# UNIWERSYTET JAGIELLOŃSKI
## W KRAKOWIE

## Wydział Biologii i Nauk o Ziemi

### Instytut Nauk o Środowisku

# Molecular basis of adaptation in the bank vole

## Mateusz Konczal

Kraków 2015

*Nothing in Biology Makes Sense Except in the Light of Evolution*

Theodosius Dobzhansky

**Molecular basis of adaptation in the bank vole**

**Author:**  Mateusz Konczal

**Supervisor:**  dr hab. Wiesław Babik

**Reviewers:**  prof. dr hab. Jarosław Burczyk

prof. dr hab. Marek Konarzewski

# TABLE OF CONTENTS

# SUMMARY

Selection experiments combined with genome or transcriptome resequencing represent a promising approach for advancing our understanding of the genetic basis of adaptive evolution. Here, I investigated a bank vole (*Myodes [=Clethrionomys] glareolus*) selection experiment, with four lines selected for aerobic capacity, four lines selected for predatory behavior and four unselected control lines. I developed transcriptomic resources for this non-model species and assessed the accuracy of cost-effective pool RNA-Seq approach. This approach can be used in studies of adaptation on molecular level in organisms without high-quality reference genome assembled. Based on transcriptomic patterns of polymorphism and divergence in the selection experiment I drew conclusions about initial response to selection on molecular level in small populations. I tried to understand how directional selection and genetic drift shape allele frequency and expression dynamics, and what are the consequences of that for studying and understanding adaption in natural populations.

     **Chapter I** is dedicated to testing the accuracy the estimates of allele frequency obtained from RNA pools sequenced using Illumina technology. While sequencing has become relatively inexpensive, the preparation of a large number of sequencing libraries remains costly and technically challenging. Pooling samples is then an attractive alternative, which is often applied in genome resequencing. Unfortunately, high-quality reference genomes are still lacking for many organisms, which complicates population genomic studies. For such species genome-wide information from entire populations may be obtained by sequencing pooled transcriptomes. The study reported in Chapter I shows that the pooled RNA-Seq approach is a reliable and cost-effective strategy for obtaining genome-wide information about potentially functionally relevant variation, provided that high-quality transcriptome assembly and stringent SNP-calling and filtering criteria based on individual sequencing are used. This result is based on data from 10 vole liver transcriptomes, sequenced both as an individually barcoded libraries and as a pool. Liver reference transcriptome was assembled *de novo* and around 24 thousands high-quality single nucleotide polymorphisms (SNPs) were identified. Allele frequency estimates were highly correlated between individual and pooled samples, indicating that pooled RNA-Seq exhibits an accuracy comparable with pooled genome resequencing. In this part of the thesis I also proposed a strategy which may be used to filter out SNPs potentially associated with transcriptome-specific errors.

This strategy was applied for studying molecular-level response to artificial selection. In **Chapter II** I investigated liver and heart transcriptomes of bank voles selected for increased aerobic metabolism. Samples were collected from four selected and four control lines at 13th generation of selection. At that point, selectively bred voles had 48% higher maximum rates of oxygen consumption than did the control ones. Liver and heart transcriptomes were assembled and annotated to evaluate molecular-level response to selection in genes which are expressed in those organs. The two transcriptomes contained transcripts of more than 18,000 and 11,000 known protein coding genes, respectively. Analyses of both transcriptomes allowed to identify 172,246 SNPs, a vast majority of them being located in putatively protein-coding sequences. For each SNP allele frequencies were calculated within each line, and mean pairwise distances between lines ($F_{ST}$ distance) were used to estimate genetic differentiation between selectively bred and control lines. I did not find evidence for separate clustering of the four lines selected for aerobic performance and the four control lines. To compare allele frequency changes with expectations from drift I performed pedigree-based simulations. Again, I did not find evidence for mechanisms other than drift driving differentiation between treatments. On the other hand, expression analyses showed that liver transcriptome-wide pattern of expression distinguished selected from controls lines, and over 300 genes were found to be differentially expressed between treatments in heart or liver. Hence, these results support the hypothesis that initial molecular-level response to selection occurs primarily through changes in gene expression. In order to understand molecular underpinnings of phenotypic selection, I chose genes which are plausible targets of selection. They were selected from the most differentiated genes between treatments, based on information about gene functions and the biology of the selected phenotype. Some of those genes are associated with mobilizing fats and sugars from body reserves, with stress response and mating success.

Response to selection for another trait - predatory behavior - is the topic of **Chapter III**. In this study, after 13 generations, the proportion of voles showing predatory behavior was 5 times higher in the selected lines than in unselected control lines. I sequenced transcriptomes of liver and hippocampus of the selected and control animals to investigate patterns of repeatability of selection on molecular level. Assembled hippocampus transcriptome contained 21,407 sequences of protein coding genes. Both transcriptomes allowed to identify 179,468 polymorphisms. In contrast to the results obtained for aerobic performance, here I found higher than expected repeatable differentiation in allele

frequencies. This is mainly explained by noncoding and synonymous changes, suggesting their role in the evolution of gene expression or alternative splicing. Analyses of gene expression showed that hippocampus pattern of expression of selected lines clusters separately from the control ones, and 149 genes were differentially expressed between treatments in hippocampus or in liver. The difference between aerobic and predatory lines in SNPs analyses suggests that the architecture of adaptive variation differs depending on trait. More repeatable changes in allele frequencies in predatory lines may result from selection for variants of higher frequencies in base population or stronger selection per variant. Finally, investigation of genes most differentiated between predatory and control lines points to the potential role of hunger, aggression, biological rhythm and functioning of the nervous system in shaping the response to selection.

The selection experiment was established using animals derived from a natural population and at the onset of selection the laboratory population contained an amount and spectrum of genetic variation similar to that typically found in populations of small mammals. Therefore, the results presented here are relevant to the understanding of the molecular basis of complex adaptation. They support the hypothesis that changes in gene expression play predominant role at early stages of adaptive evolution. My results suggest also that the molecular pattern and repeatability of adaptation are likely a function of the genetic architecture of the selected trait. Finally, methodological and transcriptomic resources developed and presented here can be used in future studies of molecular basis of adaptation in non-model species.

# STRESZCZENIE

Analiza eksperymentów selekcyjnych poprzez sekwencjonowanie genomów lub transkryptomów jest obiecującą strategią, która może pomóc w zrozumieniu genetycznego podłoża ewolucji adaptatywnej. W prezentowanej pracy wykorzystuję tę strategię do badania genetycznych podstaw adaptacji u nornicy rudej (*Myodes [=Clethrionomys] glareolus*). Podczas eksperymentu selekcyjnego, cztery niezależne linie selekcjonowane były na wysoką wydolność tlenową, cztery na zachowania drapieżnicze, a cztery pozostałe nieselekcjonowane linie traktowane były jako kontrola. W mojej pracy analizowałem osobniki pochodzące z 13. pokolenia selekcji. W pierwszej kolejności prezentuję zasoby transkryptomowe dla tego niemodelowego gatunku oraz szacuję dokładność estymowania frekwencji alleli w zmieszanych próbkach RNA pochodzących od kilku osobników. Na podstawie transkryptomowych wzorców zmienności i zróżnicowania pomiędzy liniami, wnioskuję o początkowej odpowiedzi na dobór na poziomie molekularnym. Próbuję zrozumieć, jak dobór kierunkowy i dryf genetyczny kształtują dynamikę zmian frekwencji alleli oraz poziomu ekspresji. Analizuję jakie są tego konsekwencje dla badania i zrozumienia adaptacji w naturalnych populacjach.

Rozdział I dotyczy testowania dokładności oszacowania frekwencji alleli ze zmieszanych próbek RNA sekwencjonowanych przy użyciu technologii Illumina. Chociaż sekwencjonowanie staje się relatywnie tanie, przygotowanie dużej liczby indywidulanie znakowanych bibliotek pozostaje kosztowne oraz trudne technicznie. Mieszanie próbek pochodzących z wielu osobników i przygotowywanie z nich jednej biblioteki, jest w takim przypadku rozsądną alternatywą, często stosowaną przy sekwencjonowaniu genomów. Niestety, dobrze złożone i anotowane referencyjne genomy nie są dostępne dla wielu gatunków, co utrudnia prowadzenie na nich ogólnogenomowych analiz. Dla takich gatunków populacyjna ogólnogenomowa informacja może być uzyskana poprzez sekwencjonowanie zmieszanych próbek transkryptomów. W rozdziale pierwszym pokazuję, że takie podejście jest wiarygodną i oszczędną strategią uzyskania ogólnogenomowej informacji o funkcjonalnie ważnej zmienności genetycznej. Należy jednak zwrócić uwagę na zapewnienie wysokiej jakości referencyjnego transkryptomu i ustalić rygorystyczne kryteria wywoływania oraz filtrowania polimorfizmów pojedynczych nukleotydów (SNPów), opartych na próbce pojedynczo zsekwencjonowanych osobników. Ta strategia testowana jest przy użyciu transkryptomów pochodzących z wątrób 10 osobników nornicy rudej. Transkryptomy zsekwencjonowane zostały indywidualnie oraz jako jedna zmieszana próbka. Referencyjny

transkryptom złożony został *de novo* i użyty do wywołania ok. 20 tys. SNPów. Frekwencje alleli były silnie skorelowane pomiędzy oszacowaniami pochodzącymi ze zmieszanej próbki a danymi uzyskanymi z indywidualnego sekwencjonowania każdego osobnika. Wyniki świadczą o tym, że szacowanie frekwencji alleli ze zmieszanych próbek transkryptomów ma podobną dokładność do analogicznych analiz opartych na sekwencjonowaniu genomów. W tym rozdziale proponuję także strategię, która może być użyta do filtrowania SNiPów potencjalnie związanych ze specyficznymi dla transkryptomów błędami.

To podejście zostało użyte w badaniach molekularnych podstaw odpowiedzi na dobór sztuczny w eksperymencie selekcyjny. W rozdziale drugim analizuję transkryptomy pochodzące z wątrób i ser nornic selekcjonowanych na wysoki metabolizm tlenowy. Próbki zostały pobrane z czterech linii selekcyjnych (linie A) i z czterech kontrolnych (linie C) w 13. pokoleniu selekcji. W tym pokoleniu, nornice z linii A miały o 48% wyższe maksymalne tempo konsumpcji tlenu w porównaniu z nornicami z linii C. Transkryptomy wątroby i serca zostały złożone i anotowane w celu użycia ich do oszacowania odpowiedzi na dobór na poziomie molekularnym. W tych transkrytpomach zidentyfikowałem ok 18 tys. (wątroba) i 11 tys. sekwencji kodujących znane białka. Analizy wszystkich zsekwencjonowanych próbek pozwoliły na zidentyfikowanie 172,246 SNPów, z których zdecydowana większość znajdowała się w sekwencjach kodujących białka. Dla każdego SNPa frekwencje alleli policzone zostały dla każdej linii i użyte do oszacowania średniego dystansu genetycznego pomiędzy liniami (dystans $F_{ST}$). Nie znalazłem dowodów na oddzielnie grupowanie się czterech linii selekcyjnych w stosunku do linii kontrolnych. Do porównania zmian frekwencji alleli z oczekiwaniami z dryfu przeprowadziłem symulacje komputerowe oparte o znane rodowody oraz oszacowane frekwencje alleli w populacji wyjściowej. Po raz kolejny nie znalazłem dowodów na działanie innych mechanizmów niż dryf, które powodują różnicowanie się linii selekcjonowanych i kontrolnych. Z drugiej strony, analizy zmian poziomu ekspresji genów pokazały, że ogólnotranskryptomowy profil ekspresji zmienił się w wątrobie, a ponad 300 genów miało istotnie różne poziomy ekspresji pomiędzy liniami A i C. Te wyniki wspierają hipotezę mówiącą o tym, że w początkowym okresie działania doboru, odpowiedź na poziomie molekularnym związana jest głównie ze zmianami w regulacji ekspresji genów. Aby zrozumieć molekularne podstawy selekcjonowanej cechy wybrałem potencjalnych kandydatów leżących u podstaw odpowiedzi na dobór. Zostały one wybrane spośród genów które wykazywały największe zróżnicowanie pomiędzy liniami A i C, używając przy tym informacji o ich molekularnej funkcji oraz wiedzy o selekcjonowanej

cesze. Geny te związane są z mobilizacją energii z zapasów organizmu, odpowiedzią na stres i sukcesem rozrodczym.

Odpowiedź na dobór na inną cechę – zachowania drapieżnicze – jest tematem trzeciego rozdziału. Po 13 pokoleniach selekcji, cztery linie selekcjonowane na zachowania drapieżnicze (linie P), miały pięć razy większą proporcję nornic prezentujących zachowania drapieżnicze, niż cztery linie kontrolne (linie C). Zsekwencjonowałem transkryptomy pochodzące z wątrób i hipokampów nornic z linii P i C, by uzyskać odpowiedź na pytania dotyczące molekularnych podstaw zachowań drapieżniczych oraz powtarzalności wzorców ewolucji na poziomie molekularnym. Złożenie transkryptomu hipokampu pozwoliło zidentyfikować 21,407 sekwencji genów kodujących białka, a oba transkryptomy zidentyfikować 179,468 SNPów. W przeciwieństwie do wyników otrzymanych dla nornic selekcjonowanych na wysoka wydolność tlenową, linie P i C różnicowały się we frekwencjach alleli bardziej niż wynikałoby to z dryfu. Efekt ten wyjaśniają głównie niekodujące i synonimowe SNPy, co sugeruje ich potencjalną rolę w ewolucji ekspresji lub alternatywnego splicingu. Analiza ekspresji genów pokazała że wzorzec ekspresji zmienił się istotnie w hipokampie a 149 genów ulegało istotnie różnej ekspresji w hipokampie lub wątrobie. Różnica pomiędzy liniami selekcjonowanymi na wysoki metabolizm tlenowy i tymi selekcjonowanymi na zachowania drapieżnicze sugeruje, że architektura genetyczna tych cech istotnie się różni. Bardziej powtarzalne zmiany frekwencji alleli w liniach drapieżniczych mogą być związane z wariantami o wyższych częstościach w populacji wyjściowej albo z silniejszym doborem działającym na pojedyncze warianty. Geny najbardziej zróżnicowane pomiędzy liniami drapieżniczymi i kontrolnymi sugerują potencjalną rolę zmian w procesach związanych z głodem, agresją, rytmem okołodobowym i funkcjonowaniem układu nerwowego w odpowiedzi na dobór na zachowanie drapieżnicze.

Studiowany eksperyment selekcyjny został założony z osobników pochodzących z naturalnej populacji, dlatego ilość i spektrum zmienności genetycznej była dobrym przybliżeniem zmienności dostępnej dla doboru w populacjach naturalnych małych ssaków. Prezentowane wyniki są więc relewantne dla zrozumienia molekularnych podstaw adaptacji. Wspierają one hipotezę zakładającą, że zmiany poziomu ekspresji genów odgrywają zasadniczą rolę w początkowych stadiach adaptatywnej ewolucji. Wyniki te sugerują także, że molekularne wzorce adaptacji i ich powtarzalność mogą być funkcją architektury genetycznej selekcjonowanej cechy. Zaprezentowane strategie badawcze oraz zasoby transkryptomowe mogą zostać użyte w dalszych badaniach molekularnych podstaw adaptacji.

# GENERAL INTRODUCTION

In the world of finite resources, some organisms will make more efficient use of them and so will leave more descendants than their less efficient relatives. If the variation between organisms in producing offspring includes a heritable component, this force, known as natural selection, pushes populations towards a phenotype that better fits the current environment (Darwin 1859; Fisher 1930). This simple process is called adaptation.

The importance of adaptive evolution is indisputable. However, it remains controversial what proportion of the entire evolutionary change at molecular level arises from natural selection (Nei 2005; Hahn 2008; Sella et al. 2009; Wagner 2010). Random mutations are the primary source of genetic variation, but the extraordinary amount of genetic variation between populations or species can result from two different mechanisms: either deterministic (selection) or stochastic (genetic drift) processes. Thus, it is critical for evolutionary biology to figure out which of them dominates evolution (Hahn 2008; Nei 2010).

Adaptation can be studied not only as a process but also as an outcome. An adaptive trait is one that increases organism's fitness in a particular environmental context and has been, hence, the target of natural selection. There is extensive debate going on regarding the prevalence of adaptive traits, the mechanisms by which they arise and the levels at which selection operates (Nowak et al. 2010; Abbott et al. 2011; Barrett and Hoekstra 2011).

There are several uncertainties concerning adaptation at the molecular level, for which recent advances in sequencing technology promise to find explanations. One of the problems concerns the source of adaptive variation. A population can adapt by the means of either new mutations or variants already present in the population (standing genetic variation) (Barrett and Schluter 2008). The source of adaptive variation determines the strength of signatures of positive selection in genomes: adaptive evolution from standing genetic variation leaves more subtle signs in genomes (soft sweeps) than adaptation due to new mutations (hard sweeps) (Hermisson and Pennings 2005; Teshima et al. 2006; Messer and Petrov 2013). This may lead to underestimation of the relative importance of adaptation at the molecular level, if natural selection utilizes mainly standing genetic variation. However, recent findings suggest that the source of adaptive variation appears to vary among evolutionary lineages. Extensive work on microorganisms has contributed to our understanding of adaptation scenarios driven by new mutations (Herring et al. 2006; Barrick et al. 2009; Tenaillon et al. 2012). On the other hand, in multicellular, sexually reproducing

species (the subject of this dissertation) standing genetic variation is the main source of variation at the initial stage of adaptive evolution (Barrett and Schluter 2008; Teotónio et al. 2009; Burke et al. 2010).

Because most of the theory on the genetics of adaptation has focused on adaptation from new mutations, many questions about the dynamics, circumstances and consequences of adaptation from standing variation remain unanswered (Barrett and Schluter 2008). They concern for instance the distribution of fitness effect sizes (Barrett and Schluter 2008),as well as parallelism and convergence at molecular level (Conte et al. 2012).

The genetic basis of parallel adaptations has intrigued researchers for years because its understanding may help to answer questions about repeatability of evolution (Stern and Orgogozo 2009; Radwan and Babik 2012). If adaptation is generally due to allele frequency changes at loci with standing variation, then evolution can proceed in parallel among derived populations experiencing similar environmental conditions (Teotónio et al. 2009). However, this parallel process should depend on population size and the architecture of genetic variation available for selection. Thus, the debate involves also the question whether, and under what circumstances, parallel evolution occurs on the level of nucleotide, gene and molecular pathway (Elmer and Mayer 2011).

Another issue concerns the role that gene expression changes play in adaptation. King and Wilson (1975) proposed that adaptive evolutionary change is largely due to changes in gene expression, and there is empirical evidence both supporting (Wray 2007; Jones et al. 2012) and contradicting this view (Hoekstra and Coyne 2007). A growing number of studies suggests that regulation of gene expression is a common source of adaptation, but performing a comprehensive comparison between expression and structural changes remains a challenging task and such comparison has been done only recently for humans on a genome-wide scale (Fraser 2013).

In order to study molecular basis of adaptation, scientists tried to identify alleles that affect adaptive phenotype (Glazier et al. 2002; Olson-Manning et al. 2012), or to scan genomes for signs of positive selection (Storz 2005; Biswas and Akey 2006). Such studies, targeting natural populations and utilizing whole-genome sequencing technology, allowed to collect interesting observations improving our understanding of genetic basis of adaptation. For example studies of stickleback fish demonstrated the importance of standing genetic variation and changes in regulatory elements for adaptive evolution (Jones et al. 2012). After the retreat of Pleistocene glaciers, marine sticklebacks adapted to newly formed freshwater

habitats, exhibiting repeatable changes in body shape, skeletal armour, pigmentation, life history and mating performance. By sequencing genomes of marine and freshwater fish researchers demonstrated that regulatory changes dominate their adaptive evolution, which occurred mainly from globally shared standing genetic variation (Jones et al. 2012). Another study focused on the relative importance of adaptive and neutral evolution (Brawand et al. 2014). By sequencing the genomes of African cichlid fish the authors demonstrated that both processes are necessary for generating new, highly diverse species in very short periods of time. African cichlids are a well-known example of rapid adaptive radiation with around 2,000 known species. Analyses of genomes from five lineages suggest that neutral processes were crucial for retaining genomic variation in cichlids, whereas selection subsequently sorted some of this variation (Brawand et al. 2014). In another study, researchers analyzed deer mice that have recently colonized light-colored soil of Nebraska Sand Hills (Linnen et al. 2013). Their strongly adaptive light coat color is composed of multiple traits and arose from many independent mutations within a single gene. This shows that even adaptive evolution of a simple Mendelian trait may proceed in complex manner.

The studies of natural populations have some limitations, due to lack of replications, complex population history and population structure. Another important issue which also needs to be emphasized is that they concern the outcome of adaptation processes, but often make conclusions on the process itself. The pattern does not necessarily reveal the process because multiple scenarios may produce similar patterns. Thus, it was argued that molecular studies of natural populations have rarely altered fundamental understanding of the relationship between evolution of a genotype and evolution of a phenotype (Rockman 2012; Travisano and Shaw 2013).

An alternative research framework that offers the opportunity to study evolutionary processes such as adaptation and genetic drift is experimental evolution (Kawecki et al. 2012). Under such approach hypotheses and theories concerning evolution are tested by the use of controlled experiments. Recently, selection experiments have gained new value, namely insight into molecular-level response to selection. Thus, evolve and resequencing (E&R, Turner et al. 2011) promises to unify two branches of genetics: molecular and populations genetics, by identifying loci contributing to adaptation (molecular genetics) and analyzing trajectories of those loci during adaptation (population genetics) (Kofler and Schlötterer 2014). In the wake of the latest sequencing technologies and the ongoing drop in DNA sequencing costs, this ultimate goal has now become reachable. Moreover, the advantages of

E&R studies - controlled conditions, selection pressure, demography and history, studying evolutionary processes in the real time, and finally the ability to replicate experiments under identical conditions - allow to overcome the limitations which are an inherent part of the studies of natural populations.

Most of the E&R studies of multicellular organisms were performed on fruit flies. Burke et al. (2010) studied flies that had been selected for accelerated development over 600 generations. They concluded that the probability of fixation of selected variants is relatively low and that selection does not readily expunge genetic variation. Other studies confirm the predominant role of complex evolutionary trajectories of selected variants as well as emphasize the extraordinary importance of standing genetic variation and epistasis (Teotónio et al. 2009; Turner et al. 2011; Huang et al. 2012; Orozco-Terwengel et al. 2012).

The E&R studies are less common for vertebrates, which usually have smaller population sizes, longer generation times and larger genomes. Therefore, much more effort on both experimental and molecular levels is required to perform such investigations. One of the few examples includes an experiment where mice were selected for voluntary activity (Swallow et al. 1998). The murine model was used to show heritability of the predisposition to engage in voluntary exercise with epistatic interactions accounting for a considerable amount of genetic variation (Kelly and Pomp 2013). However, a limitation for the conclusions drawn from this example is the fact, that it used laboratory model species for which the standing genetic variation may significantly deviate from natural populations. Other experiments on vertebrates suffer also from the nature of standing genetic variation in base populations, or have problems with achieving the required number of replicates or a satisfactory population size (Johansson et al. 2010; Kukekova et al. 2011; Chan et al. 2012).

In the present dissertation I make use of state-of-the-art sequencing technology and bioinformatic tools to study molecular-level response to selection in a unique vertebrate experiment. Bank voles (*Myodes [=Clethrionomys] glareolus*) have been selected in two ecologically important directions: for increased aerobic metabolism and for predatory behavior (Sadowska et al. 2008). The base population of the selection experiment originated from animals captured in the Niepołomice Forest, and thus represents good approximation of genetic variation segregating in natural populations of small vertebrates. After 13 generations, the voles from selected lines (in each direction) differed in the selection trait from the four control lines by more than 2 standard deviations.

Bank vole is an important organism in evolutionary, ecological, behavioral and medical studies (Kotlik et al. 2006; Nonno et al. 2006; Radwan et al. 2008; Sadowska et al. 2008; Tschirren et al. 2012), but unfortunately this small common European rodent does not have the reference genome assembled. Therefore, genomic analyses of bank vole are much more challenging than in case of model species, such as fruit flies or mice (Ekblom and Galindo 2011; Ekblom and Wolf 2014). Such situation is still common for many species (called non-model species) studied by ecologists and evolutionary biologists. This is because a large genome assembly is a non-trivial task, which requires a lot of work and resources. Promising alternative for that is now provided by transcriptome sequencing (RNA-Seq) (Mortazavi et al. 2008; Vijay et al. 2013). In Chapter I I discuss benefits of RNA-Seq for non-model species, I present the assembly of bank vole liver transcriptome and assess the accuracy of allele frequency estimation with the pooled RNA approach. This technical chapter is dedicated to test cost-effective pooled RNA-Seq approach and shows that sequenced pools of transcriptomes can be used as an alternative approach for population genomic analyses if reference genome is unavailable. I make use of these conclusions in the next two chapters, where pooled RNA-Seq is applied to study molecular-level response to selection.

In Chapter II I test the hypothesis that adaptation is mainly associated with changes in gene expression. By sequencing and analyzing liver and heart transcriptomes of voles selected for increased metabolism I found support for that. Transcriptomic-wide pattern of expression changed drastically, whereas there is no evidence for selection-driven changes in nonsynonymous SNPs. This conclusion was based on analyses of distance matrices. To compare SNPs allele frequency differentiation with expectations from drift I developed tools simulating drift and selection on known pedigree. Comparison of observed and simulated data showed that no repeatable changes in allele frequencies could be unambiguously attributed to directional selection, although the low power of these analyses limited the resulting inferences to large effect variants only.

Similar analyses were performed for hippocampus and liver transcriptomes of voles selected for predatory behavior (Chapter III). Again, lack of rapid nonsynonymous changes was accompanied by substantial expression differentiation. In contrast to the selection for aerobic performance however, for predatory lines I found repetitive changes in other than nonsynonymous classes of SNPs. They are potentially associated with expression or alternative splicing. Selection for predatory behavior seems to be associated with variants of larger effects or variants segregating in higher frequencies in the base populations. This

potential difference in the genetic architecture of those two traits may affect the repeatability of selection on molecular level. Finally, I identified candidate genes potentially underlying the selected traits. They should be prioritized as a target for future research, as potentially underlying ecologically important traits.

In the context of adaptation genomics, my findings support the hypothesis claiming that changes in gene expression play predominant role in adaptive evolution. My results suggest also that molecular pattern and repeatability of response to selection are likely a function of the genetic architecture of the selected trait. The methodological part of this dissertation provides evidence that pooled RNA-Seq approach can be widely used in adaptation studies  in non-model species and that pedigree-based simulations are a powerful method of evaluating genome-wide effects of selection. Specifically for research on adaptation in bank vole, this dissertation presents transcriptomic data and a list of candidate genes which can be used in future laboratory and field studies on adaptation in this species.

# References

Abbot P, Abe J, Alcock J, Alizon S, Alpedrinha JA, Andersson M, ..., Gardner A. 2011. Inclusive fitness theory and eusociality. Nature 471(7339): E1-E4.

Barrett RDH, Hoekstra HE. 2011. Molecular spandrels: tests of adaptation at the genetic level. Nat. Rev. Genet. 12:767–780.

Barrett RDH, Schluter D. 2008. Adaptation from standing genetic variation. Trends Ecol. Evol. 23:38–44.

Barrick JE, Yu DS, Yoon SH, Jeong H, Oh TK, Schneider D, Lenski RE, Kim JF. 2009. Genome evolution and adaptation in a long-term experiment with Escherichia coli. Nature 461:1243–1247.

Biswas S, Akey JM. 2006. Genomic insights into positive selection. Trends Genet. 22:437–446.

Brawand D, Wagner CE, Li YI, Malinsky M, Keller I, Fan S, …, Di Palma F. 2014. The genomic substrate for adaptive radiation in African cichlid fish. Nature 513: 375–381.

Burke MK, Dunham JP, Shahrestani P, Thornton KR, Rose MR, Long AD. 2010. Genome-wide analysis of a long-term evolution experiment with Drosophila. Nature 467:587–590.

Chan YF, Jones FC, McConnell E, Bryk J, Bünger L, Tautz D. 2012. Parallel selection mapping using artificially selected mice reveals body weight control loci. Current Biology 22: 794-800.

Conte GL, Arnegard ME, Peichel CL, Schluter D. 2012. The probability of genetic parapllelism and convergence in natural populations. Proc. R. Soc. B 279: 5039-5047.

Darwin C. 1859. On the origin of species by means of natural selection, or the preservation of favoured races in the struggle for life. London, UK: John Murray.

Ekblom R, Galindo J. 2011. Applications of next generation sequencing in molecular ecology of non-model organisms. Heredity *107*(1): 1-15.

Ekblom R, Wolf JBW. 2014. A field guide to whole-genome sequencing, assembly and annotation. Evol. Appl. 7(9): 1026-1042.

Elmer KR, Meyer A. 2011. Adaptation in the age of ecological genomics: insights from parallelism and convergence. Trends in ecology & evolution 26(6): 298-306.

Fisher RA. 1930. The genetical theory of natural selection. Oxford University Press.

Fraser HB. 2013. Gene expression drives local adaptation in humans. Genome research, 23(7): 1089-1096.

Glazier AM, Nadeau JH, Aitman TJ. 2002. Finding genes that underlie complex traits. Science 298:2345–2349.

Hahn MW. 2008. Toward a selection theory of molecular evolution. Evolution. 62:255–265.

Hermisson J, Pennings PS. 2005. Soft sweeps: Molecular population genetics of adaptation from standing genetic variation. Genetics 169:2335–2352.

Herring CD, Raghunathan A, Honisch C, Patel T, Applebee MK, Joyce AR, Albert TJ, Blattner FR, van den Boom D, Cantor CR, et al. 2006. Comparative genome sequencing of Escherichia coli allows observation of bacterial evolution on a laboratory timescale. Nat. Genet. 38:1406–1412.

Hoekstra HE, Coyne J a. 2007. The locus of evolution: Evo devo and the genetics of adaptation. Evolution. 61:995–1016.

Huang W, Richards S, Carbone MA, Zhu D, Anholt RR, Ayroles JF, ..., Mackay TF. 2012. Epistasis dominates the genetic architecture of Drosophila quantitative traits. Proceedings of the National Academy of Sciences 109(39): 15553-15559.

Johansson AM, Pettersson ME, Siegel PB, Carlborg Ö. 2010. Genome-wide effects of long-term divergent selection. PLoS genetics 6: e1001188.

Jones FC, Grabherr MG, Chan YF, Russell P, Mauceli E, Johnson J, …, Kingsley DM. 2012. The genomic basis of adaptive evolution in threespine sticklebacks. Nature 484:55–61.

Kawecki TJ, Lenski RE, Ebert D, Hollis B, Olivieri I, Whitlock MC. 2012. Experimental evolution. Trends Ecol. Evol. 27:547–560.

Kelly SA, Pomp D. 2013. Genetic determinants of voluntary exercise. Trends in Genetics, 29(6): 348-357.

King MC, Wilson AC. 1975. Evolution at two levels in humans and chimpanzees. Science 188: 107-116.

Kukekova AV, Trut LN, Chase K, Kharlamova AV, Johnson JL, Temnykh SV, ..., Lark KG. 2011. Mapping loci for fox domestication: deconstruction/reconstruction of a behavioral phenotype. Behavior genetics 41(4): 593-606.

Kofler R, Schlötterer C. 2014. A guide for the design of evolve and resequencing studies. Mol. Biol. Evol. 31:474–483.

Kotlík P, Deffontaine V, Mascheretti S, Zima J, Michaux JR, Searle J B. 2006. A northern glacial refugium for bank voles (Clethrionomys glareolus). Proceedings of the National Academy of Sciences, 103(40), 14860-14864.

Linnen CR, Poh YP, Peterson BK, Barrett RD, Larson JG, Jensen JD, Hoekstra HE. 2013. Adaptive evolution of multiple traits through multiple mutations at a single gene. Science, 339(6125): 1312-1316.

Messer PW, Petrov D a. 2013. Population genomics of rapid adaptation by soft selective sweeps. Trends Ecol. Evol. 28:659–669.

Mortazavi A, Williams B a, McCue K, Schaeffer L, Wold B. 2008. Mapping and quantifying mammalian transcriptomes by RNA-Seq. Nat. Methods 5:621–628.

Nei M. 2005. Selectionism and neutralism in molecular evolution. Mol. Biol. Evol. 22:2318–2342.

Nei M, Suzuki Y, Nozawa M. 2010. The neutral theory of molecular evolution in the genomic era. Annual review of genomics and human genetics.11: 265-289.

Nonno R, Di Bari MA, Cardone F, Vaccari G, Fazzi P, Dell'Omo G, ..., Agrimi U. 2006. Efficient transmission and characterization of Creutzfeldt–Jakob disease strains in bank voles. PLoS pathogens 2(2): e12.

Nowak MA, Tarnita CE, Wilson EO. 2010. The evolution of eusociality. Nature, 466(7310): 1057-1062.

Olson-Manning CF, Wagner MR, Mitchell-Olds T. 2012. Adaptive evolution: evaluating empirical support for theoretical predictions. Nat. Rev. Genet. 13:867–877.

Orozco-Terwengel P, Kapun M, Nolte V, Kofler R, Flatt T, Schlãtterer C. 2012. Adaptation of Drosophila to a novel laboratory environment reveals temporally heterogeneous trajectories of selected alleles. Mol. Ecol. 21:4931–4941.

Radwan J, Tkacz A, Kloch A. 2008. MHC and preferences for male odour in the bank vole. Ethology 114:827-833.

Radwan J, Babik W. 2012. The genomics of adaptation. Proceedings of the Royal Society B: Biological Sciences 279(1749): 5024-5028.

Rockman M V. 2012. The QTN program and the alleles that matter for evolution: All that's gold does not glitter. Evolution. 66:1–17.

Sadowska ET, Baliga-Klimczyk K, Chrząścik KM, Koteja P. 2008. Laboratory Model of Adaptive Radiation: A Selection Experiment in the Bank Vole. Physiol. Biochem. Zool. 81:627–640.

Sella G, Petrov D a., Przeworski M, Andolfatto P. 2009. Pervasive natural selection in the Drosophila genome? PLoS Genet. 5.

Stern DL, Orgogozo V. 2009. Is genetic evolution predictable?. Science 323(5915): 746-751.

Storz JF. 2005. Using genome scans of DNA polymorphism to infer adaptive population divergence. Mol. Ecol. 14:671–688.

Swallow JG, Carter PA, Garland Jr T. 1998. Artificial selection for increased wheel-running behavior in house mice. Behavior genetics. 28(3): 227-237.

Tenaillon O, Rodríguez-Verdugo A, Gaut RL, McDonald P, Bennett AF, Long AD, Gaut BS. 2012. The molecular diversity of adaptive convergence. Science 335(6067), 457-461.

Teotónio H, Chelo IM, Bradić M, Rose MR, Long AD. 2009. Experimental evolution reveals natural selection on standing genetic variation. Nat. Genet. 41:251–257.

Teshima KM, Coop G, Przeworski M. 2006. How reliable are empirical genomic scans for selective sweeps? Genome research 16(6): 702-712.

Travisano M, Shaw RG. 2013. Lost in the map. Evolution 67:305–314.

Tschirren B, Andersson M, Scherman K, Westerdahl H, Råberg L. 2012. Contrasting patterns of diversity and population differentiation at the innate immunity gene toll-like receptor 2 (tlr2) in two sympatric rodent species. Evolution 66:720-731.

Turner TL, Stewart AD, Fields AT, Rice WR, Tarone AM. 2011. Population-based resequencing of experimentally evolved populations reveals the genetic basis of body size variation in Drosophila melanogaster. PLoS Genet. 7.

Vijay N, Poelstra JW, Künstner A, Wolf JBW. 2013. Challenges and strategies in transcriptome assembly and differential gene expression quantification. A comprehensive in silico assessment of RNA-seq experiments. Mol. Ecol. 22:620–634.

Wray G. 2007. The evolutionary significance of cis-regulatory mutations. Nat. Rev. Genet. 8:206–216.

# CHAPTER I

**Accuracy of allele frequency estimation using pooled RNA-Seq**

M. Konczal, P. Koteja, M.T. Stuglik, J. Radwan, W. Babik

**Abstract**

For nonmodel organisms, genome-wide information that describes functionally relevant variation may be obtained by RNA-Seq following de novo transcriptome assembly. While sequencing has become relatively inexpensive, the preparation of a large number of sequencing libraries remains prohibitively expensive for population genetic analyses of nonmodel species. Pooling samples may be then an attractive alternative. To test whether pooled RNA-Seq accurately predicts true allele frequencies, we analysed the liver transcriptomes of 10 bank voles. Each sample was sequenced both as an individually barcoded library and as a part of a pool. Equal amounts of total RNA from each vole were pooled prior to mRNA selection and library construction. Reads were mapped onto the de novo assembled reference transcriptome. High-quality genotypes for individual voles, determined for 23 682 SNPs, provided information on 'true' allele frequencies; allele frequencies estimated from the pool were then compared with these values. 'True' frequencies and those estimated from the pool were highly correlated. Mean relative estimation error was 21% and did not depend on expression level. However, we also observed a minor effect of interindividual variation in gene expression and allele-specific gene expression influencing allele frequency estimation accuracy. Moreover, we observed strong negative relationship between minor allele frequency and relative estimation error. Our results indicate that pooled RNA-Seq exhibits accuracy comparable with pooled genome resequencing, but variation in expression level between individuals should be assessed and accounted for. This should help in taking account the difference in accuracy between conservatively expressed transcripts and these which are variable in expression level.

**Introduction**

Next-generation sequencing (NGS) technologies have resulted in enormous progress not only in the field of medicine but also in the fields of ecology and evolutionary biology. Comparative studies of natural variation at the molecular level have yielded important insights into the evolutionary history of populations, as well as the genomics of adaptation and speciation (Gilad et al. 2009; Rice et al. 2011; Radwan & Babik 2012). For example, NGS technologies have recently been instrumental in enabling findings as impressive and varied as evidence of interbreeding between modern humans and Neanderthals (Reich et al. 2010), the discovery that adaptive evolution results from standing genetic variation in the stickleback (Jones et al. 2012) and the identification of epistasis as one of the most important factors in evolution (Breen et al. 2012). Nonmodel, ecologically well-characterized organisms are being studied at a scale and with a precision unimaginable a few years ago (Ekblom & Galindo 2011). Unfortunately, high-quality reference genomes are still lacking for many organisms that are commonly used in evolutionary and ecological studies, mainly because the de novo assembly of complex genomes that include a large number of repetitive sequences remains a challenging task (Brenchley et al. 2012). In such cases, genome-wide information that describes functionally relevant variation may be obtained through RNA sequencing (RNA-Seq) that utilizes de novo reference transcriptome assembly. This approach has been broadly used in ecological genomics (Vera et al. 2008; Babik et al. 2010; Jeukens et al. 2010; Wolf et al. 2010; Salem et al. 2012).

RNA-Seq is an approach in which RNA molecules are selected, reverse-transcribed and then sequenced using an NGS platform (Mortazavi et al. 2008). Genome complexity and redundancy are reduced because only transcribed sequences are used, which enable the *de novo* assembly of entire transcripts, even when a relatively modest amount of sequence data are available (Martin & Wang 2011). It is important to note that RNA-Seq does not reduce genomic complexity randomly, but rather produces reads from regions in which a large proportion of functionally relevant variation is expected to be located (Jones et al. 2012). Such variation may be assessed and compared with known variation in genes in other organisms without requiring any pre-existing genomic information. Gene expression, alternative splicing patterns and the association of both with phenotypic traits may be also studied using RNA-Seq (Lu et al. 2010; Barbosa-Morais et al. 2012).

RNA-Seq is usually less costly than genome resequencing. However, if transcripts with low levels of expression are to be assembled, greater sequencing depth may be required, which increases the overall cost. Furthermore, the cost of preparing a large number of RNA-Seq libraries, for example from many individuals, is still prohibitively high. An attractive possible solution to this problem is sample pooling (i.e. a pooled RNA-Seq). However, meaningful inferences from pooled RNA data require that allele frequencies estimated from pooled samples adequately reflect true allele frequencies. In case of RNAseq, uncertainty about population allele frequency arises not only because of sampling finite number of individuals, but also from additional stochasticity introduced due to differences in expression level among genes or even among alleles of the same gene. It may bias allele frequency estimates drastically, and to our knowledge, the extent to which these RNAseq-specific issues bias allele frequency estimates has not been explored.

Pooling strategies using DNA samples (Pool-Seq) have been comprehensively tested (Sham et al. 2002; Futschik & Schlootterer 2010; Kim et al. 2010; Gompert & Buerkle 2011; Li 2011; Zhu et al. 2012), and they share some of the difficulties of the pooled RNA-Seq approach. In both Pool-Seq and pooled RNA-Seq approaches, the error associated with allele frequency estimates is inversely proportional to 'true' allele frequency. Several computational approaches that have been proposed to find rare variants in DNA pools and estimate their frequencies (Druley et al. 2009; Bansal 2010) could possibly be applied in pooled RNA-Seq analyses as well. Furthermore, variability introduced by technical errors (inaccurate pipetting, sequencing errors, etc.) is expected to be similar for RNA and DNA samples. However, three sources of error specific to pooled RNA-Seq have not been previously studied: (i) variation in expression levels among individuals, (ii) variation in expression levels among loci and (iii) allele-specific gene expression (Fig. 1).

Substantial differences in gene expression levels commonly occur among individuals of the same sex or developmental stage and are attributable to differences in genetic background and environment. For example, Whitehead and Crawford (2006) showed that 64% of genes are differentially expressed among individuals of the teleost fish genus Fundulus. Other studies argue that gene expression varies extensively both within and among populations (Sandberg et al. 2000; Morloy et al. 2004; Oleksiak et al. 2005; Lynch & Wagner 2008; Barbosa-Morais et al. 2012). In the pooled RNA-Seq approach, interindividual variation in expression level may bias estimates of allele frequency because different individuals contribute unequal numbers of reads. If individuals differ greatly in their expression of a given

gene, allele frequency estimates will be biased towards individuals with higher expression levels.

