1

**Figure S1.** IUPred "long" and "short" estimators of intrinsic structural disorder disagree on the relation

between GC content and the intrinsic structural disorder of junk polypeptides and novel functional

polypeptides under a random-sequence model. (*A*) Contour plot of the predicted average of IUPred

"long" disorder among novFPs (identical to fig. 2B). (*B*) Contour plot of the predicted average IUPred

"short" disorder among novFPs. (*C*) The predicted mean and standard deviation of IUPred "long"

disorder among JPs as functions of the GC content (identical to fig. 2D). (*D*) The predicted mean and

standard deviation of IUPred "short" disorder among JPs as functions of the GC content. Hatched areas

9    indicate impossible percentages of ISD, i.e. outside the 0%-100% interval. The landscapes in panels A

10   and B can be understood as the results of applying equation 2 to the curves in panels C and D,

11   respectively. As a result, the vertical "slice" of a landscape at a given GC content is a straight line

12   whose intercept and slope are respectively the mean and standard deviation associated with this GC

13   content in the corresponding bottom panel. The curve obtained by taking a horizontal "slice" where

14   there is no birth bias ($\delta = 0$) corresponds to the relation between the mean of the property among JPs,

15   i.e. the solid blue curve in the corresponding bottom panel. Since the vertical distance between contour

16   lines is inversely proportional to the vertical slope of the landscape, it is inversely proportional to the

17   standard deviation of the property among JPs, i.e. the dashed red curve in the corresponding bottom

18   panel.

19

20                            SUPPLEMENTARY METHODS

21

**22   Applying the Radon-Nikodým theorem to de novo gene birth**

23

24   This section explains why our framework fits the general setting of the branch of mathematics called

25   measure theory and its sub-branch, probability theory. We introduce some concepts from these theories

26   to clarify why the Radon-Nikodým theorem can be used to compare JPs and novFPs.

27

28   Given a set $\Omega$, measure theory provides the basic notions required to develop a self-consistent concept

29   of the "measure" or "size" of subsets of $\Omega$ (such as length, area, volume or probability). It is not always

30   possible to consistently define a measure for all subsets of $\Omega$, so that we must choose certain subsets

31   that form a structure called a $\sigma$-field (or $\sigma$-algebra). A $\sigma$-field on $\Omega$ is a set $\mathcal{F}$ whose elements are

32   subsets of $\Omega$ that meet certain conditions. The consequences of these conditions are that both $\Omega$ and the

33   empty set $\emptyset$ are elements of $\mathcal{F}$, and the combination of arbitrary elements of $\mathcal{F}$ through a finite or

34   infinite sequence of standard set operations (union, intersection, complementation, difference and

35   symmetric difference) always produces an element of $\mathcal{F}$ (Vestrup 2003a).

36

37   In our framework, the elements of $\Omega$ are all the possible polypeptides that are distinct in terms of

38   sequence and/or *cis*-regulation, and the elements of $\mathcal{F}$ are classes of polypeptides. Since the sequence

39   and *cis*-regulatory properties of a polypeptide are determined by a finite DNA sequence containing its

40   ORF, the set $\Omega$ is discrete or "countable", i.e. it is not larger than the set of all whole numbers (Komjáth

41   and Totik 2006). Because of this, we can choose $\mathcal{F}$ to be the set of all subsets of $\Omega$, which would cause

42   complications if $\Omega$ was a continuum (Vestrup 2003b). Nevertheless, we will continue the explanations

43   in the context of an arbitrary $\sigma$-field because that is how the Radon-Nikodým theorem is formulated.

44

45   A measure $\mu$ defined on a $\sigma$-field $\mathcal{F}$ of subsets of $\Omega$ is a function that assigns a number to each element

46   of $\mathcal{F}$. If $S$ is an element of $\mathcal{F}$, then $\mu(S)$ denotes the number that $\mu$ assigns to $S$. To meet the definition

47   of a measure, $\mu$ must also satisfy three other conditions: 1) $\mu(S) \geq 0$ for each $S \in \mathcal{F}$, 2) $\mu(\emptyset) = 0$,

48   where $\emptyset$ is the empty set, and 3) for any finite or infinite sequence $S_1, S_2, S_3, ...$ of non-overlapping

49   elements of $\mathcal{F}$, their union $S = S_1 \cup S_2 \cup S_3 ...$ satisfies $\mu(S) = \mu(S_1) + \mu(S_2) + \mu(S_3) + \cdots$ (Vestrup

50   2003c). The triple $(\Omega, \mathcal{F}, \mu)$ is called a measure space. If a measure $P$ defined on $\mathcal{F}$ also satisfies

51   $P(\Omega) = 1$, then $P$ is called a probability measure and $(\Omega, \mathcal{F}, P)$ is called a probability space, and they

52   are studied by probability theory.

