

Differential Retention of Pfam Domains Contributes to Long-term Evolutionary Trends

Jennifer E. James ^{1,2}, Paul G. Nelson ^{†,1} and Joanna Masel ^{*,1}

¹Department of Ecology & Evolutionary Biology, University of Arizona, Tucson, AZ

²Department of Ecology and Genetics, Uppsala University, Uppsala, Sweden

[†]Present address: Empress Therapeutics, Cambridge, MA.

*Corresponding author: E-mail: masel@arizona.edu.

Associate editor: Fabia Ursula Battistuzzi

Abstract

Protein domains that emerged more recently in evolution have a higher structural disorder and greater clustering of hydrophobic residues along the primary sequence. It is hard to explain how selection acting via descent with modification could act so slowly as not to saturate over the extraordinarily long timescales over which these trends persist. Here, we hypothesize that the trends were created by a higher level of selection that differentially affects the retention probabilities of protein domains with different properties. This hypothesis predicts that loss rates should depend on disorder and clustering trait values. To test this, we inferred loss rates via maximum likelihood for animal Pfam domains, after first performing a set of stringent quality control methods to reduce annotation errors. Intermediate trait values, matching those of ancient domains, are associated with the lowest loss rates, making our results difficult to explain with reference to previously described homology detection biases. Simulations confirm that effect sizes are of the right magnitude to produce the observed long-term trends. Our results support the hypothesis that differential domain loss slowly weeds out those protein domains that have nonoptimal levels of disorder and clustering. The same preferences also shape the differential diversification of Pfam domains, thereby further impacting proteome composition.

Key words: clade selection, protein evolution, Cope's rule, intrinsic structural disorder, protein folding, phylostratigraphy, gene families.

Introduction

The possibility that evolution could result in long-term directional change over time has enduring appeal, but there are few well-documented examples. It has proven difficult to rigorously assess purported universal tendencies toward increases in complexity and size over time (McShea 1991; Gregory 2008). Perhaps, the best-documented example of a long-term trend is that some taxonomic groups show a tendency toward increasing body size over time, aka "Cope's rule" (Cope 1885; Payne et al. 2009; Heim et al. 2015). Interestingly, this trend does not seem to be caused by directional selection slowly exerting a preference for larger individuals but rather by the differential diversification of larger body size clades (Heim et al. 2015). This is unsurprising, because directional selection among individuals of the same species works quickly, making it an unlikely explanation for such a slow directional trend.

New examples of long-term trends were recently observed in the properties of protein domains (e.g., intrinsic structural disorder [ISD]) as a function of how much time has elapsed since their birth, with considerable work performed to exclude artifactual explanations (Foy et al. 2019; James et al. 2021). We find it more plausible that some bias in the evolutionary process is responsible for

these trends, rather than that the external environments shaping birth have moved in such a consistent direction over such long timescales. One hypothesis is that, regardless of the time of birth, all domains were born with properties biased toward promoting de novo gene birth (Wilson et al. 2017), and since then have had different amounts of time to evolve away from that starting point toward properties that are more optimal for correct folding (Bucciantini et al. 2002; Debès et al. 2013; Chiti and Dobson 2017). This hypothesis is supported by the fact that the amino acid frequencies characteristic of newborn animal domains also make the expression of a random peptide more benign in *Escherichia coli* (Kosinski et al. 2022). But it is very hard to explain why directional selection on amino acid frequencies has been so slow to take full effect. We know that the evolution of amino acid frequencies can be rapid, as it is able to keep pace with the rapid evolution of nucleotide composition (Brbić et al. 2015).

Here, we consider the possibility that differential retention versus loss of protein domains with different properties might be responsible for long-term trends as a function of age. We hypothesize that domains are born with a wide range of biochemical properties. From this initial diversity, "fitter" domains, that is, those that duplicate

© The Author(s) 2023. Published by Oxford University Press on behalf of Society for Molecular Biology and Evolution.

This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

Open Access

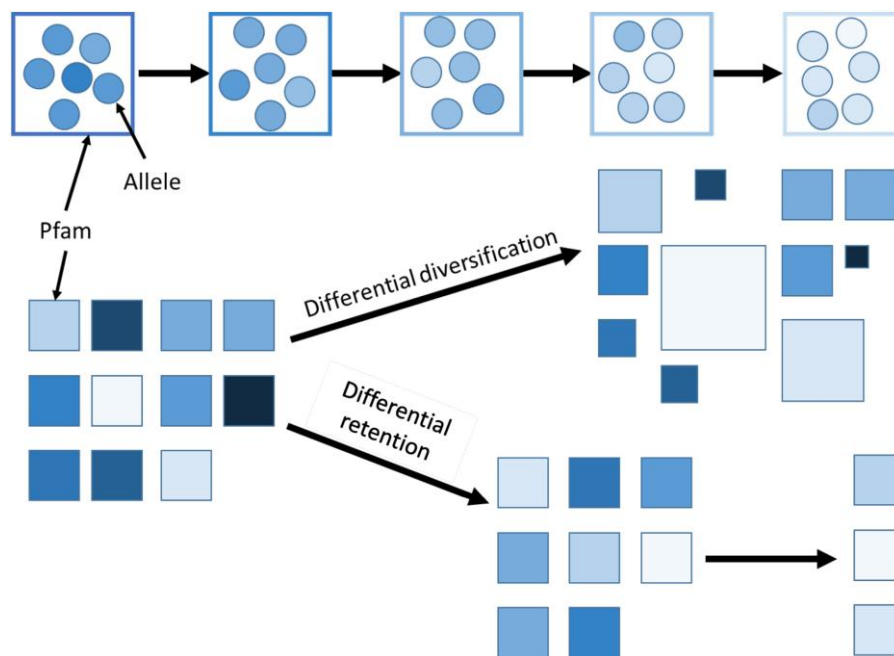


Fig. 1. Alternative mechanisms/levels of directional selection over evolutionary time. Directional evolution in a trait with cohort age, represented by a change from dark to pale, can occur either via (top) descent with modification during which selected (pale) alleles replace other alleles, or (bottom) differential retention in which an initial diversity of Pfams is differentially retained. Differential diversification (middle) combines tendencies for duplication and tendencies for loss. Note that this is a level of selection situation; a Pfam is a clade of protein-coding sequences and their component alleles, and the same figure could be interpreted with individual interests in lieu of allele interests and clade interests in lieu of Pfam interests.

more often and/or are lost less often, are more likely to have descendants in extant organisms and thus to appear in the older domain age classes.

Unlike other hypotheses, this causal hypothesis predicts that a protein property that shows a long-term trend with age will also be an indicator of loss rate. One observed trend, extending all the way back to the last universal common ancestor, is that younger domains tend to have “clustered” rather than interspersed hydrophobic acids along their primary sequence (Foy et al. 2019; James et al. 2021). A second trend, restricted to the evolutionary history of animals, is that younger domains tend to have a higher level of ISD, primarily due to trends in the frequencies of the 20 amino acids (James et al. 2021). Here, we test whether and how ISD and hydrophobic clustering are related to the loss rate of a domain.

What we are proposing is that the relevant level of selection may not be the one that selects among alternative alleles of the same gene (fig. 1, top), but instead the level that selects among alternative domains capable of fulfilling the same cellular functions as genes duplicate, differentiate, and are lost (fig. 1, bottom). Rapid loss of recently born protein-coding genes has previously been postulated (Tautz and Domazet-Lošo 2011; Palmieri et al. 2014; Schlötterer 2015). The number of domain copies also appears to be surprisingly changeable, with fairly extreme expansions and contractions observed even over short evolutionary timescales (Hahn et al. 2007; Moore and Bornberg-Bauer 2012), leaving plenty of scope for differential diversification and retention to be a biologically significant influence.

