



An imputed genotype resource for the laboratory mouse

Journal:	<i>Mammalian Genome</i>
Manuscript ID:	draft
Manuscript Type:	Original Contributions
Date Submitted by the Author:	n/a
Complete List of Authors:	Szatkiewicz, Jin; The Jackson Laboratory Beane, Glen; The Jackson Laboratory Ding, Yueming; The Jackson Laboratory Hutchins, Lucie; The Jackson Laboratory Pardo-Manuel de Villena, Fernando; University of North Carolina, Chapel Hill, Department of Genetics, Carolina Center for Genome Sciences and Lineberger Comprehensive Cancer Center Churchill, Gary; The Jackson Laboratory
Keyword:	mouse, SNP, hidden Markov model, missing data



view

An imputed genotype resource for the laboratory mouse

Jin P. Szatkiewicz¹, Glen L. Beane¹, Yueming Ding¹, Lucie Hutchins¹, Fernando Pardo-Manuel de Villena², Gary A. Churchill¹

1. The Jackson Laboratory, Bar Harbor, Maine 04609, USA. 2. Department of Genetics, Carolina Center for Genome Sciences and Lineberger Comprehensive Cancer Center, University of North Carolina, Chapel Hill, Chapel Hill, North Carolina 27599, USA.

Corresponding author:

Gary A. Churchill

The Jackson Laboratory

Bar Harbor

Maine 04609, USA

Email: gary.churchill@jax.org

Phone: 207-288-6189

Fax: 207-288-6847

Running Title: Mouse genotype resource

Key words: mouse, SNP, hidden Markov model, missing data

ABSTRACT

We have created a high-density SNP resource encompassing 7.87 million polymorphic loci across 49 inbred mouse strains of the laboratory mouse by combining data available from public databases and training a hidden Markov model to impute missing genotypes in the combined data. The strong linkage disequilibrium found in dense sets of SNP markers in the laboratory mouse provides the basis for accurate imputation. Using genotypes from eight independent SNP resources, we empirically validated the quality of the imputed genotypes and demonstrate that they are highly reliable for most inbred strains. The imputed SNP resource will be useful for studies of natural variation and complex traits. It will facilitate association study designs by providing high density SNP genotypes for large numbers of mouse strains. We anticipate that this resource will continue to evolve as new genotype data become available for laboratory mouse strains. The data are available for bulk download or query at <http://cgd.jax.org/>.

INTRODUCTION

The laboratory mouse owes much of its popularity as a model organism in biomedical research to the existence of a large collection of inbred strains that represent an immortal population of genetic clones derived by repeated brother sister mating (Lyon et al. 1996). Because mice from each strain are genetically identical it is possible to collect and combine biological data over time and space leading to a depth of phenotype characterization rarely achieved in other mammalian systems (Bogue 2003). Furthermore, the existence of a definite set of genetic differences among inbred strains allows scientists to explore the effect of genetic diversity on almost any phenotype of interest (Wade and Daly 2005). These studies require an accurate description of the level and distribution of genetic variation present among the hundreds of existing inbred strains. This is a challenging problem because the diversity between strains varies from extremely low levels found among sister substrains to very high levels found among strains derived from different species and subspecies (Petkov et al. 2004; Ideraabdullah et al. 2004; Yang et al. 2007).

Inbred strains can be classified into classical and wild-derived strains according to whether they were derived in the 20th century from a small set of founders known as “fancy” mice or derived from mice captured from natural populations more recently. Common wild-derived strains include representatives from two species *Mus spretus* and *M. musculus*, thought to have diverged almost 2 million years ago (Guenet and Bonhomme 2003). There is well over 1% sequence divergence between these species resulting in one SNP every 75bp (Ideraabdullah et al. 2004). Within the *M. musculus* species there are four subspecies,

1
2
3
4
5 *M. m. domesticus*, *M. m. castaneus*, *M. m. musculus* and *M. m. molossinus* from which
6
7 inbred strains have been derived. These subspecies are thought to have diverged 750,000
8
9 years ago (Guenet and Bonhomme 2003). There is roughly 1% divergence among
10
11 subspecies corresponding to one SNP every 150bp (Ideraabdullah et al. 2004). Recent
12
13 analysis of high density genotype data demonstrates that many wild-derived strain
14
15 genomes carry regions of intersubspecific introgression (“contamination”) from a different
16
17 subspecies (Yang et al. 2007). This analysis also confirms that classical strains are derived
18
19 from multiple subspecies but that the contribution of *M. m. domesticus* represents over
20
21 90% of the genome in most of these strains. These data are critical to interpret the results
22
23 of any mouse experiment in the proper evolutionary context and may have profound
24
25 implications for our understanding of basic biological processes such as divergence,
26
27 selection and speciation (Payseur and Hoekstra 2005; Mott 2007).
28
29
30
31
32

33 Since the genomic sequence of the C57BL/6 strain was reported (Waterston et al. 2002)
34
35 much effort has been focused on the discovery and characterization of single nucleotide
36
37 polymorphisms (SNPs) in inbred strains (Wade et al. 2002; Wiltshire et al. 2003; Yalcin et
38
39 al. 2004; Pletcher et al. 2004; Frazer et al. 2004). Early SNP discovery projects carried out
40
41 resequencing in a limited number of classical strains (Mural et al. 2002). More recently,
42
43 the NIEHS used a hybridization based strategy to discover ~8.3 million SNPs in a survey
44
45 of 15 inbred mouse strains, including four wild-derived strains representing the major
46
47 subspecies of *M. musculus* (Frazer et al. 2007). In parallel to these SNP discovery efforts
48
49 the Broad Institute of Harvard and MIT carried out genotyping of ~138,000 known SNPs
50
51 on 49 inbred mouse strains (Wade and Daly 2005) and the Wellcome-CTC genotyped 499
52
53
54
55
56
57
58
59
60

1
2
3
4 inbred strains and outbred stocks at a lower SNP density with only ~13,370 SNPs

5
6
7 (Shifman et al. 2006). Additional SNP resources are listed in Table 1.

8
9
10 In the follow we refer to the NIEHS data as high density (>7 million genotyped SNPs).

11
12 Medium density genotypes are similar in magnitude to the Broad set (>100,000 genotyped

13
14 SNPs) and low density genotypes are similar in magnitude to Wellcome-CTC set (>10,000

15
16 genotyped SNPs). Much of biomedical research involves inbred strains for which the

17
18 description of the diversity is based on low to medium density SNP panels (Liao et al.

19
20 2004; Pletcher et al. 2004; Cervino et al. 2005; Shifman et al. 2006; McClurg et al. 2007;

21
22 Payseur and Place 2007). Linkage disequilibrium (LD) among classical inbred strains is

23
24 extensive (Wade et al. 2002; Petkov et al. 2005), suggesting that we could leverage the

25
26 NIEHS data to impute genotypes at high density in a larger set of inbred mouse strains.

27
28 Achieving this goal should immediately empower hundreds of laboratories to narrow

29
30 quantitative trait loci, help design the next generation of experiments in mammalian

31
32 genetics and provide invaluable support in the field of comparative and evolutionary

33
34 genomics for the study of biological processes such as recombination, mutation and

35
36 selection (Dipetrillo et al. 2005; Siebert and Schadt 2007; Roberts et al. 2007).

37
38 We propose a method to impute genotypes at high density in strains for which only

39
40 medium or low density genotype data are available. We apply this method to create a

41
42 resource of SNP genotypes at ~7.9 million loci across 49 inbred strains by combining

43
44 existing public databases and imputing missing genotypes. The quality of the imputed

45
46 genotypes is quantified and empirically validated. We find that the imputed genotypes are

47
48 most reliable for classical strains that have at least medium density genotyping data

1
2
3
4 available. The accuracy of imputed genotypes is somewhat lower in wild-derived strains.
5
6
7 We provide a confidence score that can be used to identify those imputed genotypes that
8
9 are most reliable.
10

11 **MATERIALS AND METHODS**

12 **Data preparation**

13
14
15
16
17
18 Prior to combining databases, a multi-step quality control procedure was applied to the
19
20 original NIEHS (<http://mouse.perlegen.com/mouse/download.html>, July 2006 release) and
21
22 Broad (<http://www.broad.mit.edu/~claire/MouseHapMap>, February 2006 release) SNPs.
23
24 First, we eliminated SNPs whose reported physical locations are impossible. We compared
25
26 the genotypes and the 100-mer flanking sequences of all C57BL/6 SNPs in the data to the
27
28 published mouse genome sequence NCBI build 36. This process remapped the NCBI build
29
30 33 Broad data to NCBI build 36 coordinates and removed those SNPs that revealed any
31
32 discrepancy. A total of 75 Broad and 2,925 NIEHS SNPs were excluded. We then
33
34 identified SNPs that are present in duplicated regions. We have previously observed that
35
36 these SNPs can have very high false positive rates due to the detection of paralogous
37
38 variation at other sites (Yang et al. 2007; unpublished data). We used BLAT (Kent 2002)
39
40 under the most sensitive parameter settings to map all 25-mers centered in each SNP and
41
42 defined a duplication as a SNP for which the 25-mer map to multiple genomic locations
43
44 with less than three mismatches. We then used a sliding window to search for clustered
45
46 duplications. Whenever two duplications were found out of four consecutive SNPs, the
47
48 duplicated SNPs and any intervening SNPs were removed. A total of 448,999 SNPs
49
50 (~5.4%) were removed. Of the remaining NIEHS SNPs, 15,068 were reported to have
51
52
53
54
55
56
57
58
59
60

