

OPEN

Extinction of chromosomes due to specialization is a universal occurrence

Jason Wilson^{1*}, Joshua M. Staley² & Gerald J. Wyckoff^{1,2,3}

The human X and Y chromosomes evolved from a pair of autosomes approximately 180 million years ago. Despite their shared evolutionary origin, extensive genetic decay has resulted in the human Y chromosome losing 97% of its ancestral genes while gene content and order remain highly conserved on the X chromosome. Five 'stratification' events, most likely inversions, reduced the Y chromosome's ability to recombine with the X chromosome across the majority of its length and subjected its genes to the erosive forces associated with reduced recombination. The remaining functional genes are ubiquitously expressed, functionally coherent, dosage-sensitive genes, or have evolved male-specific functionality. It is unknown, however, whether functional specialization is a degenerative phenomenon unique to sex chromosomes, or if it conveys a potential selective advantage aside from sexual antagonism. We examined the evolution of mammalian orthologs to determine if the selective forces that led to the degeneration of the Y chromosome are unique in the genome. The results of our study suggest these forces are not exclusive to the Y chromosome, and chromosomal degeneration may have occurred throughout our evolutionary history. The reduction of recombination could additionally result in rapid fixation through isolation of specialized functions resulting in a cost-benefit relationship during times of intense selective pressure.

The human Y chromosome has lost its ability to recombine with its once homologous partner, the X chromosome, except in its pseudoautosomal regions (PARs) at the termini of the X and Y chromosomes¹⁻³. This has resulted in the majority of the Y chromosome's gene content being inherited as a unit, known as the human MSY (male-specific region of the Y chromosome)⁴. Suppression of recombination occurred at five discrete time points, probably caused by inversions, that integrated each segment into the MSY and initiated the degradative processes⁵ that resulted in wide-spread gene deletion and loss^{1,2,6,7}. These evolutionary strata show a continuum of degeneration that is highly correlated with the age of X-Y gene pairs within each stratum^{1,2,4}. The oldest of which contains only four remaining genes, including the sex-determining factor SRY⁸. The degenerative nature of the Y chromosome has led some researchers to suggest it may lose all functional genes and become extinct in as little as 5 million years⁸⁻¹⁰, an evolutionary phenomenon that has been observed in other species¹¹⁻¹³. Recent research, however, suggests that the Y chromosome has maintained a stable assortment of genes for the last 25 million years^{3,14,15} through effective purifying selection on single-copy genes¹⁶, and intrachromosomal gene conversion of ampliconic sequences¹⁷⁻²¹. Despite conflicting views on the terminal fate of the Y chromosome, functional specialization and biased gene retention²²⁻²⁴ on the Y chromosome is believed to be unique in the genome²⁵ and may have played an essential role in Y chromosome degeneration.

The remaining functional genes in the human MSY fall into three classes: X-degenerate, ampliconic, and X-transposed^{3,4,26}. The X-transposed sequences are a result of an X-to-Y transposition that occurred after the divergence of the human and chimpanzee lineages, approximately 3-4 million years ago^{4,26}. These sequences remain 99% identical to their X counterparts⁴. In contrast, the X-degenerate sequences are single-copy MSY genes that are surviving relics of the ancestral autosomes from which the sex chromosomes evolved⁴. With the notable exception of SRY, these genes are functionally coherent²⁵, and ubiquitously expressed^{1,14,17}. Their homologous X counterparts also disproportionately escape X-inactivation and are subject to stronger purifying selection

¹University of Missouri-Kansas City School of Medicine, Department of Biomedical and Health Informatics, Kansas City, 64108, Missouri, USA. ²Kansas State University College of Veterinary Medicine, Department of Diagnostic Medicine/Pathobiology, Olathe, 66061, Kansas, USA. ³University of Missouri-Kansas City School of Biological and Chemical Sciences, Department of Molecular Biology and Biochemistry, Kansas City, 64108, Missouri, USA. *email: jw84d@mail.umkc.edu

than other X-linked genes¹⁷. Thus, researchers have suggested that this class of sequences is dosage-sensitive and potentially haplolethal¹⁷. The last class of functional genes in the MSY consists of nine protein-coding gene families that have undergone various levels of amplification⁴. Unlike the ubiquitously expressed X-degenerate genes, the ampliconic gene families are expressed primarily or exclusively in the testes^{4,18} and rely on intrachromosomal gene conversion to offset the degenerative nature of the MSY^{18–21}. Surviving Y-linked genes were therefore retained through two evolutionary mechanisms: effective purifying selection on single copy dosage-sensitive genes¹⁶ and intrachromosomal gene conversion of ampliconic sequences^{17–21}.

Wide-spread gene loss accompanied by preferential retention appears to be a unique phenomenon. A review of genomic evolution, however, suggests that these trends are not unique to the Y chromosome, with the relevant literature rarely being cross-cited²⁷. Recent research suggests that the ancestral vertebrate karyotype was much larger than previously estimated, consisting of an estimated 54 chromosomes²⁸ resulting from two ancestral whole-genome duplication (WGD) events^{28–31}. The majority of genes following a WGD event are rapidly lost or pseudogenized due to loss of function mutations^{7,32–34}. This loss has also been shown to continue on a power scale^{33,35}. Consequently, a large portion of the ancestral vertebrate chromosomes has been subsequently lost through fusion in the descent of the human lineage^{28,31}, explaining the apparent haphazard gene content of most autosomes¹. Highly expressed genes³⁶, dosage-sensitive protein complexes^{34,37}, and transcriptional and developmental regulators and signal transducers, however, are preferentially retained^{30,33,38–40}. Furthermore, these genes have been maintained through purifying selection³⁷, a trend that has been observed in ubiquitously expressed genes throughout the genome^{41–45}. The factors that led to the biased retention of ubiquitously expressed single-copy genes, therefore, appear not to be restricted to the evolutionary history of the Y chromosome and have been observed in the events following large scale duplications. The biased acquisition of male-advantage traits on the Y chromosome is a subject of more considerable ambiguity in the context of genomic duplications.

Subfunctionalization has been shown to increase the likelihood a gene will be preserved in duplicate due to partial loss of function mutations in both copies⁴⁶. This targeted divergence of the duplicates may lead to differential tissue expression of the paralogs^{34,35,45,47–49} and has been proposed to occur frequently following WGD events⁵⁰. If the remaining functions are under selective constraint, the duplicates will likely remain in the population⁴⁷. A lack of genome-wide representation of subfunctionalized gene pairs, however, suggests that this may be a transition phase to neofunctionalization due to an absence of purifying selection on the redundant portions of the gene⁵¹, an evolutionary phenomenon known as the subneofunctionalization model⁵². In 2009, Wilson and Makova suggested that suppression of recombination could be thought of as a duplication event and showed X-Y genes followed similar patterns of evolution following recombination suppression as duplicated paralogs^{53,54}. Following a review of experimental data, they also concluded that the acquisition of unique expression patterns and functions might have contributed to the retention of Y-linked genes. Strong expression reduction has also been implicated in the evolution of Y genes towards testis specificity⁵⁵. The biased content of male reproductive genes on both sex chromosomes^{4,56–58}, therefore, suggests that subfunctionalization of Y-linked genes could explain the initial retention and accelerated divergence of male-advantage genes, as new evolutionary features typically bear marks of their ancestry⁵⁰.

The WGD events at the origin of the vertebrate lineage may have had significant impacts on biological complexity and evolutionary novelties of the time due to the large-scale increase in genetic redundancy³⁸. The mechanisms by which this was achieved and the selective pressures resulting in differential chromosome survival remain unknown. If the evolutionary history of the Y chromosome provides a model of genomic evolution, it would suggest that large scale duplication events allowed genes to subfunctionalize and experience periods of relaxed purifying selection through relief from pleiotropic constraints that were operating on single-gene loci⁴⁶. It has also been hypothesized that the Y chromosome's long-term fragility may be driven by short-term selective pressures⁵⁹, the most obvious of which is the accumulation of sexually antagonistic alleles in a non-recombining portion of the genome^{59–65}, a phenomenon that is supported by the transposition of male-advantage genes into the MSY from autosomes^{1,3,4,14,66–68}. The rapid evolution of male reproductive genes^{45,69–71} and the implication of inversions in local adaptation^{7,72}, however, suggest that functional isolation may become selectively favored even in the absence of sexually antagonistic traits under certain circumstances, despite the deleterious effects of reduced recombination. In order to test these hypotheses, and in lieu of the large amount of literature pertaining to expression, we analyzed the nonsynonymous to synonymous mutation rate (K_a/K_s) of 6,734 human genes with surviving mammalian orthologs in the context of their Gene Ontology (GO) annotations and chromosomal locations to determine if functional specialization and genomic isolation convey a selective advantage, respectively.

