



HAL
open science

Androgen-binding protein (Abp) evolutionary history: Has positive selection caused fixation of different paralogs in different taxa of the genus *Mus*?

Robert Karn, Golbahar Yazdanifar, Željka Pezer, Pierre Boursot, Christina
Laukaitis

► To cite this version:

Robert Karn, Golbahar Yazdanifar, Željka Pezer, Pierre Boursot, Christina Laukaitis. Androgen-binding protein (Abp) evolutionary history: Has positive selection caused fixation of different paralogs in different taxa of the genus *Mus*?. *Genome Biology and Evolution*, 2021, 10.1093/gbe/evab220 . hal-03370475v1

HAL Id: hal-03370475

<https://hal.science/hal-03370475v1>

Submitted on 8 Oct 2021 (v1), last revised 8 Aug 2022 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

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
50
51
52
53
54
55
56
57

1 *Androgen-binding protein (Abp)* evolutionary history:

2 Has positive selection caused fixation of different paralogs in different taxa of the genus *Mus*?

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
50
51
52
53
54
55
56
57

4 Robert C. Karn¹, Golbahar Yazdanifar², Željka Pezer³, Pierre Boursot⁴ and Christina M. Laukaitis⁵

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
50
51
52
53
54
55
56
57

6 ¹Butler University, Indianapolis, IN, 46208, USA; ²Department of Medicine, College of Medicine,
7 University of Arizona, Tucson, Arizona 85724 USA; ³Ruder Bošković Institute, Bijenička 54, 10000
8 Zagreb, Croatia; ⁴Institut des Sciences de l'Evolution Montpellier, Université de Montpellier, CNRS,
9 IRD, 34095 Montpellier cedex 5, France; ⁵Carle Health and Carle Illinois College of Medicine,
10 University of Illinois, Urbana-Champaign, IL 61801 USA

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

12 ²gyazdanifar@arizona.edu; ³zpezer@irb.hr; ⁴pierre.boursot@umontpellier.fr; ⁵laukaiti@illinois.edu

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

13 Corresponding author: Robert Karn, Department of Biology (emeritus), Butler University, Indianapolis,
14 IN, USA rkarn@butler.edu

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

15 Key words: Androgen-binding protein; gene family expansion; alternative paralogs; copy number variant;
16 positive selection; structural variation

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

17 Running title: Evolutionary history of the *Abp* expansion in *Mus*

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

19 **Significance:** The *Androgen-binding protein (Abp)* gene family is much larger in the house mouse than in
20 the rat and most other vertebrates. This comparative genomics study identified 206 unique *Abp* sequences
21 in six *Mus* taxa to trace their evolution through the genus. We previously suggested that salivary-
22 expressed alleles of *Abpa27* are under positive selection. We found evidence for duplication of *Abpa27* in
23 two taxa and asked whether previous selection tests were done with the assumption that the *Abpa27* genes
24 were orthologous when they were not. Our new evidence suggests that positive selection acted on
25 paralogs rather than alleles. Deletion of different *Abpa27* gene copies in different taxa, therefore, may
26 better explain the apparent positive selection than fixation of different alleles.

1
2
3 27 **Abstract**

4
5 28 Comparison of the *Androgen-binding protein (Abp)* gene regions of six *Mus* genomes provides insights
6
7 29 into the evolutionary history of this large murid rodent gene family. We identified 206 unique *Abp*
8
9 30 sequences and mapped their physical relationships. At least 48 are duplicated and thus present in more
10
11 31 than two identical copies. All six taxa have substantially elevated LINE1 densities in *Abp* regions
12
13 32 compared with flanking regions, similar to levels in mouse and rat genomes, although non-allelic
14
15 33 homologous recombination (NAHR) seems to have only occurred in *Mus musculus domesticus*.
16
17 34 Phylogenetic and structural relationships support the hypothesis that the extensive *Abp* expansion began
18
19 35 in an ancestor of the genus *Mus*. We also found duplicated *Abpa27*'s in two taxa, suggesting that
20
21 36 previously reported selection on *a27* alleles may have actually detected selection on haplotypes wherein
22
23 37 different paralogs were lost in each. Other studies reported that *a27* gene and species trees were
24
25 38 incongruent, likely because of homoplasy. However, L1MC3 phylogenies, supposed to be homoplasy-
26
27 39 free compared to coding regions, support our paralog hypothesis because the L1MC3 phylogeny was
28
29 40 congruent with the *a27* topology. This paralog hypothesis provides an alternative explanation for the
30
31 41 origin of the *a27* gene that is suggested to be fixed in the three different subspecies of *Mus musculus* and
32
33 42 to mediate sexual selection and incipient reinforcement between at least two of them. Finally, we ask why
34
35 43 there are so many *Abp* genes, especially given the high frequency of pseudogenes and suggest that relaxed
36
37 44 selection operates over a large part of the gene clusters.
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Introduction

With the development of dideoxy chain-termination DNA sequencing (Sanger, et al. 1977), the sequences of individual genes and early genomes produced a wealth of single-nucleotide polymorphism (SNP) and small insertion/deletion (indel) data. These have been invaluable for studies of molecular evolution, but since then, it has become apparent that SNPs and small indels explain only a small proportion of trait heritability (Manolio, et al. 2009; Eichler, et al. 2010). Massively parallel sequencing and microarray methods led to the discovery of small microscopic and sub-microscopic variations, structural and quantitative chromosomal rearrangements referred to collectively as structural variation (SV; Alkan, et al. 2011; Huddleston and Eichler 2016). SV includes duplications, deletions, copy-number variants (CNVs), insertions, inversions and translocations ≥ 50 bp in length (Sudmant, et al. 2015). They contribute substantially to the genetic diversity of genomes and therefore are of much interest for studies of cancer genetics, rare diseases and evolutionary genetics (Spielmann, et al. 2018).

SVs have been important in evolutionary studies, e.g. the role of CNVs in reproductive isolation, a requirement for biological speciation (Loire, et al. 2017). They are also of interest in studies of the evolution of gene families in genomes (Higuchi, et al. 2010; Young, et al. 2016; Pezer, et al. 2017; Zhou, et al. 2017; Wang, et al. 2018; Clifton, et al. 2020). Gene families are important sources of phenotypic diversification and genetic innovation and are therefore central to understanding phenotypic change and the molecular origins of heritable variation (Janousek, et al. 2016; Clifton, et al. 2020). Their early genomic origins are, however, still poorly understood because the evolutionary processes that produce tandemly arrayed genes (TAGs) in these families are obscured by the high degree of sequence identity that have arisen from recent and rapid gene duplications (Pan and Zhang 2008; Karn and Laukaitis 2009; Pavlopoulou, et al. 2010; Uriu, et al. 2021). We have studied and described the extensive *Androgen-binding protein* (*Abp*) gene family in the mouse genome (mm10; hereinafter “reference genome”). Our long-term goals are to understand the origin of this large and recently expanded gene family and to trace the evolutionary history of the expansion, including the role of SV, especially CNV, the mechanisms of duplication, and the contributions of retrotransposons (RTs).

1
2
3 72 ABPs are members of the secretoglobin (SCGB) superfamily. These small, soluble cytokine-like
4
5 73 proteins share significant amino acid sequence with uteroglobin (UG; Karn 1994; Laukaitis, et al. 2005)
6
7 74 and share the UG tertiary structure of a four-helix bundle in a boomerang configuration (Callebaut, et al.
8
9 75 2000). The first SCGB superfamily member identified was blastokinin (Krishnan and Daniel 1967),
10
11 76 which was renamed UG when it was found to be secreted in large amounts by the rabbit endometrium
12
13 77 around the time of embryo implantation (Beier 1968; reviewed in Mukherjee and Chilton 2000). The
14
15 78 precise physiological role of UG is not yet known despite the production of two mouse lines with targeted
16
17 79 disruptions of different regions of the UG-encoding gene (*SCGB1A1*; reviewed in Mukherjee, et al.
18
19 80 2007). Most of the emphasis of early work was on biochemistry and physiology, although unfortunately
20
21 81 the SCGB research field has contracted substantially since the 2000 review. SCGBs are thought to be
22
23 82 involved in autocrine, paracrine and endocrine signaling as immunomodulating agents, a role best
24
25 83 understood for the SCGB3A2 protein that reduces airway inflammation and injury in mice (Chiba, et al.
26
27 84 2006; Cai and Kimura 2015; Yoneda, et al. 2016). In the lung, the protein encoded by *SCGB1A1* is called
28
29 85 CC10 and CC16 (Club Cell Secretory Protein). This protein is upregulated after tissue damage, in
30
31 86 premature infants, and in allergy/asthma (Almuntashiri, et al. 2020), suggesting a role in detoxifying
32
33 87 harmful substances inhaled into the lungs.
34
35

36
37 88 Mouse ABPs are produced primarily by glands of the face and neck (Laukaitis, et al. 2005; Karn,
38
39 89 et al. 2014; Karn and Laukaitis 2014) and, as their name suggests, some bind male sex steroid hormones
40
41 90 with a clear specificity (Dlouhy and Karn 1983; Karn 1998; Karn and Clements 1999). Mouse salivary
42
43 91 ABPs mediate assortative mate selection based on subspecies recognition that potentially limits gene
44
45 92 exchange between subspecies where they meet (Laukaitis, et al. 1997; Talley, et al. 2001; Laukaitis and
46
47 93 Karn 2012). There is also evidence that ABP constitutes a system of incipient reinforcement across the
48
49 94 European hybrid zone where house mouse subspecies make secondary contact (Vošlajerová Bímová, et
50
51 95 al. 2011).
52

53
54 96 ABP proteins are secreted dimers composed of an alpha subunit, encoded by an *Abpa* gene,
55
56 97 connected by disulfide-bridging to a betagamma subunit, encoded by an *Abpbg* gene (reviewed in
57

1
2
3 98 Laukaitis and Karn 2012). They are a mammalian novelty with variable evolutionary histories in those
4
5 99 mammals studied to date. The last common ancestor of the eutherian and metatherian lineages 180 MYA
6
7 100 (van Rheede, et al. 2006) had both *Abpa* and *Abpbg* genes, with three each in the metatherian
8
9 101 *Monodelphis* (the opossum; Laukaitis, et al. 2008). By contrast, the entire *Abp* gene region in most
10
11 102 mammals amounts to only a single *Abpa-Abpbg* gene pair as exemplified by the ground squirrel
12
13 103 (*Spermophilus tridecemlineatus*), an outgroup for the purpose of this study. A few species, however,
14
15 104 show varying degrees of expansion of the region (e.g. opossum, cattle, horse, rabbit, house mouse and rat)
16
17 105 and evidence to date suggests that these expansions are independent of one another (Laukaitis, et al.
18
19 106 2008). The house mouse has the largest expansion with 64 total genes spanning 3 Mb (i.e. ~0.1% of the
20
21 107 total genome; Laukaitis, et al. 2008).

22
23
24 108 The *Abpa* and *Abpbg* subunit genes in the mouse genome are most commonly associated in 5'-5'
25
26 109 pairs, (<*a-bg*> or <*bg-a*>, where the arrows point in the 3' directions), with intervening sequence lengths
27
28 110 of 5-20 kb. The mouse reference genome (mm10) has 27 of these gene pairs, called "modules", with ten
29
30 111 singletons (Pezer, et al. 2017). The mouse reference genome *Abp* cluster is ten times the size of that in the
31
32 112 rat genome (rn3) which has only three modules and no singletons (Laukaitis, et al. 2008; Karn and
33
34 113 Laukaitis 2009). These two rodent *Abp* gene families appear to have expanded independently (Emes, et al.
35
36 114 2004) and the 64 *Abp* genes in the mouse genome group into five ancestral clades that may have their
37
38 115 origins as far back in the evolutionary past as the ancestor of the genus *Mus* (Laukaitis, et al. 2008).

39
40
41 116 Our previous investigations examined the expansion of the large *Abp* gene family in the genome
42
43 117 of C57BL/6, a laboratory research strain derived primarily from *Mus musculus domesticus*, to develop a
44
45 118 history of the evolutionary processes involved (Karn and Laukaitis 2009; Janousek, et al. 2013; Janousek,
46
47 119 et al. 2016). The comparative genomics study we report here used high-quality genomes representing
48
49 120 enough other mouse species to determine how extensive and how widespread the expansion was in the
50
51 121 genus *Mus*. Our goals were to resolve old issues concerning *Abp* gene duplication and to reveal emerging
52
53 122 properties of the expansion in house mice. To achieve these goals, we interrogated six recently reported
54
55 123 genomes of members of the genus *Mus* and compared their *Abp* gene families. These include: *Mus*

1
2
3 124 *musculus domesticus* (strain WSB), *M. m. musculus* (strain PWK), *M. m. castaneus* (strain CAS), *M.*
4
5 125 *spretus* (*spr*), *M. caroli* (*car*) and *M. pahari* (*pah*) as well as *Rattus norvegicus* (*Rn*). Hereinafter we will
6
7 126 refer to the subspecies of *M. m. musculus* by their strain names. We identified a total of 206 unique *Abp*
8
9 127 genes from the six *Mus* genomes that we relate to the *Abp* genes in the reference genome and to each
10
11 128 other. We determined their chromosomal positions and explored the roles of retrotransposons (RTs) in
12
13 129 their evolutionary histories. Some of the revelations from this work required us to extensively revise what
14
15 130 we thought we knew about *Abp* evolutionary history, and to discard some of it altogether.

