

The Functional Annotation of Mammalian Genomes: The Challenge of Phenotyping

Steve D.M. Brown,¹ Wolfgang Wurst,² Ralf Kühn,² and John M. Hancock¹

¹MRC Mammalian Genetics Unit, MRC Harwell, Harwell Science and Innovation Campus, Oxfordshire OX11 0RD, United Kingdom; email: s.brown@har.mrc.ac.uk, j.hancock@har.mrc.ac.uk

²Helmholtz Center Munich, German Research Center for Environmental Health, 85764 Munich, Germany; email: wurst@helmholtz-muenchen.de, ralf.kuehn@helmholtz-muenchen.de

Annu. Rev. Genet. 2009. 43:305–33

First published online as a Review in Advance on August 18, 2009

The *Annual Review of Genetics* is online at genet.annualreviews.org

This article's doi:
10.1146/annurev-genet-102108-134143

Copyright © 2009 by Annual Reviews.
All rights reserved

0066-4197/09/1201-0305\$20.00

Key Words

mouse genetics, mutagenesis, phenotype, bioinformatics, ontologies

Abstract

The mouse is central to the goal of establishing a comprehensive functional annotation of the mammalian genome that will help elucidate various human disease genes and pathways. The mouse offers a unique combination of attributes, including an extensive genetic toolkit that underpins the creation and analysis of models of human disease. An international effort to generate mutations for every gene in the mouse genome is a first and essential step in this endeavor. However, the greater challenge will be the determination of the phenotype of every mutant. Large-scale phenotyping for genome-wide functional annotation presents numerous scientific, infrastructural, logistical, and informatics challenges. These include the use of standardized approaches to phenotyping procedures for the population of unified databases with comparable data sets. The ultimate goal is a comprehensive database of molecular interventions that allows us to create a framework for biological systems analysis in the mouse on which human biology and disease networks can be revealed.

INTRODUCTION

The sequencing of the human and mouse genomes has transformed the landscape of mammalian biology, providing a fundamental underpinning to ongoing developments in functional, population, and evolutionary genomics. The mouse and human genome sequences have allowed us to determine a fairly accurate picture of the nature and content of protein-coding genes in these two mammalian genomes (70, 89). These analyses reinforce the high levels of similarity between the two genomes—approximately 99% of mouse genes have homologues in the human genome (25). Recent analyses have also revealed the extraordinary breadth of noncoding transcripts in mammalian genomes (11, 29) and have underlined that a much larger fraction of the genome is transcribed than had heretofore been realized. Despite the extensive sequence and transcript annotation that has taken place, the extent of our knowledge of the function of both coding and noncoding loci in the mammalian genome is woeful. Studies of the role and function of the various classes of noncoding transcripts are in their infancy. But even for protein-coding loci, we are some way from developing a comprehensive picture of gene function.

Undertaking a comprehensive functional annotation of the mouse genome will be pivotal to developing a systems biology of mammals—which in itself will be critical if we are to understand how transformations in genetic networks lead to disease. We need to establish through model organisms, such as the mouse, the functions and interactions of genetic loci, carry out perturbations, by mutation or other approaches, and investigate the phenotypic outcome. The ultimate goal is to create a comprehensive database of molecular interventions *in vivo*. This database will be a central, though not the only, element in developing a framework for biological systems in the mouse through which human biology and disease networks can be revealed.

The mouse is the organism of choice to undertake a comprehensive functional annotation of a mammalian genome. It has a number of

distinctive attributes that underline its role as a key model organism in functional genomic studies (42). First, an extensive genetic toolkit has been developed that allows the generation of a wide variety of mutational lesions. This adaptability in generating a range of mutant alleles underpins recent efforts to establish mutations for every gene in the mouse genome (38). Ultimately, we can expect to produce a resource with multiple mutant alleles at every mouse locus. Second, the mouse has been an important experimental organism for over a century, resulting in a depth and breadth of knowledge of developmental, physiological, behavioral, and biochemical mechanisms that underpins continuing investigations as a model for human biological processes. Third, given its relatively small size and short generation time, the mouse remains the most economically viable mammalian organism for large-scale functional genomic studies. This combination of features means that the mouse is uniquely placed for the development and analysis of models for human disease, particularly for systems that are not possible or relevant to model in other organisms.

The first step in addressing the challenge of delivering a functional annotation of the mouse genome has been to put in place programs to generate comprehensive mutant resources (4, 5): A worldwide effort is now under way to generate new libraries of mutant alleles for the bulk of mouse genes. However, this effort, which should be substantially realized in the next five years, is only the beginning. The biggest obstacle to undertaking a comprehensive determination of the function of genes remains the determination of phenotype. In this article, we review the complementary approaches that are critical for generating a rich diversity of alleles at mouse loci. However, we pay considerable attention to the many and formidable challenges that need to be met if these extensive mutant resources are to be characterized in a manner that will contribute to a profound knowledge of mammalian biological systems. There need to be considerable developments in the science of phenotyping as well as significant

Functional annotation:

gaining knowledge about the function of individual genes to provide a detailed picture of the functions of all genes

advances in technologies, infrastructure, and logistics that underpin phenotyping. We review and discuss current progress in mouse phenotyping as well as future areas for development, including the many critical informatics challenges that lie ahead.

GENE-DRIVEN MUTAGENESIS

Gene Targeting in Embryonic Stem Cells

Gene targeting allows the introduction of pre-designed, site-specific modifications into the genome of embryonic stem (ES) cells by homologous recombination (28). It is extensively used for the preplanned disruption of genes in the murine germline resulting in mutant knockout mouse strains. Since the first demonstration of homologous recombination in ES cells in 1987 (137), more than 4000 knockout mouse strains have been generated. Gene inactivation is achieved through the insertion of a selectable marker into an exon of the target gene or the replacement of one or more exons. The mutant allele is initially assembled in a specifically designed gene targeting vector such that the selectable marker is flanked at both sides with genomic segments derived from the target gene that serve as homology regions to initiate homologous recombination. Upon electroporation of such a vector into ES cells and the selection of stable integrants, clones that underwent homologous recombination can be identified through the analysis of genomic DNA using polymerase chain reaction or Southern blotting. Upon the isolation of recombinant ES cell clones, modified ES cells are injected into blastocysts to transmit the mutant allele through the germline of chimeras and to establish a mutant mouse strain.

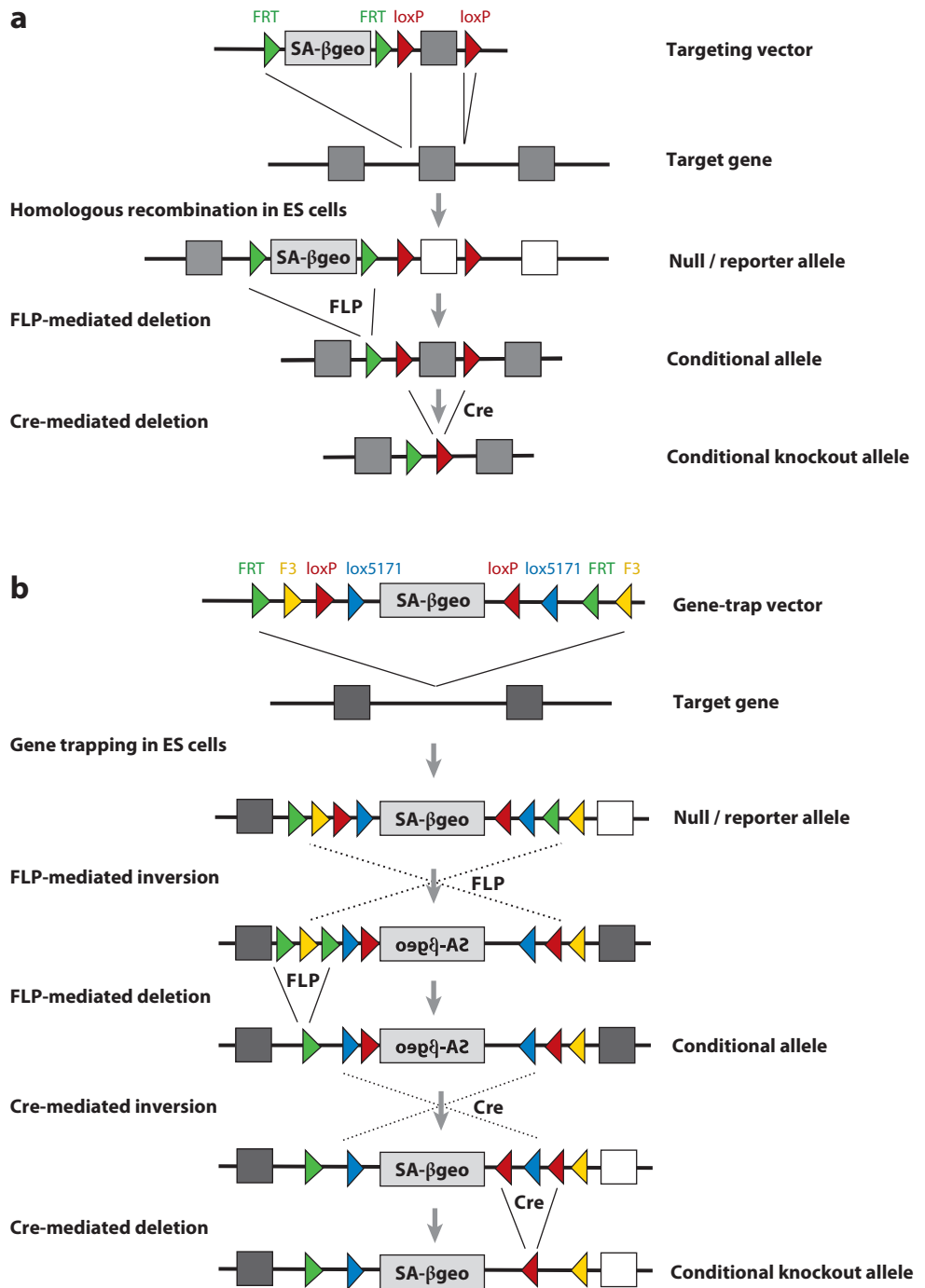
Using this classical gene targeting approach, researchers can obtain germline mutants that harbor the knockout mutation in all cells throughout development. Approximately 30% of all targeted genes are essential for embryonic development such that their inactivation leads to lethality, precluding further analysis in

adult mice. To avoid embryonic lethality and to study gene function only in specific cell types, Gu et al. (64) introduced a modified, conditional gene targeting scheme that allows researchers to restrict gene inactivation to specific cell types or developmental stages. In a conditional mutant, gene inactivation is achieved by the insertion of two 34-base pair (bp) recognition (loxP) sites of the DNA recombinase Cre into introns of the target gene such that recombination results in the deletion of loxP-flanked exons (113). Recombination and gene inactivation are achieved by crossing the strain harboring the conditional (loxP-flanked) allele with a transgenic strain expressing Cre recombinase in one or several cell types. Target gene inactivation occurs in a spatially and temporally restricted manner, according to the pattern of Cre expression. More than 100 Cre transgenic strains with tissue-specific recombinase expression have been generated covering a large range of cell types (101). In addition, gene inactivation can also be induced in adult mice using Cre transgenes activated by small molecule inducers (16, 53, 86). The generation of conditional alleles involves the same technology as the production of germline knockouts (87, 139). For a conditional gene targeting vector, a selection marker and a loxP site are inserted into one intron of the target gene while a second loxP sequence is placed into another intron (**Figure 1a**). Upon homologous recombination and germline transmission of the modified allele, the selection marker gene, flanked by FLP recombinase recognition (FRT) sites, can be removed by a cross with a strain that expresses FLP in germ cells (118). Upon removal of the selection marker, the loxP-flanked conditional allele can be bred to the required Cre transgene to obtain conditional mutants for phenotype analysis (**Figure 1a**).

Gene targeting has largely progressed in a one-by-one manner by the contributions of a large number of laboratories. As technology advanced, it became possible to produce targeted mutants more quickly and at a larger scale. Via the use of DNA-engineering strategies that rely on homologous recombination in

Gene targeting: the introduction of pre-designed, site-specific modifications into the genome of embryonic stem cells by homologous recombination

FRT: FLP recombinase recognition target



bacteria, conditional gene targeting vectors can be assembled in a high-throughput format. Starting from genomic BAC (bacterial artificial chromosome) clones, this technology, termed ET cloning or recombineering (41, 100), allows researchers to derive gene targeting vectors in a multiwell format within a few days (32). Furthermore, the modified gene targeting strategy, termed targeted trapping, allows for the isolation of recombined ES cell clones at a much higher frequency (~50%) versus that obtained with the traditional vector design. As shown in **Figure 1a**, a targeted trapping vector combines gene homology arms with a promoterless selection/reporter cassette flanked by FRT sites. The selection cassette is preceded by a splice acceptor sequence that traps the gene's mRNA such that drug resistance is acquired upon successful targeting but only rarely by random integration. It has been estimated that 60–70% of all genes are sufficiently expressed in ES cells to support mutagenesis by targeted trapping. These developments and the availability of reliable ES cell lines of C57BL/6 origin enable the pursuit of gene targeting at a large scale in the preferred genetic background. As part of the International Knockout Mouse Consortium, three international collabora-

tion projects (see <http://www.eucomm.org>, <http://www.knockoutmouse.org>, and <http://www.norcomm.org>) using a combination of gene traps, targeted trapping, and gene targeting are progressing to provide a resource of mutant ES cell clones for every gene within the next several years (38).