Results

Functional diversity of orthologs. Two primary paths to survival have occurred on the Y chromosome. Broadly expressed dosage-sensitive genes have been maintained through purifying selection¹⁶, while amplification and gene conversion have supported testis-specific genes^{17–21}. Selection for conservation through amplification and gene conversion of testis-specific sequences and the rapid evolution of male reproductive genes^{45,69–71} suggest that adaptability may be a selected phenomenon. Evolution of testis-specific functions is also believed to have preceded amplification on the Y chromosome¹⁷, suggesting subfunctionalization may have facilitated their initial retention⁴⁵ and subsequent divergence by relieving redundant portions of the genes from adaptive constraint.

To determine if this is a universal trend, we analyzed the divergence of human genes across their mammalian orthologs to provide a conservative estimate of the degree to which newly duplicated genes may diverge⁷³ following subfunctionalization. Summary statistics can be found in Supplementary Table S1. The results of our analysis suggest that a human gene's average K_a/K_s across its related orthologs and number of GO annotations are positively skewed (skew = 3.92 & 3.06, respectively) (Supplementary Fig. S1). Additionally, average K_a/K_s is

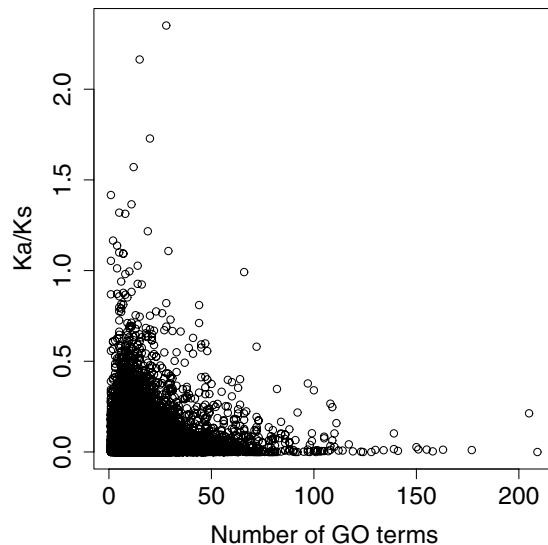


Figure 1. Ortholog K_a/K_s -GO Relationship: Average K_a/K_s values of 6,734 human genes across their related orthologs plotted against their number of GO terms. Each data point represents a human gene with surviving mammalian orthologs ($n = 6,734$). The x-axis corresponds to the number of gene ontology annotations that were found to be associated with each gene. The y-axis corresponds to a human gene's average K_a/K_s value across its related species comparisons. The results of a gamma regression analysis indicated the association was statistically significant ($p = 8.25 \times 10^{-20}$). This figure was generated using R statistical software (Version 3.5.3).

zero-inflated. This suggests that the majority of orthologous genes are under purifying selection and related to a small set of functions. As a gene's average K_a/K_s value appears to be negatively associated with its number of GO annotations (Fig. 1), a gamma-hurdle model was employed (see methods) to determine the statistical significance of this relationship. Our results suggest that a gene's average K_a/K_s decreases with increasing numbers of functional annotations ($p = 8.25 \times 10^{-20}$) (Supplementary Fig. S2), and the probability that it is entirely conserved increases ($p = 5.25 \times 10^{-9}$, odds-ratio = 1.011) (Supplementary Fig. S3).

The results of this analysis suggest that genes with high functional diversity are under more intensive purifying selection than their more functionally specific counterparts. These findings parallel those showing higher levels of purifying selection on broadly expressed essential genes throughout the genome^{41–45} as well as on the Y chromosome¹⁷ and suggest an association between the two. We conclude that functional specificity and reduced expression are associated with relaxed purifying selection, suggesting that subfunctionalization of duplicated paralogs could result in differential tissue expression^{34,35,45,47–49} and accelerated protein divergence.

Genomic isolation of functional annotations. Next, we were interested in determining if functional isolation can provide a selective advantage in the absence of sexual antagonism. The rapid evolution of male reproductive genes^{45,69–71} and the implication of inversions in local adaptation^{7,72} suggest that the localization of functionally related genes may accelerate protein divergence and facilitate adaptability. To determine if localization of genes related to a given function is associated with reduced purifying selection, we analyzed the average K_a/K_s of GO annotations related sequence comparisons with respect to the genomic distribution of their related genes. Summary statistics can be found in Supplementary Table S2. GO annotations show positively skewed distributions for their number of associated genes (skew = 34.68), number of chromosome arms they are expressed on (skew = 2.58), and average K_a/K_s (skew = 2.85) (Supplementary Fig. S4). This suggests that the majority of functional annotations we analyzed are carried out by a limited number of genes, are expressed in specific locations and under purifying selection. Their relationships with one another, however, suggest all three trends are not typically present at the same time.

A GO term's average K_a/K_s value appears to be negatively associated with the number of chromosome arms it is expressed on and its number of related genes (Fig. 2). Due to the positive skewed nature of these distributions (Supplementary Fig. S4), a gamma model with a log link was used to determine if a function's number of related genes or expressed chromosome arms is significantly associated with its average K_a/K_s . Increasing the number of genes or expressed chromosome arms related to a given function, however, increases the probability one K_a/K_s value is non-zero. Zero average K_a/K_s values in the context of this analysis, therefore, are not informative and were removed from the analysis, negating the need for a hurdle method. The two predictive variables were fit separately to determine their individual effects. The results of our analyses suggest a function's average K_a/K_s decreases with the number of chromosome arms it is expressed on ($p = 7.01 \times 10^{-19}$, Supplementary Fig. S5), however, a function's number of related genes was non-significant ($p = 0.05$, Supplementary Fig. S6). This suggests that genomic isolation is more strongly associated with relaxed purifying selection than a function's number of related genes. The non-significance of a function's number of related genes additionally suggests that low K_a/K_s values for anno-

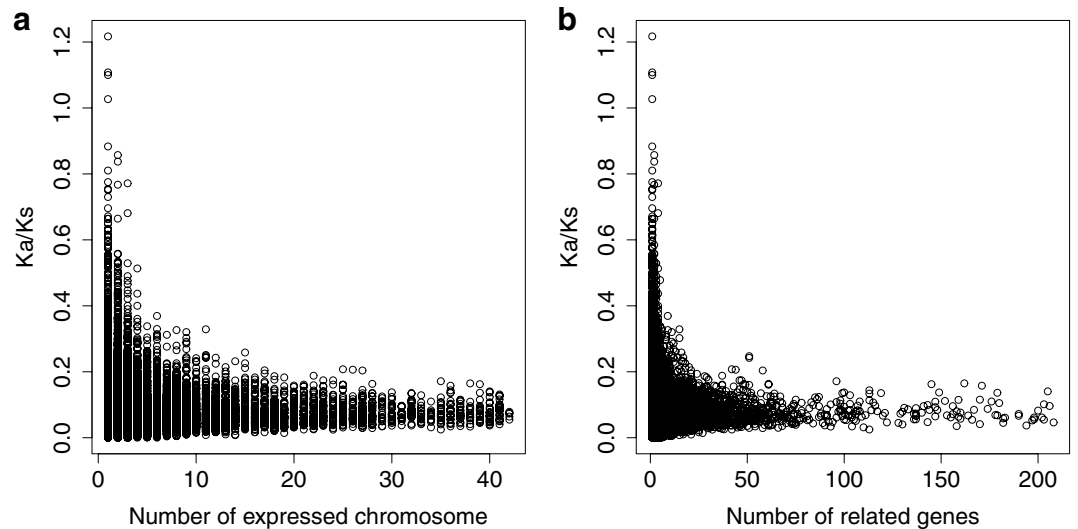


Figure 2. GO K_a/K_s Relationships: **(a)** Average K_a/K_s values of the 11,016 GO terms with unique gene sets in our dataset across each term's related genes and human sequence comparisons plotted against the number of chromosome arms a given GO term was found. **(b)** Average K_a/K_s values of the same GO terms plotted against their number of related genes (trimmed within 3 SD of the mean number of genes for clarity). Each data point represents a GO annotation ($n = 11,016$). **(a)** The x-axis corresponds to the number of chromosome arms a given GO annotation was expressed on (the number of chromosome arms containing a gene that was associated with the given GO annotation). The y-axis corresponds to a GO annotation's average K_a/K_s value across its related genes and their species comparisons. The results of a gamma regression analysis indicated the association was statistically significant ($p = 7.01 \times 10^{-19}$). **(b)** The x-axis corresponds to the number of genes a given GO annotation was associated with. The y-axis corresponds to a GO annotation's average K_a/K_s value across its related genes and their species comparisons. The results of a gamma regression analysis indicated the association was not statistically significant ($p = 0.05$). The accompanying figures were generated using R statistical software (Version 3.5.3).

tations expressed on a large number of chromosome arms cannot be attributed to convergence to the genome-wide average alone.

