

# Comparative sequencing provides insights about the structure and conservation of marsupial and monotreme genomes

Elliott H. Margulies\*, NISC Comparative Sequencing Program<sup>††</sup>, Valerie V. B. Maduro\*, Pamela J. Thomas<sup>‡</sup>, Jeffery P. Tomkins<sup>§</sup>, Chris T. Amemiya<sup>¶</sup>, Meizhong Luo<sup>||</sup>, and Eric D. Green<sup>\*†††</sup>

\*Genome Technology Branch and <sup>†</sup>NISC, National Human Genome Research Institute, National Institutes of Health, Bethesda, MD 20892; <sup>‡</sup>Clemson University Genomics Institute, Department of Genetics and Biochemistry and Life Science Studies, Clemson University, Clemson, SC 29634; <sup>§</sup>Benaroya Research Institute at Virginia Mason, Seattle, WA 98101; and <sup>¶</sup>Arizona Genomics Institute, Department of Plant Sciences, University of Arizona, Tucson, AZ 85721

Communicated by Francis S. Collins, National Institutes of Health, Bethesda, MD, November 18, 2004 (received for review August 30, 2004)

Sequencing and comparative analyses of genomes from multiple vertebrates are providing insights about the genetic basis for biological diversity. To date, these efforts largely have focused on eutherian mammals, chicken, and fish. In this article, we describe the generation and study of genomic sequences from noneutherian mammals, a group of species occupying unusual phylogenetic positions. A large sequence data set (totaling >5 Mb) was generated for the same orthologous region in three marsupial (North American opossum, South American opossum, and Australian tammar wallaby) and one monotreme (platypus) genomes. These ancient mammalian genomes are characterized by unusual architectural features with respect to G + C and repeat content, as well as compression relative to human. Approximately 14% and 34% of the human sequence forms alignments with the orthologous sequence from platypus and the marsupials, respectively; these numbers are distinctly lower than that observed with nonprimate eutherian mammals (45–70%). The alignable sequences between human and each marsupial species are not completely overlapping (only 80% common to all three species) nor are the platypus-alignable sequences completely contained within the marsupial-alignable sequences. Phylogenetic analysis of synonymous coding positions reveals that platypus has a notably long branch length, with the human–platypus substitution rate being on average 55% greater than that seen with human–marsupial pairs. Finally, analyses of the major mammalian lineages reveal distinct patterns with respect to the common presence of evolutionarily conserved vertebrate sequences. Our results confirm that genomic sequence from noneutherian mammals can contribute uniquely to unraveling the functional and evolutionary histories of the mammalian genome.

comparative genomics | genome sequencing | genome analysis | phylogenetics | mammalian evolution

Comparisons of genome sequences from evolutionarily diverse species are central to decoding the functions of vertebrate genomes (1). Of particular interest is the use of highly diverged species for detecting and characterizing sequences under purifying selection (2). Large-scale sequence comparisons have been reported for eutherian (commonly referred to as “placental”) mammals (3) or fish (4), with the most detailed studies to date emphasizing human–rodent comparisons (5, 6).

We previously described our efforts to sequence the same orthologous regions from large collections of vertebrates (7, 8) and to perform multispecies sequence comparisons (9). These analyses have helped to refine phylogenetic relationships (7), to gain insight about the mutational process (10, 11), and to reveal differences between eutherian mammals and other vertebrates (e.g., birds and fish) with respect to their utility for detecting highly conserved regions in the human genome (9). However, these studies also demonstrate that for comparative sequence

analyses, the optimal phylogenetic distances among species vary, depending on the question(s) being addressed [with the distance between humans and eutherian mammals sometimes being too close, and that between humans and birds (or fish) sometimes being too far].

Within this large phylogenetic gap between eutherian mammals and birds reside the marsupials and monotremes (12, 13). These metatherian and prototherian mammals diverged before the eutherian radiation, estimated at 185 and 200 million years ago (mya), respectively (14). Indeed, these divergence dates, as well as the origins of prototherian mammals relative to metatherian mammals, remain a source of scientific debate, in part because of insufficient molecular data (13, 15–17). Until recently, very little marsupial or monotreme DNA sequence was available in public databases. Although comparative studies involving small amounts of genomic sequence from a marsupial species [the stripe-faced dunnart (*Sminthopsis macroura*)] have been described (18), no comparisons involving large, contiguous blocks of marsupial or monotreme sequence have been reported to date.

In this article, we present the results of comparative sequence analyses involving >5 Mb of sequence from four noneutherian mammals. Specifically, we describe the features of their genomes, provide insights about their phylogenetic relationships, and reveal similarities and differences among mammalian lineages with respect to the presence of evolutionarily conserved vertebrate sequences.

