

14 Phylogenomic Approach to the Evolutionary Dynamics of Gene Duplication in Birds

CHRIS L. ORGAN

Museum of Comparative Zoology, Harvard University, Cambridge, Massachusetts

MATT RASMUSSEN

Center for Genome Research, Massachusetts Institute of Technology, Cambridge, Massachusetts

MAUDE W. BALDWIN

Museum of Comparative Zoology, Harvard University, Cambridge, Massachusetts

MANOLIS KELLIS

Center for Genome Research, Massachusetts Institute of Technology, Cambridge, Massachusetts

SCOTT V. EDWARDS

Museum of Comparative Zoology, Harvard University, Cambridge, Massachusetts

1 INTRODUCTION

New genes are thought to arise primarily through a process of gene duplication. Genes that are homologous as a result of divergence across lineages via speciation are said to be *orthologous*, whereas genes that are homologous as a result of gene duplication are *paralogous* (Li, 2006). Paralogous genes that are functionally redundant and selectively nearly neutral can result in one copy being mutated into a functionless sequence called a *pseudogene*, or they can be deleted altogether. On the other hand, some duplicated genes can be beneficial from their time of origin because of dosage effects (Kondrashov et al., 2002) and may ultimately be important for speciation. In Passeriformes (perching birds) this may be the case for growth hormone (GH) paralogs, which have undergone differential selection since their divergence (Yuri et al., 2008). In a process called *subfunctionalization*, each paralog adopts partial function of their ancestral gene (Nowak et al., 1997; Lynch and Force, 2000). Changes in gene expression immediately following duplication as a result of

subfunctionalization appear to be common (Gu et al., 2002). Duplicated genes can also diverge to produce novel functions in a process known as *neofunctionalization* (Zhang, 2003). For example, some duplicated members of the *RNaseA* gene superfamily in primates evolved a novel antibacterial function that was not present in the common ancestral gene or its descendants (Zhang et al., 1998). Through the acquisition of novel functions, gene duplication plays a vital role in generating diversity at both molecular (Eirin-Lopez et al., 2004) and organismal (Hittinger and Carroll, 2007) levels and can increase the phenotypic complexity and diversity of animals (Ohno, 1970; Maniatis and Tasic, 2002). Protein coding genes, regulatory genes, and RNA noncoding genes are all subject to gene duplication. Gene duplication occurs in all three major domains of life: Bacteria, Archaea, and Eukarya (Zhang, 2003), and comparative analyses have suggested an average origin rate of new gene duplicates on the order of 0.1 per gene per million years (Gao and Innan, 2004; Osada and Innan, 2008). However, Gao and Innan (2004) and Osada and Innan (2008) have shown that using silent-site divergence as Lynch and Connery (2000) did can seriously overestimate the rate of duplication, due to spurious similarity of paralogs on account of gene conversion. More robust rates can be estimated from phylogenetic trees in which the role of gene conversion can be properly assessed, and they have suggested rates of duplication for yeast and *Drosophila* one to three orders of magnitude lower than Lynch and Connery's rates.

The ubiquity of gene duplication and its power to generate material on which selection may act is an especially interesting topic in birds because of their uniquely structured genomes. Birds have smaller genomes than those of any other amniote group (Gregory, 2002), but they also have a substantially reduced density of active repetitive elements (Shedlock, 2006) and segmental duplications and pseudogenes (Hillier et al., 2004). Birds may also have fewer protein-coding genes in their genomes than mammals have, with roughly 18,000 in chickens compared with approximately 22,000 in humans (Hillier et al., 2004). By comparing the genomes of chicken with those of humans and pufferfish (*Fugu*), the difference in gene count can be explained by significantly less lineage-specific gene duplication along the lineage of birds. For example, gene duplication by retrotransposition appears to be rare in birds compared with mammals, which have roughly 300 times the number of retrotransposed gene duplicates than in chicken (over 15,000 compared with 51) (Hillier et al., 2004). Of the 51 duplicates detected in the chicken genome by Hillier et al. (2004), 36 appear to be pseudogenes. This is probably due in part to the abundance of long interspersed nuclear elements (LINEs) in mammal genomes, which are probably responsible for the reverse transcription of retrotransposed gene duplicates. Chicken repeat 1 (CR1) is the dominant active transposable element in chicken and other archosaurs (Shedlock et al., 2007), a feature that may account for the lack of retrotransposed gene duplicates due to this element's inability to copy polyadenylated mRNA (Haasa et al., 2001; Hillier et al., 2004).

The whole-genome sequencing of the chicken (*Gallus gallus*) provided an unprecedented window into the architecture of bird genomes (Hillier et al., 2004). This study revealed a number of important details of gene duplication. For example, most of the expansions in gene families within chicken are associated with the immune system and host defense against parasites (Ota and Nei, 1994; Nei et al., 1997). The chicken genome project found an expansion of the scavenger receptor cysteine-rich (SRCR) domain (Hillier et al., 2004), a highly conserved protein module involved in the innate immune system (Sarrias et al., 2004). Linked with the major histocompatibility complex (MHC) class I gene cluster in humans are certain olfactory receptors

which also underwent a lineage-specific expansion within the chicken. These olfactory receptor genes, related to two orthologs in human (*OR5U1* and *OR5BF1*), appear to have expanded within birds relatively recently to constitute the majority of the over 200 olfactory receptor genes in the chicken genome. The expansions of these genes could conceivably be part of the genetic mechanism linking the immune and olfactory systems with mate choice, kin recognition, and social interactions in birds (Zelano and Edwards, 2002). Other gene expansions were probably also important for key innovations in birds, such as expansions in the keratin gene family, which are the proteinaceous building blocks of feathers and therefore vital for thermoregulation, sexual display, and flight. However, even in the absence of selection (most gene duplicates are thought to be pseudogenized during evolution) gene duplication may still play a critical role in creating postmating reproductive barriers that aid in speciation (Lynch and Conery, 2000).