16
17
18 131 We first determined whether the expansion that produced this large family began in an ancestor
19
20 132 of the genus *Mus*. If that is so, we would expect that different *Mus* taxa with diverse distributions and
21
22 133 habitats would still share an evolutionary history of their *Abp* gene region expansions. The counter, or
23
24 134 null, expectation is that their *Abp* evolutionary histories would be dissimilar, even to the point of being
25
26 135 unique in some taxa as we have found in the independent rat duplication. Testing this hypothesis requires
27
28 136 answering three questions: 1) Did duplication occur primarily by $\langle a-bg \rangle$ or $\langle bg-a \rangle$ modules? 2) Does a
29
30 137 phylogeny of *Abp* modules in these six genomes support five ancestral clades (five deeply rooted,
31
32 138 monophyletic gene groups shown in Fig. 3 of Laukaitis, et al. 2008)? and 3) Is the temporal history of the
33
34 139 appearance and expansion of those clades consistent among the six taxa?

35
36
37 140 The results of this genome comparison brought us closer to an understanding of a long-standing
38
39 141 conundrum in the evolutionary history of mouse *Abp* genes, specifically that key genes expressed in
40
41 142 salivary glands and secreted into saliva have phylogenies non-congruent with the species phylogeny.
42
43 143 Karn, et al. (2002) studied the complex history of *Abpa* (later *Abpa27* or *a27*), a gene proposed to
44
45 144 participate in a sexual isolation mechanism in house mice. They observed an abnormal intron phylogeny
46
47 145 for *a27* with an unexpected topology wherein *M. musculus* is not monophyletic and its subspecies stand
48
49 146 as outgroups relative to other Palearctic species (*M. spretus* [*spr*], *M. spicilegus*, and *M. macedonicus*).
50
51 147 Could assessing the copy numbers (CN) of *a27* in the lineage of the genus *Mus* resolve this issue?

52
53
54 148 In this process, we revisited the question of how selection has influenced the expansion history of
55
56 149 the *Abp* gene family. The evolution of gene families is still poorly understood and there is sparse evidence

1
2
3 150 that an increased number of specific genes offers a selective advantage (Hastings, et al. 2009), although
4
5 151 changes (increase or decrease) in the copy number of dosage-sensitive genes can cause clear selective
6
7 152 disadvantage (reviewed in Harel and Lupski 2018). Early evolutionary studies indicated that CNVs might
8
9 153 be advantageous because the genes involved are often those that encode secreted proteins and/or are
10
11 154 enriched for “environmental” functions, including olfaction, immunity, toxin metabolism and
12
13 155 reproduction. Such genes were reported to be under positive selection because they contain higher than
14
15 156 average frequencies of nonsynonymous mutations (Johnson, et al. 2001; Nguyen, et al. 2006; Perry, et al.
16
17 157 2007; Emerson, et al. 2008; Nguyen, et al. 2008; Xue, et al. 2008; Sjodin and Jakobsson 2012). Others,
18
19 158 however, have suggested instead that a non-adaptive explanation could account for their previous
20
21 159 observations (Nguyen, et al. 2006). Finally, is it possible that these six *Abp* clusters are experiencing a
22
23 160 form of genome instability in which large blocks of genes are being gained and lost by non-allelic
24
25 161 homologous recombination (NAHR), possibly representing runaway gene duplication (Janousek, et al.
26
27 162 2016)?
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

Results and Discussion

Comparative genomics of the *Abp* gene families in the genus *Mus*

Interrogation of six wild-derived mouse genomes identified the following number of *Abp* paralogs (defined as unique DNA sequences) in each taxon (**Table 1; Fig. 1**): 11 in *Mus pahari* (*pah*), 33 in *M. caroli* (*car*), 35 in *M. spretus* (*spr*), 43 in *Mus musculus domesticus* (strain WSB), 38 in *M. m. musculus* (strain PWK), and 46 in *M. m. castaneus* (strain CAS) (**supplementary tables S1-S6**). This compares with 64 found in the reference genome (mm10, C57BL/6, B6; Laukaitis, et al. 2008) and we named the 206 *Mus Abp* sequences on those. We did this by: a) comparing their associations in *Abpa* and *Abpbg* phylogenies with the reference genes and ancestral Clades 1-5 (**Fig. 2; supplementary figs. S1, S2**), b) mapping intra-genome linear relationships relative to those of the paralogs (**Fig. 3**); and c) examining the associations of *Abpa* and *Abpbg* genes in putative modules (**supplementary tables S1-S6**). Hereinafter, specific modules will be abbreviated M followed by a number, both in italics, e.g. *M9* and *M10* for the *<bg9-a9>* and *<bg10-a10>* modules, respectively. We evaluated the numbers of *Abp* genes that could be expressed and the numbers of putative pseudogenes (**supplementary table S7**; see also **Table 1 and supplementary tables S1-S6**) and compared them to those in the reference genome (Karn, et al. 2014). The result showed high percentages of pseudogenes in the six *Mus* gene families (WSB, 58%; PWK, 47%; CAS, 50%; *spr*, 53%; *car*, 48%; and *pah*, 36%; comparable to mm10, 53%).

Copy number analysis

Initially, we attempted to estimate *Abp* gene copy number with CNVnator software (Abyzov, et al. 2011). That approach yielded suspiciously low numbers of small CNVs across the *Abp* gene regions of the six *Mus* genomes, very likely because of numerous gaps in the 1504 assemblies. Instead, we calculated CNs based on differences in read depth between *Abp* genes and putative single-copy regions (**supplementary table S8**). Each “unique” gene sequence may be present in a number of copies in a diploid organism. Those that have two copies likely have one on each chromosome (i.e. alleles), while those with copy number >2 have been duplicated and any with CN<2 have been subject to deletion. The numbers of unique *Abp* genes and pseudogenes and the inferred total gene numbers, including

1
2
3 190 duplications, are summarized in **Table 1**. The discrete numerical CN estimates from direct read-depth
4
5 191 calculations on the 206 paralogs we found in this study (**supplementary tables S1-S6**) are consistent with
6
7 192 previous analyses in the three subspecies of *M. musculus* and in *M. spretus* (Pezer, et al. 2017).

8
9 193 Altogether, 85% (40/47) of the copy number variable genes appear in the proximal region (defined as *MI*-
10
11 194 *MI2*; **supplementary tables S1-S6**) and 95% belong to ancestral Clade 2, which has the largest number
12
13 195 of paralogs. The expansion history of these regions hints at a complex sequence of events causing rapid
14
15 196 *Abp* gene expansion in the genus *Mus*.

17 197 *Technical challenges*

18
19
20 198 The genome data we analyzed are of high quality, however, problems remain: 1) the earlier 1504
21
22 199 builds were used to mine *Abp* sequences from the six genomes because they contain genes not found in
23
24 200 the later builds; 2) there seem to be assembly problems, including unexpected gene orders, in the 1504
25
26 201 builds; 3) it is not possible to determine the locations of the duplicated gene copies found in the CN
27
28 202 analysis; and 4) there are many gaps that make it difficult to estimate efficiency of gene-finding.

29
30 203 Our previous studies have shown extensive structural variation in WSB and much less in PWK,
31
32 204 CAS and *spr* (Pezer, et al. 2017). We were able to draw these conclusions because we compared the
33
34 205 sequences of multiple individuals. Here, we were not able to detect homozygous deletions given that the
35
36 206 CN calculations are based on the genome read depth of a single inbred individual from which the
37
38 207 assembly was made. In our study, if there was more than average read depth at a certain locus of the
39
40 208 assembled genome, we called it duplication (i.e. amplification), if there was less than average, we called it
41
42 209 deletion. Thus, it is reasonable to assume that the variation in this region is even higher than we can see
43
44 210 by sequencing and assembling the genome from only one individual. In order to find more paralogous
45
46 211 genes and to detect possible copy number variation in them, one would need to sequence more individuals
47
48 212 of the same taxon or population.

49
50
51 213 While the 1504 builds of these genomes provided the largest number of *Abp* genes, they were
52
53 214 mapped to the reference genome, which may have created or perpetuated assembly problems. The very
54
55 215 high levels of *Abp* sequence identity ($\geq 95\%$) and the use of short reads may have caused additional issues.

1
2
3 216 For example, the proximal and some central genes in *spr*, PWK and CAS do not share the same order as
4
5 217 they do in the reference genome, nor do any two of them share a single, alternative pattern. This is
6
7 218 especially evident in ancestral Clade 1 (i.e. *M1* and *M2* in **Fig. 3**). We suggest several possible
8
9 219 explanations for the “scrambled” appearance of the *Abp* genes in the 1504 builds: a) some of them are
10
11 220 misidentified; b) the genome builds placed them incorrectly; and/or c) small chromosomal rearrangements
12
13 221 occurred locally. The absence of a single, alternative order favors choice b): underlying assembly
14
15 222 problems caused by high sequence identity and high density of repetitive sequences.

16
17
18 223 Assembly problems are expected in genome regions containing segmental duplications (SDs)
19
20 224 because they are repeated sequences with high pairwise similarity. SDs may collapse during the assembly
21
22 225 process causing the region to appear as a single copy in the assembly when it is actually present in two
23
24 226 copies in the real genome (Morgan, et al. 2016). Moreover, individual genes and/or groups of genes may
25
26 227 appear to be out of order compared with the reference and other genomes. In some studies, genotyping of
27
28 228 sites within SDs is difficult because variants between duplicated copies (paralogous variants) are easily
29
30 229 confounded with allelic variants (Morgan, et al. 2016). Latent paralogous variation may bias
31
32 230 interpretations of sequence diversity and haplotype structure (Hurles 2002), and ancestral duplication
33
34 231 followed by differential losses along separate lineages may result in a local phylogeny that is discordant
35
36 232 with the species phylogeny (Goodman, et al. 1979). Concerted evolution may also cause difficulties if, for
37
38 233 example, local phylogenies for adjacent intervals are discordant due to nonallelic gene conversion
39
40 234 between copies (Dover 1982; Nagylaki and Petes 1982).

41
42
43 235 The annotations of these sequences were complicated because existing programs for identifying
44
45 236 orthologs between sequenced taxa (Altenhoff, et al. 2019) were not applicable to our data. The databases
46
47 237 these programs interrogate do not include many of these newly-sequenced taxa of *Mus* and also do not
48
49 238 include the complete sets of gene predictions we make here. Thus, we had to manually predict both gene
50
51 239 sequences and orthology/paralogy relationships. This is a problem facing other groups working with
52
53 240 complex gene families in other non-model organisms (Denecke, et al. 2021). Most importantly, we treated
54
55 241 the problem of orthology in our own, original way. Our conclusion is that orthology is not applicable to at

1
2
3 242 least one of the *Abpa27* paralogs, and possibly to other paralogs (*Abpa26*, *Abpbg26*, *Abpbg25*; **Fig. 5**),
4
5 243 probably due to the apparent frequencies of duplication and deletion and this is precisely the interesting
6
7 244 point of our study.

9 245 Comparison of the gene orders of the six *Mus Abp* regions with the reference genome suggests
10
11 246 perturbed synteny of many *Abp* genes (**Fig. 3**). Overall, the *proximal* region (*M1-12* with some
12
13 247 singletons) shows significant differences among the six taxa while the *distal* region (*M20-27*, singletons
14
15 248 *bg34* and *a30*) has gene orders in the six taxa much more like the same regions in the reference genome.
16
17 249 The *central* region (from singleton *a29* through *M19*, with some singletons) in WSB is unique in that it
18
19 250 includes the penultimate and ultimate duplications, shown above the blue triangle in **Fig. 3** (Janousek, et
20
21 251 al. 2013). The order of proximal and distal genes in *car* agrees relatively well with that in the reference
22
23 252 genome *Abp* region, considering how early *car* diverged from the lineage compared to *spr*, PWK and
24
25 253 CAS. In those three taxa, however, the proximal and some central genes do not share the same order as in
26
27 254 the reference genome, nor do any two of them share a single, alternative pattern, especially *M1* and *M2*.
28
29 255 We note that **Fig. 3** also shows that functional copies of either *Abpa* or *Abpbg* in some species are 1-to-1
30
31 256 orthologs to pseudogenized versions of these genes, which is interesting from the standpoint of patterns of
32
33 257 gene family evolution.

34
35
36
37 258 Gaps in the assembly might account for the perturbed synteny among different taxa. They might
38
39 259 be true deletions in the sequenced genomes when compared with the mouse reference genome used to
40
41 260 guide the genome assembly. Counting from the start position of the first *Abp* gene to the end position of
42
43 261 the last gene in the *Abp* gene region, there are in total 1933, 1115, 637, 488, 75 and 32 assembly gaps in
44
45 262 the 1504 builds of the genome for WSB, CAS, *spr*, PWK, *car* and *pah*, respectively (**supplementary fig.**
46
47 263 **S4**). We looked for reads with perfect matches that span gap positions and found that only 2-6% of gaps
48
49 264 in the *Abp* region could be closed in this way and over 50% of them are not supported by more than two
50
51 265 to four reads. At best, only a small fraction of gaps can be considered as truly missing sequences in the *de*
52
53 266 *novo* assemblies, and the majority of the gaps actually represent unknown sequences. In short, there are
54
55 267 too many gaps in this region to obtain deeper insights into its structure. This illustrates the technical

1
2
3 268 limitations related to assembly of this particular region. With deeper sequencing, additional gene and
4
5 269 intervening sequences may be identified, and the gene orders may change. For example, a significant
6
7 270 number of new *Abp* genes were found in the mouse genome from the first report (Emes, et al. 2004) to the
8
9 271 second (Laukaitis, et al. 2008).

