

# Gene birth and constraint avoidance both contribute to structural disorder in overlapping genes

S. Willis and J. Mase1\*

Department of Ecology and Evolutionary Biology, University of Arizona

\*Corresponding Author: [masel@email.arizona.edu](mailto:masel@email.arizona.edu)

## Abstract

The same nucleotide sequence can encode two protein products in different reading frames. Overlapping gene regions are known to encode higher levels of intrinsic structural disorder (ISD) than non-overlapping genes (39% vs. 25% in our viral dataset). Two explanations for elevated ISD have been proposed: that high ISD relieves the increased evolutionary constraint imposed by dual-coding, and that one member per pair was recently born *de novo* in a process that favors high ISD. Here we quantify the relative contributions of these two alternative hypotheses, as well as a third hypothesis that has not previously been explored: that high ISD might be an artifact of the genetic code. We find that the recency of *de novo* gene birth explains  $\sim 32\%$  of the elevation in ISD in overlapping regions of viral genes, with the rest attributed, by a process of elimination, to relieving constraint. While the two reading frames within a same-strand overlapping gene pair have markedly different ISD tendencies, their effects cancel out such that the properties of the genetic code do not contribute overall to elevated ISD. Same-strand overlapping gene birth events can occur in two different frames, favoring high ISD either in the ancestral gene or in the novel gene; surprisingly, most *de novo* gene birth events contained completely within the body of an ancestral gene favor high ISD in the ancestral gene (23 phylogenetically independent events vs. 1). This can be explained by mutation bias favoring the frame with more start codons and fewer stop codons.

## I. INTRODUCTION

Protein-coding genes sometimes overlap, i.e. the same nucleotide sequence encodes different proteins in different reading frames. Most of the overlapping pairs of genes that have been characterized to date are found in viral, bacterial and mitochondrial genomes, with emerging research showing that they may be common in eukaryotic genomes as well [9], [18], [25], [29].

Overlapping genes tend to encode proteins with higher intrinsic structural disorder (ISD) than those encoded by non-overlapping genes [28]. The term disorder applies broadly to proteins which, at least in the absence of a binding partner, lack a stable secondary and tertiary structure. There are different degrees of disorder: molten globules, partially unstructured proteins and random coils with regions of disorder spanning from short (less than 30 residues in length) to long. Disorder can be shown experimentally or predicted from amino acid sequences using software [13]. Rancurel et al. showed, using the latter approach, that 48% of amino acids in overlapping regions exhibit disorder, compared to only 23% in non-overlapping regions. In this work we explore three non-mutually-

exclusive hypotheses to quantify the extent to which each explains elevated ISD. Two have previously been considered: that elevated ISD in overlapping genes is a mechanism that relieves evolutionary constraint, and that elevated ISD is a holdover from the *de novo* gene birth process. We add consideration of a third, previously-unexplored hypothesis - that elevated ISD with dual-coding may be the result of an artifact of the genetic code - to the mix.

The greater evolutionary constraint on overlapping genes is usually invoked as the sole [36], [42] or at least dominant [28] explanation for their high ISD. A mutation in an overlapping region simultaneously affects both of the two (or occasionally more) genes involved in that overlap. Because  $\sim 70\%$  of mutations that occur in the third codon position are synonymous, versus only  $\sim 5\%$  and  $0\%$  of mutations in the first and second codon positions respectively [30], a mutation that is synonymous in one reading frame is highly likely to be nonsynonymous in another, so to permit adaptation, overlapping genes must be relatively tolerant of nonsynonymous changes. Any amino acid substitution that maintains disorder has a reasonable chance of being tolerated, in contrast to the relative fragility of a well-

defined three-dimensional structure. This expectation is confirmed by the higher evolutionary rates observed for disordered proteins [4].

The second hypothesis that has previously been proposed is that high ISD in overlapping genes is an artifact of the process of de novo gene birth [28]. There is no plausible path by which two non-overlapping genes could re-encode an equivalent protein sequence as overlapping; instead, an overlapping pair arises either when a second gene is born de novo within an existing gene, or when the boundaries of an existing gene are extended to create overlap [30]. In the latter case of “overprinting” [7], [17], [28], the extended portion of that gene, if not the whole gene, is born de novo [26]. One overlapping protein-coding sequence is therefore always evolutionarily younger than the other; we refer to these as “novel” versus “ancestral” overlapping genes or portions of genes. Genes may eventually lose their overlap through a process of gene duplication followed by subfunctionalization [17], enriching overlapping genes for relatively young genes that have not yet been through this process. However, gene duplication may be inaccessible to many viruses (in particular, many RNA, ssDNA, and retroviruses), due to intrinsic geometric constraints on maximum nucleotide length [6], [8], [10].

Young genes are known to have higher ISD than old genes, with high ISD at the moment of gene birth facilitating the process [41], perhaps because cells tolerate them better [37]. Domains that were more recently born de novo also have higher ISD [3], [5], [12], [21]. High ISD could be helpful in itself in creating novel function, or it could be a byproduct of a hydrophilic amino acid composition whose function is simply the avoidance of harmful protein aggregation [14], [20]. Regardless of the cause of high ISD in young genes, the “facilitate birth” hypothesis makes a distinct prediction from the constraint hypothesis, namely that the novel overlapping reading frames will tend to encode higher ISD than the ancestral overlapping reading frames. Unlike the constraint hypothesis, only novel overlapping gene regions, and not ancestral ones, are predicted to have elevated ISD.

Finally, here we also consider the possibility that the high ISD observed in overlapping genes might simply be an artifact of the genetic code [19]. We perform for the first time the appropriate control, by predicting what the ISD would be if codons were read from alternative reading frames of existing non-overlapping genes. Any DNA sequence can be read in three reading frames on

each of the two strands, for a total of 6 reading frames. We focus only on same-strand overlap, due to superior availability of reliable data on same-strand overlapping gene pairs. We classify the reading frame of each gene in an overlapping pair relative to its counterpart; if gene A is in the +1 frame with respect to gene B, this means that gene B is in the +2 frame with respect to gene A. We use the +0 frame designation just for non-overlapping genes in their original frame. If the high ISD of overlapping genes is primarily driven by the intrinsic properties of the genetic code, then we expect their ISD values to closely match those expected from translation in the +1 vs. +2 frames of non-overlapping genes.