The age of gene duplications within lineages is also of fundamental interest in the study of gene family evolution. Analysis of animal, plant, and fungal genomes has shown that the majority of duplications are recent, because most duplicated genes are thought to be pseudogenized shortly after the duplication event (Lynch and Conery, 2000). It is currently unknown if the small, streamlined genomes of birds deviate from this pattern, given the lower number of total genes, paucity of transposable elements, and highly recombinant microchromosomes found in chicken.

Here we explore the evolution of chicken gene families within the larger context of amniote evolution. Phylogenomics (Eisen, 1998) attempts systematically to reconstruct the phylogeny of gene families across multiple complete genomes so as to deduce putative protein homologies and functions as well as the specific gene duplications and losses responsible for gene family diversity. The phylogenomic approach to multigene family evolution has shown recent success in many different clades, such as vertebrates (Zmasek and Eddy, 2002; Storm and Sonnhammer, 2003; Li et al., 2006; Huerta-Cepas et al., 2007) and 16 fungi species (Wapinski et al., 2007). Recently, Rasmussen and Kellis (2007) deduced a pattern of substitution rates within and among *Drosophila* and yeast species by using machine learning to compare the gene trees in multigene families to the presumed species tree. In addition, Heger and Ponting (2007) developed a comprehensive database of orthologous and paralogous gene sets for several clades, including amniotes. Here we follow a complementary approach to characterize the dynamics and evolution of gene duplication within birds. Whereas many previous approaches to recognizing paralogs use reciprocal BLAST hits and pairwise comparison of orthologs, we apply Ensembl phylogenies and a variety of phylogenetic and comparative analyses within a clade of two dozen sequenced amniote species to analyze the dynamics of gene duplication in the complete chicken (*Gallus gallus*) genome.

2 METHODOLOGY AND COMPUTATIONAL APPROACH

For our analysis, we obtained 25,363 gene trees from Ensembl's (v50) gene tree database. Ensembl's gene trees were produced with their custom computational pipeline. Initially, family clusters were defined by single-linkage clustering of best reciprocal BLAST hits between amino acid translations of all genes within the Ensembl database (Hubbard et al., 2007). Each cluster was then processed by Ensembl to create a multiple amino acid alignment using the alignment algorithm MUSCLE

(Edgar, 2004) and a maximum likelihood phylogenetic tree using PHYML (Guindon and Gascuel, 2003). Our database of gene trees assembled from Ensembl contains 611,441 genes from 39 species, including 25 mammals, *Gallus gallus*, *Xenopus tropicalis*, five fish, six invertebrates, and *Saccharomyces cerevisiae*. There are 17,487 annotated chicken (*Gallus gallus*) genes in the Ensembl database, 13,649 (78%) of which belong to one of the 7785 gene trees.

To study gene duplication along the bird lineage, we identified the divergence between birds and mammals in each gene tree (Figures 1 and 2). This event is represented as a node that is parental to one clade of chicken-specific genes and another clade of orthologs and paralogs found only in mammals. Although we describe the chicken paralogs as chicken-specific, we do so only because chicken is the only bird in our analysis; presumably, sampling of additional genomes within Reptilia would reveal many of these genes and gene duplication to be shared with other birds and nonavian reptiles. If duplications occurred prior to the common ancestor of birds and mammals, a single family may have diversified significantly into multiple orthologs through speciation. In addition, some of these ancient gene duplications can be partially obscured when there is complete gene loss in either the bird or mammalian lineage. We were able to identify 12,094 unambiguous and datable divergences due to the bird–mammal divergence (using *Gallus* and *Homo*). For each such divergence we isolated the tree of amniote genes rooted at the amniote ancestor for further analysis. Outgroup species

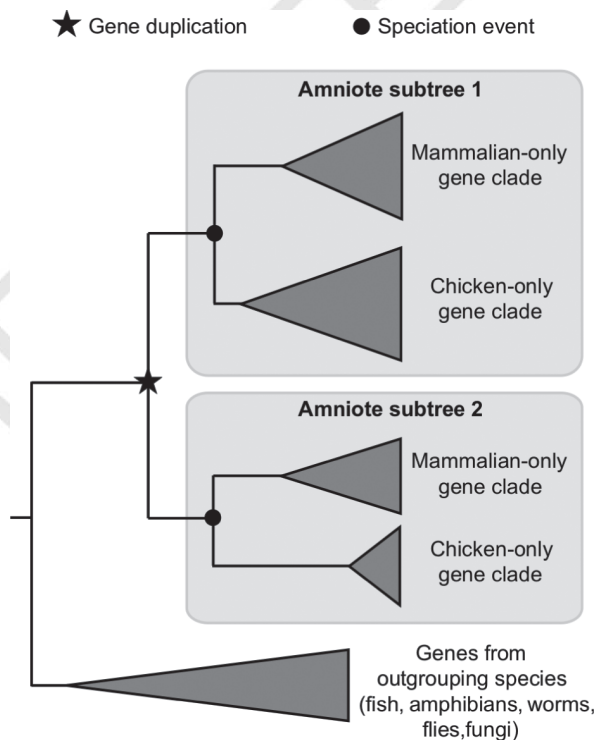


Figure 1 Phylogenetic pattern that gene duplication takes in multispecies comparisons.