1
2
3
4 same genomic location. We kept one copy of each when all genotypes were fully
5
6 consistent; otherwise, they were removed. For Broad data, we removed 2762 SNPs (~2%)
7
8 that mapped to the duplicated locations in the NIEHS data. SNPs that mapped to identical
9
10 genomic locations (redundant SNPs) were combined. During the process of combining
11
12 databases, strand orientation adjustment was done whenever necessary. Conflicting
13
14 genotypes were recoded as missing data. When more than half of the strains had discordant
15
16 genotypes, the SNP locus was excluded.
17
18
19

20 21 **Genotype Imputation** 22

23
24 We use a hidden Markov model (HMM) with left to right architecture (Figure 1) to impute
25
26 the missing genotypes. In this model, there are six hidden states ($H = 6$) representing
27
28 different haplotypes at each SNP. State transitions proceed from one SNP (columns in
29
30 figure 1) to the next according to a Markov process. The haplotypes of a strain can be
31
32 viewed as a path through the model visiting one state per SNP locus, from the first SNP to
33
34 the last on a given chromosome. Given a trained model and the genotypes of a strain, the
35
36 path decoding problem is solved by Viterbi's algorithm (Viterbi 1967). In Figure 1, the
37
38 Viterbi paths are shown as colored lines. Strains with identical path through this region are
39
40 grouped, but in general each strain will have its own unique path through the haplotype
41
42 states. States have a probabilistic output, representing the observed genotype. Missing
43
44 genotypes are imputed as the allele that is most likely to be emitted by the states along the
45
46 Viterbi path. The most probable genotype for each state is indicated in the Figure 1. For
47
48 every genotype, imputed and experimental, the posterior probability under the trained
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4 HMM serves as a confidence score, which is computed as the product of the posterior
5
6 probability of the inferred haplotype state and genotype probability given the state.
7
8

9
10 Training the HMM involves estimation of a large number of free parameters. Parameter
11
12 estimation is accomplished using the Expectation-Maximization (E-M) method as
13
14 elaborated in Churchill (1989). The convergence criterion for the E-M algorithm is set as
15
16 10^{-6} change in the log-likelihood. Initial values for the EM algorithm are sampled at least
17
18 10 times and the training run that achieves the highest likelihood is chosen.
19
20

21
22 Despite the vast amount of data (millions of genotypes) in the training sets, this is a data
23
24 poor problem. At any given SNP we have genotypes for only a small to moderate number
25
26 of strains, distinct haplotypes may not be equally represented, and the information
27
28 available in adjacent SNPs decays more or less rapidly depending on marker density and
29
30 the extent of local linkage disequilibrium. The number of parameters in the HMM is large
31
32 and it grows linearly with the number of SNPs; therefore the prior distributions can be
33
34 influential and should be chosen carefully to obtain the best results.
35
36
37

38
39 Transition and emission probabilities are assumed to follow Dirichlet prior distributions.
40
41 For state transitions, the prior density is biased towards the transitions between the same
42
43 haplotype, with probability $1 - \lambda$, and is equally distributed among the other $(H-1)$
44
45 haplotypes. A prior that favors small values of λ will encourage the use of more
46
47 information from adjacent SNPs. Emission probabilities are assumed to follow a uniform
48
49 prior distribution for the two possible alleles. The Dirichlet pseudocount method is used to
50
51 combine prior information with maximum likelihood estimates (Durbin et al. 1998).
52
53
54
55
56
57
58
59
60

1
2
3
4 A series of computational experiments was carried out to optimize the predictive accuracy
5 of the HMM. We varied both the number of haplotypes and the prior parameters and
6 assessed the accuracy of imputation by randomly masking portions of the genotype data.
7
8 When H is too small, accuracy declines but we saw little improvement for these data when
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

When H is greater than 5 or 6. In genomic regions with fewer than six distinct haplotypes, a subset of the states will typically have small marginal probabilities and are effectively unused. Based on these studies, we chose $H=6$ and a prior mean transition rate of 0.01 as optimal values for imputation in the merged NIEHS-Broad data (Figure 2).

RESULTS

Combining large scale SNP panels

In order to create the imputed genotype resource, we first merged the SNP genotypes in the NIEHS and Broad data. We refer to these data as the *merged set*. The relationship of the merged set to other SNP sets used in this study is summarized in Figure 2.

The NIEHS data (<http://mouse.perlegen.com/mouse/download.html>, July 2006 release), include 109 million genotypes on ~8.3 million SNPs spanning the 19 autosomes, the X and Y chromosomes, and the mitochondrial genome, for 11 classical and four wild-derived strains. The genotypes were generated by Perlegen Sciences using high density oligonucleotide arrays. Approximately 54% of the NIEHS SNPs have missing genotypes for one or more of the 15 strains summing to 12% incomplete genotypes. The frequency of missing data is higher in the three wild-derived strains of non-*domesticus* origin (CAST, MOLF and PWD). We removed 5.5% of SNPs from the initial NIEHS set due to potential

1
2
3
4 problems (see Materials and Methods and Table S1). The 7,804,762 remaining SNPs
5
6 spanning the autosomes and the X chromosome were used in this study.
7
8

9
10 The Broad data (<http://www.broad.mit.edu/~claire/MouseHapMap>, February 2006 release)
11
12 include over 6 million genotypes for 138,793 SNPs distributed at ~20kb intervals across
13
14 the autosomes and the X chromosome, for 38 classical and 11 wild-derived strains
15
16 representing four subspecies of *M. musculus* and two representatives of *M. spretus*.
17
18 Genotypes were generated using custom Affymetrix SNP array technology. Approximately
19
20 68% of the SNPs have missing genotypes for one or more strains summing to 8%
21
22 incomplete genotypes. The frequency of incomplete genotypes is higher for wild-derived
23
24 strains with the exception of the two wild-derived *M. m. domesticus* strains, WSB and
25
26 PERA. After mapping to NCBI build 36 coordinates and data preparation (see Materials
27
28 and Methods), 135,846 SNPs were retained for this study.
29
30
31
32

33
34 The 15 NIEHS strains are common to both datasets and the remaining 33 strains are unique
35
36 to the Broad data. Genotypes from the reference C57BL/6 genome sequence are included
37
38 in the merged data. In total there are 116 million genotypes for 7,870,134 SNPs spanning
39
40 the autosomes and the X chromosome of 49 strains.
41
42

43
44 In the process of merging datasets of dramatically different marker densities we assigned
45
46 all of the unavailable genotypes to be missing. Thus there are two types of missing
47
48 genotypes in the merged set. Experimental missing data are due to failure of a genotyping
49
50 assay. Missing data created as a result of merging the data have, for the most part, not been
51
52 directly assayed. In the merged set, 7,734,384 of the loci are assayed only in NIEHS data,
53
54 65,468 only in the Broad data and 70,282 loci have been assayed in both. There are a total
55
56
57
58
59
60

1
2
3
4 of 14,078,485 experimental missing genotypes, 255,234,672 missing genotypes created by
5
6 merging these data, and 6,318 missing values due to conflicts (see Methods). The
7
8 frequency of incomplete genotypes is 98.3% for the 33 strains unique to the Broad set and
9
10 ranges from 10% to 16% for the 15 strains common to both Broad and NIEHS data.
11
12
13

14 **The imputed genotypes**

15
16
17 We implemented a hidden Markov model (HMM) with a left to right architecture (Figure 1)
18
19 for the primary purpose of genotype imputation and for the secondary purpose of haplotype
20
21 identification. The architecture is similar to those described in Kimmel and Shamir (2005),
22
23 and in Scheet and Stephens (2006). The number of hidden states at each SNP locus and the
24
25 prior distribution of the model parameters of the HMM were optimized for genome-wide
26
27 imputation accuracy (see Materials and Methods). All 49 strains in the merged set were used
28
29 to train the model (see Materials and Methods). The posterior probability of each imputed
30
31 genotype under the trained model provides a confidence score. A total of 269,319,475
32
33 missing genotypes were imputed in the merged set (*merged-imputed set* in Figure 2), of
34
35 which 14.9%, 15.0%, and 70.1%, fall into the low, medium and high confidence score bins of
36
37 (0,0.6), (0.6,0.9) and (0.9,1), respectively (Table S2).
38
39
40
41
42
43