11 272 Finally, it is important to consider a possible role for gene conversion because if it can be ruled
12
13 273 out, then estimates of gene age are reliable. Otherwise, in addition to obscuring orthology, it can also
14
15 274 cause problems for estimating the role of natural selection (Casola and Hahn 2009) and the potential role
16
17 275 of pseudogenes as donors of new mutations for functional genes in the vicinity (Casola, et al. 2012).
18
19 276 Several previous studies suggested that gene conversion has played little if any role in the evolution of the
20
21 277 *Abp* gene region. Laukaitis, et al. (2008) discounted the explanation that nonallelic gene conversion
22
23 278 caused the low divergence of the *Abp* sequences they studied because the phylogenetic tree of ribosomal
24
25 279 protein *L23a* pseudogenes suggests that they frequently co-duplicated with *Abpa-Abpbg* gene modules.

28 280 Laukaitis and Karn (2012) pursued this at a smaller scale by analyzing the 64 *Abp* genes in the
29
30 281 reference genome with GENECONV to look for evidence of short gene conversion tracks. In the case of
31
32 282 the *Abpa* paralogs, GENECONV identified no inner (conversion between genes within alignment) and no
33
34 283 outer (conversion with genes outside alignment) fragments that were globally significant, suggesting that
35
36 284 there is no compelling evidence of gene conversion in *Abpa* paralogs. In the case of the *Abpbg* paralogs,
37
38 285 the analysis identified only one inner (*Abpbg26* and *Abpbg34*) and two outer fragments (*Abpbg5p* and
39
40 286 *Abpbg19*) that were globally significant. They also calculated the GC content of the *Abp* gene region
41
42 287 because sequences undergoing frequent gene conversion, either ectopic or allelic, are expected to become
43
44 288 GC rich (Galtier, et al. 2001; Galtier, et al. 2009). Their results showed that the average GC content in the
45
46 289 *Abp* gene region is low (~41–42%) compared with genes undergoing gene conversion, such as ribosomal
47
48 290 operons and transfer RNAs which have much higher GC contents (Galtier, et al. 2001). They concluded
49
50 291 that gene conversion has made a minimal, but not nonexistent, contribution to the evolutionary history of
51
52 292 the *Abp* gene family. Moreover, it certainly was not significant enough to have confounded the
53
54
55
56
57
58
59
60

1
2
3 293 phylogenetic inference presented by Laukaitis, et al. (2008), and it should not have adversely affected
4
5 294 their analysis of recently duplicated products.
6

7 295 Overall, despite problems in estimating the efficiency of gene finding because of non-uniformity
8
9 296 of read coverage, we propose that we missed relatively few unique gene sequences because we: a) used
10
11 297 three different methods of gene-finding to interrogate the six high-quality genomes; b) found genes with
12
13 298 disparate sequences, even in taxa diverged many MYR apart; and c) identified novel *Abp* genes without
14
15 299 counterparts in the reference genome (e.g. *bgl*, *bgJ*, *bgKp*, *aL*, *aMp* and *aN* in *car*). Similarity in the
16
17 300 numbers of protein-coding and noncoding genes found in the *car* and *pah* genomes compared with the
18
19 301 mouse and rat reference genomes (Thybert, et al. 2018) further supports this conclusion.
20
21

22 302 *Contributions of structural variants*

23
24 303 There appears to be a pattern of fewer genes in the earlier-diverging taxa (**Fig. 1** and **Fig. 3**) with
25
26 304 a large jump from *pah* (11) to *car* (33). To the extent that deletions occurred, they seem clearly to have
27
28 305 been outnumbered by duplications, but that may be an ascertainment bias because “presence” is easier to
29
30 306 identify than “absence”. Insertions and translocations are also not immediately recognizable, however
31
32 307 either could have happened locally and that might explain the “scrambled” gene orders in the Palearctics,
33
34 308 especially on the left flank of **Fig. 3**.
35
36

37 309 The *Abp* regions in the Palearctic taxa that we analyzed in this study have much higher LINE1
38
39 310 (L1) densities than their flanking regions, consistent with those in the mouse and rat genomes (Janousek,
40
41 311 et al. 2013) and the *car* and *pah* genomes (Thybert, et al. 2018), regardless of how populous the *Abp*
42
43 312 genes are in any individual taxon (**supplementary fig. S3**). By contrast, we found intragenic RTs only in
44
45 313 *car_a1bp* (a lineage-specific invasion of LIMd_A into Exon 3) and an insertion of IAP1-MM_I-int, a
46
47 314 member of the ERVK family of LTRs, into the *a30* paralogs of *spr*, WSB, PWK and CAS.
48
49

50 315 Pezer, et al. (2017) used CNVnator to show that the penultimate duplication segment (*M7-M12*
51
52 316 and *bg33-M19*) of the large-block NAHR duplication pattern in the mouse genome is seen only in WSB
53
54 317 and other inbred strains (their Fig. 6), as well as in six wild *Mus* genomes (their Fig. 2). That is consistent
55
56 318 with the derivation of the mouse genome (C57BL/6) from an *M. m. domesticus* mouse. Taken together,
57
58
59
60

1
2
3 319 these genes, along with those of the ultimate duplication (*M14-M17*) shown above the blue triangle at the
4
5 320 top of **Fig. 3** represent these NAHR duplications (Janousek, et al. 2013). By contrast, Pezer, et al. (2017)
6
7 321 found no such correspondence in the CNVnator patterns of the other three taxa (PWK, CAS and *spr*) or in
8
9 322 those of inbred strains derived from them. Here we have actual paralog data to test the idea that genes that
10
11 323 vary in copy number in the other five *Mus* taxa are not in these large, contiguous blocks (**supplementary**
12
13 324 **tables S1-S6**), despite also having enrichment in L1 sequences (**Fig. S3**). Instead, some show unique
14
15 325 duplications (e.g. *car_alap*, *albp* and *alcp*; CAS_ *a12a* and *a12b*), which are not paired in modules
16
17 326 indicating that they might have duplicated after these taxa diverged from the lineage.

20 327 *The Abp family expansion: Modules, ancestral clades and the growth of three regions*

21
22 328 We estimate that 154/206 (~75%) of unique *Abp* genes in these six taxa are pair-members of
23
24 329 modules (**supplementary tables S1-S6**). The 11 *Abp* genes of *pah* are organized into five <*bg-a*>
25
26 330 modules and only one singleton with the *Abpbg* genes on one strand and the *Abpa* genes on the other, as
27
28 331 expected of tandemly arrayed genes (TAGs). *car* has a total of 33 genes with 22 in 11 modules; *spr* has a
29
30 332 total of 35 genes with 24 in 12 modules; WSB has a total of 43 with 38 in 19 modules; PWK a total of 38
31
32 333 with 30 in 15 modules; and CAS has a total of 46 with 30 in 15 modules.

33
34 334 Because some of the new *Abp* genes we found in the six *Mus* taxa are incomplete, we could not
35
36 335 build a phylogeny using only the large intron as done previously for *Abp*'s in the mouse reference genome
37
38 336 (Laukaitis, et al. 2008). This raised the concern that positive selection on coding regions would bias the
39
40 337 phylogenetic trees. To test for this, we needed an unselected sequence common to both genes within a
41
42 338 module. Since most of the *Abp* genes are organized into modules, we searched the intra-modular
43
44 339 sequences between the 5' ends of the *Abpa* and *Abpbg* gene partners for common sequences.

45
46 340 **Supplementary table S9** shows that a LINE1, L1MC3, is nearly ubiquitous within modules of
47
48 341 the reference genome and within those of the six *Mus* taxa (93/103, 90%), suggesting that it is ancestral to
49
50 342 the genus *Mus*. It does not occur in any of the three rat *Abp* modules, however, there is an L1MC3 in the
51
52 343 single intramodular sequence of the ground squirrel, an outgroup to these taxa. That suggests that this RT
53
54 344 was lost in the rat ancestor following divergence from the rodent lineage. **Figure 4** shows a phylogeny

1
2
3 345 constructed from L1MC3 sequences in the modules of the new genes and the reference genome rooted on
4
5 346 the L1MC3 in *pah_M24* because it is the only module common to all six *Mus* taxa (**Fig. 2**;
6
7 347 **supplementary table S1**). The *pah_M24* is also the only *M24* that has an L1MC3, suggesting that this
8
9 348 RT was lost from the lineage following divergence of *pah*. In general, the module phylogeny has a
10
11 349 topology congruent with the gene (*Abpa* and *Abpbg*) phylogenies in **Fig. 2**, suggesting that the **Fig. 2**
12
13 350 phylogenies built on the genes themselves were not biased by the combination of coding regions and
14
15 351 introns that were available to use.

16
17
18 352 The module-based phylogeny we made using L1MC3 was valuable for the insights it provided
19
20 353 into the ancestral clades in the reference genome (those most deeply rooted in the *Mus* phylogeny;
21
22 354 Laukaitis, et al. 2008). **Figure 4** defines the relationships of *pah* modules to several of these ancestral
23
24 355 clades: 1) *pah_M3* and *car_M3* group with the Palearctic *M3*s on the left flank of the large ancestral
25
26 356 Clade 2 in the reference genome, while *pah_MU* is basal to the rest of that ancestral clade; 2) *M24* is the
27
28 357 sole occupant of ancestral Clade 3 and is found in all six of the *Mus* genomes (**supplementary tables S1-**
29
30 358 **S6 and Fig. 2**), however, it appears alone here because only the *M24* in *pah* has L1MC3. Comparison of
31
32 359 the two *Abp* subunit gene phylogenies in **Fig. 2** with the module phylogeny in **Fig. 4** suggests that
33
34 360 Ancestral Clade 1 is more closely related to *M3* than it is to any of the other modules in Clade 2. In fact,
35
36 361 the *bg3* clade in the *Abpbg* phylogeny groups with Clade 1, not with Clade 2 as is the case with the *a3*
37
38 362 clade. As well, the L1MC3 of *M3* has the shortest branch with Clade 1 in **Fig. 4** and *M3* lies physically
39
40 363 next to *M2* as might be expected for tandem duplication products at least when it occurred.

41
42
43 364 **Figure 2** shows that the duplication that gave birth to the ancestor of *M25* and the ancestor of
44
45 365 *M26-27-MX* occurred in an ancestor of the *Mus* lineage, prior to the divergence of *pah*, since it is older
46
47 366 than the divergence between *pah_MX* and *M26-27*. Thus, the duplication that gave rise to *M25* is older
48
49 367 than that which gave rise to *M26-27*. The duplication that gave rise to *M1-M2* (Clade 1) must also have
50
51 368 occurred previously to the divergence of *pah*, confirming the status of Clade 1 as ancestral.

52
53
54 369 In summary, Clades 1-5 are confirmed as ancestral, although clearly Clades 4 and 5 are closely
55
56 370 related. Clade 2 began expanding in the ancestor of *car* and the Palearctic taxa, and some copies survived

1
2
3 371 and were duplicated, while other paralogs died after the divergence of the Palearctics (**Fig. 2;**
4
5 372 **supplementary table S2**; see also **Fig. 3**). This clade is larger and more complex in the three subspecies
6
7 373 of *M. musculus* and seems to have been the source of most of the volatility identified when comparing the
8
9 374 *Abp* gene regions of 15 inbred strains to the mouse genome using the Mouse Paralogy Browser (Karn and
10
11 375 Laukaitis 2009).

12
13
14 376 Modules *M24*, *MX* and *MY* in *pah* (**supplementary table S2**) may represent the ancestors of the
15
16 377 entire right-flank in *car* (the segment in the mouse genome stretching from *M24* to *a30*). We did not find
17
18 378 a “classical” ancestral Clade 1 (*M1-M2*) in *pah*, because *aU*, *bgUp* and *aVp* are not in the reverse order
19
20 379 (i.e. switched strands) in relation to the other *pah* genes/modules, as Clade 1 is in the other four taxa (**Fig.**
21
22 380 **3**). One possibility, however, is that they do represent *pah* Clade 1 but the strands on the other four taxa
23
24 381 represent the outcome of an event that occurred between the divergence of *pah* and the other four, perhaps
25
26 382 during the massive genome rearrangement that followed divergence of *M. pahari* from the ancestral
27
28 383 lineage and before divergence of *M. caroli* 3-6 MYA (Thybert, et al. 2018).

29
30
31 384 The central gene region (ancestral Clade 2), is smaller and less complex in *pah*, probably only
32
33 385 represented by *M3*. However, in *car*, it is comprised of nearly 20 genes: *M3*, three *a28*-like paralogs, eight
34
35 386 genes variously related to *M21-23* and six more deeply rooted paralogs (*aL*, *aMp*, *aNp*, *bgI*, *bgJ*, and
36
37 387 *bgKp*), which likely explains the jump from 11 genes in *pah* to 33 in *car* (see above). The gene numbers
38
39 388 making up the populous and volatile central region in the *M. musculus* subspecies are consistently larger
40
41 389 than in the other three taxa. Ancestral Clade 4 (*M25*) is seen only in the Palearctic taxa, however, it had to
42
43 390 have a progenitor in the ancestor of *Mus* because it is basal to *M26* and *M27* (**Fig. 2** and **4**). So, *M25* was
44
45 391 either deleted or we failed to find it in both *pah* and *CAS*.

46
47 392 Taken together, our observations on the *Abp* gene family expansion, the modules, the Clades and
48
49 393 the growth of the three regions, provide strong support for the idea that expansion of the large reference
50
51 394 genome *Abp* family began in an ancestor of the genus *Mus*. They also suggest that most or all of the *Abp*
52
53 395 genes in these six *Mus* genomes are related as branches within one or another of the five ancestral Clades.
54
55 396 The alternative would have been independent expansions, similar to the rat *Abp* region where individual
56
57

1
2
3 397 paralogs are not orthologous with those in the genus *Mus*. Another way of thinking about this is that most
4
5 398 of the *Abps* in *Mus* have orthologs in some or all of the six taxa we studied. That suggests that they
6
7 399 evolved from a shared lineage while none of them have orthologs in the rat, which apparently had an
8
9 400 independent expansion.