## Materials and Methods

**Genomic Sequence Data Set.** Genomic segments orthologous to a 1.9-Mb region on human chromosome 7q31.3, encompassing the

Freely available online through the PNAS open access option.

Abbreviations: NISC, National Institutes of Health Intramural Sequencing Center; mya, million years ago; N.A., North American; S.A., South American; BAC, bacterial artificial chromosome; TBA, THREADED BLOCKSET ALIGNER; MCS, multispecies conserved sequence; 4D, 4-fold degenerate; SINEs, short interspersed nucleotide elements; LINEs, long interspersed nucleotide elements.

Data deposition: The sequences reported in this paper have been deposited in the GenBank database [accession nos. AC127465, AC129065, AC129066, AC129885, AC142561, AC144364, AC144365, AC144600, AC144690, AC144691, AC144755, and AC144756 (N.A. opossum); AC147869, AC147870, AC147871, AC147872, AC147873, AC147874, AC148151, and AC148214 (S.A. opossum); AC127464, AC129882, AC129883, AC129884, AC130185, AC138553, AC144363, AC144689, AC144753, AC144754, AC144788, AC146535, and AC146754 (platypus); and AC145041, AC145042, AC145183, AC145184, AC145249, AC145250, AC145407, AC145408, AC145409, and AC145841 (wallaby)]. See Table 3, which is published as supporting information on the PNAS web site for specific versions of all GenBank accession nos. used in this study.

<sup>††</sup>National Institutes of Health Intramural Sequencing Center (NISC) Comparative Sequencing Program: Leadership provided by Robert W. Blakesley, Gerard G. Bouffard, Nancy F. Hansen, Baishali Maskeri, and Jennifer C. McDowell.

<sup>†††</sup>To whom correspondence should be addressed at: National Human Genome Research Institute, National Institutes of Health, 50 South Drive, Building 50, Room 5222, Bethesda, MD 20892. E-mail: egreen@nhgri.nih.gov.

© 2005 by The National Academy of Sciences of the USA



**Table 2. Architectural features of different species' sequences**

Species	G + C content*				
	Total	Nonrepetitive sites	Synonymous 4D sites	Relative size <sup>†</sup>	Percentage repetitive <sup>‡</sup>
Human	0.384	0.369	0.432	NA	40.3
Cat	0.383	0.372	0.434	0.95	36.4
Pig	0.377	0.366	0.455	0.92	31.9
Mouse	0.401	0.391	0.479	0.90	32.6
<b>N.A. opossum</b>	<b>0.358</b>	<b>0.358</b>	<b>0.415</b>	<b>1.17</b>	<b>43.2</b>
<b>S.A. opossum</b>	<b>0.358</b>	<b>0.358</b>	<b>0.380</b>	<b>0.99</b>	<b>34.2</b>
<b>Wallaby</b>	<b>0.373</b>	<b>0.374</b>	<b>0.412</b>	<b>1.15</b>	<b>37.0</b>
<b>Platypus</b>	<b>0.459</b>	<b>0.457</b>	<b>0.642</b>	<b>0.76</b>	<b>44.9</b>
Chicken	0.412	0.407	0.423	0.44	6.0
<i>Fugu</i>	0.486	0.485	0.721	0.16	2.3

Boldface indicates the data for noneutherian mammals.

\*Fraction of G + C bases in the entire sequence (total), the nonrepetitive portion of sequence (i.e., sequence not masked by REPEATMASKER), and synonymous 4D sites (the third position of codons that can be any base and still code for the same amino acid).

<sup>†</sup>Ratio of sequence length in each species to the amount of corresponding human sequence (as defined in Table 1).

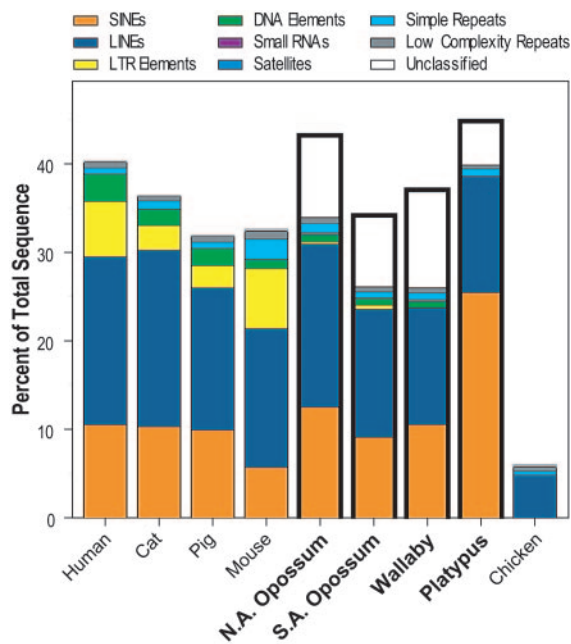
<sup>‡</sup>Percentage of sequence masked by REPEATMASKER.