Here we test the predictions of all three hypotheses, as summarized in Figure 1, and find that both the birth-facilitation and conflict-resolution hypotheses play a role. The artifact hypothesis plays no appreciable role in elevating the ISD of overlapping regions; while reading frame (+1 vs. +2) strongly affects the ISD of individual genes, each overlapping gene pair has one of each, and the two cancel out such that there is no net contribution to the high ISD found in overlapping regions. Surprisingly, novel genes are more likely to be born in the frame prone to lower ISD; this seems to be a case where mutation bias in the availability of ORFs is more important than selection favoring higher ISD.

Hypothesis	ISD Prediction
Artifact of Genetic Code	+1 Frame = +1 Controls +2 Frame = +2 Controls
Conflict Resolution	Overlapping* > Non-Overlapping *Incl. Ancestral, Controlling for Frame Effects
Facilitate Birth	Novel > Frameshifted Controls Novel > Ancestral

Fig. 1. Three non-mutually-exclusive hypotheses about why overlapping genes have high ISD. The column on the right describes the ISD patterns we would expect if the hypotheses were true. Predictions of the conflict resolution hypothesis apply to all categories of overlapping genes, including ancestral.

## II. RESULTS AND DISCUSSION

### A. Causes of elevated ISD

Because most verified gene overlaps in the literature, especially longer overlapping sequences, are in viruses [24], [28], [39], we focused on viral genomes, compiling a list of 92 verified overlapping gene pairs from 80 viral species. The mean predicted ISD of all

overlapping regions ( $0.39 \pm 0.02$ ) was higher than that of the non-overlapping genes ( $0.25 \pm 0.01$ ), confirming previous findings that overlapping genes have elevated ISD.

To test whether the elevated ISD of overlapping genes is an artifact of the genetic code, we artificially generated amino acid sequences from the +1 and +2 reading frames of 150 non-overlapping viral genes in those 80 species and calculated their ISD using IUPred re phylogenetically independent except as noted in the footnotes. [11]. Mean ISD is higher in the +2 reading frame ( $0.35 \pm 0.02$ ) than in the +1 reading frame ( $0.19 \pm 0.01$ ). While exact ISD expectations are a function of %GC content [1] and hence species-specific, this result is not specific to viruses; *Mus musculus* yields a similar gap between the +2 reading frame ISD of  $0.573 \pm 0.002$  vs. +1 reading frame ISD of  $0.393 \pm 0.002$ .

The artifact hypothesis predicts that the +1 and +2 members of the 92 verified overlapping gene pairs will follow suit. In agreement with this, the overlapping regions of genes in the +2 reading frame had higher mean ISD ( $0.48 \pm 0.03$ ) than those in the +1 reading frame ( $0.31 \pm 0.02$ ). While this provides strong evidence that frame shapes ISD as an artifact of the genetic code, average ISD across both frameshifted control groups ( $0.27 \pm 0.01$ ) is significantly lower than the ISD of all overlapping sequences ( $0.39 \pm 0.02$ ), showing that the artifact hypothesis cannot fully explain elevated ISD in the latter.

We find stronger support for the birth-facilitation hypothesis. Of the 92 verified overlapping viral gene pairs, we were able to classify the relative ages of the component genes as ancestral vs. novel for 47 pairs (Table III-G). In agreement with the predictions of the birth-facilitation process, and controlling for frame, novel genes have higher ISD than either ancestral members of the same gene pairs or artificially-frameshifted controls (Figure 2). We confirmed this using a linear mixed model (on Box-Cox-transformed data with  $\lambda = 0.4$ ), with frame (+1 vs. +2) as a fixed effect, gene type (novel vs. ancestral vs. frameshifted controls) as a fixed effect, species (to control for %GC content and other subtle sequence biases) as a random effect, and homology group (to control for phylogenetic confounding) as a random effect. Within this linear model, the prediction unique to the birth-facilitation hypothesis, namely that ISD in the overlapping regions of novel genes is higher than that in ancestral genes, is supported with  $p = 0.01$ .