44 In order to assess the quality of the imputed genotypes, we assembled a *validation set* of
45
46 genotypes reported in eight SNP resources developed independently from the NIEHS and
47
48 Broad sets (Figure 2, Table 1, Table S3). A total of 969,457 imputed genotypes could be
49
50 validated using these resources, of which 15.4%, 13.7%, and 70.8%, fall into the low
51
52 medium and high confidence score bins. The number of validated genotypes varies
53
54
55
56
57
58
59
60

1
2
3
4 substantially among inbred strains (Table 2) and seven strains (129S4/SvJae, DDK/Pas,
5
6 MAI/Pas, O20, Qsi5, ST/bJ and SEG/Pas) have no genotypes available for validation.
7
8

9
10 We compared imputed genotypes to genotypes in the validation set and conservatively
11
12 assume that discordant genotypes represent imputation errors. The overall imputation error
13
14 rate based on comparison with the validation set is 0.104 (Table 2). Error rates vary
15
16 substantially among strains and to a lesser degree across chromosomes (Table S2). Strain
17
18 specific error rates are lower for classical strains than for wild-derived strains.
19
20

21
22 Furthermore, error rates vary for the two different types of missing data (Table 2). For
23
24 imputed genotypes with high confidence scores, the error rate is 0.044 (Table 2). Among
25
26 validated genotypes with high confidence scores wild-derived strains have higher error
27
28 rates than classical strains. The NIEHS strains have higher error rates because most of the
29
30 missing genotypes are experimental.
31
32

33
34 As a consequence of high levels of divergence between the classical and wild-derived
35
36 strains, 59% of SNPs in the merged set are private to the wild-derived strains (i.e., the
37
38 genotypes of those SNPs are constant within classical strains). We estimated error rates
39
40 stratified by the status of a SNP being constant or polymorphic within the classical strains
41
42 and found that they were essentially identical (Table S4). We note that, for strains A,
43
44 DBA/2 and 129S1, the error rates within the constant SNPs (26% of the total validated
45
46 SNPs) were elevated compared to the unstratified version (Table 2). These strains have a
47
48 large number of validation genotypes available (Mural et al. 2002). This interesting pattern
49
50 and the higher experimental error rates in the NIEHS strains suggest that gene conversion
51
52
53
54
55
56
57
58
59
60 may be responsible for a large fraction of the imputation error.

1
2
3
4 To test how wild-derived strains impact the imputation accuracy among the classical
5 strains, we retrained the HMM on a subset of the merged data including only the 38
6 classical strains and estimated the imputation error rates by comparison with the validation
7 data. The impact on error rates varies across chromosomes and is most evident in
8 chromosomes where there is a substantial contribution to the classical strains from *M. m.*
9 *musculus* (Yang et al. 2007). The higher errors observed overall suggest that the inclusion
10 of the wild-derived strains improves the imputation of missing genotypes in regions where
11 some classical inbred strains are not of *M. m. domesticus* origin, without negatively
12 impacting the rest.
13
14
15
16
17
18
19
20
21
22
23
24

25
26 Based on this empirical validation, we conclude that the quality of the imputed genotypes
27 improves as information from more strains is utilized in training the HMM, and that the
28 imputation of missing genotypes is highly reliable for most strains.
29
30
31
32
33

34 **Imputation of genotypes in other mouse strains**

35
36 The trained HMM can be used to infer high density genotypes for strains that are not
37 included in the training set using the Viterbi algorithm (Viterbi, 1967). Strains NZO/HILtJ
38 and PWK/PhJ are not in the merged set but are of interest because they are founder strains
39 of the Collaborative Cross (Churchill et al. 2004). NZO is a classical inbred strain closely
40 related to strain NZB of the merged set. PWK is a wild derived *M. musculus* strain that is
41 most closely related to strain PWD in the merged set, although the overall sequence
42 divergence between PWK and PWD is greater than between any pair of classical strains.
43
44 We have assembled 140,269 genotypes and 132,862 genotypes for NZO and PWK,
45 respectively, from four independent resources (Table S3; Tim Wiltshire, personal
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4 communication) and created a medium density SNP set by retaining only genotypes that
5
6 correspond to Broad loci in the merged set. The remaining SNPs were used for validation
7
8 (Table 3). Similarly, we created low density versions of NZO and PWK SNPs by selecting
9
10 SNP loci that correspond to loci genotyped in the Wellcome-CTC study (Shifman et al,
11
12 2007). We ran the Viterbi algorithm and imputed missing genotypes for both strains at
13
14 medium and low density. At medium density, NZO imputations yield a large proportion of
15
16 high confidence SNPs; but at low density, confidence scores shift substantially downward.
17
18 At medium density about one third of the imputed PWK SNPs have low confidence and at
19
20 low density, nearly 80% of imputed SNPs fall into the lowest confidence category. The
21
22 low confidence in the PWK imputations is likely due to the presence of *M .m. musculus*
23
24 haplotypes in PWK that are not represented in the NIEHS strains. The validation accuracy
25
26 of the imputed SNPs (Table 3) follows the same pattern as the confidence scores. Although
27
28 the overall validation error rate of 38% for low density PWK imputations is unacceptably
29
30 high, but for those SNPs that achieve the highest confidence scores (>0.9), chromosome
31
32 specific error rates range from 1% to 8%. These results suggest that additional strains with
33
34 medium density SNP genotyping can be accurately imputed and that the best overall
35
36 results will be achieved for classical strains.
37
38
39
40
41
42
43
44

45 **The imputed genotype resource**

46
47
48 To generate the most accurate imputation of genotypes on the 49 inbred mouse strains, we
49
50 added the 969,457 SNP genotypes in the validation set to the merged set and created a
51
52 *combined set* of SNP genotypes (Figure 2). We then retrained the HMM using 49 strains in
53
54 the combined set and imputed the missing genotypes. The confidence score distribution of
55
56
57
58
59
60

1
2
3
4 imputed genotypes from the combined set demonstrated a slight shift towards the higher
5
6 confidence score bins, indicating increased accuracy with 14.1%, 14.7%, 71.2% of all
7
8 imputed genotypes falling into the low, medium and high confidence score bins,
9
10 respectively. In addition, we recomputed the imputation for two strains, NZO and PWK,
11
12 using all of the available genotypes. These strains were not included in the training set. The
13
14 complete set of imputed genotypes and confidence scores for 51 strains are available at
15
16 <http://cgd.jax.org/>. The data can be downloaded or queried through a MYSQL database.
17
18 The SNPs have been quality-checked, and cross-linked with other features, including
19
20 ENSEMBL annotations, GO annotations, and MGI gene and phenotype information. We
21
22 created a web interface that allows SNP retrieval filtered on features such as neighboring
23
24 genes, genomic location, SNP functional implication, CpG sites, and substitution types.
25
26
27
28
29
30

31 DISCUSSION

32
33
34 We have created a data resource of experimental and imputed SNP genotypes at a density
35
36 of 7,870,134 loci on 49 commonly used inbred mouse strains by combining data and
37
38 imputing missing genotypes from two major public SNP collections. Our results support
39
40 the hypothesis that strains with medium density genotyping can be accurately imputed to
41
42 obtain high density genotypes. Confidence scores assigned to each genotype reflect the
43
44 reliability of imputed genotypes and identify experimental genotypes that depart from
45
46 expectations based on local LD. We have demonstrated the accuracy of imputed genotypes
47
48 by comparison to experimental genotypes obtained independently.
49
50
51

52
53 Imputed genotypes are not a replacement for experimental measurements. We encourage
54
55 investigators to use this resource as an exploratory tool, but critical conclusions based on
56
57
58
59
60

1
2
3
4 imputed genotypes, should be validated. Low confidence imputations are more common in
5
6 the wild-derived strains due the limited representation of appropriate taxa in the NIEHS
7
8 strain panel. Nonetheless, the reliability of most genotypes in this resource is sufficient for
9
10 high throughput analysis and hypothesis generating.
11
12

13
14 Extensive local LD, reflecting the small number of founders and the presence of admixture
15
16 in laboratory mice, provides the essential structure that allows accurate imputation of
17
18 missing genotypes. To achieve this, the density of SNP genotyping should be sufficient to
19
20 tag most regions of local LD in the population of strains to which the imputation algorithm
21
22 is being applied. Imputation may be unreliable if a novel haplotype, not present in the high
23
24 density training data, is encountered. Furthermore, SNPs that were not identified in the
25
26 discovery process are not represented in the imputed data resource. We have previously
27
28 estimated the false negative rate in the NIEHS data to be 67% (the false negative rate in the
29
30 classical strains is 43%) but the rate is significantly higher for singleton SNPs (Yang et al.
31
32 2007). Therefore, absence of a SNP in this, or any, resource does not imply genetic
33
34 identity. False negatives are of concern when the missing SNPs are private to one strain or
35
36 to a small group of related strains. Discovery bias will significantly impact our ability to
37
38 accurately impute SNPs in inbred strains derived from diverse mouse lineages. The
39
40 solution is to carry out more SNP discovery at high density. The report of the sequence of
41
42 multiple *Drosophila* species highlights the importance and benefits of resequencing and
43
44 SNP discovery in diverse taxa (*Drosophila* 12 Genomes Consortium 2007). Similarly, it
45
46 will be particularly useful to identify lineage specific SNPs and to genotype additional
47
48 representatives of lineages with high error rates such as *M. m. castaneus* and *M. spretus*.
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5 During the data preparation phase of this project we intentionally removed SNPs from
6
7 regions with evidence of multiple copies in the C57BL/6 genome due to higher error rates
8
9 observed in these SNPs (unpublished) and there are variable repeat regions present in other
10
11 strains. Repeats and copy number variation need to be included to achieve a complete
12
13 understand of the landscape of genetic variation in the laboratory mouse.
14
15