Interlocus variation in expression levels produces enormous differences in sequencing coverage, which may cause differences in the accuracy of allele frequency estimates for different loci. In non-normalized RNA-Seq analyses, gene expression levels may differ by six orders of magnitude (Mortazavi et al. 2008). The estimated allele frequencies for genes expressed at low levels will be less accurate than those obtained for genes covered by millions of reads. This problem is known to occur in transcriptomic studies, but it has not been studied in the context of pooling.

The third major issue is allele-specific gene expression (Serre et al. 2008; Ge et al. 2009). Cis-acting regulation or epigenetic silencing may cause differential expression of a diploid individual's two alleles. Although allele-specific gene expression is a widespread phenomenon that affects the expression of 20% of genes, allele expression ratios higher than 70:30 are rather rare (Serre et al. 2008). As a result, heterozygotes can in the vast majority of cases be successfully identified, given sufficient sequencing depth (Skelly et al. 2011). However, using pooling techniques, we expect frequency estimates to be distorted for over- and underexpressed alleles.

Although both potentially attractive and inexpensive, the utility of pooled RNA-Seq may be limited by the above issues, and thus, the accuracy of the allele frequency estimates obtained from pooled data should be characterized empirically. Building on results of Pool-Seq studies, we explore here additional, RNA-Seq specific, aspects of allele frequency estimation. Our general aim is to determine how various aspects of expression level variation influence allele frequency estimation.

To examine the accuracy of allele frequency estimates obtained with a pooled RNA-Seq approach, we used bank vole (*Myodes [=Clethrionomys] glareolus*) liver transcriptomes. This rodent species is an important organism in evolutionary, ecological and behavioural studies (Sadowska et al. 2005; Radwan et al. 2008; Boratynski & Koteja 2009; Mokkonen et al. 2011; Tschirren et al. 2012). The bank vole genome is not available, and a high-quality reference genome is unlikely to become available in the near future. The bank vole thus serves as a good example of a nonmodel organism for which obtaining genome-wide data is an important but nontrivial task. RNA samples from the livers of 10 voles were sequenced to generate both individually barcoded libraries and a pooled sample. Allele frequencies were

estimated from the pool and then compared with the 'true' frequencies obtained from the individual libraries.
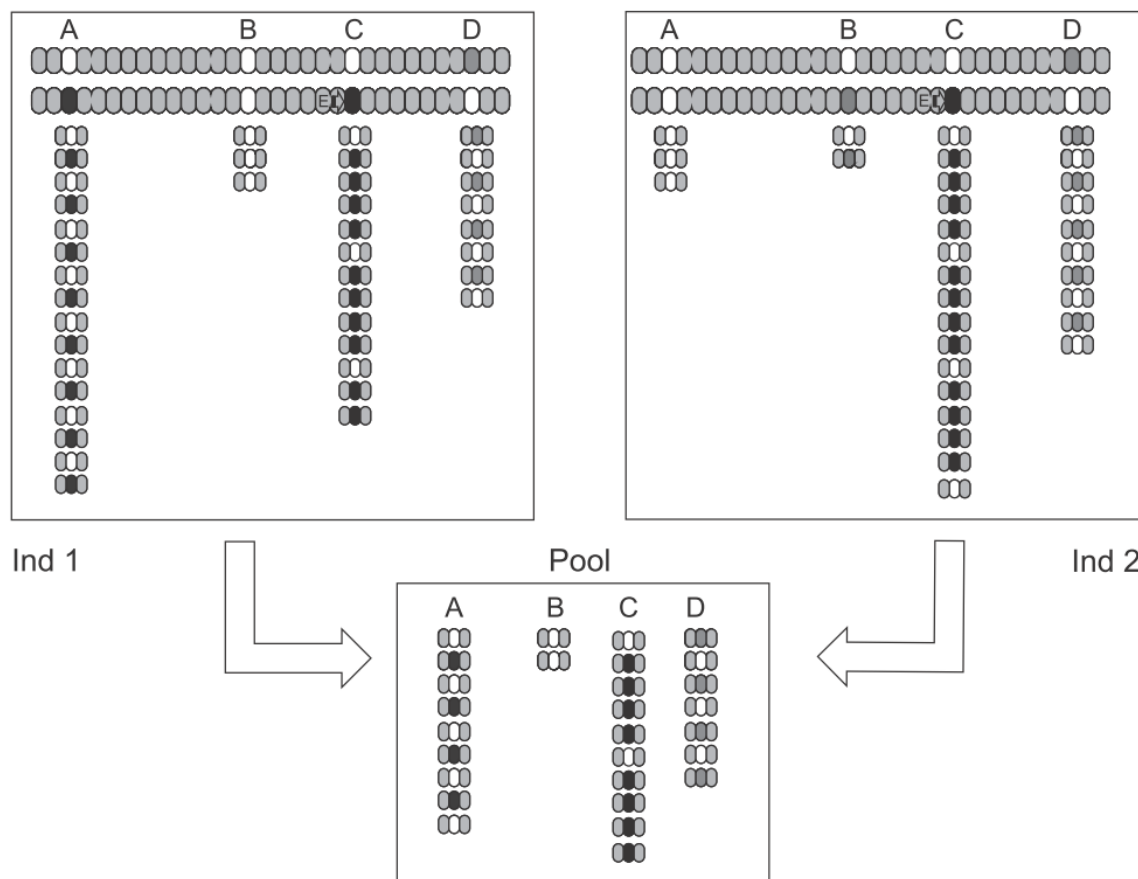


**Fig. 1** Transcriptome-specific sources of error in allele frequency estimates obtained from a pooled sample. Interindividual variation in gene expression (A), interlocus variation in gene expression (B) and allele-specific gene expression (C) are compared with a locus for which the allele frequency estimate is not biased (D).

**Materials and methods**

*Sample collection*

Liver samples were obtained from ten bank voles (*Myodes [=Clethrionomys] glareolus*) from a single control line (unselected) of an artificial selection experiment (generation 13), designed to study correlated evolution of behavioural and physiological traits (Sadowska et al. 2008). The laboratory colony was created using voles captured in the Niepołomice Forest near Krakow (Poland) in 2000. Details related to colony protocols have been provided elsewhere (Sadowska et al. 2008). The experimental protocols were approved by the I Local Ethical Committee in Krakow (decision number 99/2006).

Five male and five female voles, each 75–80 days old, were euthanized using an overdose of isoflurane (Aerrane). The animals were dissected immediately, and a small part (ca. 0.01 g) of the left liver lobe was placed in RNAlater. The samples were stored overnight at 4 °C and then frozen at -20 °C.

Total RNA was extracted using RNAzol (Molecular Research Center) in accordance with the manufacturer's instructions. Residual DNA was removed with a DNA-free Kit (Ambion). RNA concentration and quality were measured using Nanodrop and Agilent 2100 Bioanalyzer, respectively. All samples had an RNA integrity number (RIN) > 7.0, which indicated quality sufficient for poly-A selection and cDNA library preparation.

The pooled sample was prepared using an equal amount of total RNA from each individual. RNA concentration and quality in the pool were assessed as described earlier.

In the final step, the eleven RNA samples (ten individual and one pooled) were used in poly-A selection and the preparation of barcoded cDNA libraries by the Georgia Genomics Facility. Paired-end 2 x 100 bp sequencing was performed in one lane of an Illumina HiSeq 2000, a process that produced a similar number of reads from all the individually tagged samples together and from the pool.

*Reference transcriptome reconstruction*

After trimming low-quality reads using DynamicTrim (Cox et al. 2010), the Trinity assembler (version 2012-06-08) was employed to reconstruct the bank vole liver

transcriptome de novo (Grabherr et al. 2011). For computational reasons, only reads from the pool were used in the assembly. We then processed the Trinity output by merging transcripts that were probably derived from the same genomic location and subsequently produced transcriptome-based gene models (M. Stuglik, W. Babik & J. Radwan, unpublished data). In brief, in the first step of this process, we aggregated Trinity transcripts with overlapping ends using CAP3 (Huang & Madan 1999). The cut-offs for overlap length and per cent identity of the overlap were 40 bases and 99%, respectively. In the next step, we discarded contigs that were entirely contained within other sequences using CD-HIT (Li & Godzik 2006) (settings: identity 0.95 and word size 8). Finally, MegaBLAST was employed to merge all sequences that shared at least 70% of the length of the shortest sequence and had a minimum identity value of 0.96. Contigs were clustered, aligned and merged to form a single consensus sequence. The 'reference transcriptome' that results from this procedure should contain sequences from all exons from all genes that are expressed in at least one transcript and should thus correspond to an assembly of transcriptome-based gene models.

*Mapping, SNP calling and allele frequency estimation*

Because mapping algorithms take into account quality scores, we used nontrimmed reads when mapping and SNP calling. Reads that mapped onto multiple locations in the reference transcriptome were discarded. Reads were mapped onto the reconstructed reference transcriptome using Bowtie 2 (2.0.0-beta6) and employing a very sensitive alignment approach (Langmead & Salzberg 2012). The resulting bam file was post-processed using SAMtools (Li et al. 2009). SNP calling was performed separately for the 10 individually tagged samples and for the pool using mpileup in SAMtools. For SNP calling in individual samples, the default settings were applied; for SNP calling in the pool, a flat prior for the allele frequency spectrum was used.

Low-quality SNPs were filtered out of the VCF file that contained information on the individual genotypes. We excluded SNPs with individual genotypes that were based on less than five reads, and sites at which more than two variants were present. We then retained only SNPs that were reliably genotyped for all 10 individuals (Phred scores of at least 30 for SNP quality and individual genotype quality). Moreover, we discarded all contigs that contained one or more SNPs that would have led us to classify 9 or 10 of the individuals as heterozygotes. As the probability of obtaining such a sample by chance, even assuming equal

allele frequencies for both variants, is only 0.01, reads that mapped onto such contigs were most probably derived from highly similar paralogues. Such stringent filtering practices allowed us to classify the genotypes at these polymorphic sites as high-quality SNPs with known 'true' allele frequencies in the sample. In the next step, we assessed how accurately these 'true' values were reflected by the pool.

*Accuracy estimates*

For each high-quality SNP position, the number of non-reference bases was calculated ($N_O$). The expected number of non-reference bases ($N_E$) was the 'true' allele frequency estimated from individually tagged samples multiplied by the coverage at the SNP position. The accuracy of the allele frequency estimates was quantified as the relative estimation error, which was defined as the absolute value of $(N_E - N_O)/N_E$.

To quantify the effect of allele-specific expression level (ASE) on relative estimation error, we first selected contigs showing evidence of ASE using the following procedure. For each contig, one SNP with the highest number of heterozygotes (max 8 for the reasons explained earlier) was selected. Then, for each heterozygous individual, the hypothesis of equal expression of both alleles was tested (chi-squared test), using the number of reads derived from each allele. SNPs with at least 80% of heterozygotes showing P < 0.001 were considered as indicators of contigs exhibiting ASE. Mean relative estimation error was compared between ASE genes and SNPs randomly selected from the data (sampling the same number of SNPs from each MAF class as in genes with ASE), and the significance of the difference between these two groups was tested using randomization test.

To assess the effect of inaccuracy in allele frequency estimation on the results of a typical population genetic analysis, we simulated a Wright–Fisher population (N e = 10 000, u = $10^{-9}$), in which the expected distribution of allele frequencies is given by equation $\varphi(i) = 4N_e u/i$ (where $0 < i < 2N$; i is the number of copies of the derived allele) (Charlesworth & Charlesworth 2010). We estimated $F_{ST}$ under two scenarios: (i) differentiation was only due to sampling error (allele frequencies in the sample were known precisely) and (ii) differentiation was due to errors resulting from both sampling a limited number of individuals and form estimation of allele frequency from pool. In each of 10 000 simulations, we sampled one SNP from the expected distribution of allele frequencies and simulated two samples of 10

individuals each (sampling from binomial distribution with P set to population allele frequency) and calculated $F_{ST}$ according to the formula $(H_T - H_S)/H_T$ (Hartl & Clark 2006). Next, we simulated second scenario by adding estimation error caused by pooling. We replaced the sample allele frequencies by frequencies randomly drawn from our empirical results obtained from pool for the given 'true' allele frequency. We calculated $F_{ST}$ and compared $F_{ST}$ distributions between two scenarios using the Wilcoxon signed-rank test.

*Accuracy of gene expression estimation*

To estimate gene expression, we used RSEM package (Li & Dewey 2011). We performed TMM normalization (Robinson & Oshlack 2010) to account for differences in the mass of the RNA-Seq samples and thus provide a scaling parameter for each sample. This parameter was then used to calculate the fragments per kilobase of transcript per million fragments mapped (FPKM). FPKM was calculated for each transcriptome-based gene model in each sample. Accuracy was estimated for each contig with a mean FPKM value higher than one. Relative estimation error was calculated in the same way as for allele frequency. The mean expression level calculated from 10 individuals was used as the expected value, and observed values were the FPKM values calculated using the pool.

**Results**

*Reference transcriptome assembly*

A total of 194.1 million read pairs (2 x 100 bp) were obtained; the average per individual was 8.0 (SD 0.41) million pairs, and 114.1 million pairs were re-covered from the pool. Trimming resulted in the removal of 9.6% of the bases. Trimmed reads from the pool were used to assemble the liver transcriptome *de novo*; 181 698 contigs (contig length max: 16 742 bp; mean: 1111.8 bp; median: 429 bp; N50: 2662 bp) totalling 202.0 megabases were generated. Transcriptome assemblers attempt to reconstruct the sequences of all the transcripts present in the sample, which results in considerable redundancy in the assembled transcriptome – a large fraction of exons will be represented many times, reflecting their presence in multiple alternatively spliced transcripts. While such redundancy reflects biological reality, it is undesirable if one wants to construct transcriptome-based gene models in order to detect

polymorphism. We therefore further processed the results generated by Trinity using a custom pipeline that aims to produce transcriptome-based gene models, or a 'reference transcriptome'. The reference transcriptome comprised 146 758 contigs (contig length max, mean, median and N50, respectively: 16 742 bp, 702.7 bp, 353 bp and 1225 bp) and had a total length of 103.1 Mb (Table 1). These contigs represented protein and nonprotein coding sequences expressed in the liver.

**Tab. 1** Overview of the assembly of a hepatic transcriptome for bank voles (*Myodes [=Clethrionomys] glareolus*). The transcriptome was assembled using Trinity and filtered using CAP3, CD-HIT, and MegaBLAST. Statistics for the final transcriptome-based gene models are provided in the last column

|  | **TRINITY** | **CAP3** | **CD-HIT** | **MEGABLAST** |
|---|---|---|---|---|
| Min contig length | 201 | 201 | 201 | 201 |
| Max contig length | 16 742 | 16 742 | 16 742 | 16 742 |
| Mean contig length | 1 111.8 | 1 061.6 | 1 030.6 | 702.7 |
| SD contig length | 1 529.4 | 1 487.4 | 1 440.9 | 977.3 |
| Median contig length | 429 | 411 | 406 | 353 |
| N50 | 2 662 | 2 587 | 2 504 | 1 225 |
| N contigs | 181 698 | 173 496 | 171 077 | 146 758 |
| N contigs > 1kb | 51 988 | 46 524 | 44 551 | 23 512 |
| N contigs in N50 | 22 252 | 2 0648 | 20 371 | 19 101 |
| N bases in contigs | 202 007 816 | 184 191 537 | 176 303 671 | 103 123 071 |
| N bases in contigs > 1kb | 151 525 798 | 135 260 359 | 127 596 475 | 56 436 078 |

*Mapping and SNP calling*

Reads that mapped uniquely onto the reference transcriptome (83.8% of raw reads) were used to identify polymorphic sites. SNPs were identified separately in individually tagged samples and the pool. SNP calling performed on the 10 individually barcoded libraries yielded 264,310 putative SNPs and 40,277 short indels that had scaled Phred quality scores of greater than 10.

The same analysis on the pooled sample (which differed only in the use of the flat prior for the allele frequency distribution) yielded 246,122 putative SNPs and 40,621 short indels. High-quality SNPs were further analysed. We found 95 contigs that were extremely heterozygous at one or more sites (probably representing pairs of paralogues), and all SNPs

from these sequences were discarded. In total, we identified 23 682 high-quality polymorphisms within 4,128 contigs. Only 6,336 (26.8%) of high-quality SNPs within 2380 (57.7%) contigs were called from the pool. Not surprisingly, polymorphisms with rare variants (Fig. 2) and a low proportion of alternative variant reads were under-represented among the SNPs called from the pool. We found that 7% of SNPs with minor allele frequency (MAF) values of less than 0.25 and 74% of SNPs with MAF values greater than 0.25 were called from the pool.
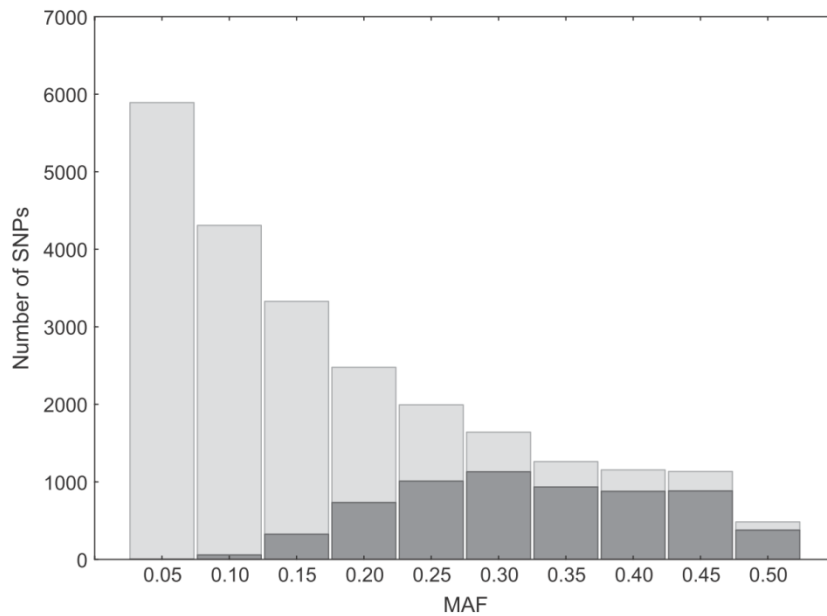


**Fig. 2** Polymorphic sites discovered in the pool. The number of identified (dark) and unidentified high-quality SNPs in the pooled sample.

*Accuracy of allele frequency estimation*

The observed and expected number of reads were strongly correlated ($R^2 = 0.96$; $P < 10^{-14}$; Fig 3). Mean relative estimation error was $0.21 \pm 0.001$ SE (median 0.16), and it was negatively correlated with the minor allele frequency (Fig. 4). Relative estimation error was relatively high for SNPs for which MAF equaled 0.05 (mean = 0.33, median = 0.25), but it decreased significantly when MAF was greater than 0.25 (mean = 0.12, median = 0.09) (Fig. 5). However, the absolute differences in frequencies did not decrease with increasing MAF

(Fig. 6). We found a very weak negative correlation between the sequencing depth for a given SNP ('SNP expression level') and the relative estimation error ($R^2 = 0.002$; $P < 10^{-11}$; Figs 5 and 7).

Relative estimation error correlated significantly with the coefficient of variation in gene expression level among individuals ($R^2 = 0.04$, $P < 0.10^{-15}$; Fig. 8). We identified 43 contigs (containing 283 SNPs) with signatures of ASE. These genes have higher relative estimation error than randomly sampled genes (mean = 0.32, $P < 0.0001$, randomization test).

$F_{ST}$ values were generally overestimated in the pool simulation ($mean_{ind} = 0.026$, $mean_{pool} = 0.033$; Wilcoxon test: $P < 10^{-15}$). Also, we observed some extreme outliers for pools (0.3% observations higher than twice maximum $F_{STind}$), which suggests that in some cases, $F_{ST}$ may be strongly overestimated due to inaccuracy in estimation of allele frequencies introduced by pooling.
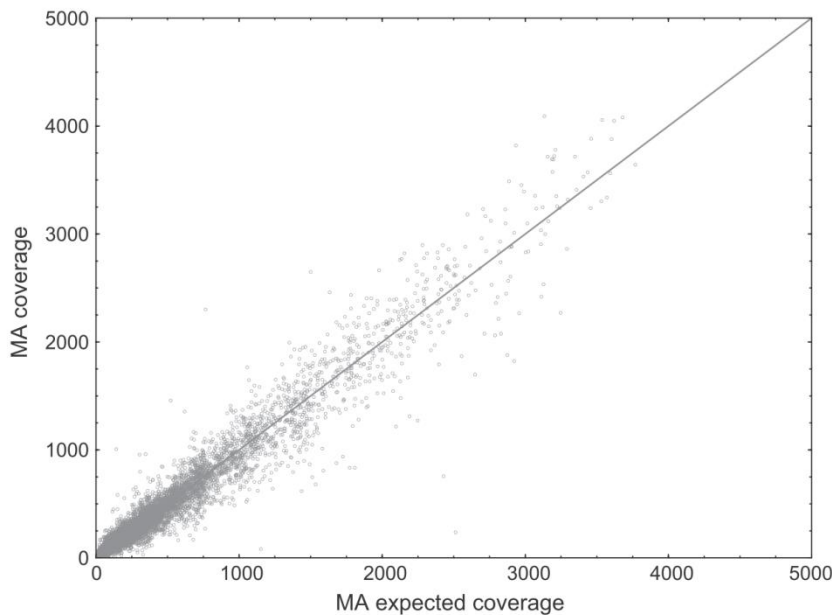


**Fig. 3** Relationship between the observed and expected frequencies of minor alleles in the pool. The observed and expected numbers of bases for minor alleles in the pool are represented for 23 682 high-quality SNPs. SNPs were originally identified during individual genotyping, and the expected numbers of minor allele bases were calculated based on allele frequency and coverage.
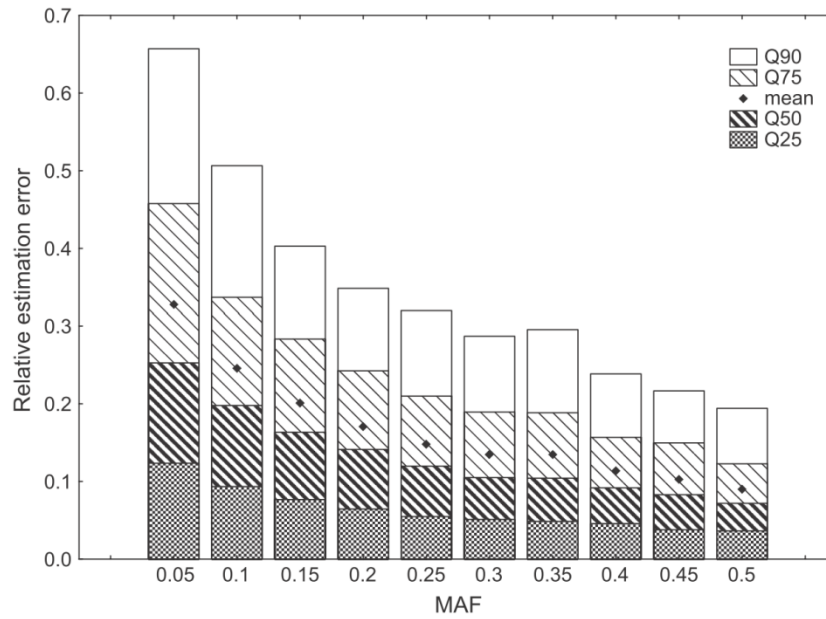
**Fig. 4** Relationship between MAF and allele frequency relative estimation error values for the pooled sample. Columns represent the 25% (Q25), 50% (Q50), 75% (Q75) and 90% (Q90) percentiles for all the relative estimation error values associated with the minor allele frequency classes.
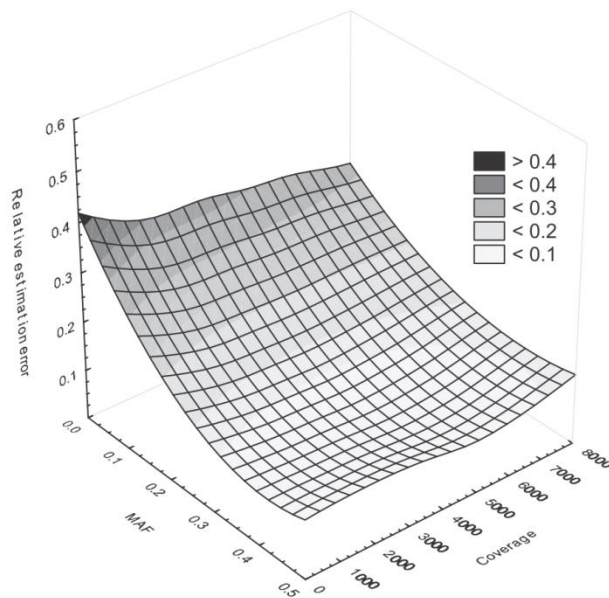


**Fig.5** Allele frequency relative estimation errors for different sequencing coverage and MAF values. The surface contours were obtained using the distance-weighted least squares method for all 23,682 high-quality SNP positions. Relative estimation error was calculated using the expected and observed number of reads of minor frequency alleles in the pooled sample.

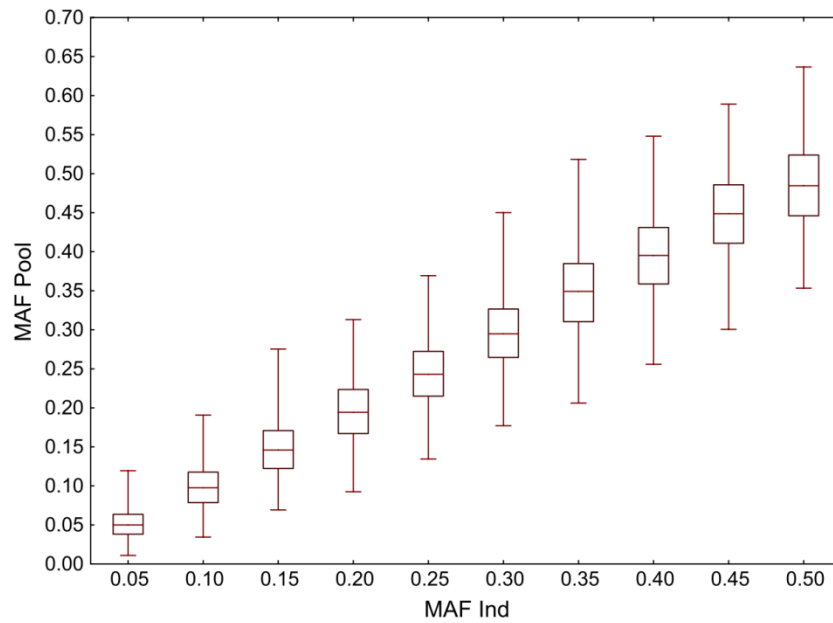**Fig. 6** Allele frequencies estimates from the pool sample within minor allele frequency classes. Boxes indicate 50% of observations, whiskers – 98% of all estimates. Horizontal lines represent medians.
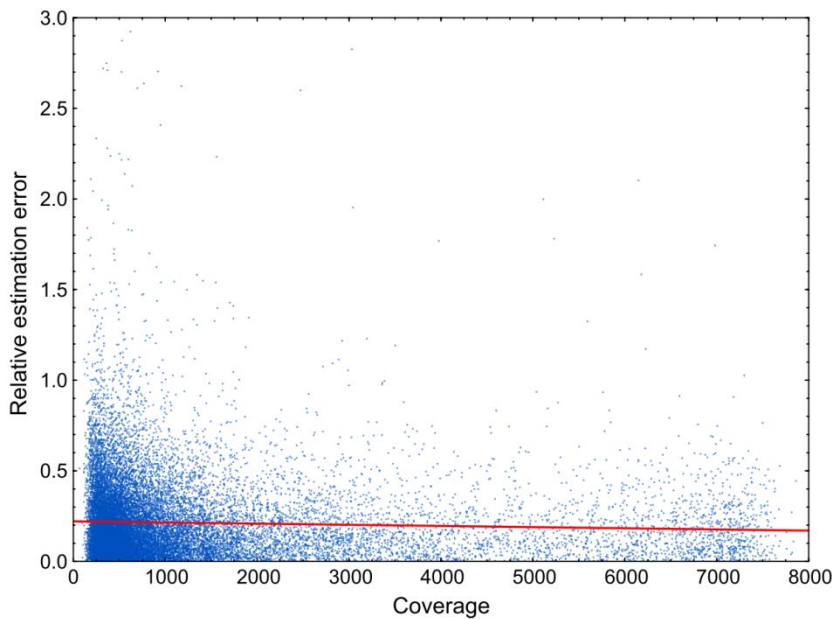


**Fig. 7** Relationship between sequencing coverage and the accuracy of the allele frequency estimates. The correlation plot includes all high-quality SNPs. The regression line is given by equation $y = 0.22 - 6 \cdot 10^{-6}x$.
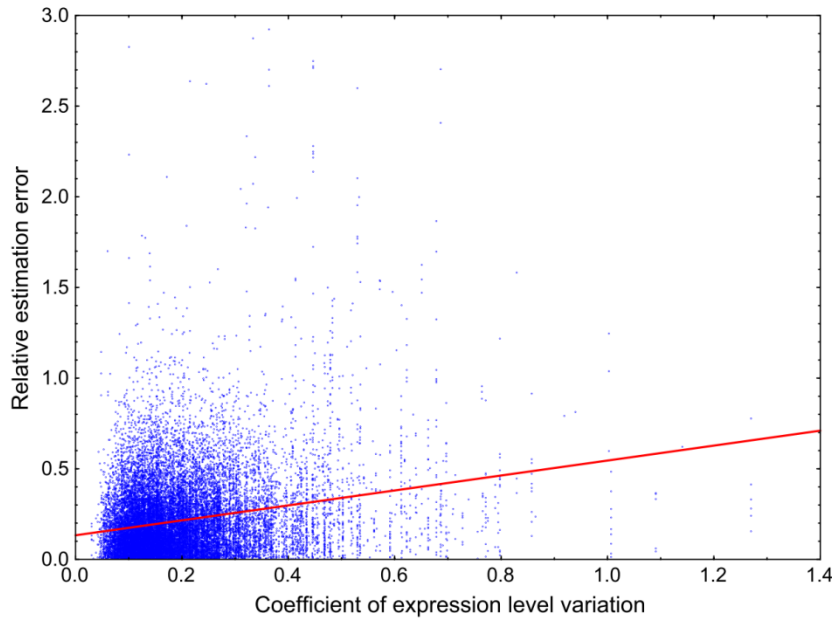
**Fig. 8** Relationship between coefficient of expression level variation and the accuracy of the allele frequency estimates. The correlation plot includes all high-quality SNPs. The regression line is given by equation y = 0.13 + 0.41x.

*Accuracy of gene expression estimation*

In total, 17 861 contigs were analysed to quantify the accuracy of gene expression estimates. Mean relative estimation error was $0.14 \pm 0.001$ SE (median 0.12). We found a significant but very weak negative correlation between mean expression level and relative estimation error ($R^2 = 0.0004$; P = 0.01). The means of expression levels calculated from the individual samples were highly correlated with those calculated from the pooled sample ($R^2 = 0.998$; P < $10^{-5}$).

**Discussion**

We used a nonmodel organism to quantitatively assess the accuracy of allele frequency estimates obtained from pooled RNA samples. Liver RNA samples of ten bank voles were sequenced both separately and as a pool. When we compared the allele frequencies estimated from the pool with the 'true' allele frequencies obtained from the individual samples, we found that the estimates from the pool were generally accurate.

We used only one pooled sample as variability introduced by technical issues should be similar for DNA and RNA pools, and the effect of such variability has been thoroughly explored for DNA pools (Barratt et al. 2002; Zhu et al. 2012). However, RNA pools differ from DNA pools in that the biological variation in the RNA pool is due to inherent differences in expression levels among genes and individuals. As a result, it is more important to examine the accuracy of frequency estimates for SNPs called from multiple genes found in a sample of individuals than to examine that of a few genes across a number of pools.

*SNP calling in the de novo assembled transcriptome*

For  nonmodel organisms, transcriptome assembly is the first, crucial step of RNA-based SNP identification (Singhal  2013). This step is challenging, as divergent alleles may be identified as separate transcripts, sequences of similar paralogues may be lumped together and chimeric transcripts may arise as artefacts of the assembly process. The effectiveness of transcriptome reconstruction has been discussed in several other studies (Bao et al. 2011; Earl et al. 2011; Martin & Wang 2011; De Wit et al. 2012; Singhal 2013), in which different sequencing strategies and assemblers were compared. We should note that, for a comprehensive test of this problem *in silico*, a high-quality reference genome is needed (Vijay et al. 2013). If no reference genome is available, we have to accept an unknown rate of false positives and subsequently test candidate SNPs in future analyses (Singhal  2013). While we recognize that there are problems related to *de novo* transcriptome assembly, we wish to emphasize that the results of our study appear to be robust to the many possible artefacts of transcriptome assembly.

First, we focused on high-quality, high-coverage SNPs that were derived from genes that were at least moderately expressed and had well-assembled transcripts. Second, by discarding SNPs called from contigs that exhibited excessive heterozygosity, we probably filtered out similar paralogues that were represented by a single transcriptome-based gene model (TGM). Heterozygosity at a biallelic locus is not expected to exceed 0.5, and, even then, we are unlikely to observe 9 or 10 heterozygotes of 10 individuals (P = 0.01); therefore, sites with high heterozygosity probably indicate that the contigs represent more than one region in the genome. By removing them from our analyses, we reduced the number of falsely positive SNPs caused by merging paralogues during assembly and reference transcriptome reconstruction. Third, in chimeric transcripts, individual SNPs were most probably properly

called; consequently, the presence of such chimeras, which may constitute a noticeable fraction of transcripts (Edgar et al. 2011), should not systematically bias our results. Taken together, these filtering steps considerably reduced the number of putative SNPs, but the numbers of retained SNPs and genes were still large. This data set provided information on the accuracy of allele frequency estimates for high-quality SNPs varying in sequence coverage and minor allele frequency.

*SNP identification in the sample pool*

Our results suggest that SNP calling from the pool remains challenging for rarer alleles. However, this problem is common to pooling approaches and has been widely discussed in genome resequencing studies focusing on improving the discovery of SNPs with rare variants (Bansal et al. 2010). Several programs dedicated to SNP calling from pools, such as PoPoolation2 (Kofler et al. 2011), vipR (Altmann et al. 2011) or Varscan (Koboldt et al. 2009), are available, but they usually require at least two pooled samples. In most experimental and case–control studies, at least two pools are compared (Sham et al. 2002), and thus, the identification of polymorphic positions and the estimation of allele frequency may be considered somewhat separate tasks. If true sample allele frequencies can be accurately estimated from pools, then existing software used to identify SNPs in DNA pools could potentially be successfully applied to RNA-Seq surveys as well (Thumma et al. 2012). Moreover, in experimental and case–control studies, the aim is to identify SNPs whose allele frequencies differ between groups. At such sites, alternative variants should occur at least an intermediate frequency in one group and thus be easily detected with available software. For example, in our study, 74% of SNPs for which MAF values were greater than 0.25 were called from the pool. SNP discovery is therefore not a limiting factor in the identification of candidate sites from pooled samples because our results support the ability of a pooled approach to identify most of the relevant genetic variation.

*Accuracy of allele frequency estimation from the pooled sample*

Many population genetic analyses require estimates of allele frequencies for comparing different natural populations, experimental treatments or phenotypic classes. Sampling a finite number of individuals from population always introduce stochasticity to these estimates,

which was studied elsewhere (Futschik & Schlotterer 2010; Buerkle & Gompert 2013). Obviously, as more individuals are sequenced from a population, allele frequencies are estimated more precisely and bias is eliminated. In some cases, however, we are not able to sample as many individuals as required (small groups/populations, laboratory colonies of vertebrates, etc.). Estimates of allele frequencies obtained from small samples have wider confidence intervals and are biased, which should be taken into consideration (Gompert & Buerkle 2011). In this study, we estimated the magnitude of additional uncertainty in estimates of allele frequencies introduced by variation in expression level in pooled RNA sample.

We found that estimates of allele frequency obtained from the RNA pool were acceptable for many purposes. The strong correlation between the observed and expected number of nonreference bases demonstrates the utility of pooled RNA samples in wide range of population genetic analyses. The correlation we found was only slightly weaker than that found in a study in which pooled DNA was used (Sham et al. 2002; Ramos et al. 2012). Moreover, almost no bias was present for SNPs with lower expression levels, and estimates of expression level obtained from the pool were accurate even for genes exhibiting moderate expression. However, it is important to note that we focused on genes that were at least moderately expressed by all individuals, and thus, extrapolating our results to genes expressed at very low levels would not be justified.

On the other hand, in our analysis, some gene and SNP categories demonstrate elevated estimation error. We found a negative correlation between MAF and relative estimation error, a result that has been observed for DNA pools as well (Guo et al. 2013). Along with SNP discovery, the low accuracy of allele frequency estimates for rare alleles remains a challenge in analyses of both DNA and RNA pools.

We found evidence that between-individual variation in expression increases estimation error only slightly but significantly: ca. 4% variation in relative estimation error can be explained by variation in expression level between individuals. Allele-specific expression also significantly influences estimates of allele frequency, but ASE seems to occur only in a minor fraction of genes (ca. 1% in our data set according to the applied criteria). These results suggest that inaccuracy in allele frequency estimation may be higher for some classes of genes, and, ideally, such genes should be identified and excluded or analysed separately. Finally, our simple simulations indicate that variation introduced by pooling

systematically increases estimates of population differentiation which may result in some false positives in outliers' analyses.

Using RNA pooling has some additional limitations, namely that a well-assembled reference transcriptome is needed. When using a pooling approach, we do not have access to individual genotypes and thus have no possibility of removing sites with excessive heterozygosity. Therefore, it is worthwhile to invest time and resources in obtaining a high-quality reference transcriptome and sequencing several individually barcoded samples to test and remove the sequences of similar paralogues. These individuals can be used to explore variation in expression level between individuals, and for assessment of ASE. If such resources are available, one can control additional sources of variation in estimating allele frequency and then pooled RNA-Seq is a reliable technique to study nonmodel organisms at the genome- and population-wide scale.

Obviously, pooled approach is not applicable to analyses, which require individual genotypes (e.g. estimating admixture coefficient or estimating linkage disequilibrium among loci). Therefore, clear arguments need to be made for using this approach for molecular ecology studies. Cost effectiveness of large-scale studies is the most obvious such case. Although sequencing itself has become relatively inexpensive, library preparation remains expensive, especially when many samples are processed. With two experimental treatments, 4 replicates within treatment and only ten individuals sampled per treatment, at least 80 libraries need to be prepared. The cost of library preparation for such a modest experiment would be $4800 (NEBNext® Ultra TM Directional RNA Library Prep Kit for Illumina®) or even up to $32,000 (Illumina TruSeq Kit (Stranded Total RNA LT)). This can be reduced ten times if samples within replicates are not barcoded. For studying many populations of nonmodel species pool RNA-seq may reduce laboratory costs drastically. If studied organisms and/or organs are very small and pooling is necessary to obtain enough material for library preparation – pooled RNA-seq is the only viable solution.

Our study tested the accuracy of allele frequency estimates obtained from RNA pools sequenced using Illumina technology. We demonstrated that pooled RNA-Seq approach is a reliable, and cost-effective strategy for obtaining genome-wide information about potentially functionally relevant variation, provided that high-quality transcriptome assembly and stringent SNP-calling and filtering criteria based on sequencing of subset of individuals are used. The lack of such filtering can result in higher inaccuracy for some categories of

transcripts, which may in turn result in a higher rate of false positives in some downstream analyses. When aforementioned prerequisites are fulfilled, the accuracy obtained is very similar to that obtained for DNA pools.

## Data Accessibility

Raw sequences: NCBI BioProject PRJNA222572; reference transcriptome, variant calling files (VCF), files containing numbers of reads and statistics, custom Python scripts: Dryad Digital Repository entry. doi: 10.5061/dryad.bh23t.

## References

Altmann A, Weber P, Quast C, et al. (2011) vipR: Variant identification in pooled DNA using R. Bioinformatics 27, i77-i84.

Babik W, Stuglik M, Qi W, et al. (2010) Heart transcriptome of the bank vole (Myodes glareolus): Towards understanding the evolutionary variation in metabolic rate. BMC Genomics 11, 390.

Bansal V (2010) A statistical method for the detection of variants from next-generation resequencing of DNA pools. Bioinformatics 26, i318-i324.

Bansal V, Libiger O, Torkamani A, Schork NJ (2010) Statistical analysis strategies for association studies involving rare variants. Nature Reviews Genetics 11, 773-785.

Bao S, Jiang R, Kwan W, et al. (2011) Evaluation of next-generation sequencing software in mapping and assembly. Journal of Human Genetics 56, 406-414.

Barbosa-Morais NL, Irimia M, Pan Q, et al. (2012) The evolutionary landscape of alternative splicing in vertebrate species. Science 338, 1587-1593.

Barratt BJ, Payne F, Rance HE, et al. (2002) Identification of the sources of error in allele frequency estimations from pooled DNA indicates an optimal experimental design. Annals of Human Genetics 66, 393-405.

Boratyński Z, Koteja P (2009) The association between body mass, metabolic rates and survival of bank voles. Functional Ecology 23, 330-339.

Breen MS, Kemena C, Vlasov PK, Notredame C, Kondrashov FA (2012) Epistasis as the primary factor in molecular evolution. Nature 490, 535-538.

Brenchley R, Spannagl M, Pfeifer M, et al. (2012) Analysis of the bread wheat genome using whole-genome shotgun sequencing. Nature 491, 705-710.

Buerkle CA, Gompert Z (2013) Population genomics based on low coverage sequencing: How low should we go? Molecular Ecology 22, 3028-3035.

Charlesworth B, Charlesworth D (2010) Elements of Evolutionary Genetics Roberts & Company Publishers.