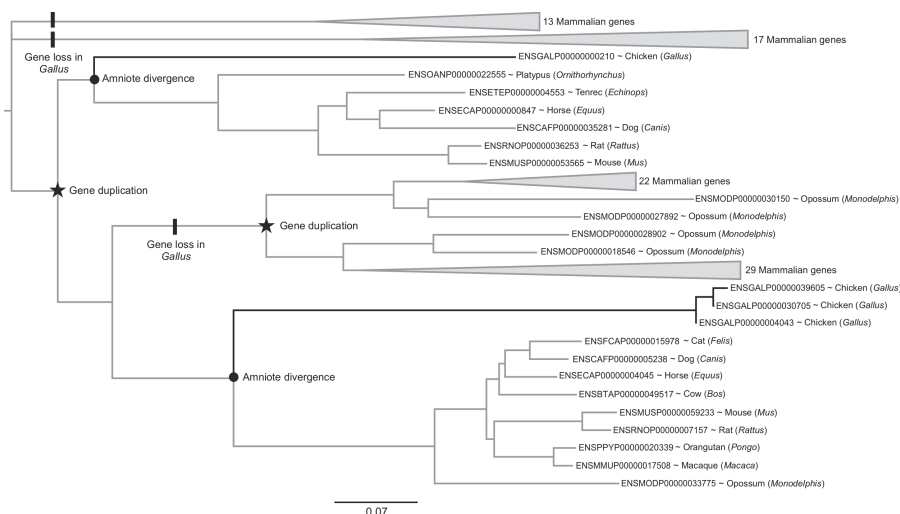


Figure 2 Olfactory receptor gene tree as defined by single-linkage clustering of best reciprocal BLAST hits between amino acid translations within the Ensembl database. It includes four chicken genes related to chicken olfactory receptor 7 (*COR7*). In Table 1 we give statistics for a gene tree including the three *Gallus COR7* family members that cluster together, as the fourth *Gallus* member split before the amniote divergence, and thus was part of a second gene tree. These three chicken genes, together with the nine mammalian paralogs that form a monophyletic group after the amniote speciation, are an example of a gene family used in various analyses throughout this chapter. The tip names are designated as Ensembl proteins.

(fish; *Fugu*) were used to help position the root of the gene tree. We expect this protocol to yield chicken gene families that differ from those circumscribed, for example, solely by reciprocal BLAST hits; sometimes our definition of chicken gene families includes genes not detected by nonphylogenetic methods, whereas in other instances our approach will miss some genes that manual inspection and curation would have revealed. On the other hand, our approach has the advantage of being objective and repeatable, and can easily be extended to study the dynamics of gene duplication in other taxa.

3 RESULTS: DYNAMICS OF CHICKEN-SPECIFIC GENE DUPLICATION

We found that the distribution of the number of paralogs per chicken gene family was heavily skewed toward small families (Figure 3). Whereas nearly 30 gene families had three members, only six families had more than nine. These figures focus only on those gene families that duplicated after the divergence of chicken from the amniote ancestor (by comparing *Gallus* with *Homo*). Thus, very ancient gene families—many potentially with large numbers of family members—do not figure into this calculation. Still, we were surprised by the small number of very large (>20) gene families in the chicken. Gene family size has been found to follow a power-law distribution in animals, a pattern thought to be produced by differential rates of pseudogenization among families (Huynen and van Nimwegen, 1998; Hughes and Liberles, 2008). Moreover,

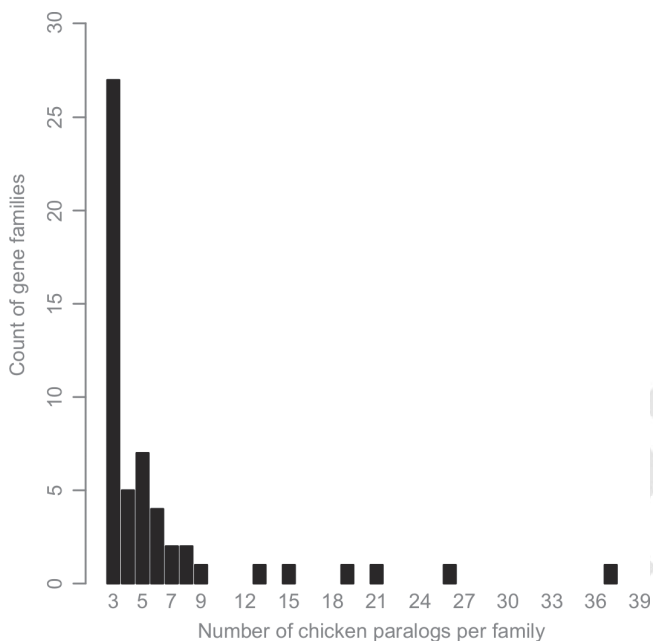


Figure 3 Count of gene families vs. the number of chicken genes in amniote gene trees.

birth–death models and purifying selection appear to account for much of the conservation seen within lineage-specific paralogs (Ota and Nei, 1994; Piontkivska et al., 2002; Piontkivska and Nei, 2003; Eirin-Lopez et al., 2004).

Substitution rates within amniote lineages are known to be quite variable, and we expect similar rate variation among paralogous members of chicken gene families. We estimated dates of divergence for each node in our chicken multigene family trees by using a penalized likelihood model, as implemented in the r8s rate analysis program (Sanderson, 2003), and by using 310 million years as an estimate for the *Gallus/Homo* common ancestor (Benton and Donoghue, 2007). The penalized likelihood model finds the optimal trade-off between maximizing the likelihood of a Poisson process for nucleotide substitution and a penalty term for rate variation between neighboring branches. The weighting of these two terms is determined by a coefficient λ , such that higher values of λ greatly penalize rate variation in favor of a clocklike model, and lower values allow a large amount of rate variation. Using a cross-validation procedure (Sanderson, 2003), we determined the optimal choice of λ from the possible values 10^{-2} to 10^3 . Within our amniote gene trees we found that over 54% of them were best explained by a λ value of 10^{-2} , indicating that the vast majority of gene families exhibit substantial rate variation among lineages (Figure 4). This rate variation probably stems from two sources: natural deviations in the clock as commonly found, for example, in phylogenetic analyses of different species; and bursts of adaptive evolution among newly evolved gene family members. Under the first hypothesis, it might be expected that larger gene families would exhibit more rate variation than would small ones, and that the incidence of rate variation would increase with family size. However, we did not find this trend when we regressed λ on gene family size