Our third hypothesis regarding the elevated ISD of

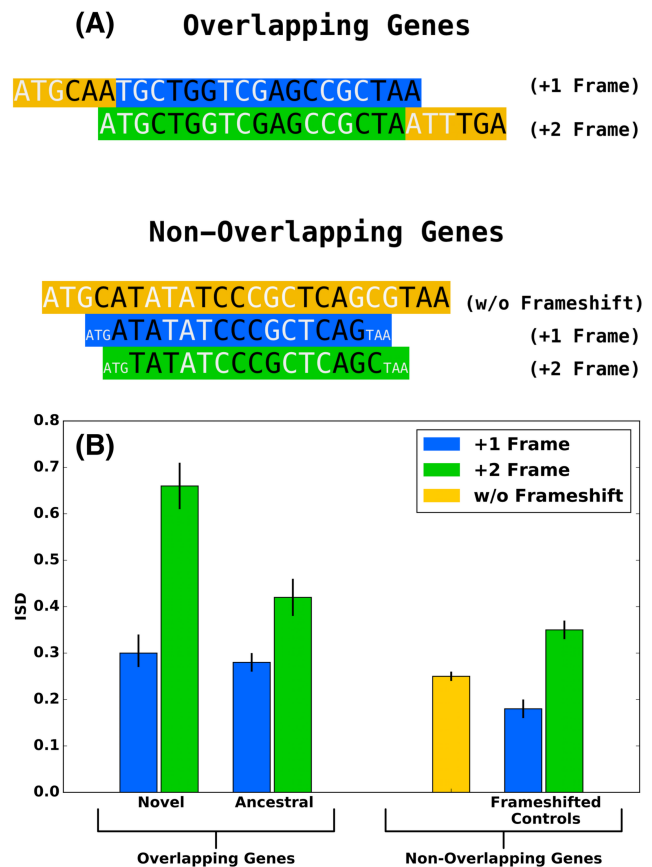


Fig. 2. Results support the birth-facilitation and conflict resolution hypotheses: novel > ancestral, and ancestral > non-genic sequences controlled for frame. (A) Data are from the overlapping sections of the 47 gene pairs whose ages could be classified, and from non-overlapping genes, and frameshifted versions of these non-overlapping genes, in the species in which the overlapping gene pairs were found. (B) While frame significantly impacts disorder content, it does not drive the high ISD of overlapping genes. Means and 66% confidence intervals were calculated from the Box-Cox transformed (with  $\lambda = 0.4$ ) means and their standard errors, and are shown here following back-transformation.

overlapping genes, that it loosens evolutionary constraint and so helps resolve conflict between paired genes, also plays a role. In agreement with this hypothesis, even ancestral overlapping sequences have higher ISD than non-overlapping genes (second vs. fourth cluster in Figure 2), and than frameshifted control sequences (see overlapping2cond cluster vs. yellow in Figure 2). In our linear model with two fixed effects and two random effects, the pairwise ancestral vs. frameshifted control comparison is supported with  $p = 0.02$ .

In contrast, a pairwise comparison between the non-overlapping genes and the mean of +1 and +2 frameshifted control versions of the same nucleotide

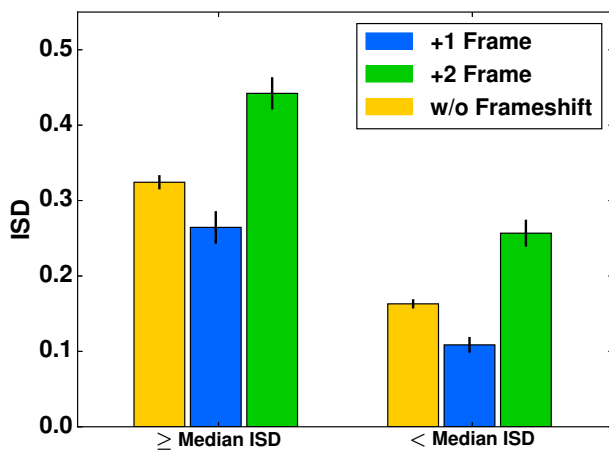


Fig. 3. The average ISD of artificially-frameshifted controls scales with the ISD of the non-overlapping genes that generated them. The 150 non-overlapping genes were separated into two groups, according to the median of the 150 original proteins. Below-median sequences have 41% GC, while above-mean sequences have 47% GC. Means and standard errors are shown.

sequences was not statistically significant ( $p = 0.85$ ; contrast statement applied to a linear model with fixed effect of actual non-overlapping gene sequence vs. +1 frameshifted version vs +2 frameshifted version vs a randomly scrambled version, using square-root transformed ISD values, and which non-overlapping gene sequence was used as a random effect). Despite the enormous effect of +1 vs. +2 reading frame, we find no support for the artifact hypothesis in explaining elevated ISD of overlapping regions. In each overlapping gene pair, there is always exactly one gene in each of the two reading frames, such that the large effects of each of the two frames cancel each other out when all overlapping genes are considered together. It is nevertheless important to control for the large effect of frame while testing and quantifying other hypotheses.

The relative magnitudes of the other two causes of high ISD in overlapping genes were quantified using our linear model with two fixed effects and two random effects. The degree to which birth facilitation elevates ISD was calculated, using a contrast statement, as half the difference between novel and ancestral genes, because exactly half of the genes are novel, and hence elevated above the “normal” ISD level of ancestral genes. Birth facilitation accounts for  $32.0\% \pm 9.7\%$  of the estimated total difference in ISD between overlapping and non-overlapping genes. We attribute the remainder, in the absence of other alternative hypotheses, to the pressure to relieve evolutionary constraint.

Note that frameshifted versions of high-ISD proteins

have higher ISD than frameshifted versions of low-ISD proteins (Figure 3). The ISD values of the two reading frames are likely linked via sharing the same %GC content, given that random sequences with higher %GC have substantially higher ISD [1]. A facilitate-birth bias toward high ISD thus imparts high %GC to overlapping genes at the time of birth, and so causes overlapping sequences to be biased toward not just high-ISD novel genes, but also high-ISD ancestral genes. This makes our estimate of the contribution of birth facilitation to elevated ISD an under-estimate. After birth, during evolution to avoid constraint, high ISD in the two reading frames is a positively correlated trait, and thus can more easily evolve in tandem.

### B. Frame of Gene Birth

Given the strong influence of frame combined with support for the facilitate-birth hypothesis, we hypothesized that novel genes would be born more often into the +2 frame (Figure 2A, green) because the intrinsically higher ISD of the +2 reading frame would facilitate high ISD in the novel gene and hence birth. Our dataset contained 41 phylogenetically independent overlapping pairs. Surprisingly, we found the opposite of our prediction: 31 of the novel genes were in the +1 frame of their ancestral counterparts, while only 10 were in the +2 frame ( $p = 10^{-3}$ , cumulative binomial distribution with trial success probability 0.5).

This unexpected result is stronger for “internal overlaps”, in which one gene is completely contained within its overlapping partner (23 +1 events vs. 1 +2 event,  $p = 3 \times 10^{-6}$ ), and is not found for “terminal overlaps”, in which the 5' end of the downstream gene overlaps with the 3' end of the upstream member of the pair (9 +1 events vs. 9 +2 events). (This double-counts a +1 event for which there were three homologous gene pairs, two of which were internal overlaps, and one of which was a terminal overlap.) Following [2], we interpret the restriction of this finding to internal overlaps as evidence that the cause of the bias applies to complete de novo gene birth, but not to the addition of a sequence to an existing gene.