Gene-Trap Mutagenesis in Embryonic Stem Cells

Gene-trap mutagenesis is based on the random integration of a gene-trap vector across the genome of ES cells and the disruption of genes through vector-specific elements. Such vectors mutate a gene at the site of insertion, provide a sequence tag for the identification of the disrupted gene, and mimic the expression of the tagged gene by a reporter gene. A typical gene-trap vector contains a promoterless reporter-selector cassette that functions by generating a fusion transcript with the endogenous gene. The most widely used β geo cassette contains an ATG-less hybrid coding region for the β -galactosidase reporter and the neomycin phosphotransferase selection marker preceded by a splice acceptor element. If the translational reading frames of the trapped transcript

Mutagenesis: the introduction of mutations into the germline of a mouse line or strain using either gene-driven or random, chemical approaches



Figure 1

Conditional gene targeting and gene-trap mutagenesis in embryonic stem (ES) cells and mice. (a) Vector for targeted trapping mutagenesis, including a gene disruption cassette (SA- β geo) flanked by FLP recombinase recognition target (FRT) sites (*green triangles*) and a loxP site (*red triangle*), inserted upstream of the exon of interest (*square*), and a second loxP site of the same orientation inserted downstream of the exon. The SA- β geo cassette consists of a splice acceptor (SA) sequence, a fusion protein of β -galactosidase and the neomycin phosphotransferase (β geo), and a polyadenylation signal sequence (not shown). Upon the isolation of drug-resistant colonies and the identification of homologous recombined ES cell clones, knockout mice harboring the null allele can be established. Upon deletion of the SA- β geo cassette by FLP-mediated recombination, a new allele for conditional mutagenesis is established. In combination with a strain expressing Cre recombinase in a specific cell type, conditional knockout mice are obtained. Transcribed exons are shown as shaded squares; exons silenced by the gene disruption cassette are shown as open squares. (b) Vector for conditional gene-trap mutagenesis that includes a SA- β geo cassette flanked by a pair of (incompatible) FRT and F3 recognition sites for FLP recombinase and a pair of (incompatible) loxP and lox5171 recognition sites for Cre recombinase (FLEEx strategy) (123). The pairs of recognition sites are placed in opposite orientation to each other. Vector integration into an intron creates a standard gene-trap (null) allele. The SA- β geo cassette can be inverted by FLP-mediated inversion and deletion between the pairs of FRT and F3 sites, creating a conditional allele. Mice harboring the conditional allele can be crossed with a strain expressing Cre recombinase in a specific cell type to generate conditional knockout mice via recombination between the pairs of loxP and lox5171 sites. Transcribed exons are shown as shaded squares; exons silenced by the gene disruption cassette are shown as open squares.

Gene trapping:

random integration of a gene-trap vector across the genome of embryonic stem cells

and the β geo cassette are in line, a fusion protein is produced that confers drug resistance to the ES cell clone. Upon vector introduction, ES cell cultures are selected for drug resistance such that only clones harboring a productive vector integration into an active gene can survive. This stringent selection scheme and the use of a single vector to hit a large number of genes are the basis for the high efficiency of the gene-trap approach because each resistant colony represents an independent integration event into a gene (55, 57, 68, 131, 153). Thus, libraries of mutant ES cell clones can be established rapidly and at low costs per mutant. The resulting collections of mutant genes provide the basis for the establishment of mutant mouse strains through germline chimeras raised from selected ES cell clones. The overall analysis of homozygous mouse mutants derived from gene-trap ES cell clones revealed obvious phenotypes and embryonic lethality at a frequency comparable to mutants generated by gene targeting, indicating that gene-trap insertions typically result in null alleles (68, 131).

Classical gene-trap vectors irreversibly modify target genes and create germline null mutations that can lead to embryonic lethality. A new type of gene-trap vector also allows conditional gene inactivation by Cre/loxP recombination (**Figure 1b**) and combines the advantages of gene-trap and conditional mutagenesis. A conditional gene-trap vector is designed like a classical vector, except that the β geo reporter-selector cassette is flanked with two pairs of loxP and two pairs of FRT sites that enable it to independently invert the cassette and provide a switch to activate or inactivate its function (122, 123). In its original orientation, the cassette disrupts the expression of the trapped gene, whereas the inverted cassette is nonmutagenic and enables gene expression. Irreversible inversion is obtained in a two-step recombination process that uses pairs of incompatible wild-type and mutant loxP/lox5171 as well as FRT/F3 sites, which must be ordered and oriented in a specific way (**Figure 1b**). In its initial configuration, the cassette disrupts the target gene and enables the

isolation of drug-resistant ES cell clones. Mice established with such an allele can be studied as null mutants and crossed to a strain expressing FLP recombinase in the germline to inactivate the mutagenic cassette by inversion, thus creating a conditional allele (**Figure 1b**). The silent gene-trap cassette can then be reactivated *in vivo* by crossing to a strain expressing Cre in a specific cell type, representing the conditional mutant for phenotyping. In addition, further modifications at the gene-trap locus in ES cells can be made to drive the expression of any foreign cDNA from the trapped gene.

Several centers within the academic research community currently run gene-trap screens that are combined in the International Gene Trap Consortium (<http://www.genetrap.org>). At present, this collection includes 380,863 ES cell clones, which represent 36% of all mouse genes (127). The sequence tags of all insertions are mapped on the Ensembl mouse genome server, providing a direct link to any gene of interest. Gene-trap mutagenesis has proven to be the most effective strategy to mutate a substantial fraction of all mouse genes. Therefore, gene trapping is an essential component of the European as well as the North American Conditional Mouse Mutagenesis initiatives for the complete mutagenesis of all mouse genes (38).

Transposon-Based Mutagenesis in Embryonic Stem Cells and Mice

Transposons are mobile genetic elements that have been used for many years as a tool for genetic analysis in prokaryotes and flies. The synthetic Tc1-like transposon *Sleeping beauty* (SB) (81) was first used in mice for gene-trap mutagenesis (36, 59) and to discover cancer genes (51, 52). To achieve insertional mutagenesis *in vivo*, gene disruption cassettes flanked with SB terminal repeats are first introduced as multi-copy transgenes into the mouse germline. Upon cross of such mice with a transgenic strain expressing SB transposase, the vector can be mobilized in somatic cells and inserted into new genomic sites. Such double-transgenic mice are highly prone to tumor development (51), and

the analysis of chromosomal integration sites in tumor cells revealed new cancer-related genes.

To achieve high rates of germline mutants through transposition of gene-trap-like vectors, the utility of the SB system is restricted by its low efficiency. To enable transposon-based genome-wide screens in mice, the insect-derived *piggyBac* (PB) system is more efficient in mammalian cells and can mobilize up to 9 kb of foreign sequence (48, 150). In double-transgenic mice with male germ cell-specific expression of PB transposase and a reporter gene flanked with PB recognition sites, one new insertion was found per gamete (48). The vector integration sites that are accessible through the PB mice database (132) show a wide chromosomal distribution and a preference for transcription units. The PB system has been further adapted to distribute a gene-trap vector together with a single loxP site across the genome of male germ cells (149). While this vector disrupts a single gene at the site of insertion, a collection of such mouse strains harboring loxP sites at various chromosomal positions provides a resource to manipulate chromosome segments (149). A Tamoxifen-inducible PB transposase fusion protein has been recently developed enabling the additional control of the timing of *in vivo* vector mobilization (27). Although the PB system was developed primarily for *in vivo* mutagenesis, it also offers new possibilities to mutagenize the genome of ES cells *in vitro* (144). Thus, gene-trap libraries established by retroviral delivery vectors may be complemented in the future by transposon-based technologies.

PHENOTYPE-DRIVEN MUTAGENESIS

Traditionally, geneticists have used phenotype-driven approaches to generate new mutant resources and to explore the relationship between gene and phenotype. This is sometimes referred to as forward genetics, but we prefer the term phenotype-driven genetics, which emphasizes the key nature of the approach. In contrast to the gene-driven approach where the gene

is the starting point, phenotype-driven mutagenesis employs efficient chemical mutagens or radiation to induce DNA lesions at random in the organism of interest. Mutagenized animals are screened for phenotypes of interest, and once the relevant phenotype has been identified, the mutation is then identified using a combination of genetic techniques, including mapping and sequencing. Mutation discovery in the phenotype-driven approach is not always straightforward. However, phenotype-driven mutagenesis has one important advantage: Unlike in the gene-driven approach, no *a priori* assumptions are made about the relationship among genes, genetic pathways, and a particular phenotype (24). The phenotype is the starting point in the discovery process, and having identified a relevant mutant phenotype, researchers characterize the underlying gene or genetic pathway. In this way, the phenotype-driven approach provides a powerful tool for identifying novel genes and genetic pathways associated with diverse physiological or disease states. For these reasons, a number of large-scale phenotype-driven mutagenesis programs have been initiated in the mouse over the past few years, focusing largely on the use of the chemical mutagen N-ethyl-N-nitrosourea (ENU) (20, 83).

ENU: A Powerful Chemical Mutagen

ENU is a chemical mutagen that acts as an alkylating agent, transferring its ethyl group to nucleophilic nitrogen or oxygen sites on the deoxyribonucleotides (8, 82, 104). The vast majority of ENU-induced variants are point mutations, but there is a nucleotide bias. Most commonly, mutations occur in A-T base pairs. Seventy to eighty-five percent of all ENU-induced nucleotide substitutions are either A-T to T-A transversions or A-T to G-C transitions (104). Conversely, the G-C to C-G transversion event is rarely seen (3). Following translation into proteins, these substitutions result in approximately 70% nonsynonymous changes, of which approximately 65% are

ENU: N-ethyl-N-nitrosourea

missense changes and the remainder are nonsense or splice mutations (82, 128).

ENU is administered as a series of intraperitoneal injections to male adult mice. Effective dose and injection regimes that optimize mutation rate vary according to genetic background (73, 82, 145). ENU acts on spermatogonial stem cells. Following ENU treatment, mice undergo a period of sterility due to depletion of differentiated spermatogonia. Following recovery of fertility, male mice are bred to produce progeny carrying ENU mutations. ENU has a high specific locus mutation rate that with optimum dose regimes is approximately 0.0015 (73). Given the high locus mutation rate, we can extrapolate that each gamete will carry ~30–50 functional mutations in coding sequences around the genome. In addition, each gamete will carry numerous other point mutations at other noncoding sequences with potential functional consequences. However, sequence-based approaches provide an increasingly clear picture of the power and spectrum of ENU mutagenesis across the genome. A number of studies indicate that ENU changes occur on average at a frequency of one mutation every 1–1.5 Mb (37, 39, 84, 95, 112, 120). In addition to its efficiency, the value of ENU lies in its ability to generate point mutations. It is possible to recover mutations demonstrating the full range of phenotypic effects from null mutations to partial loss of function (hypomorphs), gain of function (neomorphs), and dominant negatives (antimorphs). Indeed, ENU mutations in either coding or noncoding sequences will in many cases reflect more closely the kind of genetic variation found to be contributing to complex disease in humans. We return to this point below.

Dominant and Recessive Pipelines: From Mutagenesis to Gene

ENU mutagenesis can be used to uncover both dominant and recessive mutant phenotypes (23, 83). In the most simple strategy, mutagenized male mice are bred to generate G1 offspring that can be analyzed for dominant phenotypes.

To screen for recessive mutations, pedigrees are bred by intercrossing the offspring of a G1 individual (G2s) or crossing them back to the original G1, thus homozygosing mutations in a proportion of the resultant G3 offspring. Recessive mutations at specific chromosome regions can also be revealed by strategies involving crosses incorporating chromosome deletions (13, 115, 116) or balancer chromosomes (71, 85). Following identification of a new phenotype and confirmation via inheritance testing, the mutation is mapped, usually to a relatively small chromosome region of a few megabases or less and positionally cloned. Identification of the underlying mutation from the mapped chromosomal region is now relatively trivial given the annotated mouse genome sequence and the availability of high-throughput sequencing technologies. In addition, the advantage of ENU is that mutagenesis is carried out on an inbred background and the underlying mutation will represent the only coding sequence change in the mapped region—the likelihood of two ENU mutations occurring close together in the same region, for example, within 5Mb, is low (84).

ENU Mutagenesis and Phenotyping Screens

Perhaps more than any other mouse mutagenesis approach to date, the phenotype-driven nature of ENU mutagenesis pipelines has focused attention on the phenotyping platforms that are best suited for the efficient identification of mutant phenotypes. Although a single limited test may be employed to mutant mice emerging from dominant or recessive mutagenesis screens, there are economies of scale to be had from applying a battery of tests covering a variety of phenotypes and disease states. Indeed, the major ENU centers worldwide have driven much of the development in the thinking and application of mouse phenotype screens. We discuss these developments in phenotyping in detail below, but suffice it to say at this point that a wide diversity of systems have been assessed in ENU mutagenesis screens.

ENU screens have ranged broadly from developmental phenotypes to adult late-onset phenotypes, encompassing diverse systems and disease pathologies. Some of the early ENU screens applied relatively broad batteries of tests, aiming to identify large numbers of mutants across several phenotypic areas (35, 77, 92, 103, 136). Other mutagenesis pipelines have focused on particular phenotype domains or disease areas including immunological conditions (40, 142), bone (7, 76), circadian rhythms (6, 61), deafness and middle ear disease (69, 107, 126), and metabolic phenotypes (75, 79). Much effort has also been extended toward the use of recessive mutagenesis pipelines to identify developmental phenotypes (13, 72, 85). Some screens have applied challenges, such as an infection challenge, to reveal phenotypes (46). Overall, a plethora of novel phenotypes has been revealed. More importantly, a large number of these novel mutations has now been cloned, identifying novel gene function in many cases. In summary, ENU phenotype-driven screens have been a powerful tool for revealing novel gene function, and they can be expected to continue to play a role in adding to the diversity of mutant alleles from the mouse genome. Nonetheless, exciting opportunities remain for utilizing ENU in novel modalities that depend on recently developed approaches to ENU gene-driven mutagenesis.