16
17 Other important features of genetic variation, such as gene conversion and recurrent
18
19 mutations, may be missed because they are not consistent with the local LD pattern. In
20
21 fact, the patterns of errors observed in the Celera strains suggest that both processes
22
23 contribute significantly to the imputation error. This situation can be solved through
24
25 additional experimental determination of missing genotypes. Finally, genotyping errors in
26
27 the training data present a major barrier to improving the accuracy of imputed genotypes.
28
29 Identification and resolution of genotyping errors in data as extensive as these is a daunting
30
31 task. Confidence scores obtained from the HMM can suggest potential genotyping errors,
32
33 but not all genotyping errors can be detected in this way. New computational approaches to
34
35 error detection would be beneficial.
36
37
38
39

40
41 We used an empirical approach to validate imputed genotypes. Our estimated error rates
42
43 are likely to be conservative because of genotyping errors in the validation set. An
44
45 alternative method to assess imputation accuracy (Roberts et al. 2007) is to randomly mask
46
47 a proportion of the genotypes within the combined dataset and to compare the imputed
48
49 values to the masked genotypes. We found that the random masking approach provides
50
51 higher estimated accuracy compared to the empirical validation. The processes that lead to
52
53
54
55
56
57
58
59
60

1
2
3
4 missing data are likely to be more complex than the masking model, thus we prefer the
5
6 empirical estimates but acknowledge their conservative bias.
7
8

9
10 The imputed genotype resource must be viewed dynamically because the coverage of
11
12 strains, the number of SNP loci, and the accuracy of imputation will improve with
13
14 additional data from ongoing genotyping projects. In conclusion, this study reports a
15
16 method for accurate imputation of missing genotypes and its use in generating a dense map
17
18 of the genetic variation in the mouse genome. Our results support the proposal by Frazer
19
20 and coworkers (2007) that such a resource could and must be generated.
21
22
23
24
25
26

27 **ACKNOWLEDGMENTS**

28
29 This work was supported by the US National Institutes of General Medical Sciences as
30
31 part of the Center of Excellence in Systems Biology (1P50 GM076468). We thank Tim
32
33 Wiltshire for sharing genotyping data prior to its publication, Jesse Hammer and Susan
34
35 Moxley for graphics assistance.
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

REFERENCES

Abe K, Noguchi H, Tagawa K, Yuzuriha M, Toyoda A, et al. (2004) Contribution of Asian mouse subspecies *Mus musculus molossinus* to genomic constitution of strain C57BL/6J, as defined by BAC-end sequence-SNP analysis. *Genome Res* 14: 2439-2447.

Bogue MA (2003) Mouse Phenome Project: understanding human biology through mouse genetics and genomics. *J Appl Physiol* 95: 1335–1337.

Cervino AC, Li G, Edwards S, Zhu J, Laurie C, et al. (2005) Integrating QTL and high-density SNP analyses in mice to identify *Insig2* as a susceptibility gene for plasma cholesterol levels. *Genomics* 86: 505-517.

Churchill GA, Airey DC, Allayee H, Angel JM, Attie AD, et al. (2004) The Collaborative Cross, a community resource for the genetic analysis of complex traits. *Nat Genet* 36: 1133-1137.

Churchill GA (1989) Stochastic models for heterogeneous DNA sequences. *Bulletin of Mathematical Biology* 51: 79-94.

DiPetrillo K, Wang X, Stylianou L, and Pagien B (2005) Bioinformatics toolbox for narrowing rodent quantitative trait loci . *Trends Genet* 21: 684-692.

Durbin R, Eddy SR, Krogh A, and Mitchison G (1998) *Biological sequence analysis*. Cambridge University Press, Cambridge, UK.

Drosophila 12 genomes consortium. (2007) Evolution of genes and genomes on the *Drosophila* phylogeny. *Nature* 450: 203-218.

1
2
3
4
5 Frazer KA, Wade CM, Hinds DA, Patil N, Cox DR, et al.(2004) Segmental phylogenetic
6 relationships of inbred mouse strains revealed by fine-scale analysis of sequence variation
7 across 4.6 Mb of mouse genome. *Genome Res* 14: 1493-1500.

8
9
10
11 Frazer KA, Eskin E, Kang HM, Bogue MA, Hinds DA, et al. (2007) A sequence-based
12 variation map of 8.27 million SNPs in inbred mouse strain. *Nature* 448:1050-1053.

13
14
15
16 Guenet JL, and Bohomme F (2003) Wild mice: an ever-increasing contribution to a
17 popular mammalian model. *Trends Genet* 19: 24-31.

18
19
20
21 Ideraabdullah FY, de la Casa-Esperon E, Bell TA, Detwiler DA, Magnuson T, et al. (2004)
22 Genetic and haplotype diversity among wild derived mouse inbred strains. *Genome Res*
23
24
25
26
27 14: 1880-1887.

28
29
30
31 Kent WJ (2002) BLAT the BLAST-like alignment tool. *Genome Res* 12: 656-664.

32
33
34
35 Kimmel G, and Shamir R (2005) A block-free hidden Markov model for genotypes and its
36 application to disease association. *J Comput Biol* 12: 1243.

37
38
39
40 Liao G, Wang J, Guo J, Allard J, Cheng J, et al. 2004. In silico genetics: identification of a
41 functional element regulating H2-Ealpha gene expression. *Science* 306: 690-695.

42
43
44
45 Lyon MF, Rastan S, and Brown SDM. (ed.) (1996) Genetic variants and strains of the
46 laboratory mouse. 3rd edition, Oxford Univeristy Press, Oxford, UK.

47
48
49
50 McClurg P, Janes J, Wu C, Delano DL, Walker JR, et al. (2007) Genomewide association
51 analysis in diverse inbred mice: power and population structure. *Genetics* 176: 675-683.

52
53
54
55
56
57
58
59
60 Mott R (2007) A haplotype map for the laboratory mouse. *Nat Genet* 39: 1054-1056.

1
2
3
4 Mural RJ, Adams MD, Myers EW, Smith HO, Miklos GL, et al. (2002) A comparison of
5 whole-genome shotgun-derived mouse chromosome 16 and the human genome. *Science*
6
7 296: 1661-1671.
8
9

10
11 Payseur BA, and Hoekstra HE (2005) Signatures of reproductive isolation in patterns of
12 single nucleotide diversity across inbred strains of mice. *Genetics* 171: 1905-1016.
13
14

15
16 Payseur BA, and Place M (2007) Prospects for association mapping in classical inbred
17 mouse strains. *Genetics* 175: 1999-2008.
18
19

20
21 Pletcher MT, McClurg P, Batalov S, Su AI, Barnes SW, et al. 2004. Use of a dense single
22 nucleotide polymorphism map for in silico mapping in the mouse. *PLoS Biol* 2: 2159 -2169
23
24

25
26 Petkov PM, Graber JH, Churchill GA, DiPetrillo K, King BL, et al. (2005) Evidence of a
27 large-scale functional organization of mammalian chromosomes. *PLoS Genet* 1: e33.
28
29

30
31 Petkov PM, Ding Y, Cassell MA, Zhang W, Wagner G, et al. (2004) An efficient SNP
32 system for mouse genome scanning and elucidating strain relationships. *Genome Res* 14:
33 1806-1811.
34
35
36

37
38 Roberts A, McMillan L, Wang W, Parker J, Rusyn I, et al. (2007) Inferring missing
39 genotypes in large SNP panels using fast nearest-neighbor searches over sliding windows.
40
41
42 *Bioinformatics* 23: i401.
43

44
45 Roberts A, Pardo-Manuel de Villena F, Wang W, McMillan L, Threadgill DW (2007) The
46 polymorphism architecture of mouse genetic resources elucidated using genome-wide
47 resequencing data: implications for QTL discovery and systems genetics. *Mamm Genome*
48
49 18: 473-481
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4 Siebert SK, and Schadt EE (2007) Moving toward a system genetics view of disease.

5
6
7 Mamm Genome 18: 389-401.

8
9 Scheet P, and Stephens M (2006) A fast and flexible statistical model for large-scale
10
11 population genotype data: applications to inferring missing genotypes and haplotypic
12
13 phase. Am J Hum Genet 78: 129.