11 401 *The role of selection in Mus Abp gene evolution: Reconciling topologies of the gene and species trees*

13 402 Studies of selection on *Abp* genes have focused on *a27*, *bg27*, and *bg26*, the three saliva-
14
15 403 expressed paralogs because of evidence that *Abp* has a role in sexual selection between house mouse
16
17 404 subspecies (Laukaitis, et al. 1997; Talley, et al. 2001; Bímová, et al. 2005). Hwang, et al. (1997) observed
18
19 405 a high nonsynonymous/synonymous substitution ratio (dN/dS) in their *Abpa* (now *a27*) sequence data
20
21 406 from six *Mus* taxa and proposed that directional selection was a sufficient explanation of their data. They
22
23 407 envisioned the possibility of cyclical selection of certain amino acid variants that became advantageous at
24
25 408 some stage and they posited that homoplasy occurred in the phylogeny of the *Abpa* haplotypes that was
26
27 409 incongruent with the canonical phylogeny of the genus. Karn and Nachman (1999) used the HKA test
28
29 410 (Hudson, et al. 1987) to investigate patterns of DNA sequence variation at *a27* within and between
30
31 411 species of mice. Their results provided evidence that selection has shaped the evolution of *Abpa* in house
32
33 412 mice and was consistent with a recent adaptive fixation (a selective sweep) at or near *Abpa*. They also
34
35 413 calculated the ratio of nonsynonymous substitutions to synonymous substitutions on a per-site basis
36
37 414 (Ka/Ks) for the *Mus* sequences of Hwang, et al. (1997). Based on the combined observations of no
38
39 415 variation at *a27* within *M. m. domesticus* and uniformly high Ka/Ks values between species, they
40
41 416 suggested that positive directional selection has acted recently at this locus. Laukaitis, et al. (2012)
42
43 417 assessed site-specific positive selection on the coding sequences of three genes, *a27*, *bg26* and *bg27*, in
44
45 418 five *Mus* taxa using the program CODEML in the PAML package (Yang 2007). They concluded that at
46
47 419 least two (*a27*, *bg26*) of the three genes encoding the subunits of ABP dimers evolved under positive
48
49 420 selection and suggested that the third one may have also.

51
52
53 421 These selection tests were based on the assumption that the *a27* genes in the subspecies of *M.*
54
55 422 *musculus* are orthologs and thus that the studied variants were alleles. However, some genes have a

1
2
3 423 phylogeny at variance with the species phylogeny and Karn, et al. (2002) suggested that the *M. musculus*
4
5 424 taxa are not monophyletic and its subspecies are outgroups relative to other Palearctic species. Here, we
6
7 425 provide evidence that *pah* and *car* both appear to have duplications of modules related to *M27*,
8
9 426 specifically *MX* and *MY* in *pah*; as well as *M27a* (<*bg27a-a27a*>) and *M26/27b* (<*bg26-a27bp*>) in *car*
10
11 427 (**Figs. 2, 3, and 5**). These extra *M27* modules are not found in the Palearctic taxa that have their *a27*
12
13 428 topologies incongruent with that of the species phylogeny (Karn, et al. 2002). Such duplications and
14
15 429 deletions may also have occurred in the ancestor of the Palearctics, so that the copies we observe now are
16
17 430 not necessarily all orthologous. That could provide a parsimonious explanation for why the gene
18
19 431 phylogeny is incongruent with the species phylogeny. Interestingly, **Fig. 2** shows that clades *a26*, *bg25*
20
21 432 and *bg26* are also noncongruent with the species phylogeny.
22
23

24 433 Karn, et al. (2002) discussed and discarded an explanation for the incongruent gene and species
25
26 434 trees that was based on a hypothetical duplication that produced two copies of *a27* in an early ancestor(s).
27
28 435 In this view, different *a27* paralogs were fixed or lost producing very different “*a27*” sequences in *M. m.*
29
30 436 *domesticus* and *M. m. musculus* that were not orthologous. The critical point is that, if duplication of *M27*
31
32 437 and related modules led to fixation of different paralogs in *M. m. musculus* and the other Palearctic
33
34 438 species/subspecies, then the many selection tests reported for *a27*, *bg26* and other ancestral Clade 5 genes
35
36 439 (Karn and Nachman 1999; Laukaitis, et al. 2003; Laukaitis, et al. 2012) were done with the assumption
37
38 440 that they were orthologous in all the Palearctic taxa when they were not.
39
40

41 441 In this study we discovered that *pah* and *car* both have two modules that appear to be duplicated
42
43 442 ancestral versions of *M27* (**Fig. 5**). Later in *Mus* evolutionary history, random mutation could have
44
45 443 created a situation with two haplotypes segregating in a population, one haplotype having paralog *a27a*
46
47 444 with paralog *a27b* deleted and another haplotype that retained paralog *a27b* with paralog *a27a* deleted.
48
49 445 These would fit the description of “pseudoalleles” if tandem duplication had produced the two paralogs.
50
51 446 Assuming that the population gave rise to two separate migrations (as in the case of the progenitors of *M.*
52
53 447 *m. domesticus* and *M. m. musculus*), selection and/or drift could have increased the frequency of paralog
54
55 448 *a27a* in one population and conversely paralog *a27b* in the other, perhaps even to fixation in both. If
56
57
58
59
60

1
2
3 449 individual animals in the two subpopulations could sense the different salivary proteins expressed by the
4
5 450 pseudoalleles, it might have led to olfactory recognition resulting in homosubspecific selection and
6
7 451 eventually incipient reinforcement. The ancestors of the *M. m. domesticus* and *M. m. musculus* subspecies
8
9 452 made secondary contact 5000-10,000 years ago, forming what is now the European mouse hybrid zone
10
11 453 (Boursot, et al. 1993; Sage, et al. 1993). It seems clear from the literature that “*a27*”, whether orthologous
12
13 454 or paralogous in these two subspecies, mediates sexual selection and constitutes a system of incipient
14
15 455 reinforcement at the mouse hybrid zone (Vošlajerová Bímová, et al. 2011).

16
17
18 456 This emergent property highlights perhaps the most important contribution of the module trees
19
20 457 because previous explanations of the topology of these genes tended to cite homoplasy as the result of
21
22 458 strong positive selection (Hwang, et al. 1997; Karn, et al. 2010; Laukaitis, et al. 2012). One of the reasons
23
24 459 we used LIMC3 to build *Abp* module phylogenies is that L1 RTs are thought to be homoplasy-free
25
26 460 regions compared to gene regions (Verneau, et al. 1998; Semple and Steel 2002; Alexeev and Alekseyev
27
28 461 2018). Since the abnormal topology in the *a27* phylogeny is not eliminated by building a phylogeny with
29
30 462 the intramodular LIMC3 (**Fig. 4**), we conclude that it is the result of SV rather than homoplasy.
31
32
33 463 Moreover, it also shows that other genes in *M25* and *M26*, which are related by descent to *a27* as the
34
35 464 result of duplications, also have abnormal topologies. Coupled with the observation that *a27* and *bg27*
36
37 465 have duplicates in *pah* and *car*, this better supports the notion that duplication produced two copies of
38
39 466 *M27*, which then were differentially eliminated, causing one to be fixed in an ancestor of PWK and the
40
41 467 other in an ancestor of the rest of the Palearctic taxa. We feel that this explanation, rather than
42
43 468 explanations such as the occurrence of secondary genetic exchanges along the lineages leading to the
44
45 469 Palearctic taxa (Karn, et al. 2002), is more parsimonious and better fits the data we report here.

46
47 470 *Why are there so many Abp genes in the genus Mus?*

48
49
50 471 Only genes found in ancestral Clade 5 of the reference genome are expressed in salivary glands
51
52 472 and secreted into saliva, while many more genes from three of the other clades are expressed in lacrimal
53
54 473 glands and secreted into tears (Karn, et al. 2014). This led to the proposal that, early in the expansion of
55
56 474 the mouse *Abp* gene family, neo- or subfunctionalization occurred to create this clear-cut partitioning of

1
2
3 475 *Abp* expression between these glands of the face and neck. Their proposal raises the possibility that the
4
5 476 function(s) of ABP proteins in tears differs from those of ABP proteins in saliva. We have already
6
7 477 reviewed the functions of the paralogs expressed in ancestral Clade 5 above (mate recognition, sexual
8
9 478 selection, incipient reinforcement) and no other functions have yet been found for the *Abp* paralogs in
10
11 479 ancestral Clades 1-4.

12
13
14 480 So why are there so many other *Abp* genes in the members of the genus *Mus*? We have searched
15
16 481 repeatedly for other potential functions for *Abp* genes but as yet have found none (Karn and Dlouhy 1991;
17
18 482 Chung, et al. 2017). An equally important question is: Why does there appear to be so much genome
19
20 483 instability and polymorphism of CNV in naturally occurring WSB populations (Karn and Laukaitis 2009;
21
22 484 Pezer, et al. 2017), suggesting runaway gene duplication (Janousek, et al. 2016), and why is this not found
23
24 485 in the other Palearctic taxa (Pezer, et al. 2017) and this report?

25
26 486 Nguyen, et al. (2008) questioned their previous proposal (Nguyen, et al. 2006) that CNVs are
27
28 487 often retained in the human population because of their adaptive benefit. Rather, they showed that genic
29
30 488 biases of CNVs are best explained, not by positive selection, but by reduced efficiency of selection in
31
32 489 eliminating deleterious changes from the human population. We propose here that this might also apply to
33
34 490 *Abp* genes in mouse populations. They are environmental genes (*sensu* Nguyen, et al. 2006; Nguyen, et
35
36 491 al. 2008) associated with SDs and so are subject to frequent duplication, deletion and pseudogene
37
38 492 formation (Lander, et al. 2001; Waterston, et al. 2002; Gibbs, et al. 2004; Perry, et al. 2008; Karn and
39
40 493 Laukaitis 2009; Sjodin and Jakobsson 2012). This may be why so much volatility has been observed,
41
42 494 especially in the central region of the *Abp* cluster in the reference genome (Karn and Laukaitis 2009) and
43
44 495 in the six *Mus* genomes we studied here. It could also explain why >50% of *Abp* paralogs in the reference
45
46 496 genome are pseudogenes (Karn, et al. 2014) and why we found similarly high percentages of pseudogenes
47
48 497 in the six *Mus* taxa. These observations suggest that the environment of the *Abp* gene region may be more
49
50 498 permissive for duplication in the sense of retention of pseudogenized duplicates. On the one hand, this
51
52 499 might facilitate duplication followed by neo- and/or subfunctionalization (e.g. Karn, et al. 2014) while, on
53
54
55
56
57
58
59
60

1
2
3 500 the other hand, volatility may have resulted from an elevated rate of gene birth and death (Karn and
4
5 501 Laukaitis 2009; Janousek, et al. 2013; Pezer, et al. 2017).

6
7 502 *Broader evolutionary implications*

8
9 503 Given that most of the emphasis of past secretoglobin work was on biochemistry and physiology,
10
11 504 there is little information on the evolution of this superfamily, a regrettable situation that may account for
12
13 505 the lack of a holistic picture of secretoglobin function(s). We have pursued questions about the evolution
14
15 506 of the extensive *Abp* families in house mice but it seems that we and others in the field have been remiss
16
17 507 in not considering the cytokine nature of these molecules in an evolutionary context. For example, gene
18
19 508 families, such as those involved in chemosensation, reproduction, host defense and immunity, and toxin
20
21 509 metabolism that are expanded, usually as tandem duplications, in one lineage are often expanded in
22
23 510 another (Janousek, et al. 2013). Cytokines are part of the body's immune response to infection, and
24
25 511 considering that in conjunction with their detoxifying capability and with ABP's reproductive function
26
27 512 (mate recognition), we propose that there may be a connection among these seemingly diverse
28
29 513 capabilities. Karn, et al. (2014) observed that a partitioning of *Abp* expressions occurred early in the
30
31 514 evolution of the ancestral clades in the genome of the ancestor of the genus *Mus*. They suggested that
32
33 515 neo- and/or subfunctionalization was responsible for partitioning the expression of eleven *Abp* paralogs
34
35 516 from ancestral Clades 1, 2 and 3 into lacrimal gland/tears, and three *Abp* paralogs from ancestral Clade 5
36
37 517 into submandibular gland/saliva. Again considering possible multiple functions, this time in the context of
38
39 518 the large *Abp* gene family in the mouse, could it be that a small group of these secretoglobins initially
40
41 519 served in an immune and/or detoxification capacity but that one or more mutated to take on a
42
43 520 communication function in the reproductive sense of sexual selection either exclusively or in addition to
44
45 521 the earlier function? This might explain the extensive expansion of the *Abp* families not only in rodents
46
47 522 but also in a few other more distant mammals.
48
49
50
51
52
53
54
55
56
57
58
59
60

Conclusions

We identified 206 unique *Abp* gene sequences in the genomes of six taxa of the genus *Mus* and mapped their relative positions in those *Abp* clusters. Our CN estimates suggest that the total number of paralogs is closer to 300. We present evidence that the roots of the mouse reference genome expansion in the ancestor of the *Mus Abp* lineage had substantially elevated L1 densities over their *Abp* regions. Further, we suggest that previous analyses of selection on *a27* alleles fixed in different species and subspecies of *Mus* actually detected selection on different paralogs that resulted from tandem duplication (i.e., pseudoalleles) of *a27* in an ancestor of *Mus*. This alternative explanation is based on our finding that *pah* and *car* both have two modules that appear to be duplicated ancestral versions of *M27*. These alternative *a27* paralogs could have arisen as haplotypes having lost different copies through random mutation. If these became recognizable by olfaction, it could have led to sexual selection and eventually incipient reinforcement. We offer evidence for this paralog hypothesis from our module phylogenies built with L1MC3 because L1 RTs are thought to be homoplasmy-free regions compared to gene regions. Fixation of different *a27* paralogs in the subspecies by selection is consistent with; 1) different *a27* sequences fixed in each of the three subspecies of *M. musculus*; 2) finding that ABP mediated sexual selection and incipient reinforcement in previous behavioral analyses; and 3) the incongruence of the gene and species phylogenies which was explained previously by homoplasmy.