ENU Gene-Driven Screens

In a seminal paper in 2002, researchers at MRC Harwell (37) demonstrated that ENU could be effectively employed in a gene-driven approach by identifying specific ENU changes in a gene of interest. Parallel archives of DNA and frozen sperm were generated from male mutants from the Harwell ENU mutagenesis pipeline. From these, an initial archive of 2230 mice was established. Via mutation scanning of the DNA archive, ENU mutants in a gene of interest are identified and the relevant mutant mouse can then be recovered by IVF from the sperm archive. Coghill et al. (37) applied this approach

to the recovery and characterization of several mutant alleles from the Connexin 26 gene.

Key goals in applying this technique were to establish the likely mutation rates and the size of the archive that would be required to recover a given number of alleles at any locus. Given a specific locus mutation rate of 0.0015, extrapolating genome-wide, we would expect a functional change of approximately 1 in every 2 Mb of coding sequence screened. Indeed, both Coghill et al. (37) and Quwailid et al. (112) in screening a much larger archive of approximately 6000 DNAs found a rate of functional mutations close to this predicted level: The mutation rate for all mutations was in the region of 1 in 1Mb (also see discussion above). Moreover, the probabilities of finding n or more mutant alleles of a gene in varying numbers of DNAs may be calculated from F1 offspring of male mutagenized mice (see **Figure 2**) (37). An archive of 5000 DNAs provides a very high frequency (>99%) of identifying a single allele, and the probability of recovering four or more alleles is significant (around 90%). A number of centers have established ENU DNA and sperm archives (3, 95, 120, 134), and overall, the numbers of alleles identified in screens to date bears

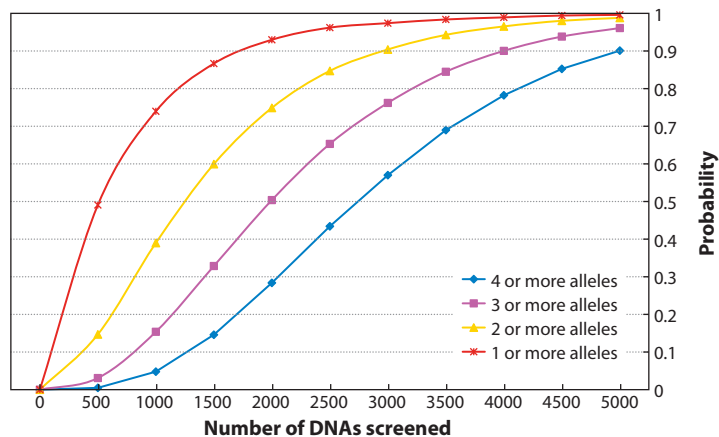


Figure 2

Summed probabilities of finding n or more mutant alleles of a gene in varying numbers of DNAs from F1 offspring of ENU mutagenized male mice. Adapted from Reference 37. Curves were calculated using BINOM version 1.72 (courtesy of J. Ott), assuming a mutation rate of 0.0015 mutations per locus that was corrected for a mutation detection rate of 90% (i.e., 0.00135 mutants/locus).

out these calculations. It can be extrapolated that an archive of 10,000 mice would provide 1 ENU change every 100 bp on average across the mouse genome.

In parallel, a number of groups have also demonstrated the utility of the ENU gene-driven approach in ES cell lines (33, 99), particularly in recovering large allelic series at loci of interest (143). The advantages here are the ability to optimize dose without giving consideration to animal welfare and, most importantly, the ability to generate large libraries of mutations without extensive animal breeding.

The utility of ENU gene-driven screens is thrown into fresh perspective by the success and scope of the recent plethora of genome-wide association studies in humans (26, 50, 56, 117, 138). These investigations are revealing a host of potential regions and genes that may underlie common, complex human disease states such as diabetes, Crohn's disease, obesity, autoimmune disease, and others. Providing a variety of mouse mutant alleles both in coding and non-coding sequences surrounding these loci will be imperative in order to validate loci and their involvement in the disease process. Generating null alleles via gene targeting will be an important first step in terms of assessing the function of these candidate loci. However, generating a variety of point mutations in both coding sequences and in surrounding promoter and regulatory sequences is likely to recapitulate better the genetic variation that segregates in the human population and that contributes to disease phenotypes. Screening of ENU archives can provide the necessary allelic variants. Moreover, as the throughput of sequencing technology improves and costs plummet, it is not fanciful to consider a major initiative to undertake the comprehensive genome resequencing of ENU archives. The aim would be to short-circuit the somewhat laborious process of mutation scanning and provide an *in silico* process of scanning for relevant mutations in any sequence of interest that can then be followed by the trivial and rapid recovery of the relevant mutant mouse from the sperm archive.

Although complete resequencing of all the thousands of DNAs in the current archive may be too ambitious (and too costly), a current alternative is to consider applying present developments in selection technologies (1, 74, 105, 111) toward, for example, capturing and resequencing all coding exons or some other selection of sequences. However, ultimately, the aim must be to generate complete genome sequences for all animals from the ENU archives, providing a unique and unparalleled database and genetics resource. Effectively, the user would be able to recover mutant alleles for any sequence within the mouse genome, including not only coding sequences but also point mutations within enhancers and other regulatory motifs as well as transcribed noncoding sequences. This resource would transform the functional annotation of the mouse and human genomes and provide a new modality for the functional assessment of any mammalian sequence.

Transposons and Phenotype-Driven Screens

Above, we discuss the utility of transposon-based systems (48, 49) for mouse mutagenesis. We can expect this approach to make a distinctive contribution to phenotype-driven screens. The efficient generation of mutant mice by transposon-based mutagenesis followed by their analysis using comprehensive phenotype screens has the potential to uncover a large range of novel mutant phenotypes. Most importantly, the identification of the underlying locus will be relatively trivial.

THE CHALLENGES OF PHENOTYPING: DELIVERING FUNCTIONAL ANNOTATION FROM THE MOUSE MUTANT RESOURCE

What is a Phenotype?

Phenotype comprises a complex biological output of gene allele, the genetic background,

environment, and, importantly, the test utilized to measure the phenotype. This can be visualized as a three-dimensional matrix (Figure 3) consisting of data points that reflect in one dimension the genetic context and orthogonally the environment in which the phenotype has been measured (22). We need to assess phenotypes in many different mutations carried on diverse genetic backgrounds if we are to explore fully phenotype space. At the same time, phenotype outcome will be affected by the environmental conditions in which animals were raised and in which the phenotype test has been carried out. We can expect many environmental parameters to have a significant impact upon phenotype measurements (151) including, for example, home cage environment, environmental enrichment, diet, and the pathogenic load on the animal. To some extent, these parameters can be controlled, but their obvious impact on phenotypic outcome identifies them as critical data components for phenotype databases (see discussion below). The third dimension of our phenotype data matrix is the phenotype test itself. For each allele, genetic background, and environmental condition, multiple tests or even variants of tests can be carried out. Different variants of a phenotype test may measure different phenotypic parameters, and this richness or granularity in phenotyping must be recognized and taken into account in the three-dimensional data sets that are generated. The major challenge of mouse phenotyping is to populate this three-dimensional matrix, delivering a unified and integrated data set. This is indeed a phenomenal goal, especially if we were to consider an early ambition to phenotype 20,000 mouse lines on two genetic backgrounds, requiring 100 phenotype tests, each measuring five parameters. Further consider that these tests are carried out under a minimum of five environmental conditions. Thus, we need to accrue 5×10^7 phenotype data points. However, even this limited goal represents a very minimal data set for detailed downstream systems analysis.

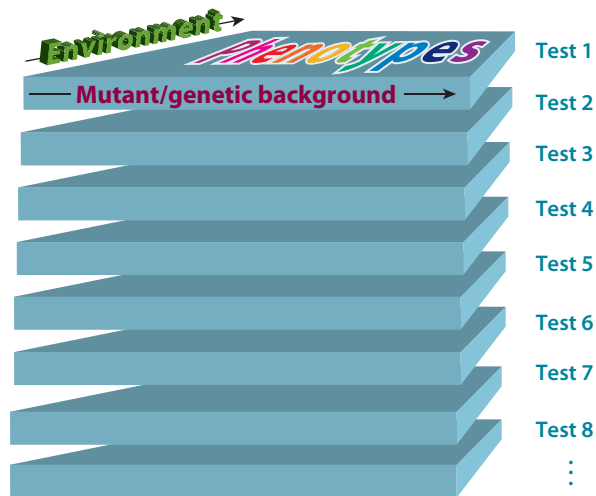


Figure 3

Phenotype data comprises a complex three-dimensional matrix of outputs reflecting multiple phenotype tests, the genetic lesion, genetic background, and the environmental conditions under which the tests were undertaken. Populating this matrix is the key challenge that faces mouse geneticists.

What is a Phenotype Test?

This may seem a trite query, but within this simple question lies a fundamental quandary about the nature of the mouse phenotyping project, its aims, and its early directions. Phenotypes can be measured at the molecular (transcriptomic, proteomic, and metabolomic), the cellular, and physiological levels. In addition, analysis at these various hierarchies can be carried out at different developmental stages and different time windows within adults, thereby encompassing phenotypes from early-onset disease to age-related conditions. If we add these complexities to our phenotyping programs, the size of the data sets that need to be accrued significantly increases.

Enormous progress in mouse phenotyping has been made at all levels. Most recently, considerable emphasis has been placed on the development and application of mouse phenotyping platforms for observational, physiological, and biochemical screens to study diverse features such as dysmorphology, neurological, metabolic, cardiovascular, behavioral, and sensory disease states. The strength of these

phenotype platforms is that they determine phenotypic parameters that may be mapped to human disease features. However, as discussed below, in many cases these phenotype platforms offer only modest throughput. In contrast, molecular phenotypic analysis (for example, transcriptomic arrays) is more immediately scalable to the analysis of thousands of mutant lines. However, algorithms that would allow for the accurate relation of transcriptomic or proteomic data to disease conditions are currently lacking. This reflects our lack of understanding of genetic systems and their networks—the very problem that we are trying to solve. For this reason, early large-scale efforts in mouse phenotyping have focused their efforts on acquiring substantive data sets from observational, physiological, and biochemical platforms. Nevertheless, it will be important to populate and integrate phenotype databases with molecular phenotyping data as it becomes available if we are to make a profound analysis of genetic systems from the molecular to the organismal.

The Importance of Standardization

Given the scale of the data sets that need to be acquired, mouse mutant phenotyping is being carried out across multiple centers. The distributed nature of the mouse phenotyping effort underlines the need for standardization. As discussed above, there is enormous potential for variation in the use of phenotype tests as well as disparities in environmental conditions that may introduce unwanted variability into phenotype outcomes. Indeed, considerable discussion has ensued as to the benefits of standardization in contrast to a more systematic but inevitably laborious approach to analyzing genetic and environmental variables and detecting gene-environment interactions (108, 114, 151, 152). For example, it has recently been proposed that standardization of the environment may reduce the external validity of phenotyping data contributing to poor reproducibility and that systematic environmental heterogenization approaches be used to minimize spurious and conflicting findings. However, to

carry out phenotype analyses in this way on any comprehensive scale is impractical given the scale of the endeavor. A contrasting view is that standardization ensures that strain ranking (for example, between controls and mutants) measured among centers is consistent (even though absolute values may vary), and this has been the focus of several labs' efforts (91).

Overall, standardization of phenotyping procedures will be critical if we are to populate unified and integrated databases with comparable data sets. Data sets must help answer questions such as, Is the phenotype of allele 1 (measured in center 1 with respect to the background control) the same or similar to the phenotype of allele 2 (measured in center 2 with respect to the same background control)? Or, for example, Is the phenotype of allele 3 (measured in center 3) different from the phenotype of allele 4 (measured in center 4)? In the latter case, we need to be sure that differences do not reflect trivial variations in test procedure or environmental conditions that may result in unexpected outcomes in test output. The use of robust, validated phenotyping platforms generating data that is comparable across time and place will be integral to populating unified phenotype databases. The need for standardization in mouse phenotyping, as well as standardization in data acquisition, storage, and annotation, reflects a wider concern in developing a systems biology for any organism (18, 30).

Standardization of Mouse Phenotyping Procedures: The Importance of Standard Operating Procedures and Environmental Standards

We need to adopt standard operating procedures (SOPs) that deliver comparable phenotype outputs across time and place if we are to undertake a comprehensive functional annotation of the mouse genome. Developing robust procedures requires us to assess how the operation of the tests along with associated environmental conditions, such as cage environment or diet, impact the test output. A landmark study

by Crabbe and colleagues (43) examined the reproducibility of a battery of neurobehavioral tests across three laboratories. Each laboratory carried out a series of six tests on a panel of several inbred strains. Despite rigorous efforts to standardize protocols, test apparatuses, and environmental conditions, significant variations in phenotype output were found among the centers. These results may support the idea that standardization of phenotyping platforms, particularly for behavior, is fraught with difficulty. More likely, the poor reproducibility arose because of unrecognized factors in test procedure or environment that were not taken into account, which underlines the need to undertake deep and extensive analyses of the sources of variation contributing to test output. One recent study (31) investigated the effects of a wide variety of handling and housing test procedures [such as cage density, diet, gender, length of fast, site of bleed (retro-orbital versus tail), timing, and anesthesia] on biochemical, hematological, and metabolic/endocrine parameters. Although many parameters yielded no significant effects on phenotype output, other minor changes in procedures had dramatic impacts.