14
15
16 Shifman S, Bell JT, Copley RR, Taylor MS, Williams RW, et al. (2006) A high-resolution
17
18 single nucleotide polymorphism genetic map of the mouse genome. PLoS Biol 4: e395.

19
20
21 Viterbi A J (1967) Error bounds for convolutional codes and an asymptotically optimum
22
23 decoding algorithm, IEEE Trans Information Theory 13: 260-269.

24
25
26 Wade CM, Kulbokas EJ 3rd, Kirby AW, Zody MC, Mullikin JC, et al. (2002) The mosaic
27
28 structure of variation in the laboratory mouse genome. Nature 420: 574-578.

29
30
31 Wade CM, and Daly MJ (2005) Genetic variation in laboratory mice. Nat Genet 37: 1175-
32
33 1180.

34
35
36 Waterston RH, Lindblad-Toh K, Birney E, Rogers J, Abril JF, et al. (2002) Initial
37
38 sequencing and comparative analysis of the mouse genome. Nature 420: 520-562.

39
40
41 Wiltshire T, Pletcher MT, Batalov S, Barnes SW, Tarantino LM, et al. 2003. Genome-wide
42
43 single-nucleotide polymorphism analysis defines haplotype patterns in mouse. Proc Natl
44
45 Acad Sci USA 100: 3380-3385.

46
47
48 Yalcin B, Fullerton J, Miller S, Keays DA, Brady S, et al. (2004) Unexpected complexity
49
50 in the haplotypes of commonly used inbred strains of laboratory mice. Proc Natl Acad Sci.
51
52 USA 101: 9734-9739.

1
2
3
4 Yang H, Bell TA, Churchill GA, Pardo-Manuel de Villena F (2007) On the subspecific
5
6 origin of the laboratory mouse. Nat Genet 39: 1100-1107.
7
8
9
10

11 **WEB SITE REFERENCES**

12 [13 http://www.broad.mit.edu/~claire/MouseHapMap/](http://www.broad.mit.edu/~claire/MouseHapMap/)

14 [15 http://mouse.perlegen.com/](http://mouse.perlegen.com/)

16 [17 http://www.well.ox.ac.uk/mouse/INBREDS/](http://www.well.ox.ac.uk/mouse/INBREDS/)

18 [19 http://www.sanger.ac.uk/modelorgs/mouse.shtml/](http://www.sanger.ac.uk/modelorgs/mouse.shtml/)

20 [21 http://www.ensembl.org/](http://www.ensembl.org/)

22 [23 http://www.ncbi.nlm.nih.gov/SNP/](http://www.ncbi.nlm.nih.gov/SNP/)

24 [25 http://snp.gnf.org/](http://snp.gnf.org/)

26 [27 http://phenome.jax.org/](http://phenome.jax.org/)

28 [29 http://stt.gsc.riken.jp/msm/](http://stt.gsc.riken.jp/msm/)

30 [31 http://mousesnp.roche.com/](http://mousesnp.roche.com/)

32 [33 http://cgd.jax.org/](http://cgd.jax.org/)

Table 1. A summary of SNP resources. For each SNP resource, the name, the number of SNPs, the number of the strains reported, the genotyping technology used, and a literature citation are shown.

Resource	No. SNP	No. Strain	Genotyping technology	Web site	Reference
NIEHS	8,272,574	16	Perlegen Science oligonucleotide arrays	mouse.perlegen.com	Frazer et al. (2007)
Broad	138,608	49	Affymetrix SNP array	www.broad.mit.edu/~claire/MouseHapMap	Wade and Daly (2005)
Celera	2,093,327	5	Whole genome shotgun assembly	www.celera.com	Mural et al. (2002)
GNF	9,594	48	PCR and DNA sequencing	snp.gnf.org	Wiltsire et al. (2003)
TJL	1,638	144	Amplifluor technology	aretha.jax.org	Petkov et al. (2004)
Japanese MSM	210,160	1	BAC-end sequence analysis	stt.gsc.riken.jp/msm/	Abe et al. (2004)
Wellcome-CTC	13,348	499	Illumina SNP array	www.well.ox.ac.uk/mouse/INBREDS	Shifman et al. (2006)
Sanger	748,723	7	whole genome shotgun, clone based sequencing	www.ensembl.org; www.sanger.ac.uk	Waterston et al. (2002)
db-SNP	11,814	7	Central repository	www.ncbi.nlm.nih.gov/SNP	
Rosetta/Merck	12,473	62	Illumina	www.rii.com	Cervino et al. (2005)
Roche	214,706	20	Sequencing	mousesnp.roche.com/	Liao et al (2004)

1
2
3 **Table 2.** Estimated strain specific error rates in the merged-imputed set. The table
4 provides (1) the total number of imputed missing genotypes, the percentage of imputed
5 genotypes that were missing due to merging datasets; (2) the total number of validated
6 imputed genotypes, the error rate of validated imputed genotypes, the error rates for
7 imputed genotypes that were experimental missing data and for imputed genotypes that
8 were missing due to merging datasets; (3) the number and error rates of validated
9 imputed genotypes with confidence scores greater than 0.9. The table is divided to
10 indicate classical (top) versus wild derived (bottom) strains. NIEHS strains are listed first
11 within these categories.
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table 2.

Strains	ALL MISSING DATA IMPUTED		VALIDATED IMPUTED				VALIDATED IMPUTED CONFIDENCE >0.9			
	# missing genotype imputed	% missing data due to merge	# validated imputed	error for ALL missing data	error for experimental missing	error for missing due to merge	# validated imputed conf>0.9	error for ALL missing data	error for experimental missing	error for missing due to merge
DBA/2J	887413	0	70678	0.156	0.156	-	50734	0.077	0.077	-
A/J	738284	0	66222	0.160	0.160	-	47485	0.078	0.078	-
129S1/SvImJ	860002	0	26846	0.168	0.168	-	19284	0.089	0.089	-
C3H/HeJ	921492	0	1337	0.147	0.147	-	944	0.058	0.058	-
BTBR.T<+>tf/J	845230	0	1256	0.143	0.143	-	926	0.054	0.054	-
FVB/NJ	852011	0	1212	0.152	0.152	-	850	0.069	0.069	-
KK/HLJ	845405	0	1209	0.168	0.168	-	816	0.082	0.082	-
NOD/LtJ	871145	0	1190	0.155	0.155	-	828	0.070	0.070	-
AKR/J	836285	0	1168	0.170	0.170	-	809	0.079	0.079	-
NZW/LacJ	846575	0	1160	0.147	0.147	-	812	0.074	0.074	-
BALB/cByJ	795150	0	7	0.000	0.000	-	7	0.000	0.000	-
129X1/SvJ	7740291	98.3%	485452	0.076	0.060	0.076	364089	0.030	0.018	0.030
SM/J	7753890	98.3%	10695	0.108	0.102	0.109	7108	0.048	0.041	0.048
NZB/BINJ	7745598	98.3%	10456	0.103	0.098	0.103	7075	0.040	0.021	0.040
NON/LtJ	7747432	98.3%	10447	0.086	0.113	0.086	7081	0.027	0.036	0.027
SJL/J	7743333	98.3%	10424	0.089	0.086	0.089	7060	0.031	0.007	0.032
CBA/J	7742711	98.3%	10423	0.077	0.039	0.078	7276	0.022	0.019	0.022
CE/J	7745164	98.3%	10421	0.124	0.132	0.124	6892	0.058	0.056	0.059
BUB/BnJ	7742202	98.3%	10416	0.090	0.067	0.090	7111	0.030	0.028	0.030
LG/J	7743879	98.3%	10413	0.094	0.123	0.093	7098	0.038	0.038	0.038
SWR/J	7742958	98.3%	10407	0.097	0.111	0.097	7069	0.041	0.034	0.041
LP/J	7742766	98.3%	10387	0.093	0.085	0.093	7096	0.030	0.030	0.030
C58/J	7741022	98.3%	10379	0.098	0.123	0.098	7131	0.034	0.024	0.034
C57BR/cdJ	7742389	98.3%	10366	0.096	0.059	0.097	7195	0.034	0.000	0.035
I/LnJ	7743740	98.3%	10363	0.095	0.101	0.095	7092	0.039	0.014	0.040
PL/J	7741981	98.3%	10332	0.087	0.076	0.088	7092	0.029	0.009	0.029
RIIS/J	7743061	98.3%	10331	0.104	0.131	0.104	6862	0.039	0.040	0.039
C57L/J	7743799	98.3%	7996	0.111	0.105	0.111	5266	0.033	0.008	0.034
MA/MyJ	7743694	98.3%	7973	0.108	0.119	0.107	5232	0.036	0.042	0.035
SEA/GnJ	7743682	98.3%	7971	0.098	0.078	0.099	5272	0.024	0.000	0.025
C57BLKS/J	7741314	98.3%	7958	0.084	0.092	0.084	5668	0.026	0.026	0.026
DBA/1J	7742347	98.3%	7926	0.103	0.110	0.103	5292	0.026	0.020	0.027
CAST/EiJ	1170222	0	1877	0.295	0.295	-	603	0.250	0.250	-
MOLF/EiJ	1137192	0	1764	0.190	0.190	-	887	0.106	0.106	-
WSB/EiJ	876918	0	1121	0.241	0.241	-	667	0.186	0.186	-
PWD/Ph	1211132	0	520	0.212	0.212	-	238	0.134	0.134	-
MSM/Ms	7750801	98.3%	79953	0.133	0.165	0.133	45522	0.052	0.122	0.052
PERA/EiJ	7748601	98.3%	10498	0.199	0.196	0.199	6487	0.130	0.134	0.130
SPRET/EiJ	7765710	98.3%	8909	0.219	0.238	0.218	5084	0.193	0.175	0.194
C2ECHII/EiJ	7749890	98.3%	7789	0.114	0.167	0.112	4846	0.070	0.131	0.067
JF1/Ms	7751232	98.3%	3205	0.070	0.119	0.068	1966	0.033	0.082	0.030
Genome-wide	215,077,943		969,457	0.104	0.158	0.090	686,852	0.044	0.078	0.036

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

Table 3. Estimated error rates for imputed genotypes using SNP genotyping data at different densities.