The unexpected prevalence of +1 gene births, despite birth facilitation favoring +2, can be explained by mutation bias. One artifact of the genetic code is that +1 frameshifts yield more start codons and fewer stop codons, and hence fewer and shorter ORFs [2]. In our control set of 150 non-overlapping viral genes, we confirm that stop codons are more prevalent in the +2 frame (1 per 11 codons) than the +1 frame (1 per 14),

decreasing the mean ORF length, and that start codons are more prevalent in the +1 frame (1 per 27 codons) than the +2 frame (1 per 111). Our results are consistent with a major role for mutational availability in shaping adaptive evolution [34], [35], [43].

The preponderance of novel genes in the +1 frame further demonstrates the need to control for frame when testing hypotheses. Ancestral genes are more frequently in the +2 frame with elevated ISD, while the depressed ISD of the +1 frame lowers the ISD of the novel. As a result, when frame is not considered, ancestral and novel overlapping sequences encode very similar levels of disorder ( $0.41 \pm 0.03$  vs.  $0.42 \pm 0.04$ , respectively), making it easy to miss the evidence for the facilitate-birth hypothesis.

### III. MATERIALS AND METHODS

Scripts and data tables used in this work may be accessed at: [https://github.com/MaselLab/Willis\\_Masel\\_Overlapping\\_Genes\\_Structural\\_Disorder\\_Explained](https://github.com/MaselLab/Willis_Masel_Overlapping_Genes_Structural_Disorder_Explained)

#### A. Overlapping Viral Genes

A total of 102 viral same-strand overlapping gene pairs were collected from the literature [27], [28], [31]–[33], [40]. Of these, ten were discarded because one or both of the genes involved in the overlap were not found in the ncbi databases, either because the accession number had been removed, or because the listed gene could not be located. This left 92 gene pairs for analysis from 80 different species, spanning 33 viral families. Six of these pairs were ssDNA, five were retroviruses, while the remaining 81 were RNA viruses: 7 dsRNA, 61 positive sense RNA and 13 negative sense RNA.

#### B. Relative Gene Age

For 39 of the remaining 92 gene pairs available for analysis, the identity of the older vs. younger member of the pair had been classified in the literature [22], [28], [31], [32] via phylogenetic analysis. There was disagreement in the literature regarding the TGBp2/TGBp3 overlap; we followed [22] rather than [28].

We also used the relative levels of codon bias to classify the relative ages of members of each pair. Because all of the overlapping genes are from viral genomes, we can assume that they are highly expressed, leading to a strong expectation of codon bias in general. Novel genes are expected to have lower codon bias than ancestral genes due to evolutionary inertia [27].

For each viral species, codon usage data [23], [44] was used to calculate a relative synonymous codon usage (RSCU) value for each codon [15]:

$$\text{RSCU}_i = \frac{X_i}{\frac{1}{n} \sum_{i=1}^n X_i}$$

where  $X_i$  is the number of occurrences of codon  $i$  in the viral genome, and  $1 \leq n \leq 6$  is the number of synonymous codons which code for the same amino acid as codon  $i$ . The relative adaptedness value ( $w_i$ ) for each codon in a viral species was then calculated as:

$$w_i = \frac{\text{RSCU}_i}{\text{RSCU}_{\max}}$$

where  $\text{RSCU}_{\max}$  is the RSCU value for the most frequently occurring codon corresponding to the amino acid associated with codon  $i$ .

The codon adaptation index (CAI) was then calculated for the overlapping portion of each gene. The CAI is defined as the geometric mean of the relative adaptedness values:

$$\text{CAI} = \left( \prod_{i=1}^L w_i \right)^{\frac{1}{L}}$$

where  $L$  is the number of codons in the overlapping portion of the gene, excluding ATG and TGG codons. This exclusion is because ATG and TGG are the only codons that code for their respective amino acids and so their relative adaptedness values are always 1, thereby introducing no new information. To ensure sufficient statistical power to differentiate between CAI values, we did not analyze CAI for gene pairs with overlapping sections less than 200 nucleotides long.

Within each overlapping pair, we provisionally classified the gene with the higher CAI value as ancestral and the gene with lower CAI value as novel. We used a Mann-Whitney U Test to determine the statistical significance of the difference in CAI values for each gene pair, and chose a p-value cutoff of 0.035 after analyzing a receiver operating characteristic (ROC) plot (Figure 4A). The combined effects of our length threshold and p-value cutoff are illustrated in Figure 4B.

Of the 19 gene pairs whose ancestral vs. novel classification was obtained both by statistically significant CAI differences and by phylogenetics, there was one for which the CAI classification contradicted the phylogenetics. The exception was the p104/p130 overlap in the Providence virus. This overlap is notable