We focus on protein domains rather than on complete genes. We take domain annotations from the Pfam database (El-gebali et al. 2019) and refer to domains simply as Pfams. Protein-coding domains are sometimes thought of as functional units of the protein sequence that are able to fold independently (Ponting and Russell 2002). However, Pfams are annotated by HMMER3 on the basis of evolutionary rather than functional relatedness (Finn et al. 2011) and as such are fundamental units of protein sequence homology. This makes them easier to work with than genes, which commonly contain multiple different Pfams of a variety of ages (Bagowski et al. 2010; Bornberg-Bauer and Albà 2013).

We infer loss rates by maximizing the likelihood of the observed presence/absence of at least one Pfam instance across all species in a phylogeny. As well as this central focus on total loss rates, leading to differential retention over time (fig. 1, bottom), we perform a complementary analysis of differential diversification (fig. 1, middle) by taking the mean number of distinct instances of Pfams within a species. Both presence/absence data and number-of-instances data are vulnerable to both false positives and false negatives during homology detection. Our previous age assignment procedure (James et al. 2021) used a keyword-based approach to remove Pfams that were likely incorrectly annotated in the focal genome due to contamination (Salzberg 2017; Lu and Salzberg 2018; Breitwieser et al. 2019), a problem that if uncorrected will result in nonsensical loss rates for affected Pfams. James et al. (2021) also removed Pfams that exhibited an implausible phylogenetic

distribution, which could have resulted from either contamination or horizontal gene transfer. Here, we go further, reducing the rate of false negatives by checking all six reading frames of contigs for unannotated Pfam instances and then reconsidering whether the phylogenetic distribution of the Pfam is implausible. Secondly, we infer total loss rates jointly with a rate of false positive hits. Finally, we consider the hypothesis of homology detection bias while interpreting results.

We ask whether and how ISD and clustering scores are related to the rate of total loss of a Pfam in a lineage and also to the number of Pfam instances per species. We analyze animal lineages only, because the ISD trends as a function of Pfam age found by [James et al. \(2021\)](#) were specific to animals.

Results

We used the phylostratigraphy data set of [James et al. \(2021\)](#), restricted to the 6,841 Pfams annotated in at least 2 out of our 343 animal species and already subjected to a number of filters used by [James et al. \(2021\)](#). We scanned all six reading frames of intergenic regions to find species that contained unannotated instances of these Pfams. We also performed additional quality controls to remove likely viral or other contaminants on the basis of incoherent distribution across the tree (see Materials and Methods), resulting in a data set of 6,700 Pfams. We then estimated the rate of total loss of each Pfam over the animal species in our data set (see Materials and Methods). During this maximum likelihood (ML) procedure, we iteratively removed Pfam presence data for clades (mostly single species) where they are more likely to be false positive hits or horizontal gene transfer events.

The median inferred loss rate of a Pfam in the animal phylogeny was 0.0009/My (first and third quartiles of 0.00021 and 0.0035), that is 0.9 losses per 100,000 years. Note that this is not the rate at which a single instance of a Pfam is lost, but rather the rate at which all instances of a Pfam are lost from a lineage. Unsurprisingly, Pfams with a higher mean number of instances per genome tend to have a lower rate of total loss ([supplementary fig. S1, Supplementary Material](#) online, adjusted $R^2 = 0.30$, $P < 1e-16$).

Pfams With a Mean ISD of 0.18 Are Least Often Lost

Pfams with high ISD have high loss rates (adjusted $R^2 = 0.023$, slope = 1.35, $P = 4e-36$; data and loess regression shown in [fig. 2A](#)). This is the predicted direction in order for their relationship to explain the animal phylostratigraphy trend ([fig. 3A](#)). Note that the Box–Cox transformation of loss rates makes the units of the slope in these linear models hard to interpret. Interestingly, this relationship is not just nonlinear, but also non-monotonic (loess regression in [fig. 2A](#)). A regression model that includes a breakpoint at ISD = 0.18 (95% confidence interval [CI]: 0.15–0.21), representing the minimum loss rate, explains substantially more of the variance in loss rate among Pfams (adjusted $R^2 = 0.030$, slopes -2.00 and 3.86), and fits the data significantly better than

a model without a breakpoint ($P = 6e-12$), better than a model with zero slope for ISD values higher than the breakpoint ($P = 1e-45$) and better than a model with zero slope for ISD values lower than the breakpoint ($P = 0.0002$).

The existence of a breakpoint rather than a monotonic relationship contradicts predictions from the homology detection bias hypothesis. That hypothesis expects an overestimation of loss rates for high ISD Pfams whose homologs are more difficult to detect. While this could drive the overall relationship, we found in which Pfams with high ISD have high loss rates, it cannot explain why loss rates also rise with low ISD < 0.18 .

Our results are instead consistent with selection acting to preferentially remove Pfams at either end of the spectrum of ISD. In further support of this, there is less variability in ISD within older Pfam cohorts ([fig. 4A](#), slope = -4.1 to 5 ISD/My, adjusted $R^2 = 0.88$, $P = 3e-22$, using linear regression weighted by the number of Pfams per age cohort) than within younger Pfam cohorts. Ancient Pfams (those born over 2,000 Ma) that are still around today have a mean ISD of 0.21 (median 0.20), at the upper edge of the 95% CI for the breakpoint representing minimum loss rate. These two observations further support the hypothesis that differential retention drives the long-term trend in ISD as a function of evolutionary age observed in [James et al. \(2021\)](#) and [Foy et al. \(2019\)](#). The fact that ancient Pfams may have ISD a little above the breakpoint might be noise in the data or it might indicate that the trend of falling ISD with age (shown in [fig. 4A](#)) has not yet fully saturated.

Pfams With Moderately Interspersed Hydrophobic Residues Are Least Often Lost

The hydrophobic residues of younger domains are more clustered (i.e., have higher clustering) than those of old domains ([James et al. 2021](#)). This directional trend could be explained most simply if Pfams with higher clustering were differentially lost. However, as with ISD, our data strongly support a non-monotonic relationship between clustering and loss rate ([fig. 2B](#), losses/My is two-parameter Box–Cox transformed prior to analyses). The best-fitting linear model for clustering incorporated a breakpoint with a clustering value of 0.81 (95% CI: 0.77–0.85) corresponding to the lowest loss rate. This model has an adjusted R^2 of 0.020 and fits the data significantly better than a model without a breakpoint ($P = 1e-25$), than a model with zero slope for the lower clustering values ($P = 8e-28$), and than a model with zero slope for the higher clustering values ($P = 2e-29$). Clustering is not known to affect homology detection, but, even if it did, it is again not clear how homology detection bias could produce a non-monotonic relationship.