We expect that the majority of functions related to a large number of genes and expressed throughout the genome are higher-level ontology functions. Genes that are beneficial in increased dosage, however, are preferentially retained following duplication events^{73,74}. Thus, large scale duplications may have resulted in the stability of higher-level functions, while relieving more redundant duplicates from adaptive conflict. We conclude that functional isolation is associated with relaxed purifying selection on the genes related to that function, potentially through relief from background selection acting on more highly conserved linked sites⁷⁵. This finding parallels the accelerated evolution of Y-linked genes following recombination suppression and suggests isolation of functions may accelerate sequence divergence of their related genes through relaxation of purifying selection. These findings provide only a modest estimate of the extent to which protein functions may diverge in isolation when recombination is suppressed or following a WGD event when genetic redundancy is at its peak.

Potential retention of functionally related haplogroups. A GO annotation's number of associated genes also appears to increase exponentially with the number of chromosome arms it is expressed on (Fig. 3). This relationship was fit using gamma regression and a log link, the results of which were highly significant ($p < 0.0005$, Supplementary Fig. S7). For a GO annotation's number of related genes to increase in this manner, the genes on a given chromosome arm must be moderately functionally related. This suggests that the retention of genes following large-scale duplication events may operate at the haplogroup level, a trend that is predicted due to the dosage-sensitive nature of protein complexes. Chromosomes enriched with blocks of functionally related genes that are beneficial in increased dosage would show the highest levels of gene retention. Thus, the functional coherence of the Y chromosome could be attributed to a lower content of functionally related haplogroups that were beneficial in increased dosage on the ancestral autosomes.

Biased retention of orthologs on existing chromosomes. The ancestral vertebrate chromosomes displayed substantial differences in gene number, potentially as a result of more significant gene deletion and loss on chromosomes with a smaller number of resulting genes²⁸. This has led to speculation of systematic biases in the deletion of duplicates on a subset of chromosomes following rediploidization, which may have resulted in the chromosome's eventual loss²⁸. We were interested in determining if this systematic bias could be attributed to the gene content of the pre-duplicated chromosomes from which they were derived. The results of our chromosome analysis show human chromosome arms have normally distributed numbers of orthologous genes (Shapiro-Wilk 0.97, $p = 0.34$) and average K_a/K_s values (Shapiro-Wilk 0.98, $p = 0.72$) (Supplementary Table S3 and Fig. S8), and

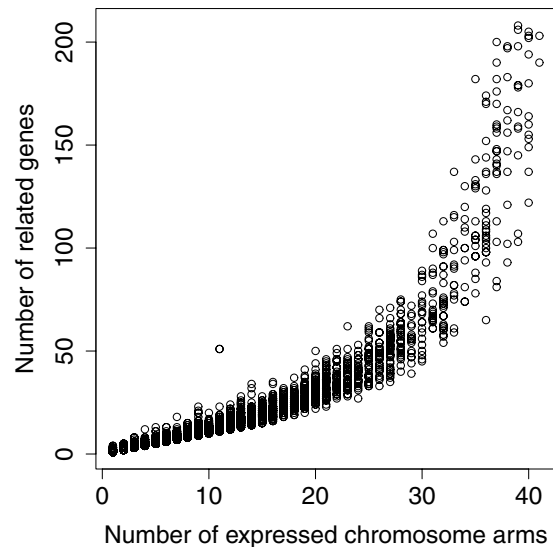


Figure 3. GO Chromosome Arm and Gene Relationship: Number of genes related to each GO term plotted against the number of chromosome arms it was found (trimmed to within 3 SD of the mean number of genes for clarity). Each data point represents a GO annotation ($n = 11,020$). The x-axis corresponds to the number of chromosome arms a given GO annotation was expressed on (the number of chromosome arms containing a gene that was associated with the given GO annotation). The y-axis corresponds to the number of genes a given GO annotation was associated with. The results of a gamma regression analysis indicated the association was statistically significant ($p < 0.0005$). This figure was generated using R statistical software (Version 3.5.3).

that a chromosome arm's number of genes and GO annotations are linearly related ($p = 5.68 \times 10^{-40}$, adjusted $R^2 = 0.985$, Supplementary Figs. 9 and 10). Density of orthologs on existing chromosome arms, however, was found to be non-normally distributed (Shapiro-Wilk 0.912, $p = 0.002$). This is due to biased ancestral gene conservation on a subset of chromosomes. Despite differential average K_a/K_s rates, selection at the chromosome arm level since the divergence of mammals does not appear to influence the number or density of orthologs on a given chromosome (Supplementary Fig. S11).

In contrast, we found that arms of chromosomes that have retained large clusters of genes resulting from the ancestral WGD events contain a disproportionate number of mammalian orthologs in our analysis. Approximately 35% of genes still exist in duplicate copies³⁹, and several paralogs resulting from WGD events^{31,32} (also known as ohnologs) were retained as quadruplicates²⁸. These include clusters containing the four HOX regions on chromosomes 2, 7, 12, and 17, as well as the MHC region on chromosome 6 containing ohnologs on chromosomes 1, 9, and 19 that are a result of single pre-duplicated regions²⁸. The gene content of chromosomes 14 and 15 have also been shown to be almost entirely derived from individual pre-duplicated chromosomes²⁸. The arms of these chromosomes show some of the higher levels of mammalian ortholog retention in our analysis, and chromosomes 17 and 19 have the highest ortholog densities in the genome (Table 1 and Supplementary Fig. S12).

To determine the extent to which these chromosomes remain functionally related, aside from their conserved gene families, we performed hierarchical clustering of the chromosome arms based on a weighted frequency (see methods) for each GO annotation on a given arm in our dataset. The results of our dendrogram (Fig. 4) indicate several trends in the functional relationships between genes of chromosome arms. The two top-level clusters appear to be differentiated based on the number of related functions retained on the chromosome arms (Table 1), a result that was expected given clustering with Euclidean distance. We additionally found lower level clustering of chromosome arms that include both the ancestral HOX and MHC regions. These include the clustering of chromosome 2q, 12q, and 17q, as well as 6pq, 19pq, and 9q. This suggests that the ancestral WGD events have had a profound impact on the retention and organization of mammalian orthologs throughout the human genome.

As stated earlier, specific classes of genes are preferentially retained following WGD events. Chromosomes that have maintained a large portion of their ancestral genes are therefore a result of the gene content and functional annotations of the pre-duplicated chromosomes from which they were derived. It has been hypothesized that the specialization of the Y chromosome is a result of the number of functional genes initially present on the ancestral autosomes⁵⁷, a hypothesis supported by the low functional gene density of the X chromosome^{2,58}. In our present analysis, we also found low ortholog density on both the X chromosome, and chromosomes that are orthologous to the chicken Z chromosome (chromosomes 5, 9, and 18⁵⁸). However, we did not find as significant a disparity in ortholog density between the X chromosome and the human genome-wide average (2.33) relative to overall gene density comparisons⁵⁸, suggesting that ancestral gene density may be less heavily influenced by the invasion of interspersed repeats. Biased gene deletion following the ancestral WGD events resulting in low gene density on a subset of chromosomes suggests that several chromosomes were pre-adapted to specialize similar to the sex chromosomes. Furthermore, chromosomal rearrangements would be under less negative selection in these regions.

