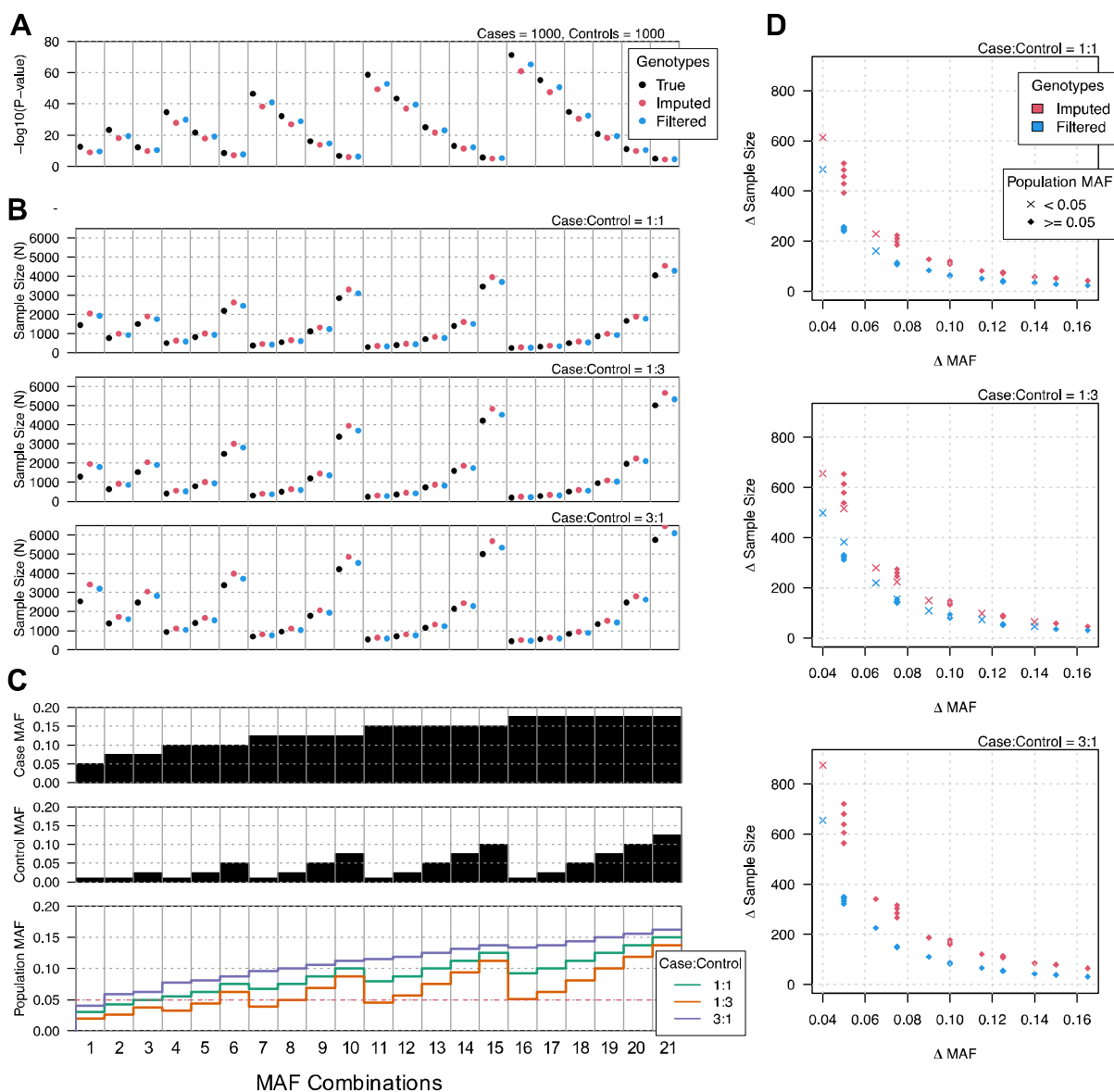
7 *Imputation errors reduce statistical power*

8 After characterizing genotyping errors introduced by imputation, we simulated the impact of
9 these errors on case-control association tests to determine best practices for study designs
10 involving low-pass imputed genotypes. Given a specific genotype detected by high coverage
11 WGS within a specific 0.01 population MAF interval, imputation errors were characterized as the
12 probability of imputing either a homozygous reference, heterozygous, or homozygous alternate
13 genotype. Overall, the imputation process led to decreased significance levels, suggesting that
14 imputation errors may cause statistical significance to be lost for certain experimental
15 configurations (**Fig 6A**). To investigate the loss of statistical significance in the context of study
16 design we performed a power analysis with a focus on the number of samples required to reach
17 sufficient power at 0.8 (**Fig 6B**). Specifically, we tested 21 case and control population MAF
18 combinations over three different case-control ratios and used the MAF of the entire population
19 to simulate imputation errors (**Fig 6C**). Case and control genotypes were based on Hardy-
20 Weinberg equilibrium and were calculated using each population's MAF (Methods). Results
21 showed that experimental configurations for which the difference between case and control was
22 smallest required the largest sample sizes. The required samples sizes were also higher for
23 scenarios with non 1:1 case-control ratios, indicating that asymmetric population sizes increase
24 the impact of imputation errors.

25

26 Importantly, the increased imputation accuracy afforded from quality filtered data translated to
27 reduced sample size requirements for achieving sufficient power (**Fig 6B**). Also, the difference

1 in MAFs between cases and controls was strongly predictive. MAF differences of 0.05 required
2 > 500 additional samples for unfiltered imputed genotypes to achieve the same power as the
3 true genotypes, while MAF differences of 0.1 required approximately 100 additional samples
4 (**Fig 6C**). Similar to the higher number of total samples required for asymmetric populations,
5 sample requirements for imputed data were also increased. Finally, the additional samples
6 required for quality filtered imputed data is slightly greater than 50% of the additional samples
7 required for unfiltered data (**Fig 6C**). Together, these results indicate the importance of
8 considering imputation error in the role of study design. Importantly, the impact from imputation
9 errors is inversely proportional to the effect size of the association.