A clustering value of 1 corresponds to a random distribution of hydrophobic residues, while the optimal clustering value of 0.81 represents the interspersion of hydrophobic residues. Similar to the ISD results, this breakpoint value corresponds closely to the average clustering score of the most ancient Pfams (over 2,000 My old), which

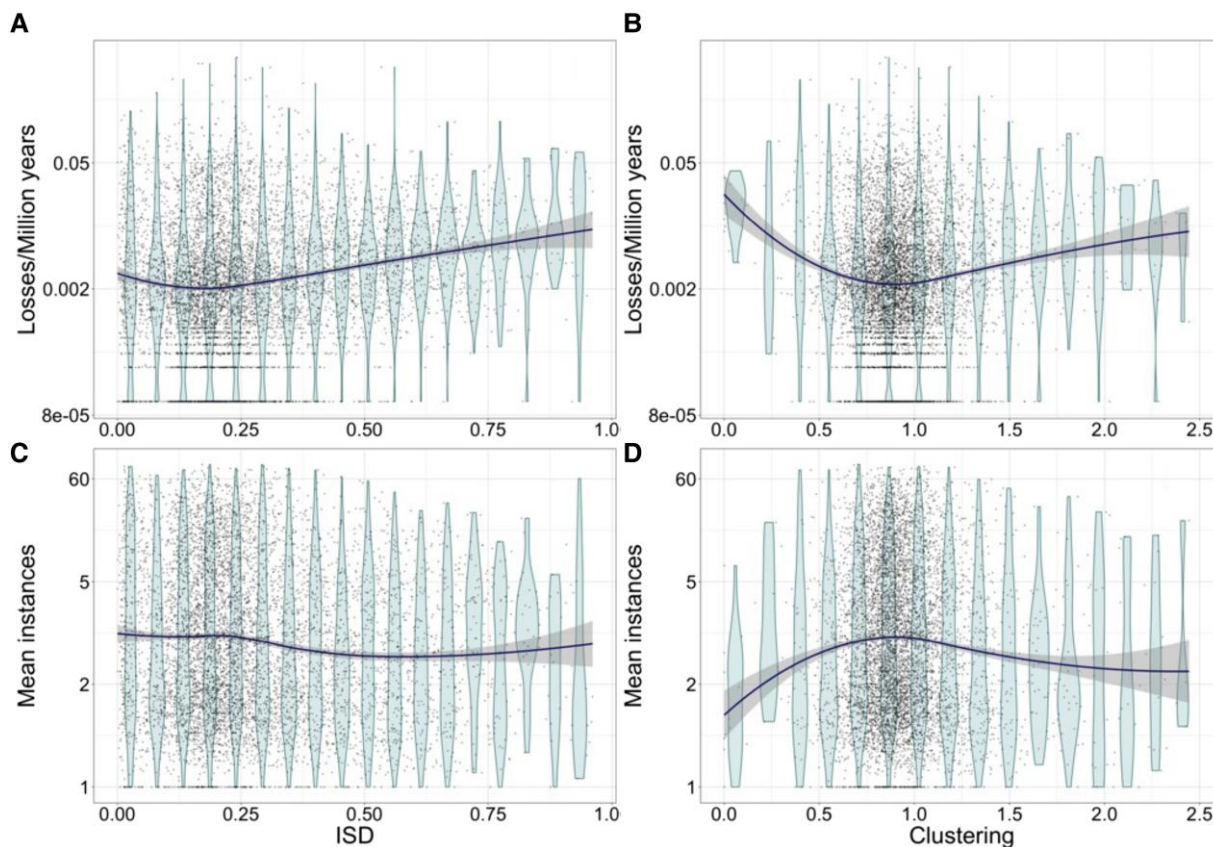


Fig. 2. Pfam loss rates and copy number depend, often non-monotonically, on ISD and clustering values. The lowest loss rates are seen for Pfams with a mean ISD of 0.18 (A) and mean clustering value of 0.85 (B). Pfams with a low ISD have a higher number of instances per genome, in a more monotonic relationship (C). Pfams with a mean clustering value of 0.85 have the highest number of instances per genome. Each point represents a Pfam. Violin plots are a visual guide only, based on dividing the data into groups of equal range of ISD (A, C) or clustering (B, D). Loess regression curves are shown in dark, with 95% CIs shown as shading. In (A) and (B), horizontal bands of points near the bottom represent 0, 1, 2, etc. observed loss events, while in (C) and (D) horizontal bands represent Pfams present in only a single copy per genome. In (B) and (D), for better visualization, the x-axis has been truncated at 2.5.

have mean and median clustering scores of 0.89 and 0.88, respectively. The fact that this is higher than the upper end of our 95% CI for the breakpoint suggests that the trend in clustering values with age has not yet saturated. The y-intercept of [figure 3B](#) shows that Pfams are born with clustering scores a little above 1. Because Pfams tend to be born with relatively random amino acid placement corresponding to a clustering score of approximately 1 that is so much higher than the optimum value of 0.81, differential Pfam retention can explain the progressive drop in clustering scores with Pfam age. In further support of the hypothesis that differential loss is responsible for the trend in clustering with age, there is more variability in clustering within younger Pfam cohorts ([fig. 4B](#), slope of standard deviation in clustering with age = -5.5×10^{-5} , adjusted $R^2 = 0.60$, $P = 2e-10$, using linear regression weighted by number of Pfams per age cohort).

Results Are Robust to Pfam Age as a Confounding Factor

When ascertaining the causal effects of ISD and clustering on loss rates, Pfam age has the potential to be a

confounding factor in our analyses, for two reasons. Firstly, if older Pfams also have lower ISD and less clustered hydrophobic amino acids ([James et al. 2021](#)) for reasons other than differential loss due to these properties, this could potentially drive a spurious relationship between Pfam properties and loss rates. As expected, we observe lower loss rates for older Pfams in our data (adjusted $R^2 = 0.050$, $P = 2e-76$, [supplementary fig. S2, Supplementary Material](#) online). Secondly and more subtly, Pfams of identical age share the same speciation events, which enable Pfam losses to be observed. This results in phylogenetic structure in the data. A further complication is that the fact that older Pfams have survived to be observed might create ascertainment bias toward estimated loss rates that are lower than the true loss propensity.

To control both for any confounding properties of Pfams and for phylogenetic structure, we used mixed-effect linear regression models, including age as a random rather than a quantitative effect. Note that each branch of the tree corresponds to a possible Pfam age, making age a discrete quantitative variable. Empirically, we found that random effects models are better-supported model choices: for example, for models of loss rate and ISD, $R^2 = 0.33$ if we include age