The asserted correlation between genome size and repeat content (4, 26) prompted us to investigate the amount and composition of repetitive elements within each species' sequence. Because repetitive sequences in noneutherian mammals have not been fully characterized, this analysis first required assembling repeat libraries for each marsupial and monotreme species (see *Materials and Methods*). Fig. 1 shows a summary of the content and types of repeats in each species' sequence, with data from several other vertebrates provided for comparison. Note the considerable variation in total repeat content among these species and the lack of correlation with genome size

(relative to human; see Table 2). Specifically, the orthologous platypus genomic region is smaller than the human region yet contains a larger proportion of repetitive sequences; similarly, the wallaby genomic region is larger than the human region yet contains a smaller proportion of repetitive sequences. Another finding is the relatively large proportion of short interspersed nucleotide elements (SINEs) in the platypus sequence (27, 28), markedly different from other vertebrate sequences. The latter is consistent with the PCR-based identification of an abundant SINE repeat within monotreme genomes (J. A. M. Graves and P. J. Kirby, personal communication).

The overall G + C content is similar among the three marsupial sequences (35.8–37.3%; see Table 2), which is slightly lower than that of the orthologous human genomic region (38.4%). In contrast, the overall G + C content of the platypus sequence is notably high (45.9%), more like that seen with the orthologous *Fugu* genomic region (48.6%). A similarly high G + C content for platypus is seen in the nonrepetitive sites and at synonymous 4D sites (see Table 2). Examining the distribution of G + C content in 1-kb windows across the noneutherian sequences reveals the same general trends (see the supporting information).

**Multispecies Sequence Comparisons.** Analyses of a multisequence alignment generated by using data from 27 vertebrates revealed notable patterns of sequence conservation. For example, the fraction of the human sequence forming alignments with non-primate eutherian mammals is typically 45–70% (Fig. 2A) (7); these alignments include both neutrally evolving and functionally constrained portions of the sequence. This fraction of alignable sequence is significantly lower for the noneutherian mammals (14–34%), with the decrease mostly reflecting fewer alignments within nonannotated regions (i.e., those reflecting sequences not thought to be genes or repeats). A substantially larger amount of noneutherian sequence could be aligned to the human sequence by generating a true multisequence alignment with the program TBA (20) as opposed to simple pair-wise alignments (Fig. 2A, purple bars). In the case of eutherian mammals (where no such difference is seen), it is thought that both pair-wise and multisequence alignments contain virtually all neutrally evolving sequence (5). However, with the noneutherian mammals, the dramatic difference likely reflects a larger amount of neutrally evolving sequence within the multisequence alignment; it re-



**Fig. 1.** Comparison of the content and types of repetitive elements among different species' sequences. Sequences from the orthologous regions of the indicated species' genomes were analyzed by REPEATMASKER, allowing detection and quantification of the indicated types of repetitive elements. The data for the noneutherian mammals are highlighted for emphasis. SINEs, short interspersed nucleotide elements; LINES, long interspersed nucleotide elements.





further evidence for the considerable divergence of monotremes relative to both the marsupial and eutherian mammals (16). The synonymous substitution rates we calculated for the mouse and rat sequences are similar to the genome-wide estimates (5, 6), whereas that for the chicken sequence is substantially lower than the genome-wide estimate (29). The latter is likely attributable to differences in the methods and assumptions used and/or characteristics of the respective data sets (i.e., pair-wise whole-genome analyses vs. multisequence targeted analyses).

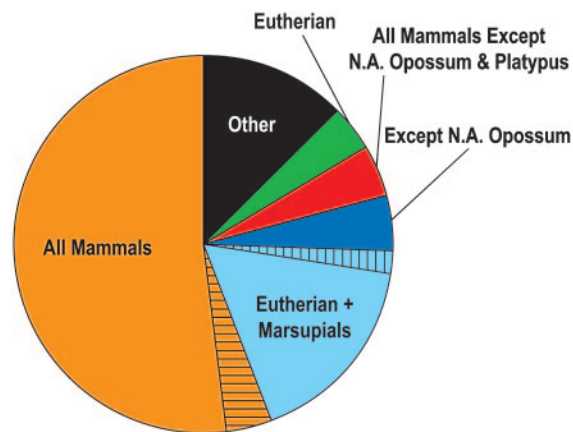
These findings reinforce the distinct phylogenetic positions of marsupials and monotremes within the vertebrate and mammalian radiations (12, 13). In addition, the simultaneous examination of alignment and branch length properties of each species' sequence compared to human (Fig. 2*B*) reveals a clear grouping of the marsupials at an intermediate position between the eutherian mammals and birds, consistent with the purported phylogenetic relationships. In contrast, the grouping of platypus and chicken in this analysis is surprising based on the significant evolutionary distance thought to separate these species (30, 31).