1
 2 **Fig 6: Impact of imputation errors on case-control GWAS.** **A)** Significance of case-control
 3 GWAS at multiple MAFs. The true genotypes are represented by black circles, where the
 4 frequency of heterozygous and homozygous variants follow Hardy–Weinberg equilibrium. Red
 5 circles represent the outcomes of significance testing on imputed genotypes, while blue circles
 6 represent outcomes after filtering imputed genotypes. Note, decreases in significance were due
 7 to estimates of errors introduced during the process of imputation. Imputation errors were
 8 modeled according to the probability of a given genotype being imputed as any other genotype
 9 at any stated MAF. **B)** Power analysis of significance testing for case-control GWAS of true and
 10 imputed genotypes. Y axis shows required samples size to reach a statistical power of 0.80.
 11 Each individual plot shows different case-control ratios. Power was calculated for a 2x2 chi-
 12 square test for significance level 5×10^{-8} , where effect size was calculated as Cohen's w . **C)**
 13 Case and control MAFs used for each significance test analysis and the combined population
 14 allele frequency for each case and control configuration. **D)** Additional samples required to

1 reach sufficient power for imputed genotypes. Delta sample size is the difference between
2 required sample sizes for true genotypes and imputed or quality filtered imputed genotypes.
3 Delta MAF is the difference in MAFs between cases and controls. Delta MAF is proportional to
4 effect size.
5

6 **Discussion**

7 Genotype imputation is a valuable tool for filling in missing genotypes and improving power to
8 detect genome wide associations. It also provides an opportunity to combine samples
9 genotyped on different platforms into a single analysis (Marchini et al. 2007; Uh et al. 2012;
10 Hayward et al. 2019). Ultimately, this increases the amount of available data by allowing
11 datasets to be reused in larger studies (Ho and Lange 2010; Zhuang et al. 2010). However,
12 imputation techniques have gone a step further and now facilitate genotyping of poor quality or
13 low-pass sequencing data, which often lack sufficient coverage for genotyping software to
14 assign confident calls. Imputation reference panels and algorithms provide the additional
15 statistical support required to assign genotypes to individual samples (Rubinacci et al. 2021;
16 Wasik et al. 2021). However, since the samples used to impute genotypes are different from
17 those that undergo low-pass WGS, imputation has the potential to introduce genotyping errors
18 and biases. Therefore, we investigated the impact of imputation on genotyping in dogs, an
19 important genetic system whose unique population structure, defined by strict breeding
20 programs over hundreds of years (Ostrander et al. 2017), may uniquely influence imputation
21 accuracy. To characterize imputation errors and provide best practices for analyzing low-pass
22 imputation data in dogs, we tested imputation for 97 high converge WGS dog samples with
23 approximately 1x coverage per sample. Imputation errors were detected by comparing the
24 imputed genotypes to genotypes determined using high coverage WGS data. We analyzed
25 these errors in the context of genotype quality filtering, imputation error biases, the role of
26 MAFs, and their impact on case-control analyses.

27

1 Our analysis of case-control association tests of imputed genotypes provides the necessary
2 information to outline best practice guidelines for working with low-pass imputed genotypes. We
3 show herein that the most important factor to consider is the expected allele frequencies in both
4 cases and controls for any potentially associated markers. This is important, as the most
5 common error observed was associated with imputing a heterozygous genotype as
6 homozygous reference, leading to an overall reduction in observed effect size and therefore
7 requiring more samples to reach sufficient power. Importantly, these reduced associations were
8 most pronounced when the allele frequency difference between cases and controls was small (\leq
9 0.05). Therefore, best practices in study design for low-pass imputation are to investigate
10 genetic associations with medium to large effect sizes. Alternatively, if effect sizes are likely
11 small, investigators need to consider increased sample sizes, balance between case and
12 control populations, and the role of quality filtering for improving overall accuracy.

13
14 Quality filtering imputed genotypes is achieved by removing sites where $\geq 5\%$ of genotypes
15 have maximum genotype probabilities (GP) < 0.9 . This threshold was chosen as it is robust to
16 small shifts in genotype failure rate, which causes only a small change in the total number of
17 sites removed. The genotype failure rate threshold of 5% is well-calibrated to optimally reduce
18 imputation errors at a GP threshold of 0.9. We demonstrate that quality filtering is able to
19 improve IQSs by a value of ~ 0.04 and reduce required increases in sample sizes for sufficient
20 power by as much as 50%. Essentially, these improvements are achieved by removing $\sim 20\%$
21 of sites that together share a disproportionate number of the imputation errors. In addition,
22 removing sites with MAFs < 0.05 removes those with comparatively lower imputation accuracies
23 and sites that aren't found in high coverage WGS datasets.

24

1 After both quality and MAF filtering ~ 7M SNVs remain, with SNVs found approximately every
2 360 bp. Despite removing the majority of SNV markers through filtering sufficient numbers of
3 SNVs remain for association analyses. Typically, association studies within breeds demarcate
4 regions on the order of one Mb (Sutter et al. 2004; Lindblad-Toh et al. 2005; Vaysse et al. 2011;
5 Karlsson et al. 2013), whereas across breeds, the scale of LD is approximately 10 – 100 kb
6 (Karlsson et al. 2007; Parker et al. 2017). In addition, currently available canine DNA genotyping
7 arrays contain just over 710,000 markers (Axiom™ Canine HD Array), one tenth the total
8 number of markers available from low-pass imputation after quality and MAF filtering. Therefore,
9 the benefit of increased genotyping accuracy using filtering likely exceeds the cost incurred from
10 reduced marker density.

11
12 An additional consideration in performing GWASs using imputed data with small effect sizes is
13 the required increase of both the case and control populations to reach sufficient power. For
14 filtered genotypes with population MAF differences of 5%, approximately 250 additional samples
15 with equal proportions of cases and controls are required to reach sufficient power. Importantly,
16 as the MAF difference between cases and controls decreases, the required increase to sample
17 sizes appears to grow exponentially. This is perhaps linked to increased rates of imputation
18 errors at MAFs < 0.05.

19
20 Finally, for best practices study design, balance between case and control populations should
21 also be addressed. The impact on power from unbalanced case and control populations is most
22 prominent at low MAF differences between cases and controls. When the ratio between cases
23 and controls favors either population, there is an overall loss of power compared to when the
24 ratio is even. For example, at a case-control MAF difference of 5%, where the ratio of cases to
25 controls is either 1:3 or 3:1, > 300 additional samples are required to reach sufficient power,

1 whereas if the ratio is 1:1, only 250 additional samples would be required. This is because a
2 lower number of total associated alleles, as opposed to proportion of alleles, within either cases
3 or controls, increases the likelihood that the accumulation of imputation errors can cause
4 statistical significance to be lost.