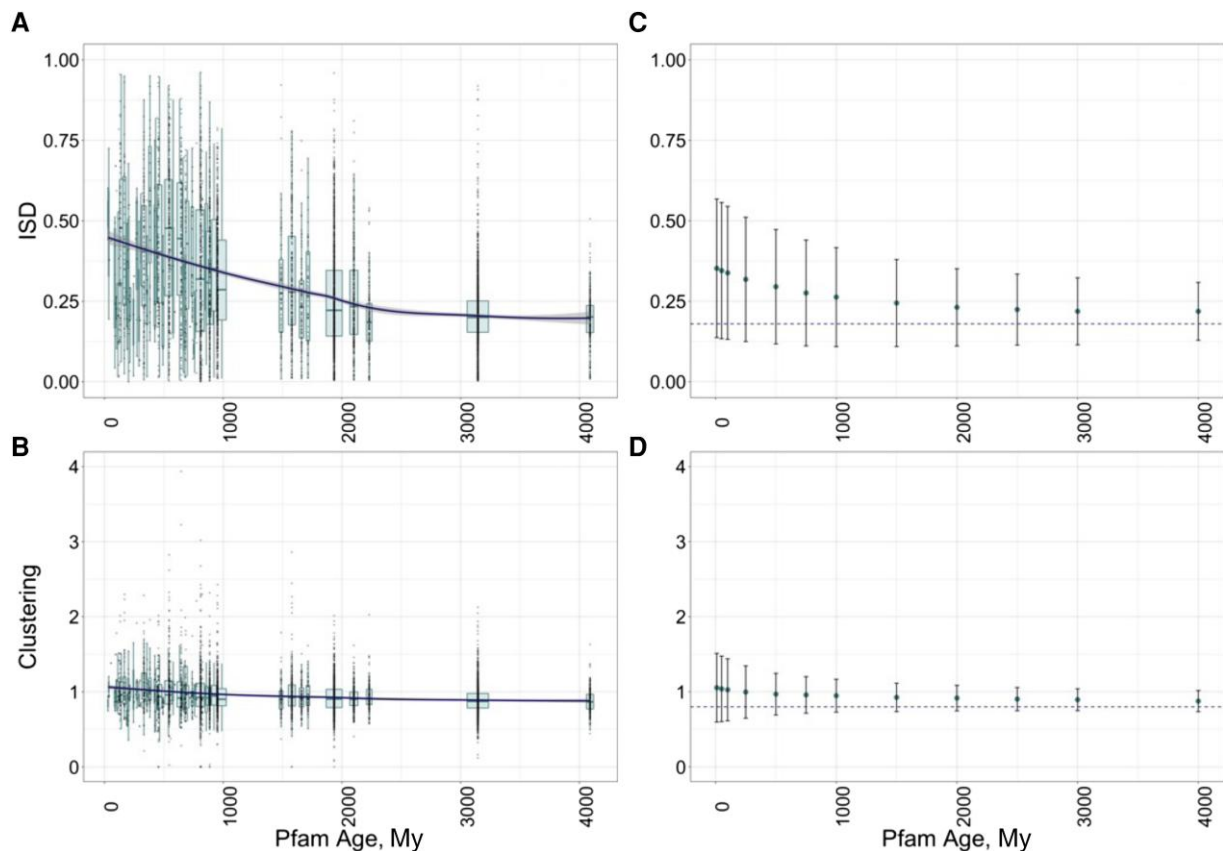


Fig. 3. Observed phylostratigraphy trends in ISD (A) and Clustering (B) values occur on similar timescales as in our simulations of differential loss rates (C and D). In (A) and (B), each box represents all Pfams of a particular age category, with the widths of the boxes proportional to the square root of the number of Pfams in the group. Boxplots show the median, upper, and lower quartiles of the data, while the whiskers represent the largest and smallest data value above and below the interquartile range multiplied by 1.5. Grey points show the values for individual Pfams. The dark line is the loess regression slope of the effect of age on ISD (A) and clustering (B), with the 95% CI shown as shading. In (C) and (D), we plot the mean of ISD (C) and clustering (D), with error bars showing the standard deviation. Timepoints all use the same simulations and so are not independent, that is, the number of domains in our simulations decreases with time. Dashed lines indicate the values of ISD and clustering corresponding to the lowest loss rates. Loss rates for the simulations are those inferred from piecewise linear regression from [figure 2A and B](#).

as a random effect, and $R^2 = 0.058$ if we include age as a linear predictor, where R^2 here is the conditional R^2 describing the variance explained by the entire model. Akaike information criterion (AIC) scores for the respective models are 23,580 and 24,432. Similarly, for models of loss rate and clustering, a model that includes age as a random effect model has $R^2 = 0.34$, while a fixed effect model has $R^2 = 0.05$. AIC values are 23,603 and 24,476, respectively.

The non-monotonic relationship between loss rate and ISD remains statistically supported in the random effect model ($P = 0.0003$), with the model having a marginal R^2 , which represents the variance explained in loss rate just by the fixed effect of ISD, of 0.0030, and a slope of -1.17 for values of ISD below 0.18, and 0.64 for values above. The non-monotonic relationship with clustering also remained supported in the random effect model ($P = 0.02$), with a marginal R^2 of 0.00054, with a slope of -0.24 for values of clustering below 0.81 and 0.097 for values of clustering above 0.81. While controlling for Pfaam age causes marginal R^2 values to drop dramatically, the causal direction is unclear. That is differences in ISD and clustering might drive changes in loss rates as a function

of Pfaam age, in which case controlling for age would remove a good portion of the genuine biological effect. However, the statistical significance we find indicates that Pfaam age alone is insufficient to be the sole driver of the dependence of loss rates on ISD and clustering.

Differences in Loss Rates Are Small Enough to Produce Slow Change

Even without controlling for phylogenetic structure/age, the effect sizes for loss and clustering on ISD are clearly small. In most contexts, this would cast doubt on their biological relevance. However, our motivating question is to discover what mechanism might have a small enough effect size such that its steady action could continue for billions of years before saturating. We, therefore, next use simulations to ask whether the effects of ISD and clustering on loss are of approximately the right magnitude to explain the relationship between these variables and Pfaam age, as observed in [James et al. \(2021\)](#).

We simulate a set of Pfams initialized with the distribution of ISD or clustering values observed in young animal

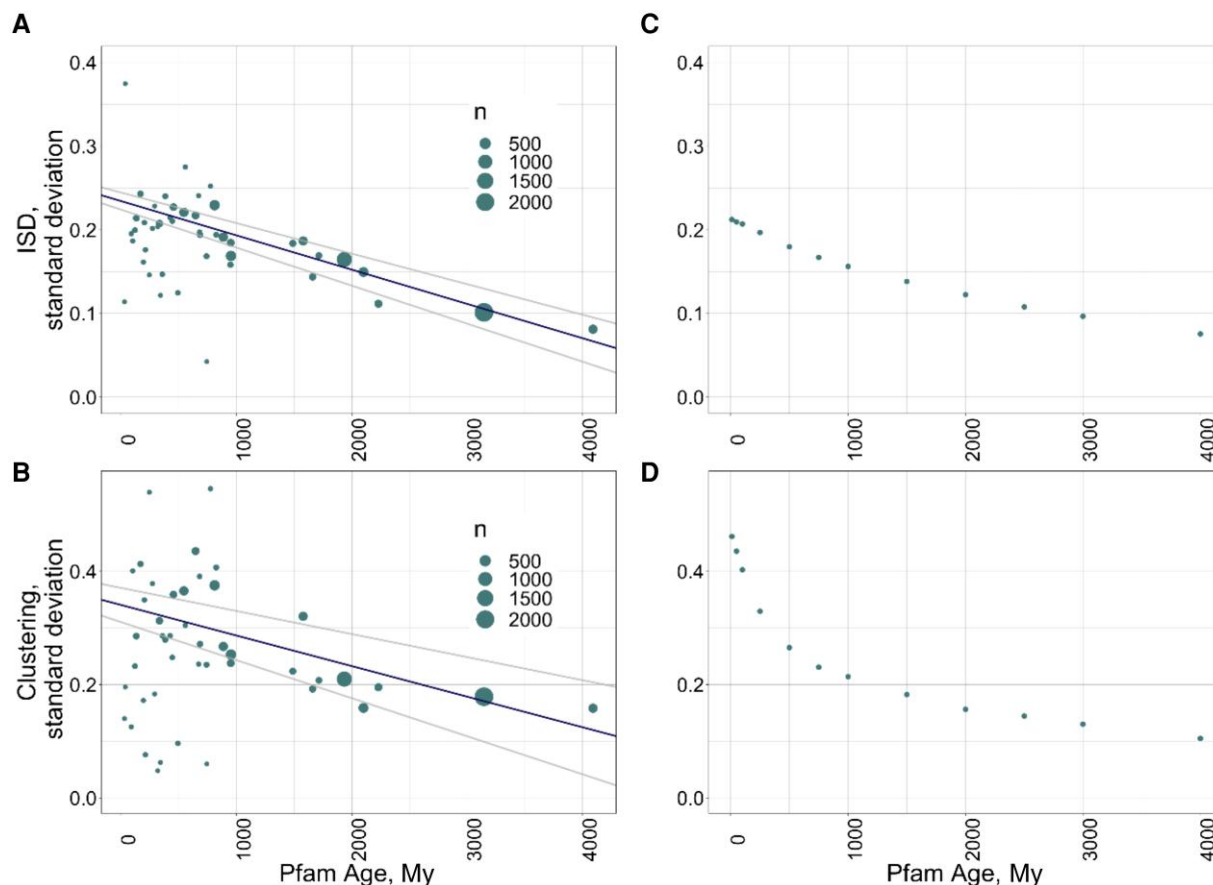


Fig. 4. Younger Pfam cohorts have higher variance within their ISD (A) and clustering (B) distributions, a trend whose timescale is captured in our simulations (C, D). (A) and (B) Point size is scaled by n = number of Pfams in that category. Lines show the weighted linear regressions (weighted by number of observations per category), and their 95% CIs. In (A), for ISD, weighted $R^2 = 0.60$, $P = 2e-10$. In (B), for clustering, weighted $R^2 = 0.60$, $P = 2e-10$. In (C) and (D), the standard deviation of ISD (C) and clustering (D) values are plotted for cohorts of simulated Pfams.