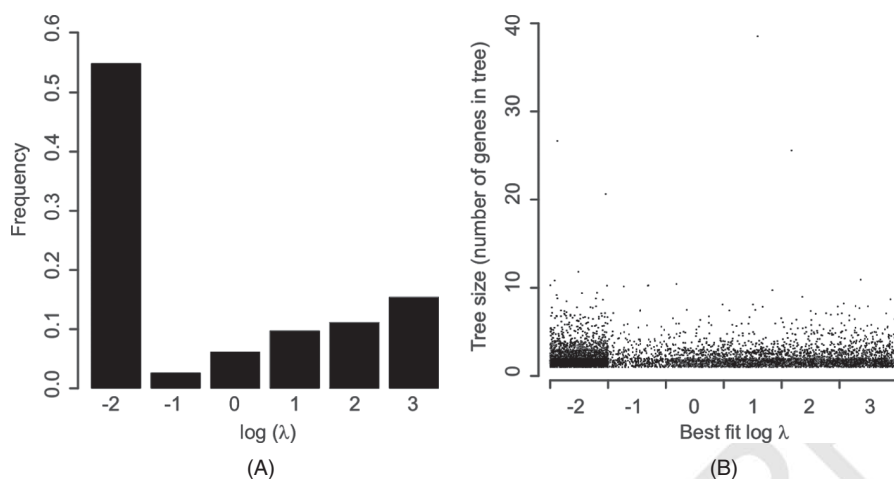


Figure 4 Branch length rate variation (molecular clock; λ) for amniote gene trees (25 mammals and chicken). **(A)** Distribution of branch length rate variation binned in different λ values. Low values of λ represent highly variable rates, and high values of λ represent clocklike rates. **(B)** The best fit (by likelihood criterion) λ by tree size. r^2 , relating tree size with $\log \lambda$, is equal to 0.0002843 ($p = 0.08$, $n = 6,901$). Values of λ were altered for display purposes to illustrate density.

(Figure 4). Many of the gene families best explained by a $\log \lambda$ value of -2 showed levels of divergence that suggests duplication since the Cenozoic era, especially since the Neogene period. For this reason we suspect that much of the rate variation among gene family members may in fact be due to adaptive bursts, because generation time effects among different lineages of birds are expected to influence rate variation only for those gene families that duplicated prior to the chicken's divergence from other lineages. Of the other categories of substitution rate variation among gene family members, the class best explained by $\lambda = \log(3)$ was the next most common. Substitution rates among gene family members in this category are fairly clocklike.

The ages of gene duplications in chicken are distributed exponentially, with most duplications occurring recently (Figure 5). This pattern is consistent with previous analyses, regardless of whether synonymous substitutions or phylogenetic analyses are used (Lynch and Conery, 2000), and suggests that assuming a relatively constant rate of gene duplication, most genes are pseudogenized or eliminated from the genome soon after duplication. However, these results are also consistent with widespread gene conversion between paralogs (Gao and Innan, 2004; Osada and Innan, 2008), in which the duplication event between highly similar sequences would be older than direct sequence comparisons would suggest. With the full genome sequence of only one bird, it is impossible to untangle the relative contributions of gene loss or gene conversion in producing the skewed age distribution of gene duplications. Moreover, this pattern could also suggest that concerted evolution might be more common among chicken gene families than in other groups. For example, concerted evolution among major histocompatibility complex (MHC) paralogs in birds is thought to occur more frequently, and over a shorter time scale, than in mammals (Hess and Edwards, 2002),

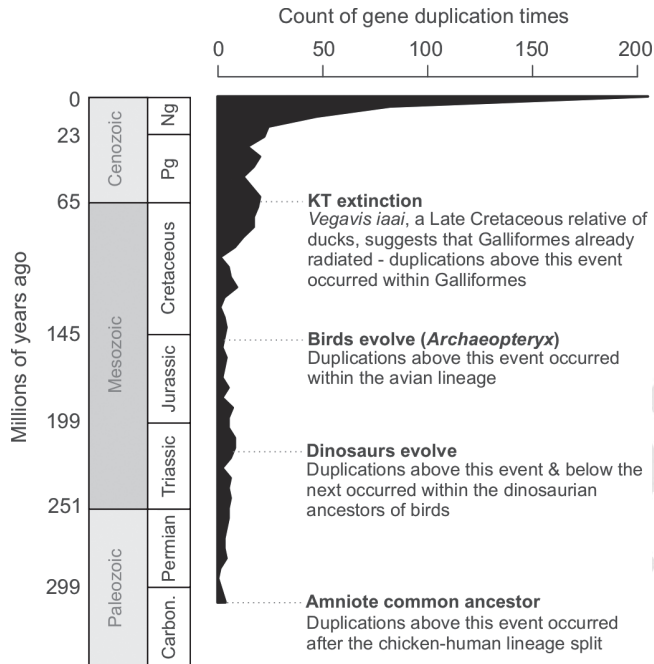


Figure 5 Depth of paralogs on the lineage leading to *Gallus gallus* that evolved after the divergence between *Gallus* and *Homo*. Pivotal events during the evolution of this lineage are noted on the figure. Pg, Paleogene; Ng, Neogene.

and the phylogenetic scale over which MHC orthologs can be identified may be smaller in birds than in mammals (but see Burri et al., 2008).