5

6 There are two strategies for developing reference panels, the first is to use a population with
7 closely matched ancestry to that of the group under study, and the second is to use as many
8 samples as possible. While not evaluated in the context of low-pass sequencing imputation,
9 analysis of DNA arrays shows that reference panels matched to the population of interest
10 outperform diverse reference panels of similar sizes (Mitt et al. 2017; Zhou et al. 2017; Bai et al.
11 2019; Yoo et al. 2019). This would suggest that larger reference panels are preferable as long as
12 they contain sufficient representation of the study population. However, the addition of diverse
13 samples to a reference panel can decrease imputation accuracy at low MAFs, where the
14 magnitude of this effect varies according to the population being studied (Bai et al. 2019). These
15 observed effects were for initial population-matched reference panels of ~100 samples, with
16 additional diverse samples increasing the reference panels to over 860 samples (Bai et al.
17 2019). Other analyses compared reference panels of > 1500 samples to the Haplotype
18 Reference Consortium (HRC) reference panel (<http://www.haplotype-reference-consortium.org/>)
19 (McCarthy et al. 2016; Mitt et al. 2017; Zhou et al. 2017; Yoo et al. 2019), which consists of
20 32,611 samples, indicating the potential for increased resolution in human studies compared to
21 canine studies, which used a panel of just 676 samples from 91 breeds (Piras et al. 2020).
22 Altogether, at MAFs > 0.05, human imputation studies conducted using DNA array genotypes
23 show non-reference concordance rates > 97.5% and mean r^2 values > 0.95 (Mitt et al. 2017;
24 Zhou et al. 2017; Yoo et al. 2019). By comparison, prior to filtering non-reference concordance
25 rates for low-pass sequence imputation in dogs were at ~ 95% and had mean r^2 values between

1 0.90 and 0.94 (**Fig 4C**), highlighting the potential gains that can be achieved from improved
2 canine reference panels.

3
4 An important initiative that may help address short comings in available high coverage WGS
5 samples in dogs is the Dog 10K project which aims to achieve 10,000 modest-high coverage
6 dog genomes representing an array of canine genetic diversity (Ostrander et al. 2019). Similar
7 to the current dog reference panel, the initial phase of the Dog 10K project prioritizes collecting
8 samples from as many modern breeds as possible. Alternatively, human data suggests that for
9 imputation purposes it is better to use samples from the population of interest. Whether this
10 same strategy is preferable in dogs is unknown.

11
12 Many mapping studies in dogs focus on traits that segregate across breeds, as breeds sharing
13 recent common ancestry likely share the same genetic underpinnings for any given trait (Parker
14 et al. 2017). While multiple breeds are often included within a single analysis, creating many
15 breed-specific reference panels is not feasible. As large haplotype blocks are shared between
16 many breeds and clades, perhaps a large reference panel representing a greater number of
17 breeds could provide even higher levels of imputation accuracy than a breed-specific reference
18 panel. This idea is supported by an array-based imputation analyses that tested the imputation
19 accuracy for a group of poodles with three different reference panel configurations. Results
20 showed a composite panel of poodles and non-poodles outperformed the poodle only and the
21 non-poodle reference panels (Friedenberg and Meurs 2016). A key finding from our analysis,
22 was that SNV discovery in individual dogs was similar between reference panel and non-
23 reference panel breeds (**Fig 2C**), and while for imputation accuracy rates reference panel
24 breeds scored highest, several non-reference panel breeds outperformed the majority of
25 reference panel breeds (**Fig 5A**). This was likely due to the fact that many of breeds from the

1 test samples belonged to clades represented in the reference panel. However, since the original
2 VCF for the reference panel was unavailable, breeds were identified from matching IDs across
3 databases and breed membership for 122 dogs could not be determined. Also, many of the
4 samples in the reference panel were either village dogs or other canid species. In our test
5 samples, of the 31 dogs from 14 breeds not previously associated with a clade (**Supplemental**
6 **table S3**), five were terrier breeds that likely belong to the primary terrier clade and two were
7 spaniel breeds that can be assigned to the spaniel clade. Since breeds of European origin are
8 heavily represented in our reference panel, and most breeds with no clade assignment are of
9 European ancestry, it is likely that many of the non-clade assigned test sample haplotypes are
10 at least partially represented in the reference panel. Supporting this idea of partial haplotype
11 representation providing accurate imputation in dogs is the observation of high levels of
12 imputation accuracy shared between mixed and pure breeds (Hayward et al. 2019). As more
13 high coverage WGS samples become available through initiatives such as Dog 10K, optimal
14 reference panel designs can be constructed.

15
16 One additional improvement in imputation accuracy may derive from the choice of imputation
17 algorithm. Currently, available tools for imputation of low-pass WGS data include STITCH
18 (Davies et al. 2016), Beagle (Browning and Browning 2016), GeneImp (Spiliopoulou et al.
19 2017), GLIMPSE (Rubinacci et al. 2021), and loimpute (Wasik et al. 2021). Our analyses used
20 loimpute, as it had already been implemented with a dog reference panel and used by the
21 canine genomics community (Piras et al. 2020). However, if low-pass sequencing imputation
22 approaches in dogs are going to improve, other algorithms need to be appropriately assessed
23 with accuracy across a range of MAFs. Currently, human studies demonstrate that GLIMPSE
24 outperforms the other algorithms in terms of both accuracy and required computations
25 resources (Rubinacci et al. 2021). The largest differences were observed for variants with MAFs

1 < 1%. For common alleles, imputation accuracy was similarly high between almost all
2 algorithms.

3
4 Our work provides the first in-depth analysis of low-pass WGS and imputation in canine
5 genomics and can act as a road map for analysis in other non-human species. By comparing
6 genotypes imputed from downsampled reads to a high coverage truth set, we have been able to
7 rigorously investigate the nature of imputation errors and their biases. We were able to optimize
8 filtering strategies to improve accuracy rates and also demonstrate the impact imputation errors
9 have on case-control GWAS. Our results inform a series of best practices guidelines and
10 demonstrated the utility of this quickly evolving resource for future analyses. Altogether,
11 widespread adoption of low-pass sequencing and imputation within the canine genomics field,
12 together with investment in developing improved reference panels, will lead to more high-
13 powered analyses and successful discovery of genotype-phenotype associations.

14

15 **Methods**

16 *Sample selection*

17 Samples were selected on the basis of whether they belonged to known breeds, were absent
18 from the reference panel, and had mean sequence coverage levels > 15x. Non-publicly
19 available samples were not used in the reference panel and therefore provide an accurate test
20 of imputation performance. Gencove Inc. provided a list of sample IDs, most of which were
21 matched to known samples within the Plassais et al. (2019) dataset, which was used to
22 represent variant population frequencies. Breed names were based on annotated records and
23 clade membership was based on previously published results (Parker et al. 2017). Breeds with

1 no recorded clade membership were assigned to a clade based on their phenotype, historical
2 information or phylogenetic clustering in Plassais et al. (2019).

3

4 *Variant calling and imputation*

5 Sample reads were mapped to CanFam3.1 using BWA-mem (Li 2013). Variant calling was
6 performed using GATK4 best practices (McKenna et al. 2010). Base quality score recalibration
7 and duplicate marking was applied to each sample (DePristo et al. 2011; Van der Auwera et al.
8 2013), and haplotypcaller was used for variant discovery (Poplin et al. 2017). Average
9 coverage was estimated using Samtools depth tool (Li et al. 2009). To simulate low-pass
10 sequencing, BAM files were downsampled to approximately 1x coverage using the
11 DownsampleSam tool from GATK4. To obtain the correct coverage level the parameter “-p” was
12 set as the sample’s mean coverage divided by one. Downsampled BAMs were converted to
13 fastq files using samtools “fastq” function and were uploaded to Gencove, Inc. using the
14 Gencove command line interface (CLI). Imputation was performed using loimpute as part of
15 Gencove’s imputation pipeline with the “Dog low-pass v2.0” configuration (Piras et al. 2020;
16 Wasik et al. 2021). Imputed genotypes were received from Gencove as a VCF for each
17 individual. Individual VCFs were split according to chromosome. Each sample’s genotypes and
18 genotyping statistics were merged to create a single dataset for each chromosome that
19 contained all individuals. This task was performed using the program extract_genotype_wg.R
20 which was written in R and used the vcfR package (Team 2013; Knaus and Grunwald 2017).