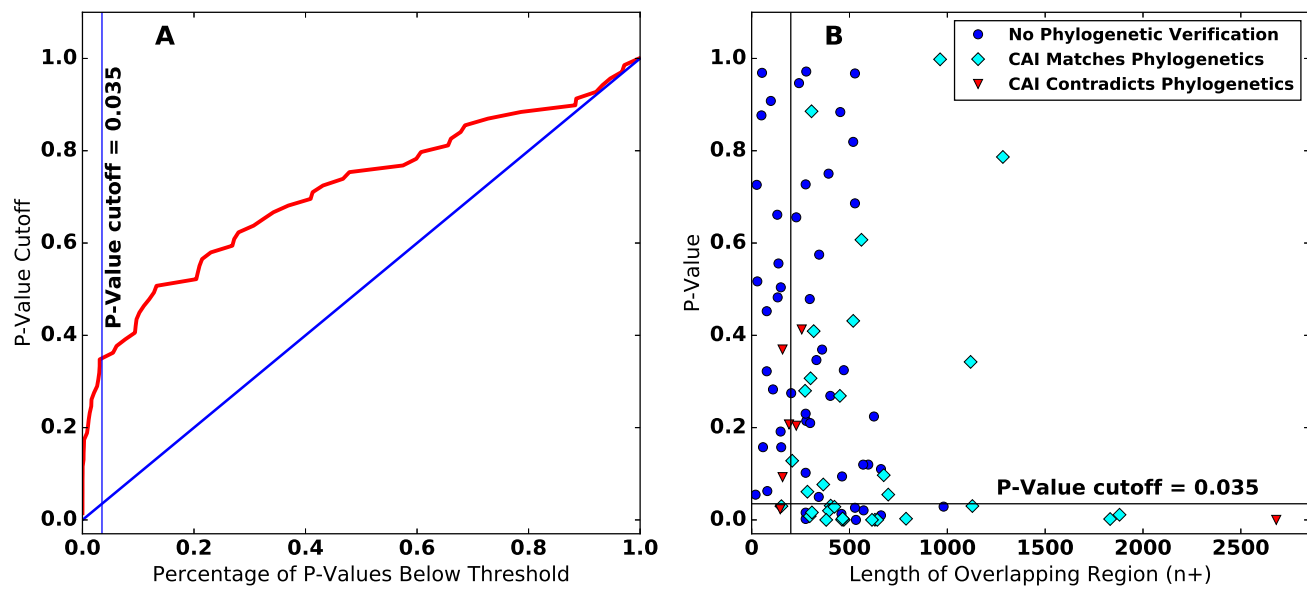


Fig. 4. Statistical classification of relative ages. (A) The receiver operating characteristic plot for determining which member of an overlapping gene pair has higher CAI, and is hence presumed to be ancestral. Only genes with an overlapping region of at least 200 nucleotides were included. (B) A plot of the p-values vs. the length of the overlapping regions of the 91 gene pairs for which codon usage data were available. The vertical line shows the overlapping length cutoff of 200 nucleotides

because the ancestral member of the pair was acquired through horizontal gene transfer, which renders codon usage an unreliable predictor of relative gene ages [27]. We therefore used the phylogenetic classification and disregarded the CAI results.

In total, we were able to classify ancestral vs. novel status for 47 overlapping gene pairs (Figure 5).

### C. Artificially-Frameshifted Viral Controls

150 non-overlapping control genes were compiled from the genomes of the viruses where the 92 overlapping gene pairs were found; matching for species helps ensure that results are not affected by %GC content or other idiosyncrasies of nucleotide composition. We removed one or two nucleotides immediately after the start codon and two or one nucleotides immediately before the stop codon in order to generate +1 and +2 frameshifted controls, respectively.

### D. Artificially-Frameshifted *mus musculus* Controls

A second set of frameshifted controls was generated from 22,778 protein-coding genes from the *Mus musculus* genome acquired from <http://uswest.ensembl.org>. 136 genes were excluded because they contained at least one N (unknown nucleotide) in their sequence, and an

additional 49 genes were excluded because their length was not a multiple of three, leaving 22,593 genes. In order to ensure independent datapoints, one gene was selected at random from each of the 10,664 gene families annotated by [41], for frameshifting as above.

### E. Homology Groups

Treating each gene as an independent datapoint is a form of pseudoreplication, because homologous genes can share properties via a common ancestor rather than via independent evolution. This problem of phylogenetic confounding can be corrected for by using gene family as a random effect term in a linear model [41], and by counting each gene birth event only once.

We constructed a pHMMer (<http://hmmer.org/>) database including all overlapping regions, non-overlapping genes and artificially-frameshifted controls. After an all-against-all search, sequences that were identified as homologous, using an expectation value threshold of  $10^{-4}$ , were provisionally assigned the same homology group ID. These provisional groups were used to determine which gene birth events were unique. Two pairs were considered to come from the same gene birth event when both the ancestral and the overlapping sequence were classified as homologous. We also used published phylogenetic analysis to classify the TGBp2/TGBp3 overlap as two birth events

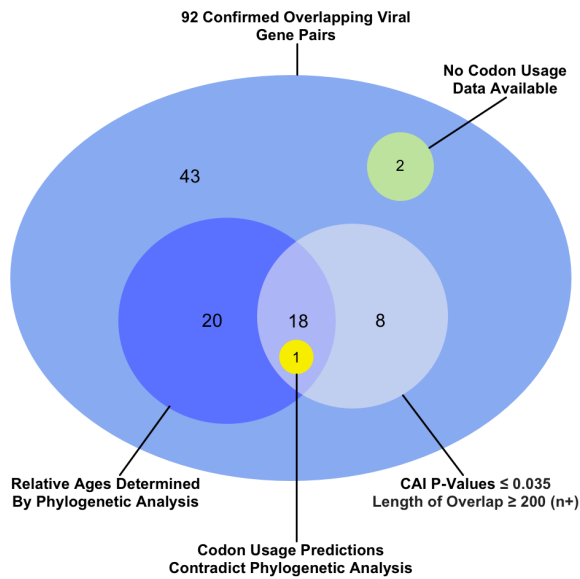


Fig. 5. A breakdown of the curated list of 92 confirmed overlapping gene pairs for which sequence data were available and how the relative ages of the genes were classified. Each gene pair is counted within only one of the subtotals shown.

(one occurring *Virgaviridae*, the other occurring in *Alpha-* and *Betaflexiviridae*) [22].

Some homologous pairs had such dissimilar protein sequences that ISD values were essentially independent. We therefore manually analyzed sequence similarity within each homology group using the Geneious [16] aligner with free end gaps, using Blosum62 as the cost matrix. The percent similarity using the Blosum62 matrix with similarity threshold 1 was then used as the criterion for whether a protein sequence would remain in its homology group for the ISD analysis. We used  $\geq 50\%$  protein sequence similarity as the threshold to assign a link between a pair, and then used single-link clustering to assign protein sequences to 561 distinct homology groups.