Cox MP, Peterson DA, Biggs PJ (2010) SolexaQA: At-a-glance quality assessment of Illumina second-generation sequencing data. BMC Bioinformatics 11, 485.

De Wit P, Pespeni MH, Ladner JT, et al. (2012) The simple fool's guide to population genomics via RNA-Seq: An introduction to high-throughput sequencing data analysis. Molecular Ecology Resources 12, 1058-1067.

Druley TE, Vallania FLM, Wegner DJ, et al. (2009) Quantification of rare allelic variants from pooled genomic DNA. Nature Methods 6, 263-265.

Earl D, Bradnam K, St. John J, et al. (2011) Assemblathon 1: A competitive assessment of de novo short read assembly methods. Genome Research 21, 2224-2241.

Edgar RC, Haas BJ, Clemente JC, Quince C, Knight R (2011) UCHIME improves sensitivity and speed of chimera detection. Bioinformatics 27, 2194-2200.

Ekblom R, Galindo J (2011) Applications of next generation sequencing in molecular ecology of non-model organisms. Heredity 107, 1-15.

Feder AF, Petrov DA, Bergland AO (2012) LDx: Estimation of Linkage Disequilibrium from High-Throughput Pooled Resequencing Data. PLoS ONE 7, e48588.

Futschik A, Schlötterer C (2010) The next generation of molecular markers from massively parallel sequencing of pooled DNA samples. Genetics 186, 207-218.

Ge B, Pokholok DK, Kwan T, et al. (2009) Global patterns of cis variation in human cells revealed by high-density allelic expression analysis. Nature Genetics 41, 1216-1222.

Gilad Y, Pritchard JK, Thornton K (2009) Characterizing natural variation using next-generation sequencing technologies. Trends in Genetics 25, 463-471.

Gompert Z, Buerkle CA (2011) A hierarchical bayesian model for next-generation population genomics. Genetics 187, 903-917.

Grabherr MG, Haas BJ, Yassour M, et al. (2011) Full-length transcriptome assembly from RNA-Seq data without a reference genome. Nature Biotechnology 29, 644-652.

Hartl D, Clark A (2006) Principles of Population Genetics, Fourth Edition Sinauer Associates, Inc.

Huang X, Madan A (1999) CAP3: A DNA sequence assembly program. Genome Research 9, 868-877.

Jeukens J, Renaut S, St-Cyr J, Nolte AW, Bernatchez L (2010) The transcriptomics of sympatric dwarf and normal lake whitefish (Coregonus clupeaformis spp., Salmonidae) divergence as revealed by next-generation sequencing. Molecular Ecology 19, 5389-5403.

Jones FC, Grabherr MG, Chan YF, et al. (2012) The genomic basis of adaptive evolution in threespine sticklebacks. Nature 484, 55-61.

Kim SY, Li Y, Guo Y, et al. (2010) Design of association studies with pooled or un-pooled next-generation sequencing data. Genetic Epidemiology 34, 479-491.

Koboldt DC, Chen K, Wylie T, et al. (2009) VarScan: Variant detection in massively parallel sequencing of individual and pooled samples. Bioinformatics 25, 2283-2285.

Kofler R, Pandey RV, Schlötterer C (2011) PoPoolation2: Identifying differentiation between populations using sequencing of pooled DNA samples (Pool-Seq). Bioinformatics 27, 3435-3436.

Langmead B, Salzberg SL (2012) Fast gapped-read alignment with Bowtie 2. Nature Methods 9, 357-359.

Li B, Dewey CN (2011) RSEM: Accurate transcript quantification from RNA-Seq data with or without a reference genome. BMC Bioinformatics 12.

Li H (2011) A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. Bioinformatics 27, 2987-2993.

Li H, Handsaker B, Wysoker A, et al. (2009) The Sequence Alignment/Map format and SAMtools. Bioinformatics 25, 2078-2079.

Li W, Godzik A (2006) Cd-hit: A fast program for clustering and comparing large sets of protein or nucleotide sequences. Bioinformatics 22, 1658-1659.

Lu T, Lu G, Fan D, et al. (2010) Function annotation of the rice transcriptome at single-nucleotide resolution by RNA-seq. Genome Research 20, 1238-1249.

Lynch VJ, Wagner GP (2008) Resurrecting the role of transcription factor change in developmental evolution. Evolution 62, 2131-2154.

Martin JA, Wang Z (2011) Next-generation transcriptome assembly. Nature Reviews Genetics 12, 671-682.

Mokkonen M, Kokko H, Koskela E, et al. (2011) Negative frequency-dependent selection of sexually antagonistic alleles in Myodes glareolus. Science 334, 972-974.

Morloy M, Molony CM, Weber TM, et al. (2004) Genetic analysis of genome-wide variation in human gene expression. Nature 430, 743-747.

Mortazavi A, Williams BA, McCue K, Schaeffer L, Wold B (2008) Mapping and quantifying mammalian transcriptomes by RNA-Seq. Nature Methods 5, 621-628.

Oleksiak MF, Roach JL, Crawford DL (2005) Natural variation in cardiac metabolism and gene expression in Fundulus heteroclitus. Nature Genetics 37, 67-72.

Radwan J, Babik W (2012) The genomics of adaptation. Proceedings of the Royal Society B: Biological Sciences 279, 5024-5028.

Radwan J, Tkacz A, Kloch A (2008) MHC and preferences for male odour in the bank vole. Ethology 114, 827-833.

Ramos E, Levinson BT, Chasnoff S, et al. (2012) Population-based rare variant detection via pooled exome or custom hybridization capture with or without individual indexing. BMC Genomics 13, 683.

Reich D, Green RE, Kircher M, et al. (2010) Genetic history of an archaic hominin group from Denisova cave in Siberia. Nature 468, 1053-1060.

Rice AM, Rudh A, Ellegren H, Qvarnström A (2011) A guide to the genomics of ecological speciation in natural animal populations. Ecology Letters 14, 9-18.

Robinson MD, Oshlack A (2010) A scaling normalization method for differential expression analysis of RNA-seq data. Genome Biology 11, r25.

Sadowska ET, Baliga-Klimczyk K, Chrzascik KM, Koteja P (2008) Laboratory model of adaptive radiation: A selection experiment in the bank vole. Physiological and Biochemical Zoology 81, 627-640.

Sadowska ET, Labocha MK, Baliga K, et al. (2005) Genetic correlations between basal and maximum metabolic rates in a wild rodent: Consequences for evolution of endothermy. Evolution 59, 672-681.

Salem M, Vallejo RL, Leeds TD, et al. (2012) RNA-seq identifies SNP markers for growth traits in rainbow trout. PLoS ONE 7, e36264.

Sandberg R, Yasuda R, Pankratz DG, et al. (2000) Regional and strain-specific gene expression mapping in the adult mouse brain. Proceedings of the National Academy of Sciences of the United States of America 97, 11038-11043.

Serre D, Gurd S, Ge B, et al. (2008) Differential allelic expression in the human genome: A robust approach to identify genetic and epigenetic Cis-acting mechanisms regulating gene expression. PLoS Genetics 4, 2.

Sham P, Bader JS, Craig I, O'Donovan M, Owen M (2002) DNA pooling: A tool for large-scale association studies. Nature Reviews Genetics 3, 862-871.

Singhal S (2013) De novo transcriptomic analyses for non-model organisms: An evaluation of methods across a multi-species data set. Molecular Ecology Resources 13, 403-416.

Skelly DA, Johansson M, Madeoy J, Wakefield J, Akey JM (2011) A powerful and flexible statistical framework for testing hypotheses of allele-specific gene expression from RNA-seq data. Genome Research 21, 1728-1737.

Thumma BR, Sharma N, Southerton SG (2012) Transcriptome sequencing of Eucalyptus camaldulensis seedlings subjected to water stress reveals functional single nucleotide polymorphisms and genes under selection. BMC Genomics 13, 364.

Tschirren B, Andersson M, Scherman K, Westerdahl H, Råberg L (2012) Contrasting patterns of diversity and population differentiation at the innate immunity gene toll-like receptor 2 (tlr2) in two sympatric rodent species. Evolution 66, 720-731.

Vera JC, Wheat CW, Fescemyer HW, et al. (2008) Rapid transcriptome characterization for a nonmodel organism using 454 pyrosequencing. Molecular Ecology 17, 1636-1647.

Vijay N, Poelstra JW, Künstner A, Wolf JBW (2013) Challenges and strategies in transcriptome assembly and differential gene expression quantification. A comprehensive in silico assessment of RNA-seq experiments. Molecular Ecology 22, 620-634.

Wolf JBW, Bayer T, Haubold B, et al. (2010) Nucleotide divergence vs. gene expression differentiation: Comparative transcriptome sequencing in natural isolates from the carrion crow and its hybrid zone with the hooded crow. Molecular Ecology 19, 162-175.

Zhu Y, Bergland AO, González J, Petrov DA (2012) Empirical validation of pooled whole genome population re-sequencing in Drosophila melanogaster. PLoS ONE 7, e41901.

# CHAPTER II

# Initial molecular-level response to artificial selection for increased aerobic metabolism occurs primarily through changes in gene expression

M. Konczal, W. Babik, J. Radwan, E.T. Sadowska, P. Koteja

**Abstract**

Experimental evolution combined with genome or transcriptome resequencing (Evolve and Resequence) represents a promising approach for advancing our understanding of the genetic basis of adaptation. Here, we applied this strategy to investigate the effect of selection on a complex trait in lines derived from a natural population of a small mammal. We analyzed the liver and heart transcriptomes of bank voles (*Myodes (=Clethrionomys) glareolus*) that had been selected for increased aerobic metabolism. The organs were sampled from 13th generation voles; at that point, the voles from four replicate selected lines had 48% higher maximum rates of oxygen consumption than those from four control lines. At the molecular level, the response to selection was primarily observed in gene expression: over 300 genes were found to be differentially expressed between the selected and control lines and the transcriptome-wide pattern of expression distinguished selected lines from controls. No evidence for selection-driven changes of allele frequencies at coding sites was found: no SNP changed frequency more than expected under drift alone and frequency changes aggregated over all SNPs did not separate selected and control lines. Nevertheless, among genes which showed highest differentiation in allele frequencies between selected and control lines we identified, using information about gene functions and the biology of the selected phenotype, plausible targets of selection; these genes, together with those identified in expression analysis have been prioritized for further studies. Because our selection lines were derived from a natural population, the amount and the spectrum of variation available for selection probably closely approximated that typically found in populations of small mammals. Therefore our results are relevant to the understanding of the molecular basis of complex adaptations occurring in natural vertebrate populations.

**Introduction**

One of the central goals of evolutionary biology is to understand the genetic mechanisms by which organisms evolve new, adaptive phenotypes under natural selection and thus diverge phenotypically (Stapley et al. 2010; Butlin et al. 2012). Despite decades of research, detecting and deciphering the molecular changes underlying adaptation remain challenging tasks to which researchers have applied various approaches, such as study of candidate genes, genome-wide scans for positive selection or experimental evolution (Sabeti et al. 2007; Garland and Rose 2009; Stapley et al. 2010; Barrett and Hoekstra 2011; Fournier-Level et al. 2011). Recently, however, the Evolve and Resequence (E&R) approach has been gaining popularity. E&R studies provide better control over confounding factors than other approaches and allow investigators to choose the traits under selection (Turner et al. 2011; Kawecki et al. 2012). This approach involves genetic analyses of populations of organisms that are either adapted to specific, experimentally controlled ambient conditions, or that are selected for increased performance with respect to a specific behavioral, life-history, or morpho-physiological trait. Such studies have helped to answer questions concerning adaptation (Tenaillon et al. 2012; Soria-Carrasco et al. 2014), the importance of new mutations (Burke et al. 2010) and the genomic patterns of a recent response to selection (Johansson et al. 2010; Pettersson et al. 2013).

The sources of adaptive variation appear to vary among evolutionary lineages. For example, extensive work on microorganisms has contributed to our understanding of adaptation scenarios that are driven by selection acting on new mutations: a substantial number of de novo mutations are expected during the course of experiments in such organisms as a result of the large population sizes involved (Herring et al. 2006; Barrick et al. 2009; Tenaillon et al. 2012). However, in multicellular, sexually reproducing species (the subject of the present study), standing genetic variation is the main source of variation at the initial stages of adaptive evolution (Barrett and Schluter 2008; Burke et al. 2010). The E&R approach has been used to comprehensively and successfully investigate some traits in D. melanogaster (Teotónio et al. 2009; Burke et al. 2010; Turner et al. 2011); for example, Burke and colleagues (2010) studied flies that had been selected for accelerated development for 600 generations. They concluded that the probability of fixation of selected variants is relatively low and that selection does not readily expunge genetic variation (Burke et al. 2010). Subsequent studies using fruit flies have confirmed the complex evolutionary trajectories of

selected variants and emphasized the importance of epistatic interactions (Huang et al. 2012; Orozco-Terwengel et al. 2012).

The genetic response to selection at the early stages of adaptation is less well understood in vertebrates, which usually have smaller population sizes (Johansson et al. 2010; Chan et al. 2012; Pettersson et al. 2013). Interestingly, however, the results from the few experiments performed thus far contrast with those obtained from D. melanogaster. For instance, in lines of chickens that had been selected for high and low body mass, Johansson and colleagues (2010) observed many genetic regions with fixed differences; likewise, signals of classical, hard selective sweeps were detected in a mouse line that had been selected for high body mass (Chan et al. 2012). In these studies, the high number of regions detected that were presumably under the influence of divergent selection (50 in chickens and 67 in mice) suggests that the initial phase of selection substantially increases divergence between lines while simultaneously reducing polymorphism within them. However, these results may reflect the differences in experimental setup rather than true contrast between vertebrates and Drosophila. Both vertebrate experiments utilized crossed inbred lines as a base population. In such cases long haplotype blocks are present at the beginning of the experiment and they may be fixed rapidly in small experimental populations, mimicking the effects of hard sweeps. Such situations are less likely in nature, where at the early stage of adaptation standing genetic variation is subject to selection. To understand the basis of adaptive processes occurring in the wild, it is therefore crucial to conduct selection experiments that control for the effect of genetic drift and utilize lines derived from natural populations. Although selection experiments on non-model organisms were always possible to perform, in practice they were rarely undertaken, partly because until recently uncovering molecular genetic mechanisms of the evolution in non-model organisms was often not possible. This has changed with the advent of high throughput sequencing (Schlötterer et al. 2014).

In the present study, we used high-throughput transcriptome sequencing to test whether recent, intense selection acting over multiple generations in mammalian populations would result in repeatable changes in the frequencies of variants in protein-coding genes and/or patterns of gene expression. This study was performed using an experimental evolution model system, with four lines of bank voles (*Myodes (=Clethrionomys) glareolus*), selectively bred for high swim-induced aerobic metabolism (A lines) and four unselected control lines (C lines; Sadowska et al. 2008). The experiment has been established as a tool for testing hypotheses concerning correlated evolution of aerobic locomotor performance and basal

metabolic rates, which is believed to have been a crucial element in evolution of terrestrial vertebrates (literature cited in Sadowska et al. 2005, 2008). Thus, the model is likely to illuminate many eco-physiological questions concerning physiological genomics and the evolution of endothermy (Nespolo et al. 2011; Pérusse et al. 2013). The swim-induced maximum rate of oxygen consumption differed significantly between the selected and control lines already in generation 2 (Sadowska et al. 2008), and in generation 13 it was 48% higher in A line voles than in C line voles (mean ± SD: 5.32 ± 0.64 ml O2/min vs. 3.59 ± 0.57 ml O2/min, respectively; Chrząścik et al. 2014, Stawski et al. 2015; see also Supplementary materials 1.3). Voles from the A lines (also referred to as "selected" lines) differed significantly from control voles not only in the trait directly under selection, but also in their basal metabolic rate and a number of other behavioral and morpho-physiological traits (Supplementary materials 1.3). This experiment presented a unique opportunity to study the genetic basis of the response to selection in mammals thanks to a combination of several factors: i) selection could operate on the natural genetic variation directly derived from a wild population, ii) known pedigrees allowed for the exact calculation of drift expectations, iii) the trait under selection was complex and ecologically important, and iv) the replicated lines provided an appropriate system to study the role of drift in phenotypic and genetic differentiation.

The eight lines (four selected and four control) were sequenced using a pooled RNA-Seq approach (Konczal et al. 2014). We used transcriptome analysis as a convenient way to determine whether the response to selection at the molecular level was dominated by gene expression or structural changes. King and Wilson (1975) proposed that adaptive evolutionary change is largely due to changes in gene expression, and there is empirical evidence from genetic mapping and interspecies comparisons that both supports (Wray 2007; Jones et al. 2012) and contradicts this view (Hoekstra and Coyne 2007). A recent study of patterns of polymorphism and divergence in murid rodents suggested that most of adaptive changes appear in regulatory regions. On the other hand, wider regions of reduced diversity around exons than around conserved noncoding elements may be interpreted as a result of substantially larger effects of adaptive substitutions  (Halligan et al. 2013). However, it is unclear whether rapid adaptation from standing genetic variation produces similar patterns. To address the question of the relative importance of coding mutations versus changes in expression levels, we studied the transcriptomes of two organs: the heart, which plays a crucial role in an organism's aerobic capacity (Bye et al. 2008), and the liver, which, as a

central metabolic organ (Malarkey et al. 2005, Konczal et al. 2014), was a promising target for investigations of the molecular mechanisms that were responsible for the increased basal metabolism observed in selected lines. The scale of the project limited the possibility of detecting significant responses to selection in allele frequencies in coding regions or gene expression only to loci of large effects. However, we could still infer the role of many loci of small effect if selection changes allele frequencies in coding regions or gene expression of many genes in replicable way (across the four selected and four control lines). In such case, the aggregate effect of these changes should result in multi-dimensional differentiation of selected lines from controls (Turchin et al. 2012), although covariances of allele frequencies, resulting from between-population component of LD (Linkage Disequilibrium; Ohta 1982), may weaken this effect (Storz 2005, Le Corre and Kremer 2012).

We identified over 300 differentially expressed genes that are associated with diverse molecular functions; many of these functions appeared to be highly relevant to the phenotypic response to selection for increased aerobic metabolism. This result, combined with significant clustering of genome-wide transcriptional profiles, highlights the role of rapid changes in gene expression at the early stages of adaptive evolution. In contrast, allele frequency changes in coding sequences appear to play, at best, a minor role: the differences observed in the allele frequencies between the selected and control lines could be entirely explained by drift and the aggregate effect of allele frequency changes does not separate selected lines from controls. Nevertheless, among the genes that showed the highest differentiation in allele frequencies, we identified, on the basis of their molecular function, a set of candidates, which may possibly contribute to phenotypic changes between the selected and control lines. These genes should be prioritized as a target for future research.

**Results**

*Single nucleotide polymorphisms*

From each sample, an average of 37.1 (± 10.6 SD) million 1 x 100 bp reads were obtained; of these, 75.9% were uniquely mapped to the bank vole liver and heart reference transcriptomes (Tab. 1). After several steps of data filtering (see Methods), we identified 172,246 SNPs. The vast majority of identified variants were found in putative protein-coding genes, with an

average of 3.95 and 3.48 SNPs per kb in open reading frames (ORFs) and untranslated regions (UTRs) of SNP-containing contigs (Tab. 2).

**Table 1.** Overview of the Assembly of Bank Vole Transcriptomes

|  | Liver | Heart |
| --- | --- | --- |
| No. of genes | 146,758 | 252,281 |
| No. of genes >1kb | 23,512 | 24,825 |
| N50 gene length (bp) | 1,225 | 650 |
| No. of genes within N50 | 19,101 | 47,439 |
| No. of genes with likely CDS | 18,050 | 11,110 |
| N50 of genes with likely CDS | 3,296 | 3,081 |
| No. of bases (Mb) | 103.1 | 134.9 |

Note. – N50, 50% of the assembly length is in genes of the length of N50 bp or longer; genes, TGMs contain both coding and noncoding sequences; genes with likely CDS, genes containing successfully annotated ORFs.

**Table 2.** Overview of SNPs Used for Analyses

| | |
| --- | --- |
| No. of SNPs | 172,246 |
| No. of genes with SNPs | 15,043 |
| No. of nonsynonymous SNPs | 22,963 |
| No. of synonymous SNPs | 44,844 |
| No. of UTR-located SNPs | 71,657 |
| No. of SNPs in noncoding genes | 32,782 |

To estimate effective population sizes, the mean inbreeding coefficient (F) was calculated from pedigree for each of the four selected (A) and four control (C) lines in each generation. The degree of inbreeding increased slightly faster in the selected lines, probably reflecting a subtle difference in the breeding scheme between the selected and control lines (see Methods). The mean effective population size (Ne) was about 16.4% lower in the selected than in the control lines (56.1 vs. 67.1; $p = 0.06$, t-test; Fig. 1A). To evaluate the effect of differences in Ne between lines on the amount of genetic variation we examined the allele frequency spectra (Fig. 1B). Specifically we calculated for each line the number of such

SNPs which were polymorphic in the entire dataset but showed little or no variation (minor allele frequency, MAF <0.05) within the line. An ANCOVA was used to examine how well Ne (covariate) and treatment (selected vs control lines) explained the number of such SNPs. We found a significant effect of Ne ($F(1,5) = 6.92$, p=0.047), but no effect of treatment ($F(1, 5) = 0.14$, p = 0.72; Fig. 1D).

For each SNP, $F_{ST}$ values were calculated between all pairs of lines. Mean pairwise FST distances did not reveal any clustering of selected or control lines ($F(1, 6) = 0.97$, p = 1.00, randomization test; Fig. 1C); and variation among selected lines (calculated as a mean distance to centroids) was slightly, but non-significantly higher than that between control lines ($F(1, 6) = 1.17$, p = 0.32; ANOVA). The two control lines (C1, C3) with the largest effective population sizes were least distant from each other, suggesting that drift played the dominant role in the differentiation of allele frequencies among lines.

Additionally, a principal components analysis (PCA) was performed to look for correlated changes in allele frequencies in various subsets of SNPs; such changes could reflect the response of multiple genes to the same selection pressure. None of the eight PCs clearly differentiated between selected and control lines (Fig. S2.1).

In the next step, folded allele frequency spectra were compared both between selection regimes and with expectations generated from simulations of genetic drift over the course of the experiment. Forward simulations were performed using known pedigrees; for the initial allele frequency spectrum, these simulations used the average spectrum calculated from control lines.

The allele frequency spectra were less skewed in the simulated data than in the observed data (Fig. 1B), which could have been caused by two effects: bias in the estimation of the initial allele frequency spectrum or selection against slightly deleterious alleles. We assessed the overall effect of deleterious alleles by comparing the allele frequency spectra of synonymous and nonsynonymous sites. Minor allele frequencies were lower for nonsynonymous SNPs than for synonymous SNPs (synonymous median MAF = 0.091, nonsynonymous median MAF = 0.068; p = 10-16; KS [Kolmogorov Smirnov] tests), indicating the presence of purifying selection. For synonymous sites, the difference in the percentage of rare variants between simulated and observed sites was 2.7% ; in contrast, the difference was 7.4% for nonsynonymous sites.
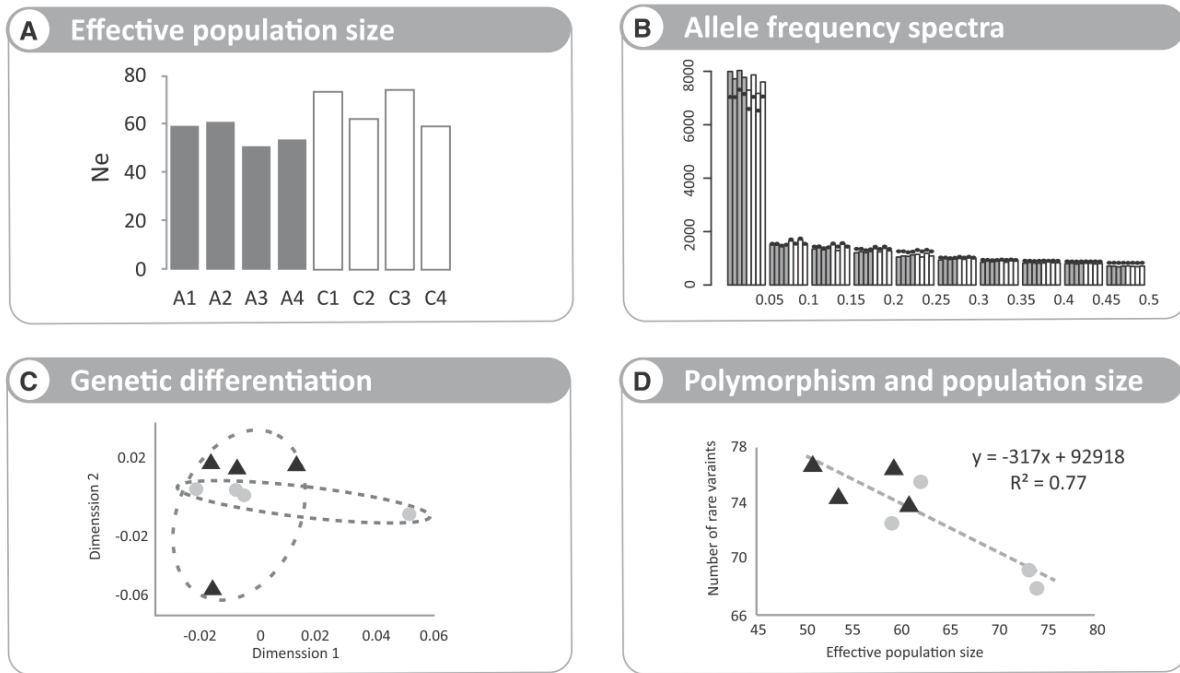
**Fig. 1** Effect of selection and population size on allele frequency changes in the bank vole selection experiment. A. Effective population sizes (Ne) of selected (grey) and control (white) lines, calculated from pedigrees for 13 generations of the selection experiment. B. Folded allele frequency spectra for selected (grey) and control (white) lines, compared with pedigree-based simulations (dots – expectations from simulations ). C. Multidimensional scaling plot (MDS) of genetic distances (pairwise FST) between selected and control lines. Triangles represent selected lines; circles represent controls. The MDS plot drawn using mean pairwise FST values calculated for all SNPs. D. Regression of the number of SNPs with rare variants (minor allele frequency < 0.05) on effective population size. Number of rare variants in thousands.

If the same alleles contribute to the response to selection in all lines, SNP frequencies in the selected lines should diverge from those in the control lines to a greater degree than expected under neutrality. To investigate whether such effect occurred, we identified variants that had ranges of allele frequencies non-overlapping between the selected and the control lines (3233 (1.88%) SNPs in 1873 genes). The number of SNPs with non-overlapping allele frequencies was significantly lower than expected from drift simulations (p = 0.01; randomization test), but this effect was not significant for subsets of synonymous (1.97%, p=0.33, randomization test) or nonsynonymous SNPs (1.96%, p=0.78, randomization test). For each of these sites, the minimum allele frequency difference between the sixteen possible A-C comparisons (diffStat) was used as a composite statistic (Turner et al. 2011). The distribution of diffStat values did not differ between the data and drift simulations, and we did not observe overrepresentation of high diffStat values (Fig. S2.2).

The relatively small population sizes decrease the population recombination rates, which may cause entire long haplotypes to drift. To control for the effect of linkage within genes we used the following procedure. First, we generated 1000 datasets consisting of SNPs sampled randomly one per gene. Then, for each dataset the number of SNPs with the ranges of allele frequencies non-overlapping between selection and control lines was calculated. Finally, we recorded the proportion of datasets in which the number of SNPs with non-overlapping allele frequencies was higher than expected under drift (upper 10% of the distribution from simulations). None of the datasets fell into the upper 10% of the distribution, and the relative number of differentiated SNPs was slightly lower than expected from simulations ($p < 10^{-50}$; t test). In coding regions  this effect was mostly explained by nonsynonymous sites ($p < 10^{-50}$, t test), while the fraction of synonymous SNPs with non-overlapping allele frequencies closely followed drift expectations ($p=0.15$; t test; Fig. S2.3).

To get some insight about the power to detect variants under selection in our experiment, we performed pedigree-based simulations of selection. These simulations were used to estimate the probability of obtaining non-overlapping allele frequencies between the selected and control lines, depending on the strength of selection and initial allele frequency. With increasing selective advantage the probability of obtaining non-overlapping frequencies increased considerably ($s = 0.05 – 5.5\%$; $s = 0.2 – 41.6\%$, averaged over the range of initial allele frequencies Fig. S2.4). The probability of obtaining non-overlapping frequencies after 13 generations was highest when the favored allele initially segregated at an intermediate frequency  (initial frequency $0.05 – 5.8\%$; $0.5 – 29.9\%$; $0.9 – 4.0\%$, averaged over the range of selection coefficients). This probably reflects the fact, that rare positively selected variants will often be lost due to drift in some of selected lines. Similarly, selected variants at high initial frequencies will often become fixed in at least some control lines.

Genes that harbored differentiated SNPs (diffStat > 0) had a higher density of polymorphisms ($p < 10^{-6}$, randomization test) which in turn exhibited more equal allele frequencies ($p = 10^{-12}$; KS test). This effect was present for nonsynonymous SNPs ($p = 1.6 \times 10^{-5}$, randomization test), but not for synonymous ($p=0.38$, randomization test), what  may mean either that highly polymorphic genes are more likely to be targets of selection, or that they are more likely to differentiate by drift because of their effective neutrality.

To explore whether some of the genes with differentiated nonsynonymous SNPs (Tab. S1) were somehow associated with phenotypic differences between selected and control

lines, we investigated their functions using relevant databases, and the most intriguing cases are described in Discussion.

Overall, these results did not provide evidence that selection for increased maximum metabolic rate caused allele frequency changes at coding SNPs. The changes in allele frequencies that we did observe can be explained by the actions of two other evolutionary forces, namely drift and purifying selection, that acted in the same way in all lines.

*Gene expression*

To determine differences in expression levels between the selected and control lines, we investigated all expressed genes with at least 10 mapped reads and performed ordination of the lines using a multidimensional scaling analysis that was based on estimates of pairwise similarity in expression levels. In contrast with the SNP results, this analysis found that the selected lines and control lines clustered separately for the liver samples; for heart samples clustering was not significant (liver: p=0.002; heart: p=0.384; randomization tests; Fig. 2A, 2B). Thus, it appears that similar changes in gene expression in the most important metabolic organ, the liver, might have occurred in all selected lines, distinguishing them from controls.

In the heart samples, 79 genes were differentially expressed between selected and control lines (52 genes were overexpressed and 20 were underexpressed in selected lines; false discovery rate (FDR) = 0.05; Fig. 2C). Many more genes were differentially expressed in the liver (278 genes at FDR = 0.05; 123 genes were overexpressed and 155 were underexpressed in selected lines; Fig. 2D). We annotated 110 differentially expressed genes (28 in heart and 82 in liver), all putatively protein coding (Tab. S2, S3). As an additional assessment of these differentially expressed genes we performed t-test on FPKM (Fragments Per Kilobase of transcript per Million mapped reads) values and calculated the proportion of genes with FPKM values non-overlapping between the selected and control lines. 63% of coding genes differentially expressed in liver and 61% in heart showed statistically significant result of the t-test (p<0.05) and 55% and 63% of them respectively had non-overlapping expression level values. The molecular functions of some of these 110 genes are described below.
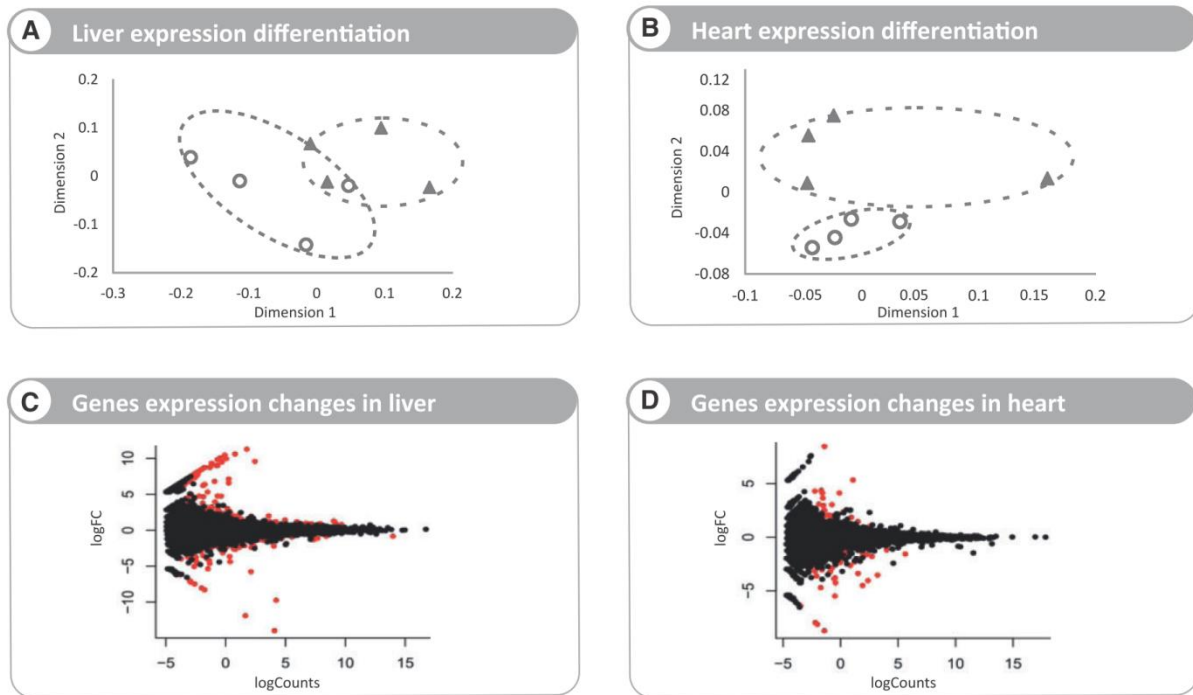
**Fig. 2** Effect of selection on expression changes in the bank vole selection experiment. A. B. - Multidimensional scaling plot (MDS) of transcriptome distances (in terms of the biological coefficient of variation (BCV)) between selected and control lines of the bank vole experiment. Triangles represent selected lines; circles represent controls. The MDS plots were drawn using the expression level values of 91,760 liver and 108,656 heart genes. C. D. - Log-fold-change expression versus log abundance of gene expression in liver and heart samples. Gene expression data are TMM-normalized. Genes that qualified as significantly differentially expressed (FDR 0.05) are in red.

## Discussion

*Differentiation at the molecular level*

Using replicate selected and control lines derived from a natural population of a small mammal, we experimentally quantified the responses to artificial selection at the molecular level. The trait investigated here, maximum metabolic rate during exercise, likely has a complex genetic basis (Hagberg et al. 2011; Roth et al. 2012; Pérusse et al. 2013; Wolfarth et al. 2014) and since the genomic basis of evolutionary change in complex traits is still poorly understood (Rockman, 2012), our results are of interest from a broad evolutionary perspective. This experiment thus addresses a question of general evolutionary and physiological relevance by applying the strict criteria associated with the design of E&R studies (replicates and use of control populations). In doing so, it has provided insight into the genetic patterns of adaptation that arise from standing genetic variation in populations of

mammals. Here, we clearly observed that the effects of artificial selection were visible at both the phenotypic and molecular levels. However, while the artificial selection applied in this experiment resulted in reproducible changes in the expression levels of many genes, it did not cause appreciable changes in allele frequencies at coding SNPs, which were instead influenced predominately by drift.

The importance and contribution of expression changes and coding mutations to adaptation has long been a topic of great interest (King and Wilson 1975; Hoekstra and Coyne 2007; Wray 2007; Stern and Orgogozo 2008; Fraser 2013). Recent findings in human populations suggest that adaptation in regulatory elements, likely affecting gene expression, is 10 times more frequent than in protein-coding parts of the genome (Fraser 2013). Similar evidence has been obtained from diverse taxa: since the split between marsupials and placental mammals, many more new regulatory elements than coding exons have emerged to differentiate the two groups (Mikkelsen et al. 2007), and in the evolution of rodents, most adaptive mutations have occurred in regulatory elements rather than in protein-coding exons (Halligan et al. 2013). Our results are consistent with these findings in showing that, during a relatively short period of selection in small populations, the pattern of expression of multiple genes can change rapidly and in a reproducible manner. Our two main observations—that more than 300 genes were differentially expressed between selected and control lines and that changes in allele frequencies were caused predominantly by drift—support the hypothesis that changes in gene expression, rather than changes in allele frequencies of coding regions, play a central role in adaptation. Other genetic analyses of rodent selection experiments found hundreds of differentially expressed genes by either eQTL investigations (Kelly et al. 2012, Kelly et al. 2014) or by comparing expression profiles between treatments (Bye et al. 2008, Roberts et al. 2013). These observations suggest that expression analyses may be among the most promising strategies to identify the molecular basis of phenotypic differences.

The contrast between expression and frequency changes of coding alleles may be explained by selection acting mostly on alleles in regulatory elements. However other factors cannot be ignored here. Gene expression can be thought of as a first-order phenotype. Because many different SNPs, possibly located in many different genes may affect expression level of a particular gene it may be easier to observe the effects of selection on gene expression levels than on the frequency changes of individual alleles. Thus significant and reproducible changes in expression levels may result from the combined effects of a number of subtle (and not necessary repeatable among selected lines) allele frequency changes in regulatory elements.

For many genes expression is essentially a polygenic trait and as such might be less prone to drift. Therefore random variation between lines in gene expression might be lower than in SNPs, increasing statistical power to detect subtle changes. On the other hand, if the artificial selection in this experiment had resulted in subtle allele frequency changes in many coding SNPs, then this pattern should have been detected by the multidimensional scaling analysis: the selected lines would have formed clusters separated from the controls. However, we did not observe such clustering for the allele frequency data; it was only seen in the gene expression data. This suggest that the changes in gene expression, but not repeatable changes in coding variants, underlie the observed response to selection.

Multidimensional scaling and additional analyses showed that the differences in allele frequencies in coding regions between selected and control lines were driven mainly by genetic drift. These results do not necessarily mean that selection does not affect variation in coding regions; however, they are not compatible with a scenario in which widespread positive selection on coding genes shapes the genomic patterns of polymorphism within, and divergence between, selection regimes. The fact that we did not find SNPs differentiated more than expected under drift effectively rules out the possibility that any genetic variant of large effect was repeatedly selected for. Therefore, we suspect that if positive selection affected the coding sequences in our experiment, it acted on a limited number of variants that provided a small-to-moderate fitness advantage. In this respect, our study contrasts with other selection experiments performed on vertebrates which showed large genomic regions being fixed for alternative variants between treatments (Johansson et al. 2010; Chan et al. 2012; Pettersson et al. 2013). This difference may be a consequence of differences in the genetic architecture of the traits investigated or result from differences in experimental setup, in particular the origin and genetic makeup (eg. presence of linkage disequilibria) of the base population. However, many previous studies relied on observations of reduced genetic diversity as evidence for the effects of selection and, in doing so, may have suffered from the confounding effects of genetic drift. For example, (Johansson et al. 2010) selected chicken lines for high body mass and interpreted the decrease in heterozygosity as reflecting the operation of selection. However, the high-body-mass selected line had an effective population size that was reduced by around 10% (44.5 vs. 49.3). Our experimental design allowed us to ascribe reduced polymorphism in the selected lines to their reduced effective population sizes, and we were able to show that even such minor differences can significantly affect polymorphism.

Our study differed from many other E&R studies (especially these performed on vertebrates) in the nature of the standing genetic variation available at the onset of the experiment. We directly utilized genetic variation that was segregating in a natural population. This has not been the case in many other experiments, in which source populations were created by crossing inbred or isofemale lines (Johansson et al. 2010; Chan et al. 2012; Orozco-Terwengel et al. 2012; Turner and Miller 2012). As a result, experimental populations may not have adequately reflected the standing genetic variation available for positive selection in natural populations. For example, inbred lines are likely to have been cleaned of large-effect recessive deleterious mutations but to have fixed many slightly deleterious ones. When inbred lines are crossed, slightly deleterious mutations become common, and the initial allele frequency spectrum is expected to depart from that observed in nature. In nature most of deleterious variants are rare and the shape of allele frequency spectrum depends on effective population size. Here, we inferred that negative selection is an important force that might shape allele frequencies, even in populations of small Ne.

However, a small Ne is a limitation inherent to E&R studies in vertebrates, and because of this, our study had limited power to detect the effects of selection on SNPs. Pedigree-based selection simulations demonstrated that, due to the effect of drift in relatively small experimental populations, only strongly selected (s~ 0.2) variants  segregating at appreciable frequencies in the base population can be detected with high probability. Therefore the effective size of experimental populations has critical consequences for the E&R approach. Several theoretical studies have examined the effect of population size on analyses of artificial selection, and all of them have found that Ne is a crucial factor that influences the power of such analyses (Kofler and Schlötterer 2013; Baldwin-Brown et al. 2014; Kessner D and Novembre J, unpublished data, http://dx.doi.org/10.1101/005892, last access August 15, 2014). Specifically, Baldwin-Brown and colleagues (2014) argue that, to localize causative SNPs with at least 80% success, researchers should use a population size of 1000 diploid individuals. This is obviously not feasible for most laboratory experiments involving vertebrates, and therefore only variants with large effects can be detected with high probability (Baldwin-Brown et al. 2014). The same situation is however observed in nature — many vertebrate populations are small, having effective population sizes comparable to those reported here, which makes distinguishing effects of drift and selection a challenging task (Palstra and Ruzzante 2008). Additionally, population recombination rate is low in small populations, which increases the rate of false positives because drift affects entire long

haplotypes and leads to correlated allele frequency changes in multiple SNPs. We partially controlled for the effect of linkage by sampling one SNP per gene. However, this problem needs to be considered in future E&R studies.