	Number of Genes	Number of GO Terms	Density	Avg K_a/K_s
1p	394	3057	3.19	0.14
1q	354	3096	2.82	0.17
2q	295	2710	1.99	0.14
11q	290	2426	3.55	0.17
17q	289	2630	4.97	0.15
5q	272	2401	2.05	0.14
12q	258	2369	2.64	0.12
15q	231	2130	2.78	0.13
7q	229	2352	2.31	0.16
10q	228	2172	2.43	0.13
3q	218	1950	2.03	0.14
6p	200	1682	3.34	0.16
14q	197	1881	2.19	0.12
3p	197	2035	2.17	0.12
2p	193	2015	2.06	0.14
6q	192	1754	1.73	0.15
4q	188	2099	1.34	0.14
9q	178	1857	1.87	0.15
19q	170	1782	5.24	0.15
8q	168	1685	1.68	0.15
19p	154	1362	5.88	0.14
11p	136	1451	2.55	0.13
Xq	128	1322	1.35	0.13
16p	122	1029	3.32	0.14
13q	121	1156	1.25	0.14
16q	120	1296	2.24	0.13
17p	115	1457	4.58	0.12
7p	114	1213	1.90	0.12
22q	112	1399	3.13	0.13
20q	109	1296	3.00	0.12
12p	99	1026	2.79	0.17
8p	89	1085	1.97	0.15
18q	89	1127	1.44	0.14
Xp	80	894	1.31	0.13
4p	74	878	1.48	0.13
10p	67	790	1.68	0.11
5p	64	787	1.31	0.15
9p	62	865	1.44	0.19
20p	58	838	2.06	0.15
21q	49	542	1.41	0.14
18p	29	455	1.57	0.12
Yq	1	10	0.02	0.09
Yp	1	7	0.10	0.16

Table 1. Chromosome arms summary statistics: Includes number of genes with surviving orthologs, number of GO terms, density of genes or orthologs/Mb, and average K_a/K_s of all genes and their species comparisons on a given chromosome arm. Sorted by number of genes.

Discussion

Since its discovery, the perceived functional importance of the Y chromosome has grown exponentially within the scientific community and now may provide further insight into chromosomal evolution following the ancestral WGD events at the origin of the vertebrate lineage. Our present analysis, in conjunction with existing literature, has shown that evolutionary trends believed to be unique to the Y chromosome are observed in the events following large-scale duplications and are still present in mammalian ortholog comparisons. These include higher levels of purifying selection on functionally diverse, ubiquitously expressed genes^{41–45}, as well as reduced purifying selection on genomically isolated protein functions. The biased distribution of ancestral mammalian genes on chromosomes primarily derived from single pre-duplication chromosomes additionally suggests that gene retention was dependent on the gene content of the ancestral chromosome from which they were derived, and this retention may persist over long periods of evolutionary time. The conservation of the functionally coherent,

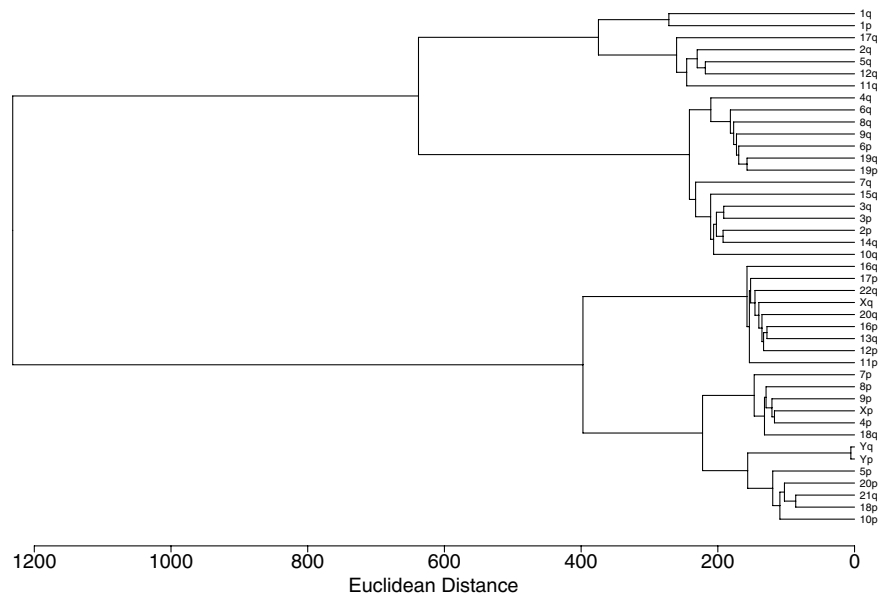


Figure 4. Chromosome Arm Dendrogram: Hierarchical clustering of the chromosome arms based on functional relatedness of their genes. Hierarchical clustering of the chromosome arms based on functional relatedness of their genes. Chromosome arm functionality was obtained by deriving weighted frequencies for each GO annotation within our dataset. Frequencies were weighted based on the specificity of the GO annotation to a given chromosome arm. Distance measure is Euclidean. Multi-node clusters are joined based on minimum increase in within group variance.

potentially haplolethal X-degenerate sequences through purifying selection¹⁶, the rapid evolution of ampliconic sequences expressed primarily in the testes¹⁴, and the pseudogenization and loss of redundant sequences are consistent with a large-scale duplication event. This suggests that the Y chromosome may serve as a model for chromosome evolution following a large-scale duplication event.

Examining suppression of recombination on the Y chromosome in the light of large-scale duplication events has essential implications for karyotypic evolution at the onset of the vertebrate lineage. It has been proposed that the ancestral WGD events contributed to the proliferation of vertebrates during the Cambrian period due to the increase in genetic variation and tolerance to environmental conditions^{38,76}. Recent research suggests that all extant vertebrate karyotypes are descendants of an ancestral marine chordate consisting of 17 chromosome pairs that underwent two successive WGDs²⁸. Rapid loss through the fusion of seven chromosomes between duplications and the loss of an additional five chromosomes following the second duplication resulted in an ancestral Amniota karyotype of 49 chromosomes with highly differential gene content²⁸. The smaller size of extant genomes, therefore, suggests a consistent pattern of karyotype reduction following the ancestral WGD events²⁸, and speciation rates have been shown to be strongly correlated with chromosomal evolution rates⁷⁷.

Duplication events should occur at a fitness cost, and an optimal gene copy number should exist⁷³. Duplicate genes, therefore, would be subjected to three potential evolutionary fates: retention of genes that are beneficial in increased dosage, inactivation of genes that are harmful in increased dosage, and a period of neutral evolution of redundant sequences. Consequently, the observed gene retention on a given chromosome could be attributed to its density of genes that were beneficial in an increased dosage and subsequently retained through purifying selection. Loss of chromosomes due to widespread gene inactivation of detrimental duplicates, however, should have occurred early and ubiquitously, contributing little to the evolutionary novelties and speciation observed at the time. Considering Y chromosome evolution as a potential model for other chromosomes following a large-scale duplication event and high selective pressure would suggest an alternative hypothesis: chromosomal rearrangements may have protected large regions of the genome from gene flow, which allowed isolated genes to diverge until complete reproductive barriers were formed⁷⁸.

The specialization of SRY as the sex-determining factor appears to have played a significant role in X-Y divergence, as its emergence is correlated with the first stratification event that reduced recombination between the neo-sex chromosomes⁷⁹. Single gene sex determination alone should not select for recombination suppression⁸⁰. However, the presence of gonadal dysgenesis in XY individuals with an SRY deletion⁸¹ and sterility of XX individuals containing an inactivated copy of SRY^{82–87} suggests that multiple genes are required to produce fertile offspring. The accumulation of sexually antagonistic alleles in a non-recombining portion of the genome could have provided a sufficient selective advantage that outweighed the deleterious effects of reduced recombination due to their synergistic effects on fertility. Despite Mueller's ratchet being implicated in the early stages of Y chromosome degeneration⁸⁸, genetic decay due to strong positive selection resulting in hitchhiking events is believed to be responsible for its extensive divergence and continued degeneration^{88–91}. This is supported by the stepwise repression of recombination¹ and the correlation of Y-degeneration with levels of female promiscuity in related

species comparisons^{3,14,15,92}. This suggests that strong positive selection on mutually beneficial alleles at linked sites may drive recombination suppression to become selectively favored.

Ohno's original model of genetic evolution suggesting that newly duplicated genes would be functionally redundant and able to escape purifying selection⁹³ has mixed empirical evidence^{54,73}. The events of large-scale duplications, however, create an environment in which newly duplicated genes or complexes that are beneficial in increased dosage may be retained through purifying selection, while the remainder of duplicates would show a continuum of redundancy. The scale of such duplications would allow a small subset of genes to achieve a beneficial mutation. If one mutation resulted in a novel function or further specialized a gene towards one of its respective functions, the likelihood of the gene being retained would increase⁴⁶. Our analysis has also shown that increasing its functional specificity may relax purifying selection, resulting in further divergence. In the event that this new, beneficial mutation occurred on a highly redundant chromosome, the additional reduction of purifying selection due to isolation of a function in an environment with little to no background selection may selectively specialize the chromosome. The survival of the remaining neutrally evolving sequences on that chromosome would depend on their acquisition of functionally related beneficial mutations. If the chromosome bearing this specialized function captured a pair of alleles that together significantly increased the organism's fitness, selection for recombination suppression may result in an inversion becoming prevalent in the population. As observed on the Y chromosome, the resulting chromosome would now contain a complex of genes maintained through purifying selection, as well as a subset of specialized genes that are rapidly evolving resulting in a period of extensive divergence from its homologous counterpart.