### F. ISD Prediction

We used IUPred [11] to calculate ISD values for each sequence. Following [41], before running IUPred, we excised all cysteines from each amino acid sequence, because of the uncertainty about their disulphide bond status and hence entropy [38]. Whether cysteine forms a disulphide bond depends on whether it is in an oxidizing or reducing environment. IUPred implicitly, through the selection of its training data set, assumes most cysteines are in disulphide bonds, which may or may not be accurate for our set of viral proteins.

Because cysteines have large effects on ISD (in either direction) depending on disulphid status and hence introduce large inaccuracies, cysteines were dropped from consideration altogether.

IUPred assigns a score between 0 and 1 to each amino acid. To calculate the ISD of an overlapping region, IUPred was run on the complete protein (minus its cysteines), then the average score was taken across only the pertinent subset of amino acids. These sequence-level ISD values were transformed using a Box-Cox transform. The optimal value of  $\lambda$  for the combined ancestral, novel and artificially-framedshifted control group data was 0.41, which we rounded to 0.4.

### G. Statistical models

Linear mixed models were generated using the lmer and gls functions contained in the nlme and lme4 R packages. In our main model, frame, gene designation (ancestral vs novel vs non-genic control), species, and homology group terms were retained in the model with  $p = 5 \times 10^{-20}$ ,  $1 \times 10^{-6}$ ,  $9 \times 10^{-19}$ , and  $2 \times 10^{-3}$  respectively. To determine the statistical significance of each effect, we used the anova function in R to compare nested models.

Accession Number	Organism	Ancestral Gene	Novel Gene	Overlap Length (n+)	Novel Frame
NC.001401	Adeno-Associated Virus 2	VP2	AAP	615	+1
NC.004285	Aedes Albopictus Dengovirus	NS1	NS2	1119	+1
NC.001467	African Cassava Mosaic Virus	AL1	AC4	423	+1
NC.009896	Akabane Virus <sup>1</sup>	N	NSs	276	+1
NC.001749	Apple Stem Grooving Virus	MP	Polyprotein	963	+1
NC.001719	Arctic Ground Squirrel Hepatitis Virus <sup>2</sup>	P	L	1284	+1
NC.003481	Barley Stripe Mosaic Virus <sup>3,4,5</sup>	TGBp2	TGBp3	191	+1
NC.003680	Barley Yellow Dwarf Virus <sup>6,7</sup>	P5	MP	465	+1
NC.005041	Blattella Germanica Dengovirus	NS-1	ORF4	789	+1
NC.001927	Bunyamwera Virus <sup>1</sup>	N	NSs	306	+1
NC.001658	Cassava Common Mosaic Virus <sup>4,5</sup>	TGBp2	TGBp3	152	+1
NC.001427	Chicken Anemia Virus	VP2	Apoptin	366	+1
NC.003688	Cucurbit Aphid-Born Yellowing Virus <sup>6,7,8</sup>	CP	P5	572	+1
NC.005899	Dendrolimus Punctatus Tetravirus <sup>6</sup>	p71	p17	381	+1
NC.016561	Hepatitis B <sup>2</sup>	P	L	1128	+1
NC.003608	Hibiscus Chlorotic Ringspot Virus <sup>6</sup>	Coat	p25	675	+1
NC.003608	Hibiscus Chlorotic Ringspot Virus	Replicase	p23	630	+1
NC.004730	Indian Peanut Clump Virus	P14	P17	158	+1
KR732417	Influenza A Virus H5N1	PB1	PB1-F2	273	+1
NC.009025	Israel Acute Paralysis Virus Of Bees	ORF2	ORFx	285	+1
NC.003627	Maize Chlorotic Mottle Virus	Coat	p31	451	+1
NC.001498	Measles Virus <sup>9</sup>	P	C	561	+1
NC.005339	Mossman Virus <sup>9</sup>	P	C	459	+1
NC.008311	Murine Norovirus	VP1	VF1	642	+1
NC.001633	Mushroom Bacilliform Virus	ORF1	Vpg-protease	533	+1
NC.001718	Porcine Parvovirus	Capsid	SAT	207	+1
NC.001747	Potato Leafroll Virus	P0	P1	661	+1
NC.003725	Potato Mop-Top Virus <sup>3,4</sup>	TGBp2	TGBp2	146	+1
NC.003768	Rice Dwarf Virus	Pns12	OP-ORF	276	+1
NC.003771	Rice Ragged Stunt Virus	P4b	Replicase	981	+1
NC.004718	SARS Coronavirus	Nucleocapsid	Protein I	297	+1
NC.003809	Spinach Latent Virus	Replicase	2b	308	+1
NC.003448	Striped Jack Nervous Necrosis Virus	Protein A	B2	228	+1
NC.001366	Theiler's Virus	L	L*	471	+1
NC.002199	Tupaia Paramyxovirus <sup>9</sup>	P	C	462	+1
NC.003743	Turnip Yellow Virus <sup>6,7,8</sup>	CP	ORF5	528	+1
NC.001409	Apple Chlorotic Leaf Spot Virus	CP	MP	317	+2
NC.001719	Arctic Ground Squirrel Hepatitis Virus	P	Capsid Precursor	158	+2
NC.001719	Arctic Ground Squirrel Hepatitis Virus	P	X	256	+2
NC.003532	Cymbidium Ringspot Virus	MP	p19	519	+2
NC.003093	Indian Citrus Ringspot Virus	CP	NABP	301	+2
NC.004178	Infectious Bursal Disease Virus <sup>10</sup>	VP2	VP5	404	+2
NC.001915	Infectious Pancreatic Necrosis Virus <sup>10</sup>	VP2	VP5	395	+2
NC.001990	Nudaurelia Capensis Beta Virus <sup>6</sup>	CP	Replicase	1832	+2
NC.014126	Providence Virus	p104	p130	2681	+2
NC.004366	Tobacco Bushy Top Virus	MP	RNP	698	+2
NC.004063	Turnip Yellow Mosaic Virus	Replicase	MP	1880	+2