Pfams and subject them to differential loss at rates given by our fitted regression models with breakpoints. This produces directional trends in ISD and clustering on the pertinent timescale (fig. 3C and D). The relationship between ISD (fig. 3A) and clustering (fig. 3B) is visually extremely similar to that generated in our simulations (fig. 3C and D). Our simple single-branch single loss event simulation ignores many complexities of the actual phylogenetic tree but is sufficient to show that the effect sizes of ISD and clustering on loss rates are of approximately the right magnitude to explain phylostratigraphy trends.

ISD and Clustering Are Related to the Mean Number of Pfam Instances

The mean number of instances of the Pfam per genome is an alternative metric of evolutionary success. Not only is it likely an indirect proxy of total loss rate (supplementary fig. S1, Supplementary Material online), but it also shapes proteome composition in its own right, independently of its effect on total loss rate (fig. 1, middle vs. bottom). The addition of substantially novel protein features through domain duplication

followed by significant divergence (Nasir et al. 2014) is subsumed within our metric of instance number.

Unfortunately, the mean number of instances per genome is expected to depend on genome annotation quality, which is likely to vary across the animal species in our data set. We have at least partly addressed this problem by restricting analyses to genomes deemed “Complete.”

In agreement with our loss rate results, Pfams with a higher mean number of instances per genome have lower levels of ISD (fig. 2C, $R^2 = 0.0069$, $P = 7e-11$). Pfams with a higher mean number of instances per genome also have low levels of hydrophobic clustering (adjusted $R^2 = 0.0013$, $P = 0.0017$). A nonlinear model for the effect of ISD on mean number instances is only marginally supported over a linear model ($P = 0.049$ in model comparison with linear relationship, breakpoint = 0.57, 95% CI: 0.42, -0.73), unlike the clear case for a nonlinear relationship with loss rates. However, for hydrophobic clustering, we again find support for a nonlinear relationship (fig. 2D, adjusted $R^2 = 0.0099$), with a highly supported breakpoint at the same value of clustering (0.85, 95% CI: 0.79–0.90). A model with a clustering breakpoint fits the data significantly better than one without ($P = 1e-13$).

Discussion

The persistence of trends in biochemical properties of protein domains (ISD and hydrophobic clustering) over such strikingly long timespans has been an unresolved puzzle. We hypothesized that at birth, Pfam domains have a wide range of properties and are then differentially retained versus lost as a function of these properties. Our empirical results provide a striking confirmation of the predictions of this hypothesis. Specifically, optimal retention is achieved with somewhat low ISD values close to 0.18, and with interspersions of hydrophobic amino acids corresponding to a clustering value of 0.8. The oldest cohort of Pfams has mean ISD and clustering close to these values, and there is less variance within old Pfam cohorts than within young Pfam cohorts. Simulations show that the effect sizes of differential retention as a function of ISD and clustering are of the right order of magnitude (i.e., small enough) to explain how phylostratigraphy trends can persist for such long evolutionary times without saturating. Conditional on retention, similar values of ISD and clustering are associated with a higher mean number of domain instances per genome, magnifying the impact of these biases on proteome-wide composition.

This study implements high-throughput phylostratigraphy; while we implemented a number of novel measures to reduce the frequency of false positive and false negative homologs, our data set will not be free of error. We focus on Pfam domains, for which homology detection is performed using the more sensitive program HMMER rather than BLASTP (Finn et al. 2011; El-gebali et al. 2019). Concerns about homology detection bias usually focus on its potential for causing some rapidly evolving Pfams to be annotated as too young (Moyers and Zhang 2017; Weisman et al. 2020). However, Pfam age is not a central focus of our analysis. The potential issue here is rather that homology detection bias might explain why Pfams with high ISD appear to be lost more often and have fewer instances per genome (a similar bias for clustering has not been documented). However, the known difficulty in detecting homology given higher ISD cannot explain our observed non-monotonic dependence of loss rates and instance number on ISD.

To detect as many true homologs as possible, instead of relying on existing genome annotations, we scan all six reading frames for hits. Each Pfam was as a result included in a mean of 7.35 and a median of 4 new animal species, out of 343. However, HMMER can produce false positives (Pearson et al. 2017); while the use of curated Pfam seeds might reduce this problem, it will not completely eliminate it. We minimized impact by co-inferring false positive (or horizontally transferred) hits, with a higher rate for those found in our 6-frame scan, at the time of our inference of loss rate. Overall, these data quality measures resulted in a substantial decrease in the estimated mean loss rate of approximately 14%. Another key quality control was to remove some Pfams altogether as potential contaminants not from the genomes under study. Our data set

and methodology represent a new standard for work on protein domain evolution, which, while imperfect, has resulted in a higher-quality data set for further analysis. Nevertheless, given the residual likely presence of false positive and false negative Pfam hits, we expect our results on the total loss of a Pfam to be more robust to genome annotation issues than the number of Pfam instances per genome. However, we note that our instance number and loss rate results qualitatively agree.

While protein domains no doubt experience idiosyncratic selection pressures associated with their specific structures and functions within organisms, our findings suggest a more general statistical tendency for Pfams with ISD and clustering values close to the “optimum” to be differentially retained in the long term. A candidate cause for this tendency stems from the fact that the same biophysical properties that promote correct folding also promote toxic aggregation, posing a trade-off dilemma for proteins in general (Bucciantini et al. 2002; Chiti and Dobson 2017). Foy et al. (2019) proposed that young genes tend to address this dilemma via a more “primitive” strategy, with high ISD representing a sacrifice of folding but safety from misfolding. The slow evolution of hydrophobic interspersions was proposed to eventually provide an alternative method for avoiding toxic aggregation that allows for tighter protein folds (Foy et al. 2019). Our findings suggest that the emergence of a tight protein fold might instead depend on good luck in the properties of the newborn protein sequence from which it derives. This idea is supported by lattice protein simulations designed to capture frustration between folding and aggregation, in which there is an abundance of lower fitness peaks and where sequences born with high hydrophobicity (the lattice protein equivalent of low ISD) are more likely to find a high fitness peak (Bertram and Masel 2020).

The current study provides evidence in support of the predictions made by the differential retention hypothesis. These predictions are not made by the hypothesis of directionality in descent with modification, but nor is this hypothesis ruled out by our findings. Biased descent with modification might be operating simultaneously, in either the same or the opposite direction to the forces documented here. While the evolution of amino acid frequencies can be rapid enough to keep pace with changes in nucleotide composition (Brbić et al. 2015), epistatic effects might lead to significant deceleration beyond a certain point (Vieira-Silva and Rocha 2008). Vertebrate species with higher effective population size tend to evolve higher ISD (Weibel et al. 2020), which is the opposite direction to what would be needed to produce the phylostratigraphy trends under the assumption that most ancestral lineages had high population size. While ISD was not directly assessed, long-term directional trends in amino acid frequencies via descent with modification have been inferred during early evolution (Groussin et al. 2013), as well as during the more recent evolution of cold tolerance (Fontanillas et al. 2017; Lecocq et al. 2021).

However, while selection among alternative alleles of the same domain (i.e., descent with modification) is often the default explanation for molecular evolution trends (fig. 1, top), we provide evidence for selection acting on a higher level, among protein domains, highlighting the importance of differential retention and diversification of domains as an evolutionary force in protein sequence evolution (fig. 1, middle and bottom). Selection acting on multiple different levels to affect a trait does not necessarily result in conflict between the levels (Lewontin 1970). This higher level is analogous to the dependence of speciation and extinction rates on the traits of species (see, e.g., Aguilée et al. 2018). An interesting parallel to our findings is the work of Heim et al. (2015), who found that long-term trends do not seem to be caused by directional selection, but rather by the differential diversification of certain clades. Here, we find that the same broad category of explanation applies to long-term trends in protein properties as a function of age since they originated.