The probability that a new inversion captures an advantageous haplotype can be high⁷²; however, for an inversion to become fixed when sexually antagonistic alleles are not present the selective advantage would have to strongly outweigh the negative fitness consequences of reduced recombination⁹⁴. In 1973, Leigh Van Valen showed that the probability of extinction of a population was constant over time and suggested an evolutionary arms race where survival is dependent on a population's ability to adapt to changing selective pressures⁹⁵. During times of intense selective pressure, selection for rapid fixation of a highly advantageous haplotype may have driven recombination suppression to become selectively favored due to the reduced effective population size and increased fixation rate. Recombination suppression events, such as inversions, in the absence of sexual antagonism, would have markedly different evolutionary consequences. This is due in part to inversions only reducing recombination in heterozygotes⁷. If the inversion is driven to fixation, recombination would resume between the new homologous chromosomes. In isolated populations, this divergence from the ancestral chromosome may have been sufficient to create a reproductive barrier, such as in the divergence of ancestral *Equus* populations⁹⁶. As evidenced by the Y chromosome, recombination suppression can also occur progressively⁸⁰ and may be related to continued selection for newly introduced, functionally related, beneficial alleles. Subsequent inversions resulting from extended periods of intense selective pressure on the associated functions would continue to drive the degeneration of the chromosome through successive hitchhiking events, increasing its long-term fragility.

For the Y chromosome, or any significantly degraded chromosome to go extinct, its functions would need to be replaced elsewhere or no longer under selective constraint to prevent fitness consequences²⁴. Relaxation of selection at the locally adapted sites (e.g., predator/prey adaptation), however, would render the genes functionally inert, and selection for recombination resumption between ubiquitously expressed genes would result in fusion events becoming selectively favored. Thus, prolonged strong selection for specialization would have driven a subset of ancestral chromosomes to extinction. As this pertains to the fate of the human Y chromosome, continued selection for localization of sexually antagonistic traits, reduction of female promiscuity resulting in less intensive sexual selection, and its recent stability may suggest it is here to stay. Its continued survival in species still experiencing strong sexual selection, however, may be suspect.

It is worth noting that an apparent contradiction in this logic is the ZW sex-determining chromosomes in avian lineages in which females are the heterogametic sex. Similar to the evolution of the Y chromosome, suppression of recombination in the W chromosome has resulted in significant degeneration of its ancestral gene content⁹⁷. Those that remain functionally active have been shown to be ubiquitously expressed and are believed to be essential to both sexes⁹⁷. In contrast, the W chromosome lacks genes coding for female-advantage traits⁹⁷, suggesting that selection for specialization has not resulted in its degeneration. To reduce the recombination between mutually beneficial, sexually antagonistic alleles, however, an inversion can occur in either chromosome⁸⁰. DMRT1 has also been implicated as the sex-determining locus in the ZW system and is present only on the Z chromosome^{98,99}, suggesting testes development may function through a dosage-dependent mechanism. In conjunction with a lack of dosage-compensation observed in the ZW system¹⁰⁰, the degeneration of the W chromosome is still a result of male-driven positive selection, only on the opposite chromosome.

The chicken Z has been found to be orthologous to portions of human chromosomes 5, 9, and 18 while the human X is orthologous to chicken chromosomes 1 and 4^{58,97}. Due to a lack of structural similarity with their respective orthologous regions, researchers have suggested the chicken Z and human X chromosome were not predisposed to become sex chromosomes and their low gene density is a result of convergent evolution⁵⁸. Chromosomes 5, 9, and 18, however, show similar levels of ortholog density as the X chromosome in our analysis, as well as, some of the most substantial differences in mammalian ortholog content between the arms of individual chromosomes. This phenomenon is most likely a result of the arms being derived from different ancestral chromosomes that underwent fusion events²⁸. Chromosome 9p also contains the ortholog of the believed avian sex-determining locus DMRT1 and is functionally related to the short arm of the X chromosome in our analysis. This may indicate that the convergent evolution of the sex chromosomes is a result of the differential fusion of ancestral chromosomes containing sex-related genes, and the process of sex chromosome evolution further lowered overall gene density to their current state. However, further research needs to be conducted to determine the significance of this relationship.

The results of our analysis in the context of existing literature present a model by which chromosomes, and therefore populations, rapidly evolved at the onset of the vertebrate lineage. The large-scale duplication events allowed a subset of genes to subfunctionalize, thereby reducing pleiotropic constraints and accelerating evolutionary rates. The isolation of these genes on redundant chromosomes further relieved purifying selection, resulting in a period of rapid chromosomal evolution and divergence due to specialization. If this divergence alone did not create a reproductive barrier, the chromosome's eventual loss due to a change in adaptive pressures would have resulted in differential karyotypes of isolated populations. Thus, the extinction of chromosomes due to specialization is not unique to the Y chromosome, or sex chromosomes in general.

Methods

In the present analysis, we examined the divergence of human genes from their mammalian orthologs with respect to their GO annotations and chromosomal locations. For mammalian genes, the probability a newly arisen nonsynonymous mutation is fixed, relative to what is expected under neutrality, is resolved by the strength of selection¹⁰¹. The nonsynonymous to synonymous mutation rate (K_a/K_s ratio) is commonly utilized as a measurement of this selective strength, with low values suggesting strong purifying selection and high values indicating relaxed purifying selection and/or positive selection¹⁰¹. By definition, human orthologs are identical by descent¹⁰² with at least one other species in our analysis. The use of mammalian ortholog comparisons in conjunction with GO annotations, therefore, allowed us to analyze how gene protein functions influence the patterns of selection that led to the differential divergence of ancestral mammalian genes.

Data collection. Divergence data was collected from The Searchable Prototype Experimental Evolutionary Database (SPEED)¹⁰³. SPEED contains orthologous sequence comparisons of nine species including human (*Homo sapiens*), chimp (*Pan troglodytes*), rhesus macaque (*Macaca mulatta*), mouse (*Mus musculus*), rat (*Rattus norvegicus*), dog (*Canis familiaris*), cow (*Bos Taurus*), opossum (*Momodelphis domestica*), and chicken (*Gallus gallus*) as a true outgroup. Methodology on the identification of orthologous groups and calculation of divergence data can be found elsewhere¹⁰³.

Orthologous sequence pairs were queried from SPEED for genetic summary information and their related K_a/K_s values. Data cleaning was performed using PySpark (Spark version 2.3.1, Python version 2.7.10). Sequence comparisons containing a K_s value of zero or less due to computational error were removed from the analysis. Any sequence comparison set with unusually low K_s values were removed, as they gave spuriously high K_a/K_s values. Where multiple comparisons existed, divergence data inconsistencies were resolved by computing a zero-corrected harmonic mean; therefore, more significant weight was given to conservative estimates¹⁰⁴, and comparisons containing at least one zero K_a/K_s value were assigned a K_a/K_s value of zero. Lastly, sequence comparisons that did not include a human comparison with an associated gene name and chromosomal location were excluded. Our resulting dataset included a total of 68,006 comparisons across 10,849 genes.

Gene ontology information was collected from the European Bioinformatics Institute¹⁰⁵. The most recent version of human gene ontology annotations (9/19/19) was downloaded and joined to their respective genes. The dataset included 19,395 genes and 18,211 GO terms. The validity of GO terms with IEA (Inferred from Electronic Annotation) evidence codes has been questioned due to their inferential nature¹⁰⁶. The quality of IEA terms, however, has significantly improved and rival those inferred by curators¹⁰⁷. To alleviate potentially biased numbers of GO annotations on well-studied genes, IEA terms were also retained. After joining with the ortholog dataset and removing genes that lacked annotation, our final dataset included 6,734 annotated genes across 14,121 GO terms. IEA, IDA (inferred from direct assay), ISS (Inferred from Sequence or structural Similarity), IBA (Inferred from Biological aspect of Ancestor), IMP (Inferred from Mutant Phenotype) and TAS (Traceable Author Statement) evidence codes were the primary methods of annotation in our dataset at 30.2%, 20.96%, 13.2%, 12.5%, 10.3%, and 7.3%, respectively.