Finally, we propose that the roles of cytokines in immune response to infection and detoxification be considered in further work on the evolutionary histories of the secretoglobins generally and the ABPs specifically. A unique aspect of ABPs in the secretoglobin superfamily evolution is that some diverse mammal species have *Abp* families that expanded extensively. Genes involved in adaptation and functional innovation are subject to frequent duplication, deletion and pseudogene formation and prevalent among rapidly evolving genes are those involved in immunity, reproduction, chemosensation and toxin metabolism. Considering possible multiple functions in the context of the large *Abp* gene family in the mouse, we ask if a small group of them initially served in an immune and/or detoxification capacity but that one or more mutated to take on a communication function in the reproductive sense of sexual

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
50
51
52
53
54
55
56
57
58
59
60

550 selection. That might explain the extensive expansion of *Abp* families not only in rodents but in a few
551 other more distant mammals.

Downloaded from <https://academic.oup.com/gbe/advance-article/doi/10.1093/gbe/evab220/6377336> by guest on 07 October 2021

553 **Materials and Methods**

554 *Genomic sequences*

555 Paired-end Illumina HiSeq2000 sequencing data were obtained from the European Nucleotide
556 Archive (ENA) as fastq files for WSB (*Mus musculus domesticus*, ENA Accession number ERS076380);
557 PWK (*M. m. musculus*, ERS076378); CAST (*M. m. castaneus*, ERS076381); and *M. spretus*
558 (ERS076388, ERS138732). The 1504 builds of the genomes of WSB, PWK, CAST and *M. spretus* were
559 obtained from Sanger Mouse Genomes project (<ftp://ftp-mouse.sanger.ac.uk/> (Keane, et al. 2011; Lilue, et
560 al. 2018)) as were *M. caroli* and *M. pahari* (Thybert, et al. 2018). Although the later *de novo* assemblies
561 uncoupled from the reference genome had better overall statistics than the 1504 builds, more unique *Abp*
562 sequences were identified on chromosome 7 (chromosome 1 in *pah*) in the 1504 builds (Thybert, et al.
563 2018). Because they yielded the largest number of *Abp* genes from each genome and the most
564 parsimonious set of gene assignments to chromosomes, we used the gene predictions and coordinates of
565 the 1504 builds.

566 *Data mining genomes for Abp sequences*

567 We employed the BLASTn function in DNA Workbench v3.1.0 ([https://www.ncbi.nlm.nih.gov/](https://www.ncbi.nlm.nih.gov/tools/gbench/)
568 [tools/gbench/](https://www.ncbi.nlm.nih.gov/tools/gbench/)), hmmersearch in HMMER/ 3.1 (<https://github.com/EddyRivasLab/hmmer>), and
569 Exonerate/2.2.0 (<https://www.ebi.ac.uk/about/vertebrate-genomics/software/exonerate>) to search for
570 sequences similar to *Abp* genes identified previously in rodents. We then searched again with the newly
571 identified sequences but did not obtain additional genes. The sequences we identified by these methods
572 were searched manually for start and stop codons and for donor and acceptor intron splice sites. We also
573 verified the flanking genes *Scn1b* and *Uba2* (formerly *Uble1b*). Once a mouse taxon *Abp* gene sequence
574 was identified bioinformatically, it was verified by designing a set(s) of primers, amplifying it in genomic
575 DNA, and sequencing it (UAGC core facility, University of Arizona) as previously reported (Laukaitis, et
576 al. 2005). The DNA samples used as PCR templates for these six genus *Mus* taxa were obtained from
577 Jackson Laboratory (Bar Harbor, ME). We aligned the sequences we obtained in the laboratory with the
578 corresponding data-mined sequences using DNAsis max/3.0 (MiraiBio, San Francisco, CA). This allowed

1
2
3 579 us to verify and/or correct the data-mined sequences and complete some that were not full-length.

4
5 580 **supplementary tables S1-S6** contain the gene descriptions by taxon with their GenBank number;

6
7 581 **supplementary table S7** shows their gene/pseudogene status.

8
9 582 *Sequence coverage and calculation of gene copy numbers (CN)*

10
11 583 We first attempted to estimate copy number (CN) of *Abp* genes using CNVnator software
12
13 584 (Abyzov, et al. 2011), but due to numerous gaps in the *Abp* regions of the 1504 build assemblies (see
14
15 585 below), this yielded suspiciously low numbers of small CNVs across the *Abp* gene regions of the six
16
17 586 mouse genomes. Hence, we calculated the CNs based on differences in read depth between *Abp* genes and
18
19 587 supposedly single-copy regions. With samtools (Li, et al. 2009) we extracted the number of reads mapped
20
21 588 to each *Abp* gene and calculated the coverage as (read count/gene length) x (average read length). In the
22
23 589 same way we also calculated the coverage for the 1000 randomly chosen regions of 2 kb in length of each
24
25 590 *Abp*-containing chromosome where CNVnator did not reveal CNVs. Regions with less than 10% of reads
26
27 591 of low mapping quality (defined as MAPQ < 20) were chosen for calculating the average coverage for
28
29 592 single-copy sequences (**supplementary table S8**). We derived a diploid copy number for each *Abp* gene
30
31 593 by dividing the coverage of the *Abp* gene with the average coverage for single-copy sequences. Very few
32
33 594 reads in any of the *Abp* genes were of low quality (MAPQ<20). In the case of two sequences of *car_a28a*
34
35 595 and *b* (**supplementary table S2**), multiple mapping locations could have inflated their apparent coverage.
36
37 596 This made only a minor contribution (<1% of the 206 genes we found) to CN determination. We have
38
39 597 shown that the average GC content in the *Abp* gene region is in line with the genome average (Karn and
40
41 598 Laukaitis 2009). Hence, we assumed that GC bias is not an issue for read-depth based inference of copy
42
43 599 number in the *Abp* region.

44
45 600 *Analysis of gaps in the assemblies*

46
47 601 To examine the possibility that some of the gaps represent truly missing sequences rather than
48
49 602 undefined difficult to assemble regions, we looked for read evidence to dismiss gap placement in the
50
51 603 assembly. For that reason, we deleted all gaps (N's) in the reference chromosome where the *Abp* cluster
52
53 604 resides in order to create a new reference chromosome without gaps. We next performed mapping and
54
55
56
57
58
59
60

1
2
3 605 subsequent processing steps with the gap-free reference as described above. We searched for reads that
4
5 606 span gap positions and are either ≥ 80 nucleotides (nt) long or have mapping quality > 0 . The reads were
6
7 607 also required to exactly match the gap-free reference over the entire read length and to have at least 8 nt
8
9 608 of sequence around the predicted gap position (i.e. that the gap position is not at the very end of the read).

11 609 *Retrotransposon content*

12
13
14 610 The data table “rmsk” was obtained from the UCSC FTP server for C57Bl/6 and the Sanger
15
16 611 Mouse Genomes project for six mouse taxa (see above). Data for LINEs was extracted (Wicker, et al.
17
18 612 2007; Kapitonov and Jurka 2008). Sliding windows of 100 kb with 10 kb steps were created across each
19
20 613 genome assembly using the “bedtools makewindows” command of bedtools. The number of bases within
21
22 614 each window that is covered by LINEs was calculated using bedtools coverage ([https://github.com/](https://github.com/arq5x/bedtools)
23
24 615 [arq5x/bedtools](https://github.com/arq5x/bedtools)). When gaps were present within the assembly, two coordinate systems were created: one
25
26 616 before the gap removal and one after the gap removal, and positions of LINEs and *Abp* genes were
27
28 617 converted between these ([https:// github.com/ucscGenomeBrowser/kent](https://github.com/ucscGenomeBrowser/kent)). The density plots with gaps
29
30 618 removed were made using the ggplot2 package in R.

31
32
33 619 The rmsk RT data table was obtained from the Sanger Mouse Genomes project for six mouse taxa
34
35
36 620 (<ftp://ftp-mouse.sanger.ac.uk/> Keane, et al. 2011; Lilue, et al. 2018). Positions of L1MC3 Elements were
37
38
39 621 extracted from the rmsk table and were filtered using the *Abp* intra-module coordinates. L1MC3
40
41
42 622 sequences were then obtained using “bedtools getfasta” command of bedtools
43
44
45 623 (<https://github.com/arq5x/bedtools>). Those intra-module sequences negative for L1MC3 when searched in
46
47
48 624 that manner were searched again by aligning L1MC3 sequences from other modules and, in some cases
49
50
51 625 this revealed the retrotransposon in the intra-module sequences.

54 626 *Data analysis*

1
2
3 627 We assigned exons and introns to the verified and/or corrected DNA sequences of the six taxa of
4
5 628 *Mus musculus* by aligning them with the known exon and intron sequences of four *Abpa* and four *Abpbg*
6
7 629 genes from the mouse genomes (*a2, a27, a24, a27, bg2, bg7, bg24 and bg27*). The donor and acceptor
8
9 630 splice sites were identified and the exons were assembled into putative mRNAs and translated *in silico*.
10
11 631 From the translations, we identified each gene as either a potentially expressed gene or as a pseudogene if
12
13 632 it had either a disruption in the coding region and/or a noncanonical splice site (Emes, et al. 2004). **Tables**
14
15 633 **S1-S6** show the disruptions for the putative pseudogenes. MAFFT was used to align the *Abp* gene
16
17 634 sequences from the genus *Mus* and the mouse and rat reference genomes, IQtree (<http://www.iqtree.org>;
18
19 635 Trifinopoulos, et al. 2016) was used to build maximum likelihood phylogenetic trees that were visualized
20
21 636 with FigTree v1.4.3 (<http://tree.bio.ed.ac.uk/software/figtree>). Initially, we built trees with the larger
22
23 637 intron b, that lies between Exons 2 and 3, in order to avoid bias caused by selection (Laukaitis, et al.
24
25 638 2008). Comparisons with trees constructed with the full genes (ATG to the stop codon) showed
26
27 639 essentially the same topologies and allowed us to include partial sequences lacking most or all of intron b.
28
29 640 Bootstrap values (1000 repetitions) were obtained with the MAFFT ultrafast bootstrap approximation.
30
31 641 L1MC3 RTs from the intramodular regions were aligned and used for producing MAFFT and IQTree.
32
33
34
35
36

642 **Acknowledgements**

37 643
38
39 644 This work was supported by the National Cancer Institute at the National Institutes of Health
40
41 645 [grant number P30 CA23074] for laboratory infrastructure. The authors gratefully acknowledge Willie
42
43 646 Swanson for helpful discussions and Miloš Macholán for additional recommendations.
44
45
46
47
48

49 648 **Data Availability Statement**

50 649 All sequence data are released into GenBank and their Accession numbers are listed in
51
52 650 **supplementary tables S1-6, supplementary material** online.
53
54
55
56
57
58
59
60

1
2
3 **652 Author Contributions**
4

5 653 RCK, GY and CML conceived of the project, mined the *Abp* and *LIMC3* sequence data, designed primers
6
7 654 and sequenced the genes and built phylogenies. ZP did the *Abp* module alignments, the copy number
8
9 655 analyses and the gap analyses. PB and RCK assessed the evolutionary forces acting on *Abp* orthologs vs.
10
11 656 paralogs. All the authors participated in writing the manuscript.
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