<sup>1</sup>N/NSs overlaps share  $\geq 50\%$  sequence similarity

<sup>2</sup>P/L overlap predicted homologous in HMMer run

<sup>3</sup>TGBp2 genes share  $\geq 50\%$  protein sequence similarity

<sup>4</sup>TGBp2 genes predicted homologous (Morozov and Solovyev, 2003)

<sup>5</sup>TGBp3 genes predicted homologous (Morozov and Solovyev, 2003)

<sup>6</sup>Ancestral genes predicted homologous in HMMer run

<sup>7</sup>Novel genes predicted homologous in HMMer run

<sup>8</sup>Novel genes predicted homologous in HMMer run

<sup>9</sup>Novel genes predicted homologous in HMMer run

<sup>10</sup>Ancestral VP2 genes share  $\geq 50\%$  protein sequence similarity



#### IV. ACKNOWLEDGMENTS

We thank S. Foy and B. Wilson for programming assistance, M. Cordes and R. Neme for comments on the manuscript, and D. Karlin for pointing us to issues with the TGBp2/TGBp3 overlapping genes. This work was supported by the National Institutes of Health (R01 GM104040) and the John Templeton Foundation (60814).

#### REFERENCES

- [1] Annamária F Ángyán, András Perczel, and Zoltán Gáspári. Estimating intrinsic structural preferences of de novo emerging random-sequence proteins: Is aggregation the main bottleneck? *FEBS Letters*, 586(16):2468–2472, 2012.
- [2] Robert Belshaw, Oliver G Pybus, and Andrew Rambaut. The evolution of genome compression and genomic novelty in rna viruses. *Genome Research*, 17(10):1496–1504, 2007.
- [3] Erich Bornberg-Bauer and M Mar Alba. Dynamics and adaptive benefits of modular protein evolution. *Current Opinion in Structural Biology*, 23(3):459–466, 2013.
- [4] Celeste J Brown, Sachiko Takayama, Andrew M Campen, Pam Vise, Thomas W Marshall, Christopher J Oldfield, Christopher J Williams, and A Keith Dunker. Evolutionary rate heterogeneity in proteins with long disordered regions. *Journal of Molecular Evolution*, 55(1):104–110, 2002.
- [5] Marija Buljan, Adam Frankish, and Alex Bateman. Quantifying the mechanisms of domain gain in animal proteins. *Genome Biology*, 11(7):R74, 2010.
- [6] José A Campillo-Balderas, Antonio Lazcano, and Arturo Becerra. Viral genome size distribution does not correlate with the antiquity of the host lineages. *Frontiers in Ecology and Evolution*, 3:143, 2015.
- [7] Joseph J Carter, Matthew D Daugherty, Xiaojie Qi, Anjali Bheda-Malge, Gregory C Wipf, Kristin Robinson, Ann Roman, Harmit S Malik, and Denise A Galloway. Identification of an overprinting gene in merkel cell polyomavirus provides evolutionary insight into the birth of viral genes. *Proceedings of the National Academy of Sciences*, 110(31):12744–12749, 2013.
- [8] Nicola Chirico, Alberto Vianelli, and Robert Belshaw. Why genes overlap in viruses. *Proceedings of the Royal Society of London B: Biological Sciences*, 277(1701):3809–3817, 2010.
- [9] Wen-Yu Chung, Samir Wadhawan, Radek Szklarczyk, Sergei Kosakovsky Pond, and Anton Nekrutenko. A first look at arfome: dual-coding genes in mammalian genomes. *PLoS Computational Biology*, 3(5):e91, 2007.
- [10] John M Coffin, Stephen H Hughes, and Harold E Varmus. Principles of retroviral vector design. 1997.
- [11] Zsuzsanna Dosztányi, Veronika Csizmók, Péter Tompa, and István Simon. The pairwise energy content estimated from amino acid composition discriminates between folded and intrinsically unstructured proteins. *Journal of Molecular Biology*, 347(4):827–839, 2005.
- [12] Diana Ekman and Arne Elofsson. Identifying and quantifying orphan protein sequences in fungi. *Journal of Molecular Biology*, 396(2):396–405, 2010.
- [13] François Ferron, Sonia Longhi, Bruno Canard, and David Karlin. A practical overview of protein disorder prediction methods. *Proteins: Structure, Function, and Bioinformatics*, 65(1):1–14, 2006.
- [14] Scott G Foy, Benjamin A Wilson, Matthew HJ Cordes, and Joanna Masel. Progressively more subtle aggregation avoidance strategies mark a long-term direction to protein evolution. *bioRxiv*, 176867, 2017.
- [15] Dan Graur. *Molecular and Genome Evolution*, chapter 4, pages 140–141. Sinauer Associates Inc., first edition, 2016.
- [16] Matthew Kearse, Richard Moir, Amy Wilson, Steven Stones-Havas, Matthew Cheung, Shane Sturrock, Simon Buxton, Alex Cooper, Sidney Markowitz, Chris Duran, et al. Geneious basic: an integrated and extendable desktop software platform for the organization and analysis of sequence data. *Bioinformatics*, 28(12):1647–1649, 2012.
- [17] Paul K Keese and Adrian Gibbs. Origins of genes: “big bang” or continuous creation? *Proceedings of the National Academy of Sciences*, 89(20):9489–9493, 1992.
- [18] Dae-Soo Kim, Chi-Young Cho, Jae-Won Huh, Heui-Soo Kim, and Hwan-Gue Cho. Evog: a database for evolutionary analysis of overlapping genes. *Nucleic Acids Research*, 37(suppl\_1):D698–D702, 2008.
- [19] Erika Kovacs, Peter Tompa, Karoly Liliom, and Lajos Kalmar. Dual coding in alternative reading frames correlates with intrinsic protein disorder. *Proceedings of the National Academy of Sciences*, 107(12):5429–5434, 2010.
- [20] Zhirong Liu and Yongqi Huang. Advantages of proteins being disordered. *Protein Science*, 23(5):539–550, 2014.
- [21] Andrew D Moore and Erich Bornberg-Bauer. The dynamics and evolutionary potential of domain loss and emergence. *Molecular Biology and Evolution*, 29(2):787–796, 2011.
- [22] Sergey Yu Morozov and Andrey G Solovyev. Triple gene block: modular design of a multifunctional machine for plant virus movement. *Journal of General Virology*, 84(6):1351–1366, 2003.
- [23] Y. Nakamura, T. Gojobori, and T. Ikemura. Codon usage tabulated from the international dna sequence databases: status for the year 2000. *Nucleic Acids Research*, 28:292, 2000.
- [24] Tomohiro Nakayama, Satoshi Asai, Yasuo Takahashi, Oto Maekawa, and Yasuji Kasama. Overlapping of genes in the human genome. *International Journal of Biomedical Science: IJBS*, 3(1):14, 2007.
- [25] Anton Nekrutenko, Samir Wadhawan, Paula Goetting-Minesky, and Kateryna D Makova. Oscillating evolution of a mammalian locus with overlapping reading frames: an xlas/alex relay. *PLoS Genetics*, 1(2):e18, 2005.
- [26] Rafik Neme and Diethard Tautz. Phylogenetic patterns of emergence of new genes support a model of frequent de novo evolution. *BMC Genomics*, 14(1):117, 2013.
- [27] Angelo Pavesi, Gkikas Magiorkinis, and David G Karlin. Viral proteins originated de novo by overprinting can be identified by codon usage: application to the gene nursery of deltaretroviruses. *PLoS Computational Biology*, 9(8):e1003162, 2013.
- [28] Corinne Rancurel, Mahvash Khosravi, A Keith Dunker, Pedro R Romero, and David Karlin. Overlapping genes produce proteins with unusual sequence properties and offer insight into de novo protein creation. *Journal of Virology*, 83(20):10719–10736, 2009.
- [29] Sebastien Ribrioux, Adrian Brünger, Birgit Baumgarten, Klaus Seuwen, and Markus R John. Bioinformatics prediction of overlapping frameshifted translation products in mammalian transcripts. *BMC Genomics*, 9(1):122, 2008.
- [30] Niv Sabath, Dan Graur, and Giddy Landan. Same-strand overlapping genes in bacteria: compositional determinants of phase bias. *Biology Direct*, 3(1):36, 2008.
- [31] Niv Sabath, Andreas Wagner, and David Karlin. Evolution of