21

22 *Assessing imputation accuracy*

23 Imputation accuracy was assessed by comparing imputed genotypes to high coverage WGS
24 genotypes. This was made possible by identifying sites shared across both datasets. Sites were

1 considered shared if position, reference allele, and alternate allele were identical. Importantly,
2 all multiallelic sites in all VCFs were split into biallelic states using the bcftools “norm” function
3 with the “-m -” parameter to ensure all potential allelic combinations were matched (Li 2011).
4 Sites where all samples were homozygous for the reference allele were removed from the
5 analysis. A genotype for a given individual at a particular site was considered concordant if the
6 imputed genotype was identical to the genotype determined using high coverage WGS.
7 Imputation errors were those that were not identical between the high coverage WGS dataset
8 and imputed dataset for a given individual at a given site. The total number of concordant
9 genotypes per sample was calculated using the program “venn_filter_wg.R”. Once filtering
10 thresholds were determined (below), imputation accuracies were determined according to MAF
11 intervals of 0.01 using the program “gt_by_af.R”.

12

13 *Determining filtering thresholds*

14 The imputation process provides genotype probabilities as a measure of confidence regarding
15 whether a call is homozygous reference, heterozygous, or homozygous alternate. The max
16 genotype probability (GP) is the level of confidence for the imputed genotype. We therefore
17 investigated the relationship between GP and the proportion of concordant genotypes between
18 our imputed and high coverage WGS datasets in order to determine optimal strategies for
19 filtering imputed variants. Low confidence or failed genotypes were identified according to their
20 GP values and variant sites were filtered out if the number of low confidence genotypes was
21 above a given threshold. (**Supplemental Fig. S1A**). Filtering strategies were evaluated
22 according to the number of remaining genotypes after filtering and the proportion of these
23 genotypes that were concordant with the truth set. These values were further investigated in
24 terms of the true positive rate (TPR), false positive rate (FPR), false discovery rate (FDR), and
25 the number of genotypes kept after filtering (**Supplemental Fig. S1B**).

1

2 *Simulation of imputation errors*

3 Imputation errors were simulated using probabilities derived from the observed fraction of any
4 given genotype that was incorrectly imputed. Error probabilities were also grouped according to
5 population MAFs. For example, for sites with MAFs between 0.03 and 0.04, 5% of heterozygous
6 sites may be imputed as homozygous reference, whereas for sites with MAFs between 0.10 and
7 0.11, 3% of heterozygous sites may be imputed as homozygous reference. This strategy is
8 used to capture the MAF impact on imputation accuracy and error rates. The imputed
9 genotypes, I , were stored as a 1 x 3 matrix, consisting of the imputed genotype counts for
10 homozygous reference, heterozygous, and homozygous alternate genotypes. The values for I
11 were calculated as $I = GP$, where G is the counts for the true starting genotypes stored in a 1 x
12 3 matrix and P is a 3 x 3 matrix of the probabilities a true genotype is imputed as any other
13 genotype. In P , rows represent the true starting genotype and columns represent the imputed
14 genotypes. Importantly, the values used for P depend on the population MAF and whether the
15 genotypes were filtered for quality. A total of 100 P matrices were defined, with one matrix for
16 each 0.01 MAF interval between 0 and 0.5 for both quality filtered and unfiltered imputation error
17 rates.

18

19 *Association analyses and statistical power calculations*

20 Case-control association analyses were performed as a chi square test on a 2 x 2 contingency
21 table, measuring the association between cases and the presence of the minor allele.
22 Association tests were carried out on the true genotypes, imputed genotypes, and quality
23 filtered imputed genotypes. Power was calculated using the “pwr” package in R (Champely
24 2020), which uses Cohen’s w to calculate effect size. Probabilities for the null hypothesis were

1 calculated as if the minor allele was evenly distributed across both cases and controls.
2 Significance levels were set at 5×10^{-8} and power was calculated across a variety of case and
3 control MAFs for all population sizes between 100 and 7000. Required sample sizes for
4 sufficiently powered analyses were identified as the lowest sample size that could achieve a
5 power level of 0.8 or greater for a particular case-control analysis.

6

7 *Software and data analysis*

8 All original software used here can be obtained from the following URL:

9 <https://github.com/NHGRI-dog-genome>

10

11 **Data Access**

12 Test samples used in this analysis are deposited under the BioProject accession
13 PRJNA648123. Individual accessions for each sample are recorded in Supplemental Table S1.

14

15 **Competing Interests**

16 All authors declare no competing interests and that the presented work is original.

17

18 **Acknowledgements**

19 All authors acknowledge Gencove, Inc. for their assistance and guidance throughout the project.

20 R.M.B., A.C.H., D.T.W., and E.A.O. were supported by the Intramural Program of the National

21 Human Genome Research Institute at the National Institutes of Health. Support was also

22 provided by the National Key R&D Program of China (2019YFA0707101), Key Research

23 Program of Frontier Sciences of the CAS (ZDBS-LY-SM011), and Innovative Research Team

1 (in Science and Technology) of Yunnan Province (202005AE160012). G.D.W. is supported by
2 the Youth Innovation Promotion Association of CAS. We especially thank dog owners who have
3 provided samples from their dogs for these and other studies.

4

5 **Supplemental Tables**

6 Supplemental Table S1: Sample-specific metadata.

7 Supplemental Table S2: Breed composition of datasets.

8 Supplemental Table S3: Breed clade membership.

9 Supplemental Table S4: Breeds not included in previous phylogenetic analyses and have no
10 clade membership.

11 Supplemental Table S5: Mean chromosome 38 coverage levels at variant sites.

12 Supplemental Table S6: Imputation accuracy rates for individual samples measured as non-
13 reference concordance.

14

15 **References**

16 Ali MB, Evans JM, Parker HG, Kim J, Pearce-Kelling S, Whitaker DT, Plassais J, Khan QM,

17 Ostrander EA. 2020. Genetic analysis of the modern Australian labradoodle dog breed
18 reveals an excess of the poodle genome. *PLoS Genet* **16**: e1008956.