1
2
3 658 **References**
4
5

- 6 659 Abyzov A, Urban AE, Snyder M, Gerstein M. 2011. CNVnator: an approach to discover, genotype, and
7 660 characterize typical and atypical CNVs from family and population genome sequencing. *Genome Res*
8 661 21:974-984.
- 10 662 Alexeev N, Alekseyev MA. 2018. Combinatorial Scoring of Phylogenetic Trees and Networks Based on
11 663 Homoplasy-Free Characters. *J Comput Biol* 25:1203-1219.
- 13 664 Alkan C, Coe BP, Eichler EE. 2011. Genome structural variation discovery and genotyping. *Nat Rev*
14 665 *Genet* 12:363-376.
- 16 666 Almuntashiri S, Zhu Y, Han Y, Wang X, Somanath PR, Zhang D. 2020. Club Cell Secreted Protein
17 667 CC16: Potential Applications in Prognosis and Therapy for Pulmonary Diseases. *J Clin Med* 9.
- 18 668 Altenhoff AM, Glower NM, Dessimoz C. 2019. Inferring Orthology and Paralogy. In: Anisimova M,
19 669 editor. *Evolutionary Genomics: Statistical and Computational Methods*. New York, NY: Humana Press.
20 670 p. 149-176.
- 22 671 Beier HM. 1968. Uteroglobin: a hormone-sensitive endometrial protein involved in blastocyst
23 672 development. *Biochim Biophys Acta* 160:289-291.
- 25 673 Bímová B, Karn RC, Pialek J. 2005. The role of salivary androgen-binding protein in reproductive
26 674 isolation between two subspecies of house mouse: *Mus musculus musculus* and *Mus musculus*
27 675 *domesticus*. *Biological Journal of the Linnean Society* 84:349-361.
- 29 676 Boursot P, Auffray J-C, Britton-Davidian J, Bonhomme F. 1993. The evolution of house mice. *Annu Rev*
30 677 *Ecol Syst* 24:119-152.
- 31 678 Cai Y, Kimura S. 2015. Secretoglobin 3A2 Exhibits Anti-Fibrotic Activity in Bleomycin-Induced
32 679 Pulmonary Fibrosis Model Mice. *PLoS One* 10:e0142497.
- 34 680 Callebaut I, Poupon A, Bally R, Demaret JP, Housset D, Delettre J, Hossenlopp P, Mornon JP. 2000. The
35 681 uteroglobin fold. *Ann N Y Acad Sci* 923:90-112.
- 37 682 Casola C, Conant GC, Hahn MW. 2012. Very low rate of gene conversion in the yeast genome. *Mol Biol*
38 683 *Evol* 29:3817-3826.
- 40 684 Casola C, Hahn MW. 2009. Gene conversion among paralogs results in moderate false detection of
41 685 positive selection using likelihood methods. *J Mol Evol* 68:679-687.
- 43 686 Chevret P, Veyrunes F, Britton-Davidian J. 2005. Molecular phylogeny of the genus *Mus*
44 687 (Rodentia:Murinae) based on mitochondrial and nuclear data. *Biol J Linn Soc* 84:417-427.
- 45 688 Chiba Y, Kurotani R, Kusakabe T, Miura T, Link BW, Misawa M, Kimura S. 2006. Uteroglobin-related
46 689 protein 1 expression suppresses allergic airway inflammation in mice. *Am J Respir Crit Care Med*
47 690 173:958-964.
- 49 691 Chung AG, Belone PM, Bimova BV, Karn RC, Laukaitis CM. 2017. Studies of an Androgen-Binding
50 692 Protein Knockout Corroborate a Role for Salivary ABP in Mouse Communication. *Genetics* 205:1517-
51 693 1527.
- 53 694 Clifton BD, Jimenez J, Kimura A, Chahine Z, Librado P, Sanchez-Gracia A, Abbassi M, Carranza F,
54 695 Chan C, Marchetti M, et al. 2020. Understanding the Early Evolutionary Stages of a Tandem
55
56
57
58
59
60

- 1
2
3 696 *Drosophilamelanogaster*-Specific Gene Family: A Structural and Functional Population Study. *Mol Biol*
4 697 *Evol* 37:2584-2600.
- 5
6 698 Denecke S, Rankic I, Driva O, Kalsi M, Luong NBH, Buer B, Nauen R, Geibel S, Vontas J. 2021.
7 699 Comparative and functional genomics of the ABC transporter superfamily across arthropods. *BMC*
8 700 *Genomics* 22:553.
- 9
10 701 Dlouhy SR, Karn RC. 1983. The tissue source and cellular control of the apparent size of androgen
11 702 binding protein (Abp), a mouse salivary protein whose electrophoretic mobility is under the control of
12 703 sex-limited saliva pattern (Ssp). *Biochem Genet* 21:1057-1070.
- 13
14 704 Dover G. 1982. Molecular drive: a cohesive mode of species evolution. *Nature* 299:111-117.
- 15
16 705 Eichler EE, Flint J, Gibson G, Kong A, Leal SM, Moore JH, Nadeau JH. 2010. Missing heritability and
17 706 strategies for finding the underlying causes of complex disease. *Nat Rev Genet* 11:446-450.
- 18
19 707 Emerson JJ, Cardoso-Moreira M, Borevitz JO, Long M. 2008. Natural selection shapes genome-wide
20 708 patterns of copy-number polymorphism in *Drosophila melanogaster*. *Science* 320:1629-1631.
- 21
22 709 Emes RD, Riley MC, Laukaitis CM, Goodstadt L, Karn RC, Ponting CP. 2004. Comparative evolutionary
23 710 genomics of androgen-binding protein genes. *Genome Res* 14:1516-1529.
- 24
25 711 Galtier N, Duret L, Glemin S, Ranwez V. 2009. GC-biased gene conversion promotes the fixation of
26 712 deleterious amino acid changes in primates. *Trends Genet* 25:1-5.
- 27
28 713 Galtier N, Piganeau G, Mouchiroud D, Duret L. 2001. GC-content evolution in mammalian genomes: the
29 714 biased gene conversion hypothesis. *Genetics* 159:907-911.
- 30
31 715 Gibbs RA, Weinstock GM, Metzker ML, Muzny DM, Sodergren EJ, Scherer S, Scott G, Steffen D,
32 716 Worley KC, Burch PE, et al. 2004. Genome sequence of the Brown Norway rat yields insights into
33 717 mammalian evolution. *Nature* 428:493-521.
- 34
35 718 Goodman M, Czelusniak J, William Moore G, Romero-Herrera AE, Matsuda G. 1979. Fitting the gene
36 719 lineage into its species lineage, a parsimony strategy illustrated by cladograms constructed from globin
37 720 sequences. *Systematic Zoology* 28:132-163.
- 38
39 721 Guy L, Kultima JR, Andersson SG. 2010. genoPlotR: comparative gene and genome visualization in R.
40 722 *Bioinformatics* 26:2334-2335.
- 41
42 723 Harel T, Lupski JR. 2018. Genomic disorders 20 years on-mechanisms for clinical manifestations. *Clin*
43 724 *Genet* 93:439-449.
- 44
45 725 Hastings PJ, Lupski JR, Rosenberg SM, Ira G. 2009. Mechanisms of change in gene copy number. *Nat*
46 726 *Rev Genet* 10:551-564.
- 47
48 727 Higuchi DA, Cahan P, Gao J, Ferris ST, Poursine-Laurent J, Graubert TA, Yokoyama WM. 2010.
49 728 Structural variation of the mouse natural killer gene complex. *Genes Immun* 11:637-648.
- 50
51 729 Huddleston J, Eichler EE. 2016. An Incomplete Understanding of Human Genetic Variation. *Genetics*
52 730 202:1251-1254.
- 53
54 731 Hudson RR, Kreitman M, Aguade M. 1987. A test of neutral molecular evolution based on nucleotide
55 732 data. *Genetics* 116:153-159.
- 56
57 733 Hurles M. 2002. Are 100,000 "SNPs" useless? *Science* 298:1509; author reply 1509.

- 1
2
3 734 Hwang JM, Hofstetter JR, Bonhomme F, Karn RC. 1997. The microevolution of mouse salivary
4 735 androgen-binding protein (ABP) paralleled subspeciation of *Mus musculus*. *J Hered* 88:93-97.
5
6 736 Janousek V, Karn RC, Laukaitis CM. 2013. The role of retrotransposons in gene family expansions:
7 737 insights from the mouse *Abp* gene family. *BMC Evol Biol* 13:107.
8
9 738 Janousek V, Laukaitis CM, Yanchukov A, Karn RC. 2016. The Role of Retrotransposons in Gene Family
10 739 Expansions in the Human and Mouse Genomes. *Genome Biol Evol* 8:2632-2650.
11
12 740 Johnson ME, Viggiano L, Bailey JA, Abdul-Rauf M, Goodwin G, Rocchi M, Eichler EE. 2001. Positive
13 741 selection of a gene family during the emergence of humans and African apes. *Nature* 413:514-519.
14
15 742 Kapitonov VV, Jurka J. 2008. A universal classification of eukaryotic transposable elements implemented
16 743 in Repbase. *Nat Rev Genet* 9:411-412; author reply 414.
17
18 744 Karn RC. 1994. The mouse salivary androgen-binding protein (ABP) alpha subunit closely resembles
19 745 chain 1 of the cat allergen Fel dI. *Biochem Genet* 32:271-277.
20
21 746 Karn RC. 1998. Steroid binding by mouse salivary proteins. *Biochem Genet* 36:105-117.
22
23 747 Karn RC, Chung AG, Laukaitis CM. 2014. Did androgen-binding protein paralogs undergo neo- and/or
24 748 Subfunctionalization as the *Abp* gene region expanded in the mouse genome? *PLoS One* 9:e115454.
25
26 749 Karn RC, Clements MA. 1999. A comparison of the structures of the alpha:beta and alpha:gamma dimers
27 750 of mouse salivary androgen-binding protein (ABP) and their differential steroid binding. *Biochem Genet*
28 751 37:187-199.
29
30 752 Karn RC, Dlouhy SR. 1991. Salivary androgen-binding protein variation in *Mus* and other rodents. *J*
31 753 *Hered* 82:453-458.
32
33 754 Karn RC, Laukaitis CM. 2009. The mechanism of expansion and the volatility it created in three
34 755 pheromone gene clusters in the mouse (*Mus musculus*) genome. *Genome Biol Evol* 1:494-503.
35
36 756 Karn RC, Laukaitis CM. 2014. Selection shaped the evolution of mouse androgen-binding protein (ABP)
37 757 function and promoted the duplication of *Abp* genes. *Biochemical Society Transactions* 42:in press.
38
39 758 Karn RC, Nachman MW. 1999. Reduced nucleotide variability at an androgen-binding protein locus
40 759 (*Abpa*) in house mice: evidence for positive natural selection. *Mol Biol Evol* 16:1192-1197.
41
42 760 Karn RC, Orth A, Bonhomme F, Boursot P. 2002. The complex history of a gene proposed to participate
43 761 in a sexual isolation mechanism in house mice. *Mol Biol Evol* 19:462-471.
44
45 762 Karn RC, Young JM, Laukaitis CM. 2010. A candidate subspecies discrimination system involving a
46 763 vomeronasal receptor gene with different alleles fixed in *M. m. domesticus* and *M. m. musculus*. *PLoS*
47 764 *One* 5.
48
49 765 Keane TM, Goodstadt L, Danecek P, White MA, Wong K, Yalcin B, Heger A, Agam A, Slater G,
50 766 Goodson M, et al. 2011. Mouse genomic variation and its effect on phenotypes and gene regulation.
51 767 *Nature* 477:289-294.
52
53 768 Krishnan RS, Daniel JC, Jr. 1967. "Blastokinin": inducer and regulator of blastocyst development in the
54 769 rabbit uterus. *Science* 158:490-492.
55
56 770 Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, Devon K, Dewar K, Doyle M,
57 771 FitzHugh W, et al. 2001. Initial sequencing and analysis of the human genome. *Nature* 409:860-921.

- 1
2
3 772 Laukaitis C, Karn RC. 2012. Recognition of subspecies status mediated by androgen-binding protein
4 773 (ABP) in the evolution of incipient reinforcement on the European house mouse hybrid zone. In:
5 774 Macholan M, Munclinger P, Baird SJ, Pialek J, editors. *Evolution of the House Mouse*. Cambridge, UK:
6 775 Cambridge University Press.
- 8 776 Laukaitis CM, Critser ES, Karn RC. 1997. Salivary Androgen-Binding Protein (Abp) Mediates Sexual
9 777 Isolation in *Mus Musculus*. *Evolution* 51:2000-2005.
- 11 778 Laukaitis CM, Dlouhy SR, Emes RD, Ponting CP, Karn RC. 2005. Diverse spatial, temporal, and sexual
12 779 expression of recently duplicated androgen-binding protein genes in *Mus musculus*. *BMC Evol Biol* 5:40.
- 14 780 Laukaitis CM, Dlouhy SR, Karn RC. 2003. The mouse salivary androgen-binding protein (ABP) gene
15 781 cluster on chromosomes 7: characterization and evolutionary relationships. *Mamm Genome* 14:679-691.
- 17 782 Laukaitis CM, Heger A, Blakley TD, Munclinger P, Ponting CP, Karn RC. 2008. Rapid bursts of
18 783 androgen-binding protein (Abp) gene duplication occurred independently in diverse mammals. *BMC Evol*
19 784 *Biol* 8:46.
- 21 785 Laukaitis CM, Mauss C, Karn RC. 2012. Congenic strain analysis reveals genes that are rapidly evolving
22 786 components of a prezygotic isolation mechanism mediating incipient reinforcement. *PLoS One* 7:e35898.
- 24 787 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R, Genome
25 788 Project Data Processing S. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics*
26 789 25:2078-2079.
- 28 790 Lilue J, Doran AG, Fiddes IT, Abrudan M, Armstrong J, Bennett R, Chow W, Collins J, Collins S,
29 791 Czechanski A, et al. 2018. Sixteen diverse laboratory mouse reference genomes define strain-specific
30 792 haplotypes and novel functional loci. *Nat Genet* 50:1574-1583.
- 32 793 Loire E, Tusso S, Caminade P, Severac D, Boursot P, Ganem G, Smadja CM. 2017. Do changes in gene
33 794 expression contribute to sexual isolation and reinforcement in the house mouse? *Mol Ecol* 26:5189-5202.
- 35 795 Manolio TA, Collins FS, Cox NJ, Goldstein DB, Hindorff LA, Hunter DJ, McCarthy MI, Ramos EM,
36 796 Cardon LR, Chakravarti A, et al. 2009. Finding the missing heritability of complex diseases. *Nature*
37 797 461:747-753.
- 39 798 Morgan AP, Holt JM, McMullan RC, Bell TA, Clayshulte AM, Didion JP, Yadgary L, Thybert D, Odom
40 799 DT, Flicek P, et al. 2016. The Evolutionary Fates of a Large Segmental Duplication in Mouse. *Genetics*
41 800 204:267-285.
- 43 801 Mukherjee A, Chilton BE. 2000. The uteroglobin/clara cell protein family. *Annals of the New York*
42 802 *Academy of Sciences* 932:1-356.
- 45 803 Mukherjee AB, Zhang Z, Chilton BS. 2007. Uteroglobin: a steroid-inducible immunomodulatory protein
46 804 that founded the Secretoglobin superfamily. *Endocr Rev* 28:707-725.
- 48 805 Nagylaki T, Petes TD. 1982. Intrachromosomal gene conversion and the maintenance of sequence
49 806 homogeneity among repeated genes. *Genetics* 100:315-337.
- 51 807 Nguyen DQ, Webber C, Hehir-Kwa J, Pfundt R, Veltman J, Ponting CP. 2008. Reduced purifying
52 808 selection prevails over positive selection in human copy number variant evolution. *Genome Res* 18:1711-
53 809 1723.
- 55 810 Nguyen DQ, Webber C, Ponting CP. 2006. Bias of selection on human copy-number variants. *PLoS*
54 811 *Genet* 2:e20.
- 56 812 Pan D, Zhang L. 2008. Tandemly arrayed genes in vertebrate genomes. *Comp Funct Genomics*:545269.