It will be imperative to develop a strong evidence base for the reproducibility of phenotype platforms across diverse systems, requiring the researching and cataloging of the many variables that may confound test standardization. An understanding of the impact of environmental conditions, particularly the impact of environmental enrichment, will be critical to ensure test reproducibility. This need is well exemplified by recent studies that have examined the effect of environmental enrichment on the behavioral analysis of inbred strains. Wolfner and colleagues (148) found that environmental enrichment had no significant effect on the reproducibility of data acquired in replicate studies, and concluded that the introduction of enrichment has little impact upon the standardization of results. They also found that enrichment conditions did not affect the strain ranking observed in the tests utilized. This study examined only two inbred strains and used a limited battery of tests. In contrast, a study

carried out by the EUMORPHIA project (see below) examined a larger number of inbred strains with a more extensive battery of behavioral tests (141). A complex picture emerges showing that for some phenotypes strain ranking is clearly affected by the enrichment conditions, whereas for other tests no effect on strain ranking is seen.

These results underline the complexities and potential pitfalls as we aim to develop and utilize SOPs to populate unified phenome databases. Nevertheless, these have not prevented a major drive to define a comprehensive set of standardized and validated SOPs for wide use across the mouse genetics community. EUMORPHIA, a consortium of European laboratories, recently completed a program to develop standardized phenotyping platforms. The consortium comprised 18 research institutes working on establishing and validating new phenotyping methods. EUMORPHIA has developed a new robust primary screening protocol, EMPReSS (21, 91). EMPReSS incorporates more than 150 SOPs and associated annexes and appendixes, many validated on a cohort of inbred strains across a number of laboratories. Importantly, as discussed above, consideration was given to the impact of animal handling and environmental conditions (31, 141). EMPReSS SOPs are available for all the major body systems and include SOPs for generic approaches such as imaging, pathology, and gene expression. Ultimately, we need to enlarge the resource of SOPs to provide an even more comprehensive catalog of standardized and validated screens, bringing further comparability and reproducibility to phenotyping platforms.

Phenotyping Pipelines: Development of New Platforms and Hierarchies of Tests

There is a continuing need to develop new phenotyping platforms, first to assess phenotype domains for which tests are currently not available and second to satisfy the demand for faster and cheaper phenotyping approaches.

Phenotyping pipeline:

phenotyping tests using predefined SOPs to obtain a standard set of phenotype data about a mouse line or strain

Phenotype test development continues apace with a number of new modalities being added to the classic batteries. A few examples are useful.

In the areas of neurological and behavioral disorders, there is an ongoing requirement to develop new tests that explore novel domains. For example, a new test of motor coordination—the MoRaG (mouse reaching and grasping) test—was recently described (140) and can be used to assess the skilled movements of reaching and grasping previously thought to be confined to primate lineages. The MoRaG test is a useful tool to uncover motor deficits in complex neurological, neuromuscular, neurodegenerative, and behavioral disorders. Within neurology, there has been considerable discussion of the platforms used to assess pain in the mouse (147). Phenotyping development is also addressing the need for models of complex human behavioral disorders. Standardized behavioral assays have now been developed to examine mice's preference for social interactions with novel conspecifics (97). Procedures to quantify sociability and social preference are important for assessing social avoidance tendencies and identifying models of autism.

Combining or juxtaposing existing phenotyping tools is also advantageous to bring a new dimension to uncovering novel phenotypes. For example, combining MRI with high-resolution episcopic microscopy provides a valuable high-resolution, high-throughput tool for accurately determining phenotypes in transgenic and mutant embryos (110). Phenotyping tools such as these have wide application ranging from the discovery of cardiac development disorders (121) to more esoteric phenotypes such as left-right asymmetry.

Phenotyping tests are not always used in isolation. Indeed, as we aim for a broad determination of mouse phenotypes, the use of batteries of tests, or phenotyping pipelines, will become increasingly common. As new SOPs are developed and validated, they need to be incorporated within phenotyping pipelines in a way that ensures the impact of preceding tests on subsequent assays is minimized (94)

and that intertest intervals are appropriate (109). This is particularly true for neurological and behavioral tests. In the EUMORPHIA program, considerable attention was paid to defining an appropriate order of tests within the neurobehavioral battery (21, 91; also see <http://empress.har.mrc.ac.uk/>). Age is also an important consideration in the design of phenotyping pipelines, and it needs to be recognized that certain phenotypes, such as bone mineralization deficits (7), may emerge only at later time points. Factoring in test age and test order will be critical if we are to develop robust, efficient pipelines and move beyond the application of individual, standardized phenotyping tests to pipelines that offer a more comprehensive assessment of gene function.

The development of phenotyping pipelines or batteries of tests brings into focus the considerable variation in the sophistication and speed of different phenotyping tests—a factor that has underlined the adoption of hierarchical approaches to phenotype assessment. Given the large numbers of animals that need to be analyzed in phenotype screens, there has been considerable emphasis on the development of batteries of tests for primary, broad-based, high-throughput screens that provide an initial, but superficial, phenotype assessment. Subsequently, more sophisticated secondary or tertiary tests may be applied depending on the outcome of the primary screen. The advantages of a hierarchical approach to screening has stimulated the development of a number of primary screening pipelines. The SHIRPA (SmithKline Beecham, Harwell, Imperial College, Royal London Hospital Phenotype Assessment) test is a simple and rapid primary phenotyping tool (119) that provides a semiquantitative assessment of muscle and lower motorneuron, spinocerebellar, sensory, neuropsychiatric, and autonomic function. The test protocol is adapted from early work by Irwin (80) on phenotype assessment of pharmacological and toxicological responses. The simple design of SHIRPA, employing relatively unsophisticated equipment and allied to the range of phenotypes assessed, has resulted in its

wide usage. SHIRPA was a critical component of the primary screen used to identify a variety of novel disease models in one of the first large-scale phenotype-driven ENU mutagenesis programs (103). Since its introduction, SHIRPA has evolved and been improved, and a modified version of the protocol is often utilized (92). The hierarchical screening approach has also been adopted for pipelines in specific phenotyping domains such as behavior (44, 45).

However, it is important to emphasize that there is no simple relationship between test sophistication and throughput. Increasingly, we find that technically sophisticated tests, such as MRI, are being used as a primary screen (110, 121). If we are to increase the breadth and depth of the phenotype space that can be captured in a primary screen, then it will be necessary to continue this trend and reverse the traditional relationship between test sophistication and throughput. This will require investment in new technology that enables us to ally speed and sophistication while realizing potential cost benefits.

Bringing Technology Developments to Bear: Speeding Phenotyping Pipelines

Complex behavioral assays, including, for example, circadian rhythm analyses, have been a focus of technology development and are already benefiting from automation in data capture (6, 97). The increase in throughput that this brings allows an increasing number of behavioral phenotypes to be assessed at the primary level. Moreover, assays that are traditionally primary in scope, such as blood analyses (75), are also benefiting from improved technologies and automation. For example, clinical chemistry screens have been revolutionized by the Luminex system, which enables multiplex analysis of a wide range of blood proteins (47).

However, some areas of phenotyping, i.e., techniques that require direct human intervention, will continue to be refractory to dramatic shifts in speed, even when more sophisticated techniques are brought to bear. This is especially true for pathology analysis.

Comprehensive pathological analysis is vital for identifying and validating disease models (124). Yet, the process from necropsy to sectioning, staining, and interpretation is highly labor intensive and depends on direct observational skills. Not much prospect of automating or streamlining pathology analysis appears in the foreseeable future. However, one aspect that can improve the accessibility and standardized recording of pathology results is the development of tools that record results using standard vocabularies, such as the MPATH pathology ontology (124), and export them in formats that can be imported into phenotype databases. One such tool that has recently been published is MoDIS (133), and it is likely that this will become an area of intensive research over the next few years.

Genes and Environment: Employing Challenges in Phenotype Platforms

Despite the emphasis on controlling environment as we standardize test output and populate integrated databases, it will also be increasingly important to determine the effects of specific and disease-related challenges on phenotype outcome. Challenges that are critical to understanding disease development and will be a focus of future efforts include infection, diet, and exercise (for a discussion of the influence of environment on phenotype outcome, see above). Indeed, it will be important to explore in more depth how the richness of the environment impacts disease development, particularly in terms of behavioral and neurodegenerative diseases, given existing evidence that environmental enrichment can influence disease progression (54, 102).

RECENT SIGNIFICANT PROGRESS IN LARGE-SCALE MOUSE PHENOTYPING

Mouse Phenome Project

Inbred strains play a number of critical roles in mouse genetics, including providing uniform genetic backgrounds on which to generate

Ontology: a computational structure for describing data that consists of two elements (standardized terms and standard relationships between those terms)

Mouse clinic: a center with a wide range of facilities for mouse phenotyping

Large-scale/high-throughput phenotyping: the acquisition of phenotype information, using a predefined battery of tests on the phenotypic characteristics of mouse strains or lines

EUMODIC:
European Union
Mouse Disease Clinic

mutant resources. For example, the new mutant libraries from the International Mouse Mutagenesis Consortium are being generated on a C57BL/6N background. Inbred lines also display a wide range of phenotypic variation that can be utilized in genetic crosses, recombinant inbred crosses, or consomic lines to reveal genetic loci underpinning the variant phenotypes (34, 42, 63, 93). It is therefore important to determine systematically the phenotype of key inbred strains, thereby providing baseline data for mutant resources as well as cataloging the inherent variation between inbred strains that may be utilized in genetic studies.

In 2000, the Jackson Laboratory in Bar Harbor, Maine, initiated the Mouse Phenome Project to undertake the systematic phenotypic characterization of a group of 40 mouse inbred strains (14). The strains were chosen on the basis of their popularity and utility in mouse genetics, and they included not only household strains, for example, C57BL/6, but also a number of wild inbred strains, such as *Mus castaneus*, which were employed because of their genetic divergence from the standard laboratory strains. The Mouse Phenome Project was not a directed effort involving one or a few laboratories. Rather, the program was initiated as a distributed effort, whereby investigators were encouraged to contribute phenotyping data on inbred strains to a central database. A large amount of data has been acquired from diverse sources, including internal programs from the Jackson Laboratory and other laboratories worldwide, and deposited in the Mouse Phenome Database (15). This database provides a variety of tools to download and view phenome data, offering a useful basis for strain selection when carrying out diverse mouse genetic experiments.

Development of Mouse Clinics

Some diverse data sets on mouse inbred strains and mutants are emerging from centers with a breadth and depth of expertise in mouse phenotyping, so-called mouse clinics (58). The

emergence of mouse clinics will enhance the scientific community's ability to address the challenge of large-scale phenotyping. It would be unrealistic to expect mouse clinics to undertake or have expertise in all of the more sophisticated physiological, biochemical, and developmental platforms. Rather, it is envisaged that these clinics will sit at the hub of a network of centers, each of which has close links to a clinic and can offer a range of specialist phenotype platforms not available at the actual clinic. There are important scientific, economic, and translational arguments for the development of the clinic concept and its associated networks. First, the scale of the phenotyping effort required to develop a comprehensive functional annotation of the mouse genome is beyond any individual laboratory. Instead, a global network of mouse clinics is required, each able to undertake a broad-based phenotyping pipeline and work together toward common integrated goals. Second, the clinic concept, bringing diverse platforms to efficient phenotyping pipelines, provides economies of scale. Moreover, clinics are efficiently scalable as we aim to tackle the characterization of increased numbers of mutant lines. Finally, as centers for mouse biology and phenotyping, mouse clinics will be able to offer numerous opportunities for the identification, application, and translation of preclinical disease models. For example, the clinics will be in an ideal position to work with diverse partners investigating the efficacy and phenotypic outcomes of small molecules or other therapeutic interventions on appropriate disease models.

EUMODIC

In 2007, the EUMODIC (European Mouse Disease Clinic) program was initiated. The program is a major pilot program for large-scale phenotyping, aiming to undertake the broad-based phenotyping of 500 mouse mutant lines from the EUComm project. EUMODIC comprises four mouse clinics (Helmholtz Zentrum, Munich in Germany, ICS in France, as well as MRC Harwell and the

Sanger Institute in the United Kingdom), each of which is carrying out a primary screen of a proportion of the mutants. The primary screen, EMPReSSslim, comprises a subset of the EMPReSS protocols organized into two pipelines of tests carried out on mice aged 9-15 weeks and encompassing a wide diversity of disease systems (Figure 4). A cohort of 10 age-matched females and 10 age-matched males enters each of the pipelines. A proportion of mutants with interesting phenotypes will undergo more detailed secondary/tertiary phenotyping at other centers within the EUMODIC consortium. All

phenotype data will be made publicly available through the EuroPhenome database (see Reference 90 and below). The EUMODIC program will provide a stern test of the logistics and infrastructure required for large-scale, broad-based phenotype screens and forms a foundation for the continued development of improved approaches that will allow us to scale mouse phenotyping from hundreds to thousands of mutants. Moreover, it will enable us to assess the merits of distributed networks of clinics and associated specialist phenotyping partners.