An alternative explanation for the lack of considerable changes in allele frequencies is that adaptation is due to different variants in different lines. Repeatability of adaptation is however surprisingly high on the gene level both in experimental evolution experiments and natural populations (Conte et al. 2012; Tenaillon et al. 2012; Martin and Orgogozo 2013). Because of that observation and the fact that initial standing genetic variation was similar in all selected lines derived from a single base population, many SNPs initially in moderate frequencies should be repetitively selected. Because the number of repetitively selected coding SNPs was probably modest, and the power to detect them was limited, we attempted to identify potential candidates by exploring the molecular functions of differentiated genes. We also carried out a similar analysis on genes with significantly different expression levels. This strategy is often used in experimental selection surveys (Bye et al. 2008; Kelly et al. 2012; Roberts et al. 2013; Kelly et al. 2014) allowing to pinpoint the most promising candidates for future investigations. Below, we very briefly discuss the molecular processes associated with these plausible candidates.

*Molecular function of plausible candidates*

To assess which biological pathways have possibly changed in response to selection, we investigated the genes that had been identified as having non-overlapping allele frequencies between selected and control lines (despite the overall lack of support for a role of selection in allele frequency changes, some variants may nevertheless be weakly selected for) and those that were differentially expressed in at least one organ. We refer to these genes as "plausible candidates" and list them in Tables S1, S2, and S3. None GO category was significantly overrepresented relative to all GO categories (FDR < 0.05). We argue, though, that some of these plausible candidates are more likely than others to explain some phenotypic changes. We narrowed down the list of candidates based on their functions and present the most interesting genes below.

The stromal interaction molecule 1 (*STIM1*) gene showed highest differences in allele frequency between the selected and control lines, i.e. harbored nonsynonymous SNPs with the

highest diffStat values. STIM1 senses exhaustion of Ca2+ in the endoplasmic reticulum and activates an ion channel in the plasma membrane, causing continuous influx of the extracellular Ca2+ (Kurosaki and Baba 2010). Heterozygous mutations in human *STIM1* cause tubular aggregate myopathy (Bohm et al. 2013) and sotormorken syndrome (Misceo et al. 2014). In tubular aggregate myopathy all patients were characterized by mild and slowly progressive lower limb muscle weakness causing frequent falls and running difficulties (Bohm et al. 2013), which suggests that mutations in *STIM1* may play an important role in swimming performance. Another gene of great interest is that of glycogen phosphorylase (*PYGL*). The physiological role of this liver phosphorylase is to ensure constant supply of glucose for extrahepatic tissues by catalyzing the rate-limiting step in glycogenolysis (Newgard et al. 1989; Bollen et al. 1998). Nonsynonymous mutations in human PYGL cause glycogen storage disease type VI. In substantial number of patients with such disease mild hypotonia, delayed motor development and muscle weakness and cramps were observed (Beauchamp et al. 2007). Interestingly, another nonsynonymous SNP that showed significant differences between selected and control lines was located in the gene that encodes the glycogen-debranching enzyme AGL, which acts together with PYGL to mobilize glucose from glycogen reserves. Mutations in human *AGL* cause Glycogen Storage Disease type III affecting calves and peroneal muscles (Lucchiari et al. 2007).

The gene characterized by the highest number of differentiated nonsynonymous SNPs was *MYO18B*, encoding unconventional myosin XVIIIb. Previous studies have demonstrated the important role of this gene in myocardic structures (Ajima et al. 2008), as well as its contribution to cognitive phenotypes (Purcell et al. 2009; Ludwig et al. 2013).

Another interesting gene is insulin-like growth factor 2 (*IGF2*), its expression increases in response to endurance training and extent of this change differs between humans with highest and lowest improvement in aerobic capacity (Keller et al. 2011). Next gene with an interesting function is the one that encodes fibroblast growth factor 21, which stimulates glucose uptake in adipocytes and plays a critical role in the regulation of lipid homeostasis (Badman et al. 2007). We also identified changes in other genes involved in lipid metabolism (e.g. *ABCG1,CYP17A,APOB,LIPA,APOA1,APOA2,CYP4A14*), the formation and proper functioning of the heart (e.g., *XIRP2, KDM4A, JPH2*), and stress responses (e.g., *IRGM, DELE, PARP, HSP70, HSP105*). All these genes may be involved in response to selection for aerobic performance.

In liver tissue, we found significant differences in expression of the gene that encodes retinoblastoma-like protein-2 (RBL2). RBL2 acts as a transcriptional repressor of the enzymes DNMT3A and DNMT3B, which catalyze the transfer of methyl groups to specific CpG structures in DNA, a process called DNA methylation (Benetti et al. 2008). Also on the list of differentially expressed liver genes are the genes encoding heterogeneous nuclear ribonucleoprotein H2 (HNRPH2), which plays an important role in pre-mRNA processing (Alkan et al. 2006), and methyl-CpG-binding domain protein 4 (MBD4), which takes part in the active demethylation process (Roloff et al. 2003). Additionally, one of the genes whose allele frequencies differed the most between selected and control lines was that coding for lysine-specific demethylase 4A (KDM4A), which plays a central role in modifying the "histone code" (Tan et al. 2011). Taken together, these observations suggest that genes associated with epigenetic changes might represent important targets of selection.

One of the most significant changes in expression level was observed for the gene that encodes aphrodisin—a protein that transports pheromones that stimulate copulatory behavior (Briand et al. 2004, Stopková et al. 2010). Genes coding for aphrodisin-like proteins in bank voles may be used in chemical communication among individuals and thus may play an important role in aggression, dominance, and mate choice (Stopková et al. 2010). Changes in expression of this gene are interesting in the context of differences in reproductive success between the selected and control lines. Already in previous generations of the selection experiment we observed that voles from the selected lines produced litters sooner after the mating (Koteja et al. 2010). Also, in generation 12 and 13 (parents and siblings, respectively, of the voles used in transcriptome analysis), the proportion of mated pairs that produced offspring was significantly higher in the selected than in the control lines (generation 12 – selected: 93.2%, control: 70.1%, p = 0.011; generation 13 – selected: 92.9%, control: 68.2%, p = 0.010; GLIMMIX procedure in SAS 9.3). It is tempting to speculate that changes in the expression of aphrodisin may have been the underlying mechanism.

**Conclusions**

We characterized, through transcriptome sequencing, the response to selection for increased aerobic metabolism in lines derived from a natural population of the bank vole. We showed that the initial response to selection occurs mainly via changes in gene expression. After applying a rigorous control for the effect of drift, no repeatable changes in allele frequencies

at coding SNPs could be unambiguously attributed to directional selection. These results differ from a handful of previous analyses of selection experiments in birds and mammals, in which signals of multiple selective sweeps were detected by resequencing of genomes. Because our selection lines were derived from a natural population, the amount and spectrum of variation available for selection probably closely approximates these typically found in populations of small mammals. Therefore our results are relevant to the understanding of the molecular basis of complex adaptations occurring in vertebrate populations. By combining transcriptome analyses, information about gene functions, and knowledge about selected traits and phenotypes, we identified genes and pathways that could be the targets of selection for increased aerobic metabolism. To further investigate the patterns uncovered here, novel methods that combine knowledge from both population genetics and molecular biology should be developed and exploited in order to effectively characterize the candidate genes that were identified during this experiment.

## Materials and methods

### Selection experiment

This study was performed using individuals from the 13th generation of a laboratory colony of bank voles (*Myodes (=Clethrionomys) glareolus*) that was subjected to selection for improved aerobic metabolism. The rationale for the selection experiment as well as detailed breeding and selection protocols are described elsewhere (Sadowska et al. 2008; Supplementary materials 1.1, 1.2). Briefly, the colony was founded using approximately 320 voles captured in 2000 and 2001 in the Niepołomice Forest in southern Poland. For 6-7 generations, the animals were bred randomly, and the colony was used for quantitative genetic analyses of metabolic rates (Sadowska et al. 2005). In 2004, a multidirectional selection experiment was established. In the A lines analyzed here, the selection criterion was the maximum mass-independent (residual from regression) 1-min rate of oxygen consumption achieved during 18 min of swimming. The swim test was conducted at 38ºC so that no thermoregulatory burden was imposed, and animals were tested when they were around 75 to 85 days old (see Supplementary materials 1.1 and 1.2 for details of the protocol and results of the selection).

*Estimating inbreeding effective population size*

To explore breeding differences among lines, individual inbreeding coefficients were calculated for each line using the R package "pedigree". Changes in inbreeding over time were calculated as:

$$\Delta F_i = \frac{F_i - F_{i-1}}{1 - F_{i-1}}$$

where Fi is the mean inbreeding coefficient in generation i (Falconer and MacKay 1996). The effective population size Ne was calculated for each line according to the formula:

$$N_e = \frac{1}{2 \; x \; \Delta F}$$

where ΔF is the mean change in inbreeding over time.

*Sampling, RNA extraction, and sequencing*

For the transcriptome analysis, five males and five females of 75 to 80 days in age were sampled from each line; each individual came from a different family. These individuals had not been previously used in the swimming trials or for any other specific measurements (except routine measurements of body mass). Voles were euthanized by being placed one by one in a jar containing isoflurane (Aerane®) fumes; this process took place between 8.00 a.m. and 2.00 p.m. The animals were then weighed, and a small part of their left liver lobes and hearts were immediately excised and placed in RNAlater (Sigma). Samples were stored overnight at 4°C and then frozen at -20°C.

Total RNA was extracted using RNAzol® (Molecular Research Center) in accordance with the manufacturer's instructions. RNA concentration and quality were measured using Nanodrop and Agilent 2100 Bioanalyzer, respectively. All samples had an RNA Integrity Number (RIN) higher than 7.0 and were thus suitable for poly-A selection and cDNA library preparation.

For each organ, we prepared one pooled sample per line—using equal amounts of total RNA from each individual—for a total of 16 samples. Residual DNA was removed from

pooled samples using a DNA-free Kit (Ambion®). RNA quality and concentration following the DNAse treatment were assessed as described above.

Poly-A selection, reverse transcription, and the preparation of barcoded cDNA libraries with the TrueSeq RNA kit were performed by the Georgia Genomics Facility, USA. Liver samples from one control line (C3) were pair-end (2 x 100bp) sequenced on an Illumina HiSeq 2000 and used for reference transcriptome construction (Konczal et al. 2014). For the remaining 15 pools, single-end (1 x 100 bp) sequencing was performed. The reads were deposited in Sequence Read Archive (Bioproject PRJNA267038).

*Reference transcriptome reconstruction and annotation*

We first trimmed low-quality reads using DynamicTrim, removed adaptors with Cutadapt, and removed reads shorter than 20bp with LengthSort (Cox et al. 2010; Grabherr et al., 2011; Martin 2011). As references, we used a previously assembled liver transcriptome (Konczal et al. 2014) and the heart reference transcriptome generated by Kaczyńska and colleagues (unpublished), which had been assembled for other purposes. The transcriptomes were processed by merging transcripts that were likely derived from the same genomic locations. This produced transcriptome-based gene models (TGMs), which we refer to here as "genes" (Stuglik et al. 2014).

We did not assemble one transcriptome from pooled reads of two organs to avoid the problem of potential redundancy of the reference transcriptome. Transcriptome complexity negatively affects de novo assembly and TGMs reconstruction (Vija et al. 2013), and it is known that most alternative splicing occurs between organs (Wang et al. 2008). Redundancy of the reference transcriptome has serious implications for SNP calling, because reads mapping equally well to multiple locations are filtered out during this procedure. Because we assembled transcriptomes of each organ separately, this problem is reduced to within-organ splice variant variation and even if it occurs for one transcriptome, SNPs may be still identified using reference transcriptome from the other organ (see below).

TGMs were annotated using Trinotate software. Trinotate makes use of a number of methods for functional annotation (e.g., homology search to Swissprot database, protein domain identification, protein signal prediction) of likely coding regions (likely CDSs). Likely CDSs were identified using a pipeline implemented in Trinity, but this approach did

not successfully annotate all of them. Non-annotated genes represent either errors, fast evolving genes, or genes whose homologs are not present in the Swissprot database.

*Mapping and identification of SNPs*

Filtered reads were mapped to the reference transcriptomes using Bowtie2 (Langmead and Salzberg 2012). For liver samples of line C3, we subsampled reads to obtain a comparable number of single-end sequences. Reads mapped into multiple locations were removed from analyses. We mapped all reads from both organs together to the liver and heart transcriptomes to increase accuracy of allele frequency estimation.

SNP calling was performed in two steps. First, we identified SNPs with samtools (mpileup with options: -Q 10, -E), which is dedicated to diploid genomes (Li et al. 2009). SNPs that contained more than two variants in samtools output were discarded. In the second step, we applied PoPoolation2 (Kofler, et al. 2011) to filter data and estimate allele frequencies. Only SNPs with a minimum of 10x coverage in each sample and a minimum of three reads that supported minor allele were considered. Additional to multi-allelic SNPs removal, two procedures were applied to identify and exclude similar paralogs that had been assembled into single genes: i) we removed most polymorphic genes (more than five SNPs per 100 bp using a minimum of 10x coverage) and ii) we discarded all genes that contained SNPs with an excess of observed heterozygotes or had BLASTN hits with E-value < 10-150 to such genes. This procedure was based on those developed for the individually sequenced liver transcriptomes of 10 voles (Konczal et al. 2014) and the heart transcriptomes of 20 voles (Kaczyńska et al. unpublished)—these studies excluded SNPs for which more than 8/10 or 14/20 samples, respectively, were heterozygotes. Using custom python scripts, SNPs were classified as being synonymous, nonsynonymous, UTRs, or localized in putative non-protein-coding genes.

The above procedure was performed for both liver and heart transcriptomes and yielded highly overlapping sets of SNPs. The differences resulted from differences in reference transcriptomes (lack of a SNP in one transcriptome may be caused by incompletely assembled or non-assembled genes or by splice variants which were not collapsed into a single gene model). To remove the redundancy in the SNP dataset we preformed the following procedure. First we clustered liver and heart SNPs-containing genes using

reciprocal blast searches (BLASTN hits with E-value < 10-100 and > 99% identity). Genes which did not form clusters were apparently expressed in one organ only and were retained (liver: 5,786 genes, heart: 1,759 genes). Genes with significant hits in the other transcriptome were reduced using the criterion of completeness From clusters with one to one relation (containing a single sequence from each transcriptome; 6082 clusters) we retained the longer one. The relation one to many (939 clusters) was mainly caused by fragmentation of the gene in one of the assemblies, therefore we retained SNPs from the transcriptome with the single assembled sequence. For clusters containing >1 sequence from liver and > 1 from heart (many to many, 200 clusters) we included in analyses sequences from this transcriptome in which the total length of sequences was larger. SNPs identified in thus selected genes were used for all analyses. For genes that contained at least one SNP, FST was calculated using PoPoolation2. FST was calculated for each SNP using the formula $FST = (\pi T - \pi W)/\pi T$. Mean FST values between each pair of lines were calculated, and this matrix of pairwise FST was used to test: i) whether the extent of variation among lines within treatments differed between selected and control lines and ii) whether selected and control lines cluster separately. Multivariate homogeneity of group dispersion was tested using betadisper{vegan}, (Oksanene et al. 2013) followed by an ANOVA. To test for separate clustering of selected and control lines we calculated the ratio of between treatment to within treatment variance using adonis{vegan}and assessed its statistical significance through 1000 randomizations. Randomized matrices of mean FST were obtained by shuffling pairwise FST values for each gene independently. The original mean pairwise FST matrix was visualized using nonmetric multidimensional scaling.

*Simulations of allele frequency distribution under drift and positive selection*

To obtain the allele frequency distributions that would be expected under drift, we performed forward drift simulations on known pedigrees. Simulations were performed separately for allele frequency spectra derived from all, synonymous and nonsynonymous SNPs.

The simulations were divided into four parts and were repeated 10 million times (steps 2-4):

1. Estimation of the initial allele frequency distribution. As we did not know the allele frequencies in the ancestral population, we had to estimate them using data from the control lines. For each SNP, we calculated the mean allele frequency from the four control lines. If

control lines diverge mainly due to drift (a reasonable assumption for most polymorphisms), such averages are unbiased estimates of allele frequencies in the ancestral population, which may then be used to reconstruct the allele frequency spectrum in the ancestral population.

2. Simulation of the genotypes of "generation 0" individuals. We simulated the genotype of each individual in the ancestral population by randomly choosing one initial allele frequency (p0) from the set of frequencies estimated in step 1. Then, for each individual, we sampled from a binomial distribution with n = 2 and p = p0 (n - number of draws, p - probability of success), thus obtaining the number of allele copies (0, 1, 2) for each individual.

3. Simulation of the effect of drift on known pedigrees. Based on known pedigrees, we simulated genotypes for each individual by randomly choosing one chromosome from each of the parents. We then obtained genotypes for 10 individuals (that were selected for sequencing for each line) and calculated allele frequencies.

4. Simulation of pooling and sequencing error. Pooling and sequencing cause inaccuracy in allele frequency estimation. Therefore, we decided to add relevant variation to the simulated allele frequencies using the relative errors of allele frequency estimation that had been previously calculated (Konczal et al. 2014). For a given MAF class, a gamma function was fitted to the distribution of experimentally obtained relative errors. Then, one value of estimation error was randomly chosen from the fitted gamma distribution and incorporated into the simulation results.

To asses power to detect selected variant given its selective advantage and initial frequency we used the approach similar as in drift simulations. Several initial allele frequencies (f = 0.05, 0.1, 0.25, 0.5, 0.75, 0.9) and four different values of selection advantage (s = 0, 0.05, 0.1, 0.2) were used. In the course of pedigree based simulations, in selectively bred lines advantageous allele was passed from heterozygote parents to offspring with higher probability (0.5 + ½ s) than the alternative variant (0.5 – ½ s). For each combination of f and s we performed 100000 iterations and recorded the fraction iterations with diffStat > 0.

Scripts used for simulations are available at http://www.molecol.eko.uj.edu.pl.

*SNP analyses, polymorphism, and divergence between lines/selection regimes*

To study the effect of differences in Ne between lines on the amount of genetic variation we examined the allele frequency spectra. Specifically we calculated for each line the number of such SNPs which were polymorphic in the entire dataset but showed little or no variation (minor allele frequency, MAF <0.05) within the line. We used an ANCOVA in which Ne was a covariate and treatment (control vs. selected) was a fixed effect; the interaction between the two was included in order to check the assumption that the slopes were homogeneous between treatments. The interaction was not significant ($F(1, 4) = 0.84$, $p = 0.41$). As a consequence, we used a simple model without interactions to study the general effect of Ne and treatment on allele frequency spectra.

To study differentiation between the selected and control lines, we investigated repeatable changes in allele frequencies. SNPs with frequencies that were either always higher or always lower in selected lines as compared to control lines (non-overlapping allele frequencies) were considered to be potential targets of selection (plausible candidates). For each such site, we calculated the diffStat statistic, which is the smallest difference in allele frequency between selected and control lines (Turner et al. 2011). The distribution of the number of unlinked candidate SNPs was estimated by sampling one SNP per gene 1000 times. We then sampled the same number of SNPs (the number of genes with at least one SNP) from simulated pedigrees 600 times; in these simulations, drift was the only evolutionary force in operation. We subsequently compared the two sets of results. The difference between these two distributions should reveal the genome-wide effects of selection. These analyses were performed on the set of all SNPs, as well as separately for each class of SNPs (synonymous, nonsynonymous, UTR, noncoding).

The biological functions and molecular processes associated with the differentiated genes were studied using custom scripts and Gowinda software (Kofler and Schlötterer 2012).

*Estimation and comparison of gene expression levels*

To identify differentially expressed genes, we mapped reads onto reference transcriptomes with bowtie and used the EdgeR Bioconductor and RSEM packages (Robinson et al. 2010). The matrix of expected counts over all samples was used for EdgeR analyses. Only genes for which the sum of expected counts over all samples was higher than 10 were counted. Using

the standard EdgeR procedure, we normalized counts for library size and RNA composition. We performed multidimensional scaling (BCV method, EdgeR package) over all genes to analyze general expression patterns within tissues. We also estimated dispersion and calculated exact tests for genes that were differentially expressed between control and selected lines. The FDR was calculated as per Benjamini and Hochberg (1995).

The GO terms associated with the differentially expressed genes were investigated with GOrilla software (Eden et al. 2009).

To statistically test for separate clustering of transcriptional profiles of selected and control lines  we developed a procedure analogous to that used for the FST matrix. We used table of expression values (FPKM, TMM normalized) which included only transcripts with the total FPKM > 1. For this table we calculated distance matrix (dist() function) and the ratio of between treatment to within treatment variance (adonis{vegan}). The statistical significance of this ratio was assessed through 1000 randomizations. Randomized matrices of mean gene expression distances were obtained by shuffling expression values of individual gene between lines. Differences between lines in genome-wide transcriptional profiles were visualized with multidimensional scaling (plotMDS{edgeR}).

## Acknowledgments

## References

Ajima R, Akazawa H, Kodama M, Takeshita F, Otsuka A, Kohno T, Komuro I, Ochiya T, Yokota J. 2008. Deficiency of Myo18B in mice results in embryonic lethality with cardiac myofibrillar aberrations. Genes to cells 13: 987-999.

Alkan S, Martincic K, Milcarek C. 2006. The hnRNPs F and H2 bind to similar sequences to influence gene expression. Biochem. J 393: 361-371.

Badman MK, Pissios P, Kennedy AR, Koukos G, Flier JS, Maratos-Flier E. 2007. Hepatic fibroblast growth factor 21 is regulated by PPARα and is a key mediator of hepatic lipid metabolism in ketotic states. Cell metabolism 5: 426-437.

Baldwin-Brown JG, Long AD, Thornton KR. 2014. The power to detect quantitative trait Loci using resequenced, experimentally evolved populations of diploid, sexual organisms. Molecular biology and evolution 31(4): 1040-1055.

Barrett RDH, Hoekstra HE. 2011. Molecular spandrels: tests of adaptation at the genetic level. Nature Reviews Genetics 12: 767-780.

Barrett RDH, Schluter D. 2008. Adaptation from standing genetic variation. Trends in Ecology and Evolution 23: 38-44.

Barrick JE, Yu DS, Yoon SH, Jeong H, Oh TK, Schneider D, Lenski RE, Kim JF. 2009. Genome evolution and adaptation in a long-term experiment with Escherichia coli. Nature 461: 1243-1247. doi: 10.1038/nature08480

Benjamini Y, Hochberg Y. 1995. Controlling the false discovery rate: a practical and powerful approach to multiple testing. Journal of the Royal Statistical Society. Series B (Methodological). 289-300.

Beauchamp NJ, Taybert J, Champion MP, Layet V, Heinz-Erian P, Dalton A, Tanner MS, Pronicka E, Sharrard MJ. 2007. High frequency of missense mutations in glycogen storage disease type VI. Journal of inherited metabolic disease 30(5): 722-734.

Benetti R, Gonzalo S, Jaco I, Muñoz P, Gonzalez S, Schoeftner S, Murchison E, Andl T, Chen T, Klatt P. 2008. A mammalian microRNA cluster controls DNA methylation and telomere recombination via Rbl2-dependent regulation of DNA methyltransferases. Nature structural & molecular biology 15: 268-279.

Bohm J, Chevessier F, Maues De Paula A, Koch C, Attarian S, Feger C, Hantai D, Laforet P, Ghorab K, Vallat J.-M, et. al. 2013. Constitutive activation of the calcium sensor STIM1 causes tubular-aggregate myopathy. Am. J. Hum. Genet. 92: 271-278.

Bollen M, Keppens S, Stalmans W. 1998. Specific features of glycogen metabolism in the liver. Biochem. J 336: 19-31.

Briand L, Trotier D, Pernollet J-C. 2004. Aphrodisin, an aphrodisiac lipocalin secreted in hamster vaginal secretions. Peptides 25: 1545-1552.

Burke MK, Dunham JP, Shahrestani P, Thornton KR, Rose MR, Long AD. 2010. Genome-wide analysis of a long-term evolution experiment with Drosophila. Nature 467: 587-U111. doi: 10.1038/nature09352

Butlin R, Debelle A, Kerth C, Snook RR, Beukeboom LW, Castillo Cajas RF, Diao W, Maan ME, Paolucci S, Weissing FJ, van de Zande L, Hoikkala A, Geuverink E, Jennings J, Kankare M, Knott KE, Tyukmaeva VI, Zoumadakis C, Ritchie MG, Barker D, Immonen E, Kirkpatrick M, Noor M, Garcia MC, Schmitt T, Schilthuizen M. 2012.

What do we need to know about speciation? Trends in Ecology and Evolution 27: 27-39.

Bye A, Langaas M, Høydal MA, Kemi OJ, Heinrich G, Koch LG, Britton SL, Najjar SM, Ellingsen Ø, Wisløff U. 2008. Aerobic capacity-dependent differences in cardiac gene expression. Physiological genomics 33: 100-109.

Chan YF, Jones FC, McConnell E, Bryk J, Bünger L, Tautz D. 2012. Parallel selection mapping using artificially selected mice reveals body weight control loci. Current Biology 22: 794-800.

Chrzascik KM, Sadowska ET, Rudolf A, Koteja P. 2014. Learning ability in bank voles selected for high aerobic metabolism, predatory behaviour and herbivorous capability. Physiology & Behavior 135:143-151 (doi: 10.1016/j.physbeh.2014.06.007).

Conte GL, Arnegard ME, Peichel CL, Schluter D. 2012. The probability of genetic parallelism and convergence in natural populations. *Proceedings of the Royal Society B: Biological Sciences*, 279(1749):5039-5047.

Cox MP, Peterson DA, Biggs PJ. 2010. SolexaQA: At-a-glance quality assessment of Illumina second-generation sequencing data. BMC Bioinformatics 11: 485.

Eden E, Navon R, Steinfeld I, Lipson D, Yakhini Z. 2009. GOrilla: a tool for discovery and visualization of enriched GO terms in ranked gene lists. BMC bioinformatics 10: 48.

Falconer DS, Mackay TFC. 1996. Introduction to quantitative genetics. Harlow (United Kingdom): Addison Wesley Longman.

Fournier-Level A, Korte A, Cooper MD, Nordborg M, Schmitt J, Wilczek AM. 2011. A map of local adaptation in Arabidopsis thaliana. Science 334: 86-89.

Fraser HB. 2013. Gene expression drives local adaptation in humans. Genome research 23: 1089-1096.

Garland Jr T, Rose MR (eds). 2009. Experimental evolution: concepts, methods, and applications of selection experiments. University of California Press, Berkeley

Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L, Raychowdhury R, Zeng Q, Chen Z, Mauceli E, Hacohen N, Gnirke A, Rhind N, Di Palma F, Birren BW, Nusbaum C, Lindblad-Toh K, Friedman N, Regev A. 2011. Full-length transcriptome assembly from RNA-Seq data without a reference genome. Nature Biotechnology 29: 644-652.

Hagberg JM, Rankinen T, Loos RJF, Pirusse L, Roth SM, Wolfarth B, Bouchard C. 2011. Advances in exercise, fitness, and performance genomics in 2010. Medicine and Science in Sports and Exercise 43: 743-752.

Halligan DL, Kousathanas A, Ness RW, Harr B, Eöry L, Keane TM, Adams DJ, Keightley PD. 2013. Contributions of protein-coding and regulatory change to adaptive molecular evolution in murid rodents. PLoS genetics 9: e1003995.

Herring CD, Raghunathan A, Honisch C, Patel T, Applebee MK, Joyce AR, Albert TJ, Blattner FR, van den Boom D, Cantor CR, Palsson BO. 2006. Comparative genome sequencing of Escherichia coli allows observation of bacterial evolution on a laboratory timescale. Nat Genet 38: 1406-1412.

Hoekstra HE, Coyne JA. 2007. The locus of evolution: evo devo and the genetics of adaptation. Evolution 61: 995-1016.

Huang W, Richards S, Carbone MA, Zhu D, Anholt RRH, Ayroles JF, Duncan L, Jordan KW, Lawrence F, Magwire MM, Warner CB, Blankenburg K, Han Y, Javaid M, Jayaseelan J, Jhangiani SN, Muzny D, Ongeri F, Perales L, Wu YQ, Zhang Y, Zou X, Stone EA, Gibbs RA, Mackay TFC. 2012. Epistasis dominates the genetic architecture of Drosophila quantitative traits. Proceedings of the National Academy of Sciences of the United States of America 109: 15553-15559.

Johansson AM, Pettersson ME, Siegel PB, Carlborg Ö. 2010. Genome-wide effects of long-term divergent selection. PLoS genetics 6: e1001188.

Jones FC, Grabherr MG, Chan YF, Russell P, Mauceli E, Johnson J, Swofford R, Pirun M, Zody MC, White S, Birney E, Searle S, Schmutz J, Grimwood J, Dickson MC, Myers RM, Miller CT, Summers BR, Knecht AK, Brady SD, Zhang H, Pollen AA, Howes T, Amemiya C, Baldwin J, Bloom T, Jaffe DB, Nicol R, Wilkinson J, Lander ES, Di Palma F, Lindblad-Toh K, Kingsley DM. 2012. The genomic basis of adaptive evolution in threespine sticklebacks. Nature 484: 55-61.

Kawecki TJ, Lenski RE, Ebert D, Hollis B, Olivieri I, Whitlock MC. 2012. Experimental evolution. Trends in Ecology and Evolution 27: 547-560.

Keller P, Vollaard NB, Gustafsson T, Gallagher IJ, Sundberg CJ, Rankinen T, Britton SL, Bouchard C, Koch LG, Timmons JA. 2011. A transcriptional map of the impact of endurance exercise training on skeletal muscle phenotype. Journal of applied physiology 110(1): 46-59.

Kelly SA, Nehrenberg DL, Hua K, Garland T, Pomp D. 2012. Functional genomic architecture of predisposition to voluntary exercise in mice: expression QTL in the brain. Genetics 181(2): 643-654.

Kelly SA, Nehrenberg DL, Hua K, Garland T, Pomp D. 2014. Quantitative genomics of voluntary exercise in mice: transcriptional analysis and mapping of expression QTL in muscle. Physiological genomics 46(16): 593-601.

King M-C, Wilson AC. 1975. Evolution at two levels in humans and chimpanzees. Science 188: 107-116.

Kofler R, Pandey RV, Schlötterer C. 2011. PoPoolation2: Identifying differentiation between populations using sequencing of pooled DNA samples (Pool-Seq). Bioinformatics 27: 3435-3436.

Kofler R, Schlötterer C. 2013. A guide for the design of evolve and resequencing studies. Molecular biology and evolution 31(2): 474-483.

Kofler R, Schlötterer C. 2012. Gowinda: unbiased analysis of gene set enrichment for genome-wide association studies. Bioinformatics 28: 2084-2085.

Konczal M, Koteja P, Stuglik MT, Radwan J, Babik W. 2014. Accuracy of allele frequency estimation using pooled RNA-Seq. Molecular Ecology Resources 14: 381-392.

Koteja P, Chrząścik KM, Sadowska ET, Ołdakowski Ł. 2010. Correlated responses to a multi-directional selection in bank voles (Myodes glareolus): reproductive parameters. Annual Main Meeting of the Society for Experimental Biology, Prague, 30.06 - 3.07.2010. Programme and Abstract Book: 98.

Kurosaki T, Baba Y. 2010. Ca2+ signaling and STIM1. Progress in biophysics and molecular biology 103(1): 51-58.

Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. Nature Methods 9: 357-359.

Le Corre V, Kremer A. 2012. The genetic differentiation at quantitative trait loci under local adaptation. Molecular Ecology, 21(7): 1548-1566.

Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R. 2009. The Sequence Alignment/Map format and SAMtools. Bioinformatics 25: 2078-2079.

Lucchiari S, Santoro D, Pagliarani S Comi GP. 2007. Clinical, biochemical and genetic features of glycogen debranching enzyme deficiency. Acta myologica 26(1): 72.

Ludwig K, Sämann P, Alexander M, Becker J, Bruder J, Moll K, Spieler D, Czisch M, Warnke A, Docherty S. 2013. A common variant in Myosin-18B contributes to mathematical abilities in children with dyslexia and intraparietal sulcus variability in adults. Translational psychiatry 3: e229.

Malarkey DE, Johnson K, Ryan L, Boorman G, Maronpot RR. 2005. New insights into functional aspects of liver morphology. Toxicologic Pathology 33: 27-34.

Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing reads. EMBnet. journal 17: 10-12.

Martin A, Orgogozo V. 2013. The loci of repeated evolution: a catalog of genetic hotspots of phenotypic variation. Evolution 67(5): 1235-1250.

Mikkelsen TS, Wakefield MJ, Aken B, Amemiya CT, Chang JL, Duke S, Garber M, Gentles AJ, Goodstadt L, Heger A. 2007. Genome of the marsupial Monodelphis domestica reveals innovation in non-coding sequences. Nature 447: 167-177.

Misceo D, Holmgren A, Louch WE, Holme PA, Mizobuchi M, Morales RJ, De Paula AM,

Stray-Pedersen A, Lyle R, Dalhus B, et al. 2014. A dominant STIM1 mutation causes Stormorken syndrome. Hum. Mutat 35: 556-564.

Nespolo RF, Bacigalupe LD, Figueroa CC, Koteja P, Opazo JC. 2011. Using new tools to solve an old problem: the evolution of endothermy in vertebrates. Trends in ecology & evolution 26: 414-423.

Newgard CB, Hwang PK, Fletterick RJ. 1989. The Family of Glycogen Phosphorylases: Structure and Functio. Critical reviews in biochemistry and molecular biology 24: 69-99.

Ohta T. 1982. Linkage disequilibrium due to random genetic drift in finite subdivided populations. Proceedings of the National Academy of Sciences 79(6): 1940-1944.

Oksanen JF, Blanchet G, Kindt R, Legendre P, Minchin PR, O'Hara RB, Simpson GL, Solymos P, Stevens MHH, Wagner H. 2013. vegan: Community Ecology Package. R package version 2.0-10. http://CRAN.R-project.org/package=vegan

Orozco-Terwengel P, Kapun M, Nolte V, Kofler R, Flatt T, Schlãtterer C. 2012. Adaptation of Drosophila to a novel laboratory environment reveals temporally heterogeneous trajectories of selected alleles. Molecular Ecology 21: 4931-4941.

Palstra FP, Ruzzante DE. 2008. Genetic estimates of contemporary effective population size: what can they tell us about the importance of genetic stochasticity for wild population persistence? Molecular Ecology 17: 3428-3447.

Pettersson ME, Johansson AM, Siegel PB, Carlborg Ö. 2013. Dynamics of Adaptive Alleles in Divergently Selected Body Weight Lines of Chickens. G3: Genes| Genomes| Genetics 3: 2305-2312.

Purcell SM, Wray NR, Stone JL, Visscher PM, O'Donovan MC, Sullivan PF, Sklar P, Ruderfer DM, McQuillin A, Morris DW. 2009. Common polygenic variation contributes to risk of schizophrenia and bipolar disorder. Nature 460: 748-752.

Pérusse L, Rankinen T, Hagberg JM, Loos RJF, Roth SM, Sarzynski MA, Wolfarth B, Bouchard C. 2013. Advances in exercise, fitness, and performance genomics in 2012. Medicine and Science in Sports and Exercise 45: 824-831.

Roberts MD, Brown JD, Oberle LP, Heese AJ, Toedebusch RG, Wells KD, Cruthirds CL, Knouse JA, Ferreira JA, Childs TE et al. 2013. Phenotypic and molecular differences between rats selectively bred to voluntarily run high vs. low nightly distances. American Journal of Physiology-Regulatory, Integrative and Comparative Physiology 304(11): R1024-R1035.

Robinson MD, McCarthy DJ, Smyth GK. 2010. edgeR: a Bioconductor package for differential expression analysis of digital gene expression data. Bioinformatics (Oxford, England) 26: 139-140.

Rockman MV. 2012. The QTN program and the alleles that matter for evolution: All that's gold does not glitter. Evolution 66: 1-17.

Roloff TC, Ropers HH, Nuber UA. 2003. Comparative study of methyl-CpG-binding domain proteins. BMC genomics 4: 1.

Roth SM, Rankinen T, Hagberg JM, Loos RJF, Pérusse L, Sarzynski MA, Wolfarth B, Bouchard C. 2012. Advances in exercise, fitness, and performance genomics in 2011. Medicine and Science in Sports and Exercise 44: 809-817.

Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, Cotsapas C, Xie X, Byrne EH, McCarroll SA, Gaudet R. 2007. Genome-wide detection and characterization of positive selection in human populations. Nature 449: 913-918.

Sadowska ET, Baliga-Klimczyk K, Chrzascik KM, Koteja P. 2008. Laboratory model of adaptive radiation: A selection experiment in the bank vole. Physiological and Biochemical Zoology 81: 627-640.

Sadowska ET, Labocha MK, Baliga K, Stanisz A, Wróblewska AK, Jagusiak W, Koteja P. 2005. Genetic correlations between basal and maximum metabolic rates in a wild rodent: Consequences for evolution of endothermy. Evolution 59: 672-681.

Schlötterer C, Kofler R, Versace E, Tobler R, Franssen SU. 2014. Combining experimental evolution with next-generation sequencing: a powerful tool to study adaptation from standing genetic variation. Heredity: 1-10.

Soria-Carrasco V, Gompert Z, Comeault AA, Farkas TE, Parchman TL, Johnston JS, Buerkle CA, Feder JL, Bast J, Schwander T. 2014. Stick Insect Genomes Reveal Natural Selection's Role in Parallel Speciation. Science 344: 738-742.

Stapley J, Reger J, Feulner PGD, Smadja C, Galindo J, Ekblom R, Bennison C, Ball AD, Beckerman AP, Slate J. 2010. Adaptation genomics: The next generation. Trends in Ecology and Evolution 25: 705-712.

Stern DL, Orgogozo V. 2008. The loci of evolution: How predictable is genetic evolution? Evolution 62: 2155-2177.

Stopková R, Zdráhal Z, Ryba Š, Šedo O, Šandera M, Stopka P. 2010. Novel OBP genes similar to hamster Aphrodisin in the bank vole, Myodes glareolus. BMC genomics 11: 45.

Stawski C, Koteja P, Sadowska ET, Jefimow M, Wojciechowski MS. 2015. Selection for high activity-related aerobic metabolism does not alter the capacity of non-shivering thermogenesis in bank voles. Comparative Biochemistry and Physiology A. 180: 51-56.

Storz JF. 2005. INVITED REVIEW: Using genome scans of DNA polymorphism to infer adaptive population divergence. Molecular Ecology 14(3): 671-688.

Stuglik MT, Babik W, Prokop Z, Radwan J. 2014. Alternative reproductive tactics and sex-biased gene expression: the study of the bulb mite transcriptome. Ecology and Evolution 4: 623-632.

Tan M-KM, Lim H-J, Harper JW. 2011. SCFFBXO22 regulates histone H3 lysine 9 and 36 methylation levels by targeting histone demethylase KDM4A for ubiquitin-mediated proteasomal degradation. Molecular and cellular biology 31: 3687-3699.

Tenaillon O, Rodríguez-Verdugo A, Gaut RL, McDonald P, Bennett AF, Long AD, Gaut BS. 2012. The molecular diversity of adaptive convergence. Science 335: 457-461.

Teotónio H, Chelo IM, Bradić M, Rose MR, Long AD. 2009. Experimental evolution reveals natural selection on standing genetic variation. Nature Genetics 41: 251-257.

Turchin MC, Chiang CW, Palmer CD, Sankararaman S, Reich D, Hirschhorn JN, Genetic Investigation of ANthropometric Traits (GIANT) Consortium. 2012. Evidence of widespread selection on standing variation in Europe at height-associated SNPs. Nature genetics, 44(9): 1015-1019.

Turner TL, Miller PM. 2012. Investigating natural variation in Drosophila courtship song by the evolve and resequence approach. Genetics 191: 633-642.

Turner TL, Stewart AD, Fields AT, Rice WR, Tarone AM. 2011. Population-based resequencing of experimentally evolved populations reveals the genetic basis of body size variation in Drosophila melanogaster. PLoS Genetics 7: e1001336.

Vijay N, Poelstra JW, Künstner A, Wolf JB. 2013. Challenges and strategies in transcriptome assembly and differential gene expression quantification. A comprehensive in silico assessment of RNA-seq experiments. Molecular Ecology 22(3): 620-634.

Wang ET, Sandberg R, Luo S, Khrebtukova I, Zhang L, Mayr C, Kingsmore SF, Schroth GP, Burge CB. 2008. Alternative isoform regulation in human tissue transcriptomes. Nature, 456(7221): 470-476.

Wolfarth B, Rankinen T, Hagberg JM, Loos RJF, Pérusse L, Roth SM, Sarzynski MA, Bouchard C. 2014. Advances in exercise, fitness, and performance genomics in 2013. Medicine and Science in Sports and Exercise 46: 851-859.

Wray GA. 2007. The evolutionary significance of cis-regulatory mutations. Nature Reviews Genetics 8: 206-216.

**Supplementary materials**


**Supplement 1. Animal maintenance, selection protocol and effects of the selection in the first 15 generations of the selection experiment on bank voles**


*Animal maintenance and welfare*

The animals were maintained in standard plastic mouse cages (mostly opaque, polypropylene) with sawdust bedding, at a constant temperature (20±1°C) and photoperiod (16h:8h light:dark; light phase starting at 2:00am). Breeding pairs and pairs with offspring (up to 17 days old) were maintained in model 1290D cages (Tecniplast, Bugugiatte, Italy; dimensions L x W x H: 425 x 266 x 155 mm, floor area 800 cm2), equipped with a shelter (ceramic pot), additional nest material (paper towels) and cardboard tubes (environment enrichment). At the age of 17 days the animals were weaned, marked temporarily by fur clipping and kept in family groups until the age of 30-35 days. At the age of about 34 days, all individuals were marked permanently with mouse ear tags (model 10005-1; National Band and Tag, Newport, KY; mass 0.18g) and later maintained in same-sex groups of three individuals in model 1264C cages (Tecniplast, Bugugiatte, Italy; dimensions L x W x H: 267 x 207 x 140 mm; floor area 370 cm2) or up to five (usually four) individuals in the larger model 1290D cages (described above). Cages were changed every 5-14 days, depending on the number of animals in the cage, size of the cage and their cleanliness. Water and food (a "breeding type" rodent chow: 24% protein, 3% fat, 4% fibre; Labofeed H, Kcynia, Poland) was provided ad libitum. Every day all cages were visually inspected for presence of food and water or dead animals. The colony was under supervision of a qualified veterinary surgeon. During any kind of measurements, if symptoms of poor condition were observed in an animal (problems with breathing or moving, injury, etc.), it was removed from the experiment. Depending on judgment of the observer or animal care personnel, it was either allowed to recover or was euthanized. Depending on circumstances, one of three methods of euthanasia was used: exposure to a rising concentration of $CO_2$, cervical dislocation, or isoflurane inhalation (AEranne, Baxter; applied using open-drop technique).