- 1
2
3 813 Pavlopoulou A, Pampalakis G, Michalopoulos I, Sotiropoulou G. 2010. Evolutionary history of tissue
4 814 kallikreins. *PLoS One* 5:e13781.
- 5
6 815 Perry GH, Ben-Dor A, Tsalenko A, Sampas N, Rodriguez-Revenga L, Tran CW, Scheffer A, Steinfeld I,
7 816 Tsang P, Yamada NA, et al. 2008. The fine-scale and complex architecture of human copy-number
8 817 variation. *Am J Hum Genet* 82:685-695.
- 9
10 818 Perry GH, Dominy NJ, Claw KG, Lee AS, Fiegler H, Redon R, Werner J, Villanea FA, Mountain JL,
11 819 Misra R, et al. 2007. Diet and the evolution of human amylase gene copy number variation. *Nat Genet*
12 820 39:1256-1260.
- 13
14 821 Pezer Z, Chung AG, Karn RC, Laukaitis CM. 2017. Analysis of Copy Number Variation in the Abp Gene
15 822 Regions of Two House Mouse Subspecies Suggests Divergence during the Gene Family Expansions.
16 823 *Genome Biol Evol* 9.
- 17
18 824 Sage RD, Atchley WR, Capanna E. 1993. House mice as models in systematic biology. *Syst. Biol.*
19 825 42:523-561.
- 20
21 826 Sanger F, Nicklen S, Coulson AR. 1977. DNA sequencing with chain-terminating inhibitors. *Proc Natl*
22 827 *Acad Sci U S A* 74:5463-5467.
- 23
24 828 Semple C, Steel M. 2002. Tree Reconstruction from Multi-State Characters. *Advances in Applied*
25 829 *Mathematics* 28:169-184.
- 26
27 830 Sjodin P, Jakobsson M. 2012. Population genetic nature of copy number variation. *Methods Mol Biol*
28 831 838:209-223.
- 29
30 832 Spielmann M, Lupianez DG, Mundlos S. 2018. Structural variation in the 3D genome. *Nat Rev Genet*
31 833 19:453-467.
- 32
33 834 Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, Zhang Y, Ye K, Jun G,
34 835 Fritz MH, et al. 2015. An integrated map of structural variation in 2,504 human genomes. *Nature* 526:75-
35 836 81.
- 36
37 837 Talley HM, Laukaitis CM, Karn RC. 2001. Female preference for male saliva: implications for sexual
38 838 isolation of *Mus musculus* subspecies. *Evolution* 55:631-634.
- 39
40 839 Thybert D, Roller M, Navarro FCP, Fiddes I, Streeter I, Feig C, Martin-Galvez D, Kolmogorov M,
41 840 Janousek V, Akanni W, et al. 2018. Repeat associated mechanisms of genome evolution and function
42 841 revealed by the *Mus caroli* and *Mus pahari* genomes. *Genome Res* 28:448-459.
- 43
44 842 Trifinopoulos J, Nguyen LT, von Haeseler A, Minh BQ. 2016. W-IQ-TREE: a fast online phylogenetic
45 843 tool for maximum likelihood analysis. *Nucleic Acids Res* 44:W232-235.
- 46
47 844 Uriu K, Kosugi Y, Ito J, Sato K. 2021. The Battle between Retroviruses and APOBEC3 Genes: Its Past
48 845 and Present. *Viruses* 13.
- 49
50 846 van Rheede T, Bastiaans T, Boone DN, Hedges SB, de Jong WW, Madsen O. 2006. The platypus is in its
51 847 place: nuclear genes and indels confirm the sister group relation of monotremes and Therians. *Mol Biol*
52 848 *Evol* 23:587-597.
- 53
54 849 Verneau O, Catzeflis F, Furano AV. 1998. Determining and dating recent rodent speciation events by
55 850 using L1 (LINE-1) retrotransposons. *Proc Natl Acad Sci U S A* 95:11284-11289.
- 56
57 851 Vošlajerová Bimová B, Macholán M, Baird SEB, Munclinger P, Laukaitis CM, Karn RC, Luzynski K,
58 852 Tucker P, Piálek J. 2011. Reinforcement selection acting on the European house mouse hybrid zone.
59 853 *Molecular Ecology* 20:2403-2424.

- 1
2
3 854 Wang Y, Li X, Lin Y, Wang Y, Wang K, Sun C, Lu T, Zhang M. 2018. Structural Variation, Functional
4 855 Differentiation, and Activity Correlation of the Cytochrome P450 Gene Superfamily Revealed in
5 856 Ginseng. *Plant Genome* 11.
- 7 857 Waterston RH, Lindblad-Toh K, Birney E, Rogers J, Abril JF, Agarwal P, Agarwala R, Ainscough R,
8 858 Alexandersson M, An P, et al. 2002. Initial sequencing and comparative analysis of the mouse genome.
9 859 *Nature* 420:520-562.
- 11 860 White MA, Ane C, Dewey CN, Larget BR, Payseur BA. 2009. Fine-scale phylogenetic discordance
12 861 across the house mouse genome. *PLoS Genet* 5:e1000729.
- 14 862 Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capi P, Chalhoub B, Flavell A, Leroy P, Morgante M,
15 863 Panaud O, et al. 2007. A unified classification system for eukaryotic transposable elements. *Nat Rev*
16 864 *Genet* 8:973-982.
- 18 865 Xue Y, Sun D, Daly A, Yang F, Zhou X, Zhao M, Huang N, Zerjal T, Lee C, Carter NP, et al. 2008.
19 866 Adaptive evolution of UGT2B17 copy-number variation. *Am J Hum Genet* 83:337-346.
- 21 867 Yang Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol* 24:1586-1591.
- 23 868 Yoneda M, Xu L, Kajiyama H, Kawabe S, Paiz J, Ward JM, Kimura S. 2016. Secretoglobulin Superfamily
24 869 Protein SCGB3A2 Alleviates House Dust Mite-Induced Allergic Airway Inflammation in Mice. *Int Arch*
25 870 *Allergy Immunol* 171:36-44.
- 27 871 Young ND, Zhou P, Silverstein KA. 2016. Exploring structural variants in environmentally sensitive gene
28 872 families. *Curr Opin Plant Biol* 30:19-24.
- 30 873 Zhou P, Silverstein KA, Ramaraj T, Guhlin J, Denny R, Liu J, Farmer AD, Steele KP, Stupar RM, Miller
31 874 JR, et al. 2017. Exploring structural variation and gene family architecture with De Novo assemblies of
32 875 15 *Medicago* genomes. *BMC Genomics* 18:261.

877 **Table 1.** *Abpa* and *Abpbg* genes in each wild-derived mouse genome (B6 refers to mouse reference genome (mm10)).

878

Taxon	Number of genes (unique)	Number of <i>Abpa</i> genes				Number of <i>Abpbg</i> genes				Represented Clades
		Total (unique)	Number of genes (unique)	Number of pseudogenes (unique)	Lineage-specific	Total (unique)	Number of genes (unique)	Number of pseudogenes (unique)	Lineage-specific	
B6	64 (58)	30 (27)	14 (13)	16 (14)	7	34 (31)	12 (11)	22 (20)	7	1-5
WSB	79 (43)	39 (21)	21 (8)	18 (13)	NA	40 (22)	18 (10)	22 (12)	NA	1-5
PWK	41 (38)	20 (18)	13 (11)	7 (7)	1	21 (20)	10 (10)	11 (10)	1	1-5
CAS	72 (46)	36 (22)	24 (11)	12 (11)	2	36 (24)	26 (18)	10 (6)	3	1-5
<i>spr</i>	65 (35)	30 (17)	12 (6)	18 (11)	3	35 (18)	23 (11)	12 (7)	4	1-5
<i>car</i>	40 (33)	21 (17)	9 (6)	12 (11)	9	19 (16)	13 (10)	6 (6)	8	1-3,5
<i>pah</i>	11 (11)	6 (6)	3 (3)	3 (3)	4	5 (5)	4 (4)	1 (1)	3	3, 5

879 **Figure Legends**

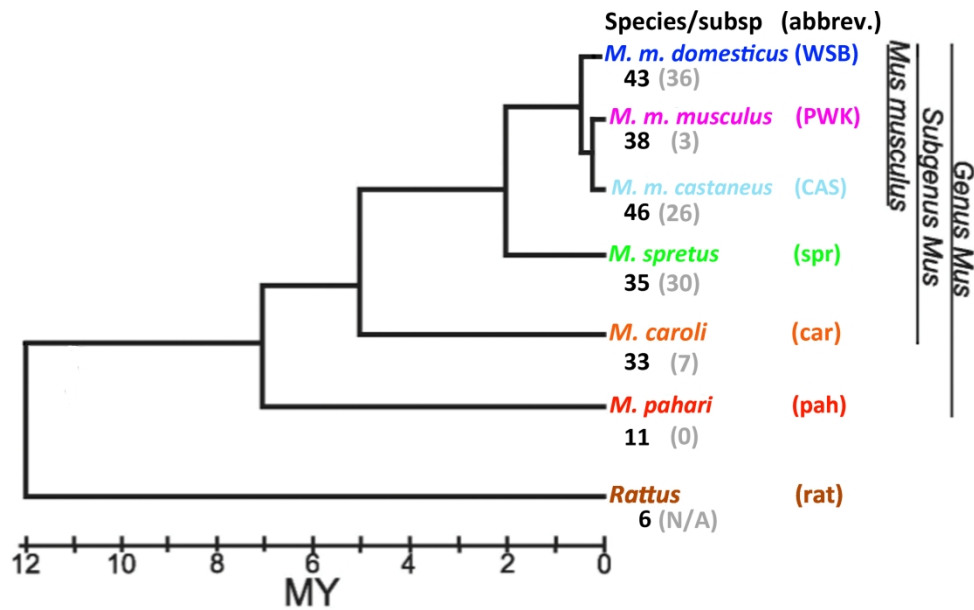
880 **Figure 1.** A canonical phylogeny of murid rodents adapted from Chevret, et al. (2005) to show the
881 divergence of *M. m. domesticus* (strain WSB) from the ancestor of *M. m. musculus* (strain PWK) and
882 *M.m. castaneus* (strain CAS) (White, et al. 2009; Keane, et al. 2011). The seven taxa are differentiated by
883 color. The black numbers under each taxon are the different gene sequences we found and the grey
884 numbers indicate total additional copies (CN) above the diploid number (the CN for each gene is given in
885 **supplementary tables S1-S6**). CN was not determined for the rat genome *Abp* region.

886 **Figure 2.** Gene phylogenies of murid rodent *Abpa* (Panel A) and *Abpbg* (Panel B) genes rooted to the
887 independent Rat clade (brown) and basal *Mus* root (gray). Paralogs from five ancestral clades (Laukaitis,
888 et al. 2008) are indicated by color-coding of branches, red (1), green (2), purple (3), yellow (4), and blue
889 (5), represented by colored bars around the periphery of the phylogeny. The taxon-specific colors of Fig.
890 1 are used for the gene names (not italicized) and genes that root more deeply than individual B6 clades
891 are named with capital letters (e.g. *pah_aW*). Bootstrap values are shown in black. See **supplementary**
892 **figures S1-S2** for parts of these trees broken out to make them easier to read.