53

54   In our framework, we define two probability measures: $P$, which represents a time average of JPs, and

55    $P_F$, which represents novFPs that functionalize in the time period considered. These measures are

56    defined on the same $\sigma$-field $\mathcal{F}$; they assign numbers to the same classes of polypeptides. Given $S$, a

57    subset of $\Omega$ which is an element of $\mathcal{F}$, the number $P(S)$ is the ratio of the time-averaged number of JPs

58    that belong to $S$ to the time-averaged total number of JPs. We can see that $P$ meets the three

59    requirements that define a measure: the ratio is never negative ($P(S) \geq 0$), the empty set contains no

60    JPs ($P(\emptyset) = 0$) and the ratio assigned to the union of several non-overlapping classes of polypeptides

61    is the sum of their individual ratios (the numerators add up and the denominator is a constant). $P$ is a

62    probability measure since $P(\Omega) = 1$, i.e. the time-averaged number of JPs that belong to $\Omega$ is precisely

63    the time-averaged total number of JPs. $P_F$ is also a probability measure: we define $P_F(S)$ as the

64    proportion of novFPs that belong to $S$. Proportions are never negative ($P_F(S) \geq 0$), the empty set

65    contains no novFPs ($P_F(\emptyset) = 0$), the proportion of novFPs belonging to the union of several non-

66    overlapping classes is the sum of the proportions belonging to each class, and the proportion of novFPs

67    belonging to the whole set $\Omega$ is $P_F(\Omega) = 1$.

68

69    Measure theory defines the notion of the integral, with respect to a measure and over a specific subset

70    of $\Omega$, of a numerical function. We use such functions to represent polypeptide properties such as length

71    and intrinsic disorder, and their integrals determine their averages among polypeptides. A function $f$

72    defined on the set $\Omega$ is a function that assigns a number $f(\omega)$ to each element $\omega$ of $\Omega$. Given a measure

73    space $(\Omega, \mathcal{F}, \mu)$, a function $f$ on $\Omega$ must have a property called $\mathcal{F}/\mathcal{B}^*$-measurability in order for its

74    integral to be well-defined. $f$ is said to be $\mathcal{F}/\mathcal{B}^*$-measurable if, for every real number $x$, there is an

75    element of $\mathcal{F}$ (called $f^{-1}((x, +\infty])$) that is exactly the set of all elements $\omega$ of $\Omega$ which satisfy

76    $f(\omega) > x$ (Vestrup 2003d). Given an $\mathcal{F}/\mathcal{B}^*$-measurable function $f$ and a subset $S$ of $\Omega$ which is an

77    element of $\mathcal{F}$, the integral of $f$ over $S$ with respect to $\mu$ is a number denoted by $\int_S f \, d\mu$. Given a

78     probability space $(\Omega, \mathcal{F}, P)$, the conditional average of $f$ "knowing" $S$ is given (Çinlar 2011) by:

$$E(f|S) = \frac{1}{P(S)} \int_S f \, dP$$

79     In particular, the (unconditional) average of $f$ is given by:

$$E(f) = E(f|\Omega) = \frac{1}{P(\Omega)} \int_\Omega f \, dP = \int_\Omega f \, dP$$

80

81     Given two measures $\mu$ and $\upsilon$ defined on the same $\sigma$-field $\mathcal{F}$ of subsets of $\Omega$, $\upsilon$ is said to be absolutely

82     continuous with respect to $\mu$ if every element $S$ of $\mathcal{F}$ which satisfies $\mu(S) = 0$ also satisfies $\upsilon(S) = 0$.

83     This relationship between $\mu$ and $\upsilon$ is also denoted by $\upsilon \ll \mu$ (Vestrup 2003e). In our framework, the

84     measure $P$ represents a time average of JPs and $P_F$ represents novFPs that functionalize in the time

85     period considered. This implies that for each novFP represented in the measure $P_F$, the JP that it was

86     immediately before functionalization is represented in the measure $P$. These two polypeptides are

87     identical because of our definition of novFPs, so they belong to exactly the same subsets of $\Omega$.

88     Therefore, if a subset $S$ of $\Omega$ is an element of $\mathcal{F}$ and never contains any JPs ($P(S) = 0$), then no novFPs

89     emerge in this subset ($P_F(S) = 0$). Thus, we have $P_F \ll P$.

90

91     The Radon-Nikodým theorem for finite measures states that given two measures $\mu$ and $\upsilon$ on $\mathcal{F}$ which

92     are both finite ($\mu(\Omega)$ and $\upsilon(\Omega)$ are finite numbers) and which satisfy $\upsilon \ll \mu$, there exists a finite-valued

93     nonnegative $\mathcal{F}/\mathcal{B}^*$-measurable function $f$ on $\Omega$ which summarizes the relationship between $\mu$ and $\upsilon$.