After completing all the measurements of the 13th generation, we discovered that the colony had been infected with Puumala hantavirus during that generation. The virus was not detected earlier because, under normal housing conditions, infection does not result in

pathology in bank voles (Bernshtein et al. 1999); some data suggest, however, that the virus may decrease vole survival under harsh winter conditions (Kallio et al. 2007). We confirmed that reproduction (litter mass and litter size during weaning), mortality, and condition (adult body mass) in the "infected" generations did not differ from the preceding, "uninfected" generations. Therefore, it is highly unlikely that infection influenced the results of the selection experiment.


*Selection protocol*

The whole selection experiment evolves lines selected in three distinct directions (Sadowska et al. 2008), but in this study we used only lines selected for high aerobic-exercise metabolism, measured as the highest 1-minute rate of oxygen consumption achieved by the voles during swimming. The measurements were performed in a positive pressure open-flow respirometric system (design 1b in Koteja 1996), similarly as described in our earlier reports (Sadowska et al. 2005). The voles swam in a 3L glass chamber partly filled with water with a drop of a shampoo for dogs (to ensure complete soaking of fur). Unlike in the earlier study, the measurements were performed at 38°C, to ensure that the increase of metabolism was solely due to locomotor activity and not due to thermoregulatory demand. The tests lasted for up to 18 minutes, unless an individual began to sink or oxygen consumption rapidly decreased.

In each generation we used two custom built computerized respirometric systems: 1) with both O2 and CO2 analysers: FC2 Oxzilla and CA2-2A (Sable systems, Las Vegas, NV), or 2) with only O2 analyser: either Applied Electrochemistry S3-A/II (AMETEK, now AEI Technologies, Pittsburgh, PA) or FC-1 (Sable Systems, Las Vegas, NV), depending on generation. Fresh air was pumped through the respirometric chamber at about 2000 ml/min (STPD), controlled with one of the following mass-flow controller systems: either ERG3000 (Beta-Erg, Warsaw, Poland) or GFC-17A (Aalborg, Orangeburg, New York) or two parallel MFS-1 (Sable Systems, Las Vegas, NV). Actual flow rate controlled by the three mass-flow systems was calibrated against the same precise glass rotameter (LO 2.5/100, Rota, Germany). Sample of excurrent air (150-200ml/min) was dried and directed to the gas analyzers. Concentration of gases was recorded every second, and the rate of oxygen consumption was calculated according to appropriate equations (Koteja 1996; modified to include information about CO2 concentration, if available). Results obtained with the first system were used to

calculate respiratory exchange rate (RQ; ratio of $CO_2$ production to $O_2$ consumption rates), and averaged RQ values were used in calculation of oxygen consumption with the second system (which did not measure $CO_2$). Raw values of oxygen consumption and $CO_2$ production calculated for each 1-sec interval were corrected for effective volume of the systems to achieve "instantaneous" rates (Bartholomew et al. 1981).

In the first three generations the respirometric measurements were performed twice, at the age of 74-86 days and again about 10 days later. The two measurements were highly repeatable and selection decision did not change markedly after consideration of the second measurement. Therefore, in further generations the measurements were performed only once (in most individuals at the age of 75-85 days). The selection criterion was the 1-minute maximum instantaneous rate of oxygen consumption adjusted (residual from ANCOVA) for body mass, sex, number of litter, litter size, age, measurement date and the type of respirometer.

In each generation we obtained offspring from at least 16 families, with 1-4 litters from a family (because average litter size in bank voles is about 4.5, multiple litters from a family were needed to allow effective selection). From each of the families 1-2 males and 1-2 females with the highest adjusted (residual) values of the maximum metabolism were selected as breeders. However, the selection was not purely of "within-family" type, because a) more individuals were selected as breeders from "good" families, in which average scores were higher than a line mean, and b) if for a given line in a particular generation an excess number of families was available, we have not selected any individuals from "poor" families, in which scores of all individuals were below the line mean. Thus, the effective population size was lower than predicted for a breeding scheme with purely within-family selection (see Results). We decided to trade-off the effective population size for an increased efficiency of selection, because we anyway breed a larger number of families per line than in typical selection experiments on laboratory rodents (e.g. Swallow et al. 1998), but have a longer generation time (due to producing multiple litters). The males and females selected as breeders were mated in pairs randomly, but with the restriction that mating between siblings and first cousins was avoided.

All the breeding and experimental protocols have been approved by the Polish State and Local Ethical Committee for Ethics in Animal Research in Kraków (decisions No. DB/KKE/PL-111/2001, 31/OP/2005, 99/2006, 21/2010 and 22/2010).

*Effects of the selection*

Significant differences between selected and control lines (P<0.001) were observed already since the second generation of selection (Fig. S1.1) and increased gradually in further generations. In generation 12, we were able to perform the measurements in only a small sample of individuals, and therefore the selection was relaxed (breeders were chosen randomly form each line). Thus, individuals used for the transcriptome analysis (gen. 13) are offspring of non-selected animals, representing a random sample of their lines. Therefore the results are not influenced by possible maternal effects associated with phenotypic selection on parental generation. In generation 13 the maximum rate of aerobic metabolism achieved during swimming was 48% higher in voles from the selected lines than in those from the unselected control lines (48% higher in A line voles than in C line voles (mean ± SD: $5.32 \pm 0.64$ ml O2/min vs. $3.59 \pm 0.57$ ml O2/min, respectively; p<0.0001). The difference continued to increase in further two generations, and at this point it is not possible determine, whether a selection limit has been already approached.
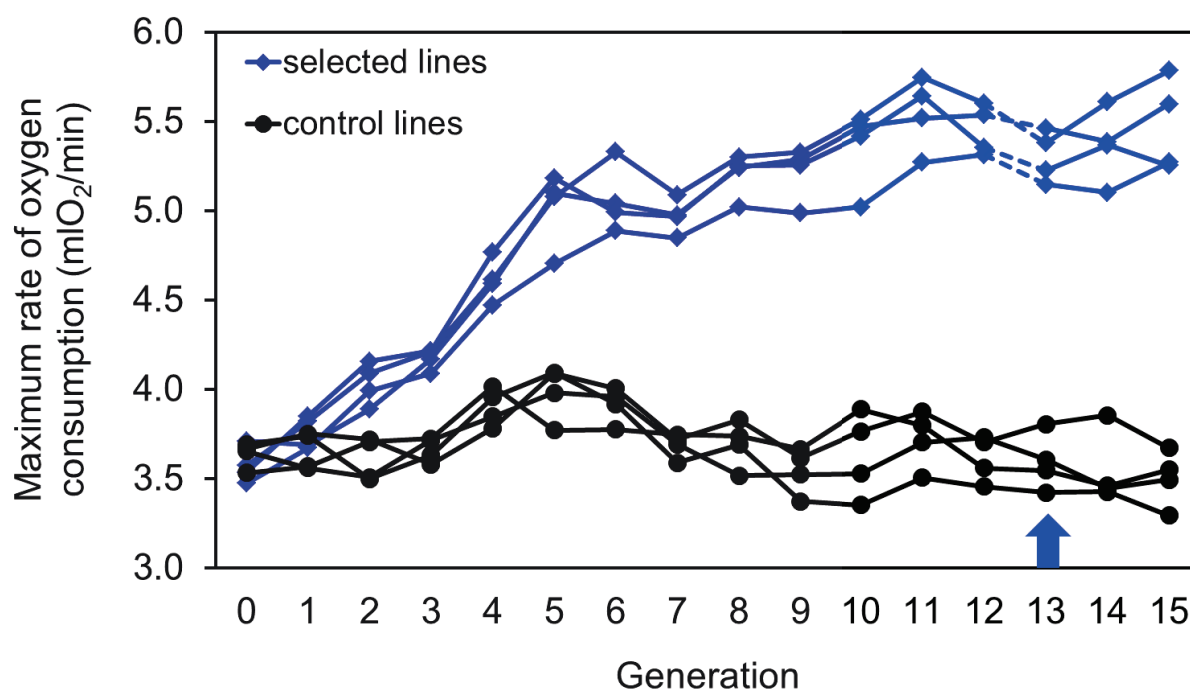


**Fig. S1.1**. Results of 15 generations of selection for the maximum rate of oxygen consumption achieved during swimming in the bank vole (replicate line means). In generation 12 the selection was relaxed (broken lines). The transcriptome analysis was performed on voles from generation 13.

In addition to the direct response to the selection we observed several correlated responses in behavioral and morpho-physiological traits, which were presented as preliminary conference reports (e.g. Koteja et al 2009, 2010, 2012, 2013, Sadowska et al 2013). Voles from the selected lines were more active in home cages, had increased basal metabolic rates, maximum forced-running and cold-induced metabolic rates, and food consumption rate, increased mass of gastrocnemius leg muscle, heart, kidney, liver and brain, and reproduced better, compared to voles from the unselected control lines. On the other hand, we found no significant difference between the lines in learning capability (Chrząścik et al. 2014), oxidative damage to lipids or proteins (Ołdakowski et al. 2012) or non-shivering thermogenesis (Stawski et al., 2015).

*References*

Bartholomew GA, Vleck D, Vleck CM. 1981. Instantaneous measurements of oxygen consumption during pre-flight warm-up and post-flight cooling in sphingid and saturnid moths. Journal of Experimental Biology 90:17-32.

Bernshtein A, Apekina N, Mikhailova T, Myasnikov YA, Khlyap L, Korotkov YS, Gavrilovskaya I. 1999. Dynamics of Puumala hantavirus infection in naturally infected bank voles (Clethrinomys glareolus). Archives of virology 144: 2415-2428.

Chrzascik KM, Sadowska ET, Rudolf A, Koteja P. 2014. Learning ability in bank voles selected for high aerobic metabolism, predatory behaviour and herbivorous capability. Physiology & Behavior 135:143-151.

Kallio ER, Voutilainen L, Vapalahti O, Vaheri A, Henttonen H, Koskela E, Mappes T. 2007. Endemic hantavirus infection impairs the winter survival of its rodent host. Ecology 88: 1911-1916.

Koteja P. 1996. Measuring energy metabolism with open flow respirometric systems: which design to choose? Functional Ecology 10:675-677.

Koteja P, Baliga-Klimczyk K, Chlad A, Chrzascik KM, Damulewicz M, Dragosz-Kluska D, Morawska-Ploskonka J, Sadowska ET. 2009. Correlated responses to a multidirectional artificial selection in the bank vole: activity, metabolism, and food consumption. Journal of Physiological Sciences 59 (Suppl. 1): 541 (XXXVI International Congress of Physiological Sciences IUPS2009, Kyoto, July 27 - August 1, 2009).

Koteja P, Chrząścik KM, Sadowska ET, Ołdakowski Ł. 2010. Correlated responses to a multi-directional selection in bank voles (Myodes glareolus): reproductive parameters. Annual Main Meeting of the Society for Experimental Biology, Prague, 30.06 - 3.07.2010. Programme and Abstract Book: 98.

Koteja P, Baliga-Klimczyk K, Chrząścik KM, Dheyongera G, Konczal M, Mait, U, Rudolf A, Stawski C, Sadowska ET. 2012. Correlated evolution of behaviour and physiology in a multidirectional artificial selection on a wild rodent, the bank vole. Annual Main Meeting of the Society for Experimental Biology, Salzburg, Austria, 28.06-2.07.2012. Programme and Abstract Book: 156-157.

Koteja P, Baliga-Klimczyk K, Chrząścik KM, Dheyongera G, Konczal M, Maiti U, Orłowska P, Rudolf A, Stawski C, Sadowska ET. 2013. Correlated responses to a multidirectional artificial selection in the bank vole: changes in organ size. IUPS 2013, 21-26.07.2013, Birmingham, UK. Abstract book: 610P.

Ołdakowski, Ł., Piotrowska, Ż., Chrząścik, K.M., Sadowska, E.T., Koteja, P. and Taylor, J.R.E. 2012. Is reproduction costly? No increase of oxidative damage in breeding bank voles. Journal of Experimental Biology 215: 1799-1805 (doi:10.1242/jeb.068452).

Sadowska ET, Labocha MK, Baliga K, Stanisz A, Wróblewska AK, Jagusik W, Koteja P. 2005. Genetic correlations between basal and maximum metabolic rates in a wild rodent: consequences for evolution of endothermy. Evolution 59: 672-681.

Sadowska ET, Baliga-Klimczyk K, Chrzascik KM, Koteja, P. 2008. Laboratory model of adaptive radiation: A selection experiment in the bank vole. Physiological and Biochemical Zoology 81: 627–640.

Sadowska ET, Chrzascik KM, Dheyongera G, Rudolf A, Koteja P. 2013. Maximum cold-induced food consumption in bank voles selected for high swim-induced aerobic capacity: Implications for the evolution of endothermy. Annual Main Meeting of the Society for Experimental Biology, 3-6.07.2013, Valencia, Spain. Abstract Book: 100.

Stawski C, Koteja P, Sadowska ET, Jefimow M, Wojciechowski MS. 2015.. Selection for high activity-related aerobic metabolism does not alter the capacity of non-shivering thermogenesis in bank voles. Comparative Biochemistry and Physiology A. 180: 51-56.

Swallow JG, Carter PA, Garland T, Jr. 1998. Artificial selection for increased wheel-running behavior in house mice. Behavior Genetics 28:227-237.
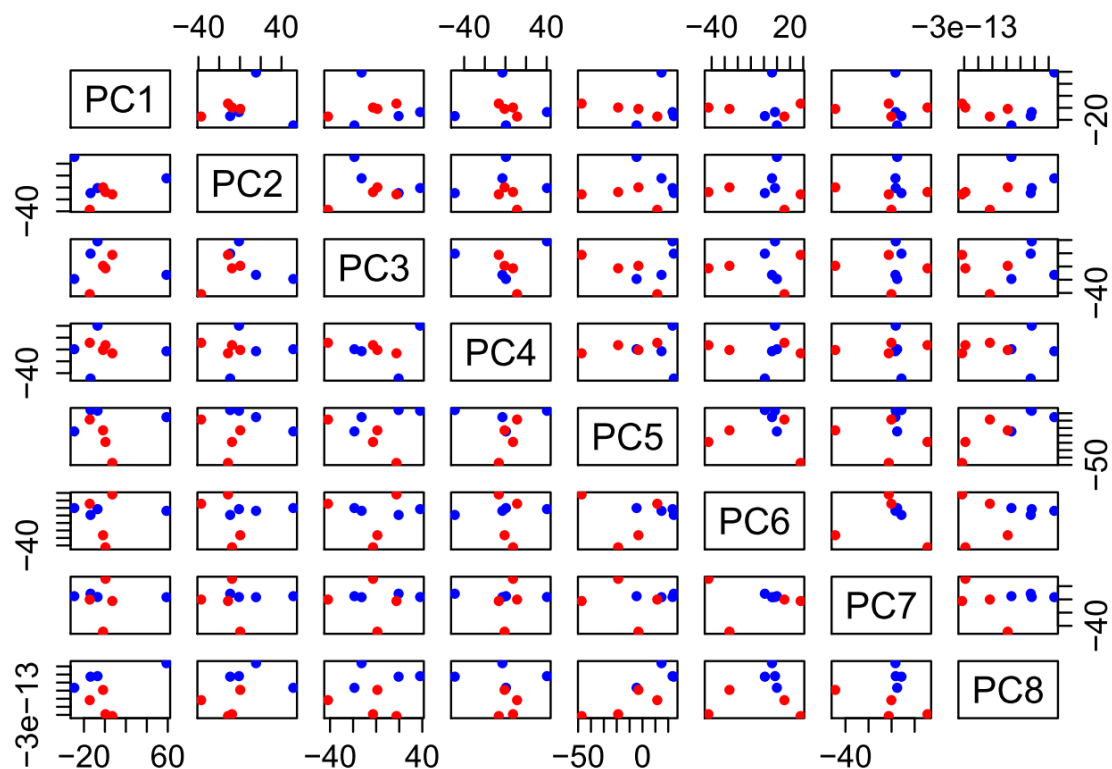
**Fig. S2.1**. PCA plot based on allele frequencies estimated from pooled transcriptomes of liver and heart samples. Blue circles represent four selected lines; red show four control lines of the selection experiment. PCA was performed to look for corelated changes in allele frequencies in various subsets of SNPs, which could reflect response of multiple genes to the same selection pressure.
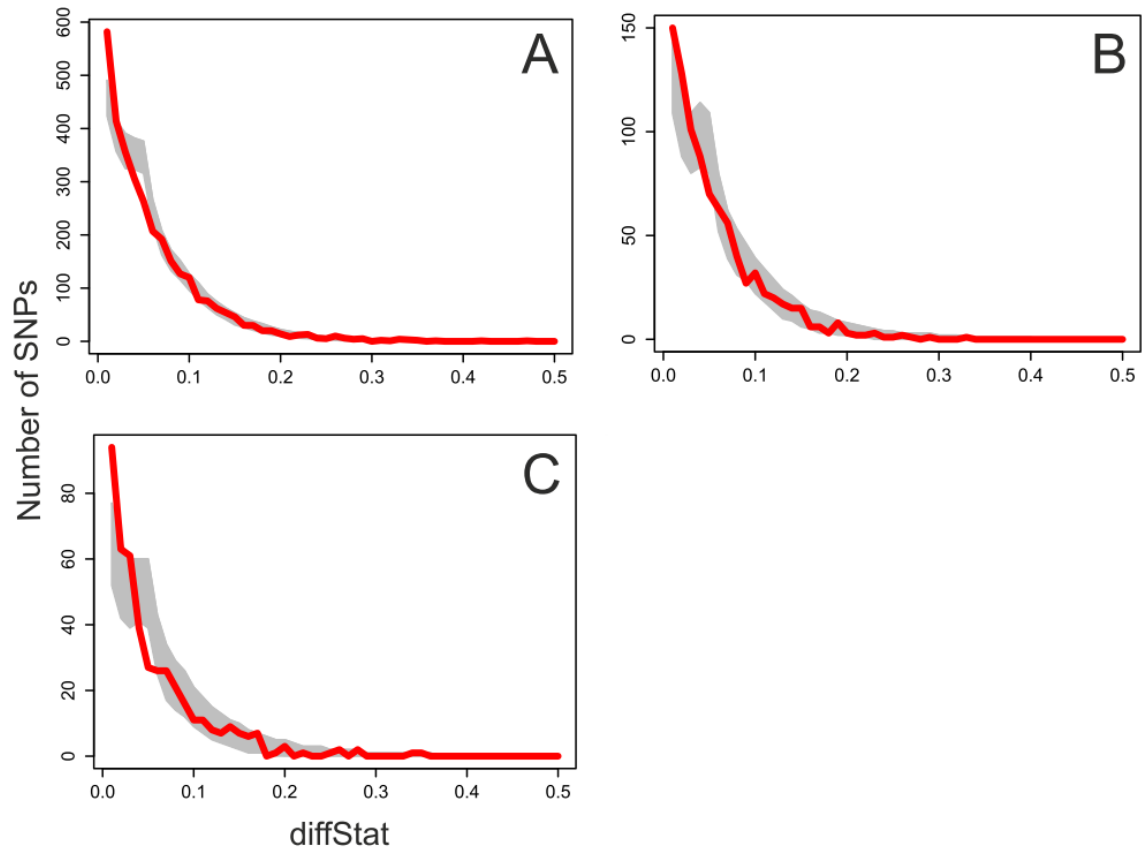
**Fig. S2.2**. Distribution of diffStat values expected from simulations (grey shaded area indicates 90% of all simulations) and observed (red line) for all (A), synonymous (B), and nonsynonymous (C) SNPs. The simulations were performed on pedigree and expected distribution was obtained from iterations showing diffStat > 0. DiffStat is the minimum allele frequency difference between selected and control lines calculated for SNPs with non-overlapping allele frequencies between treatments.
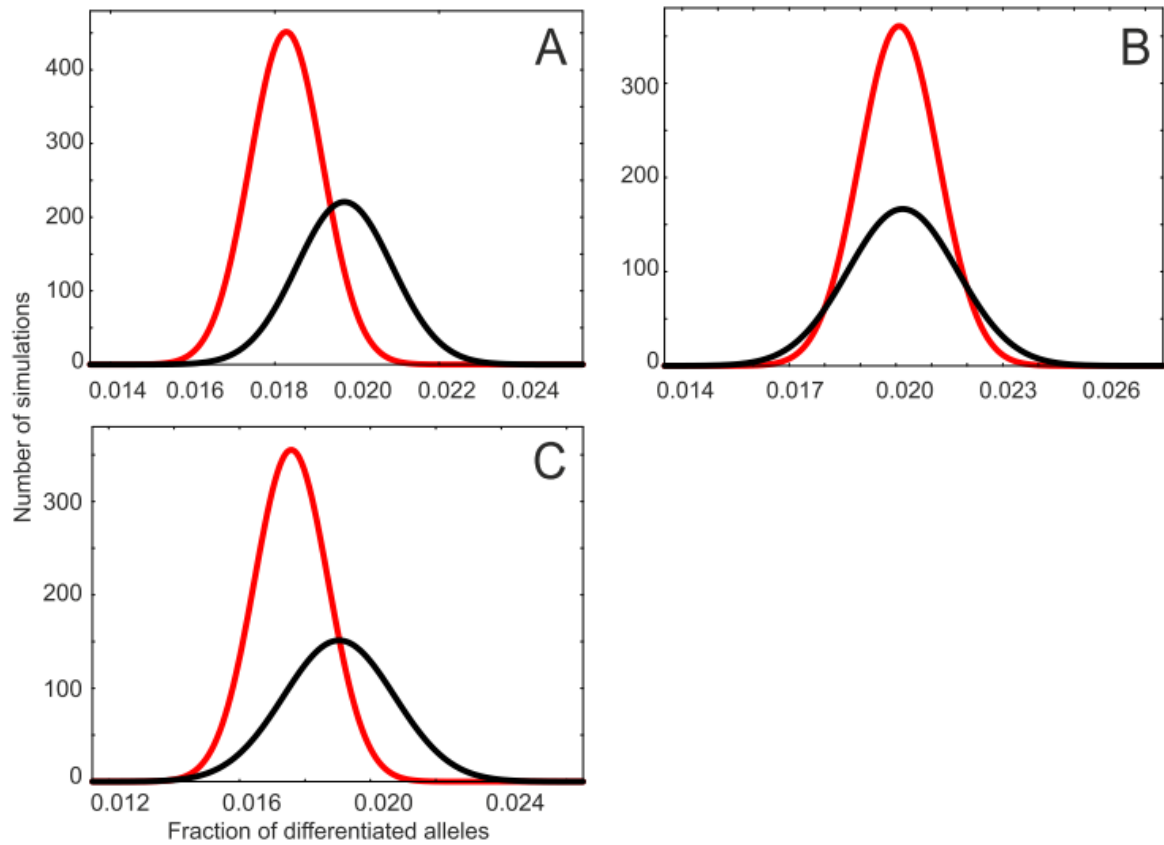
**Fig. S2.3**. Expected and observed fraction of SNPs with non-overlapping allele frequencies between selected and control lines (diffStat > 0). The plots are shown for all (A), synonymous (B) and nonsynonymous (C) SNPs. Red line represents observed variants, black – expectations from pedigree-based simulations. The distribution for observed values was obtained from 1000 datasets, consisting of SNPs sampled one per gene.
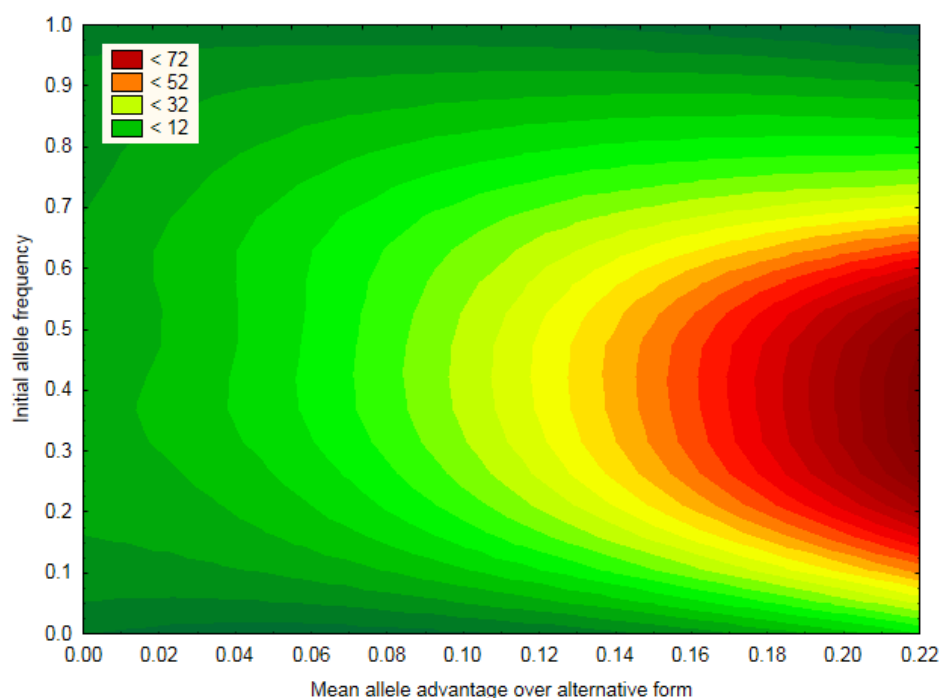
**Fig. S2.4**. Power of detection selected allele within alleles with non-overlapping allele frequencies between four selected lines and four controls. The results were obtained from pedigree-based simulations, where, in selected lines, one allele had higher probability being sampled from parents to offspring than alternative. The colors show the expected percent of SNPs with non-overlapping allele frequencies at the 13 generation of selection.

*Tables*

**Tab. S1** List of plausible candidate genes from SNP analysis. Max nonsyn DiffStat – maximum value of DiffStat (minimum difference in allele frequency between treatments) for a given gene; N of nonsyn/syn cand. – number of SNPs (nonsynonymous/synonymous) with non-overlapping allele frequencies; N of all SNPs – number of all SNPs for given gene.

| Gene | Max nonsyn DiffStat | N of nonsyn cand. | N of syn cand. | N of all SNPs |
|---|---|---|---|---|
| Stromal interaction molecule 1 | 0.28 | 1 | 0 | 5 |
| Cytosolic phospholipase A2 gamma | 0.26 | 1 | 0 | 9 |
| Lysine-specific demethylase 4A | 0.25 | 2 | 2 | 20 |
| Zinc finger FYVE domain-containing protein 9 | 0.22 | 2 | 2 | 18 |
| Glycogen phosphorylase, liver form | 0.2 | 5 | 0 | 54 |
| Putative zinc finger protein 724 | 0.2 | 3 | 0 | 6 |
| H-2 class II histocompatibility antigen, E-K alpha chain | 0.17 | 2 | 3 | 23 |
| H-2 class II histocompatibility antigen, E-K alpha chain | 0.17 | 2 | 3 | 30 |
| Pentatricopeptide repeat-containing protein 1, mitochondrial | 0.17 | 1 | 0 | 7 |
| LIM and senescent cell antigen-like-containing domain protein 2 | 0.17 | 1 | 1 | 3 |
| Apolipoprotein B-100 | 0.16 | 2 | 3 | 173 |
| Protein FAM98B | 0.16 | 1 | 0 | 9 |

| | | | | |
|---|---|---|---|---|
| Monocarboxylate transporter 12 | 0.16 | 3 | 2 | 109 |
| Xin actin-binding repeat-containing protein 2 | 0.16 | 6 | 4 | 251 |
| Steroid 17-alpha-hydroxylase/17,20 lyase | 0.15 | 1 | 1 | 14 |
| ATP-binding cassette sub-family G member 8 | 0.15 | 1 | 2 | 40 |
| Centrosomal protein of 89 kDa | 0.15 | 1 | 0 | 22 |
| Cytochrome b5 reductase 4 | 0.14 | 2 | 1 | 31 |
| G patch domain and ankyrin repeat-containing protein 1 | 0.14 | 2 | 0 | 13 |
| Interferon-induced very large GTPase 1 | 0.14 | 2 | 1 | 164 |
| Polyamine-modulated factor 1 | 0.14 | 1 | 0 | 6 |
| Poly [ADP-ribose] polymerase 9 | 0.13 | 6 | 0 | 101 |
| O-phosphoseryl-tRNA(Sec) selenium transferase | 0.13 | 2 | 0 | 3 |
| Alpha-2-macroglobulin | 0.13 | 2 | 0 | 61 |
| Putative sodium-coupled neutral amino acid transporter 10 | 0.13 | 1 | 0 | 16 |
| Immunity-related GTPase family M protein 1 | 0.12 | 5 | 0 | 17 |
| Alpha-mannosidase 2 | 0.12 | 1 | 0 | 60 |
| N-acetyltransferase 10 | 0.12 | 1 | 0 | 22 |
| Peroxisomal membrane protein PEX14 | 0.12 | 1 | 0 | 10 |
| Alanine--glyoxylate aminotransferase 2, mitochondrial | 0.12 | 1 | 4 | 42 |
| Nidogen-2 | 0.12 | 1 | 1 | 18 |
| Pleckstrin homology domain-containing family G member 3 | 0.11 | 3 | 0 | 44 |
| Uncharacterized protein C20orf194 homolog | 0.11 | 1 | 0 | 41 |
| ATP-binding cassette sub-family B member 8, mitochondrial | 0.11 | 2 | 2 | 13 |
| ERBB receptor feedback inhibitor 1 | 0.11 | 2 | 0 | 15 |
| Trinucleotide repeat-containing gene 18 protein | 0.11 | 1 | 1 | 12 |
| Nuclear RNA export factor 1 | 0.11 | 1 | 0 | 9 |
| Microtubule-associated tumor suppressor 1 homolog | 0.1 | 2 | 1 | 80 |
| Cytosolic 5'-nucleotidase III-like protein | 0.1 | 1 | 0 | 15 |
| NADP-dependent malic enzyme | 0.1 | 1 | 0 | 35 |
| Glycerol-3-phosphate acyltransferase 3 | 0.1 | 1 | 0 | 29 |
| Aldehyde dehydrogenase family 16 member A1 | 0.1 | 1 | 3 | 21 |
| Phospholipid transfer protein | 0.09 | 2 | 2 | 34 |
| Death ligand signal enhancer | 0.09 | 3 | 1 | 47 |
| Lysosomal acid lipase/cholesteryl ester hydrolase | 0.09 | 3 | 0 | 37 |
| Protein LZIC | 0.09 | 1 | 1 | 9 |
| E3 ubiquitin-protein ligase RFWD3 | 0.09 | 1 | 0 | 12 |
| Zinc finger protein with KRAB and SCAN domains 5 | 0.09 | 2 | 0 | 8 |
| MKI67 FHA domain-interacting nucleolar phosphoprotein | 0.09 | 2 | 0 | 15 |
| Junctophilin-2 | 0.09 | 1 | 1 | 45 |
| Filamin-C | 0.09 | 1 | 8 | 55 |
| Trans-1,2-dihydrobenzene-1,2-diol dehydrogenase | 0.09 | 1 | 2 | 7 |
| Glycogen debranching enzyme | 0.09 | 1 | 0 | 63 |
| Zinc finger protein 260 | 0.08 | 1 | 1 | 15 |
| ADP-ribosylation factor GTPase-activating protein 3 | 0.08 | 1 | 0 | 10 |
| Molybdenum cofactor sulfurase | 0.08 | 1 | 0 | 61 |
| Early endosome antigen 1 | 0.08 | 1 | 0 | 52 |
| Titin | 0.08 | 1 | 0 | 30 |

| | | | | |
|---|---|---|---|---|
| Transmembrane 9 superfamily member 3 | 0.08 | 1 | 0 | 19 |
| Laminin subunit gamma-2 | 0.08 | 1 | 0 | 18 |
| Transcription initiation factor TFIID subunit 2 | 0.08 | 1 | 0 | 15 |
| Centrosome-associated protein 350 | 0.08 | 1 | 0 | 15 |
| Coiled-coil domain-containing protein 66 | 0.08 | 1 | 0 | 7 |
| Protein FAN | 0.08 | 1 | 0 | 15 |
| Retrograde Golgi transport protein RGP1 homolog | 0.08 | 1 | 0 | 25 |
| EH domain-binding protein 1-like protein 1 | 0.07 | 4 | 0 | 30 |
| Basement membrane-specific heparan sulfate proteoglycan core protein | 0.07 | 1 | 4 | 45 |
| Basement membrane-specific heparan sulfate proteoglycan core protein | 0.07 | 1 | 3 | 6 |
| DNA-directed RNA polymerase III subunit RPC9 | 0.07 | 1 | 0 | 37 |
| Protein KRBA1 | 0.07 | 2 | 0 | 41 |
| Phosphoacetylglucosamine mutase | 0.07 | 1 | 0 | 8 |
| Atrophin-1 | 0.07 | 1 | 1 | 29 |
| Coiled-coil domain-containing protein 68 | 0.07 | 3 | 1 | 33 |
| Eukaryotic translation initiation factor 2-alpha kinase 1 | 0.07 | 2 | 0 | 46 |
| Cytochrome P450 3A11 | 0.07 | 2 | 1 | 7 |
| Neurogenic locus notch homolog protein 4 | 0.07 | 2 | 2 | 191 |
| Non-lysosomal glucosylceramidase | 0.07 | 2 | 0 | 36 |
| Mitochondrial assembly of ribosomal large subunit protein 1 | 0.07 | 1 | 1 | 10 |
| Thyroid transcription factor 1-associated protein 26 | 0.07 | 1 | 0 | 10 |
| Zinc finger protein 48 | 0.07 | 1 | 0 | 2 |
| Very-long-chain (3R)-3-hydroxyacyl-[acyl-carrier protein] dehydratase 3 | 0.07 | 1 | 0 | 30 |
| Coagulation factor X | 0.07 | 1 | 1 | 23 |
| Phosphoenolpyruvate carboxykinase [GTP], mitochondrial | 0.07 | 1 | 2 | 5 |
| LIM domain and actin-binding protein 1 | 0.07 | 1 | 0 | 39 |
| Zinc finger and BTB domain-containing protein 25 | 0.07 | 1 | 0 | 8 |
| Coiled-coil domain-containing protein 51 | 0.07 | 1 | 0 | 13 |
| Neuralized-like protein 2 | 0.06 | 4 | 2 | 34 |
| Acyl-coenzyme A thioesterase 2, mitochondrial | 0.06 | 3 | 0 | 8 |
| Epiplakin | 0.06 | 2 | 3 | 24 |
| Gag-Pro polyprotein | 0.06 | 2 | 1 | 63 |
| 7-alpha-hydroxycholest-4-en-3-one 12-alpha-hydroxylase | 0.06 | 2 | 0 | 36 |
| Annexin A6 | 0.06 | 1 | 3 | 23 |
| Probable glutathione peroxidase 8 | 0.06 | 1 | 0 | 17 |
| Coagulation factor V | 0.06 | 1 | 0 | 71 |
| KN motif and ankyrin repeat domain-containing protein 3 | 0.06 | 1 | 1 | 13 |
| Unconventional myosin-XVIIIb | 0.06 | 10 | 7 | 105 |
| Protein DGCR14 | 0.06 | 2 | 0 | 37 |
| Cytochrome P450 2C26 | 0.06 | 2 | 1 | 7 |
| Transmembrane protein 175 | 0.06 | 1 | 0 | 15 |
| 5'-AMP-activated protein kinase subunit gamma-2 | 0.06 | 1 | 0 | 12 |
| Inhibitor of Bruton tyrosine kinase | 0.06 | 1 | 0 | 32 |
| Dystonin | 0.05 | 3 | 6 | 82 |

| | | | | |
|---|---|---|---|---|
| Transmembrane emp24 domain-containing protein 3 | 0.05 | 2 | 0 | 19 |
| PR domain zinc finger protein 2 | 0.05 | 1 | 0 | 11 |
| Methionyl-tRNA formyltransferase, mitochondrial | 0.05 | 1 | 1 | 18 |
| Alcohol dehydrogenase 4 | 0.05 | 1 | 0 | 22 |
| Transmembrane emp24 domain-containing protein 2 | 0.05 | 1 | 0 | 8 |
| Histone-lysine N-methyltransferase SETDB2 | 0.05 | 1 | 0 | 3 |
| Regulatory solute carrier protein family 1 member 1 | 0.05 | 1 | 0 | 64 |
| Sugar phosphate exchanger 3 | 0.05 | 1 | 0 | 7 |
| Cullin-associated NEDD8-dissociated protein 2 | 0.05 | 1 | 4 | 37 |
| Plexin-B2 | 0.05 | 2 | 1 | 40 |
| Alpha-methylacyl-CoA racemase | 0.05 | 1 | 0 | 25 |
| D-3-phosphoglycerate dehydrogenase | 0.05 | 1 | 1 | 33 |
| Podocalyxin | 0.05 | 1 | 0 | 82 |
| Procollagen C-endopeptidase enhancer 1 | 0.05 | 1 | 2 | 4 |
| C-reactive protein | 0.04 | 3 | 1 | 19 |
| Lysophospholipase-like protein 1 | 0.04 | 2 | 1 | 8 |
| Zinc finger protein 791 | 0.04 | 2 | 0 | 5 |
| Cytochrome P450 3A29 | 0.04 | 1 | 0 | 10 |
| Hydroxymethylglutaryl-CoA synthase, cytoplasmic | 0.04 | 1 | 1 | 44 |
| Acyl-CoA dehydrogenase family member 9, mitochondrial | 0.04 | 1 | 2 | 15 |
| Filamin A-interacting protein 1-like | 0.04 | 1 | 0 | 11 |
| CDK5 regulatory subunit-associated protein 3 | 0.04 | 1 | 0 | 17 |
| DNA polymerase subunit gamma-1 | 0.04 | 1 | 5 | 40 |
| Conserved oligomeric Golgi complex subunit 1 | 0.04 | 1 | 0 | 39 |
| Dual specificity testis-specific protein kinase 1 | 0.04 | 1 | 5 | 25 |
| Ubiquitin carboxyl-terminal hydrolase 30 | 0.04 | 1 | 0 | 32 |
| DNA excision repair protein ERCC-8 | 0.04 | 1 | 0 | 7 |
| Cytochrome P450 2B19 | 0.04 | 2 | 1 | 84 |
| Coiled-coil and C2 domain-containing protein 1B | 0.04 | 2 | 1 | 11 |
| RRP12-like protein | 0.04 | 1 | 0 | 35 |
| Leucine-rich repeat-containing protein 28 | 0.04 | 1 | 0 | 16 |
| Cbp/p300-interacting transactivator 4 | 0.04 | 1 | 0 | 14 |
| Microtubule-associated protein 4 | 0.04 | 1 | 0 | 18 |
| Inactive serine protease PAMR1 | 0.04 | 1 | 1 | 14 |
| Zinc transporter 5 | 0.04 | 1 | 0 | 21 |
| Acyl-CoA desaturase 2 | 0.03 | 3 | 1 | 55 |
| Carcinoembryonic antigen-related cell adhesion molecule 1 | 0.03 | 3 | 1 | 45 |
| Interferon-induced protein with tetratricopeptide repeats 1 | 0.03 | 1 | 1 | 48 |
| Interferon-induced protein with tetratricopeptide repeats 1 | 0.03 | 1 | 1 | 55 |
| Cip1-interacting zinc finger protein | 0.03 | 2 | 0 | 22 |
| Valine--tRNA ligase, mitochondrial | 0.03 | 2 | 4 | 73 |
| Type 2 lactosamine alpha-2,3-sialyltransferase | 0.03 | 1 | 0 | 14 |
| Microtubule-associated protein 1A | 0.03 | 1 | 0 | 24 |
| Bifunctional epoxide hydrolase 2 | 0.03 | 1 | 2 | 17 |
| Methylcrotonoyl-CoA carboxylase subunit alpha, mitochondrial | 0.03 | 1 | 0 | 36 |