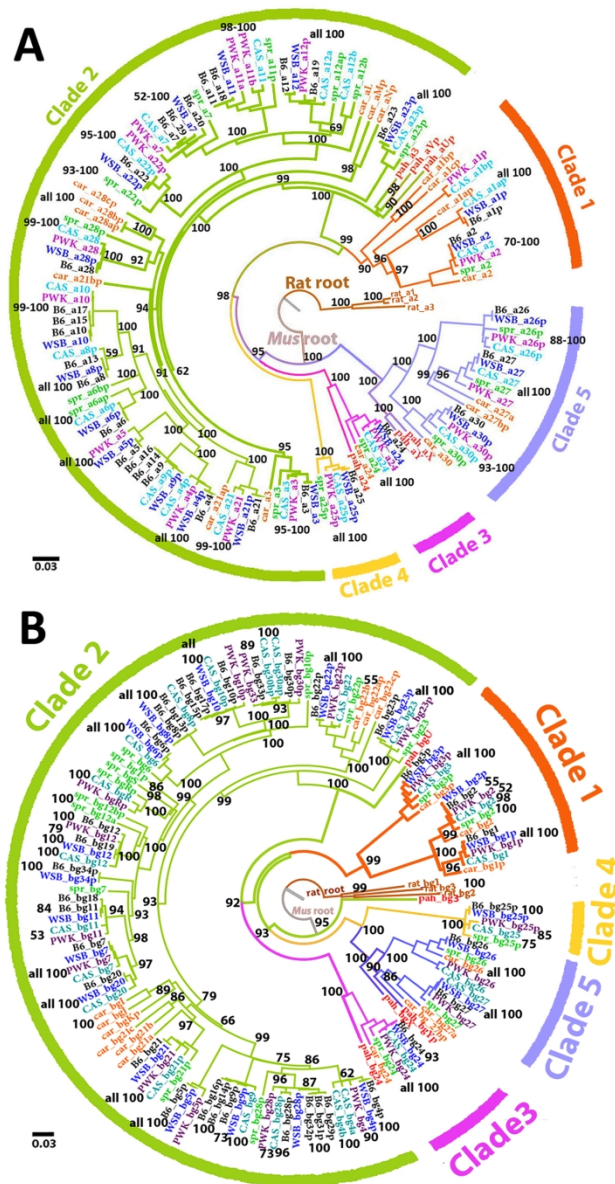
893 **Figure 3.** Relationships between *Abp* paralogs identified in six rodent taxa. Genes that could be
894 expressed are solid-filled blue (*Abpa*) or red (*Abpbg*) arrows while putative pseudogenes are unfilled
895 arrows. Taxon abbreviations with their chromosome in parentheses are shown in a phylogeny on the left.
896 A block of genes unique to the C57BL/6 (B6) reference genome is shown above a blue triangle at the top
897 of the figure and represents genes of the ultimate and penultimate duplications in the WSB lineage (Pezer,
898 et al. 2017). Orthologs between taxa are connected by bands with the colors of the Clades shown in **Fig. 2**
899 and the large Clade numbers at the top of the figure are colored the same. Starting from *pah*, genes with
900 no ortholog in the next taxon are highlighted in blue. Genes are shown in a proportional scale, however
901 regions between genes are reduced tenfold to enable clearer representation of gene relationships. Synteny
902 was plotted with *genoPlotR* package (Guy, et al. 2010).

903 **Figure 4.** Module phylogeny constructed with a LINE1, L1MC3, that is nearly ubiquitous in the intra-
904 modular sequences of gene modules in the genome mouse (mm10) and in the six *Mus* taxa. Five ancestral

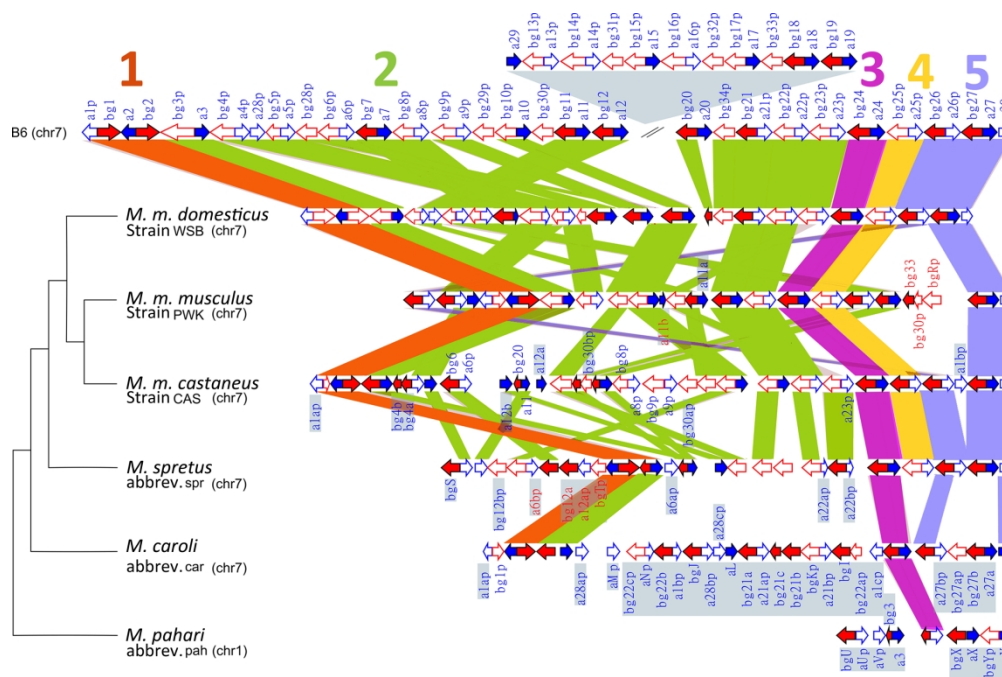
1
2
3 905 clades are labeled in red (1), green (2), purple (3, represented only by pah-M24), yellow (4), and blue (5)
4
5 906 and the gene-specific colors are the same as in **Fig. 2**. Bootstrap values >50 are shown in black.
6
7 907 **Figure 5.** Clades 4 and 5 gene and module phylogenies. Genes and modules with unusual topologies are
8
9 908 shown with red asterisks. *Abpa27* (Panel A, center) has the unexpected topology reported by Karn, et al.
10
11 909 (2002) where the PWK allele is an outgroup to the *spr* allele. The *a26* genes (Panel A, top) also have an
12
13 910 unexpected topology as do the *M27*, *M26* (Panel C) and *M25* (Panel D) modules, and *bg26* (Panel E) and
14
15 911 *bg25* (Panel F) genes. Only *a25* (Panel B) shows an expected topology.
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



A canonical phylogeny of murid rodents adapted from Chevret, et al. (2005) to show the divergence of *M. m. domesticus* (strain WSB) from the ancestor of *M. m. musculus* (strain PWK) and *M. m. castaneus* (strain CAS) (White, et al. 2009; Keane, et al. 2011). The seven taxa are differentiated by color. The black numbers under each taxon are the different gene sequences we found and the grey numbers indicate total additional copies (CN) above the diploid number (the CN for each gene is given in supplementary tables S1-S6). CN was not determined for the rat genome Abp region.

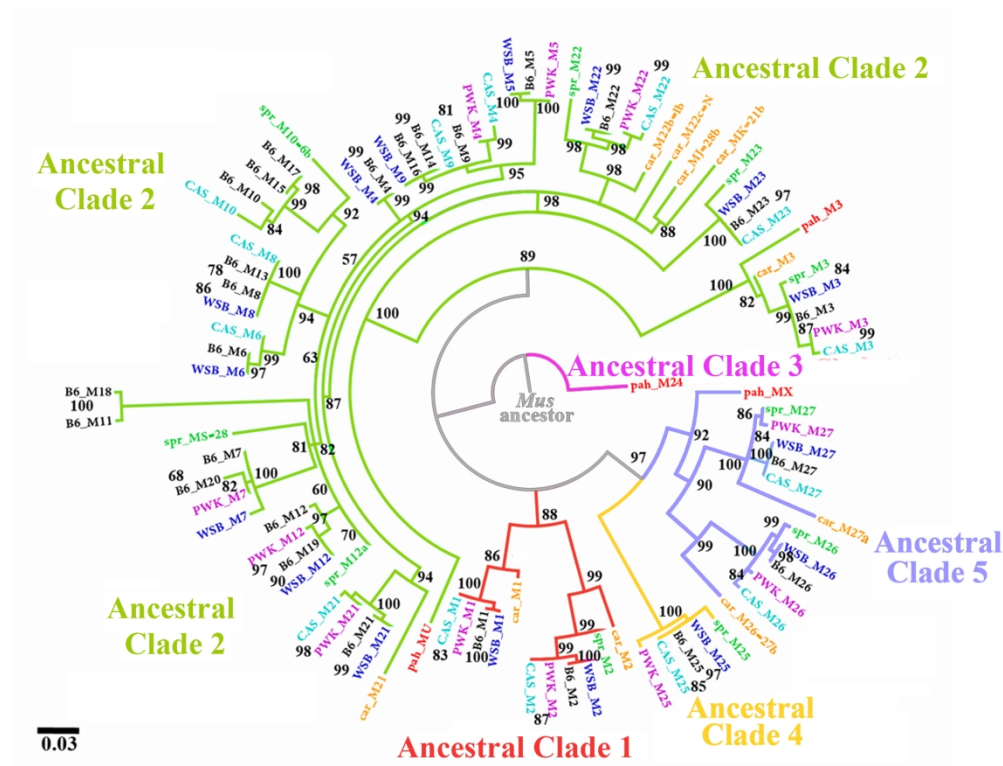


Gene phylogenies of murid rodent Abpa (Panel A) and Abpbg (Panel B) genes rooted to the independent Rat clade (brown) and basal Mus root (gray). Paralogs from five ancestral clades (Laukaitis, et al. 2008) are indicated by color-coding of branches, red (1), green (2), purple (3), yellow (4), and blue (5), represented by colored bars around the periphery of the phylogeny. The taxon-specific colors of Fig. 1 are used for the gene names (not italicized) and genes that root more deeply than individual B6 clades are named with capital letters (e.g. pah_aW). Bootstrap values are shown in black. See supplementary figures S1-S2 for parts of these trees broken out to make them easier to read.

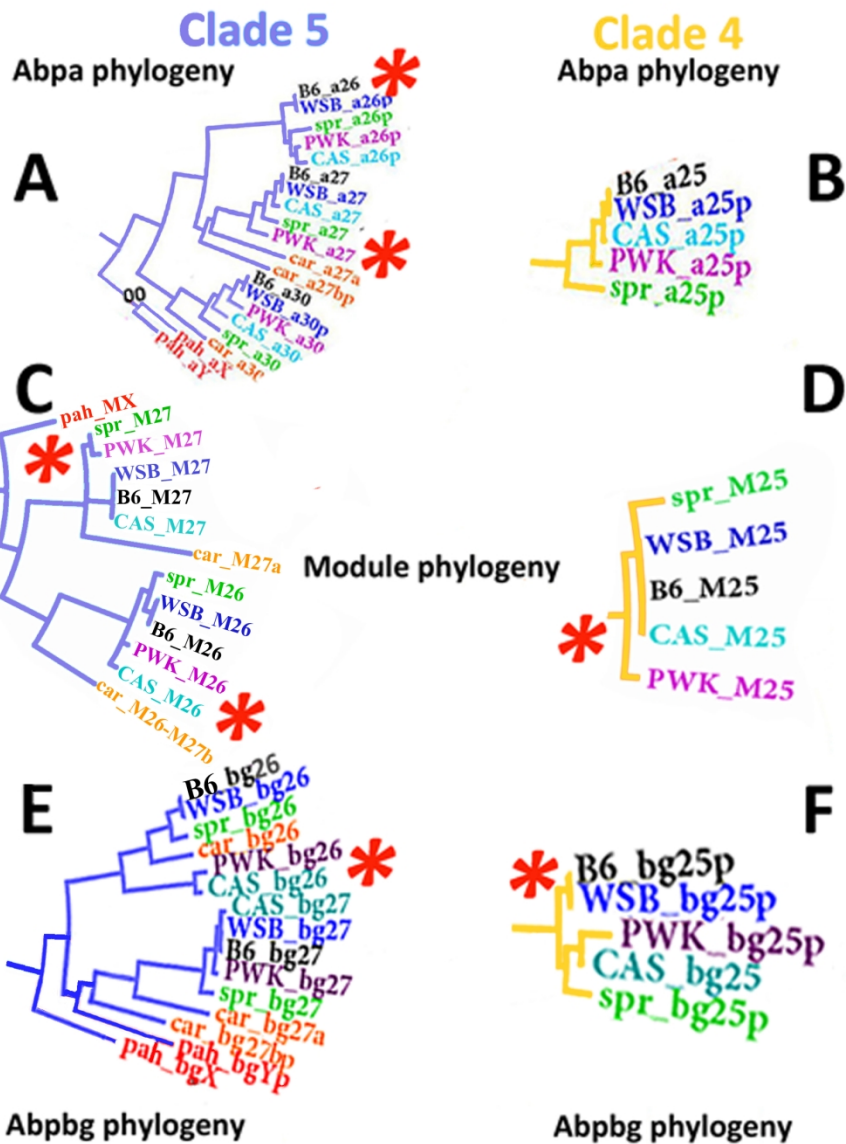


Relationships between Abp paralogs identified in six rodent taxa. Genes that could be expressed are solid-filled blue (Abpa) or red (Abpbp) arrows while putative pseudogenes are unfilled arrows. Taxon abbreviations with their chromosome in parentheses are shown in a phylogeny on the left. A block of genes unique to the C57BL/6 (B6) reference genome is shown above a blue triangle at the top of the figure and represents genes of the ultimate and penultimate duplications in the WSB lineage (Pezer, et al. 2017). Orthologs between taxa are connected by bands with the colors of the Clades shown in Fig. 2 and the large Clade numbers at the top of the figure are colored the same. Starting from pah, genes with no ortholog in the next taxon are highlighted in blue. Genes are shown in a proportional scale, however regions between genes are reduced tenfold to enable clearer representation of gene relationships. Synteny was plotted with genoPlotR package (Guy, et al. 2010).

229x152mm (300 x 300 DPI)



Module phylogeny constructed with a LINE1, L1MC3, that is nearly ubiquitous in the intra-modular sequences of gene modules in the genome mouse (mm10) and in the six *Mus* taxa. Five ancestral clades are labeled in red (1), green (2), purple (3, represented only by pah-M24), yellow (4), and blue (5) and the gene-specific colors are the same as in Fig. 2. Bootstrap values >50 are shown in black.



Clades 4 and 5 gene and module phylogenies. Genes and modules with unusual topologies are shown with red asterisks. Abpa27 (Panel A, center) has the unexpected topology reported by Karn, et al. (2002) where the PWK allele is an outgroup to the spr allele. The a26 genes (Panel A, top) also have an unexpected topology as do the M27, M26 (Panel C) and M25 (Panel D) modules, and bg26 (Panel E) and bg25 (Panel F) genes. Only a25 (Panel B) shows an expected topology.