4 EXAMPLES OF FAMILIES WITH CHICKEN-SPECIFIC DUPLICATIONS

Gene family composition is shaped both by gene gain and loss, yet as other researchers have noted (Furlong, 2005), gene family expansion is easier to detect, especially when annotation is not complete and gaps remain in recent genome builds. We examined several gene families containing lineage-specific expansions in the chicken, using the amniotic gene tree rooted at the chicken–human divergence. In Table 1 we describe the dynamics of five representative families: Toll-like receptors, hemoglobin, ovalbumin-related serpins, four subfamilies of olfactory receptors, and keratin. These families were selected for their variety in size, age, and function and because the annotation and family membership could be at least partially cross-validated with recent studies.

4.1 Toll-like Receptors

Temperley and colleagues (2008) describe the evolutionary history and chromosomal location of chicken Toll-like receptors (TLRs), a family that is part of the innate immune system and is characterized by an ancient, highly conserved pathogen-recognition domain that triggers an inflammatory response. These authors discovered that while

TABLE 1 Summary of Properties for Specific Gene Families in the Amniote–Bird Lineage^a

Family Name	Number of Amniote Paralogs	Number of Chicken Paralogs	Mean Branch Length from Tips to Duplication in Chicken	Estimated Duplication Time (Mya)	Sequence Divergence Before Chicken Duplications After Amniote Divergence	Log λ
Toll-like-receptors (<i>TLR2A</i> and <i>TLR2B</i>)	18	2	0.075	67	0.260	-2
Ovalbumin B serpins (gene <i>X</i> , gene <i>Y</i> , and <i>ovalbumin</i>)	73	3	0.141	107	0.227	-2
Ovalbumin B serpins (<i>MENT</i> , <i>serpinb10</i>)	24	2	0.183	203	0.094	-2
Hemoglobin β -globin (β^H , ρ , ϵ)	69	3	0.075	122	0.086	-2
Olfactory receptors (orthologous to <i>OR5U1</i> and <i>OR5BF1</i> in <i>Homo</i>)	511	196	0.205	281	0.320	-2 ^b
Olfactory receptors (related to <i>COR7</i>)	12	3	0.014	15	0.394	-2
Olfactory receptors (related to <i>COR 1-6</i>)	60	4	0.275	238	0.084	2
Olfactory receptors (small cluster on chromosome 1)	207	3	0.049	30	0.421	-2
β -keratin	117	117	0.635	NA ^c	NA	NA

^aNumber of amniote paralogs is defined as the number of genes across all species in the amniote tree (25 mammals and chicken; see Ensembl for exact species). Mean branch length is the average path length (in substitutions/site) between chicken paralogs (tips) and their common ancestor (the first chicken duplication). Across all families this length is on average 0.235 substitution/site. Estimated duplication time is the time of first chicken duplication in millions of years. Sequence divergence before the duplications in chicken is given in substitutions/site between the root of the amniote tree (*Gallus/Homo* common ancestor) and the first chicken duplication. λ is a molecular-rate variation parameter (low values are highly variable rates; high values are clocklike).

^bThis λ is calculated from *Homo* and *Gallus* for computational limitations due to large family size.

^cNA, cannot be computed because of chicken-only expansions.

chickens and humans both have 10 receptors, only four genes in chicken maintain one-to-one orthology with mammalian genes; much gain and loss has occurred in every lineage. In the chicken, their study suggested that a duplication event estimated at 66 million years ago (Mya) gave rise to *TLR2A* and *TLR2B*, orthologs to the single *TLR2* in mammals. Three other genes, two that duplicated in tandem (*TLR1LA* and *TLR1LB*, estimated duplication time 147 Mya), as well as *TLR15*, have no mammalian counterpart. Other mammalian members have been pseudogenized or fully lost in chicken. Using our automated phylogenetic approach, we are able to analyze one of these chicken-specific expansions: the duplication event that gave rise to *TLR2A* and *TLR2B* in chickens. We obtained the same estimate of duplication time (67 Mya) as did Temperley et al. (2008). As with all of the families that we examined in detail, there was considerable rate variation among gene lineages ($\log \lambda = -2$). The sequence divergence in the Toll-like receptor family is greater than the sequence divergence found for 50.8% of other chicken gene families. We could not analyze the other duplication in chicken, as our data set from Ensembl was missing one of the chicken-specific genes (*TLR1LB*). *TLR15*, another gene unique to birds, had no mammalian ortholog, so it also was not included in our amniote gene tree.

4.2 Ovalbumin-Related Serpins

Another gene family with documented chicken-specific expansions is the ov-serpin family, also called ovalbumin-related serpins or clade B serpins. Benarafa and Remold-O'Donnell (2005) examine the phylogenetic relationship between the chicken members (some of which function as egg-white storage proteins) and their mammalian counterparts (involved in diverse roles such as embryogenesis, inflammation regulation, and angiogenesis). The initial duplication is thought to have occurred very early in the vertebrate lineage; and, like TLRs, the family is also marked by recent lineage-specific expansions and losses. Chickens have 10 members and humans have 13 members. Three genes in chicken—*ovalbumin* and ovalbumin-like genes *X* and *Y*—are paralogs and lack a human ortholog. Another gene, with a single human ortholog, seems to have duplicated to produce the chicken genes *Serpinb10* and *MENT* (mature erythrocyte nuclear termination state-specific protein). The remaining family members from chicken each have single human orthologs. Among the two subfamilies with chicken-specific expansions, rate variation is substantial ($\log \lambda = -2$). Moreover, the sequence divergence in the subfamily containing *ovalbumin* and ovalbumin-like genes *X* and *Y* is greater than the sequence divergence found for 69.6% of other chicken gene families that have duplicated since the chicken–mammalian split. The sequence divergence in the other ovalbumin subfamily (*serpinb10* and *MENT*) is greater than the sequence divergence found for 77.6% of other chicken gene families that have duplicated since the chicken–mammalian split.