Data preparation. Single value human gene K_a/K_s rates were derived by averaging their respective K_a/K_s values across all species comparisons present in the dataset. GO annotation K_a/K_s values were obtained by averaging the K_a/K_s values of all related genes across all species comparisons. Chromosome arm K_a/K_s values were calculated by averaging the K_a/K_s values of all genes present on the respective chromosome arm across all species comparisons. Ortholog density was calculated by dividing the number of orthologs present on a given chromosome arm by arm size in Mb. Lastly, the chromosome arms and their related GO annotations were cross-tabulated to obtain the number of times a given function occurs on each arm. Due to their hierarchical nature, GO terms can be broad¹⁰⁶. This issue was addressed in a context dependent manner for each analysis. Prior to clustering the chromosome arms based on functional relatedness of their genes, the number of times each GO annotation occurred on each chromosome arm was weighted using an algorithm adapted from Martinez and Reyes-Valdés¹⁰⁸. We considered the average frequency of the i^{th} GO term among j chromosome arms as,

$$\rho_i = \frac{1}{t} \sum_{j=1}^t \rho_{ij} \quad (1)$$

and defined GO term specificity as the information that its expression provides about the identity of the chromosome arm as

$$S_i = \frac{1}{t} \left(\sum_{j=1}^t \frac{\rho_{ij}}{\rho_i} \log_2 \frac{\rho_{ij}}{\rho_i} \right) \quad (2)$$

S_i will give zero if the GO term is expressed on all chromosome arms and $\max \log_2(t)$ if the function is exclusively expressed on a single chromosome arm. We then assigned a weighted frequency for each GO term on each chromosome arm as the product of the GO term specificity and its frequency on a given chromosome arm.

$$\delta_{ij} = \rho_{ij} S_i \quad (3)$$

Thus, a higher degree of functional similarity would be found between chromosome arms if their shared functions were absent elsewhere in the genome. This method was also applied to the relationship between genes and their related GO terms. Weighted GO term counts were derived for each gene by summing the specificities of their related GO terms in equation (3). However, the weighted GO counts did not alter the distributions or significances of our ortholog regression analyses. Therefore, raw counts were used for ease of interpretability.

A primary goal of our GO annotation regression analyses was to determine if genomic representation influences the selective pressures exerted on a function. Therefore, all ontology terms were retained in this analysis in order to examine our hypothesis that large scale duplication events may relieve pleiotropic constraints in a subset of genes through increased dosage of essential functions. However, where multiple GO terms contained the same set of related genes, only one term was retained to remove redundant data points. 11,020 terms of the original 14,121 were found to have unique sets of related genes. 11,016 of these were found to have non-zero K_a/K_s values and were used in K_a/K_s regression analyses.

Statistical analysis. All regression and distribution analyses were performed in Python (see version above) using the statsmodels API. Due to the strong positive skew of several variables in our dataset, generalized linear models (GLM) were used where appropriate for hypothesis testing. Fitting positively skewed continuous data with a gamma distribution has been shown to perform comparably or outperform lognormal transformations without the need for manual manipulation of the variables¹⁰⁹ and provides a more flexible model when assumptions of ordinary least squares are violated. A log link was used to maintain a non-linear fit with interpretable coefficients while respecting the domain of the gamma function. The exponential of the coefficient for the intercept and predictor variable, therefore, represent the initial predicted outcome value and rate of change for a one-unit increase in the predictor variable, respectively. This methodology was applied to the relationship between an ortholog's number of GO terms and its K_a/K_s ratio, a GO term's number of expressed chromosome arms and its K_a/K_s ratio, a GO term's number of related genes and its K_a/K_s ratio, and a GO term's number of expressed chromosome arms and number of related genes. Where statistically meaningful zero values were present, a hurdle method was employed to counteract the calculation error introduced. This entails fitting a gamma distribution to all non-zero data, as well as a binomial distribution to the full dataset to determine the influence of the predictor variable on the probability that the dependent variable is zero¹¹⁰. This was applied to the relationship between an ortholog's K_a/K_s ratio and number of GO terms. The linear relationship of chromosome arm's number of related genes and GO annotations was fit with ordinary least squares regression without an intercept, as it was nonsensical in the given context. Normality of distributions was determined using the Shapiro-Wilk test which tests the null hypothesis that the data are normally distributed.

Hierarchical clustering of the chromosome arms based on GO annotation content was performed using the cluster package in R (Version 3.5.3). GO term counts were not scaled before distance calculation due to the homogenous nature of the variables. The distance was calculated using Euclidean distance. The linkage measure was determined by obtaining the agglomerative coefficient (amount of clustering structure found) for single, complete, average linkage and Ward's method using the `agnes()` function. For our dataset, Ward's method resulted in the highest agglomerative coefficient and was subsequently used in our clustering analysis. Therefore, multi-node clusters were joined based on minimum increase in within-group variance.

Data availability

Requests for compiled data materials pertaining to this manuscript should be addressed to J.W. Inquiries pertaining to the use of raw data from the SPEED database should be submitted to J.M.S. email address jstaley2@ksu.edu.

Code availability

Requests for code should be addressed to J.W.

Received: 14 December 2019; Accepted: 20 January 2020;

Published online: 07 February 2020

References

1. Lahn, B. T. & Page, D. C. Four evolutionary strata on the human x chromosome. *Science* **286**, 964–967 (1999).
2. Ross, M. T. *et al.* The DNA sequence of the human x chromosome. *Nature* **434**, 325–337 (2005).
3. Hughes, J. F. *et al.* Strict evolutionary conservation followed rapid gene loss on human and rhesus y chromosomes. *Nature* **483**, 82–86 (2012).
4. Skaletsky, H. *et al.* The male-specific region of the human y chromosome is a mosaic of discrete sequence classes. *Nature* **423**, 825–837 (2003).
5. Charlesworth, B. & Charlesworth, D. The degeneration of y chromosomes. *Philosophical Transactions of the Royal Society B: Biological Sciences* **355**, 1563–1572 (2000).