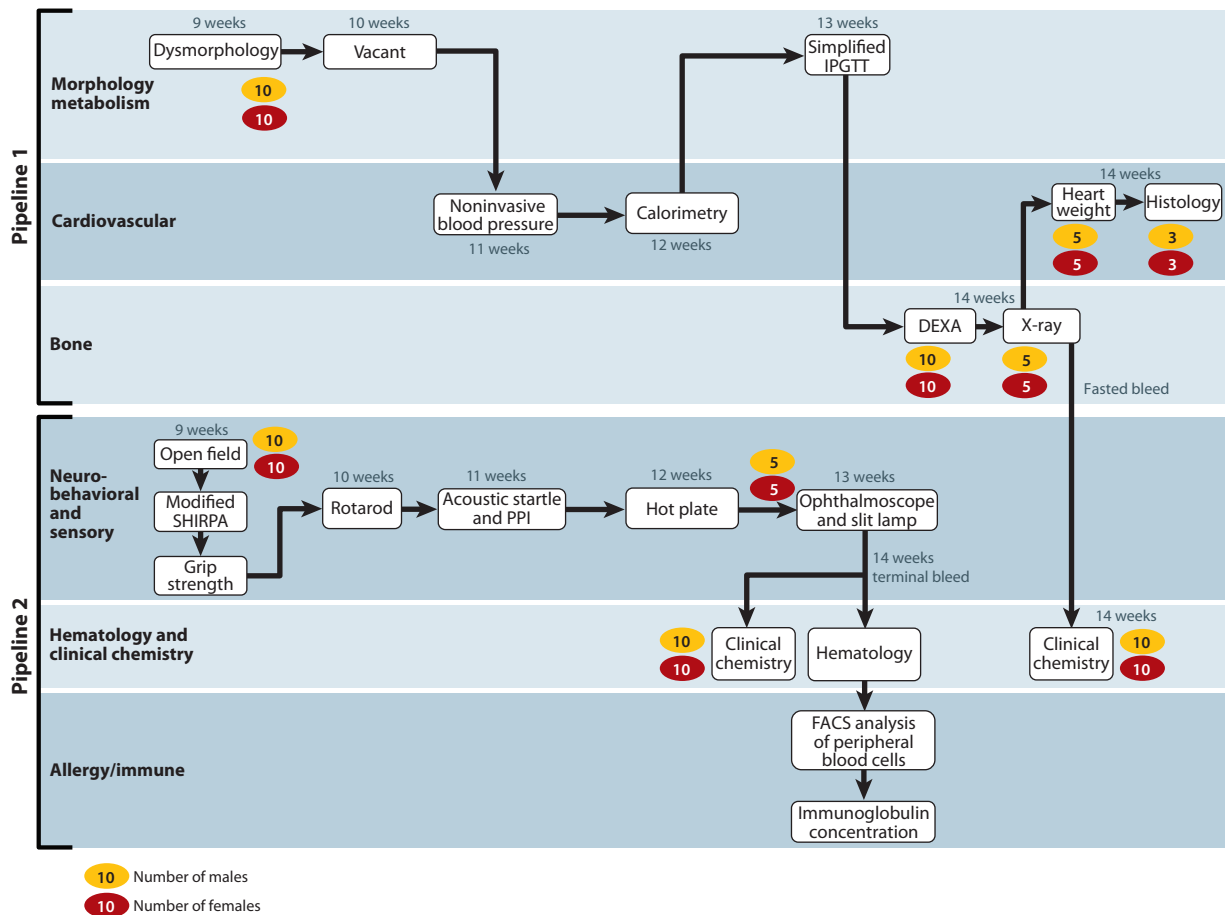


Figure 4

The EMPReSSslim phenotyping pipeline utilized by the EUMODIC consortium. EMPReSSslim comprises two pipelines covering a wide variety of body systems. Age-matched cohorts of 10 males and 10 females enter each of the pipelines at 9 weeks of age, and screening is completed by 15–16 weeks. Smaller cohorts enter some of the test platforms later in the screening process as indicated.

BIOINFORMATICS FOR HIGH-THROUGHPUT PHENOTYPING PROJECTS

Bioinformatics plays an integral part in any large-scale/high-throughput phenotyping program because of the need to track large numbers of mice, including their breeding histories, genotyping records, and progress through phenotyping pipelines, and in order to acquire, store, present, and analyze data in the best possible way. The exciting challenges for the phenotype bioinformatics community lie in two directions. First, there is a need for phenotype data and databases to make the transition from the relatively ad hoc and ill-structured state of today to one in which standards for data, operating procedures, and data transfer are consistently implemented, similar to the way that members of consortia such as caBIG (106) demand implementation of certain standards. Implementation of such standards will allow seamless integration of different phenotype data types, originating from different sources, into a single phenotype data network. This will lead to a second, as yet less well-defined set of challenges: how to analyze these disparate data sets to obtain novel insights into the origin of phenotypes and their relationship to human disease.

Phenotyping Procedures

The primary message from the EUMORPHIA project (19) was that standardization of methods was an essential prerequisite for comparability of phenotype data. This has driven a movement within the bioinformatics community to develop standards in a number of areas related to phenotype data. From an experimental perspective, the most immediately obvious of these has been the development of a standard format for the representation of SOPs for phenotyping. The first version of the EMPReSS database (62) used a simple XML (eXtensible markup language) schema (SOPML) to represent SOPs. XML is a document processing standard that describes the structure of data (17). It allows users to define their own elements and is widely used as an open

standard for exchanging information, structured documents, and data across different information systems, particularly over the Internet. XML defines the data contained within the tagged elements so that a computing application “reading” that file then “knows” what kind of information is held in it. XML data can be described, constrained, and validated within an XML environment by way of an XML schema. The original EMPReSS XML schema was designed to describe and constrain the information within the EMPReSS SOPs. It was adequate for storing the SOPs and allowing them to be regenerated in a number of file formats, but in a broader context, a more sophisticated structure is needed. As discussed above, the EUMODIC project will capture phenotyping data using the same SOPs in four major European centers. To do this, standardization is required beyond the experimental SOPs, and in particular, agreement is required on exactly which parameters for each SOP will be measured by the centers and in which units such parameters will be expressed. The SOPML schema has therefore been extended to include this information. Additionally, as part of international discussions on interoperability and integration of phenotype databases and data (96), a more general list of requirements for SOPs (in this context known as phenotyping procedures) has been drawn up to allow all participating institutions to describe all the features of their phenotyping procedures. The outcomes of a preliminary meeting to define the contents of this schema can be found at <http://tinyurl.com/2db6u9>. The schema is currently being refined as PPML (phenotyping procedure markup language). A first version for general use is planned for release later in 2009 or early in 2010 as part of a broader minimum information standard for the description of phenotyping experiments (135).

Ontological Description of Phenotypic Observations

A significant issue for the bioinformatics community is how to represent phenotype data in

XML (eXtensible markup language): a document-processing standard that describes the structure of data, widely used in computer science

a way that is comprehensible to the user, machine readable, and able to convey information at a range of depths from general summaries of the characteristics of mouse lines to detailed descriptions of phenotype data on individual mice as they come out of high-throughput phenotyping experiments. Since the advent of the gene ontology in the late 1990s (2), ontologies have been the preferred form of data representation in the bioinformatics community (12). Ontologies have the advantages of containing not only a list of defined standard terms that can be used to annotate data but also information about the relationships between terms. An anatomy ontology, for example, may contain terms such as eye, head, and visual system as well as the information that the eye is “part of” the head and “part of” the visual system. The property that a term can be linked to multiple higher level (parental) terms is known as multiple parentage, which is a particularly useful feature in biology where there are multiple points of reference. The first and best known attempt to represent mouse phenotypes is the mammalian phenotype (MP) ontology developed at the Jackson Laboratory (130). This ontology is a powerful resource that currently contains 9354 terms (version 11.03.09), which can be used to annotate abnormal phenotypes of mouse lines and individual mice. However, it is not designed to capture quantitative or qualitative data in detail. An approach to this latter problem that is gaining favor is to use the extended EQ (entity plus quality) formalism and the quality ontology (60, 129).

The EQ approach (Figure 5) represents the phenotypic character that is being observed as a combination of terms from at least two ontologies. The entity (*E*) may be an element of an anatomy ontology (such as eye) or from another relevant orthogonal ontology (which describe domains of knowledge that do not overlap). The aim is that all properties of an organism should be describable by a set of nonoverlapping ontologies so that it should never be necessary to choose between two ontologies that contain the same term, or terms with the same definition. The quality (*Q*) is an element from the

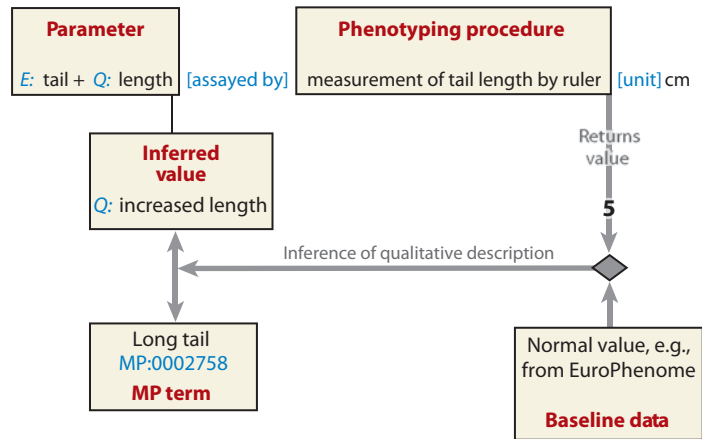


Figure 5

Representation of phenotype data using the EQ (entity plus quality) approach and inference of the mammalian phenotype (MP) terms. A parameter to be measured in a phenotyping experiment is represented by an entity (*E*) and an associated quality (*Q*), in this case the tail and its length, respectively. The measurement is carried out according to a specified assay (in this case, measurement of tail length by ruler), which provides a result in units specified by the assay. This value is compared with baseline data for the same age, sex, and strain to infer whether it is abnormal, which in this case is a tail length of 5 cm. Two qualitative descriptions are then inferred from the quantitative data: (a) a revised quality (i.e., increased length), which can be associated with the entity “tail” in the database, and (b) the MP ontology term “long tail.”

quality ontology (129). This contains properties of things that can be observed, associated by “is a” relationships with the type of property they represent, for example, red is a color. Thus, eye (*E*) + red (*Q*) describes a particular observation about an eye. However, as described here, this formalism is not sufficient to detail the results of phenotyping experiments because observations can depend on the procedure used to obtain them. Thus, a fuller and more useful approach is to combine an *E* + *Q* statement with additional information on the procedure used to obtain the value and the value of the parameter (specified in a structure such as PPML; see above) and a statement of the quantitative or qualitative value that the procedure returned with appropriate units (60). The following is an example:

tail (*E*) + length(*Q*)[tail length measured
using a ruler] = 5 cm

Additional parameters may also need to be associated with such a statement, including for example, the environmental conditions under which mice were kept. Presently, no ontology of phenotyping procedures exists. It remains to be established whether such an ontology is needed and whether its structure would need to mirror that of other ontologies, such as MP or the quality ontology, or if its relationships might differ and reflect the types of techniques used.

The two approaches to representing phenotypes using ontologies are superficially very different, and there has been some discussion as to whether they can be made compatible so that automatic conversions can be made between them. Recently, significant effort has been put into describing MP terms in EQ terms with considerable success (G.V. Gkoutos, personal communication). Although this process is not complete, if a full description of MP terms in terms of EQ could be achieved (and maintained through appropriate curation), the problem of interconversion would be essentially solved. Recent developments in the EuroPhenome database include the development of a method to generate MP terms directly from data held in the EQ format, making use of data on normal mouse lines also held in the database (10).

Beyond the relationship between MP and EQ descriptions of mouse phenotypes, more fundamental issues arise. To make maximal use of mouse phenotype data, it will be important to make the best possible linkage between mouse and human phenotypes and disease. Currently, the standard approach is to link mouse phenotypes to OMIM (Online Mendelian Inheritance in Man) (65) terms and identifications. Apart from the fact that OMIM represents only diseases showing Mendelian inheritance, the primary weakness of this approach is that disease terms as used by OMIM are not phenotypes: Human diseases are complex assemblages of phenotypic observations, and labeling of an individual with a given disease term may be a probabilistic inference based on the presence of some, but not all, of the abnormal phenotypic observations associated with the disease.

It is therefore important to be able to distinguish between the phenotypic characteristics of a particular mouse or mouse line and the human diseases with which these characteristics may be associated (125). An example is type 2 diabetes. Mouse models of type 2 diabetes are often identified using the intraperitoneal glucose tolerance test. However, type 2 diabetes in humans has many features that go well beyond glucose tolerance, and a measure is needed for how well a particular mouse model mimics the human disease in question. Such a measure may be based on the number of shared phenotypic characteristics and how similar they are between the mouse model and human disease state, perhaps also taking into account phenotypic characteristics shown by the mouse model that are not features of the human disease.

Beyond Description: The Phenotype Superhighway

As standards emerge to describe mouse phenotypes in ways that can be implemented by both database managers and Web page developers, a remaining challenge is how to establish a wider international network of phenotype databases that can provide researchers with access to all the richness of phenotype data in its various forms (66). To enable this exchange, a file format, probably XML, is needed to represent the wide variety of current and future phenotype data (96). The development of such a format will require the intensive effort of the international mouse phenotype database community. The EUMODIC community has taken some first steps in this direction. The structure of the EUMODIC project, with a single core database including data from four primary phenotyping centers and numerous secondary phenotypers, required the development of an XML schema for the transport of phenotype data (A. Mallon, H. Morgan, & J.M. Hancock, unpublished manuscript) (**Figure 6**). Further work is needed to make this format suitable for the transfer of the wide range of phenotype data, but such an aim seems achievable within the next few years.