#### Presence of Evolutionarily Conserved Sequences in Different Lineages.

The unique genomic properties of marsupials and monotremes make their sequences of particular interest for identifying and characterizing the small portion of the mammalian genome under purifying selection (5, 32, 33). We previously described an approach for using sequences from multiple vertebrates to detect evolutionarily conserved sequences in the human genome (called MCSs) and demonstrated that different species' sequences vary greatly in their relative contribution to the identification of MCSs (7–9).

Given the diverse representation of mammalian species in our sequence data set, especially with the inclusion of metatherian and prototherian sequences, we next investigated the presence of MCSs among the different mammalian lineages. For this analysis, we studied a set of 418 MCSs falling within a 571-kb portion of the targeted genomic region where there was complete sequence coverage from cat and dog (carnivores), cow and pig (artiodactyls), rat and mouse (rodents), N.A. opossum and wallaby (marsupials), and platypus (monotreme). Note that S.A. opossum sequence was not included in this analysis, so that each lineage would be represented by two species (except monotremes, where only one species was available). The presence or absence of each of the 418 MCSs in each species' sequence was determined based on whether there was a human–species sequence alignment that overlapped that MCS in the human sequence (note that virtually all such alignments reflect high levels of sequence identity). Although virtually all 58 MCSs overlapping coding regions and 46 MCSs overlapping UTRs are present in all species, the remaining noncoding MCSs show interesting patterns of conservation (Fig. 4; also see the supporting information for additional details).

Just over one-half (52%) of the human-referenced noncoding MCSs are present in all nine nonhuman mammals analyzed. These regions thus represent the most anciently constrained sequences in the mammalian lineage. An additional 3.8% of the MCSs are present in all mammals except one or both rodents; this could be due to the known high deletion rate in the rodent lineage (5) or imprecision of current MCS-detection methods. An additional 17% of MCSs are present in all mammals except monotremes, with an additional 2% present in all mammals except monotremes and both rodents. The other major combinations are MCSs in all mammals except N.A. opossum (4.5%), in all mammals except N.A. opossum and platypus (4.5%), and in all eutherian mammals (4.0%). Together, these data provide evidence for lineage specificity with respect to the presence of evolutionarily conserved sequences in the human genome.



**Fig. 4.** Lineage specificity of MCSs. The proportion of nonexonic MCSs found in the sequences of species in each category is indicated. Note that virtually all MCSs overlapping known exonic sequences are present in all mammals (data not shown). All Mammals: cat, dog, cow, pig, rat, mouse, N.A. opossum, wallaby, and platypus; Eutherian: cat, dog, cow, pig, rat, and mouse; Marsupials: N.A. opossum and wallaby; and Other: species combinations containing <2% of the analyzed MCSs (see the supporting information for the complete data set). Hashed areas of “All Mammals” reflect portions lacking one or both rodents, and hashed portions of “Eutherian + Marsupials” reflect portions lacking both rodents.

#### Discussion

Phylogenetic diversity is an important component of comparative genomic studies (8, 34). To date, the comparative sequencing of mammalian genomes largely has involved species within the eutherian radiation, each contributing relatively short branch lengths. Although short branch lengths allow for accurate sequence alignments, many species' sequences then are needed to identify those bases under purifying selection. The more diverged metatherian and prototherian mammals contribute longer branch lengths, making their sequences particularly valuable for identifying genomic regions under purifying selection, while still allowing for reliable alignments to the human sequence. The latter has been challenging with nonmammalian vertebrates, such as chicken and fish (W. Miller, personal communication).

Here, we report the large-scale generation and comparative studies of genome sequences from noneutherian mammals. This initial in-depth glimpse revealed several intriguing properties of these species' genomes. The platypus genome, which, at least for the region studied, shows: (i)  $\approx 25\%$  compression relative to the human genome; (ii) an unusually high G + C content for a mammal; (iii) a disproportionately high fraction of SINES among its repetitive sequences; (iv) a notably low fraction of human-alignable sequence (14% compared with 34% for marsupials); and (v) a markedly long branch length revealed by phylogenetic analyses. Interestingly, these last two properties of platypus are quite similar to those of chicken (see Fig. 2*B*), despite the large difference in their evolutionary distances from human [estimated at 200 versus 310 mya, respectively (12–14)]. Although the long branch length for platypus is intriguing, it was calculated by using the reversible substitution model (REV), which assumes similar nucleotide composition among analyzed sequences. Because this is not the case for platypus (Table 2), and because the synonymous 4D sites analyzed in this study might not be entirely neutrally evolving, caution should be used in making strong claims about the phylogenetic position of monotremes based on our data. Finally, it is interesting to note that the observed compression of the platypus genome (relative to human) cannot be explained fully by differences in gene or repeat content. The evolutionary events that led to this relative compression are not