4.3 Hemoglobin

Metabolic rate is an important trait that governs many organismal characters, from growth strategies to sustained physical activity. In amniotes, an elevated metabolism (endothermy) has only evolved within two extant groups, birds and mammals, although paleontologists suspect that many extinct dinosaurian lineages possessed endothermy (de Ricqlès et al., 2001; Horner et al., 2001). Whereas the typical mammalian and avian

condition is homeothermy (roughly constant body temperature), some birds, such as swifts, hummingbirds, and nightjars, are facultatively poikilothermic, a condition in which their usually elevated body temperature can vary over a wider range than that seen in mammals (e.g., Lane et al., 2004). The hemoglobin multigene family is closely associated with metabolism and the respiratory system.

Hemoglobin, a multidomain protein, has rapidly diversified with vertebrate lineages (Gribaldo et al., 2003). For example, α -globin underwent rapid duplication and deletion in mammals (Hoffmann et al., 2008). Based on an analysis of the platypus genome, which incorporated information from flanking loci, a recent model (Patel et al., 2008) proposes that the β -globin paralogs arose from a single transposition in the amniote ancestor followed by independent duplication in birds and mammals. From our data set, the β -globin paralogs β^H , ρ , and ϵ appear as a chicken-specific expansion, consistent with this model. Rate diversity is high in this family as well, and sequence evolution (point mutations) in β -globin paralogs is similar to that seen in the Toll-like receptors family (the sequence divergence in the β -globin family is greater than the sequence divergence found for 50.8% of other chicken gene families).

4.4 Olfactory Receptors

Olfaction has recently gained much recognition as an important sensory modality for birds (Nevitt et al., 2008; O'Dwyer et al., 2008; Steiger et al., 2008); Historically, birds were assumed to communicate primarily via the visual or auditory systems, but behavioral and genomic data suggest that chemosensory perception plays a larger role. The chicken genome paper remarked on the surprisingly large group of avian-specific olfactory receptors: 218 genes were identified, representing an avian expansion orthologous to the human receptors *OR5U1* and *OR5BF1*. Our bioinformatic approach detected 196 genes belonging to this subfamily; the discrepancy is perhaps in part due to different genome builds. As in other recent work (Lagerström et al., 2006), we also identified other families of olfactory receptors with small expansions in chicken (three to four genes in our analysis). These include a cluster associated with the previously identified *COR1-6* genes (chicken olfactory receptor genes) on chromosome 5, a second cluster on chromosome 10 related to *COR7* (Figure 2), and a third cluster on chromosome 1. Interestingly, the subfamilies have very different extents of sequence divergence ranging from 26.7% (family containing *COR7*) to 88.4% (family containing *COR1-6*). The latter gene family also had a more clocklike rate than other families, which together with the large sequence divergence suggests that it is among the oldest gene duplications in chicken.

4.5 Keratin

The evolution of feathers in theropod dinosaurs was a major innovation that probably provided insulation for metabolically active animals, ornamentation for display, and in one lineage transformed arms into wings (Ji et al., 1998; Zhang and Zhou, 2000; Currie and Chen, 2001; Norell et al., 2002; Sawyer and Knapp, 2003). β -Keratins differ from the keratins found in nonavian reptiles and are the basic structural elements of feathers and therefore a gene family vital for the success of birds. In the publication of the first genome draft, the International Chicken Genome Sequencing Consortium (2004) noted the large expansion of the avian-specific keratin gene family, estimated

at around 150 members. This avian keratin family, which encodes proteins forming feathers and scales (Sawyer et al., 2000), is functionally and evolutionarily distinct from the mammalian hair-specific α -keratin, which like another main component of hair, the keratin-associated proteins, have no members in the chicken genome (Wu et al., 2008). Within nonavian reptiles, β -keratins probably duplicated by retrotransposition, resulting in the loss of introns in some paralogs; in birds, all paralogs have lost introns. Unequal crossover is also thought to expand and contract keratin gene arrays in birds, resulting in a tandem organization of multiple paralogs (Toni et al., 2007). We find that the amount of sequence evolution in keratin is among the highest for any chicken gene family (greater than the sequence divergence found for 98% of other chicken gene families). This high divergence is likely a consequence of the absence of comparisons with other reptiles, but it could also be due to the adaptive significance of these proteins within birds.

5 PROSPECTS AND CONCLUSIONS

Reptilia, including birds and nonavian reptiles, is the sister group of mammals and as such holds an important phylogenetic position for shedding light on patterns of gene duplication in amniotes (Wang et al., 2006). Reptiles are arguably more diverse than mammals in many traits; with about 17,000 species (about 10,000 in birds and 7000 nonavian reptiles) they are substantially more species-rich than mammals (about 5000 species) and possess a greater diversity of sex chromosome and sex determination systems (Organ and Janes, 2008). The chicken is currently the sole member of Reptilia with a draft genome and as such provides the only point of comparison of genome dynamics between mammals and their sister group. A greater understanding of genome and multigene family dynamics in mammals will undoubtedly require greater genome sampling and characterization in Reptilia.

Gene duplication and the families they produce are vital for generating the thread with which evolution weaves new adaptations and species. We have developed a pipeline for phylogenomic analysis of gene duplication in the chicken lineage, but our approach can be applied easily to any particular clade of interest. Our approach rests on the assumption that gene orthology and paralogy are best identified through phylogenetic analysis, and we delimit chicken-specific gene duplications by an approach (Figure 1) that combines initial identification and collection of gene copies across many vertebrates that show significant sequence similarity in Ensembl, followed by phylogenetic analysis of these gene sets; identification of particular nodes in these gene trees that correspond to gene duplications, in our case the mammal–bird divergence; identification of those gene clusters that diversify from these particular nodes; and statistical analysis of the gene trees collected. Many of the duplications we have identified here as “chicken-specific” in fact will be found to have duplicated in ancestors of the chicken, since orthologs of many chicken genes will no doubt be discovered in other reptile genomes as they emerge. Nonetheless, using an approximate time scale (Figure 5) we can estimate which chicken gene paralogs might be found in upcoming reptilian genome projects based on their estimated timing of duplication relative to the divergence times of species whose genomes are being compared. Our approach has the advantage of providing an objective means of identifying chicken-specific gene duplications, but of course when conducted on a genomewide scale, it will miss some

gene family members that manual curation will identify; we have illustrated this with some specific examples (Table 1). The loss of detail for some gene families is offset by the ability to study genome-wide distributions of multigene family dynamics; both approaches are required to provide an informed view of the dynamics of multigene family evolution in birds and relatives.