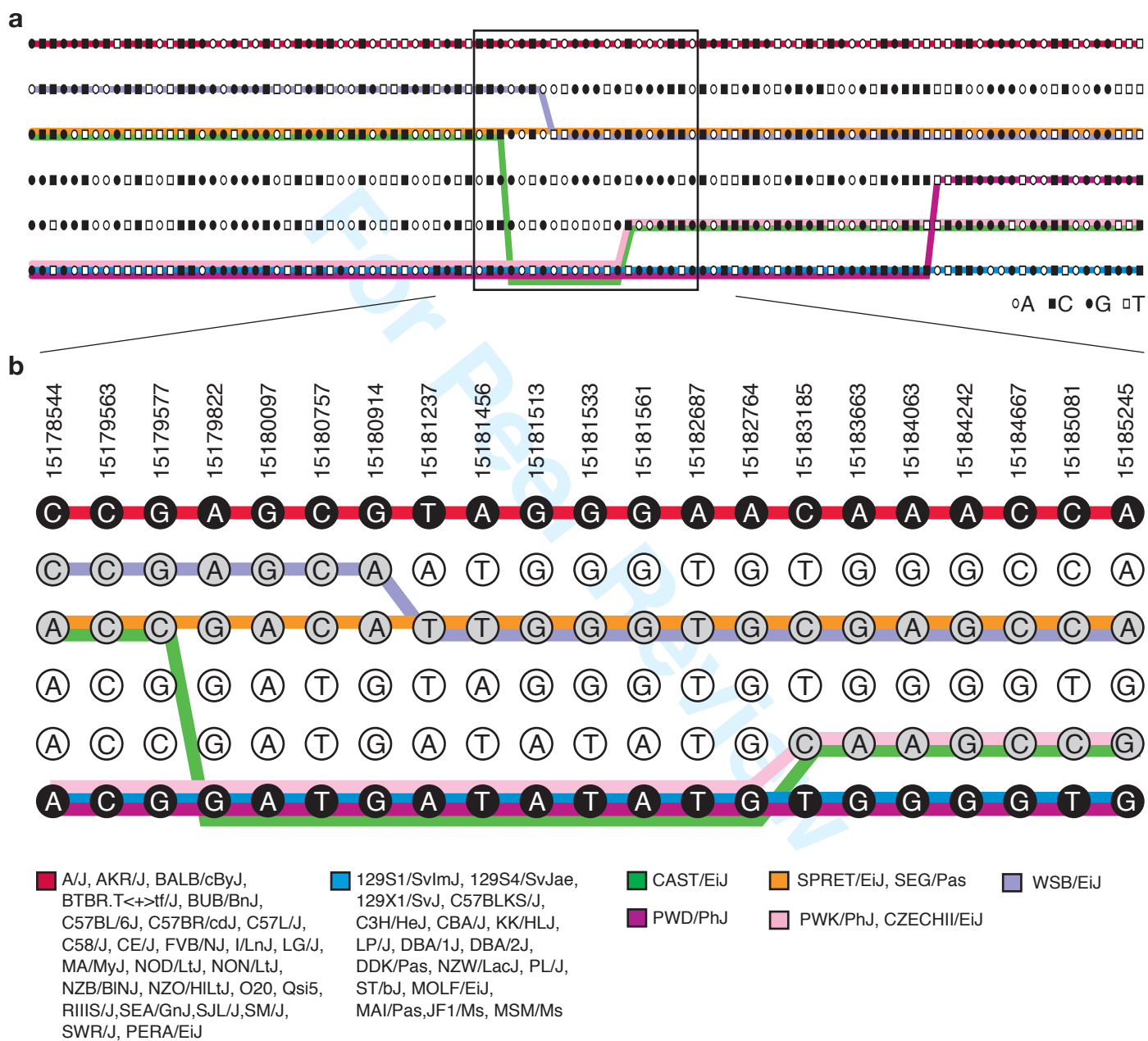
The table provides (1) the total number of genotyped SNPs; the proportion of the merged data; (2) the number of imputed genotypes, (3) the number of validated imputed genotypes and error estimate; (4) the number of validation imputed genotypes and error estimates for confidence score bins of (0,0.6), (0.6,0.9), (0.9,1).

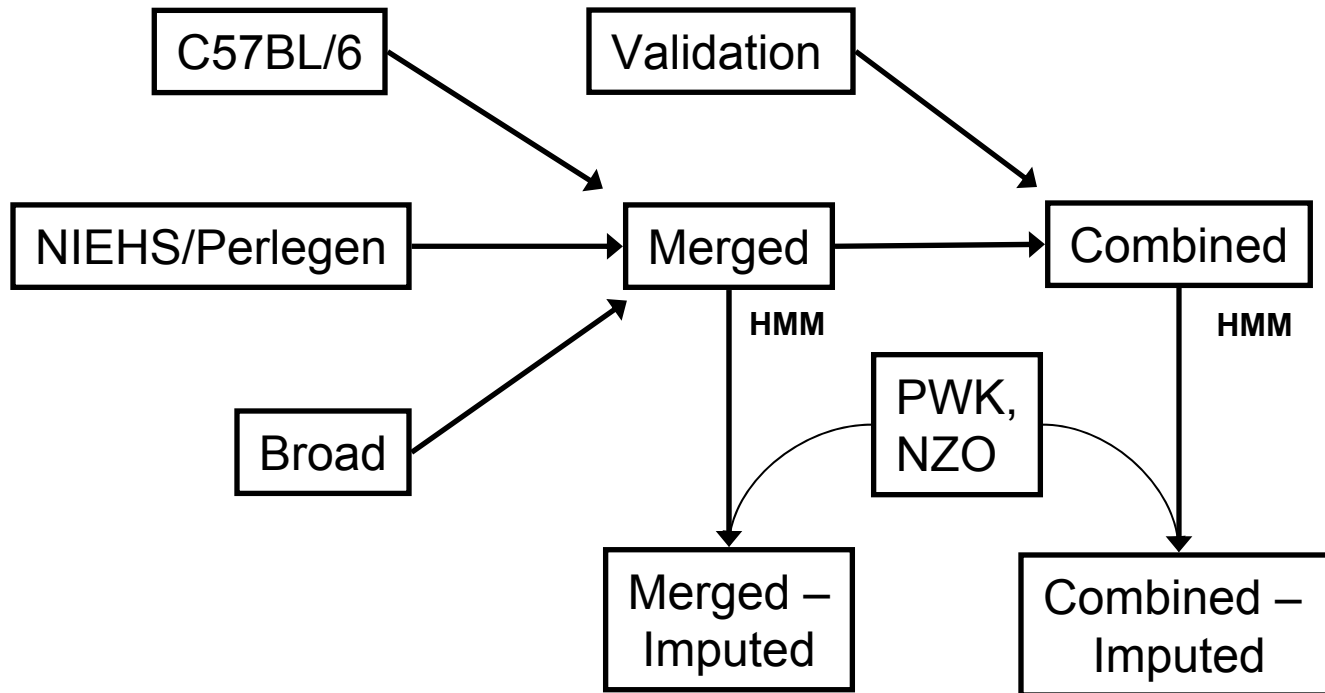
Strain	# genotyped snp	% of total snp in merged set	Number Imputed	Number validated.	Error validated	Number validated in.confidence.score.bin			error. in.confidence.score.bin		
						(0,0.6)	(0.6,0.9)	(0.9,1)	(0,0.6)	(0.6,0.9)	(0.9,1)
NZO.medium	130,387	1.7%	7,739,747	9,882	0.092	1,387	1,822	6,673	0.327	0.149	0.028
NZO.low	8,491	0.1%	7,861,643	131,778	0.159	50,220	45,532	36,026	0.336	0.080	0.011
PWK.medium	123,709	1.6%	7,746,425	9,153	0.086	2,408	1,474	5,271	0.206	0.091	0.030
PWK.low	8,085	0.1%	7,862,049	124,777	0.378	89,148	23,609	12,020	0.481	0.161	0.032

FIGURE LEGENDS

Figure 1. HMM architecture. Each SNP locus is modeled using six hidden states representing haplotypes. Each state is labeled by the nucleotide that is most likely to be observed in that haplotype. Colored lines represent the most probable haplotypes (Viterbi paths) for strains shown at the bottom. (a) A 40kb interval on chromosome 14 spanning 105 SNPs. (b) Detailed view of a 6.7 kb segment. Shading of the state nodes corresponds to the marginal state probability with darker color indicating haplotypes that are well represented in the training data.