- viral proteins originated de novo by overprinting. *Molecular Biology and Evolution*, mss179, 2012.
- [32] Aditi Shukla and Rolf Hilgenfeld. Acquisition of new protein domains by coronaviruses: analysis of overlapping genes coding for proteins n and 9b in sars coronavirus. *Virus Genes*, 50(1):29–38, 2015.
- [33] Etienne Simon-Loriere, Edward C Holmes, and Israel Pagán. The effect of gene overlapping on the rate of rna virus evolution. *Molecular Biology and Evolution*, 30(8):1916–1928, 2013.
- [34] Arlin Stoltzfus and David M McCandlish. Mutational biases influence parallel adaptation. *bioRxiv*, 114694, 2017.
- [35] Arlin Stoltzfus and Lev Y Yampolsky. Climbing mount probable: mutation as a cause of nonrandomness in evolution. *Journal of Heredity*, 100(5):637–647, 2009.
- [36] Nobuhiko Tokuriki, Christopher J Oldfield, Vladimir N Uversky, Igor N Berezovsky, and Dan S Tawfik. Do viral proteins possess unique biophysical features? *Trends in Biochemical Sciences*, 34(2):53–59, 2009.
- [37] Vyacheslav Tretyachenko, Jiří Vymětal, Lucie Bednárová, Vladimír Kopecký, Kateřina Hofbauerová, Helena Jindrová, Martin Hubálek, Radko Souček, Jan Konvalinka, Jiří Vondrášek, et al. Random protein sequences can form defined secondary structures and are well-tolerated in vivo. *Scientific Reports*, 7(1):15449, 2017.
- [38] Vladimir N Uversky and A Keith Dunker. Understanding protein non-folding. *Biochimica et Biophysica Acta (BBA)-Proteins and Proteomics*, 1804(6):1231–1264, 2010.
- [39] Vamsi Veeramachaneni, Wojciech Makalowski, Michal Galdzicki, Raman Sood, and Izabela Makalowska. Mammalian overlapping genes: the comparative perspective. *Genome Research*, 14(2):280–286, 2004.
- [40] Robert G Webster, William J Bean, Owen T Gorman, Thomas M Chambers, and Yoshihiro Kawaoka. Evolution and ecology of influenza a viruses. *Microbiological Reviews*, 56(1):152–179, 1992.
- [41] Benjamin A Wilson, Scott G Foy, Rafik Neme, and Joanna Masel. Young genes are highly disordered as predicted by the preadaptation hypothesis of de novo gene birth. *Nature Ecology & Evolution*, 1(6):0146, 2017.
- [42] Bin Xue, David Blocquel, Johnny Habchi, Alexey V Uversky, Lukasz Kurgan, Vladimir N Uversky, and Sonia Longhi. Structural disorder in viral proteins. *Chemical Reviews*, 114(13):6880–6911, 2014.
- [43] Lev Y Yampolsky and Arlin Stoltzfus. Bias in the introduction of variation as an orienting factor in evolution. *Evolution & Development*, 3(2):73–83, 2001.
- [44] Tong Zhou, Wanjun Gu, Jianmin Ma, Xiao Sun, and Zuhong Lu. Analysis of synonymous codon usage in h5n1 virus and other influenza a viruses. *Biosystems*, 81(1):77–86, 2005.