| | | | | |
|---|---|---|---|---|
| NADH dehydrogenase [ubiquinone] 1 beta subcomplex subunit 8, mitochondrial | 0.03 | 1 | 0 | 15 |
| Catenin alpha-1 | 0.03 | 1 | 0 | 41 |
| WSC domain-containing protein 2 | 0.03 | 1 | 0 | 61 |
| Adenylate cyclase type 5 | 0.03 | 1 | 1 | 12 |
| Serine/threonine-protein kinase RIO1 | 0.03 | 1 | 3 | 32 |
| Tripartite motif-containing protein 16 | 0.03 | 1 | 1 | 24 |
| Eukaryotic translation initiation factor 3 subunit A | 0.03 | 1 | 0 | 12 |
| Unconventional myosin-X | 0.03 | 1 | 3 | 69 |
| Mitotic spindle assembly checkpoint protein MAD1 | 0.03 | 1 | 0 | 14 |
| Interferon-induced guanylate-binding protein 1 | 0.03 | 1 | 1 | 10 |
| Tetratricopeptide repeat protein 32 | 0.03 | 1 | 1 | 5 |
| Complement component C7 | 0.03 | 1 | 2 | 37 |
| snRNA-activating protein complex subunit 2 | 0.03 | 1 | 0 | 2 |
| Ephexin-1 | 0.03 | 1 | 0 | 16 |
| GTPase IMAP family member 5 | 0.03 | 1 | 0 | 35 |
| A-kinase anchor protein 2 | 0.03 | 1 | 0 | 25 |
| Centrosomal protein of 44 kDa | 0.03 | 1 | 0 | 16 |
| Protein sprouty homolog 2 | 0.03 | 1 | 0 | 22 |
| Tripartite motif-containing protein 34A | 0.03 | 1 | 0 | 7 |
| Sialoadhesin | 0.03 | 1 | 0 | 6 |
| Ectonucleotide pyrophosphatase/phosphodiesterase family member 2 | 0.03 | 1 | 0 | 23 |
| Macrophage colony-stimulating factor 1 receptor | 0.03 | 1 | 0 | 16 |
| Isoaspartyl peptidase/L-asparaginase | 0.03 | 1 | 0 | 4 |
| DnaJ homolog subfamily C member 16 | 0.03 | 1 | 0 | 29 |
| Fatty-acid amide hydrolase 1 | 0.03 | 1 | 2 | 32 |
| Aminopeptidase N | 0.03 | 1 | 0 | 98 |
| Autophagy-related protein 2 homolog A | 0.03 | 1 | 0 | 10 |
| Acyl-CoA dehydrogenase family member 10 | 0.03 | 1 | 0 | 2 |
| 2-oxoglutarate and iron-dependent oxygenase domain-containing protein 3 | 0.03 | 1 | 1 | 25 |
| Cordon-bleu protein-like 1 | 0.02 | 2 | 0 | 18 |
| Fatty acid desaturase 1 | 0.02 | 2 | 2 | 35 |
| Caskin-2 | 0.02 | 2 | 0 | 30 |
| Histone H1.3 | 0.02 | 1 | 1 | 28 |
| Protein phosphatase 1 regulatory subunit 21 | 0.02 | 1 | 0 | 36 |
| Heat shock 70 kDa protein 4L | 0.02 | 1 | 0 | 5 |
| Forkhead box protein O4 | 0.02 | 1 | 0 | 2 |
| Phospholipase D3 | 0.02 | 1 | 0 | 9 |
| Probable leucine--tRNA ligase, mitochondrial | 0.02 | 1 | 0 | 48 |
| UPF0568 protein C14orf166 homolog | 0.02 | 1 | 0 | 5 |
| FAST kinase domain-containing protein 3 | 0.02 | 1 | 0 | 22 |
| RILP-like protein 1 | 0.02 | 1 | 3 | 16 |
| Disintegrin and metalloproteinase domain-containing protein 9 | 0.02 | 1 | 0 | 23 |
| Sarcosine dehydrogenase, mitochondrial | 0.02 | 1 | 1 | 54 |
| Multiple inositol polyphosphate phosphatase 1 | 0.02 | 1 | 1 | 32 |

| | | | | |
|---|---|---|---|---|
| Nicotinamide mononucleotide adenylyltransferase 1 | 0.02 | 1 | 1 | 23 |
| Bestrophin-3 | 0.02 | 1 | 0 | 3 |
| RNA polymerase II-associated protein 3 | 0.02 | 1 | 0 | 19 |
| L-selectin | 0.02 | 1 | 0 | 20 |
| Rho guanine nucleotide exchange factor 7 | 0.02 | 1 | 0 | 6 |
| Multidrug resistance-associated protein 6 | 0.02 | 1 | 1 | 71 |
| Urocanate hydratase | 0.02 | 1 | 0 | 30 |
| WD repeat and FYVE domain-containing protein 3 | 0.02 | 1 | 0 | 34 |
| H/ACA ribonucleoprotein complex subunit 4 | 0.02 | 1 | 0 | 20 |
| Activity-dependent neuroprotector homeobox protein | 0.02 | 1 | 0 | 10 |
| Aprataxin | 0.02 | 1 | 0 | 36 |
| Putative deoxyribonuclease TATDN2 | 0.02 | 1 | 0 | 4 |
| TBC1 domain family member 9B | 0.02 | 1 | 1 | 50 |
| Lipoprotein lipase | 0.02 | 1 | 1 | 68 |
| SWI/SNF-related matrix-associated actin-dependent regulator of chromatin subfamily A member 5 | 0.02 | 1 | 0 | 31 |
| Rho guanine nucleotide exchange factor 37 | 0.02 | 1 | 0 | 33 |
| Growth hormone receptor | 0.02 | 1 | 1 | 55 |
| Diphosphoinositol polyphosphate phosphohydrolase 2 | 0.02 | 1 | 0 | 5 |
| Pecanex-like protein 3 | 0.02 | 1 | 0 | 8 |
| Probable E3 ubiquitin-protein ligase HERC1 | 0.02 | 1 | 0 | 26 |
| Cytochrome c oxidase assembly protein COX14 | 0.01 | 3 | 0 | 7 |
| Complement component C9 | 0.01 | 3 | 1 | 52 |
| 14-3-3 protein theta | 0.01 | 2 | 0 | 33 |
| 2-methoxy-6-polyprenyl-1,4-benzoquinol methylase, mitochondrial | 0.01 | 2 | 0 | 18 |
| UDP-glucuronosyltransferase 3A2 | 0.01 | 2 | 0 | 32 |
| Receptor-transporting protein 3 | 0.01 | 2 | 0 | 30 |
| Macrophage-expressed gene 1 protein | 0.01 | 1 | 0 | 60 |
| DnaJ homolog subfamily C member 24 | 0.01 | 1 | 0 | 1 |
| Fibronectin type III and SPRY domain-containing protein 2 | 0.01 | 1 | 6 | 52 |
| Ser/Thr-rich protein T10 in DGCR region | 0.01 | 1 | 0 | 6 |
| Maleylacetoacetate isomerase | 0.01 | 1 | 0 | 27 |
| Serine/threonine-protein kinase Nek3 | 0.01 | 1 | 1 | 16 |
| Arylacetamide deacetylase | 0.01 | 1 | 0 | 9 |
| Activator of basal transcription 1 | 0.01 | 1 | 0 | 18 |
| Ubiquinone biosynthesis protein COQ9, mitochondrial | 0.01 | 1 | 0 | 46 |
| Transmembrane protein 106A | 0.01 | 1 | 0 | 19 |
| Glutamyl-tRNA(Gln) amidotransferase subunit A, mitochondrial | 0.01 | 1 | 1 | 47 |
| Heat shock protein HSP 90-alpha | 0.01 | 1 | 0 | 19 |
| SAYSvFN domain-containing protein 1 | 0.01 | 1 | 0 | 9 |
| Lipase maturation factor 2 | 0.01 | 1 | 0 | 36 |
| Lysosomal protective protein | 0.01 | 1 | 0 | 34 |
| Ubiquitin-protein ligase E3C | 0.01 | 1 | 0 | 6 |
| Acyl-CoA desaturase 1 | 0.01 | 1 | 0 | 45 |
| Protocadherin-12 | 0.01 | 1 | 0 | 35 |
| Platelet-activating factor acetylhydrolase 2, cytoplasmic | 0.01 | 1 | 1 | 22 |

| | | | | |
|---|---|---|---|---|
| Protein PRRC2B | 0.01 | 1 | 0 | 8 |
| Voltage-dependent calcium channel gamma-like subunit | 0.01 | 1 | 0 | 11 |
| Bifunctional purine biosynthesis protein PURH | 0.01 | 1 | 0 | 38 |
| Sulfated glycoprotein 1 | 0.01 | 1 | 0 | 12 |
| CD5 antigen-like | 0.01 | 1 | 0 | 20 |
| Solute carrier family 2, facilitated glucose transporter member 2 | 0.01 | 1 | 0 | 56 |
| HLA class I histocompatibility antigen, B-37 alpha chain | 0.01 | 1 | 0 | 3 |
| Semaphorin-4G | 0.01 | 1 | 0 | 27 |
| Interferon-activable protein 204 | 0.01 | 1 | 0 | 80 |
| Nucleolar complex protein 4 homolog | 0.01 | 1 | 0 | 5 |
| Probable cation-transporting ATPase 13A2 | 0.01 | 1 | 0 | 19 |
| Conserved oligomeric Golgi complex subunit 2 | 0.01 | 1 | 1 | 6 |
| Plectin | 0.01 | 1 | 0 | 34 |
| Leucine-rich repeat flightless-interacting protein 2 | 0.01 | 1 | 0 | 13 |
| Laminin subunit beta-2 | 0.01 | 1 | 0 | 20 |
| Vacuolar-sorting protein SNF8 | 0.01 | 1 | 0 | 5 |
| Glutathione S-transferase Mu 1 | 0.01 | 1 | 0 | 25 |
| Sodium-dependent phosphate transport protein 3 | 0.01 | 1 | 0 | 20 |
| GH3 domain-containing protein | 0.01 | 1 | 0 | 7 |
| Aldo-keto reductase family 1 member B15 | 0.01 | 1 | 0 | 4 |
| Testican-2 | 0.01 | 1 | 1 | 7 |
| Tetratricopeptide repeat protein 33 | 0.01 | 1 | 0 | 45 |

**Tab. S2** List of annotated genes from liver samples with significantly differentiated expression between selected and control lines. Overexpression in selected lines is denoted by negative values in the "Fold Change" column. Expression level values are given in log CPM (count per million reads). FDR (False Discovery Rate) calculated according to Benjamini and Hochberg 1995. The second part of the table represent results of additional verification of these genes: P values of t test performed on FPKM and whether genes have non-overlapping FPKM values (+/- means yes/ no).

| Gene name | edgeR results | | | | Tests on FPKM | |
|---|---|---|---|---|---|---|
| | Fold Change | Expression level | P value | FDR | P value (t test) | diff Stat > 0 |
| Retinoblastoma-like protein 2 | -11.90 | 1.65 | 3.9E-21 | 3.5E-16 | 7.08E-02 | + |
| Heterogeneous nuclear ribonucleoprotein H2 | 9.50 | -0.32 | 1.7E-20 | 7.6E-16 | 3.00E-02 | + |
| Methyl-CpG-binding domain protein 4 | 7.41 | -2.63 | 3.9E-10 | 2.1E-06 | 6.63E-04 | + |
| H-2 class II histocompatibility antigen, I-E beta chain | 3.95 | -0.89 | 1.1E-09 | 4.9E-06 | 1.84E-02 | + |
| Inactive serine protease 35 | 1.80 | 1.02 | 1.2E-09 | 5.1E-06 | 4.12E-03 | + |
| Apolipoprotein A-I | 1.29 | 9.03 | 9.7E-09 | 3.4E-05 | 1.02E-02 | + |
| Cytochrome P450 4A14 | -13.99 | 4.08 | 1.8E-08 | 5.5E-05 | 3.56E-01 | - |
| Protein FAM92A1 | 7.63 | -2.39 | 4.4E-08 | 1.3E-04 | 3.68E-02 | - |
| Heat shock protein 105 kDa | -2.16 | 1.57 | 9.4E-08 | 2.5E-04 | 7.71E-02 | + |
| Cytochrome P450 4A12 | -1.18 | 3.74 | 9.9E-08 | 2.5E-04 | 4.48E-04 | + |
| Testosterone 17-beta-dehydrogenase 3 | -1.20 | 3.21 | 1.3E-07 | 2.9E-04 | 8.60E-04 | + |
| Methylmalonate-semialdehyde dehydrogenase acylating, mitochondrial | -9.75 | 4.22 | 1.7E-07 | 3.5E-04 | 3.56E-01 | - |
| Intraflagellar transport protein 122 homolog | 2.34 | -0.16 | 1.8E-07 | 3.6E-04 | 2.65E-02 | + |
| Serotransferrin-B | 10.46 | -0.07 | 1.8E-07 | 3.7E-04 | 3.56E-01 | - |
| Zinc finger protein 3 | -4.01 | -2.87 | 2.0E-07 | 3.9E-04 | 6.17E-04 | + |
| HERV-K_22q11.21 provirus ancestral Gag polyprotein | -1.29 | 1.52 | 2.2E-07 | 4.2E-04 | 6.13E-03 | + |
| Probable ATP-dependent RNA helicase DHX37 | 2.07 | -0.54 | 2.9E-07 | 5.0E-04 | 1.34E-02 | + |
| Uncharacterized protein C1orf129 | -1.45 | 2.04 | 3.3E-07 | 5.6E-04 | 6.33E-04 | + |
| Serum albumin B | 10.06 | -0.46 | 3.5E-07 | 5.7E-04 | 3.56E-01 | - |
| Serum albumin B | 9.83 | -0.69 | 5.3E-07 | 8.3E-04 | 3.56E-01 | - |
| Heat shock 70 kDa protein 1A/1B | -1.50 | 3.61 | 8.0E-07 | 1.1E-03 | 5.61E-02 | + |
| Uncharacterized protein C1orf129 | -2.23 | -1.60 | 9.3E-07 | 1.2E-03 | 2.86E-04 | + |
| Glyceraldehyde-3-phosphate dehydrogenase | 9.33 | -1.18 | 1.5E-06 | 1.8E-03 | 3.56E-01 | - |
| Fibroblast growth factor 21 | 2.29 | -0.11 | 2.4E-06 | 2.6E-03 | 5.58E-02 | - |

| | | | | | | |
|---|---|---|---|---|---|---|
| Acyl-coenzyme A amino acid N-acyltransferase 2 | -1.12 | 7.69 | 2.9E-06 | 3.0E-03 | 1.98E-03 | + |
| Suppressor of cytokine signaling 2 | -1.39 | 3.11 | 3.8E-06 | 3.6E-03 | 2.74E-02 | + |
| Dedicator of cytokinesis protein 5 | 7.69 | -2.63 | 4.1E-06 | 3.7E-03 | 1.40E-01 | - |
| Cytochrome P450 4A10 | -0.94 | 5.53 | 5.9E-06 | 4.9E-03 | 1.33E-03 | + |
| Hemoglobin subunit beta-1 | 8.74 | -1.75 | 5.9E-06 | 4.9E-03 | 3.56E-01 | - |
| Alpha-1-antiproteinase 2 | 8.68 | -1.80 | 6.9E-06 | 5.6E-03 | 3.56E-01 | - |
| Cathepsin K | -3.77 | -2.48 | 8.5E-06 | 6.5E-03 | 4.69E-02 | - |
| Poly ADP-ribose polymerase 14 | -1.28 | 3.34 | 1.0E-05 | 7.7E-03 | 5.64E-02 | + |
| Apolipoprotein A-II | 8.54 | -1.94 | 1.0E-05 | 7.7E-03 | 3.56E-01 | - |
| Kynureninase | -0.85 | 8.40 | 1.1E-05 | 7.8E-03 | 2.81E-03 | + |
| Alpha-fetoprotein | 8.51 | -1.97 | 1.1E-05 | 7.9E-03 | 3.56E-01 | - |
| Interferon-induced helicase C domain-containing protein 1 | -1.85 | -1.12 | 1.2E-05 | 8.1E-03 | 6.63E-04 | + |
| Alpha-actinin-4 | -2.43 | 1.10 | 1.5E-05 | 9.6E-03 | 2.65E-02 | - |
| Retrovirus-related Pol polyprotein LINE-1 | -1.05 | 1.15 | 1.6E-05 | 9.8E-03 | 2.83E-03 | + |
| Elongation factor 1-alpha, oocyte form | 8.01 | -2.42 | 1.7E-05 | 1.0E-02 | 3.37E-01 | - |
| Perilipin-2 | 0.94 | 6.83 | 1.7E-05 | 1.0E-02 | 2.00E-02 | + |
| Ras-related protein Rab-30 | 1.01 | 2.94 | 1.8E-05 | 1.1E-02 | 2.01E-03 | + |
| Hemoglobin subunit alpha-1 | 8.31 | -2.16 | 1.9E-05 | 1.1E-02 | 3.56E-01 | - |
| Adenylate kinase isoenzyme 4, mitochondrial | 1.09 | 3.66 | 1.9E-05 | 1.2E-02 | 4.25E-02 | - |
| Mucolipin-1 | 1.34 | 0.04 | 2.1E-05 | 1.2E-02 | 2.36E-03 | + |
| Protein GPR108 | 2.02 | -0.68 | 2.3E-05 | 1.3E-02 | 2.63E-02 | - |
| EMI domain-containing protein 1 | 6.75 | -3.19 | 2.4E-05 | 1.4E-02 | 1.16E-01 | - |
| Leucine-rich repeat and transmembrane domain-containing protein 1 | 2.09 | 0.55 | 2.5E-05 | 1.4E-02 | 5.37E-02 | - |
| Aryl hydrocarbon receptor nuclear translocator 2 | -1.27 | 2.84 | 2.5E-05 | 1.4E-02 | 5.58E-02 | - |
| Elongation factor 1-alpha, oocyte form | 6.22 | -1.89 | 2.6E-05 | 1.4E-02 | 3.48E-01 | - |
| Potassium-transporting ATPase alpha chain 1 | -2.59 | 0.68 | 2.8E-05 | 1.5E-02 | 1.29E-01 | - |
| Pentraxin fusion protein | 8.11 | -2.34 | 3.2E-05 | 1.7E-02 | 3.56E-01 | - |
| Cytosolic phospholipase A2 gamma | -0.94 | 3.71 | 3.3E-05 | 1.7E-02 | 9.71E-03 | + |
| BSD domain-containing protein 1 | 1.01 | 2.95 | 3.9E-05 | 1.9E-02 | 8.03E-03 | + |
| EH domain-containing protein 1 | -7.12 | -2.99 | 3.9E-05 | 1.9E-02 | 1.59E-01 | - |
| Tenascin-X | -0.95 | 3.24 | 4.3E-05 | 2.1E-02 | 6.32E-03 | + |
| Intraflagellar transport protein | 1.20 | 2.14 | 4.5E-05 | 2.1E-02 | 2.40E-02 | + |

| | | | | | | |
|---|---|---|---|---|---|---|
| 122 homolog | | | | | | |
| Ropporin-1-like protein | -2.64 | -1.54 | 4.6E-05 | 2.1E-02 | 1.62E-01 | + |
| Lysine-specific demethylase 2A | 7.99 | -2.47 | 4.6E-05 | 2.2E-02 | 3.56E-01 | - |
| Glutathione S-transferase theta-1 | 0.70 | 8.05 | 5.3E-05 | 2.4E-02 | 3.21E-03 | + |
| Kielin/chordin-like protein | 1.75 | -1.01 | 5.3E-05 | 2.4E-02 | 3.38E-03 | + |
| Retinol-binding protein 4 | 7.92 | -2.53 | 5.6E-05 | 2.5E-02 | 3.56E-01 | - |
| NADP-dependent malic enzyme | 0.89 | 7.17 | 5.8E-05 | 2.5E-02 | 3.53E-02 | - |
| Meiosis-specific with OB domain-containing protein | 2.38 | -1.28 | 5.9E-05 | 2.6E-02 | 6.60E-02 | - |
| Homeobox protein cut-like 2 | -1.09 | 0.92 | 6.4E-05 | 2.7E-02 | 8.47E-03 | + |
| Protein FAM92A1 | 6.32 | -3.46 | 6.4E-05 | 2.7E-02 | 1.41E-01 | - |
| Aphrodisin | -0.84 | 13.98 | 6.6E-05 | 2.7E-02 | 2.12E-02 | + |
| Suppressor of cytokine signaling 2 | -1.61 | 0.73 | 8.3E-05 | 3.3E-02 | 1.07E-02 | + |
| Diacylglycerol O-acyltransferase 2-like protein 6 | 3.13 | -2.69 | 8.4E-05 | 3.3E-02 | 4.06E-03 | + |
| Cytochrome P450 4A2 | -1.11 | 3.90 | 8.4E-05 | 3.3E-02 | 2.99E-02 | - |
| Transcription factor E2F2 | -0.94 | 1.68 | 9.2E-05 | 3.6E-02 | 9.41E-03 | + |
| Fatty acid-binding protein 2, liver | 7.74 | -2.70 | 9.7E-05 | 3.7E-02 | 3.56E-01 | - |
| UDP-glucuronosyltransferase 3A2 | -0.90 | 3.21 | 9.9E-05 | 3.7E-02 | 7.40E-03 | + |
| Protein phosphatase 1 regulatory subunit 3G | -1.31 | 2.23 | 1.0E-04 | 3.7E-02 | 1.12E-01 | - |
| Bcl-2-associated transcription factor 1 | -7.49 | -2.59 | 1.1E-04 | 4.1E-02 | 3.56E-01 | - |
| Zinc finger protein 226 | 2.83 | -1.40 | 1.1E-04 | 4.1E-02 | 8.71E-02 | - |
| Heat shock 70 kDa protein 1A/1B | -1.28 | 3.92 | 1.2E-04 | 4.2E-02 | 1.06E-01 | - |
| Fatty acid desaturase 2 | 0.78 | 9.68 | 1.2E-04 | 4.2E-02 | 3.02E-02 | + |
| Rho guanine nucleotide exchange factor 4 | 1.28 | 0.34 | 1.2E-04 | 4.3E-02 | 1.68E-02 | - |
| Rho guanine nucleotide exchange factor 4 | 3.07 | -1.48 | 1.3E-04 | 4.5E-02 | 8.33E-02 | - |
| HEAT repeat-containing protein 4 | -6.21 | -3.80 | 1.4E-04 | 4.7E-02 | 4.78E-02 | - |
| Insulin-like growth factor-binding protein 2 | 2.16 | 0.17 | 1.5E-04 | 4.9E-02 | 2.15E-01 | - |
| Zinc finger and BTB domain-containing protein 47 | 0.90 | 1.30 | 1.5E-04 | 4.9E-02 | 7.84E-04 | + |

**Tab. S3** List of annotated genes from heart samples with significantly differentiated expression between selected and control lines. Overexpression in selected lines is denoted by negative values in the "Fold Change" column. Expression level values are given in log CPM (count per million reads). FDR (False Discovery Rate) calculated according to Benjamini and Hochberg 1995. The second part of the table represent results of additional verification of these genes: P values of t test performed on FPKM and whether genes have non-overlapping FPKM values (+/- means yes/ no).

| Gene name | edgeR results | | | | Tests on FPKM | |
|---|---|---|---|---|---|---|
| | Fold Change | Expression level | P value | FDR | P value (t test) | diffStat > 0 |
| Aphrodisin | -2.63 | 0.85 | 1.2E-14 | 1.3E-09 | 2.25E-02 | + |
| Centromere protein S | -2.88 | -0.75 | 7.5E-14 | 4.1E-09 | 1.33E-03 | + |
| Synaptonemal complex protein 3 | -3.21 | -0.90 | 6.8E-11 | 1.5E-06 | 1.56E-02 | + |
| Fetuin-B | -1.94 | 0.03 | 1.7E-09 | 2.3E-05 | 9.36E-03 | + |
| Interferon-induced guanylate-binding protein 2 | -8.73 | -1.41 | 3.6E-09 | 4.3E-05 | 2.35E-01 | - |
| Transmembrane protein 114 | -1.48 | 0.73 | 3.9E-09 | 4.3E-05 | 1.07E-03 | + |
| Aphrodisin | -2.37 | -0.03 | 9.1E-09 | 8.3E-05 | 3.32E-02 | + |
| Heat shock-related 70 kDa protein 2 | -0.76 | 3.04 | 2.3E-07 | 1.6E-03 | 1.94E-05 | + |
| Arachidonate 12-lipoxygenase, leukocyte-type | -3.53 | 3.21 | 2.6E-07 | 1.7E-03 | 2.88E-01 | - |
| Apolipoprotein L3 | -8.12 | -2.02 | 2.8E-07 | 1.7E-03 | 2.65E-01 | - |
| Semaphorin-4F | -1.31 | 0.29 | 3.5E-07 | 2.0E-03 | 1.81E-02 | + |
| Serum albumin | -1.79 | 2.73 | 3.8E-07 | 2.1E-03 | 8.04E-02 | + |
| Geminin coiled-coil domain-containing protein 1 | 8.46 | -1.40 | 1.4E-06 | 5.3E-03 | 3.47E-01 | - |
| Rhophilin-2 | -0.95 | 1.92 | 2.9E-06 | 9.3E-03 | 3.45E-03 | + |
| Interferon-induced protein with tetratricopeptide repeats 1 | -3.79 | -0.48 | 3.1E-06 | 9.5E-03 | 2.92E-01 | - |
| Phenylalanine-4-hydroxylase | -3.17 | -2.24 | 5.2E-06 | 1.3E-02 | 3.62E-02 | + |
| Thyrotropin-releasing hormone receptor | 2.04 | -0.58 | 6.2E-06 | 1.5E-02 | 4.27E-02 | - |
| Ubiquitin-like protein ISG15 | -4.05 | 2.38 | 6.3E-06 | 1.5E-02 | 3.51E-01 | - |
| Guanylate-binding protein 4 | -2.30 | 1.23 | 7.5E-06 | 1.7E-02 | 2.36E-01 | + |
| Apolipoprotein C-III | -2.97 | -2.72 | 8.3E-06 | 1.8E-02 | 4.09E-04 | + |
| Cytochrome P450 2A3 | -2.11 | -0.39 | 1.0E-05 | 2.0E-02 | 1.14E-01 | - |
| Immunoglobulin lambda-like polypeptide 1 | 3.07 | -1.90 | 1.3E-05 | 2.2E-02 | 4.69E-02 | - |
| Interferon-induced protein with tetratricopeptide repeats 3 | -2.02 | 1.68 | 1.2E-05 | 2.2E-02 | 1.51E-01 | + |
| Neuron-specific protein family member 2 | -1.18 | 4.01 | 1.3E-05 | 2.2E-02 | 4.51E-02 | + |
| Soluble lamin-associated protein of 75 kDa | -2.04 | -1.58 | 1.6E-05 | 2.6E-02 | 1.20E-02 | - |
| 2'-5'-oligoadenylate synthase- | -3.38 | 1.52 | 2.0E-05 | 3.2E-02 | 3.39E-01 | - |

like protein 1

| | | | | | | |
|---|---|---|---|---|---|---|
| Beta-2-glycoprotein 1 | -2.55 | -2.13 | 3.0E-05 | 4.4E-02 | 4.94E-02 | + |
| Retinaldehyde-binding protein 1 | 5.33 | 1.09 | 3.3E-05 | 4.6E-02 | 3.67E-01 | - |

# CHAPTER III

## Molecular basis of predatory behavior in bank voles

M. Konczal, W. Babik, P. Orłowska-Feuer, J. Radwan, E.T. Sadowska,

P. Koteja

**Abstract**

Changes in foraging strategies are thought to be essential in many adaptive radiations, yet little is known about the nature of genetic variation underlying such behaviors. Here, we investigated genetic basis of the response to selection for predatory behavior using a unique laboratory model of vertebrate adaptive radiation. After 13 generations of selection for increased frequency of predatory behavior, the proportion of bank voles (*Myodes [=Clethrionomys] glareolus*) showing predatory behavior was 5 times higher in 4 replicated selected lines than in 4 control lines. We used RNAseq to analyze hippocampus and liver transcriptomes of voles from the selected and control lines and found that selection resulted in repeatable changes in allele frequencies and gene expression. We did not however find evidence for the role of nonsynonymous polymorphisms in response to selection, suggesting that most single nucleotide polymorphisms (SNPs) which did respond, affect gene expression or alternative splicing. Expression analyses showed that substantial number of genes were differentially expressed between treatments (149). Together with results of the SNP analyses, this suggests that the initial response to selection for predatory behavior appears strongest in regulatory regions of the genome, what support hypothesis that changes in gene expression play predominant role at early stages adaptive evolution. Repeatable changes in SNP allele frequencies, the pattern not observed in similar analyses of voles selected for maximum metabolism rate, suggests that architecture of adaptive variation differs between traits. We suspect that alleles at higher frequencies or of larger effects have responded to selection for predatory behavior than for aerobic performance, which caused higher repeatability between replicates. Finally, we characterized genes with the largest differences between predatory and control lines. They are associated with hunger, aggression, biological rhythm and functioning of the nervous system.

**Introduction**

Many different strategies have been used to study the molecular basis of adaptive changes (Stapley et al. 2010). We are however still far from reaching a consensus about the role of natural selection vs. other processes in shaping genetic variation (Hahn 2008; Wagner 2008; Sella et al. 2009; Nei 2013). Important issues concerning the process of adaptation, such as the role of standing genetic variation, changes in gene expression or effect sizes of adaptive variants, remain controversial as well (Orr 2005, Barrett and Schluter 2008, Rockman 2012, Fraser 2013). Hence a major research program, in which researchers attempt to decipher the genetic architecture of adaptive traits, to establish links between genotype, phenotype and fitness and test hypotheses about the action of selection at the genomic level (Ellegren and Sheldon 2008, Dalziel et al. 2009, Radwan and Babik 2012, but see Travisano and Shawn 2013). Such studies contributed to better understanding of the molecular basis of morphological (e.g. height in humans, armor plates in sticklebacks, fur color in deer mice [Hoekstra et al. 2006, Frazer et al. 2009, Jones et al. 2012]) or physiological (performance, adaptation to high altitude [Storz et al. 2007, Yi et al. 2010, Konczal et al. 2015]) traits but we still know little about the genomic architecture of natural variation in complex behaviors (Boake et al. 2002, Bendesky and Bargmann 2011, Weber et al. 2013). It is unfortunate because during the initial phase of adaptation, response to selection may be primarily observed in behavioral traits (Mayr 1959, Blomberg et al. 2003, Garland and Rose 2009). Identification of genomic basis of variation in such traits (both underlying genes and genetic architecture) is therefore crucial for mechanistic and evolutionary understanding of behavioral adaptation.

From ecological and evolutionary perspective, one of the most intriguing behaviors is predation (Curio 1976, Barbosa and Castellanos 2005, Ishii and Shimada 2010, Ritchie et al. 2012). Predation is an ecological factor of almost universal importance for regulating ecosystems and sustaining biodiversity (Ritchie and Johnson 2009). At the organismal level it is associated with hunger, activity, searching behavior and prey recognition/selection (Curio 1976). Motivation and ability for predation may have serious consequences for survival and reproductive success (Eisenberg and Leyhausen 1972, Curio 1976), but the genomic basis of variation in predatory behavior is largely unknown. Sequenced genomes of some predatory species revealed genome-wide signs of positive selection, however the identification of genes underlying predatory behavior *per se* is next to impossible using such comparative genomic approach (Zhan et al. 2013, Cho et al. 2013).

Another approach however, experimental evolution, allows to identify the genotype-phenotype link, to test the role of selection in shaping genetic variation and to study evolution in real time (Garland and Rose 2009, Kawecki et al. 2012). The advantage of experimental evolution is the ability to focus on a well-defined selected trait, which is measured in experimentally controlled ambient conditions. Direct and correlated responses to selection can then be distinguished and separated (Garland and Rose 2009). Well-designed selection experiments minimize also the effect of complex demography and other historical factors, which may elevate the rate of false positives in scans for genomic signatures of adaptation in natural populations (Akey 2009).

In this study we employ experimental evolution and high throughput sequencing to get insight into the early stages of predatory behavior evolution. Specifically, we investigate the transcriptome-wide response to artificial selection for increased predatory behavior in a small mammal in an attempt to better understand the genomic architecture of variation affecting this trait. We analyzed bank vole (*Myodes [=Clethrionomys] glareolus*) selection experiment, with four lines selected for predatory behavior (referred to as selected lines or P lines) and four control lines (C lines; Sadowska at al. 2008). Bank voles are omnivore rodents with a very diverse diet that includes both invertebrates and green plant tissues, allowing selection to operate on feeding-related traits (Petrusewicz 1983; Wereszczyńska et al. 2007; Sadowska et al. 2008). In the experiment selection was applied based on intensity of predatory behavior towards crickets. This trait was measured as the time to catch a live cricket which was placed in a cage together with a fasted vole as described by Sadowska et al. (2008). After 13 generations of selection, the proportion of voles attacking crickets was 5 times higher in the selected P lines than in unselected control lines (Chrząścik et al. 2014).

Transcriptomes of two tissues – liver and brain (hippocampus) – were sequenced and compared between the selected and control lines using a cost-effective pooled RNA-Seq approach (Konczal et al. 2014). We focus not only on gene expression changes, but also on differences in allele frequencies at expressed parts of the genome to assess the impact of selection on genome-wide patterns of polymorphism and divergence. Such approach is feasible because transcriptome represents substantial part of the functional genome and RNA-seq is thus a convenient way to perform genome-wide studies of experimental evolution for species without the available reference genome (Konczal et al. 2014). Hippocampus was chosen, as a brain structure involved in feeding-related behaviors (Tracy et al. 2001) and as the structure where motivation and memory are coordinated to guide behavior (Kennedy and

Shapiro 2009). We also analyzed the liver transcriptome, because many genes expressed in this organ provide information about allele frequency changes in loci which are not expressed in the hippocampus. Moreover, if evolution of predatory behavior is accompanied by changes in overall physiology (e.g. hunger level, metabolism or stress) then gene expression changes may also be expected in liver.

Our specific questions focused on several aspects of molecular-level response to selection. First, we asked whether selection decreases genetic variation more than drift alone. Second, we were interested in repeatable changes associated with response to the same selection pressure. In particular, we wanted to compare repeatability of changes in SNPs allele frequencies with changes in genes expression. Third, we assessed whether the number of SNPs differentiated between treatments is higher than expected under drift and whether such differentiated SNPs are overrepresented at nonsynonymous sites. Fourth, we suspected that selection has affected more gene expression in hippocampus that in liver tissue. To test this hypothesis, we compared expression changes between both organs. Finally, we wanted to identify and characterize genes most differentiated between the P and C lines - candidates for molecular targets of selection for predatory behavior.

## Results

*Transcriptome sequencing, assembly and annotation*

We used 80.2 mln (M) of paired-end reads (2 x 100bp) from one control line (C3) for the reconstruction of the hippocampus transcriptome. *De novo* assembly resulted in 219,886 transcripts, which were then reduced to 153,677 transcriptome-based gene models (referred to as genes; Tab. 1). Of these 28,743 were identified as putatively protein-coding, and 21,407 (74.5%) were successfully annotated to 13,305 known genes deposited in SwissProt database (some genes were fragmented in the transcriptome assembly).

For the remaining 7 lines (3 C lines and 4 P lines) we obtained altogether 250 M of 100 bp single-end reads from the hippocampus transcriptomes (35.6 M ± (SD) 15.3 M per sample). These reads, together with subsampled sequences from the C3 line (35 M single-end reads) were used to compare hippocampus transcriptomes between the selected and control lines.

For analyses of liver transcriptomes we used the previously published liver reference transcriptome (Konczal et al. 2014, Konczal et al. 2015). We sequenced liver transcriptomes of P lines (127 M single-end 100 bp reads, 31.9 M ± (SD) 3.7 per sample) and compared them with transcriptomes of control lines from the same generation of the selection experiment (Konczal et al. 2015).

**Tab. 1** Basic statistics of the hippocampus reference transcriptome.

| | |
|---|---|
| No. of genes | 153,677 |
| No. of genes >1kb | 29,267 |
| N50 gene length | 1,598 |
| No. of genes within N50 | 18,518 |
| No. of putative protein coding genes | 28,743 |
| Total length (Mb) | 122 |

Note. – N50, 50% of the assembly length is in genes of the length of N50 bp or longer; genes, TGMs contain both coding and noncoding sequences.

*Polymorphism within the selection experiment*

Using reads from both organs we identified 179,468 SNPs, which were grouped into four classes: nonsynonymous, synonymous, UTR-located and noncoding (Tab. 2). The SNPs were localized in 15,580 genes, 11,076 of which were putatively protein coding. In accordance with expectations, allele frequency spectra differed between the classes of SNPs and they were the most skewed for nonsynonymous variants, indicating the presence of purifying selection (Fig. 1). To test an effect of selection and effective population size (calculated from the available pedigree) on genetic variation within lines we counted number of polymorphic sites (minor allele frequency > 0.05) within each line. The number of such sites was mainly affected by effective population size (p=0.037, F(1,5)=7.97, ANCOVA), while treatment (C vs. P) did not affect polymorphism (p=0.123, F(1,5)=3.45, ANCOVA).
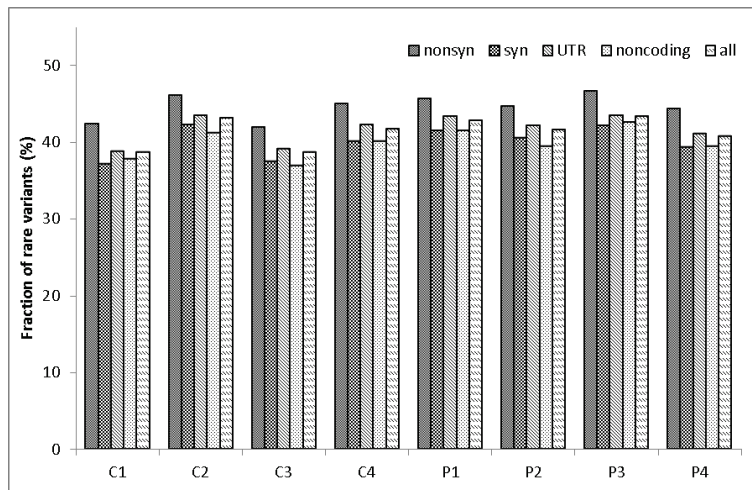
**Fig. 1** Fraction of SNPs with rare or absent alternative variant (minor allele frequency < 0.05) for a given line. Allele frequency spectra are more skewed for nonsynonymous sites than for any other class of SNPs indicating the role of purifying selection.

**Tab. 2** Number of SNPs identified in the transcriptomes of bank voles derived from selection experiment

| | |
|---|---|
| No. of SNPs | 179,468 |
| No. of genes with SNPs | 15,580 |
| No. of nonsynonymous SNPs | 21,708 |
| No. of synonymous SNPs | 44,102 |
| No. of UTR-located SNPs | 82,422 |
| No. of SNPs in noncoding genes | 31,236 |

*Repeatable allele frequency changes*

To test whether selection results in repeatable changes of allele frequencies we investigated pairwise $F_{ST}$ distances between all lines. Ordination of the matrix of mean pairwise $F_{ST}$ did not reveal any clustering of selected or control lines (Fig 2A; p=0.999; randomization test) and only separated one control line (C2) from the 7 other lines. However, when we sampled 500 SNPs with the highest mean pairwise $F_{ST}$ (i.e. showing the highest overall differentiation), the selected lines clustered apart from controls (Fig 2B; p=0.018; randomization test). This suggests, that selection causes repeatable changes in allele frequencies in selected lines. Additional evidence for this was observation, that the number of SNPs with allele frequencies non-overlapping between selected and control lines among these 500 SNPs was higher than expected by chance (p=0.02, Chi-Square test with Yates correction).
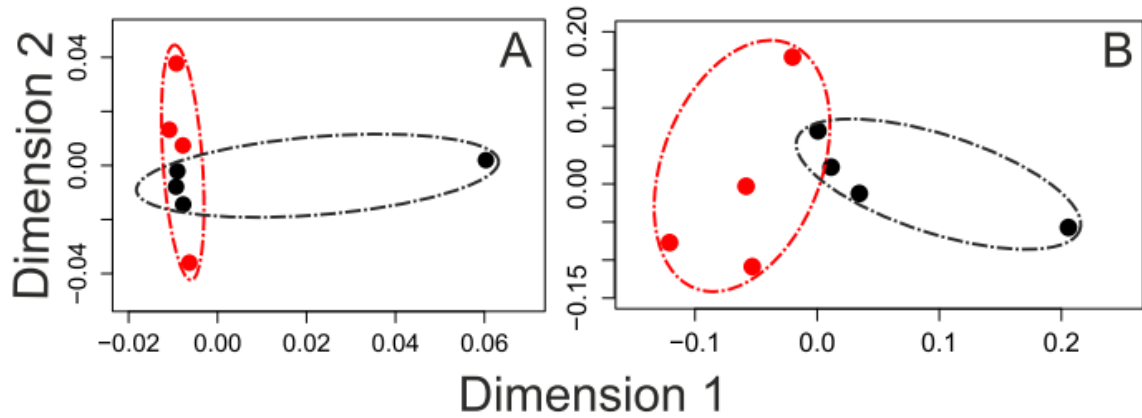
**Fig. 2** Genetic differentiation of predatory (red) and control (black) lines of the bank vole selection experiment. Multidimensional scaling was performed on the matrices of pairwise $F_{ST}$ distances between lines calculated using all SNPs (A) and for top 500 SNPs with the highest mean pairwise $F_{ST}$, i.e. showing the most overall differentiation among lines (B).

Above results provide evidence for repeatable allele frequency differentiation. To compare this effect with expectations produced from pedigree-based simulations, and to see, whether this repeatable differentiation is associated with a particular class of SNPs, such as nonsynonymous polymorphisms, we investigated SNPs which have non-overlapping allele frequencies between selected and control lines (3,715 SNPs (2.07%) located in 2,050 genes). Within this dataset we found 419 (1.93%, in 338 genes) nonsynonymous, 965 (2.19%, in 696 genes) synonymous, 1682 UTR-located (2.04 % in 962 genes) and 649 (2.08% in 407 genes) noncoding SNPs. The observed number of SNPs with non-overlapping allele frequencies was significantly higher than expected from pedigree-based simulations (p=0.02) and mostly synonymous SNPs were responsible for this effect: while the number of synonymous SNPs was higher than expected (p=0.01), the fraction of nonsynonymous (p=0.52), UTR-located (p=0.16) and non-coding (p=0.25) candidates did not differ from expectations generated by simulations.

The relatively small population sizes translate into low population recombination rate, which may cause entire long haplotypes to drift. To control for this effect we sampled one SNP per gene and compared results with drift simulations. Contrary to the analysis involving all SNPs we found slightly fewer differentiated SNPs than expected from simulations ($f_{obs}$=1.99%, $f_{exp}$=2.00% ; p=0.003, t-test). The number of differentiated SNPs was lower for SNPs localized in coding genes (nonsynonymous: $f_{obs}$=1.82%, $f_{exp}$=1.94%; p=$10^{-48}$, synonymous: $f_{obs}$=1.98% $f_{exp}$=2.03%, p=$10^{-12}$, UTR-located: $f_{obs}$=1.93%, $f_{exp}$=1.99%, p=$10^{-16}$; t-test) while it was higher in noncoding sequences ($f_{obs}$=2.15%, $f_{exp}$=2.03%, p=$10^{-33}$, t-test).