Figure 2. Relationships among SNP data sources. The sequence of C57BL/6J provides reference genotypes for all SNPs in this study. The NIEHS data on 15 strains and Broad data on 48 strains were combined to create a merged set of experimental SNPs. The HMM was trained on the merged set and used to create the merged-imputed set. The validation set of experimental SNPs was assembled and curated from sources listed in Supplemental Table S3 and compared to the merged-imputed set in our validation study. The validation set was then combined with the merged SNPs to create the combined set of experimental SNPs. The HMM was retrained on the combined set to generate the combined-imputed set. SNP data from strains PWK and NZO were threaded through both of the trained HMMs. The results from threading with the merged-imputed model were used in our validation study. The PWK and NZO SNP data were not used in the training of either HMM.





1
2
3 **SUPPLEMENTARY MATERIALS for:**
4
5
6
7

8 **An imputed genotype resource for the laboratory mouse**
9

10
11
12
13
14
15 Jin P. Szatkiewicz¹, Glen L. Beane¹, Yueming Ding¹, Lucie Hutchins¹, Fernando Pardo-
16 Manuel de Villena², Gary A. Churchill¹
17
18
19

20
21
22
23 1. The Jackson Laboratory, Bar Harbor, Maine 04609, USA. 2. Department of Genetics,
24 Carolina Center for Genome Sciences and Lineberger Comprehensive Cancer Center,
25 University of North Carolina, Chapel Hill, Chapel Hill, North Carolina 27599, USA.
26
27
28
29
30
31

32
33 **Corresponding author:**
34

35 Gary A. Churchill
36

37 The Jackson Laboratory
38

39 Bar Harbor
40

41
42 Maine 04609, USA
43

44 Email: gary.churchill@jax.org
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Table S1. Summary of NIEHS, Broad, merged, and validation data sets by chromosome.

Chromosome	NIEHS Data			Broad Data			NIEHS & Broad		Merged Data				Validation Data
	NIEHS_original: No. SNP	NIEHS_clean: No. SNP	NIEHS_clean: SNP.density/100kb	Broad_original: No. SNP	Broad_clean: No. SNP	Broad_clean: SNP.density/100kb	No. shared-SNPs	No. shared-SNP eliminated	Merged: No. SNP	Merged: SNP.density/100kb	Total.missing.genotype	Rate.of.missing.genotypes	No. genotypes for SNPs in Merged
1	725,333	689,491	349.9	12096	11,858	6	6,531	9	694,809	352.6	23,719,388	0.697	92,427
2	536,890	519,405	285.4	10438	10,287	5.7	5,020	5	524,667	288.3	18,074,289	0.703	59,530
3	529,628	505,357	316.1	9665	9,492	5.9	4,948	9	509,892	318.9	17,516,977	0.701	70,360
4	496,883	472,442	304.7	8666	8,530	5.5	4,539	8	476,425	307.3	16,322,244	0.699	66,798
5	520,144	493,182	324.5	8473	8,292	5.5	4,578	8	496,888	326.9	16,945,501	0.696	63,587
6	545,009	506,246	338.6	7321	7,157	4.8	3,851	7	509,545	340.8	17,426,831	0.698	54,962
7	453,850	401,803	276.8	8132	7,866	5.4	3,933	3	405,733	279.6	13,812,494	0.695	61,438
8	460,703	441,256	334.1	7865	7,733	5.9	4,070	9	444,910	336.8	15,245,102	0.699	59,940
9	373,007	357,954	288.7	7253	7,133	5.8	3,514	2	361,571	291.6	12,306,386	0.695	47,836
10	415,666	396,905	305.4	4710	4,628	3.6	2,405	2	399,126	307.1	13,718,592	0.701	38,392
11	263,447	254,467	208.9	7709	7,542	6.2	2,978	3	259,028	212.7	8,903,563	0.701	40,682
12	419,456	393,188	326.4	6438	6,341	5.3	3,412	3	396,114	328.8	13,564,253	0.699	54,140
13	424,242	397,099	329.2	6333	6,181	5.1	3,346	4	399,930	331.6	13,676,989	0.698	22,833
14	371,159	342,158	276.0	7562	7,430	6.0	3,801	4	345,783	278.9	11,792,367	0.696	58,473
15	347,159	334,579	323.3	6353	6,246	6.0	3,359	5	337,461	326.1	11,527,682	0.697	47,585
16	314,405	303,027	308.4	4505	4,422	4.5	2,366	5	305,078	310.5	10,453,923	0.699	38,640
17	291,718	263,955	277.3	5280	5,056	5.3	2,587	3	266,421	279.9	9,058,647	0.694	41,347
18	298,383	288,853	318.3	5157	5,102	5.6	2,682	7	291,266	321.0	9,969,336	0.699	14,130
19	233,252	220,556	359.7	3416	3,358	5.5	1,883	0	222,031	362.1	7,589,845	0.698	25,332
X	247,019	222,839	134.6	1236	1,192	0.7	575	0	223,456	135.0	7,695,066	0.703	11,025
Genome	8,267,353	7,804,762	299.3	138,608	135,846	5.2	70,378	96	7,870,134	301.8	269,319,475	0.698	969,457

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49

Table S2 Estimated error rates by chromosome. The table provides the total number of imputed genotypes, the percent total imputed genotypes that fall within confidence score bins of (0,0.6), (0.6,0.9), (0.9,1) , the number imputed genotypes that are available in other resources and are used for validation, the overall error rate as the proportion of the agreed genotypes of all comparisons, the percent validation genotypes that fall within confidence core bins of (0,0.6), (0.6,0.9), (0.9,1) and the corresponding error rate for each bin.

chromosome	Rate.missing.genotypes	Total Imputed				Validated Imputed							
		#.total.imputed.genotypes	%.total.imputed.with.conf (0,0.6)	%.total.imputed.with.conf (0.6,0.9)	%.total.imputed.with.conf (0.9,1)	#.genotypes.validated	overall.error	%.validated.with.conf (0,0.6)	error.conf (0,0.6)	%.validated.with.conf (0.6,0.9)	error.conf.(0.6,0.9)	%.validated.with.conf (0.9,0.1)	error.conf.(0.9,1)
1	0.697	23,719,388	15.3	15.3	69.4	92,427	0.106	15.4	0.331	13.2	0.158	71.4	0.047
2	0.703	18,074,289	13.8	13.0	73.3	59,530	0.090	14.4	0.343	12.4	0.163	73.2	0.028
3	0.701	17,516,977	14.7	13.2	72.0	70,360	0.105	14.8	0.335	12.3	0.171	72.9	0.047
4	0.699	16,322,244	15.1	13.9	71.0	66,798	0.099	15.0	0.328	12.7	0.178	72.3	0.038
5	0.696	16,945,501	14.5	14.3	71.2	63,587	0.122	15.4	0.347	13.7	0.194	70.8	0.059
6	0.698	17,426,831	15.7	15.1	69.2	54,962	0.105	15.4	0.334	14.6	0.147	70.1	0.046
7	0.695	13,812,494	14.8	15.5	69.7	61,438	0.107	15.3	0.331	15.4	0.129	69.3	0.052
8	0.699	15,245,102	13.3	13.3	73.4	59,940	0.090	14.5	0.304	13.5	0.135	72.0	0.039
9	0.695	12,306,386	15.1	14.9	70.0	47,836	0.105	14.9	0.348	12.9	0.165	72.2	0.044
10	0.701	13,718,592	13.8	15.0	71.2	38,392	0.097	15.2	0.350	15.2	0.127	69.6	0.035
11	0.701	8,903,563	16.1	14.6	69.2	40,682	0.079	14.1	0.318	10.3	0.149	75.5	0.024
12	0.699	13,564,253	17.1	16.6	66.3	54,140	0.112	18.5	0.318	15.0	0.134	66.5	0.049
13	0.698	13,676,989	13.7	15.3	71.0	22,833	0.097	16.5	0.287	15.5	0.116	68.0	0.047
14	0.696	11,792,367	12.1	14.6	73.3	58,473	0.079	11.9	0.339	12.1	0.105	76.0	0.034
15	0.697	11,527,682	15.1	15.8	69.1	47,585	0.110	15.8	0.339	12.6	0.148	71.6	0.053
16	0.699	10,453,923	15.1	14.6	70.2	38,640	0.113	16.6	0.310	14.1	0.188	69.3	0.050
17	0.694	9,058,647	18.2	16.6	65.2	41,347	0.131	18.4	0.365	15.2	0.173	66.4	0.056
18	0.699	9,969,336	15.8	15.1	69.2	14,130	0.121	19.3	0.319	19.0	0.141	61.7	0.053
19	0.698	7,589,845	15.3	17.5	67.2	25,332	0.107	16.0	0.333	14.0	0.142	70.0	0.048
X	0.703	7,695,066	14.8	21.1	64.0	11,025	0.140	22.2	0.375	26.2	0.151	51.6	0.033
Genome	0.698	269,319,475	14.9	15.0	70.1	969,457	0.104	15.4	0.333	13.7	0.153	70.8	0.044

Table S3. Strain coverage summary of the source data.