After this, it should be possible for individual databases to return subsets of the data they hold to query or analysis software located anywhere in the world, thus opening up an entire new field of phenotype bioinformatics. To realize this vision, mechanisms must be established to enable the programmatic recall of data. Currently, two main ways of achieving this exist, both of which fall under the general heading of Web services: (a) direct programmatic access to databases by SQL queries (e.g., remote procedure call) and (b) the use of an application programming interface, which allows access to underlying data using a set of predefined functions (e.g., service-oriented architecture). A recent questionnaire (67) investigating the implementation of Web-based access to mouse databases carried out by the Coordination and Sustainability of International Mouse Informatics Resources (<http://www.casimir.org.uk>; <http://www.i-mouse.org/>) consortium indicated that, although the penetration of Web-service implementation is still relatively low (40%), a significant number of databases (33%) are considering adopting Web services in the foreseeable future, suggesting that the infrastructure for such a phenotype data network is coming into existence. A number of higher level entities will probably be necessary before ready access to all mouse phenotype data resources is available—in particular, some form of registry regarding which types of data each database holds and what sorts of Web services are available to recall or process the data. Some Web-service integration projects, such as TAVERNA (78) and BIOMOBY (146), currently host their own registries, but an open resource that could be easily accessed by anyone wishing to design their own Web services-based analysis platform would be a welcome addition.

...AND BEYOND!

Currently, it is not possible to carry out complex queries on phenotype data from different sources and extract novel information from them. We have only glimpses of the kinds of applications that will emerge once the phenotype

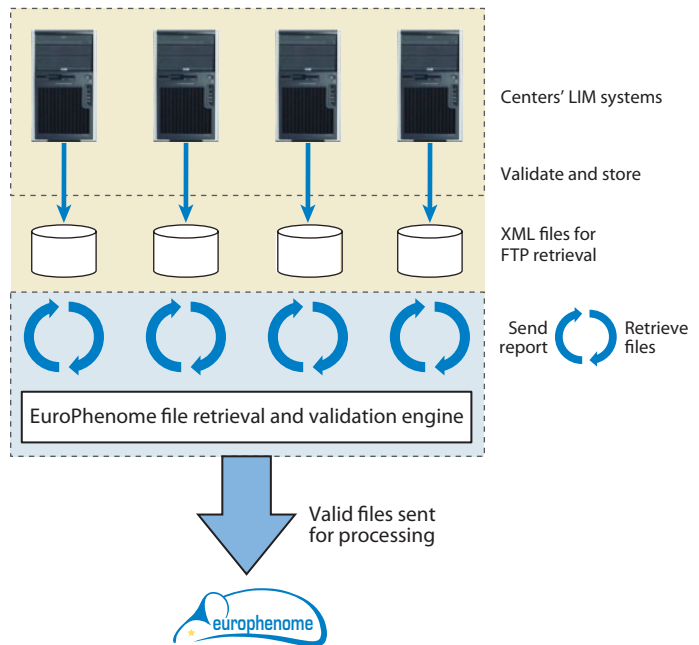


Figure 6

The EUMODIC data-collection process. Data collection takes place in individual phenotyping centers according to common standard operating procedures using the individual laboratories' information systems. These systems export the data in a commonly accepted XML format. These files are validated locally and placed on file transfer protocol servers. These servers are accessed by the central site, uploaded, validated again, and processed. The data are then loaded into the central database EuroPhenome.

data network has started to emerge. Many of the techniques that will be applied to these data sets are likely to be well-established data-mining techniques such as cluster and correlation analysis, although such approaches will need to be enhanced to take into account the large amount of qualitative and categorical data in the phenotype data set, much of which is described using ontologies. An approach that may yield significant dividends is Inductive Logic Programming (98), which is a form of machine learning that aims to learn logical rules describing the classification of a particular subset of the data. This could lead to the learning of new features of the phenotype data set. Unsupervised, as well as supervised, modes of machine learning are also likely to play a significant role. Borrowing from experiences in sequence analysis and genomics, rapid techniques for finding similarity between

phenotype data sets, e.g., a PhenoBLAST algorithm that can take into account various types of data rather than the single data type represented by sequences, may need to be developed. The phenotypes an organism exhibits are ultimately functions of the underlying genetic, signaling, and metabolic networks that are instantiated in the various tissues of the living organism [and during its development (9)]. Thus, it will be important to relate phenotype data to underlying genetic networks both to better understand the origins of

the phenotypes and to understand other gene products and metabolites that may be involved in producing them. Finally, a failing of much early genomic analysis resulted from not taking into account phylogenetic relationships between species. Hence, approaches for the analysis of phenotype data must be developed to take phylogeny into account (88). The next ten years will see exciting developments in the field of phenotype bioinformatics that will make it unrecognizable to those of us working in the area today.

SUMMARY POINTS

1. The laboratory mouse is a key organism in the functional annotation of mammalian genomes because of its ease of use, short generation time, and genetic toolkit.
2. There are now international efforts to generate a wide variety of different kinds of mutations in the mouse genome.
3. A concomitant effort is now under way to generate phenotype information on mutant mouse lines.
4. Standardization of phenotyping approaches will be essential to this project to allow comparability of phenotype data across lines, and the first steps in this direction have been taken.
5. A major requirement for this project will be the establishment of new international mouse clinics and the expansion of existing ones to provide the needed capacity.
6. Advances in phenotype bioinformatics will also be required, ranging from the establishment of databases through the design of advanced description frameworks using ontologies and data-exchange formats leading to a phenotype semantic web. The first steps in this direction have also been taken.

FUTURE ISSUES

1. The first phase of functional annotation of the mouse genome is the completion of the international project to generate mutations of every gene in the mouse genome.
2. Determining the phenotypic outcomes of each of these mutations will be a much greater challenge, and it will be necessary to develop the necessary phenotyping pipelines and clinic facilities to carry this out.
3. The bioinformatic challenges associated with obtaining phenotype data are also complex and manifold and will require the development of new approaches to data representation, exchange, and analysis.

DISCLOSURE STATEMENT

The authors are not aware of any affiliations, memberships, funding, or financial holdings that might be perceived as affecting the objectivity of this review.

ACKNOWLEDGMENTS

We thank Georgios Gkoutos and Ann-Marie Mallon for their comments on the manuscript. This work was supported by the Medical Research Council, United Kingdom.

LITERATURE CITED

1. Albert TJ, Molla MN, Muzny DM, Nazareth L, Wheeler D, et al. 2007. Direct selection of human genomic loci by microarray hybridization. *Nat. Methods* 4:903–5
2. Ashburner M, Ball CA, Blake JA, Botstein D, Butler H, et al. 2000. Gene ontology: tool for the unification of biology. *Nat. Genet.* 25:25–29
3. Augustin M, Sedlmeier R, Peters T, Huffstadt U, Kochmann E, et al. 2005. Efficient and fast targeted production of murine models based on ENU mutagenesis. *Mamm. Genome* 16:405–13
4. Austin CP, Battey JF, Bradley A, Bucan M, Capecchi M, et al. 2004. The knockout mouse project. *Nat. Genet.* 36:921–24
5. Auwerx J, Avner P, Baldock R, Ballabio A, Balling R, et al. 2004. The European dimension for the mouse genome mutagenesis program. *Nat. Genet.* 36:925–27
6. Bacon Y, Ooi A, Kerr S, Shaw-Andrews L, Winchester L, et al. 2004. Screening for novel ENU-induced rhythm, entrainment and activity mutants. *Genes Brain Behav.* 3:196–205
7. Barbaric I, Perry MJ, Dear TN, Rodrigues Da Costa A, Salopek D, et al. 2007. An ENU-induced mutation in the *Ankrd11* gene results in an osteopenia-like phenotype in the mouse mutant Yoda. *Physiol. Genomics* 32:311–21
8. Barbaric I, Wells S, Russ A, Dear TN. 2007. Spectrum of ENU-induced mutations in phenotype-driven and gene-driven screens in the mouse. *Environ. Mol. Mutagen.* 48:124–42
9. Bard J. 2007. Systems developmental biology: the use of ontologies in annotating models and in identifying gene function within and across species. *Mamm. Genome* 18:402–11
10. Beck T, Morgan H, Blake A, Wells S, Hancock JM, Mallon A-M. 2009. Practical application of ontologies to annotate and analyse large scale raw mouse phenotype data. *BMC Bioinform.* 10(Suppl. 5):S2
11. Birney E, Stamatoyannopoulos JA, Dutta A, Guigo R, Gingeras TR, et al. 2007. Identification and analysis of functional elements in 1% of the human genome by the ENCODE pilot project. *Nature* 447:799–816
12. Bodenreider O, Stevens R. 2006. Bio-ontologies: current trends and future directions. *Brief Bioinform.* 7:256–74
13. Bogani D, Warr N, Elms P, Davies J, Tymowska-Lalanne Z, et al. 2004. New semidominant mutations that affect mouse development. *Genesis* 40:109–17
14. Bogue MA, Grubb SC. 2004. The mouse phenome project. *Genetica* 122:71–74
15. Bogue MA, Grubb SC, Maddatu TP, Bult CJ. 2007. Mouse phenome database (MPD). *Nucleic Acids Res.* 35:D643–49
16. Branda CS, Dymecki SM. 2004. Talking about a revolution: the impact of site-specific recombinases on genetic analyses in mice. *Dev. Cell* 6:7–28
17. Bray T, Paoli J, Sperberg-McQueen CM. 1998. *Extensible markup language (XML) 1.0: W3C recommendation 10-February-1998*. <http://www.w3.org/TR/1998/REC-xml-19980210>
18. Brazma A, Krestyaninova M, Sarkans U. 2006. Standards for systems biology. *Nat. Rev. Genet.* 7:593–605
19. Brown S, Lad H, Green E, Gkoutos G, Gates H, et al. 2006. EUMORPHIA and the European Mouse Phenotyping Resource for Standardised Screens (EMPreSS). In *Standards of Mouse Model Phenotyping*, ed. MM Hrabé de Angelis, P Chambon, S Brown, pp. 311–20. Hoboken, NJ: Wiley

20. Brown SD, Balling R. 2001. Systematic approaches to mouse mutagenesis. *Curr. Opin. Genet. Dev.* 11:268–73
21. Brown SD, Chambon P, de Angelis MH. 2005. EMPReSS: standardized phenotype screens for functional annotation of the mouse genome. *Nat. Genet.* 37:1155
22. Brown SD, Hancock JM, Gates H. 2006. Understanding mammalian genetic systems: the challenge of phenotyping in the mouse. *PLoS Genet.* 2:e118
23. Brown SD, Hardisty RE. 2003. Mutagenesis strategies for identifying novel loci associated with disease phenotypes. *Semin. Cell Dev. Biol.* 14:19–24
24. Brown SD, Nolan PM. 1998. Mouse mutagenesis-systematic studies of mammalian gene function. *Hum. Mol. Genet.* 7:1627–33
25. Brown SDM, Hancock JM. 2006. The mouse genome. In *Vertebrate Genomes*, ed. JN Volff, pp. 33–45. Basel: S Karger AG
26. Burton PR, Clayton DG, Cardon LR, Craddock N, Deloukas P, et al. 2007. Association scan of 14500 nonsynonymous SNPs in four diseases identifies autoimmunity variants. *Nat. Genet.* 39:1329–37
27. Cadinanos J, Bradley A. 2007. Generation of an inducible and optimized piggyBac transposon system. *Nucleic Acids Res.* 35:e87
28. Capecchi MR. 1989. The new mouse genetics: altering the genome by gene targeting. *Trends Genet.* 5:70–76
29. Carninci P, Hayashizaki Y. 2007. Noncoding RNA transcription beyond annotated genes. *Curr. Opin. Genet. Dev.* 17:139–44
30. Cassman M. 2005. Barriers to progress in systems biology. *Nature* 438:1079
31. Champy MF, Selloum M, Piard L, Zeitler V, Caradec C, et al. 2004. Mouse functional genomics requires standardization of mouse handling and housing conditions. *Mamm. Genome* 15:768–83
32. Chan W, Costantino N, Li R, Lee SC, Su Q, et al. 2007. A recombineering-based approach for high-throughput conditional knockout targeting vector construction. *Nucleic Acids Res.* 35:e64
33. Chen Y, Yee D, Dains K, Chatterjee A, Cavalcoli J, et al. 2000. Genotype-based screen for ENU-induced mutations in mouse embryonic stem cells. *Nat. Genet.* 24:314–17
34. Churchill GA, Airey DC, Allayee H, Angel JM, Attie AD, et al. 2004. The Collaborative Cross, a community resource for the genetic analysis of complex traits. *Nat. Genet.* 36:1133–37
35. Clark AT, Goldowitz D, Takahashi JS, Vitaterna MH, Siepka SM, et al. 2004. Implementing large-scale ENU mutagenesis screens in North America. *Genetica* 122:51–64
36. Clark KJ, Geurts AM, Bell JB, Hackett PB. 2004. Transposon vectors for gene-trap insertional mutagenesis in vertebrates. *Genesis* 39:225–33
37. Coghill EL, Huggill A, Parkinson N, Davison C, Glenister P, et al. 2002. A gene-driven approach to the identification of ENU mutants in the mouse. *Nat. Genet.* 30:255–56
38. Collins FS, Rossant J, Wurst W. 2007. A mouse for all reasons. *Cell* 128:9–13
39. Concepcion D, Seburn KL, Wen G, Frankel WN, Hamilton BA. 2004. Mutation rate and predicted phenotypic target sizes in ethylnitrosourea-treated mice. *Genetics* 168:953–59
40. Cook MC, Vinuesa CG, Goodnow CC. 2006. ENU-mutagenesis: insight into immune function and pathology. *Curr. Opin. Immunol.* 18:627–33
41. Copeland NG, Jenkins NA, Court DL. 2001. Recombineering: a powerful new tool for mouse functional genomics. *Nat. Rev. Genet.* 2:769–79
42. Cox RD, Brown SD. 2003. Rodent models of genetic disease. *Curr. Opin. Genet. Dev.* 13:278–83
43. Crabbe JC, Wahlsten D, Dudek BC. 1999. Genetics of mouse behavior: interactions with laboratory environment. *Science* 284:1670–72
44. Crawley JN. 2003. Behavioral phenotyping of rodents. *Comput. Med.* 53:140–46
45. Crawley JN, Paylor R. 1997. A proposed test battery and constellations of specific behavioral paradigms to investigate the behavioral phenotypes of transgenic and knockout mice. *Horm. Behav.* 31:197–211
46. Crozat K, Georgel P, Rutschmann S, Mann N, Du X, et al. 2006. Analysis of the MCMV resistome by ENU mutagenesis. *Mamm. Genome* 17:398–406
47. de Jager W, te Velthuis H, Prakken BJ, Kuis W, Rijkers GT. 2003. Simultaneous detection of 15 human cytokines in a single sample of stimulated peripheral blood mononuclear cells. *Clin. Diagn. Lab. Immunol.* 10:133–39