6. Lemaitre, C. *et al.* Footprints of inversions at present and past pseudoautosomal boundaries in human sex chromosomes. *Genome Biology and Evolution* **1**, 56–66 (2009).
7. Kirkpatrick, M. How and why chromosome inversions evolve. *PLoS Biology* **8** (2010).
8. Graves, J. A. M. Sex chromosome specialization and degeneration in mammals. *Cell* **124**, 901–914 (2006).
9. Aitken, R. J. & Graves, J. A. M. Human spermatozoa: The future of sex. *Nature* **415**, 963 (2002).
10. Graves, J. A. The rise and fall of SRY. *Trends in Genetics* **18**, 259–264 (2002).
11. Arakawa, Y., Nishida-Umehara, C., Matsuda, Y., Sutou, S. & Suzuki, H. X-chromosomal localization of mammalian y-linked genes in two XO species of the ryukyu spiny rat. *Cytogenetic and Genome Research* **99**, 303–309 (2002).
12. Just, W. *et al.* Absence of sry in species of the vole *Ellobius*. *Nature Genetics* **11**, 117 (1995).
13. Kuroiwa, A., Ishiguchi, Y., Yamada, F., Shintaro, A. & Matsuda, Y. The process of a y-loss event in an XO/XO mammal, the ryukyu spiny rat. *Chromosoma* **119**, 519–526 (2010).
14. Hughes, J. F. *et al.* Chimpanzee and human y chromosomes are remarkably divergent in structure and gene content. *Nature* **463**, 536–539 (2010).
15. Hughes, J. F. *et al.* Conservation of y-linked genes during human evolution revealed by comparative sequencing in chimpanzee. *Nature* **437**, 100 (2005).
16. Rozen, S., Marszalek, J. D., Alagappan, R. K., Skaletsky, H. & Page, D. C. Remarkably little variation in proteins encoded by the y chromosome's single-copy genes, implying effective purifying selection. *American Journal of Human Genetics* **85**, 923–928 (2009).
17. Bellott, D. W. *et al.* Mammalian y chromosomes retain widely expressed dosage-sensitive regulators. *Nature* **508**, 494–499 (2014).
18. Rozen, S. *et al.* Abundant gene conversion between arms of palindromes in human and ape y chromosomes. *Nature* **423**, 873–876 (2003).
19. Teitz, L. S., Pyntikova, T., Skaletsky, H. & Page, D. C. Selection has countered high mutability to preserve the ancestral copy number of y chromosome amplicons in diverse human lineages. *American Journal of Human Genetics* **103**, 261–275 (2018).
20. Marais, G. A. B., Campos, P. R. A. & Gordo, I. Can intra-y gene conversion oppose the degeneration of the human y chromosome? a simulation study. *Genome Biology and Evolution* **2**, 347–357 (2010).
21. Connallon, T. & Clark, A. G. Gene duplication, gene conversion and the evolution of the y chromosome. *Genetics* **186**, 277–286 (2010).
22. Bachtrog, D. A dynamic view of sex chromosome evolution. *Current Opinion in Genetics & Development* **16**, 578–585 (2006).
23. Kaiser, V. B., Zhou, Q. & Bachtrog, D. Nonrandom gene loss from the drosophila miranda neo-y chromosome. *Genome Biology and Evolution* **3**, 1329–1337 (2011).
24. Bachtrog, D. Y chromosome evolution: emerging insights into processes of y chromosome degeneration. *Nature reviews. Genetics* **14**, 113–124 (2013).
25. Lahn, B. T. & Page, D. C. Functional coherence of the human y chromosome. *Science (New York, N.Y.)* **278**, 675–680 (1997).
26. Page, D. C., Harper, M. E., Love, J. & Botstein, D. Occurrence of a transposition from the x-chromosome long arm to the y-chromosome short arm during human evolution. *Nature* **311**, 119–123 (1984).
27. Connallon, T., Débarre, F. & Li, X.-Y. Linking local adaptation with the evolution of sex differences. *Philosophical Transactions of the Royal Society B: Biological Sciences* **373** (2018).
28. Sacerdot, C., Louis, A., Bon, C., Berthelot, C. & RoestCrollius, H. Chromosome evolution at the origin of the ancestral vertebrate genome. *Genome Biology* **19**, 166 (2018).
29. Dehal, P. & Boore, J. L. Two rounds of whole genome duplication in the ancestral vertebrate. *PLoS Biology* **3** (2005).
30. Putnam, N. H. *et al.* The amphioxus genome and the evolution of the chordate karyotype. *Nature* **453**, 1064–1071 (2008).
31. Nakatani, Y., Takeda, H., Kohara, Y. & Morishita, S. Reconstruction of the vertebrate ancestral genome reveals dynamic genome reorganization in early vertebrates. *Genome Research* **17**, 1254–1265 (2007).
32. Wolfe, K. H. Yesterday's polyploids and the mystery of diploidization. *Nature Reviews Genetics* **2**, 333 (2001).
33. Blomme, T. *et al.* The gain and loss of genes during 600 million years of vertebrate evolution. *Genome Biology* **7**, R43 (2006).
34. Innan, H. & Kondrashov, F. The evolution of gene duplications: classifying and distinguishing between models. *Nature Reviews Genetics* **11**, 97–108 (2010).
35. Lynch, M. & Conery, J. S. The evolutionary fate and consequences of duplicate genes. *Science* **290**, 1151–1155 (2000).
36. Gout, J.-F., Kahn, D., Duret, L. & Consortium, P. P.-G. The relationship among gene expression, the evolution of gene dosage, and the rate of protein evolution. *PLoS Genetics* **6**, 1000944 (2010).
37. Birchler, J. A. & Veitia, R. A. The gene balance hypothesis: From classical genetics to modern genomics. *The Plant Cell* **19**, 395–402 (2007).
38. Van de Peer, Y., Maere, S. & Meyer, A. The evolutionary significance of ancient genome duplications. *Nature Reviews Genetics* **10**, 725–732 (2009).
39. Singh, P. P., Arora, J. & Isambert, H. Identification of ohnolog genes originating from whole genome duplication in early vertebrates, based on synteny comparison across multiple genomes. *PLoS Computational Biology* **11**, 1004394 (2015).
40. Makino, T. & McLysaght, A. Positionally biased gene loss after whole genome duplication: Evidence from human, yeast, and plant. *Genome Research* **22**, 2427–2435 (2012).
41. Tu, Z. *et al.* Further understanding human disease genes by comparing with housekeeping genes and other genes. *BMC Genomics* **7**, 31 (2006).
42. Zhang, L. & Li, W.-H. Mammalian housekeeping genes evolve more slowly than tissue-specific genes. *Molecular Biology and Evolution* **21**, 236–239 (2004).
43. Yang, J., Gu, Z. & Li, W.-H. Rate of protein evolution versus fitness effect of gene deletion. *Molecular Biology and Evolution* **20**, 772–774 (2003).
44. Duret, L. & Mouchiroud, D. Determinants of substitution rates in mammalian genes: Expression pattern affects selection intensity but not mutation rate. *Molecular Biology and Evolution* **17**, 68–070 (2000).
45. Torgerson, D. G. & Singh, R. S. Rapid evolution through gene duplication and subfunctionalization of the testes-specific α 4 proteasome subunits in drosophila. *Genetics* **168**, 1421–1432 (2004).
46. Lynch, M. & Force, A. The probability of duplicate gene preservation by subfunctionalization. *Genetics* **154**, 459–473 (2000).
47. Wagner, A. The fate of duplicated genes: loss or new function? *BioEssays* **20**, 785–788 (1998).
48. Dermitzakis, E. T. & Clark, A. G. Differential selection after duplication in mammalian developmental genes. *Molecular Biology and Evolution* **18**, 557–562 (2001).
49. Li, W.-H., Yang, J. & Gu, X. Expression divergence between duplicate genes. *Trends in Genetics* **21**, 602–607 (2005).
50. Conant, G. C. & Wolfe, K. H. Turning a hobby into a job: How duplicated genes find new functions. *Nature Reviews Genetics* **9**, 938–950 (2008).
51. Rastogi, S. & Liberles, D. A. Subfunctionalization of duplicated genes as a transition state to neofunctionalization. *BMC Evolutionary Biology* **5**, 28 (2005).
52. He, X. & Zhang, J. Rapid subfunctionalization accompanied by prolonged and substantial neofunctionalization in duplicate gene evolution. *Genetics* **169**, 1157–1164 (2005).
53. Wilson, M. A. & Makova, K. D. Evolution and survival on eutherian sex chromosomes. *PLoS Genetics* **5** (2009).
54. Zhang, P., Gu, Z. & Li, W.-H. Different evolutionary patterns between young duplicate genes in the human genome. *Genome Biology* **4**, R56 (2003).