Overall, comparisons between observed and simulated allele frequencies do not provide evidence for larger than expected differentiation in nonsynonymous SNPs. Depending on the applied strategy of data analysis, noncoding or synonymous variants differed between treatments more than expected form drift alone. These SNPs may be involved in alternative splicing and regulation of gene expression or be linked to causative variants.

*Candidate genes for predatory behavior*

To identify genes most differentiated between treatments we sorted SNPs with non-overlapping allele frequencies by diffStat value. DiffStat is the difference in allele frequency between a selected line with the highest frequency and a control line with the lowest (or vice versa). The distributions of diffStat values were compared with these generated from pedigree-based simulations (Fig. 3). The observed distributions generally follow expectations, but for noncoding variants we observed larger fraction of SNPs with diffStat > 0.2, than in any other class of SNPs (Fig. 3D). Because the probability of obtaining the diffStat value > 0.2 by chance, as assessed by pedigree-based simulations was very low ($4.5 \times 10^{-4}$) therefore we considered genes harboring 94 such SNPs as promising candidates and investigated their molecular functions. Some of these genes are discussed below, and the full list is provided in Supplementary materials (Tab S1).
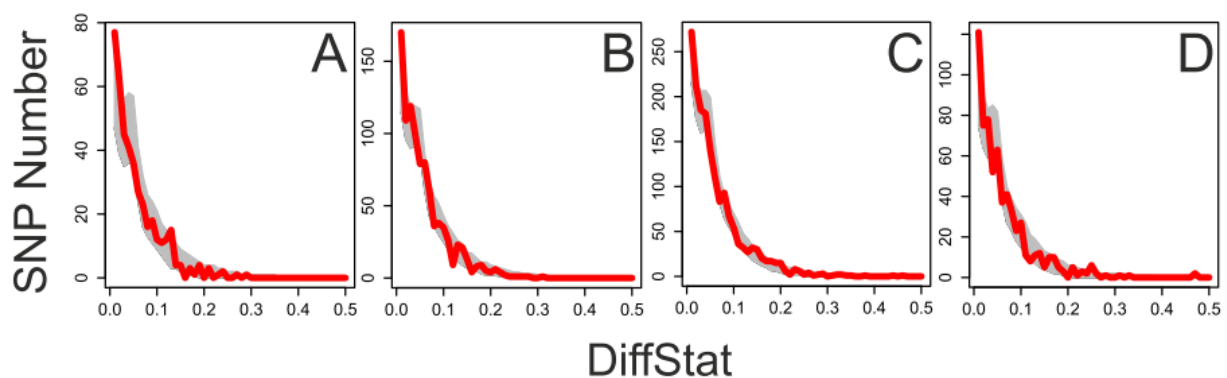


**Fig. 3** Distribution of allele frequency differences between predatory and control lines. The distributions show number of SNPs with given diffStat values. The DiffStat value is the smallest difference in allele frequency between a P- and a C- line, so for SNPs having overlapping allele frequencies between treatments diffStat =0; such SNPs are not included in the plots. Shaded area indicates 90% of all simulations, and red line represents distribution observed for nonsynonymous (A), synonymous (B), UTR-located (C) and noncoding (D) SNPs.

*Changes in gene expression*

To determine differences in gene expression level between predatory and control lines, we mapped reads from the liver and hippocampus to the respective transcriptomes, and compared expression between the P and C lines. We investigated all expressed genes with at least 10 mapped reads. Multidimensional scaling separated selected lines from controls for the hippocampus (p=0.012) but not for the liver (p=0.286). The same pattern was observed in analyses limited to 500 genes showing the highest variation among all samples (hippocampus: p=0.016; liver: p=0.289, Fig. 4). On the other hand, the number of genes with statistically significant differences (FDR < 0.05) in expression between the P and C lines was higher in the liver (90) than in the hippocampus (59) (Tab S2, Tab S3). Candidate genes potentially associated with predatory behavior are described in Discussion.
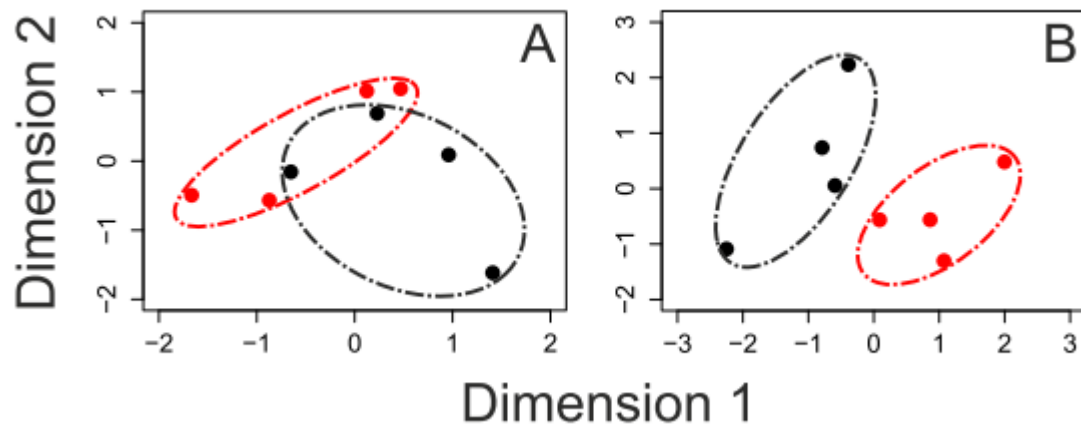


**Fig. 4** Expression differentiation of predatory (red) and control (black) lines of the bank vole selection experiment. Multidimensional scaling plots were drawn from top 500 genes with the largest variation in expression, treating all lines as a one group, for liver (A) and hippocampus (B) samples. Distances on the plot can be interpreted as leading log2-fold-change.

**Discussion**

In this study we used replicate selected and control lines derived from a natural population of an omnivore rodent to experimentally quantify the molecular level response to selection for predatory behavior. The effect of selection is manifested as consistent allele frequency and gene expression differences between the selected and control lines. The pattern of line clustering based on allele frequencies distances suggests enrichment in SNPs associated with repeatable response to selection. For some SNP classes differentiation of allele frequencies between selected and control lines was higher than expected from drift, which provided

evidence that in many cases the same alleles are selected for in the P lines. We did not find however any evidence for overrepresentation of nonsynonymous polymorphisms among selected variants or for larger than expected from drift loss of variation in selected lines. Because changes in gene expression were much more pronounced in the hippocampus than in the liver we suspect that selection is mainly associated with genes involved in functioning of the nervous system.

*Changes in allele frequencies*

Two lines of evidence indicate that selection caused repeatable allele frequency changes in the lines selected for predatory behavior. First, the P lines group separately from the C lines in ordination analysis based on the 500 SNPs most differentiated among the lines (regardless of the selective regime). Second, the number of SNPs with allele frequencies nonoverlapping between the P and C lines was higher than expected under drift; this effect was driven by synonymous polymorphisms. It is difficult to propose an explanation other than the effect of selection, for the pattern revealed by the ordination analysis. The excess of SNPs highly differentiated between the P and C lines is also suggestive, but potential limitations of the analytical approach need to be considered. The original analysis did not control for linkage between sites. When one SNP per gene was sampled, in an attempt to control for linkage, the number of differentiated SNPs higher than expected under drift was obtained only for noncoding sites, while the overall number differentiated SNPs was slightly lower than expected. However the control for linkage based on sampling one SNP per gene is far from ideal, because this approach may be affected by nonrandom distribution of differentiated SNPs across genes. If differentiated SNPs are localized in more polymorphic genes, their number will be underestimated and if they are located in less polymorphic genes their number will be overestimated.

In model species more efficient and realistic control of the effect of linkage is possible with information about haplotypes and recombination rate (Kessner et al. 2013), which is not available for our system. One potential solution in our case could be using the standard coalescent (Wakeley 2009) to simulate for each gene haplotypes in the base population. Such approach requires information about polymorphism of individual genes, which is available, and assumptions about demographic history and recombination rate (Hudson 2002). Expectations of the number of highly differentiated SNPs under neutrality can then be obtained through simulations of haplotype drift on the known pedigree. Drift

simulations could assume free recombination between and no recombination within genes, which appears reasonable given the time scale of the experiment and effective population sizes.

Overall, the comparison of SNP differentiation among lines robustly demonstrated repeatable changes in allele frequencies between selection regimes, but showed that polymorphisms affecting protein sequences are underrepresented among differentiated SNPs. Instead polymorphisms in other classes of sites, synonymous or noncoding (depending on the analytical approach used) are overrepresented, which suggests the importance of changes affecting gene expression or alternative splicing. Thus the initial response to selection for predatory behavior appears strongest in regulatory regions of the genome, supporting King and Wilson's hypothesis (King and Wilson 1975) about predominant role of gene expression changes in adaptive evolution. It is now becoming clear that selection often acts in distributed fashion on the expression of many genes, what was supported by studies performed on yeast, mice or humans (Fraser 2010, Fraser 2013, Halligan et al. 2014) as well as by analyses of bank voles selected for aerobic capacity (Konczal et al. 2015).

*Repeatable changes in gene expression*

Expression analyses showed that the overall pattern of gene expression changed in hippocampus, while we did not observe transcriptome-wide effect of selection in liver (Fig. 4). Hippocampus is a major component of the brain and it receives information from each of the sensory modalities and projects widely throughout the brain (Swanson 1983). Hippocampus also plays an important role in learning, memory, motivation and motor behavior (Morris and Hagan 1983, Tracy et al. 2001) and is one of the few regions in the adult mammalian brain that can generate new nerve cells (Gage 2000; Rhodes et al. 2003). For these reasons we suspected, that it plays an important role in response to selection for predatory behavior. The results of transcriptome-wide clustering do not mean that there is no selection-driven expression changes in liver. Contrary, we found some interesting genes with significant changed expression in this organ (see below). But the multidimensional analyses should be rather interpreted as a support for notion that aggregate effect of expression changes in many genes is more pronounced in hippocampus than in liver.

Precisely deciphering molecular changes in hippocampus requires additional studies focusing on individual data and candidate genes. First, effect sizes of candidate genes should

be evaluated at the individual level of many voles, to increase statistical power and to estimate amount of genetic variation explaining by candidate genes. Second, adaptive significance of candidate genes should be assessed in natural populations, using association studies or with the common garden experiments, to further test their role in evolution of predatory behavior. This is especially important, because laboratory conditions, epistasis and other cofounding factors may seriously affect identification of candidate genes (Carabbe et al. 1999; MacKay 2014).

*Behavioral vs. physiological adaptation*

The selection experiment we have analyzed is designed as a model of adaptive radiation with lines selected not only for predatory behavior, but also for high aerobic metabolism during swimming (A lines) and for the ability to maintain body weight on poor quality herbivorous diet (H lines) (Sadowska et al. 2008). This system provides unique opportunity for comprehensive comparison of the response to selection for different, ecologically important, traits, while minimizing the influence of confounding factors. All lines were derived from the single base population, have similar effective population sizes, and equal number of replicate lines are selected in each direction. Selected lines differ from unselected controls only with respect to the selection procedure. Thus, if molecular-level response to selection differs between treatments, it results most likely from differences in genetic architecture of traits under selection.

Previous RNAseq analysis of the A lines, selected for high metabolism, revealed that initial molecular level response to selection occurs primarily via changes in gene expression (Konczal et. al 2015). Unlike in the present study, no evidence was found that selection caused repeatable allele frequency changes in the A lines. Additional analysis of A lines, based on top 500 most differentiated SNPs mirroring the analysis applied to P lines here (see above), confirms the lack of evidence for repeatable allele frequency changes (p=0.81, randomization test, Fig. 5). Potentially, this difference may be result of differences in allele frequency spectra of 500 SNPs with the highest overall $F_{ST}$. If they represent variants of lower frequencies in base population for A lines, they may not respond to selection in repeatable manner due high probability of being lost by drift. However, the overall expected heterozygosity ($H_T$) calculated for the control lines did not differ between the two sets of SNPs (p=0.91, Wilcoxon test), indicating no difference in the initial allele frequency spectra. Thus, this qualitative difference between the A and P lines suggests that architecture of

genetic variation differs between the selected traits, which appear to result from differences in selection pressure acting on common variants in these two selection regimes. We conclude, that response to selection for predatory behavior appears to be associated with a smaller number of loci with higher individual selection coefficients, because alleles that segregated at intermediate frequencies in the base population show rapid and repeatable frequency changes. On the contrary, selection for high aerobic metabolism in the A lines has resulted in changes in expression of many genes, while allele frequency changes were not repeatable indicating smaller per locus selection coefficients or repeatable changes in in parts of the genome not expressed in liver or heart.
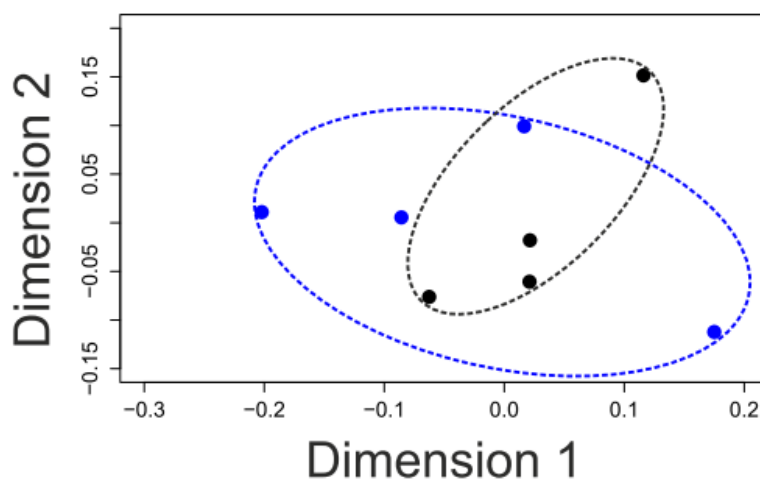


**Fig. 5** Genetic differentiation of selected (blue) and control (black) lines of bank vole selection experiment. The selection was applied for increasing maximum metabolism rate during swimming. Multidimensional scaling plot was performed on the matrices of pairwise $F_{ST}$ distances between lines calculated using top 500 SNPs with the highest mean pairwise $F_{ST}$, i.e. showing the most overall differentiation among lines. The plot is shown to compare genetic differentiation in lines selected for aerobic performance with lines selected for predatory behavior (Fig. 2B).

Our comparison of response to selection in P and A lines is relevant in a wider context of evolution of physiological and behavioral traits. A long standing idea in evolutionary biology is that "behavior evolves first" (Mayr 1958, Garland and Rose 2009) and behavioral shift may then alter the selective environment of other traits and drive their evolution. Phylogenetic analyses support this hypothesis providing evidence that compared to other traits behavior is relatively labile evolutionarily, which may be attributed to frequent adaptive changes (Blomberg et al. 2003). Also many selection experiments demonstrated that variation in behavior has a substantial genetic component. Mice were successfully selected for voluntary wheel running (Swallow et al. 1998), honeybees responded to selection for

specialization in foraging for pollen versus nectar (Page et al. 1995), mice were selectively bred for high and low activity in an open-field area (DeFries et al. 1978); many other examples of successful selection experiments which targeted behavior are reviewed in Garland and Rose (2009). However, the genetic architecture of the natural heritable variation underlying evolutionary change in behavior has rarely been traced all the way to the DNA level in vertebrates. One example is the prairie vole, in which specific motif of DNA upstream of a gene that encodes arginine vasopressine receptor protein strongly influence pair-bonding behavior (Hammock and Young 2005). Studies on laboratory mice found single chromosomal regions influencing avoidance behavior, exploration and activity (Turri et al. 2001). On the other hand, study of five human personality traits, concluded that many small-effect genes influence behavior (Terracciano et al. 2008). Because of few available examples we are still far away from the understanding of general properties of genetic variation in behavioral traits. However, if genetic architecture of behavioral traits differs from that of physiological or morphological traits, in particular if response to selection is associated with variants typically segregating at higher frequencies as suggested by our results, then genomic signals of selection may be more subtle or even indistinguishable from the neutral background in natural populations (Hermisson and Pennings 2005, Messer and Petrov 2013). On the contrary, replicated lines derived from the same base population and selected for the same behavior should show higher repeatability at the molecular level, because common variants are less prone to the loss by drift in any of replicated lines. Thus, studies focusing on more traits in multiple species are needed, to evaluate whether general properties of genetic architecture indeed differ between behavioral and physiological traits.

*Candidate genes*

Below, we present promising candidates for future studies; they were selected from genes showing largest differences between the predatory and control lines based on molecular functions potentially associated with biology of the selected trait. The minimum criteria for inclusion in the set of candidate genes were: at least one SNP with diffStat > 0.2 or significant (FDR <0.05) differences in expression and mean expression > 1 FPKM (to filter out genes with overall low expression). Such stringent criteria require repeatable changes in all 4 lines (all candidate SNPs and most of genes with differentiated expression) which may be considered independent replicates, thus should allow to identify genetic changes of the large effect on phenotype.

We found evidence for notion, that repeatable allele frequency changes are observed mainly in regulatory elements. An additional support for the effect of selection on allele frequencies in regulatory regions is provided by the observation of a larger fraction of SNPs with diffStat > 0.2 among noncoding SNPs (Fig. 2D). The noncoding SNPs can be located in noncoding RNAs, transcripts which were not automatically annotated and in partially assembled genes. To get insight into molecular function of noncoding candidate SNPs (diffStat > 0.2) we performed manual annotation of respective genes. We blasted them against the mouse genome and found that many of these sequences represent 3' untranslated regions (3'UTRs) or regions immediately downstream of genes, probably expressing expanded 3'UTRs in bank vole or unannotated transcribed regions in mouse (Tab. S1). Existing assembly strategies often fragment long 3'UTRs (Shenker et al. 2015) and some 3'UTRs may express separately from the associated protein coding sequences to which they are normally linked (Mercer et al. 2011). Allele frequency changes in such sequences may be caused by either linkage to causative variants in noncoding regions (coding nonsynonymous changes were investigated), or may be functionally important *per se*. The 3'UTRs and downstream sequences affect the expression of eukaryotic genes by regulation of mRNA translation, stability and subcellular localization (Kuersten and Goodwin 2003). 3'UTRs undergone a massive expansion during metazoan evolution and some of them are highly conserved within the mammalian genome (Siepel et al. 2005, Sandberg et al. 2008). These observations suggest that 3'UTRs have assumed an increasingly important role in the evolution of the eukaryotic genomes, thus they may be important target of selection, associated with pre- and post-transcriptional regulation of expression and evolution of alternative splicing. Although functional importance of candidate SNPs localized in 3'UTRs is more challenging to assess comparing with nonsynonymous changes, they should be proposed for future investigations, together with SNPs located in coding sequences, as a variants potentially underlying genetic variance in predatory behavior.

Two SNPs with the highest diffStat values (0.47) were localized in a noncoding transcript, which was manually annotated as 3'UTR region of cAMP-specific 3',5'-cyclic phosphodiesterase 4D (*PDE4D*). PDE4D acts as antidepressant in both animals and humans via enhancement of cAMP signaling in the brain (Zhang 2009). Mice deficient in PDE4D displayed memory enhancement and increased hippocampal neurogenesis (Li et. al 2011). Moreover, cAMP can act as a hunger signal in several tissues such as liver and muscle where glucagon can promote glycogen breakdown by activating cAMP signaling (Jiang and Zhang

2003). Changes in 3'UTR of *PDE4D* may thus be linked to functionally important variants, or cause changes in gene expression or alternative splicing of *PDE4D* (Li et al. 2011).

Some changes of gene expression in liver also suggest that a hunger level may play an important role in the response to selection. G0/G1 switch gene 2 (*G0S2*), a gene which was upregulated during chronic fasting in mouse (Zandbergen et. al 2005), showed lower expression in livers of voles from P lines. Another gene (Protein phosphatase 1 regulatory subunit 3G, *PPP1R3G*) with differences of expression in liver was reported as upregulated during fasting and downregulated after feeding in mouse (Luo et al. 2011). Finally, the leptin hormone serves as a mediator of the adaptation to fasting, and regulation of feeding and energy balance (Ahima and Flier 2000). Normal liver tissue does not express leptin, but leptin receptor (Otte et al. 2004), which expression was upregulated in predatory lines.

The motivation to catch a prey may be associated not only with a hunger but also with an aggression level. We found overexpression in liver of testosterone 17-beta-dehydrogenase 3 (*HSD17B3*), the gene encoding an enzyme which favors the reduction of androstenedione to testosterone, and thus may be potentially associated with an aggression level (Nelson nad Chiavegatto 2001).

Another potential factor associated with the selected trait are changes in nervous system. Differentiated SNPs were localized in an intron of calsyntenin 2 (*CLSTN2*), which is associated with episodic memory in humans (Preuschhof et al. 2010). Changes in expression of the hippocampus gene *NDRG4-A* may be associated with preservation of spatial learning and the resistance to neuronal cell death caused by stress (Yamamoto et al. 2011). Also genes reported as influencing mental disorders may play an important role in the evolution of the nervous system. One of the most differentiated SNPs was localized in proline dehydrogenase 1 (*PRODH*), a gene associated with cognitive dysfunctions in humans (Kempf et al. 2008). Other candidate SNPs are localized in such genes as Ubiquitin Carboxyl-Terminal Hydrolase 15 (*USP15*), associated with Parkinson disease (Cornelissen et al. 2014), or Phospholipase A2 (*PLA2*), gene linked to schizophrenia and autism (Bell et al. 2004).

We found also molecular signs of selection in genes associated with activity and circadian rhythm – differentiated expression of Aryl hydrocarbon receptor nuclear translocator 2 (*ARNT2*) in liver and differentiated allele frequencies in delta-aminolevulinate synthase 1 (*ALAS1*) intron SNPs. These findings suggest, that future phenotypic studies of

general activity and biological rhythm may show differences between predatory and control lines.

Finally, changes of neurotransmitters may also affect motivation and ability to successfully involve in predation. Candidate SNPs were detected in pyridoxine-5'-phosphate oxidase (*PNPO*); gene encoding enzyme that catalyzes the rate limiting step in the synthesis of pyridoxal 5'-phosphate, which is an important cofactor in biosynthesis of many neurotransmitters including GABA (Petroff 2002). Interestingly, one of the differentially expressed genes between treatments in hippocampus was 3'UTR of gamma-aminobutyric acid B receptor 1 (*GABBR1*). GABA is the main inhibitory neurotransmitter in the central nervous system; its actions are mediated by GABAB receptors. We thus suspect, that changes associated with GABAergic signaling may be responsible for the evolution of predatory behavior.

## Conclusions

Selection for predatory behavior affected allele frequencies in multiple genes and overall pattern of hippocampal gene expression, while it did not affect polymorphism within lines and overall gene expression pattern in liver. Our results indicate that response to selection for predatory behavior is associated with variants of relatively large effect and/or variants which were segregating at relatively high frequencies in the base population. These results are qualitatively different from those obtained in a similar analysis of lines selected for maximum metabolism, in which no repeatable changes of allele frequencies were detected. We thus hypothesize that the nature of genetic variation available for selection acting on behavioral and physiological traits may be fundamentally different. Combination of selection experiment and transcriptome sequencing allowed us to select candidate genes, potentially underlying predatory behavior in bank vole. These candidate genes are associated with hunger, aggression, biological rhythm and functioning of the nervous system. Overall, the description of the genetic architecture and genes underlying predatory behavior in bank voles is relevant to advance general understanding of evolution of ecologically important traits.

## Materials and methods

### *Selection experiment*

This study was performed using individuals from 13[th] generation of a laboratory colony of bank vole (*Myodes (=Clethrionomys) glareolus*) that was subject to selection for predatory behavior. Detailed information about the animal maintenance and welfare, selection protocols, and direct effects of selection is presented by Sadowska et al. (2008) and Chrząścik et al. (2014). Briefly, the colony was reared from about 320 voles captured in the Niepołomice Forest in southern Poland (Sadowska et al. 2008). For 6-7 generations, the animals were bred randomly, and then multidimensional selection experiment was established. In P lines analyzed here (four independent lines), the selection criterion was a time to catch live cricket (*Grillus assimillis*). The voles were fasten for 10-12h, then cricket was placed in each cage and then presence of the cricket was checked after 0.5, 1, 3, 6, and 10 min. The tests were repeated two, three or four times for each individual, depending on generation and human resources. After 13 generations of selection, the proportion of voles attacking crickets was 5 times higher in the selected P lines than in unselected control lines (Fig 6; Chrząścik et al. 2014).
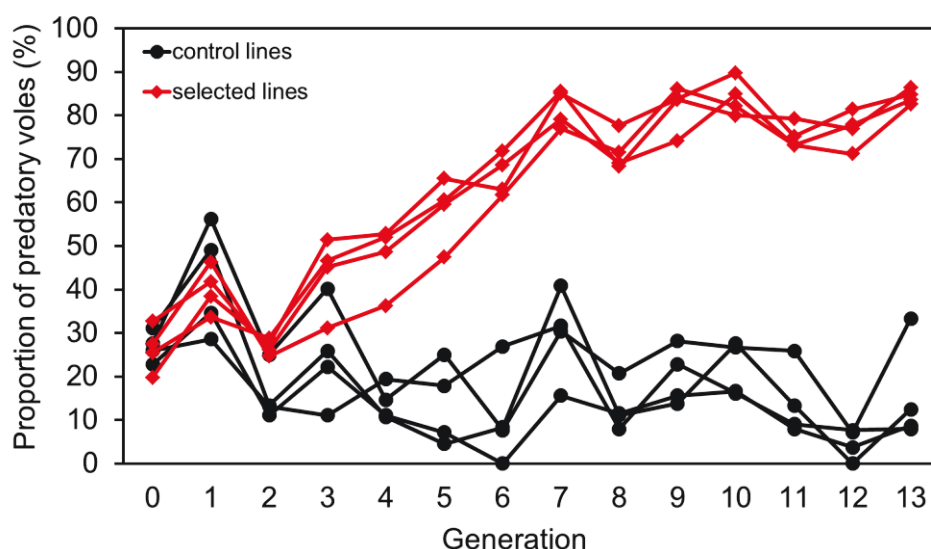


**Fig. 6** Results of 13 generations of selection for the predatory behavior in the bank vole (replicate line means).

*Sampling, sequencing and data filtering*

Five females and five males of 75-80 days in age were sampled from each line; each individual came from a different family. The individuals were previously used only for routine measurements of body mass. Voles were euthanized by being placed individually in a jar containing isofulrane (Aerane®) fumes. After that small part of the left liver lobe and entire hippocampus were excised and immediately placed in RNAlater (Sigma). Tissues were collected between 8.00 am and 2.00 pm. Samples were stored overnight at 4°C and then frozen at -20°C. Hippocampus samples were collected from individuals from four selected and four control lines. Liver samples were collected only from the P lines. For the C lines we used previously reported liver transcriptomes, obtained from tissue collected at the same time and using identical procedure (Konczal et al. 2015).

Total RNA was extracted with RNAzol® (MRC); RNA concentration and quality were measured with Nanodrop and Agilent 2100 Bioanalyzer. All samples had RNA Integrity Number higher than 7.0. Then, for each organ, we prepared one pooled sample per line using equal amounts of total RNA from each individual. Residual DNA was removed from pooled samples using DNA-free Kit (Ambion®).

Preparation of barcoded cDNA libraries with TrueSeq RNA kit was performed by Georgia Genomic Facility, USA. Hippocampus sample from one control line (C3) was paired-end sequenced (2 x 100bp) and used for reference transcriptome reconstruction. For the remaining 11 pools, single-end (1 x 100 bp) sequencing was performed.

*Hippocampus reference transcriptome reconstruction and annotation*

Pair end reads were trimmed with DynamicTrim (Cox et al. 2010) and used for the reconstruction of bank vole hippocampus transcriptome *de novo* with Trinity assembler (version 2013-02-15 with –REDUCE option; Grabherr et al. 2011). We then processed the Trinity output by merging transcripts that were probably derived from the same genomic location and subsequently produced transcriptome-based gene models, which we refer to here as "genes" (Stuglik et al. 2014).

Putative coding sequences were identified using the pipeline implemented in Trinity and they were annotated using Trinotate software and homology search to Swissprot database. For candidate genes which could not be annotated automatically, we attempted manual annotation using blast searches against the mouse genome.

*SNP analyses*

Single-end reads were trimmed with DynamicTrim (Cox et al. 2010) and adaptors were removed with Cutadapt (Martin 2011). We subsampled also single-end reads from pair-end reads, to obtain comparable amount of data for all lines and organs. Reads were mapped to the reference transcriptomes using Bowtie2 (Langmead and Salzberg 2012) and we considered only reads with mapping quality > 20 and positions with base quality > 20 phred. SNP calling was performed with samtools (Li et al. 2009), Popoolation2 (Kofler et al. 2011) and custom scripts as described in detail elsewhere (Konczal et al. 2015).

$F_{ST}$ distances were calculated for each SNP with PoPoolation2 (Kofler et al. 2011). To test for separate clustering of selected and control lines we calculated the ratio of between treatment to within treatment variance using adonis {vegan} (Oksanen et al. 2013) and assessed its statistical significance through 1000 randomizations. Randomized matrices of mean $F_{ST}$ were obtained by shuffling pairwise $F_{ST}$ values for each SNP independently.

*Simulations of allele frequency differentiation under drift*

To obtain the rate of allele frequency differentiation that would be expected under drift, we performed forward drift simulations on known pedigrees. Simulations were performed separately for allele frequency spectra derived from all, synonymous, nonsynonymous, UTR-located and noncoding SNPs using scripts available from http://www.molecol.eko.uj.edu.pl.

*Expression analyses*

To identify differentially expressed genes, we mapped reads to the reference transcriptomes with bowtie and used Trinity pipeline with EdgeR Bioconductor and RSEM (Gentleman et al. 2004, Grabherr et al. 2011, Li and Dewey 2011). Only genes for which the sum of expected counts over all samples was higher than 10 were used for analyses.

To statistically test for separate clustering of transcriptional profiles of selected and control lines we used similar strategy to that for $F_{ST}$. We used table of expression values (FPKM, TMM normalized) and calculated distance matrix (dist() function) followed by calculation of the ratio of between treatment to within treatment variance. The statistical significance of this ratio was assessed through 1000 randomizations. Differences between lines in genome-wide transcriptional profiles were visualized with multidimensional scaling (plotMDS {edgeR}).

**References**

Ahima, R. S., & Flier, J. S. (2000). Leptin. *Annual review of physiology*, *62*(1), 413-437.

Akey, J. M. (2009). Constructing genomic maps of positive selection in humans: Where do we go from here?. *Genome research*, *19*(5), 711-722.

Barbosa, P., & Castellanos, I. (Eds.). (2005). *Ecology of predator-prey interactions*. Oxford University Press.

Barrett, R. D., & Schluter, D. (2008). Adaptation from standing genetic variation. *Trends in Ecology & Evolution*, *23*(1), 38-44.

Bell, J. G., MacKinlay, E. E., Dick, J. R., MacDonald, D. J., Boyle, R. M., & Glen, A. C. A. (2004). Essential fatty acids and phospholipase A 2 in autistic spectrum disorders. *Prostaglandins, leukotrienes and essential fatty acids*,*71*(4), 201-204.

Blomberg, S. P., Garland, T., & Ives, A. R. (2003). Testing for phylogenetic signal in comparative data: behavioral traits are more labile. *Evolution*, *57*(4), 717-745.

Boake, C. R., Arnold, S. J., Breden, F., Meffert, L. M., Ritchie, M. G., Taylor, B. J., ... & Moore, A. J. (2002). Genetic tools for studying adaptation and the evolution of behavior. *The American Naturalist*, *160*(S6), S143-S159.

Cornelissen, T., Haddad, D., Wauters, F., Van Humbeeck, C., Mandemakers, W., Koentjoro, B., ... & Vandenberghe, W. (2014). The deubiquitinase USP15 antagonizes Parkin-mediated mitochondrial ubiquitination and mitophagy.*Human molecular genetics*, ddu244.

Crabbe, J. C., Wahlsten, D., & Dudek, B. C. (1999). Genetics of mouse behavior: interactions with laboratory environment. *Science*, *284*(5420), 1670-1672.

Cho, Y. S., Hu, L., Hou, H., Lee, H., Xu, J., Kwon, S., ... & Ko, J. (2013). The tiger genome and comparative analysis with lion and snow leopard genomes. *Nature communications*, *4*.

Chrząścik, K. M., Sadowska, E. T., Rudolf, A., & Koteja, P. (2014). Learning ability in bank voles selected for high aerobic metabolism, predatory behaviour and herbivorous capability. *Physiology & behavior*, *135*, 143-151.

Cox, M. P., Peterson, D. A., & Biggs, P. J. (2010). SolexaQA: At-a-glance quality assessment of Illumina second-generation sequencing data. *BMC bioinformatics*, *11*(1), 485.

Curio, E. (1976). *Ethology of predation*. Springer-Verlag.

Dalziel, A. C., Rogers, S. M., & Schulte, P. M. (2009). Linking genotypes to phenotypes and fitness: how mechanistic biology can inform molecular ecology. *Molecular ecology*, *18*(24), 4997-5017.

DeFries, J. C., Gervais, M. C., & Thomas, E. A. (1978). Response to 30 generations of selection for open-field activity in laboratory mice. *Behavior genetics*, *8*(1), 3-13.

Eisenberg, J.F., Leyhausen, P. (1972). The phylogenesis of predatory behavior in mammals. *Z Tierpsychol*, 30:59-93.

Ellegren, H., & Sheldon, B. C. (2008). Genetic basis of fitness differences in natural populations. *Nature*, *452*(7184), 169-175.

Fraser, H. B., Moses, A. M., & Schadt, E. E. (2010). Evidence for widespread adaptive evolution of gene expression in budding yeast. *Proceedings of the National Academy of Sciences*, *107*(7), 2977-2982.

Fraser, H. B. (2013). Gene expression drives local adaptation in humans. *Genome research*, *23*(7), 1089-1096.

Frazer, K. A., Murray, S. S., Schork, N. J., & Topol, E. J. (2009). Human genetic variation and its contribution to complex traits. *Nature Reviews Genetics*, *10*(4), 241-251.

Gage, F. H. (2000). Mammalian neural stem cells. *Science*, *287*(5457), 1433-1438.

Garland, T., & Rose, M. R. (Eds.). (2009). *Experimental evolution: concepts, methods, and applications of selection experiments*. Berkeley, CA, USA: University of California Press.

Gentleman, R. C., Carey, V. J., Bates, D. M., Bolstad, B., Dettling, M., Dudoit, S., ... & Zhang, J. (2004). Bioconductor: open software development for computational biology and bioinformatics. *Genome biology*, *5*(10), R80.

Grabherr, M. G., Haas, B. J., Yassour, M., Levin, J. Z., Thompson, D. A., Amit, I., ... & Regev, A. (2011). Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nature biotechnology*, *29*(7), 644-652.

Hahn, M. W. (2008). Toward a selection theory of molecular evolution. *Evolution*, *62*(2), 255-265.

Halligan, D. L., Kousathanas, A., Ness, R. W., Harr, B., Eöry, L., Keane, T. M., ... & Keightley, P. D. (2013). Contributions of protein-coding and regulatory change to adaptive molecular evolution in murid rodents. *PLoS genetics*, *9*(12), e1003995.

Hammock, E. A., & Young, L. J. (2004). Functional microsatellite polymorphism associated with divergent social structure in vole species. *Molecular Biology and Evolution*, *21*(6), 1057-1063.

Hermisson, J., & Pennings, P. S. (2005). Soft sweeps molecular population genetics of adaptation from standing genetic variation. *Genetics*, *169*(4), 2335-2352.

Hoekstra, H. E., Hirschmann, R. J., Bundey, R. A., Insel, P. A., & Crossland, J. P. (2006). A single amino acid mutation contributes to adaptive beach mouse color pattern. *Science*, *313*(5783), 101-104.

Hudson, R. R. (2002). Generating samples under a Wright–Fisher neutral model of genetic variation. *Bioinformatics*, *18*(2), 337-338.

Ishii, Y., & Shimada, M. (2010). The effect of learning and search images on predator–prey interactions. *Population ecology*, *52*(1), 27-35.

Jiang, G., & Zhang, B. B. (2003). Glucagon and regulation of glucose metabolism. *American Journal of Physiology-Endocrinology And Metabolism*,*284*(4), E671-E678.

Jones, F. C., Grabherr, M. G., Chan, Y. F., Russell, P., Mauceli, E., Johnson, J., ... & Broad Institute Genome Sequencing Platform & Whole Genome Assembly Team. (2012). The genomic basis of adaptive evolution in threespine sticklebacks. *Nature*, *484*(7392), 55-61.

Kawecki, T. J., Lenski, R. E., Ebert, D., Hollis, B., Olivieri, I., & Whitlock, M. C. (2012). Experimental evolution. *Trends in ecology & evolution*, *27*(10), 547-560.

Kempf, L., Nicodemus, K. K., Kolachana, B., Vakkalanka, R., Verchinski, B. A., Egan, M. F., ... & Meyer-Lindenberg, A. (2008). Functional polymorphisms in PRODH are associated with risk and protection for schizophrenia and fronto-striatal structure and function. *PLoS genetics*, *4*(11), e1000252.

Kennedy, P. J., & Shapiro, M. L. (2009). Motivational states activate distinct hippocampal representations to guide goal-directed behaviors. *Proceedings of the National Academy of Sciences*, *106*(26), 10805-10810.

Kessner, D., Turner, T. L., & Novembre, J. (2013). Maximum likelihood estimation of frequencies of known haplotypes from pooled sequence data. *Molecular biology and evolution*, *30*(5), 1145-1158.

King, M. C., & Wilson, A. C. (1975). Evolution at two levels in humans and chimpanzees. *Science*, *188*(4184), 107-116.

Kofler, R., Pandey, R. V., & Schlötterer, C. (2011). PoPoolation2: identifying differentiation between populations using sequencing of pooled DNA samples (Pool-Seq). *Bioinformatics*, *27*(24), 3435-3436.

Konczal, M., Koteja, P., Stuglik, M. T., Radwan, J., & Babik, W. (2014). Accuracy of allele frequency estimation using pooled RNA-Seq. *Molecular ecology resources*, *14*(2), 381-392.

Konczal M., Babik W., Radwan J., Sadowska E.T., Koteja P. (2015) Initial molecular-level response to artificial selection for increased aerobic metabolism occurs primarily via changes in gene expression. *Molecuar Biology and Evolution, doi: 10.1093/molbev/msv038*

Kuersten, S., & Goodwin, E. B. (2003). The power of the 3′ UTR: translational control and development. *Nature Reviews Genetics*, *4*(8), 626-637.

Langmead, B., & Salzberg, S. L. (2012). Fast gapped-read alignment with Bowtie 2. *Nature methods*, *9*(4), 357-359.

Li, Y. F., Cheng, Y. F., Huang, Y., Conti, M., Wilson, S. P., O'Donnell, J. M., & Zhang, H. T. (2011). Phosphodiesterase-4D knock-out and RNA interference-mediated knock-down enhance memory and increase hippocampal neurogenesis via increased cAMP signaling. *The Journal of Neuroscience*,*31*(1), 172-183.

Li, B., & Dewey, C. N. (2011). RSEM: accurate transcript quantification from RNA-Seq data with or without a reference genome. *BMC bioinformatics*, *12*(1), 323.

Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., ... & Durbin, R. (2009). The sequence alignment/map format and SAMtools. *Bioinformatics*, *25*(16), 2078-2079.

Luo, X., Zhang, Y., Ruan, X., Jiang, X., Zhu, L., Wang, X., ... & Chen, Y. (2011). Fasting-induced protein phosphatase 1 regulatory subunit contributes to postprandial blood glucose homeostasis via regulation of hepatic glycogenesis. *Diabetes*, *60*(5), 1435-1445.

Martin, M. (2011). Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet. journal*, *17*(1), pp-10.

Mayr, E. 1958. Behavior and systematics. In A. Roe and G. G. Simpson, eds. Behavior and Evolution. New Haven, CT: Yale University Press.

Mercer, T. R., Wilhelm, D., Dinger, M. E., Solda, G., Korbie, D. J., Glazov, E. A., ... & Mattick, J. S. (2011). Expression of distinct RNAs from 3′ untranslated regions. *Nucleic acids research*, *39*(6), 2393-2403.

Messer, P. W., & Petrov, D. A. (2013). Population genomics of rapid adaptation by soft selective sweeps. *Trends in ecology & evolution*, *28*(11), 659-669.

Morris, R. G. M., & Hagan, J. J. (1983). Hippocampal electrical activity and ballistic movement. *Neurobiology of the hippocampus*, 321-331.

Nei, M. (2013). *Mutation-driven evolution*. Oxford University Press.

Nelson, R. J., & Chiavegatto, S. (2001). Molecular basis of aggression. *Trends in neurosciences*, *24*(12), 713-719.

Oksanen, J., Blanchet, F. G., Kindt, R., Legendre, P., Minchin, P. R., O'Hara, R. B., ... & Imports, M. A. S. S. (2013). Package 'vegan'. *Community ecology package, version*, *2*(9).

Orr, H. A. (2005). The genetic theory of adaptation: a brief history. *Nature Reviews Genetics*, *6*(2), 119-127.

Otte, C., Otte, J. M., Strodthoff, D., Bornstein, S. R., Fölsch, U. R., Mönig, H., & Kloehn, S. (2004). Expression of leptin and leptin receptor during the development of liver fibrosis and cirrhosis. *Experimental and clinical endocrinology & diabetes: official journal, German Society of Endocrinology [and] German Diabetes Association*, *112*(1), 10-17.

Page, R. E., Waddington, K. D., Hunt, G. J., & Fondrk, M. K. (1995). Genetic determinants of honey bee foraging behaviour. *Animal behaviour*, *50*(6), 1617-1625.

Petroff, O. A. (2002). Book Review: GABA and glutamate in the human brain. *The Neuroscientist*, *8*(6), 562-573.