19 Awano T, Johnson GS, Wade CM, Katz ML, Johnson GC, Taylor JF, Perloski M, Biagi T,

20 Baranowska I, Long S. 2009. Genome-wide association analysis reveals a SOD1 mutation
21 in canine degenerative myelopathy that resembles amyotrophic lateral sclerosis.

22 *Proceedings of the National Academy of Sciences* **106**: 2794-2799.

- 1 Bai WY, Zhu XW, Cong PK, Zhang XJ, Richards JB, Zheng HF. 2019. Genotype imputation and
2 reference panel: a systematic evaluation on haplotype size and diversity. *Brief Bioinform*
3 doi:10.1093/bib/bbz108.
- 4 Benjelloun B, Boyer F, Streeter I, Zamani W, Engelen S, Alberti A, Alberto FJ, BenBati M,
5 Ibbelbachyr M, Chentouf M et al. 2019. An evaluation of sequencing coverage and
6 genotyping strategies to assess neutral and adaptive diversity. *Mol Ecol Resour* **19**: 1497-
7 1515.
- 8 Boyko AR, Quignon P, Li L, Schoenebeck JJ, Degenhardt JD, Lohmueller KE, Zhao K, Brisbin A,
9 Parker HG, vonHoldt BM et al. 2010. A simple genetic architecture underlies
10 morphological variation in dogs. *PLoS Biol* **8**: e1000451.
- 11 Browning BL, Browning SR. 2016. Genotype Imputation with Millions of Reference Samples. *Am*
12 *J Hum Genet* **98**: 116-126.
- 13 Champely S. 2020. pwr: Basic functions for Power Analysis.
- 14 Davies RW, Flint J, Myers S, Mott R. 2016. Rapid genotype imputation from sequence without
15 reference panels. *Nature genetics* **48**: 965.
- 16 DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis AA, Del Angel G,
17 Rivas MA, Hanna M. 2011. A framework for variation discovery and genotyping using
18 next-generation DNA sequencing data. *Nature genetics* **43**: 491.
- 19 Friedenberg SG, Meurs KM. 2016. Genotype imputation in the domestic dog. *Mamm Genome*
20 **27**: 485-494.

- 1 Friedrich J, Strandberg E, Arvelius P, Sanchez-Molano E, Pong-Wong R, Hickey JM, Haskell MJ,
2 Wiener P. 2019. Genetic dissection of complex behaviour traits in German Shepherd
3 dogs. *Heredity (Edinb)* **123**: 746-758.
- 4 Hayward JJ, Castelhana MG, Oliveira KC, Corey E, Balkman C, Baxter TL, Casal ML, Center SA,
5 Fang M, Garrison SJ et al. 2016. Complex disease and phenotype mapping in the
6 domestic dog. *Nat Commun* **7**: 10460.
- 7 Hayward JJ, White ME, Boyle M, Shannon LM, Casal ML, Castelhana MG, Center SA, Meyers-
8 Wallen VN, Simpson KW, Sutter NB et al. 2019. Imputation of canine genotype array
9 data using 365 whole-genome sequences improves power of genome-wide association
10 studies. *PLoS Genet* **15**: e1008003.
- 11 Ho LA, Lange EM. 2010. Using public control genotype data to increase power and decrease
12 cost of case-control genetic association studies. *Hum Genet* **128**: 597-608.
- 13 Karlsson EK, Baranowska I, Wade CM, Salmon Hillbertz NH, Zody MC, Anderson N, Biagi TM,
14 Patterson N, Pielberg GR, Kulbokas EJ, 3rd et al. 2007. Efficient mapping of mendelian
15 traits in dogs through genome-wide association. *Nat Genet* **39**: 1321-1328.
- 16 Karlsson EK, Sigurdsson S, Ivansson E, Thomas R, Elvers I, Wright J, Howald C, Tonomura N,
17 Perloski M, Swofford R et al. 2013. Genome-wide analyses implicate 33 loci in heritable
18 dog osteosarcoma, including regulatory variants near CDKN2A/B. *Genome Biol* **14**: R132.
- 19 Knaus BJ, Grunwald NJ. 2017. vcfR: a package to manipulate and visualize variant call format
20 data in R. *Mol Ecol Resour* **17**: 44-53.

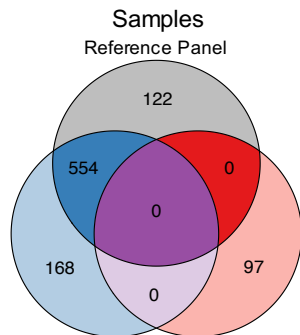
- 1 Li H. 2011. A statistical framework for SNP calling, mutation discovery, association mapping and
2 population genetical parameter estimation from sequencing data. *Bioinformatics* **27**:
3 2987-2993.
- 4 Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
5 *arXiv preprint arXiv:13033997*.
- 6 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,
7 Genome Project Data Processing S. 2009. The Sequence Alignment/Map format and
8 SAMtools. *Bioinformatics* **25**: 2078-2079.
- 9 Lin P, Hartz SM, Zhang Z, Saccone SF, Wang J, Tischfield JA, Edenberg HJ, Kramer JR, A MG,
10 Bierut LJ et al. 2010. A new statistic to evaluate imputation reliability. *PLoS One* **5**:
11 e9697.
- 12 Lindblad-Toh K Wade CM Mikkelsen TS Karlsson EK Jaffe DB Kamal M Clamp M Chang JL
13 Kulbokas EJ, 3rd Zody MC et al. 2005. Genome sequence, comparative analysis and
14 haplotype structure of the domestic dog. *Nature* **438**: 803-819.
- 15 Marchini J, Howie B, Myers S, McVean G, Donnelly P. 2007. A new multipoint method for
16 genome-wide association studies by imputation of genotypes. *Nat Genet* **39**: 906-913.
- 17 Martin AR, Atkinson EG, Chapman SB, Stevenson A, Stroud RE, Abebe T, Akena D, Alemayehu
18 M, Ashaba FK, Atwoli L et al. 2021. Low-coverage sequencing cost-effectively detects
19 known and novel variation in underrepresented populations. *Am J Hum Genet*
20 doi:10.1016/j.ajhg.2021.03.012.