55. Cortez, D. *et al.* Origins and functional evolution of y chromosomes across mammals. *Nature* **508**, 488–493 (2014).
56. Mueller, J. L. *et al.* Independent specialization of the human and mouse x chromosomes for the male germline. *Nature genetics* **45**, 1083–1087 (2013).
57. Saifi, G. M. & Chandra, H. S. An apparent excess of sex- and reproduction-related genes on the human x chromosome. *Proceedings of the Royal Society B: Biological Sciences* **266**, 203–209 (1999).
58. Bellott, D. W. *et al.* Convergent evolution of chicken z and human x chromosomes by expansion and gene acquisition. *Nature* **466**, 612–616 (2010).
59. Blackmon, H. & Brandvain, Y. Long-term fragility of y chromosomes is dominated by short-term resolution of sexual antagonism. *Genetics* **207**, 1621–1629 (2017).
60. Rice, W. R. The accumulation of sexually antagonistic genes as a selective agent promoting the evolution of reduced recombination between primitive sex chromosomes. *Evolution* **41**, 911–914 (1987).
61. Rice, W. R. Evolution of the y sex chromosome in Animals Y chromosomes evolve through the degeneration of autosomes. *BioScience* **46**, 331–343 (1996).
62. Charlesworth, D. & Charlesworth, B. Sex differences in fitness and selection for centric fusions between sex-chromosomes and autosomes. *Genetical Research* **35**, 205–214 (1980).
63. van Doorn, G. S. & Kirkpatrick, M. Turnover of sex chromosomes induced by sexual conflict. *Nature* **449**, 909–912 (2007).
64. Matsumoto, T. & Kitano, J. The intricate relationship between sexually antagonistic selection and the evolution of sex chromosome fusions. *Journal of Theoretical Biology* **404**, 97–108 (2016).
65. Charlesworth, B. The evolution of chromosomal sex determination and dosage compensation. *Current Biology* **6**, 149–162 (1996).
66. Saxena, R. *et al.* The DAZ gene cluster on the human y chromosome arose from an autosomal gene that was transposed, repeatedly amplified and pruned. *Nature Genetics* **14**, 292 (1996).
67. Lahn, B. T. & Page, D. C. Retroposition of autosomal mRNA yielded testis-specific gene family on human y chromosome. *Nature Genetics* **21**, 429–433 (1999).
68. Bhowmick, B. K., Satta, Y. & Takahata, N. The origin and evolution of human ampliconic gene families and ampliconic structure. *Genome Research* **17**, 441–450 (2007).
69. Zhou, Q. & Bachtrog, D. Sex-specific adaptation drives early sex chromosome evolution in drosophila. *Science (New York, N.Y.)* **337**, 341–345 (2012).
70. Wyckoff, G. J., Li, J. & Wu, C.-I. Molecular evolution of functional genes on the mammalian y chromosome. *Molecular Biology and Evolution* **19**, 1633–1636 (2002).
71. Wyckoff, G. J., Wang, W. & Wu, C.-I. Rapid evolution of male reproductive genes in the descent of man. *Nature* **403**, 304 (2000).
72. Kirkpatrick, M. & Barton, N. Chromosome inversions, local adaptation and speciation. *Genetics* **173**, 419–434 (2006).
73. Kondrashov, F. A., Rogozin, I. B., Wolf, Y. I. & Koonin, E. V. Selection in the evolution of gene duplications. *Genome Biology* **3**, research0008.1 (2002).
74. Kondrashov, F. A. & Koonin, E. V. A common framework for understanding the origin of genetic dominance and evolutionary fates of gene duplications. *Trends in Genetics* **20**, 287–290 (2004).
75. McVicker, G., Gordon, D., Davis, C. & Green, P. Widespread genomic signatures of natural selection in hominid evolution. *PLOS Genetics* **5**, 1000471 (2009).
76. Crow, K. D. & Wagner, G. P. What is the role of genome duplication in the evolution of complexity and diversity? *Molecular Biology and Evolution* **23**, 887–892 (2006).
77. Bush, G. L., Case, S. M., Wilson, A. C. & Patton, J. L. Rapid speciation and chromosomal evolution in mammals. *Proceedings of the National Academy of Sciences of the United States of America* **74**, 3942–3946 (1977).
78. Livingstone, K. & Rieseberg, L. Chromosomal evolution and speciation: a recombination-based approach: Research review. *New Phytologist* **161**, 107–112 (2003).
79. Foster, J. W. & Graves, J. A. An SRY-related sequence on the marsupial x chromosome: implications for the evolution of the mammalian testis-determining gene. *Proceedings of the National Academy of Sciences of the United States of America* **91**, 1927–1931 (1994).
80. Charlesworth, D., Charlesworth, B. & Marais, G. Steps in the evolution of heteromorphic sex chromosomes. *Heredity* **95**, 118 (2005).
81. Hughes, J. F. & Rozen, S. Genomics and genetics of human and primate y chromosomes. *Annual Review of Genomics and Human Genetics* **13**, 83–108 (2012).
82. Chapelle, Adl, Tippet, P. A., Wetterstrand, G. & Page, D. Genetic evidence of x-y interchange in a human XX male. *Nature* **307**, 170–171 (1984).
83. Page, D. C., Chapelle, Adl & Weissenbach, J. Chromosome y-specific DNA in related human XX males. *Nature* **315**, 224–226 (1985).
84. Pepene, C. E., Coman, I., Mihu, D., Militaru, M. & Duncea, I. Infertility in a new 46, XX male with positive SRY confirmed by fluorescence in situ hybridization: a case report. *Clinical and Experimental Obstetrics & Gynecology* **35**, 299–300 (2008).
85. Anik, A., Çatlı, G., Abacı, A. & Böber, E. 46,XX male disorder of sexual development: A case report. *Journal of Clinical Research in Pediatric Endocrinology* **5**, 258–260 (2013).
86. Wang, T., Liu, J. H., Yang, J., Chen, J. & Ye, Z. Q. 46, XX male sex reversal syndrome: a case report and review of the genetic basis. *Andrologia* **41**, 59–62 (2009).
87. Wu, Q.-Y. *et al.* Clinical, molecular and cytogenetic analysis of 46, XX testicular disorder of sex development with SRY-positive. *BMC Urology* **14**, 70 (2014).
88. Bachtrog, D. The temporal dynamics of processes underlying y chromosome degeneration. *Genetics* **179**, 1513–1525 (2008).
89. Kuroki, Y. *et al.* Comparative analysis of chimpanzee and human y chromosomes unveils complex evolutionary pathway. *Nature Genetics* **38**, 158 (2006).
90. Bachtrog, D. Evidence that positive selection drives y-chromosome degeneration in drosophila miranda. *Nature Genetics* **36**, 518 (2004).
91. Rice, W. R. Genetic hitchhiking and the evolution of reduced genetic activity of the y sex chromosome. *Genetics* **116**, 161–167 (1987).
92. Dorus, S., Evans, P. D., Wyckoff, G. J., Choi, S. S. & Lahn, B. T. Rate of molecular evolution of the seminal protein gene SEMG2 correlates with levels of female promiscuity. *Nature Genetics* **36**, 1326–1329 (2004).
93. Ohno, S. *Evolution by Gene Duplication* (Springer-Verlag, 1970).
94. Charlesworth, B., Betancourt, A. J., Kaiser, V. B. & Gordo, I. Genetic recombination and molecular evolution. *Cold Spring Harbor Symposia on Quantitative Biology* **74**, 177–186 (2009).
95. L., V. A. N. V. A. L. E. N. A new evolutionary law. *Evol Theory* **1**, 1–30 (1973).
96. Renaud, G. *et al.* Improved de novo genomic assembly for the domestic donkey. *Science Advances* **4**, eaaq0392 (2018).
97. Bellott, D. W. *et al.* Avian w and mammalian y chromosomes convergently retained dosage-sensitive regulators. *Nature Genetics* **49**, 387–394 (2017).
98. Smith, C. A. *et al.* The avian z-linked gene DMRT1 is required for male sex determination in the chicken. *Nature* **461**, 267–271 (2009).
99. MarshallGraves, J. A. Sex determination: Birds do it with a z gene. *Nature* **461**, 177–178 (2009).

100. Vicoso, B. & Bachtrog, D. Progress and prospects toward our understanding of the evolution of dosage compensation. *Chromosome Research* **17** (2009).
101. Wyckoff, G. J., Malcom, C. M., Vallender, E. J. & Lahn, B. T. A highly unexpected strong correlation between fixation probability of nonsynonymous mutations and mutation rate. *Trends in Genetics* **21**, 381–385 (2005).
102. Koonin, E. V. Orthologs, paralogs, and evolutionary genomics. *Annual Review of Genetics* **39**, 309–338 (2005).
103. Vallender, E. J., Paschall, J. E., Malcom, C. M., Lahn, B. T. & Wyckoff, G. J. SPEED: a molecular-evolution-based database of mammalian orthologous groups. *Bioinformatics (Oxford, England)* **22**, 2835–2837 (2006).
104. Ferger, W. F. The nature and use of the harmonic mean. *Journal of the American Statistical Association* **26**, 36–40 (1931).
105. Huntley, R. P. *et al.* The GOA database: gene ontology annotation updates for 2015. *Nucleic Acids Research* **43**, D1057–1063 (2015).
106. duPlessis, L., Škunca, N. & Dessimoz, C. The what, where, how and why of gene ontology—a primer for bioinformaticians. *Briefings in Bioinformatics* **12**, 723–735 (2011).
107. Škunca, N., Altenhoff, A. & Dessimoz, C. Quality of computationally inferred gene ontology annotations. *PLOS Computational Biology* **8**, e1002533 (2012).
108. Martínez, O. & Reyes-Valdés, M. H. Defining diversity, specialization, and gene specificity in transcriptomes through information theory. *Proceedings of the National Academy of Sciences* **105**, 9709–9714 (2008).
109. Gustavsson, S., Fagerberg, B., Sallsten, G. & Andersson, E. M. Regression models for log-normal data: Comparing different methods for quantifying the association between abdominal adiposity and biomarkers of inflammation and insulin resistance. *International Journal of Environmental Research and Public Health* **11**, 3521–3539 (2014).
110. Tong, E. N. C., Mues, C. & Thomas, L. A zero-adjusted gamma model for mortgage loan loss given default. *International Journal of Forecasting* **29**, 548–562 (2013).

Acknowledgements

The authors would like to thank Christine Malcom and Neil Miller for help in preparing and reviewing this manuscript. This work was supported in part by Bionexus KC.

Author contributions

J.W., J.M.S. and G.J.W. planned the project. J.W. and J.M.S. performed data collection. J.W. performed data and statistical analysis. J.W. generated all figures. J.W. and G.J.W. wrote the paper. All authors reviewed the manuscript.

Competing interests

The authors declare no competing financial interests.

Additional information

Supplementary information is available for this paper at <https://doi.org/10.1038/s41598-020-58997-2>.

Correspondence and requests for materials should be addressed to J.W.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2020