Petrusewicz, K. (1983). Ecology of the bank vole. *Acta Theriologica, 28*, 1-242.

Preuschhof, C., Heekeren, H. R., Li, S. C., Sander, T., Lindenberger, U., & Bäckman, L. (2010). KIBRA and CLSTN2 polymorphisms exert interactive effects on human episodic memory. *Neuropsychologia*, *48*(2), 402-408.

Radwan, J., & Babik, W. (2012). The genomics of adaptation. *Proceedings of the Royal Society B: Biological Sciences*, *279*(1749), 5024-5028.

Rhodes, J. S., van Praag, H., Jeffrey, S., Girard, I., Mitchell, G. S., Garland Jr, T., & Gage, F. H. (2003). Exercise increases hippocampal neurogenesis to high levels but does not improve spatial learning in mice bred for increased voluntary wheel running. *Behavioral neuroscience*, *117*(5), 1006.

Ritchie, E. G., & Johnson, C. N. (2009). Predator interactions, mesopredator release and biodiversity conservation. *Ecology letters*, 12(9), 982-998.

Ritchie, E. G., Elmhagen, B., Glen, A. S., Letnic, M., Ludwig, G., & McDonald, R. A. (2012). Ecosystem restoration with teeth: what role for predators?. *Trends in Ecology & Evolution*, *27*(5), 265-271.

Rockman, M. V. (2012). The QTN program and the alleles that matter for evolution: all that's gold does not glitter. *Evolution*, *66*(1), 1-17.

Sadowska, E. T., Baliga-Klimczyk, K., Chrząścik, K. M., & Koteja, P. (2008). Laboratory model of adaptive radiation: a selection experiment in the bank vole. *Physiological and Biochemical Zoology*, *81*(5), 627-640.

Sandberg, R., Neilson, J. R., Sarma, A., Sharp, P. A., & Burge, C. B. (2008). Proliferating cells express mRNAs with shortened 3'untranslated regions and fewer microRNA target sites. *Science*, *320*(5883), 1643-1647.

Sella, G., Petrov, D. A., Przeworski, M., & Andolfatto, P. (2009). Pervasive natural selection in the Drosophila genome?. *PLoS genetics*, *5*(6), e1000495.

Shenker, S., Miura, P., Sanfilippo, P., Lai, E. C., Carnes, J., Lerch, M., ... & Chang, H. Y. (2015). IsoSCM: improved and alternative 3′ UTR annotation using multiple change-point inference. *rna*, *21*(1), 1-13.

Siepel, A., Bejerano, G., Pedersen, J. S., Hinrichs, A. S., Hou, M., Rosenbloom, K., ... & Haussler, D. (2005). Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. *Genome research*, *15*(8), 1034-1050.

Stapley, J., Reger, J., Feulner, P. G., Smadja, C., Galindo, J., Ekblom, R., ... & Slate, J. (2010). Adaptation genomics: the next generation. *Trends in ecology & evolution*, *25*(12), 705-712.

Storz, J. F., Sabatino, S. J., Hoffmann, F. G., Gering, E. J., Moriyama, H., Ferrand, N., ... & Nachman, M. W. (2007). The molecular basis of high-altitude adaptation in deer mice. *PLoS Genetics*, *3*(3), e45.

Stuglik, M. T., Babik, W., Prokop, Z., & Radwan, J. (2014). Alternative reproductive tactics and sex-biased gene expression: the study of the bulb mite transcriptome. *Ecology and Evolution*, *4*(5), 623-632.

Swallow, J. G., Carter, P. A., & Garland Jr, T. (1998). Artificial selection for increased wheel-running behavior in house mice. *Behavior genetics*, *28*(3), 227-237.

Swanson, L. W. (1983). The hippocampus and the concept of the limbic system. *Neurobiology of the Hippocampus*, 3-19.

Terracciano, A., Sanna, S., Uda, M., Deiana, B., Usala, G., Busonero, F., ... & Costa, P. T. (2010). Genome-wide association scan for five major dimensions of personality. *Molecular psychiatry*, *15*(6), 647-656.

Tracy, A. L., Jarrard, L. E., & Davidson, T. L. (2001). The hippocampus and motivation revisited: appetite and activity. *Behavioural brain research*, *127*(1), 13-23.

Travisano, M., & Shaw, R. G. (2013). Lost in the map. *Evolution*, *67*(2), 305-314.

Turri, M. G., Datta, S. R., DeFries, J., Henderson, N. D., & Flint, J. (2001). QTL analysis identifies multiple behavioral dimensions in ethological tests of anxiety in laboratory mice. *Current Biology*, *11*(10), 725-734.

Wagner, A. (2008). Neutralism and selectionism: a network-based reconciliation. *Nature Reviews Genetics*, *9*(12), 965-974.

Wakeley, J. (2009). Coalescent theory: an introduction (Vol. 1). Greenwood Village, Colorado: Roberts & Company Publishers.

Weber, J. N., Peterson, B. K., & Hoekstra, H. E. (2013). Discrete genetic modules are responsible for complex burrow evolution in Peromyscus mice. *Nature*, *493*(7432), 402-405.

Wereszczynska, A. M., Nowakowski, W. K., Nowakowski, J. K., & Jedrzejewska, B. (2007). Is food quality responsible for the cold-season decline in bank vole density? Laboratory experiment with herb and acorn diets. *Folia zoologica*, *56*(1), 23-32.

Yamamoto, H., Kokame, K., Okuda, T., Nakajo, Y., Yanamoto, H., & Miyata, T. (2011). NDRG4 protein-deficient mice exhibit spatial learning deficits and vulnerabilities to cerebral ischemia. *Journal of Biological Chemistry*, *286*(29), 26158-26165.

Yi, X., Liang, Y., Huerta-Sanchez, E., Jin, X., Cuo, Z. X. P., Pool, J. E., ... & Cao, Z. (2010). Sequencing of 50 human exomes reveals adaptation to high altitude. *Science*, *329*(5987), 75-78.

Zandbergen, F., Mandard, S. X., Escher, P., Tan, N. X., Patsouris, D., Jatkoe, T., ... & Kersten, S. (2005). The G0/G1 switch gene 2 is a novel PPAR target gene. *Biochem. J*, *392*, 313-324.

Zhan, X., Pan, S., Wang, J., Dixon, A., He, J., Muller, M. G., ... & Bruford, M. W. (2013). Peregrine and saker falcon genome sequences provide insights into evolution of a predatory lifestyle. *Nature genetics*, *45*(5), 563-566.

Zhang, H. T. (2009). Cyclic AMP-specific phosphodiesterase-4 as a target for the development of antidepressant drugs. *Current pharmaceutical design*, *15*(14), 1688-1698.

**Supplementary materials**

**Tab. S1** SNPs with the largest differences in allele frequencies (diffStat > 0.2) between predatory and control lines. Gene annotation was performed either with Trinotate and SwissProt database, or manually (M) with mouse genome. SNP class nonORF indicates SNPs localized in genes without likely protein coding sequences and analyzed as noncoding. Column Ma show major allele and next columns present its frequencies in selection experiment lines. DiffStat is a minimum allele frequency difference between selected and control lines.

| Gene annotation | Ensembl symbol of mouse gene | SNP class | Contig name | Pos | Ma | C1 | C2 | C3 | C4 | P1 | P2 | P3 | P4 | Diff Stat |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Phosphodiesterase 4D, cAMP-specific, 3'UTR and upstrem UTR (M) | ENSMUSG00000021699 | nonORF | comp129965 | 3630 | C | 0.43 | 0.23 | 0.53 | 0.48 | 1.00 | 1.00 | 1.00 | 1.00 | **0.47** |
| Phosphodiesterase 4D, cAMP-specific, 3'UTR and upstrem UTR (M) | ENSMUSG00000021699 | nonORF | comp129965 | 3631 | A | 0.41 | 0.23 | 0.53 | 0.48 | 1.00 | 1.00 | 1.00 | 1.00 | **0.47** |
| ATP-binding cassette sub-family A member 9 | ENSMUSG00000030249 | UTR | Contig_88975 | 255 | T | 0.37 | 0.22 | 0.25 | 0.38 | 0.84 | 0.98 | 0.95 | 0.94 | **0.46** |
| Ubiquitin carboxyl-terminal hydrolase 15 | ENSMUSG00000020124 | UTR | Contig_69084 | 3026 | A | 0.83 | 0.80 | 0.95 | 1.00 | 0.35 | 0.34 | 0.29 | 0.36 | **0.44** |
| Calsyntenin 2, 3'UTR (M) | ENSMUSG00000032452 | UTR | comp138591 | 6290 | A | 0.84 | 0.73 | 0.92 | 0.86 | 0.28 | 0.35 | 0.10 | 0.00 | **0.38** |
| SLAIN motif-containing protein 2 | ENSMUSG00000036087 | UTR | comp130355 | 2814 | G | 0.31 | 0.35 | 0.33 | 0.05 | 0.95 | 0.92 | 0.86 | 0.70 | **0.35** |
| RAB6B, member RAS oncogene family (M) | ENSMUSG00000032549 | UTR | comp79801 | 1198 | A | 0.50 | 0.52 | 0.47 | 0.28 | 0.90 | 0.96 | 0.86 | 0.88 | **0.34** |
| SLAIN motif-containing protein 2 | ENSMUSG00000036087 | UTR | comp130355 | 2819 | T | 0.79 | 0.84 | 0.79 | 0.95 | 0.09 | 0.21 | 0.46 | 0.42 | **0.33** |
| SLAIN motif-containing protein 2 | ENSMUSG00000036087 | UTR | comp130355 | 3222 | A | 0.69 | 0.82 | 0.65 | 0.90 | 0.23 | 0.05 | 0.32 | 0.21 | **0.33** |
| Delta-aminolevulinate synthase 1, intron (M) | ENSMUSG00000032786 | nonORF | Contig_37392 | 146 | A | 0.91 | 0.76 | 0.78 | 0.90 | 0.26 | 0.35 | 0.26 | 0.43 | **0.33** |

| Protein | Gene ID | Type | Contig | Pos | Base | | | | | | | | | Value |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **SLAIN motif-containing protein 2** | ENSMUSG00000036087 | UTR | comp130355 | 3227 | C | 0.69 | 0.82 | 0.65 | 0.89 | 0.26 | 0.05 | 0.33 | 0.22 | **0.32** |
| **CB1 cannabinoid receptor-interacting protein 1** | ENSMUSG00000044629 | UTR | comp95622 | 2868 | T | 0.51 | 0.38 | 0.29 | 0.43 | 0.87 | 0.89 | 0.84 | 0.83 | **0.32** |
| **Probable RNA-binding protein 19** | ENSMUSG00000097368 | UTR | comp114111 | 1344 | A | 0.46 | 0.45 | 0.48 | 0.22 | 0.79 | 0.89 | 0.89 | 0.84 | **0.31** |
| **Kinesin-like protein KIFC3** | ENSMUSG00000031788 | syn | comp138183 | 3816 | C | 0.17 | 0.39 | 0.32 | 0.40 | 0.75 | 0.84 | 0.71 | 0.78 | **0.31** |
| **-** | - | nonORF | Contig_95669 | 2514 | T | 0.37 | 0.43 | 0.53 | 0.30 | 0.92 | 0.84 | 0.95 | 0.86 | **0.31** |
| **BMP-binding endothelial regulator protein (M)** | ENSMUSG00000031963 | nsyn | comp88485 | 2811 | T | 0.52 | 0.45 | 0.36 | 0.00 | 1.00 | 1.00 | 0.81 | 0.88 | **0.29** |
| **SLAIN motif-containing protein 2** | ENSMUSG00000036087 | UTR | comp130355 | 4030 | C | 0.35 | 0.14 | 0.37 | 0.13 | 0.86 | 0.89 | 0.70 | 0.66 | **0.29** |
| **BMP-binding endothelial regulator protein** | ENSMUSG00000031963 | UTR | comp88485 | 3000 | T | 0.46 | 0.45 | 0.48 | 0.14 | 1.00 | 1.00 | 0.81 | 0.77 | **0.29** |
| **Hippcalin (M)** | ENSMUSG00000028785 | UTR | comp67098 | 1412 | G | 0.59 | 0.31 | 0.45 | 0.65 | 1.00 | 0.95 | 0.94 | 1.00 | **0.29** |
| **Proline dehydrogenase 1, mitochondrial** | ENSMUSG00000003526 | syn | comp124278 | 841 | C | 0.30 | 0.40 | 0.35 | 0.29 | 0.68 | 0.72 | 0.81 | 0.75 | **0.28** |
| **Mitogen-activated protein kinase 4** | ENSMUSG00000024558 | UTR | comp128832 | 989 | G | 0.88 | 0.81 | 0.66 | 0.72 | 0.27 | 0.24 | 0.36 | 0.38 | **0.28** |
| **Delta-aminolevulinate synthase 1, intron (M)** | ENSMUSG00000032786 | nonORF | Contig_37392 | 166 | T | 0.88 | 0.78 | 0.79 | 0.95 | 0.29 | 0.40 | 0.19 | 0.50 | **0.28** |
| **Copine-2** | ENSMUSG00000034361 | UTR | comp91950 | 3166 | G | 0.75 | 0.85 | 0.78 | 0.74 | 0.04 | 0.46 | 0.43 | 0.43 | **0.28** |
| **UPF0258 protein KIAA1024-like homolog** | ENSMUSG00000050875 | nsyn | comp116305 | 815 | C | 0.85 | 0.75 | 1.00 | 0.69 | 0.42 | 0.34 | 0.32 | 0.35 | **0.27** |

| Gene | Ensembl ID | Type | Contig | Pos | Allele | | | | | | | | | Value |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **SHC-transforming protein 2** | ENSMUSG00000020312 | syn | Contig_112868 | 82 | A | 0.46 | 0.33 | 0.37 | 0.41 | 0.97 | 0.73 | 0.73 | 0.94 | **0.27** |
| **ATP-binding cassette sub-family A member 9** | ENSMUSG00000030249 | UTR | Contig_88975 | 552 | A | 0.34 | 0.04 | 0.08 | 0.32 | 0.86 | 0.66 | 0.61 | 0.89 | **0.27** |
| **V-Set And Transmembrane Domain Containing 2A, 3'UTR (M)** | ENSMUSG00000048834 | nonORF | comp126038 | 3276 | G | 0.42 | 0.19 | 0.07 | 0.36 | 0.73 | 0.68 | 0.88 | 0.83 | **0.26** |
| **Membrane magnesium transporter 2, upstrem DNA (M)** | ENSMUSG00000048497 | nonORF | comp91416 | 286 | G | 0.24 | 0.27 | 0.41 | 0.24 | 0.80 | 0.86 | 0.67 | 0.84 | **0.26** |
| **PQ-loop repeat-containing protein 1** | ENSMUSG00000034006 | syn | comp128921 | 7738 | A | 0.96 | 1.00 | 0.95 | 0.89 | 0.63 | 0.63 | 0.11 | 0.49 | **0.26** |
| **ADNP homeobox protein 2** | ENSMUSG00000051149 | UTR | comp90500 | 4586 | G | 0.81 | 1.00 | 0.94 | 0.79 | 0.53 | 0.51 | 0.17 | 0.04 | **0.26** |
| **Peptide chain release factor 1-like, mitochondrial** | ENSMUSG00000019774 | UTR | comp2146 | 1614 | C | 0.42 | 0.45 | 0.40 | 0.18 | 0.75 | 0.84 | 0.71 | 1.00 | **0.26** |
| **Interferon gamma receptor 1** | ENSMUSG00000020009 | UTR | comp121809 | 445 | G | 0.28 | 0.46 | 0.41 | 0.38 | 0.80 | 0.72 | 0.74 | 0.92 | **0.26** |
| **Mannosyl-oligosaccharide 1,2-alpha-mannosidase IB** | ENSMUSG00000008763 | UTR | Contig_82607 | 4753 | T | 0.75 | 0.83 | 0.75 | 0.72 | 0.38 | 0.30 | 0.46 | 0.23 | **0.26** |
| **V-Set And Transmembrane Domain Containing 2A, 3'UTR (M)** | ENSMUSG00000048834 | nonORF | comp60800 | 322 | C | 0.42 | 0.06 | 0.19 | 0.24 | 0.78 | 0.67 | 0.82 | 0.82 | **0.25** |
| **-** | - | nonORF | Contig_90579 | 1474 | A | 0.29 | 0.19 | 0.16 | 0.43 | 0.68 | 0.79 | 0.80 | 0.71 | **0.25** |
| **Pyridoxine-5'-phosphate oxidase** | ENSMUSG00000018659 | syn | comp117077 | 262 | G | 0.30 | 0.64 | 0.68 | 0.75 | 1.00 | 1.00 | 1.00 | 1.00 | **0.25** |

| Gene | Ensembl ID | Type | Contig | Pos | Allele | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Integrin Alpha 7 Chain (M) | ENSMUSG00000025348 | UTR | comp131887 | 6216 | G | 1.00 | 1.00 | 1.00 | 1.00 | 0.46 | 0.63 | 0.57 | 0.75 | **0.25** |
| Semaphorin-4G (M) | ENSMUSG00000025207 | UTR | Contig_79198 | 3726 | A | 0.74 | 0.85 | 0.61 | 0.61 | 0.32 | 0.34 | 0.31 | 0.36 | **0.25** |
| Ring finger protein 165, intron (M) | ENSMUSG00000025427 | nonORF | comp131856 | 2751 | C | 0.68 | 0.70 | 0.65 | 0.59 | 0.95 | 1.00 | 1.00 | 1.00 | **0.25** |
| Neurofibromin 2, 3'UTR (M) | ENSMUSG00000009073 | nonORF | comp69409 | 888 | G | 0.66 | 0.57 | 0.64 | 0.62 | 1.00 | 0.91 | 0.98 | 1.00 | **0.25** |
| Steroid 5-alpha-reductase, 3'UTR and upstream DNA (M) | ENSMUSG00000021594 | nonORF | comp107686 | 1152 | A | 0.93 | 0.90 | 0.95 | 0.94 | 0.43 | 0.65 | 0.60 | 0.54 | **0.25** |
| Myeloid differentiation primary response gene 88, 3'UTR (M) | ENSMUSG00000032508 | nonORF | Contig_89210 | 286 | T | 0.39 | 0.36 | 0.25 | 0.25 | 0.91 | 0.78 | 0.64 | 0.76 | **0.25** |
| Neurofibromin 2, 3'UTR (M) | ENSMUSG00000009073 | nonORF | comp69409 | 911 | T | 0.69 | 0.53 | 0.66 | 0.55 | 1.00 | 0.93 | 1.00 | 1.00 | **0.24** |
| Zinc finger protein 295 | ENSMUSG00000046962 | UTR | comp96589 | 2320 | T | 0.68 | 0.73 | 0.69 | 0.72 | 0.20 | 0.21 | 0.44 | 0.44 | **0.24** |
| Transmembrane protein 120A | ENSMUSG00000039886 | nsyn | comp67874 | 425 | T | 0.41 | 0.41 | 0.48 | 0.42 | 0.83 | 0.91 | 0.74 | 0.72 | **0.24** |
| Testican-2, 3'UTR (M) | ENSMUSG00000058297 | nsyn | comp106451 | 61 | T | 0.20 | 0.36 | 0.48 | 0.35 | 0.84 | 0.83 | 0.72 | 0.79 | **0.24** |
| Protein-tyrosine kinase 2-beta | ENSMUSG00000059456 | syn | Contig_82228 | 496 | C | 0.49 | 0.55 | 0.64 | 0.57 | 0.88 | 1.00 | 0.91 | 0.91 | **0.24** |
| PILR Alpha Associated Neural Protein (M) | ENSMUSG00000030329 | UTR | comp103257 | 734 | C | 0.50 | 0.25 | 0.50 | 0.48 | 0.78 | 0.90 | 0.74 | 0.75 | **0.24** |
| Uncharacterized protein KIAA1211-like | - | UTR | comp113945 | 172 | G | 0.15 | 0.60 | 0.46 | 0.55 | 0.84 | 0.94 | 0.91 | 0.90 | **0.24** |
| Integrin Alpha 7 Chain (M) | ENSMUSG00000025348 | UTR | comp131887 | 6218 | G | 1.00 | 1.00 | 1.00 | 1.00 | 0.50 | 0.66 | 0.59 | 0.76 | **0.24** |

| Gene | Ensembl ID | Type | Contig | Pos | Allele | | | | | | | | | Score |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| RAB6B, member RAS oncogene family (M) | ENSMUSG00000032549 | UTR | comp79801 | 1930 | T | 0.48 | 0.60 | 0.50 | 0.34 | 0.92 | 0.95 | 0.84 | 0.90 | **0.24** |
| Transmembrane Emp24 Protein Transport Domain Containing 4, 3'UTR (M) | ENSMUSG00000004394 | nonORF | Contig_56898 | 1329 | T | 0.94 | 1.00 | 0.94 | 0.94 | 0.52 | 0.70 | 0.65 | 0.59 | **0.24** |
| Plastin-1 | ENSMUSG00000049493 | UTR | Contig_80587 | 317 | A | 0.24 | 0.33 | 0.17 | 0.11 | 0.88 | 0.88 | 0.57 | 0.88 | **0.24** |
| N-acetylglucosamine-1-phosphotransferase subunits alpha/beta | ENSMUSG00000035311 | nsyn | comp137737 | 2603 | G | 0.94 | 1.00 | 0.92 | 0.97 | 0.69 | 0.69 | 0.25 | 0.52 | **0.23** |
| Receptor expression-enhancing protein 5 | ENSMUSG00000005873 | UTR | comp135739 | 549 | G | 0.81 | 1.00 | 0.93 | 0.90 | 0.53 | 0.58 | 0.54 | 0.57 | **0.23** |
| Yippee-Like 2 (M) | ENSMUSG00000018427 | UTR | Contig_4209 | 1638 | C | 0.93 | 1.00 | 0.96 | 0.92 | 0.67 | 0.69 | 0.61 | 0.65 | **0.23** |
| Plexin-D1 | ENSMUSG00000030123 | UTR | comp3292 | 828 | G | 0.45 | 0.31 | 0.33 | 0.36 | 0.71 | 0.68 | 0.77 | 0.71 | **0.23** |
| Angiomotin-like protein 2 | ENSMUSG00000032531 | UTR | comp104477 | 3147 | A | 0.33 | 0.21 | 0.44 | 0.43 | 0.71 | 0.94 | 0.87 | 0.67 | **0.23** |
| Serine protease HTRA1 | ENSMUSG00000006205 | syn | comp78712 | 1041 | C | 0.49 | 0.52 | 0.54 | 0.56 | 1.00 | 0.79 | 0.83 | 0.85 | **0.23** |
| H-2 class II histocompatibility antigen, E-K alpha chain | - | syn | Contig_75434 | 605 | A | 0.70 | 0.12 | 0.40 | 0.71 | 1.00 | 1.00 | 0.94 | 1.00 | **0.23** |
| Gamma-aminobutyric acid receptor subunit gamma-2 | ENSMUSG00000020436 | UTR | comp125841 | 972 | C | 0.45 | 0.60 | 0.41 | 0.62 | 1.00 | 0.87 | 1.00 | 0.85 | **0.23** |
| Golgi SNAP receptor complex member 1 | ENSMUSG00000010392 | UTR | Contig_2352 | 1107 | T | 0.96 | 1.00 | 1.00 | 0.89 | 0.61 | 0.66 | 0.60 | 0.62 | **0.23** |
| Carboxypeptidase D | ENSMUSG00000020841 | UTR | Contig_66587 | 1238 | C | 1.00 | 0.90 | 1.00 | 0.88 | 0.64 | 0.65 | 0.61 | 0.53 | **0.23** |

| Gene | Ensembl ID | Type | Contig | Position | Allele | | | | | | | | | Value |
|------|-----------|------|--------|----------|--------|---|---|---|---|---|---|---|---|-------|
| **kallikrein-related peptidase 9, 3'UTR and upstrem DNA (M)** | ENSMUSG00000047884 | nonORF | comp110212 | 2106 | G | 0.32 | 0.60 | 0.27 | 0.50 | 0.94 | 0.97 | 0.84 | 0.83 | **0.23** |
| **D16Ertd472e (M)** | ENSMUSG00000022864 | nonORF | Contig_85509 | 75 | C | 0.27 | 0.36 | 0.32 | 0.33 | 0.62 | 0.85 | 0.59 | 0.83 | **0.23** |
| **-** | - | nonORF | Contig_74825 | 846 | T | 0.95 | 1.00 | 1.00 | 0.87 | 0.61 | 0.64 | 0.62 | 0.46 | **0.23** |
| **SLAIN motif-containing protein 2** | ENSMUSG00000036087 | UTR | comp130355 | 3281 | G | 0.79 | 0.88 | 0.69 | 0.86 | 0.17 | 0.14 | 0.46 | 0.40 | **0.23** |
| **Long-chain-fatty-acid--CoA ligase ACSBG1** | ENSMUSG00000032281 | syn | comp87279 | 1043 | C | 0.35 | 0.13 | 0.42 | 0.41 | 0.78 | 0.98 | 0.64 | 0.67 | **0.22** |
| **SHC-transforming protein 3, 3'UTR (M)** | ENSMUSG00000021448 | syn | comp138222 | 3272 | G | 0.84 | 0.74 | 0.89 | 0.76 | 0.44 | 0.52 | 0.46 | 0.47 | **0.22** |
| **Transient receptor potential cation channel subfamily M member 4 (M)** | ENSMUSG00000038260 | syn | comp5057 | 1345 | A | 0.76 | 0.82 | 0.72 | 0.70 | 0.44 | 0.48 | 0.00 | 0.41 | **0.22** |
| **Serpin B6** | ENSMUSG00000060147 | syn | Contig_96291 | 4880 | T | 0.88 | 0.96 | 1.00 | 0.89 | 0.42 | 0.66 | 0.65 | 0.33 | **0.22** |
| **Glutamate receptor 1, intron (M)** | ENSMUSG00000020524 | UTR | comp134072 | 3100 | G | 0.48 | 0.40 | 0.52 | 0.63 | 1.00 | 0.92 | 1.00 | 0.85 | **0.22** |
| **CD5 antigen-like** | ENSMUSG00000015854 | UTR | Contig_60567 | 1174 | C | 0.52 | 0.42 | 0.57 | 0.68 | 0.90 | 0.99 | 0.94 | 0.93 | **0.22** |
| **-** | - | nonORF | Contig_4861 | 974 | C | 0.50 | 0.38 | 0.29 | 0.21 | 0.72 | 0.72 | 1.00 | 0.76 | **0.22** |
| **Interferon gamma receptor 1** | ENSMUSG00000020009 | nsyn | comp121809 | 800 | G | 0.53 | 0.48 | 0.59 | 0.40 | 0.84 | 0.90 | 0.80 | 0.92 | **0.21** |
| **Cytochrome P450 26A1** | ENSMUSG00000024987 | syn | Contig_67458 | 1590 | C | 0.94 | 0.95 | 0.79 | 0.95 | 0.48 | 0.58 | 0.58 | 0.12 | **0.21** |
| **FCH Domain Only 1 (M)** | ENSMUSG00000070000 | UTR | comp119681 | 1264 | T | 0.37 | 0.36 | 0.44 | 0.57 | 0.89 | 0.93 | 0.78 | 0.89 | **0.21** |

| Gene | Ensembl ID | Type | Contig | Pos | Allele | | | | | | | | | Value |
|------|-----------|------|--------|-----|--------|---|---|---|---|---|---|---|---|-------|
| **Ridine-Cytidine Kinase 2 (M)** | ENSMUSG00000026558 | UTR | comp4443 | 2380 | C | 0.45 | 0.47 | 0.35 | 0.29 | 0.74 | 0.70 | 0.80 | 0.68 | **0.21** |
| **Glutamate receptor 1, intron (M)** | ENSMUSG00000020524 | UTR | comp134072 | 2746 | T | 0.57 | 0.35 | 0.70 | 0.58 | 1.00 | 0.94 | 1.00 | 0.91 | **0.21** |
| **-** | - | UTR | Contig_30676 | 712 | T | 0.92 | 1.00 | 1.00 | 1.00 | 0.70 | 0.71 | 0.69 | 0.69 | **0.21** |
| **-** | - | nonORF | comp91703 | 509 | C | 0.96 | 0.80 | 0.85 | 0.83 | 0.56 | 0.43 | 0.50 | 0.59 | **0.21** |
| **neurofibromin 2, 3'UTR (M)** | ENSMUSG00000009073 | nonORF | comp69409 | 1089 | A | 0.63 | 0.72 | 0.69 | 0.59 | 1.00 | 0.93 | 1.00 | 1.00 | **0.21** |
| **Insulin-Like Growth Factor 11 (M)** | ENSMUSG00000022790 | nsyn | comp136218 | 2490 | G | 0.65 | 0.73 | 0.83 | 0.69 | 0.22 | 0.44 | 0.18 | 0.43 | **0.21** |
| **CDK5 regulatory subunit-associated protein 3** | ENSMUSG00000018669 | nsyn | Contig_50394 | 762 | C | 0.52 | 0.68 | 0.74 | 0.68 | 0.95 | 0.99 | 1.00 | 0.98 | **0.21** |
| **Ras-related protein Rab-13** | ENSMUSG00000027935 | syn | comp72964 | 241 | C | 0.72 | 0.67 | 0.78 | 0.79 | 1.00 | 1.00 | 1.00 | 1.00 | **0.21** |
| **Deoxynucleotidyltransf erase terminal-interacting protein 1** | ENSMUSG00000017299 | syn | comp57960 | 513 | G | 0.41 | 0.61 | 0.41 | 0.68 | 0.97 | 0.93 | 1.00 | 0.89 | **0.21** |
| **Ral GTPase-activating protein subunit beta** | ENSMUSG00000027652 | syn | comp3196 | 5159 | C | 0.55 | 0.55 | 0.61 | 0.42 | 1.00 | 0.86 | 0.82 | 0.95 | **0.21** |
| **Liver carboxylesterase 1** | ENSMUSG00000071047 | syn | Contig_788 | 91 | C | 0.69 | 0.70 | 0.76 | 0.68 | 1.00 | 1.00 | 0.97 | 1.00 | **0.21** |
| **Nesprin-1** | ENSMUSG00000096054 | syn | Contig_96388 | 3389 | C | 0.77 | 0.89 | 0.74 | 0.93 | 0.48 | 0.26 | 0.42 | 0.53 | **0.21** |
| **GTP-binding nuclear protein Ran** | ENSMUSG00000029430 | UTR | comp115810 | 1495 | G | 0.74 | 0.44 | 0.56 | 0.70 | 1.00 | 1.00 | 0.95 | 0.99 | **0.21** |
| **Coagulation factor V** | ENSMUSG00000026579 | UTR | Contig_91316 | 7330 | C | 0.79 | 0.71 | 0.71 | 0.59 | 1.00 | 1.00 | 1.00 | 1.00 | **0.21** |
| **Kallikrein-related peptidase 9, 3'UTR and upstrem DNA (M)** | ENSMUSG00000047884 | nonORF | comp110212 | 2535 | C | 0.38 | 0.50 | 0.38 | 0.54 | 0.85 | 0.78 | 0.75 | 0.82 | **0.21** |

| | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Neurofibromin 2, 3'UTR (M)** | ENSMUSG00000009073 | nonORF | comp69409 | 320 | G | 0.67 | 0.60 | 0.62 | 0.63 | 1.00 | 0.88 | 1.00 | 1.00 | **0.21** |
| - | - | nonORF | Contig_4861 | 977 | T | 0.51 | 0.37 | 0.28 | 0.22 | 0.75 | 0.72 | 0.83 | 0.75 | **0.21** |

**Tab. S2** Genes differentially expressed (FDR < 0.05) between predatory and control lines in hippocampus. Presented are only genes having moderate to high expression (mean FPKM > 1), which were annotated to known genes in SwissProt database, or manually annotated to using mouse genome (M). For each line FPKM value TMM-normalized is given and false discovery rate value is provided.

| Gene annotation | Ensembl symbol of mouse gene | Contig name | C1 | C2 | C3 | C4 | P1 | P2 | P3 | P4 | FDR |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 3'UTR of Gamma-aminobutyric acid B receptor 1 (M) | ENSMUSG00000024462 | comp125314 | 15.35 | 11.47 | 14.66 | 14.62 | 22.62 | 23.45 | 19.46 | 18.27 | 1.93E-02 |
| Hemoglobin subunit beta-1 | ENSMUSG00000052305 | comp61009 | 0 | 0 | 40.85 | 0 | 0 | 0 | 0 | 0 | 9.34E-03 |
| Protein NDRG4-A | ENSMUSG00000036564 | comp82808 | 6.03 | 6.4 | 5.4 | 6.21 | 3.84 | 4.23 | 3.9 | 3.23 | 6.85E-03 |
| - | - | comp115007 | 2.82 | 3.99 | 2.59 | 3.35 | 5.25 | 4.87 | 5.95 | 6.5 | 6.93E-03 |
| Hemoglobin subunit alpha-1 | ENSMUSG00000069919 | comp55968 | 0 | 0 | 33.12 | 0 | 0 | 0 | 0 | 0 | 1.61E-02 |
| Death-associated protein 1 | ENSMUSG00000039168 | comp80742 | 2.29 | 2.6 | 2.57 | 2.44 | 3.69 | 3.31 | 4.18 | 4.36 | 1.54E-02 |
| Peptidylprolyl isomerase (cyclophilin)-like 2 (M) | ENSMUSG00000022771 | comp106583 | 1.38 | 1.14 | 1.07 | 0.68 | 3.77 | 2.54 | 2.1 | 1.82 | 1.94E-02 |
| Sin3A associated protein (M) | ENSMUSG00000024260 | comp102628 | 1.06 | 0.9 | 0.75 | 1.54 | 1.81 | 2.21 | 3.61 | 2 | 2.22E-02 |

| | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **SH3/ankyrin domain gene 1 (M)** | ENSMUSG00000038738 | comp113420 | 2.32 | 3.13 | 1.53 | 1.72 | 1.19 | 1.09 | 0.72 | 1.02 | 2.04E-02 |
| **Interferon-induced GTP-binding protein Mx2** | ENSMUSG00000023341 | comp66523 | 1.08 | 0.22 | 0.53 | 0.68 | 1.75 | 1.75 | 3.06 | 1.81 | 5.31E-04 |
| **F-box/WD repeat-containing protein 8** | ENSMUSG00000032867 | comp80437 | 1.77 | 2.04 | 1.55 | 1.51 | 1.19 | 0.84 | 0.67 | 0.82 | 7.22E-03 |
| **-** | - | comp65388 | 0.32 | 0 | 0.4 | 0.69 | 1.29 | 1.95 | 3.49 | 1.31 | 7.22E-03 |

**Tab S3** Genes differentially expressed (FDR < 0.05) between predatory and control lines in liver. Presented are only genes having moderate to high expression (mean FPKM > 1), which were annotated to known genes in SwissProt database, or manually annotated to using mouse genome (M). For each line FPKM value TMM-normalized is given and false discovery rate value is provided

| Gene annotation | Ensembl symbol of mouse gene | Contig name | C1 | C2 | C3 | C4 | P1 | P2 | P3 | P4 | FDR |
|---|---|---|---|---|---|---|---|---|---|---|---|
| NADP-dependent malic enzyme | ENSMUSG00000032418 | Contig_69956 | 65.15 | 51.29 | 34.37 | 78.22 | 32.43 | 32.49 | 20.11 | 17.36 | 1.94E-02 |
| G0/G1 switch gene 2 (M) | ENSMUSG00000009633 | Contig_41592 | 47.31 | 49.81 | 55.77 | 58.93 | 12.48 | 23.86 | 29.49 | 27.11 | 3.96E-03 |
| Ribosome binding protein 1 (M) | ENSMUSG00000027422 | Contig_88901 | 17.9 | 21.79 | 22.32 | 19.81 | 33.72 | 36 | 33.49 | 40.91 | 3.69E-02 |
| - | - | Contig_88895 | 12.15 | 18.39 | 17.86 | 11.67 | 25.29 | 27.58 | 34 | 36.27 | 1.67E-02 |
| Phospholipase A2 | ENSMUSG00000056220 | Contig_47128 | 24.59 | 69.27 | 3.58 | 0.44 | 0.13 | 1.93 | 1.26 | 2.12 | 4.89E-02 |
| - | - | Contig_42802 | 0.84 | 0.32 | 3.33 | 1.22 | 12.35 | 18.9 | 21.03 | 19.84 | 2.05E-13 |
| Testosterone 17-beta-dehydrogenase 3 | ENSMUSG00000033122 | Contig_79364 | 4.42 | 7.05 | 6.72 | 5.02 | 14.68 | 10.25 | 11.51 | 10.42 | 1.62E-02 |
| - | - | Contig_56236 | 0 | 32.39 | 9.39 | 0 | 0 | 0.07 | 0 | 0 | 1.70E-02 |
| - | - | Contig_55080 | 2.84 | 0.28 | 1.76 | 1.49 | 6.14 | 4.73 | 11.88 | 7.94 | 1.36E-03 |
| Poly [ADP-ribose] polymerase 14 | ENSMUSG00000034422 | Contig_96742 | 2.93 | 1.47 | 2.13 | 2.19 | 9.19 | 8.38 | 5.58 | 2.21 | 2.15E-02 |
| IQ motif containing G (M) | ENSMUSG00000035578 | Contig_83221 | 5.12 | 7.06 | 5.88 | 3.5 | 1.66 | 1.79 | 3.29 | 2.1 | 4.25E-02 |
| - | - | Contig_24041 | 1.44 | 3.14 | 2.42 | 2.82 | 5.38 | 5.55 | 4.47 | 5.14 | 4.08E-02 |
| - | - | Contig_41830 | 0.99 | 0.14 | 1.07 | 1.47 | 3.6 | 5.49 | 6.29 | 2.45 | 6.59E-04 |
| Protein phosphatase 1 regulatory subunit 3G | ENSMUSG00000050423 | Contig_66518 | 1.38 | 1.57 | 1.23 | 1.41 | 6.15 | 2.52 | 2.65 | 4.45 | 2.56E-03 |
| - | - | Contig_46518 | 1.24 | 1.93 | 0.97 | 0.99 | 4.08 | 3.97 | 4.25 | 2.96 | 1.53E-04 |
| - | - | Contig_42514 | 1.22 | 0.17 | 0.49 | 1.56 | 5.5 | 3.46 | 5.22 | 1.37 | 1.94E-02 |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| - | - | Contig_44016 | 1.73 | 0 | 0.63 | 0.44 | 3.77 | 2.95 | 5.02 | 3.86 | 3.96E-03 |
| - | - | Contig_101471 | 16.59 | 0 | 0.33 | 0 | 0 | 0 | 0.06 | 0 | 2.76E-02 |
| - | - | Contig_3442 | 1.49 | 1.1 | 1.59 | 1.11 | 2.88 | 3.32 | 1.83 | 2.93 | 1.62E-02 |
| - | - | Contig_194 | 1.49 | 0.58 | 0.49 | 0.61 | 1.55 | 2.36 | 3.78 | 4.28 | 3.96E-03 |
| - | - | Contig_81892 | 5.72 | 2.75 | 3.15 | 1.39 | 0.26 | 0.79 | 0.39 | 0.67 | 9.49E-06 |
| - | - | Contig_43818 | 1.31 | 0.62 | 0.37 | 0 | 2.56 | 2.85 | 2.86 | 3.72 | 6.97E-03 |
| Transcriptional repressor CTCFL | ENSMUSG0000007049 | Contig_65729 | 0.51 | 1.03 | 1.1 | 0.91 | 2.49 | 2.41 | 3.28 | 1.82 | 1.94E-02 |
| - | - | Contig_83445 | 0.88 | 0.96 | 0.77 | 0.45 | 3.36 | 2.53 | 1.67 | 2.86 | 2.56E-03 |
| - | - | Contig_56233 | 0 | 9.48 | 2.96 | 0 | 0 | 0.01 | 0 | 0 | 1.06E-02 |
| - | - | Contig_65480 | 0.67 | 0.49 | 0.61 | 0.46 | 1.77 | 2.16 | 1.92 | 2.88 | 2.58E-02 |
| Protease, serine 35 (M) | ENSMUSG00000033491 | Contig_51709 | 2.41 | 2.39 | 1.55 | 2.54 | 0.78 | 0.38 | 0.54 | 0.04 | 1.59E-04 |
| Aryl hydrocarbon receptor nuclear translocator 2 | ENSMUSG00000015709 | Contig_1448 | 1.16 | 0.6 | 0.85 | 0.59 | 2.07 | 2.06 | 1.41 | 1.48 | 2.20E-02 |
| - | - | Contig_103697 | 0.48 | 0.53 | 0.38 | 0.47 | 3.21 | 2.62 | 1.37 | 1.04 | 1.94E-02 |
| Transient receptor potential cation channel, subfamily C, member 1 (M) | ENSMUSG00000032839 | Contig_34437 | 0.91 | 0.75 | 0.58 | 0.59 | 2.29 | 2.03 | 1.34 | 1.39 | 2.90E-02 |
| - | - | Contig_67720 | 0.75 | 0.36 | 0.48 | 0.62 | 1.47 | 1.83 | 2.36 | 1.68 | 3.01E-03 |
| - | - | Contig_97564 | 0.09 | 0.2 | 0.86 | 0.3 | 1.94 | 3.21 | 1.04 | 1.81 | 3.59E-02 |
| - | - | Contig_27010 | 0.4 | 0.13 | 0.55 | 0.67 | 1.95 | 2.45 | 1.51 | 0.97 | 1.93E-02 |
| Leptin receptor (M) | ENSMUSG00000057722 | Contig_1474 | 0.43 | 0.5 | 0.48 | 0.15 | 2.51 | 2.12 | 0.86 | 1.56 | 1.59E-04 |
| - | - | Contig_65292 | 0 | 0 | 0 | 0 | 1.73 | 2.4 | 2.22 | 1.87 | 5.38E-11 |

# ACKNOWLEDGMENTS

**STATEMENTS OF CO-AUTHORS**