- 1 McCarthy S, Das S, Kretzschmar W, Delaneau O, Wood AR, Teumer A, Kang HM, Fuchsberger C
2 Danecek P, Sharp K et al. 2016. A reference panel of 64,976 haplotypes for genotype
3 imputation. *Nat Genet* **48**: 1279-1283.
- 4 McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler D,
5 Gabriel S, Daly M. 2010. The Genome Analysis Toolkit: a MapReduce framework for
6 analyzing next-generation DNA sequencing data. *Genome research* **20**: 1297-1303.
- 7 Mitt M, Kals M, Parn K, Gabriel SB, Lander ES, Palotie A, Ripatti S, Morris AP, Metspalu A, Esko T
8 et al. 2017. Improved imputation accuracy of rare and low-frequency variants using
9 population-specific high-coverage WGS-based imputation reference panel. *Eur J Hum*
10 *Genet* **25**: 869-876.
- 11 Nosková A, Bhati M, Kadri NK, Crysanto D, Neuenschwander S, Hofer A, Pausch H. 2021.
12 Characterization of a haplotype-reference panel for genotyping by low-pass sequencing
13 in Swiss Large White pigs. *bioRxiv*.
- 14 Ostrander EA, Wang GD, Larson G, vonHoldt BM, Davis BW, Jagannathan V, Hitte C, Wayne RK,
15 Zhang YP, Dog KC. 2019. Dog10K: an international sequencing effort to advance studies
16 of canine domestication, phenotypes and health. *Natl Sci Rev* **6**: 810-824.
- 17 Ostrander EA, Wayne RK, Freedman AH, Davis BW. 2017. Demographic history, selection and
18 functional diversity of the canine genome. *Nat Rev Genet* **18**: 705-720.
- 19 Parker HG, Dreger DL, Rimbault M, Davis BW, Mullen AB, Carpintero-Ramirez G, Ostrander EA.
20 2017. Genomic Analyses Reveal the Influence of Geographic Origin, Migration, and
21 Hybridization on Modern Dog Breed Development. *Cell Rep* **19**: 697-708.

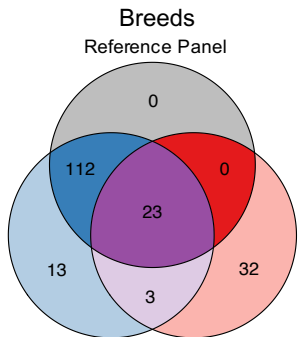
- 1 Piras IS, Bleul C, Siniard A, Wolfe AJ, De Both MD, Hernandez AG, Huentelman MJ. 2020.
2 Association of Common Genetic Variants in the CPSF7 and SDHAF2 Genes with Canine
3 Idiopathic Pulmonary Fibrosis in the West Highland White Terrier. *Genes (Basel)* **11**.
4 Plassais J, Kim J, Davis BW, Karyadi DM, Hogan AN, Harris AC, Decker B, Parker HG, Ostrander
5 EA. 2019. Whole genome sequencing of canids reveals genomic regions under selection
6 and variants influencing morphology. *Nat Commun* **10**: 1489.
7 Poplin R, Ruano-Rubio V, DePristo M, Fennell T, Carneiro M, Van der Auwera G, Kling D,
8 Gauthier L, Levy-Moonshine A, Roazen D et al. 2017. Scaling accurate genetic variant
9 discovery to tens of thousands of samples. doi:10.1101/201178. bioRxiv.
10 Rubinacci S, Ribeiro DM, Hofmeister RJ, Delaneau O. 2021. Efficient phasing and imputation of
11 low-coverage sequencing data using large reference panels. *Nat Genet* **53**: 120-126.
12 Shannon LM, Boyko RH, Castelhana M, Corey E, Hayward JJ, McLean C, White ME, Abi Said M,
13 Anita BA, Bondjengo NI et al. 2015. Genetic structure in village dogs reveals a Central
14 Asian domestication origin. *Proc Natl Acad Sci U S A* **112**: 13639-13644.
15 Snelling WM, Hoff JL, Li JH, Kuehn LA, Keel BN, Lindholm-Perry AK, Pickrell JK. 2020. Assessment
16 of Imputation from Low-Pass Sequencing to Predict Merit of Beef Steers. *Genes (Basel)*
17 **11**.
18 Spiliopoulou A, Colombo M, Orchard P, Agakov F, McKeigue P. 2017. GenImp: Fast Imputation
19 to Large Reference Panels Using Genotype Likelihoods from Ultralow Coverage
20 Sequencing. *Genetics* **206**: 91-104.

- 1 Sutter NB, Eberle MA, Parker HG, Pullar BJ, Kirkness EF, Kruglyak L, Ostrander EA. 2004.
2 Extensive and breed-specific linkage disequilibrium in *Canis familiaris*. *Genome Res* **14**:
3 2388-2396.
- 4 Team RC. 2013. R: A language and environment for statistical computing.
- 5 Uh HW, Deelen J, Beekman M, Helmer Q, Rivadeneira F, Hottenga JJ, Boomsma DI, Hofman A,
6 Uitterlinden AG, Slagboom PE et al. 2012. How to deal with the early GWAS data when
7 imputing and combining different arrays is necessary. *Eur J Hum Genet* **20**: 572-576.
- 8 Van der Auwera GA, Carneiro MO, Hartl C, Poplin R, Del Angel G, Levy-Moonshine A, Jordan T,
9 Shakir K, Roazen D, Thibault J. 2013. From FastQ data to high-confidence variant calls:
10 the genome analysis toolkit best practices pipeline. *Current protocols in bioinformatics*
11 **43**: 11.10. 11-11.10. 33.
- 12 Vaysse A, Ratnakumar A, Derrien T, Axelsson E, Rosengren Pielberg G, Sigurdsson S, Fall T,
13 Seppala EH, Hansen MS, Lawley CT et al. 2011. Identification of genomic regions
14 associated with phenotypic variation between dog breeds using selection mapping. *PLoS*
15 *Genet* **7**: e1002316.
- 16 Wasik K, Berisa T, Pickrell JK, Li JH, Fraser DJ, King K, Cox C. 2021. Comparing low-pass
17 sequencing and genotyping for trait mapping in pharmacogenetics. *BMC Genomics* **22**:
18 197.
- 19 Yoo SK, Kim CU, Kim HL, Kim S, Shin JY, Kim N, Yang JSW, Lo KW, Cho B, Matsuda F et al. 2019.
20 NARD: whole-genome reference panel of 1779 Northeast Asians improves imputation
21 accuracy of rare and low-frequency variants. *Genome Med* **11**: 64.

- 1 Zhou W, Fritsche LG, Das S, Zhang H, Nielsen JB, Holmen OL, Chen J, Lin M, Elvestad MB, Hveem
- 2 K et al. 2017. Improving power of association tests using multiple sets of imputed
- 3 genotypes from distributed reference panels. *Genet Epidemiol* **41**: 744-755.
- 4 Zhuang JJ, Zondervan K, Nyberg F, Harbron C, Jawaid A, Cardon LR, Barratt BJ, Morris AP. 2010.
- 5 Optimizing the power of genome-wide association studies by using publicly available
- 6 reference samples to expand the control group. *Genet Epidemiol* **34**: 319-326.
- 7