The study of long-term directional change in evolution, and the level of selection at which it occurs, has historically focused on species-level traits such as morphology. But there are several advantages of conducting such studies on proteins. Firstly, because Pfam domains are easily identified using HMMER, they are highly tractable for use in the study of higher levels of selection. In particular, it is easier to count how many instances there are of a Pfam in a genome than it is to identify all species in a clade.

Secondly, it is impossible to distinguish between, for example, high extinction rates versus low speciation rates from observations of extant species distributions alone (Louca and Pennell 2020). An inherent difficulty is that we cannot observe extinct species once they are gone; fossil data are required to provide additional information (Mongiardino Koch et al. 2021). However, when a domain undergoes complete loss in one lineage, we can directly infer that loss through comparison to other lineages in which it was not lost. Sophisticated methods developed to study speciation and extinction rates (Fitzjohn 2012) could be adapted to take advantage of this during the study of domain duplication and loss rates. Future studies could examine not just how loss rates depend on trait values but also account for systematic variation in duplication and loss rates among clades, or over time. The study of protein domain diversification and loss is interesting not only to advance our understanding of the evolution of the proteome but also as a powerful new study system for those interested more broadly in the causes of long-term directional trends in evolution, and in levels of selection.

Materials and Methods

Pfam Selection

We begin with the Pfam phylostratigraphy data set of James et al. (2021). Briefly, for a species to be included in this data set, a “Complete” annotated genome had to be available from the National Center for Biotechnology

Information (NCBI), and the species had to be present in Timetree (Hedges et al. 2006). We further restricted this data set to those Pfams that were present in a minimum of 2 animal species, and for which both ISD and clustering scores could be calculated, resulting in a data set of 6,841 Pfams, and their distributions over a phylogenetic tree of 343 animal species. These Pfams had previously been subject to two quality control filtering steps: a keyword-based screening designed to exclude possible contaminants and an additional z-score screen to exclude Pfams with species distributions indistinguishable from chance. Genes and Pfams annotated as mitochondrial were also already excluded from our starting database.

Pfam ages were initially taken from James et al. (2021), based on a species tree of 435 eukaryote species, with additional resolution among the two most ancient Pfam cohorts taken from Weiss et al. (2016). While we calculate loss rates only within the animal species tree, our analysis includes Pfams that originated prior to animals and that are, therefore, also found in other taxa.

False Negatives

To detect potentially missing Pfam instances in unannotated genes, including those missed due to frameshift sequencing errors, we translated all regions not annotated as coding sequences (CDS) in all six frames. For this purpose, we downloaded annotated genomes of the 343 animal species in our data set from NCBI, with access dates by species varying between May and July 2019. We removed sequences annotated as “CDS” and any sequences of “N” more than 3 nucleotides long. We then translated the resulting nucleic acid sequences in all six reading frames, not requiring a start codon for initiation but ending each polypeptide at stop codons, similar to the procedure implemented by Deutekom et al. (2019), and ran amino acid sequences through InterProScan version 5.33–72.0. We record all novel Pfam hits for each species. This resulted in each Pfam being included in a mean of 7.35 and a median of 4 new animal species. Although this six-frame search procedure generates false positives at an elevated rate relative to annotated Pfams, these false positive rates remain modest (see below). To avoid inflation of mean instances/genome by pseudogenes, we only correct putative false negatives in cases where the species in question would otherwise have zero instances of the Pfam.

False Positives

Excluding Contaminant Pfams

As in James et al. (2021), we exclude any Pfam that was both present in fewer than half the species in the original data set (i.e., not restricted to animals) and had a species distribution that was indistinguishable from chance. Such Pfams are likely to be contaminants. Briefly, for each number of species n , we performed 20,000 simulations in which we placed Pfams in n randomly selected species. We then inferred the number of total losses by parsimony and used the mean and variance for the

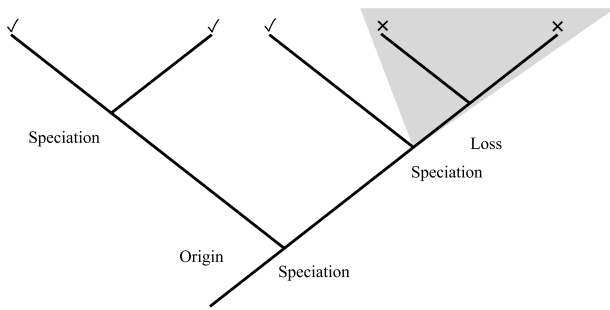


Fig. 5. Inferring loss events for a Pfam by Dollo parsimony. Dollo parsimony assumes a Pfam arose exactly once, on the branch basal to the monophyletic clade that includes all species with the Pfam, and is lost on the branches that explain the current distribution with the fewest number of losses. Species with at least one instance of the Pfam are indicated with checks and those without by crosses. We initialize the Pfam loss rate as $1/\text{sum of unshaded branch lengths}$, that is, all unshaded branches are counted by index k , whereas the most basal-shaded branch is counted by index j in equation (1). We then refine this coarse Dollo parsimony-based estimate through ML techniques that allow for false positive/horizontally transferred instances (eq. 2).

resulting distribution of simulated numbers of total losses to calculate z scores for each empirically observed distribution of n Pfam instances. To remove Pfams with species distributions indistinguishable from chance, James et al. (2021) excluded Pfams with z scores < -2 . Here, we used an additional screen where we excluded any Pfams that had z scores < -2 either with or without the intergenic hits flagged as possible false negatives as described above. This resulted in the further removal of 141 Pfams, leaving a data set of 6,700 Pfams.

Excluding False Positive Pfam Presence While Estimating Loss Rates by Maximum Likelihood

We jointly estimate the total loss rate of each Pfam by ML, together with the identity of false positive or horizontally transferred hits/instances. First consider the case of no false positive or horizontally transferred hits. Using branch lengths obtained from Timetree, we assume that loss events are independent, making the probability of loss of Pfam i with loss rate λ_i per My on branch j of length t_j equal to $1 - \exp(-t_j\lambda_i)$. The likelihood that Pfam i is completely lost on branches j and retained in at least one copy on branches k , given the total loss rate λ_i , is:

$$L_{\lambda}(\lambda_i) = \prod_k \exp(-t_k\lambda_i) \prod_j (1 - \exp(-t_j\lambda_i)) \quad (1)$$

We first provisionally assign losses by assuming a Pfam arose on the branch basal to the monophyletic clade that includes all species with the Pfam and is lost on the branches that explain the current distribution with the fewest number of losses, that is we use Dollo parsimony (see fig. 5). We then use Newton's method to find the loss rate that maximizes the log likelihood, initialized at the loss rate calculated as the number of losses/total branch length of the Pfam tree.

Because multiple loss events can occur on the same branch, the estimated ML loss rate estimate is generally very close but not identical to the initialization value (parsimony mean and median: 0.0058, 0.00088, ML mean and median: 0.0067, 0.00094). Occasionally, false positives/horizontal transfer events that are not yet accounted for trigger runaway inference of absurdly high ML loss rates at this stage. When the ML loss rate > 1 per My, we, therefore, reset the loss rate to the initialization value originally estimated under parsimony, with associated likelihood, prior to proceeding.

Next, to remove false positive/horizontally transferred Pfam presence data, we use a recursive empirical Bayesian procedure. We calculated two prior probabilities: f_c , that the annotated presence of a Pfam in a species or clade is a false positive, and f_i , that a Pfam presence discovered only by our six reading frame scan of intergenic sequences is a false positive. These were initialized at 10^{-6} and 0.2, respectively.