48. Ding S, Wu X, Li G, Han M, Zhuang Y, Xu T. 2005. Efficient transposition of the *piggyBac* (PB) transposon in mammalian cells and mice. *Cell* 122:473–83
49. Drabek D, Zagoraiou L, deWit T, Langeveld A, Roumpaki C, et al. 2003. Transposition of the *Drosophila* hydei Minos transposon in the mouse germline. *Genomics* 81:108–11
50. Duerr RH, Taylor KD, Brant SR, Rioux JD, Silverberg MS, et al. 2006. A genome-wide association study identifies IL23R as an inflammatory bowel disease gene. *Science* 314:1461–63
51. Dupuy AJ, Akagi K, Largaespada DA, Copeland NG, Jenkins NA. 2005. Mammalian mutagenesis using a highly mobile somatic *Sleeping Beauty* transposon system. *Nature* 436:221–26
52. Dupuy AJ, Jenkins NA, Copeland NG. 2006. Sleeping beauty: a novel cancer gene discovery tool. *Hum. Mol. Genet.* 15(Spec. No. 1):R75–79
53. Feil R, Wagner J, Metzger D, Chambon P. 1997. Regulation of Cre recombinase activity by mutated estrogen receptor ligand-binding domains. *Biochem. Biophys. Res. Commun.* 237:752–57
54. Fischer A, Sananbenesi F, Wang X, Dobbin M, Tsai LH. 2007. Recovery of learning and memory is associated with chromatin remodelling. *Nature* 447:178–82
55. Floss T, Wurst W. 2002. Functional genomics by gene-trapping in embryonic stem cells. *Methods Mol. Biol.* 185:347–79
56. Frayling TM, Timpson NJ, Weedon MN, Zeggini E, Freathy RM, et al. 2007. A common variant in the FTO gene is associated with body mass index and predisposes to childhood and adult obesity. *Science* 316:889–94
57. Friedrich G, Soriano P. 1991. Promoter traps in embryonic stem cells: a genetic screen to identify and mutate developmental genes in mice. *Genes Dev.* 5:1513–23
58. Gailus-Durner V, Fuchs H, Becker L, Bolle I, Brielmeier M, et al. 2005. Introducing the German Mouse Clinic: open-access platform for standardized phenotyping. *Nat. Methods* 2:403–4
59. Geurts AM, Wilber A, Carlson CM, Lobitz PD, Clark KJ, et al. 2006. Conditional gene expression in the mouse using a *Sleeping Beauty* gene-trap transposon. *BMC Biotechnol.* 6:30
60. Gkoutos GV, Green ECJ, Mallon A-M, Hancock JM, Davidson D. 2005. Using ontologies to describe mouse phenotypes. *Genome Biol.* 6:R8
61. Godinho SI, Maywood ES, Shaw L, Tucci V, Barnard AR, et al. 2007. The after-hours mutant reveals a role for Fbxl3 in determining mammalian circadian period. *Science* 316:897–900
62. Green ECJ, Gkoutos GV, Lad HV, Blake A, Weekes J, Hancock JM. 2005. EMPReSS: European mouse phenotyping resource for standardized screens. *Bioinformatics* 21:2930–31
63. Gregorova S, Divina P, Storchova R, Trachtulec Z, Fotopulosova V, et al. 2008. Mouse consomic strains: exploiting genetic divergence between *Mus m. musculus* and *Mus m. domesticus* subspecies. *Genome Res.* 18:(3)509–15
64. Gu H, Marth JD, Orban PC, Mossman H, Rajewsky K. 1994. Deletion of a DNA polymerase beta gene segment in T cells using cell type-specific gene targeting. *Science* 265:103–6
65. Hamosh A, Scott AF, Amberger J, Bocchini C, Valle D, McKusick VA. 2002. Online Mendelian Inheritance in Man (OMIM), a knowledgebase of human genes and genetic disorders. *Nucleic Acids Res.* 30:52–55
66. Hancock JM, Mallon A-M. 2007. Phenobabelomics—mouse phenotype data resources. *Brief Funct. Genomics Proteomics* 6:292–301
67. Hancock JM, Schofield PN, Chandras C, Zouberakis M, Aidinis V, et al. 2008. CASIMIR: Coordination and Sustainability of International Mouse Informatics Resources. *Proc. IEEE Int. Conf. Bioinform. Bioeng., 8th, Athens*. Piscataway, NJ: IEEE. doi:10.1109/BIBE.2008.4696712
68. Hansen J, Floss T, Van Sloun P, Fuchtbauer EM, Vauti F, et al. 2003. A large-scale, gene-driven mutagenesis approach for the functional analysis of the mouse genome. *Proc. Natl. Acad. Sci. USA* 100:9918–22
69. Hardisty-Hughes RE, Tateossian H, Morse SA, Romero MR, Middleton A, et al. 2006. A mutation in the F-box gene, Fbxo11, causes otitis media in the Jeff mouse. *Hum. Mol. Genet.* 15:3273–79
70. Hayashizaki Y, Carninci P. 2006. Genome Network and FANTOM3: assessing the complexity of the transcriptome. *PLoS Genet.* 2:e63
71. Hentges KE, Justice MJ. 2004. Checks and balancers: balancer chromosomes to facilitate genome annotation. *Trends Genet.* 20:252–59

72. Herron BJ, Lu W, Rao C, Liu S, Peters H, et al. 2002. Efficient generation and mapping of recessive developmental mutations using ENU mutagenesis. *Nat. Genet.* 30:185–89
73. Hitotsumachi S, Carpenter DA, Russell WL. 1985. Dose-repetition increases the mutagenic effectiveness of N-ethyl-N-nitrosourea in mouse spermatogonia. *Proc. Natl. Acad. Sci. USA* 82:6619–21
74. Hodges E, Xuan Z, Balija V, Kramer M, Molla MN, et al. 2007. Genome-wide in situ exon capture for selective resequencing. *Nat. Genet.* 39:1522–27
75. Hough TA, Nolan PM, Tshipouri V, Toye AA, Gray IC, et al. 2002. Novel phenotypes identified by plasma biochemical screening in the mouse. *Mamm. Genome* 13:595–602
76. Hough TA, Polewski M, Johnson K, Cheeseman M, Nolan PM, et al. 2007. Novel mouse model of autosomal semidominant adult hypophosphatasia has a splice site mutation in the tissue nonspecific alkaline phosphatase gene *Akp2*. *J. Bone Miner. Res.* 22:1397–407
77. Hrabe de Angelis MH, Flawinkel H, Fuchs H, Rathkolb B, Soewarto D, et al. 2000. Genome-wide, large-scale production of mutant mice by ENU mutagenesis. *Nat. Genet.* 25:444–47
78. Hull D, Wolstencroft K, Stevens R, Goble C, Pocock MR, et al. 2006. Taverna: a tool for building and running workflows of services. *Nucleic Acids Res.* 34:W729–32
79. Inoue M, Sakuraba Y, Motegi H, Kubota N, Toki H, et al. 2004. A series of maturity onset diabetes of the young, type 2 (MODY2) mouse models generated by a large-scale ENU mutagenesis program. *Hum. Mol. Genet.* 13:1147–57
80. Irwin S. 1968. Comprehensive observational assessment: Ia. A systematic, quantitative procedure for assessing the behavioral and physiologic state of the mouse. *Psychopharmacologia* 13:222–57
81. Ivics Z, Hackett PB, Plasterk RH, Izsvak Z. 1997. Molecular reconstruction of Sleeping Beauty, a *Tc1*-like transposon from fish, and its transposition in human cells. *Cell* 91:501–10
82. Justice MJ, Carpenter DA, Favor J, Neuhauser-Klaus A, Hrabe de Angelis M, et al. 2000. Effects of ENU dosage on mouse strains. *Mamm. Genome* 11:484–88
83. Justice MJ, Noveroske JK, Weber JS, Zheng B, Bradley A. 1999. Mouse ENU mutagenesis. *Hum. Mol. Genet.* 8:1955–63
84. Keys DA, Clark TG, Flint J. 2006. Estimating the number of coding mutations in genotypic- and phenotypic-driven N-ethyl-N-nitrosourea (ENU) screens. *Mamm. Genome* 17:230–38
85. Kile BT, Hentges KE, Clark AT, Nakamura H, Salinger AP, et al. 2003. Functional genetic analysis of mouse chromosome 11. *Nature* 425:81–86
86. Kuhn R, Schwenk F, Aguet M, Rajewsky K. 1995. Inducible gene targeting in mice. *Science* 269:1427–29
87. Kwan KM. 2002. Conditional alleles in mice: practical considerations for tissue-specific knockouts. *Genesis* 32:49–62
88. Mabee PM, Ashburner M, Cronk Q, Gkoutos GV, Haendel M, et al. 2007. Phenotype ontologies: the bridge between genomics and evolution. *Trends Ecol. Evol.* 22:345–50
89. Maeda N, Kasukawa T, Oyama R, Gough J, Frith M, et al. 2006. Transcript annotation in FANTOM3: mouse gene catalog based on physical cDNAs. *PLoS Genet* 2:e62
90. Mallon AM, Blake A, Hancock JM. 2008. EuroPhenome and EMPReSS: online mouse phenotyping resource. *Nucleic Acids Res.* 36:D715–18
91. Mandillo S, Tucci V, Holter SM, Meziane H, Banchaabouchi MA, et al. 2008. Reliability, robustness, and reproducibility in mouse behavioral phenotyping: a cross-laboratory study. *Physiol. Genomics* 34:243–55
92. Masuya H, Inoue M, Wada Y, Shimizu A, Nagano J, et al. 2005. Implementation of the modified-SHIRPA protocol for screening of dominant phenotypes in a large-scale ENU mutagenesis program. *Mamm. Genome* 16:829–37
93. Martin A, Collin GB, Asada Y, Varnum D, Nadeau JH. 1999. Susceptibility to testicular germ-cell tumours in a 129.MOLF-Chr 19 chromosome substitution strain. *Nat. Genet.* 23:237–40
94. McIlwain KL, Merriweather MY, Yuva-Paylor LA, Paylor R. 2001. The use of behavioral test batteries: effects of training history. *Physiol. Behav.* 73:705–17
95. Michaud EJ, Culiati CT, Klebig ML, Barker PE, Cain KT, et al. 2005. Efficient gene-driven germ-line point mutagenesis of C57BL/6J mice. *BMC Genomics* 6:164
96. Mouse Phenotype Database Integration Consortium. 2007. Integration of mouse phenome data resources. *Mamm. Genome* 18:157–63