A

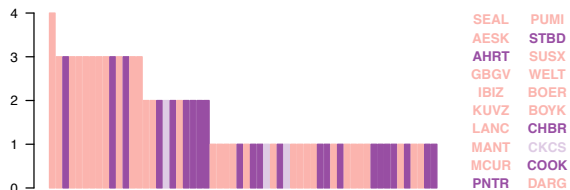
Plassais et al. (2019) Test Samples

B

Plassais et al. (2019) Test Samples

C**Plassais et al (2019)**

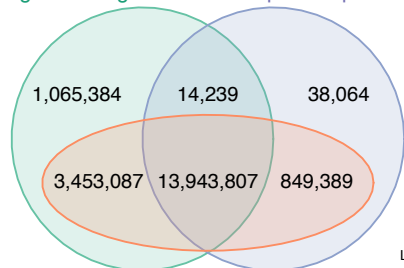
Individual samples (n)

Reference Panel**Test Samples**

Dog Breeds

A**SNVs**

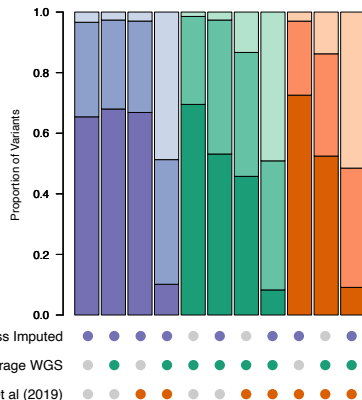
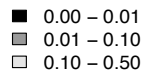
High Coverage WGS Low-pass Imputation



Plassais et al (2019)

B

MAF Range



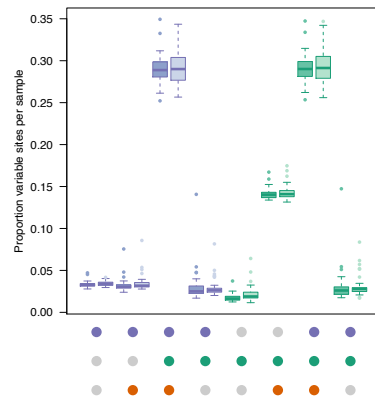
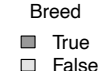
Low-pass Imputed

High Coverage WGS

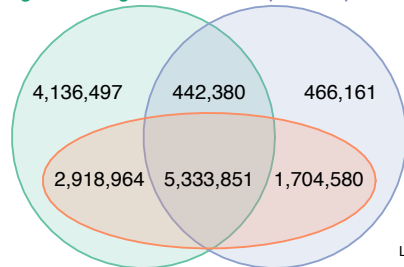
Plassais et al (2019)

C

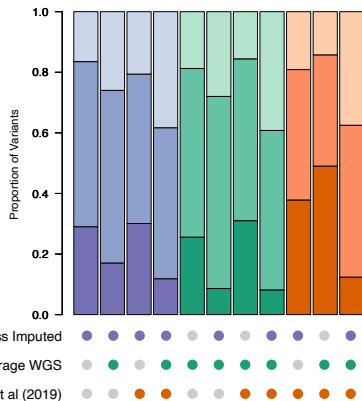
Reference Breed

**Indels**

High Coverage WGS Low-pass Imputation



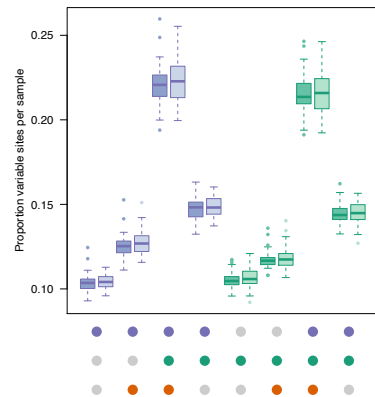
Plassais et al (2019)

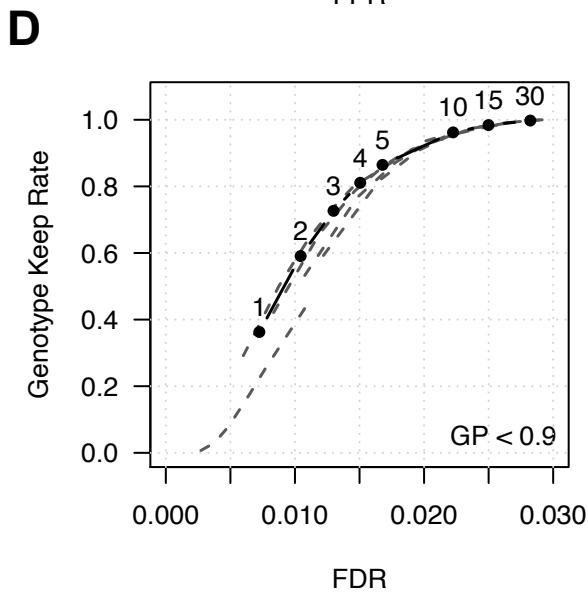
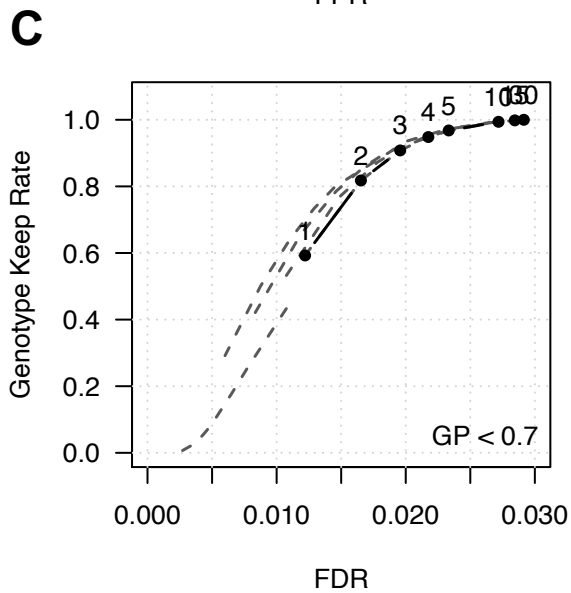
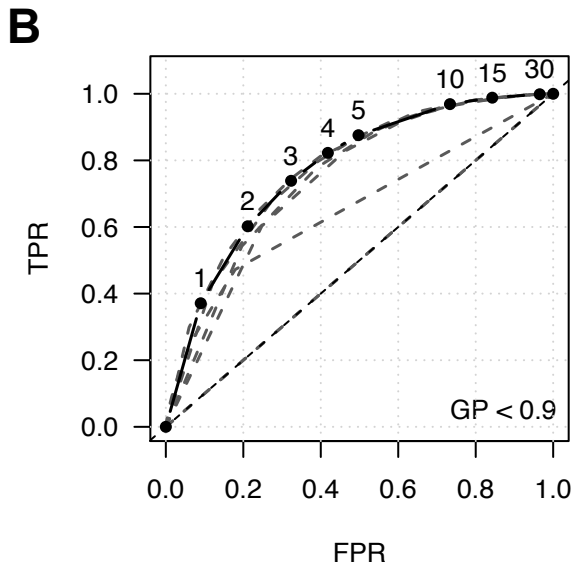
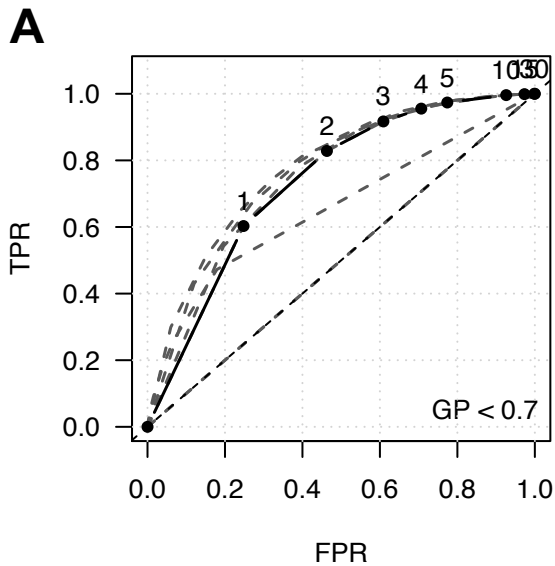