If P_C is the set of species annotated as having the Pfam from CDS observations, and P_i is the set of species annotated as having the Pfam only by us from the InterProScan intergenic search, the likelihood of observing the data if the species sets F_C and F_i are false positives is:

$$L = L_{\lambda}(\lambda_i) f_c^{|F_C|} (1 - f_c)^{|P_C - F_C|} f_i^{|F_i|} (1 - f_i)^{|P_i - F_i|} \quad (2)$$

where λ_i is re-estimated by ML for each species set F_C and F_i , and absolute magnitude notation indicates the number of species in the set.

To find the sets F_C and F_i that maximize the likelihood, we first group Pfam presence observations into the largest possible monophyletic clades where all species putatively possess at least one copy of the Pfam of interest. We then begin by “flipping” each such monophyletic clade, that is designating all species in a clade as being false positive hits or horizontal transfer events, and determining whether the likelihood improves. The false positive designation that improves the likelihood the most is kept, and then the process is repeated. In these subsequent iterations, we also test whether flipping a false positive back to a true positive improves likelihood. The process ends when no single flip improves likelihood. After applying this process to all Pfams, we use the resulting set of false positive exclusion to estimate the two false positive probabilities. Using this as a revised empirical prior, we then repeat the procedure until no further change in the false positive rates is observed. Our method rapidly converges within four iterations, to a false positive rate of 0.00021 for Pfams previously annotated as coding and 0.031 for Pfams annotated here by us based on the presence in putatively intergenic regions. We calculate the mean number of Pfam instances per animal species within the set of species annotated as true positives at the end of the procedure described above.

Our procedures to include presence data when Pfams were previously unannotated, and to exclude presence data as described above, had a noticeable quantitative impact on our data. The combined effect of these quality control measures caused the mean and median loss rate

to fall from 0.0075 and 0.00091 to 0.0065 and 0.00090. The mean and median number of Pfam instances per species rose from 5.9 and 2.2 to 6.2 and 2.4.

For some young Pfams (appearing after divergence with plants), our procedures to reduce false negative and false positive errors in presence annotation within the animal species tree led us to re-estimate Pfam age.

Intrinsic Structural Disorder

We calculated predicted ISD using IUPred 2 (Dosztányi et al. 2005; Mészáros et al. 2018), which scores each amino acid between 0 and 1, with a low score indicating low disorder. In order to obtain a single estimate for each specific Pfam instance, we averaged over all amino acids. Per Pfam scores were then calculated by finding the mean over all copies of the Pfam in our data set (i.e., animal species with more Pfam instances had higher weight).

Clustering

Briefly, we calculated clustering as a kind of normalized index of dispersion: the variance in the number of the most hydrophobic amino acids (Leu, Ile, Val, Phe, Met, and Trp) among blocks of six consecutive amino acids compared to the mean number, normalized so as to be comparable across Pfams of different lengths and hydrophobicities. If the length of a Pfam was not a multiple of 6, we took the average of all possible amino acid frames of 6, truncating the ends appropriately. For more detail, see Foy et al. (2019) and James et al. (2021). A clustering value of 1 indicates the random distribution of hydrophobic residues, with higher values indicating that hydrophobic residues tend to be arranged in clusters along the sequence and lower values indicating that hydrophobic residues are more evenly distributed along the sequence than would be expected by chance. As for ISD, per Pfam scores were then calculated as the mean over all instances of each Pfam in our animal data set.

Statistics

All plotting and statistical analyses were performed in R. Linear regression models that incorporate breakpoints were implemented using the segmented package. We used the packages lme4, MuMIn, and segmented to conduct models, and ggplot2 to create plots. Prior to linear regression, we transformed the loss rate using an optimized two-parameter Box–Cox transform ($\lambda = 0.0054$, $\lambda_2 = 6.1e-06$) calculated using the R package geoR. We transformed the mean instance number using a Box–Cox transform ($\lambda = -0.61$). We transform after rather than before taking the mean because the natural units for scoring impact on the proteome (fig. 1 middle) are linear. Histograms comparing the distribution of untransformed versus transformed loss rates, in addition to Q–Q plots illustrating the requirement for transformation in our linear regression models with ISD, are shown in supplementary figure S3, Supplementary Material online.

Throughout this work, we report the adjusted R^2 values for our linear regression model results, as implemented in

the programming language R, which uses the Wherry formula (Wherry 1931). This R^2 accounts for the number of explanatory variables included in the model. Nested models were compared using the anova function in base R. We conducted mixed effects models in R using the lmer package and calculated corresponding pseudo R^2 values using the R package MumIn. This package returns marginal R^2 values, which can be interpreted as the variance explained by fixed effects in the model, and conditional R^2 values, which can be interpreted as the variance explained by both the fixed and random effects in the model (Nakagawa and Schielzeth 2013).

Simulations

We performed simulations to assess the biological relevance of our inferred loss rates. We initialize these simulations by sampling, with replacement, the Pfams that we estimate to be younger than 100 My old, to produce an initial set of 5,000 Pfams, each represented by its ISD and clustering values. We then assign each Pfam the loss rate predicted from either its ISD value, or its clustering value, on the basis of our corresponding regression model. For each Pfam, we then sample the time to a single loss from an exponential distribution, with scale parameter = 1/loss rate. We then track the mean and standard deviation of the population of survivors over time.

Supplementary Material

Supplementary data in the form of three supplementary figures are available at *Molecular Biology and Evolution* online. Pfam loss rate and age data, R scripts required to replicate data analysis, and python simulation code are all available at <https://github.com/j-e-james/DomainLossRate>.

Acknowledgments

We thank Cecile Ané for helpful discussions that prompted a more complex project that we abandoned in favor of this simpler approach. We thank Nathan Aviles for helpful discussions about fancier statistical approaches that we did not use. We also thank all of the participants at the 2019 “Workshop on Long-Term Trends in Evolution” for stimulating discussions, and Alan Love for discussions about predictivism. Both this work and the workshop were primarily supported by the John Templeton Foundation (60814). This work was also supported by the National Institutes of Health (GM-104040), and in its final stages, by the John Templeton Foundation (62220) and a Wenner-Gren post-doctoral stipendium.

References

Aguilée R, Gascuel F, Lambert A, Ferriere R. 2018. Clade diversification dynamics and the biotic and abiotic controls of speciation and extinction rates. *Nat Commun.* 9(1):1–13.