Strain	Type	complete sequence											
			NIEHS	Broad	Celera	Wellcome-CTC	GNF	Rosetta/Merck	TJL	db-SNP	Sanger	Jap-MSM	
129S1/SvImJ	classical		X	X	X	X	X	X	X	X	X	X	
129S4/SvJae	classical			X									
129X1/SvJ	classical			X	X	X	X	X	X	X	X	X	
A/J	classical		X	X	X	X	X	X	X	X	X	X	
AKR/J	classical		X	X		X	X	X	X				
BALB/cByJ	classical		X	X		X	X	X	X				
BTBR.T<+>tf/J	classical		X	X		X	X	X	X				
BUB/BnJ	classical			X		X	X	X	X				
C3H/HeJ	classical		X	X		X	X	X	X				
C57BL/6J	classical	X	X	X	X	X	X	X	X	X	X	X	X
C57BLKS/J	classical			X			X	X	X				
C57BR/cdJ	classical			X			X	X	X				
C57L/J	classical			X			X	X	X				
C58/J	classical			X			X	X	X				
CAST/EiJ	<i>M. m. castaneus</i>		X	X		X	X	X	X				
CBA/J	classical			X			X	X	X				
CE/J	classical			X			X	X	X				
CZECHII/EiJ	<i>M. m. musculus</i>			X			X	X	X				
DBA/1J	classical			X			X	X	X				
DBA/2J	classical	X	X	X	X	X	X	X	X	X	X	X	X
DDK/Pas	classical			X									
FVB/NJ	classical		X	X		X	X	X	X				
I/LnJ	classical			X			X	X	X				
JF1/Ms	<i>M. m. molossinus</i>			X			X	X	X				
KK/HLJ	classical		X	X			X	X	X				
LG/J	classical			X			X	X	X				
LP/J	classical			X			X	X	X				
MA/MyJ	classical			X			X	X	X				
MAI/Pas	<i>M. m. molossinus</i>			X			X	X	X				
MOLF/EiJ	<i>M. m. molossinus</i>		X	X			X	X	X				
MSM/Ms	<i>M. m. molossinus</i>			X			X	X	X	X	X	X	X
NOD/LtJ	classical		X	X			X	X	X	X	X	X	X
NON/LtJ	classical			X			X	X	X				
NZB/BINJ	classical			X			X	X	X				
NZO/HILtJ	classical			X			X	X	X				
NZW/LacJ	classical	X	X	X			X	X	X				
O20	classical			X									
PERA/EiJ	<i>M. m. domesticus</i>			X			X	X	X				
PL/J	classical			X			X	X	X				
PWD/PhJ	<i>M. m. musculus</i>		X	X			X	X	X				
PWK/PhJ	<i>M. m. musculus</i>						X	X	X				
Qsi5	classical			X									
RIIIS/J	classical			X			X	X	X				
SEA/GnJ	classical			X			X	X	X				
SEG/Pas	<i>M. spretus</i>			X									
SJL/J	classical			X			X	X	X				
SM/J	classical			X			X	X	X				
SPRET/EiJ	<i>M. spretus</i>			X			X	X	X				
ST/bJ	classical			X			X	X	X				
SWR/J	classical			X			X	X	X				
WSB/EiJ	<i>M. m. domesticus</i>		X	X			X	X	X				

Table S4 Error rates stratified by the status of a SNP being polymorphic (snp.type.1) or constant (snp.type.2) within the classical strains. For each type of SNPs, the table provides (1) the total number and error rates for all validated imputed genotypes (2) the number and error rates for validated imputed genotypes with confidence score greater than 0.9.

Strains	VALIDATED IMPUTED				VALIDATED IMPUTED & CONFIDENCE > 0.9			
	SNP.TYPE.1		SNP.TYPE.2		SNP.TYPE.1		SNP.TYPE.2	
	count	error rate	count	error rate	count	error rate	count	error rate
DBA/2J	56113	0.098	14565	0.380	43469	0.044	7265	0.274
A/J	53584	0.099	12638	0.420	41805	0.046	5680	0.315
129S1/SvImJ	21055	0.110	5791	0.380	16347	0.058	2937	0.263
C3H/HeJ	1067	0.112	270	0.289	776	0.040	168	0.143
BTBR.T<+>tf/J	1008	0.125	248	0.214	754	0.042	172	0.105
FVB/NJ	958	0.125	254	0.252	684	0.053	166	0.139
KK/HLJ	960	0.145	249	0.257	674	0.079	142	0.099
NOD/LtJ	956	0.130	234	0.261	681	0.053	147	0.150
AKR/J	935	0.144	233	0.275	674	0.073	135	0.111
NZW/LacJ	911	0.132	249	0.205	655	0.066	157	0.108
BALB/cByJ	5	0.000	2	0.000	5	0.000	2	0.000
129X1/SvJ	451412	0.067	34040	0.193	342223	0.025	21866	0.106
SM/J	9238	0.114	1457	0.072	6017	0.049	1091	0.040
NZB/BINJ	9010	0.107	1446	0.080	6008	0.038	1067	0.047
NON/LtJ	8996	0.089	1451	0.070	5993	0.026	1088	0.031
SJL/J	8985	0.090	1439	0.084	6006	0.029	1054	0.043
CBA/J	8976	0.079	1447	0.063	6162	0.021	1114	0.026
CE/J	8974	0.129	1447	0.093	5838	0.059	1054	0.053
BUB/BnJ	8974	0.092	1442	0.074	6033	0.030	1078	0.033
LG/J	8967	0.096	1446	0.078	6014	0.038	1084	0.040
SWR/J	8960	0.100	1447	0.079	5989	0.041	1080	0.038
LP/J	8948	0.095	1439	0.079	6024	0.029	1072	0.036
C58/J	8938	0.104	1441	0.058	6018	0.035	1113	0.030
C57BR/cdJ	8929	0.103	1437	0.052	6071	0.035	1124	0.028
I/LnJ	8928	0.099	1435	0.070	6025	0.040	1067	0.037
PL/J	8898	0.089	1434	0.075	6005	0.029	1087	0.031
RIIS/J	8916	0.108	1415	0.082	5815	0.039	1047	0.039
C57L/J	7479	0.113	517	0.072	4885	0.034	381	0.024
MA/MyJ	7450	0.109	523	0.086	4863	0.036	369	0.024
SEA/GnJ	7451	0.097	520	0.121	4911	0.023	361	0.042
C57BLKS/J	6832	0.091	1126	0.044	4779	0.025	889	0.031
DBA/1J	7412	0.103	514	0.113	4927	0.025	365	0.052
CAST/EiJ	1652	0.281	225	0.396	535	0.243	68	0.309
MOLF/EiJ	1538	0.189	226	0.204	760	0.101	127	0.134
WSB/EiJ	902	0.239	219	0.247	560	0.193	107	0.150
PWD/Ph	494	0.209	26	0.269	229	0.135	9	0.111
MSM/Ms	39414	0.093	40539	0.173	23999	0.039	21523	0.067
PERA/EiJ	9053	0.200	1445	0.190	5545	0.130	942	0.131
SPRET/EiJ	7687	0.192	1222	0.390	4440	0.167	644	0.373
CZECHII/EiJ	6699	0.100	1090	0.195	4232	0.064	614	0.111
JF1/Ms	3103	0.069	102	0.118	1911	0.031	55	0.073

Supplementary web site references

<http://www.broad.mit.edu/~claire/MouseHapMap>

<http://mouse.perlegen.com>

<http://www.well.ox.ac.uk/mouse/INBREDS>

<http://www.sanger.ac.uk/modelorgs/mouse.shtml>

<http://www.ensemble.org>

<http://www.ncbi.nlm.nih.gov/SNP>

<http://snp.gnf.org>

<http://phenome.jax.org/>

<http://stt.gsc.riken.jp/msm/>

<http://mousesnp.roche.com/>

<http://cgd.jax.org/>