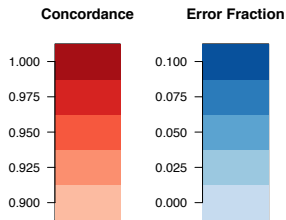
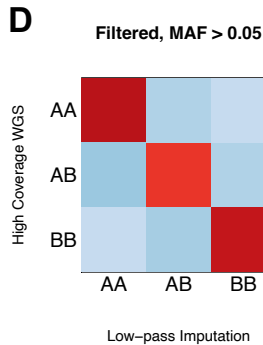
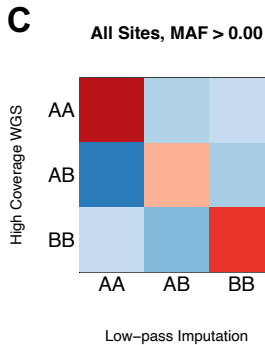
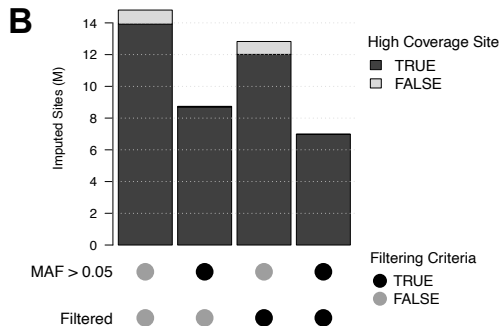
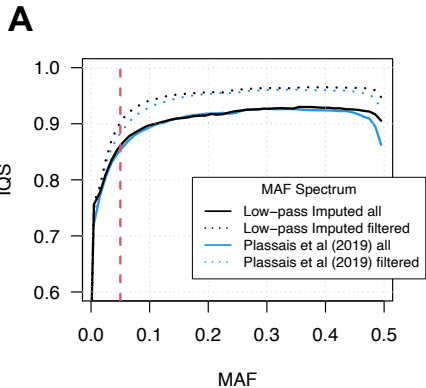
Low-pass Imputed

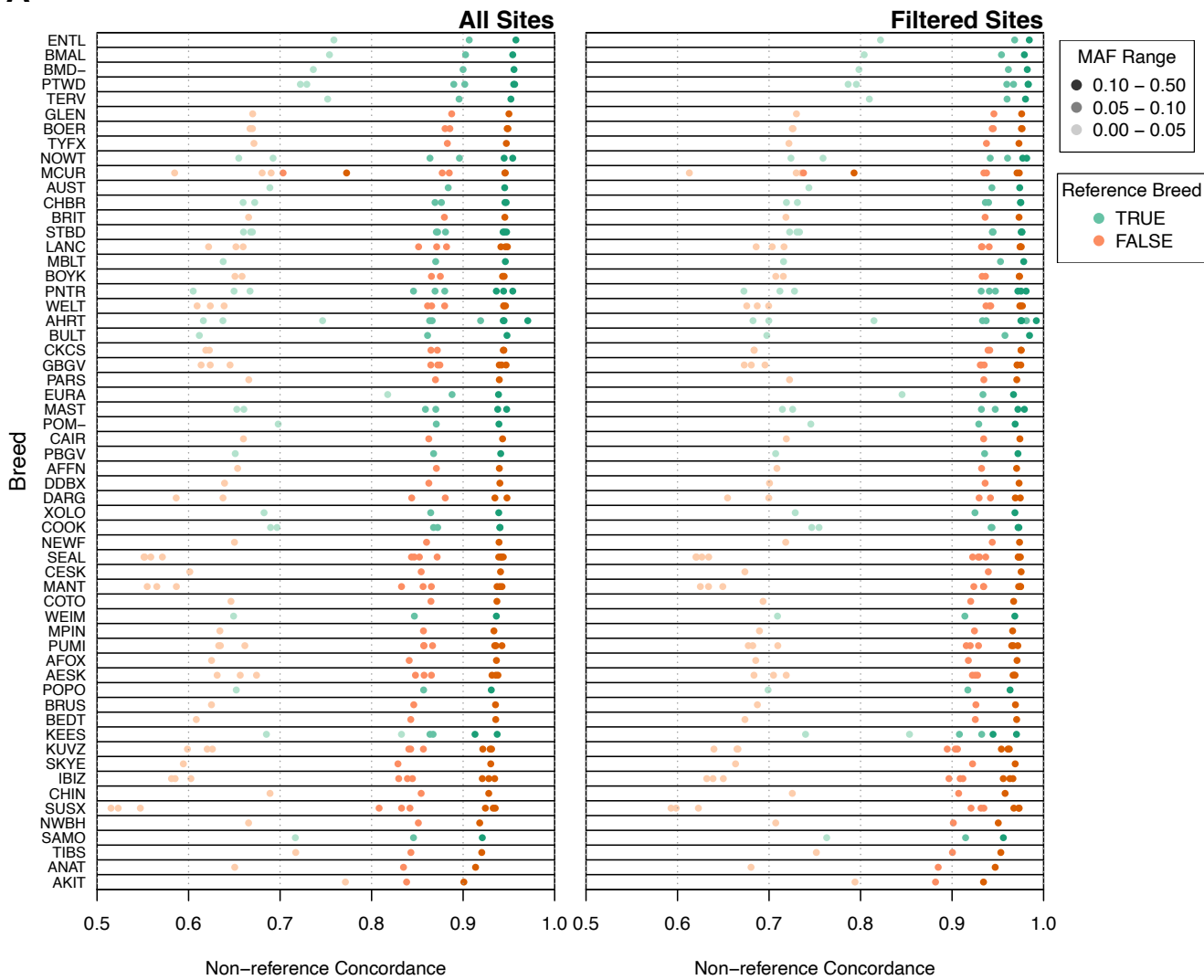
High Coverage WGS

Plassais et al (2019)







A**B**