Phylogenomic approaches such as those presented here have only just begun to provide a window into the dynamics and importance of gene duplication within organisms. For example, nonprotein coding RNA paralogs are dispersed throughout the chicken genome; this, along with an unusual paucity of nonprotein coding RNA pseudogenes, suggests that they may not undergo the same processes of duplication (unequal crossover and retrotransposition) that characterize protein coding genes (Hillier et al., 2004). Currently available data are insufficient to address this and other hypotheses, because as of the time of this writing the genome of only one reptile species has been sequenced. But progress is quickly being made with the imminent release of the zebra finch (*Taeniopygia guttata*) and anole lizard (*Anolis carolinensis*) genomes. An increase in the number of genomes will permit more detailed quantitative comparison of the evolutionary dynamics of gene duplication in amniotes and other lineages and will help clarify the role of these gene duplications in organismal diversification.

REFERENCES

- Benarafa C, Remold-O'Donnell E. 2005. The ovalbumin serpins revisited: perspective from the chicken genome of clade B serpin evolution in vertebrates. *Proc Nat Acad Sci USA* 102(32):11367–11372.
- Benton MJ, Donoghue PCJ. 2007. Paleontological evidence to date the tree of life. *Mol Biol Evol* 24(1):26–53.
- Burri R, et al. 2008. Evolutionary patterns of MHC class II B in owls and their implications for the understanding of avian MHC evolution. *Mol Biol Evol* 25:1180–1191.
- Currie PJ, Chen P-J. 2001. Anatomy of *Sinosauropteryx prima* from Liaoning, northeastern China. *Can J Earth Sci* 38(12):1705–1727.
- de Ricqlès AJ, et al. 2001. The bone histology of basal birds in phylogenetic and ontogenetic perspectives. In Gauthier J, Gall LF (eds.), *New Perspectives on the Origin and Early Evolution of Birds: Proceedings of the International Symposium in Honor of John H. Ostrom*. New Haven, CT: Peabody Museum of Natural History, pp. 411–426.
- Edgar RC. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res* 32:1792–1797.
- Eirin-Lopez JM, et al. 2004. Birth-and-death evolution with strong purifying selection in the histone H1 multigene family and the origin of orphon H1 genes. *Mol Biol Evol* 21(10):1992–2003.
- Eisen JA. 1998. Phylogenomics: improving functional predictions for uncharacterized genes by evolutionary analysis. *Genome Res* 8:163–167.
- Furlong RF. 2005. Insights into vertebrate evolution from the chicken genome sequence. *Genome Biol* 6(2).
- Gao L-Z, Innan H. 2004. Very low gene duplication rate in the yeast genome. *Science* 306:1367–1370.
- Gregory TR. 2002. A bird's-eye view of the C-value enigma: genome size, cell size, and metabolic rate in the class Aves. *Evolution* 56(1):121–130.

- Gribaldo S, et al. 2003. Functional divergence prediction from evolutionary analysis: a case study of vertebrate hemoglobin. *Mol Biol Evol* 20(11):1754–1759.
- Gu Z, et al. 2002. Rapid divergence in expression between duplicate genes inferred from microarray data. *Trends Genet* 18(12):609–613.
- Guindon S, Gascuel O. 2003. A simple, fast, and accurate algorithm to estimate large phylogenies by maximum likelihood. *Syst Biol* 52(5):696–704.
- Haasa NB, et al. 2001. Subfamilies of CR1 non-LTR retrotransposons have different 50 UTR sequences but are otherwise conserved. *Gene* 265:175–183.
- Heger A, Ponting CP. 2007. Evolutionary rate analyses of orthologs and paralogs from 12 *Drosophila* genomes. *Genome Res* 17(12):1837–1849.
- Hess CM, Edwards SV. 2002. The evolution of the major histocompatibility complex in birds. *Bioscience* 52(5):423–431.
- Hillier LW, et al. 2004. Sequence and comparative analysis of the chicken genome provide unique perspectives on vertebrate evolution. *Nature* 432(7018):695–716.
- Hittinger CT, Carroll SB. 2007. Gene duplication and the adaptive evolution of a classic genetic switch. *Nature* 449:677–681.
- Hoffmann FG, et al. 2008. Rapid rates of lineage-specific gene duplication and deletion in the alpha-globin gene family. *Mol Biol Evol* 25(3):591–602.
- Horner JR, et al. 2001. Comparative osteology of some embryonic and perinatal archosaurs: developmental and behavioral implications for dinosaurs. *Paleobiology* 27(1):39–58.
- Hubbard TJP, et al. 2007. Ensembl 2007. *Nucleic Acids Res* 35(database issue):D610–D617.
- Huerta-Cepas J, et al. 2007. The human phylome. *Genome Biol* 8:R109.
- Hughes T, Liberles DA. 2008. The power-law distribution of gene family size is driven by the pseudogenisation rate's heterogeneity between gene families. *Gene* 414(1–2):85–94.
- Huynen MA, van Nimwegen E. 1998. The frequency distribution of gene family sizes in complete genomes. *Mol Biol Evol* 15(5):583–589.
- Ji Q, et al. 1998. Two feathered dinosaurs from northeastern China. *Nature* 393:753–761.
- Kondrashov FA, et al. 2002. Selection in the evolution of gene duplications. *Genome Biol* 3(2):1–9.
- Lagerström MC, et al. 2006. The G protein-coupled receptor subset of the chicken genome. *PLoS Comput Biol* 2(6):493–507.
- Lane JE, et al. 2004. Daily torpor in free-ranging whip-poor-wills (*Caprimulgus vociferus*). *Physiol Biochem Zool* 77:297–304.
- Li W-H. 2006. *Molecular Evolution*. Sunderland, MA: Sinauer Associates.
- Li H, et al. 2006. TreeFam: a curated database of phylogenetic trees of animal gene families. *Nucleic Acids Res* 34:D572–D580.
- Lynch M, Conery JS. 2000. The evolutionary fate and consequences of duplicate genes. *Science* 290:1151–1155.
- Lynch M, Force A. 2000. The probability of duplicate-gene preservation by subfunctionalization. *Genetics* 154:459–473.
- Maniatis T, Tasic B. 2002. Alternative pre-mRNA splicing and proteome expansion in metazoans. *Nature* 418(6894):236–243.
- Nevitt GA, et al. 2008. Evidence for olfactory search in wandering albatross, *Diomedea exulans*. *Proc Nat Acad Sci* 105(12):4576–4581.
- Norell MA, et al. 2002. “Modern” feathers on a non-avian dinosaur. *Nature* 416:36–37.
- Nowak MA, et al. 1997. Evolution of genetic redundancy. *Nature* 388(6638):167–171.