97. Moy SS, Nadler JJ, Perez A, Barbaro RP, Johns JM, et al. 2004. Sociability and preference for social novelty in five inbred strains: an approach to assess autistic-like behavior in mice. *Genes Brain Behav.* 3:287–302
98. Muggleton SH. 1991. Inductive logic programming. *New Gener. Comput.* 8:295–318
99. Munroe RJ, Bergstrom RA, Zheng QY, Libby B, Smith R, et al. 2000. Mouse mutants from chemically mutagenized embryonic stem cells. *Nat. Genet.* 24:318–21
100. Muylers JP, Zhang Y, Stewart AF. 2001. Techniques: recombinogenic engineering—new options for cloning and manipulating DNA. *Trends Biochem. Sci.* 26:325–31
101. Nagy A, Mar L. 2001. Creation and use of a Cre recombinase transgenic database. *Methods Mol. Biol.* 158:95–106
102. Nithianantharajah J, Barkus C, Murphy M, Hannan AJ. 2008. Gene–environment interactions modulating cognitive function and molecular correlates of synaptic plasticity in Huntington’s disease transgenic mice. *Neurobiol. Dis.* 29:490–504
103. Nolan PM, Peters J, Strivens M, Rogers D, Hagan J, et al. 2000. A systematic, genome-wide, phenotype-driven mutagenesis program for gene function studies in the mouse. *Nat. Genet.* 25:440–43
104. Noveroske JK, Weber JS, Justice MJ. 2000. The mutagenic action of N-ethyl-N-nitrosourea in the mouse. *Mamm. Genome* 11:478–83
105. Okou DT, Steinberg KM, Middle C, Cutler DJ, Albert TJ, Zwick ME. 2007. Microarray-based genomic selection for high-throughput resequencing. *Nat. Methods* 4:907–9
106. Oster S, Langella S, Hastings S, Ervin D, Madduri R, et al. 2007. caGrid 1.0: an enterprise grid infrastructure for biomedical research. *J. Am. Med. Inform. Assoc.* 15:138–49
107. Parkinson N, Hardisty-Hughes RE, Tateossian H, Tsai HT, Brooker D, et al. 2006. Mutation at the *Evi1* locus in Junbo mice causes susceptibility to otitis media. *PLoS Genet.* 2:e149
108. Paylor R. 2009. Questioning standardization in science. *Nat. Methods* 6:253–54
109. Paylor R, Spencer CM, Yuva-Paylor LA, Pieke-Dahl S. 2006. The use of behavioral test batteries, II: effect of test interval. *Physiol. Behav.* 87:95–102
110. Pielec G, Geyer SH, Szumska D, Schneider J, Neubauer S, et al. 2007. microMRI-HREM pipeline for high-throughput, high-resolution phenotyping of murine embryos. *J. Anat.* 211:132–37
111. Porreca GJ, Zhang K, Li JB, Xie B, Austin D, et al. 2007. Multiplex amplification of large sets of human exons. *Nat. Methods* 4:931–36
112. Quwailid MM, Hugill A, Dear N, Vizor L, Wells S, et al. 2004. A gene-driven ENU-based approach to generating an allelic series in any gene. *Mamm. Genome* 15:585–91
113. Rajewsky K, Gu H, Kuhn R, Betz UA, Muller W, et al. 1996. Conditional gene targeting. *J. Clin. Invest.* 98:600–3
114. Richter SH, Garner JP, Wurbel H. 2009. Environmental standardization: cure or cause of poor reproducibility in animal experiments? *Nat. Methods* 6:257–61
115. Rinchik EM. 1991. Chemical mutagenesis and fine-structure functional analysis of the mouse genome. *Trends Genet.* 7:15–21
116. Rinchik EM, Carpenter DA, Selby PB. 1990. A strategy for fine-structure functional analysis of a 6- to 11-centimorgan region of mouse chromosome 7 by high-efficiency mutagenesis. *Proc. Natl. Acad. Sci. USA* 87:896–900
117. Rioux JD, Xavier RJ, Taylor KD, Silverberg MS, Goyette P, et al. 2007. Genome-wide association study identifies new susceptibility loci for Crohn disease and implicates autophagy in disease pathogenesis. *Nat. Genet.* 39:596–604
118. Rodriguez CI, Buchholz F, Galloway J, Sequerra R, Kasper J, et al. 2000. High-efficiency deleter mice show that FLPe is an alternative to Cre-loxP. *Nat. Genet.* 25:139–40
119. Rogers DC, Fisher EM, Brown SD, Peters J, Hunter AJ, Martin JE. 1997. Behavioral and functional analysis of mouse phenotype: SHIRPA, a proposed protocol for comprehensive phenotype assessment. *Mamm. Genome* 8:711–13
120. Sakuraba Y, Sezutsu H, Takahasi KR, Tsuchihashi K, Ichikawa R, et al. 2005. Molecular characterization of ENU mouse mutagenesis and archives. *Biochem. Biophys. Res. Commun.* 336:609–16

121. Schneider JE, Bose J, Bamforth SD, Gruber AD, Broadbent C, et al. 2004. Identification of cardiac malformations in mice lacking Ptdsr using a novel high-throughput magnetic resonance imaging technique. *BMC Dev. Biol.* 4:16
122. Schnütgen F, De-Zolt S, Van Sloun P, Hollatz M, Floss T, et al. 2005. Genome-wide production of multipurpose alleles for the functional analysis of the mouse genome. *Proc. Natl. Acad. Sci. USA* 102:7221–26
123. Schnütgen F, Doerflinger N, Calleja C, Wendling O, Chambon P, Ghyselink NB. 2003. A directional strategy for monitoring Cre-mediated recombination at the cellular level in the mouse. *Nat. Biotechnol.* 21:562–65
124. Schofield PN, Bard JB, Booth C, Boniver J, Covelli V, et al. 2004. Pathbase: a database of mutant mouse pathology. *Nucleic Acids Res.* 32:D512–15
125. Schofield PN, Rozell B, Gkoutos GV. 2008. Towards a disease ontology. In *Anatomy Ontologies for Bioinformatics: Principles and Practice*, ed. A Burger, D Davidson, R Baldock, pp. 119–32. London: Springer-Verlag
126. Schwander M, Sczaniecka A, Grillet N, Bailey JS, Avenarius M, et al. 2007. A forward-genetics screen in mice identifies recessive deafness traits and reveals that pejkakin is essential for outer hair cell function. *J. Neurosci.* 27:2163–75
127. Skarnes WC, von Melchner H, Wurst W, Hicks G, Nord AS, et al. 2004. A public gene trap resource for mouse functional genomics. *Nat. Genetics*: 543–44
128. Skopek TR, Walker VE, Cochrane JE, Craft TR, Cariello NF. 1992. Mutational spectrum at the Hprt locus in splenic T cells of B6C3F1 mice exposed to N-ethyl-N-nitrosourea. *Proc. Natl. Acad. Sci. USA* 89:7866–70
129. Smith B, Ashburner M, Rosse C, Bard J, Bug W, et al. 2007. The OBO Foundry: coordinated evolution of ontologies to support biomedical data integration. *Nat. Biotechnol.* 25:1251–55
130. Smith CL, Goldsmith CA, Eppig JT. 2005. The Mammalian Phenotype Ontology as a tool for annotating, analyzing and comparing phenotypic information. *Genome Biol.* 6:R7
131. Stanford WL, Cohn JB, Cordes SP. 2001. Gene-trap mutagenesis: past, present and beyond. *Nat. Rev. Genet.* 2:756–68
132. Sun LV, Jin K, Liu Y, Yang W, Xie X, et al. 2008. PBmice: an integrated database system of piggyBac (PB) insertional mutations and their characterizations in mice. *Nucleic Acids Res.* 36:D729–34
133. Sundberg JP, Sundberg BA, Schofield P. 2008. Integrating mouse anatomy and pathology ontologies into a phenotyping database: tools for data capture and training. *Mamm. Genome* 19:413–19
134. Takahasi KR, Sakuraba Y, Gondo Y. 2007. Mutational pattern and frequency of induced nucleotide changes in mouse ENU mutagenesis. *BMC Mol. Biol.* 8:52
135. Taylor CF, Field D, Sansone SA, Aerts J, Apweiler R, et al. 2008. Promoting coherent minimum reporting guidelines for biological and biomedical investigations: the MIBBI project. *Nat. Biotechnol.* 26:889–96
136. Thauung C, West K, Clark BJ, McKie L, Morgan JE, et al. 2002. Novel ENU-induced eye mutations in the mouse: models for human eye disease. *Hum. Mol. Genet.* 11:755–67
137. Thomas KR, Capocchi MR. 1987. Site-directed mutagenesis by gene targeting in mouse embryo-derived stem cells. *Cell* 51:503–12
138. Todd JA, Walker NM, Cooper JD, Smyth DJ, Downes K, et al. 2007. Robust associations of four new chromosome regions from genome-wide analyses of type 1 diabetes. *Nat. Genet.* 39:857–64
139. Torres RM, Kühn R. 1997. *Laboratory Protocols for Conditional Gene Targeting*. Oxford: Oxford Univ. Press. 167 pp.
140. Tucci V, Achilli F, Blanco G, Lad HV, Wells S, et al. 2007. Reaching and grasping phenotypes in the mouse (*Mus musculus*): a characterization of inbred strains and mutant lines. *Neuroscience* 147:573–82
141. Tucci V, Lad HV, Parker A, Polley S, Brown SD, Nolan PM. 2006. Gene-environment interactions differentially affect mouse strain behavioral parameters. *Mamm. Genome* 17:1113–20
142. Vinuesa CG, Cook MC, Angelucci C, Athanasopoulos V, Rui L, et al. 2005. A RING-type ubiquitin ligase family member required to repress follicular helper T cells and autoimmunity. *Nature* 435:452–58
143. Vivian JL, Chen Y, Yee D, Schneider E, Magnuson T. 2002. An allelic series of mutations in Smad2 and Smad4 identified in a genotype-based screen of N-ethyl-N-nitrosourea-mutagenized mouse embryonic stem cells. *Proc. Natl. Acad. Sci. USA* 99:15542–47

144. Wang W, Lin C, Lu D, Ning Z, Cox T, et al. 2008. Chromosomal transposition of *PiggyBac* in mouse embryonic stem cells. *Proc. Natl. Acad. Sci. USA* 105:9290–95
145. Weber JS, Salinger A, Justice MJ. 2000. Optimal N-ethyl-N-nitrosourea (ENU) doses for inbred mouse strains. *Genesis* 26:230–33
146. Wilkinson MD, Links M. 2002. BioMOBY: an open source biological web services proposal. *Brief Bioinform.* 3:331–41
147. Wilson SG, Mogil JS. 2001. Measuring pain in the (knockout) mouse: big challenges in a small mammal. *Behav. Brain Res.* 125:65–73
148. Wolfer DP, Litvin O, Morf S, Nitsch RM, Lipp HP, Wurbel H. 2004. Laboratory animal welfare: cage enrichment and mouse behavior. *Nature* 432:821–22
149. Wu S, Ying G, Wu Q, Capecchi MR. 2007. Toward simpler and faster genome-wide mutagenesis in mice. *Nat. Genet.* 39:922–30
150. Wu SC, Meir YJ, Coates CJ, Handler AM, Pelczar P, et al. 2006. *piggyBac* is a flexible and highly active transposon as compared to *Sleeping Beauty*, *Tol2*, and *Mos1* in mammalian cells. *Proc. Natl. Acad. Sci. USA* 103:15008–13
151. Wurbel H. 2001. Ideal homes? Housing effects on rodent brain and behavior. *Trends Neurosci.* 24:207–11
152. Wurbel H. 2002. Behavioral phenotyping enhanced—beyond (environmental) standardization. *Genes Brain Behav.* 1:3–8
153. Wurst W, Rossant J, Prideaux V, Kownacka M, Joyner A, et al. 1995. A large-scale gene-trap screen for insertional mutations in developmentally regulated genes in mice. *Genetics* 139:889–99



Contents

Genetic and Epigenetic Mechanisms Underlying Cell-Surface Variability in Protozoa and Fungi <i>Kevin J. Verstrepen and Gerald R. Fink</i>	1
Regressive Evolution in <i>Astyanax</i> Cavefish <i>William R. Jeffery</i>	25
Mimivirus and its Virophage <i>Jean-Michel Claverie and Chantal Abergel</i>	49
Regulation Mechanisms and Signaling Pathways of Autophagy <i>Congcong He and Daniel J. Klionsky</i>	67
The Role of Mitochondria in Apoptosis <i>Chunxin Wang and Richard J. Youle</i>	95
Biom mineralization in Humans: Making the Hard Choices in Life <i>Kenneth M. Weiss, Kazubiko Kawasaki, and Anne V. Buchanan</i>	119
Active DNA Demethylation Mediated by DNA Glycosylases <i>Jian-Kang Zhu</i>	143
Gene Amplification and Adaptive Evolution in Bacteria <i>Dan I. Andersson and Diarmaid Hughes</i>	167
Bacterial Quorum-Sensing Network Architectures <i>Wai-Leung Ng and Bonnie L. Bassler</i>	197
How the Fanconi Anemia Pathway Guards the Genome <i>George-Lucian Moldovan and Alan D. D'Andrea</i>	223
Nucleomorph Genomes <i>Christa Moore and John M. Archibald</i>	251
Mechanism of Auxin-Regulated Gene Expression in Plants <i>Elisabeth J. Chapman and Mark Estelle</i>	265
Maize Centromeres: Structure, Function, Epigenetics <i>James A. Birchler and Fangpu Han</i>	287

The Functional Annotation of Mammalian Genomes: The Challenge of Phenotyping <i>Steve D.M. Brown, Wolfgang Wurst, Ralf Kühn, and John Hancock</i>	305
Thioredoxins and Glutaredoxins: Unifying Elements in Redox Biology <i>Yves Meyer, Bob B. Buchanan, Florence Vignols, and Jean-Philippe Reichheld</i>	335
Roles for BMP4 and CAM1 in Shaping the Jaw: Evo-Devo and Beyond <i>Kevin J. Parsons and R. Craig Albertson</i>	369
Regulation of Tissue Growth through Nutrient Sensing <i>Ville Hietakangas and Stephen M. Cohen</i>	389
Hearing Loss: Mechanisms Revealed by Genetics and Cell Biology <i>Aniel A. Dror and Karen B. Avraham</i>	411
The Kinetochore and the Centromere: A Working Long Distance Relationship <i>Marcin R. Przewloka and David M. Glover</i>	439
Multiple Roles for Heterochromatin Protein 1 Genes in <i>Drosophila</i> <i>Danielle Vermaak and Harmit S. Malik</i>	467
Genetic Control of Programmed Cell Death During Animal Development <i>Barbara Conradt</i>	493
Cohesin: Its Roles and Mechanisms <i>Kim Nasmyth and Christian H. Haering</i>	525
Histones: Annotating Chromatin <i>Eric I. Campos and Danny Reinberg</i>	559
Systematic Mapping of Genetic Interaction Networks <i>Scott J. Dixon, Michael Costanzo, Charles Boone, Brenda Andrews, and Anastasia Baryshnikova</i>	601

Errata

An online log of corrections to *Annual Review of Genetics* articles may be found at <http://genet.annualreviews.org/errata.shtml>