- Bagowski CP, Bruins W, Velthuis AJ. 2010. The nature of protein domain evolution: shaping the interaction network. *Curr Genomics*. **11**:368–376.
- Bertram J, Masel J. 2020. Evolution rapidly optimizes stability and aggregation in lattice proteins despite pervasive landscape valleys and mazes. *Genetics* **214**:1047–1057.
- Bornberg-Bauer E, Albà MM. 2013. Dynamics and adaptive benefits of modular protein evolution. *Curr Opin Struct Biol*. **23**(3):459–466.
- Brbic M, Warnecke T, Kriško A, Supek F. 2015. Global shifts in genome and proteome composition are very tightly coupled. *Genome Biol Evol*. **7**(6):1519–1532.
- Breitwieser FP, Perlea M, Zimin AV, Salzberg SL. 2019. Human contamination in bacterial genomes has created thousands of spurious proteins. *Genome Res*. **29**(6):954–960.
- Bucciantini M, Giannoni E, Chiti F, Baroni F, Taddei N, Ramponi G, Dobson CM, Stefani M. 2002. Inherent toxicity of aggregates implies a common mechanism for protein misfolding diseases. *Nature* **416**(6880):507–511.
- Chiti F, Dobson CM. 2017. Protein misfolding, amyloid formation, and human disease: a summary of progress over the last decade. *Annu Rev Biochem*. **86**(1):27–68.
- Cope ED. 1885. On the evolution of the vertebrata, progressive and retrogressive. *Am Nat*. **19**(2):140–148.
- Debès C, Wang M, Caetano-Anollés G, Gräter F. 2013. Evolutionary optimization of protein folding. *PLoS Comput Biol*. **9**(1):e1002861.
- Deutecom ES, Vosseberg J, Van Dam TJP, Snel B. 2019. Measuring the impact of gene prediction on gene loss estimates in eukaryotes by quantifying falsely inferred absences. *PLoS Comput Biol*. **15**(8):1–15.
- Dosztányi Z, Csizmók V, Tompa P, Simon I. 2005. The pairwise energy content estimated from amino acid composition discriminates between folded and intrinsically unstructured proteins. *J Mol Biol*. **347**(4):827–839.
- El-Gebali S, Mistry J, Bateman A, Eddy SR, Potter SC, Qureshi M, Richardson LJ, Salazar GA, Smart A, Sonhammer ELL, et al. 2019. The Pfam protein families database in 2019. *Nucleic Acids Res*. **47**(D1):D427–D432.
- Finn RD, Clements J, Eddy SR. 2011. HMMER web server: interactive sequence similarity searching. *Nucleic Acids Res*. **39**(Web Server issue):29–37.
- FitzJohn RG. 2012. Diversitree: comparative phylogenetic analyses of diversification in R. *Methods Ecol Evol*. **3**(6):1084–1092.
- Fontanillas E, Galzitskaya OV, Lecompte O, Lobanov MY, Tanguy A, Mary J, Girguis PR, Hourdez S, Jollivet D. 2017. Proteome evolution of deep-sea hydrothermal vent alvinellid polychaetes supports the ancestry of thermophily and subsequent adaptation to cold in some lineages. *Genome Biol Evol*. **9**(2):279–296.
- Foy SG, Wilson BA, Bertram J, Cordes MHJ, Masel J. 2019. A shift in aggregation avoidance strategy marks a long-term direction to protein evolution. *Genetics* **211**:1345–1355.
- Gregory TR. 2008. Evolutionary trends. *Evol: Educ Outreach*. **1**(3):259–273.
- Groussin M, Boussau B, Gouy M. 2013. A branch-heterogeneous model of protein evolution for efficient inference of ancestral sequences. *Syst Biol*. **62**(4):523–538.
- Hahn MW, Han MV, Han SG. 2007. Gene family evolution across 12 *Drosophila* genomes. *PLoS Genet*. **3**(11):2135–2146.
- Hedges SB, Dudley J, Kumar S. 2006. Timetree: a public knowledgebase of divergence times among organisms. *Bioinformatics* **22**(23):2971–2972.
- Heim NA, Knope ML, Schaal EK, Wang SC, Payne JL. 2015. Cope's rule in the evolution of marine animals. *Science* **347**(6224):867–870.
- James JE, Willis SM, Nelson PG, Weibel C, Kosinski LJ, Masel J. 2021. Universal and taxon-specific trends in protein sequences as a function of age. *ELife* **10**:1–23.
- Kosinski LJ, Aviles NR, Gomez K, Masel J. 2022. Random peptides rich in small and disorder-promoting amino acids are less likely to be harmful. *Genome Biol Evol*. **14**(6):1–14.
- Lecocq M, Groussin M, Gouy M, Brochier-Armanet C. 2021. The molecular determinants of thermoadaptation: Methanococcales as a case study. *Mol Biol Evol*. **38**(5):1761–1776.
- Lewontin RC. 1970. Units of selection. *Annu Rev Ecol Syst*. **1**:1–18.
- Louca S, Pennell MW. 2020. Extant timetrees are consistent with a myriad of diversification histories. *Nature* **580**(7804):502–505.
- Lu J, Salzberg SL. 2018. Removing contaminants from databases of draft genomes. *PLoS Comput Biol*. **14**(6):1–13.
- McShea DW. 1991. Complexity and evolution: what everybody knows. *Biol Philos*. **6**(3):303–324.
- Mészáros B, Erdős G, Dosztányi Z. 2018. IUPred2a: context-dependent prediction of protein disorder as a function of redox state and protein binding. *Nucleic Acids Res*. **46**(Web Server issue):W329–W337.
- Mongiardino Koch N, Garwood RJ, Parry LA. 2021. Fossils improve phylogenetic analyses of morphological characters. *Proc R Soc Biol Sci*. **288**(1950):20210044.
- Moore AD, Bornberg-Bauer E. 2012. The dynamics and evolutionary potential of domain loss and emergence. *Mol Biol Evol*. **29**(2):787–796.
- Moyers BA, Zhang J. 2017. Further simulations and analyses demonstrate open problems of phylostratigraphy. *Genome Biol Evol*. **9**(6):1519–1527.
- Nakagawa S, Schielzeth H. 2013. A general and simple method for obtaining R^2 from generalized linear mixed-effects models. *Methods Ecol Evol*. **4**(2):133–142.
- Nasir A, Kim KM, Caetano-Anollés G. 2014. Global patterns of protein domain gain and loss in superkingdoms. *PLoS Comput Biol*. **10**(1):e1003452.
- Palmieri N, Kosiol C, Schlötterer C. 2014. The life cycle of *Drosophila* orphan genes. *ELife* **3**:1–21.
- Payne JL, Boyer AG, Brown JH, Finnegan S, Kowalewski M, Krause RA, Lyons SK, McLain CR, Novack-Gottshall PM, Smith FA, et al. 2009. Two-phase increase in the maximum size of life over 3.5 billion years reflects biological innovation and environmental opportunity. *Proc Natl Acad Sci U S A*. **106**(1):24–27.
- Pearson WR, Li W, Lopez R. 2017. Query-seeded iterative sequence similarity searching improves selectivity 5–20-fold. *Nucleic Acids Res*. **45**(7):e46.
- Ponting CP, Russell RR. 2002. The natural history of protein domains. *Annu Rev Biophys Biomol Struct*. **31**(1):45–71.
- Salzberg SL. 2017. Horizontal gene transfer is not a hallmark of the human genome. *Genome Biol*. **18**(1):1–5.
- Schlötterer C. 2015. Genes from scratch—the evolutionary fate of de novo genes. *Trends Genet*. **31**(4):215–219.
- Tautz D, Domazet-Lošo T. 2011. The evolutionary origin of orphan genes. *Nat Rev Genet*. **12**(10):692–702.
- Vieira-Silva S, Rocha EPC. 2008. An assessment of the impacts of molecular oxygen on the evolution of proteomes. *Mol Biol Evol*. **25**(9):1931–1942.
- Weibel CA, Wheeler AL, James JE, Willis SM, Masel J. 2023. A new codon adaptation metric predicts vertebrate body size and tendency to protein disorder. *BioRxiv*:2023.03.02.530449. <https://doi.org/10.1101/2023.03.02.530449>
- Weisman CM, Murray AW, Eddy SR. 2020. Many, but not all, lineage-specific genes can be explained by homology detection failure. *PLoS Biol*. **18**(11):1–24.
- Weiss MC, Sousa FL, Mrnjavac N, Neukirchen S, Roettger M, Nelson-Sathi S, Martin WF. 2016. The physiology and habitat of the last universal common ancestor. *Nat Microbiol*. **1**(9):1–8.
- Wherry RJ. 1931. A new formula for predicting the shrinkage of the coefficient of multiple correlation. *Ann Math Stat*. **2**(4):440–457.
- Wilson BA, Foy SG, Neme R, Masel J. 2017. Young genes are highly disordered as predicted by the preadaptation hypothesis of de novo gene birth. *Nat Ecol Evol*. **1**(6):0146.