- O'Dwyer TW, et al. 2008. Examining the development of individual recognition in a burrow-nesting procellariiform, the Leach's storm-petrel. *J Exp Biol* 211(3):337–340.
- Ohno S. 1970. *Evolution by Gene Duplication*. New York: Springer-Verlag.
- Organ CL, Janes DE. 2008. Evolution of sex chromosomes in Sauropsida. *Integrative Comparative Biol* 48(4):512–519.
- Osada N, Innan H. 2008. Duplication and gene conversion in the *Drosophila melanogaster* genome. *PLoS Genet* 4(12):e1000305.
- Ota T, Nei M. 1994. Divergent evolution and evolution by the birth-and-death process in the immunoglobulin V-H gene family. *Mol Biol Evol* 11(3):469–482.
- Patel VS, et al. 2008. Platypus globin genes and flanking loci suggest a new insertional model for beta-globin evolution in birds and mammals. *BMC Biol* 6(34).
- Piontkivska H, Nei M. 2003. Birth-and-death evolution in primate MHC class I genes: divergence time estimates. *Mol Biol Evol* 20(4):601–609.
- Piontkivska H, et al. 2002. Purifying selection and birth-and-death evolution in the histone H4 gene family. *Mol Biol Evol* 19(5):689–697.
- Rasmussen MD, Kellis M. 2007. Accurate gene-tree reconstruction by learning gene- and species-specific substitution rates across multiple complete genomes. *Genome Res* 17(12):1932–1942.
- Sanderson MJ. 2003. r8s: inferring absolute rates of molecular evolution and divergence times in the absence of a molecular clock. *Bioinformatics* 19(2):301–302.
- Sarrias MR, et al. 2004. The scavenger receptor cysteine-rich (SRCR) domain: an ancient and highly conserved protein module of the innate immune system. *Crit Rev Immunol* 24(1):1–37.
- Sawyer RH, Knapp LW. 2003. Avian skin development and the evolutionary origin of feathers. *J Exp Zool B* 298B(1):57–72.
- Sawyer RH, et al. 2000. The expression of beta (β) keratins in the epidermal appendages of reptiles and birds. *Am Zool* 40(4):530–539.
- Shedlock AM. 2006. Phylogenomic investigation of CR1 LINE diversity in reptiles. *Syst Biol* 55(6):902–911.
- Steiger SS, et al. 2008. Avian olfactory receptor gene repertoires: evidence for a well-developed sense of smell in birds?. *Proc R Soc B* 275(1649):2309–2317.
- Storm CEV, Sonnhammer ELL. 2003. Comprehensive analysis of orthologous protein domains using the HOPS database. *Genome Res* 13:2353–2362.
- Temperley ND, et al. 2008. Evolution of the chicken Toll-like receptor gene family: a story of gene gain and gene loss. *BMC Genom* 9(62).
- Toni M, et al. 2007. Hard (beta-) keratins in the epidermis of reptiles: composition, sequence, and molecular organization. *J Proteome Res* 6(9):3377–3392.
- Wang Z, et al. 2006. Tuatara (*Sphenodon*) genomics: BAC library construction, sequence survey, and application to the DMRT gene family. *J Hered* 97(6):541–548.
- Wapinski I, et al. 2007. Natural history and evolutionary principles of gene duplication in fungi. *Nature* 449:54–61.
- Wu DD, et al. 2008. Molecular evolution of the keratin associated protein gene family in mammals, role in the evolution of mammalian hair. *BMC Evol Biol* 8(241).
- Yuri T, et al. 2008. Duplication of accelerated evolution and growth hormone gene in passerine birds. *Mol Biol Evol* 25(2):352–361.
- Zelano B, Edwards SV. 2002. An Mhc component to kin recognition and mate choice in birds: predictions, progress, and prospects. *Am Nat* 160:S225–S237.

- Zhang J. 2003. Evolution by gene duplication: an update. *Trends Genet* 18(6):292–298.
- Zhang F, Zhou Z. 2000. A primitive enantiornithine bird and the origin of feathers. *Science* 290(5498):1955–1959.
- Zhang J, et al. 1998. Positive Darwinian selection after gene duplication in primate ribonuclease genes. *Proc Nat Acad Sci USA* 98:3708–3713.
- Zmasek CM, Eddy SR. 2002. RIO: analyzing proteomes by automated phylogenomics using resampled inference of orthologs. *BMC Bioinf* 3:14.