94     Specifically, $\upsilon$ can be constructed by integrating $f$ with respect to $\mu$; for each element $S$ of $\mathcal{F}$, we have

95     $\upsilon(S) = \int_S f \, d\mu$ (Vestrup 2003e). In our framework, this theorem applies to the measures $P$ and $P_F$

96     since they are both finite ($P(\Omega) = P_F(\Omega) = 1$) and $P_F \ll P$. Therefore, there exists a finite-valued

97     nonnegative $\mathcal{F}/\mathcal{B}^*$-measurable function $\hat{r}$ on $\Omega$ (a polypeptide property) such that for each element $S$ of

98     $\mathcal{F}$, we have $P_F(S) = \int_S \hat{r}\, dP$. Because of the definition of the conditional average (Çinlar 2011), we

99     have $P_F(S) = P(S) \times E(\hat{r}|S)$ and thus:

$$E(\hat{r}|S) = \frac{P_F(S)}{P(S)}$$

100     where $E(\hat{r}|S)$ is the average of $\hat{r}$ among JPs that belong to the class $S$. This provides an interpretation

101     of the polypeptide property $\hat{r}$: its average among JPs that belong to a given class of polypeptides ($S$) is

102     the ratio of the frequency of this class among novFPs ($P_F(S)$) to its frequency among JPs ($P(S)$). Since

103     a class of polypeptides may be arbitrarily small and may even contain only one JP, the value of $\hat{r}$ for a

104     single polypeptide is the factor by which its frequency changes from JPs to novFPs. We can deduce

105     from the above equation that the average of $\hat{r}$ among JPs is $E(\hat{r}) = 1$, since:

$$E(\hat{r}) = E(\hat{r}|\Omega) = \frac{P_F(\Omega)}{P(\Omega)} = \frac{1}{1} = 1$$

106

107     The function $f$ defined from two measures $v \ll \mu$ by the Radon-Nikodým theorem has a useful

108     property: for every $\mathcal{F}/\mathcal{B}^*$-measurable function $g$, its integral with respect to $v$ over any element $S$ of $\mathcal{F}$

109     is given by $\int_S g\, dv = \int_S f g\, d\mu$ (Vestrup 2003e). In our framework, this property translates to

110     $\int_S q\, dP_F = \int_S q\hat{r}\, dP$ for any polypeptide property $q$. By the definition of the conditional average

111     (Çinlar 2011), we thus have:

$$P_F(S) \times E_F(q|S) = P(S) \times E(q\hat{r}|S)$$

$$E_F(q|S) = \frac{P(S)}{P_F(S)} \times E(q\hat{r}|S)$$

$$E_F(q|S) = \frac{E(q\hat{r}|S)}{E(\hat{r}|S)}$$

112     where $E_F(q|S)$ is the average of $q$ among novFPs that belong to $S$. If we choose $S = \Omega$ and use the fact

113     that $E(\hat{r}) = 1$, we obtain:

$$E_F(q) = E(q\hat{r})$$

114     which shows how the relationship between $q$ and $\hat{r}$ among JPs determines the average of $q$ among

115     novFPs. From this equation, our main mathematical results can be derived using the universal

116     properties of averages, variances, covariances, etc. without further need for the basic concepts of

117     measure theory.

118

119     **Interpreting the coskewness of three variables**

120

121     To facilitate the interpretation of the coskewness of three variables $cosk(x,y,z) = \frac{E(\Delta x \Delta y \Delta z)}{\sigma(x)\sigma(y)\sigma(z)}$, where

122     $\Delta x = x - E(x)$ , consider the standard score $Z(x) = \frac{\Delta x}{\sigma(x)}$ which has a mean of 0 and a variance of 1.

$$cosk(q,\lambda,f) = E\big(Z(q)Z(\lambda)Z(f)\big)$$

123     Since $E(x\,y) = E(x){\times}E(y) + cov(x,y)$ , we obtain:

$$cosk(q,\lambda,f) = E\big(Z(q)\big){\times} E\big(Z(\lambda)Z(f)\big) + cov\big(Z(q),Z(\lambda)Z(f)\big)$$

124     Because $E\big(Z(q)\big) = 0$ , we obtain:

$$cosk(q,\lambda,f) = cov\big(Z(q),Z(\lambda)Z(f)\big)$$

125     Because of the definition of coskewness, its value does not change when we swap any two of the three

126     variables:

$$cosk(q,\lambda,f) = cov\big(Z(q),Z(\lambda)Z(f)\big) = cov\big(Z(\lambda),Z(q)Z(f)\big) = cov\big(Z(f),Z(\lambda)Z(q)\big)$$

127

128     Now consider the fact that $E(Z(x)Z(y)) = cov(Z(x),Z(y)) = \rho(x,y)$. Put in words, the Pearson

129     correlation coefficient is the mean of the product of the standard scores of two variables, while

130     coskewness is the covariance of this same product with the standard score of a third variable.

131 Therefore, roughly speaking, coskewness is a measure of how any of the three variables linearly affects

132 the correlation between the two others.

133

134 <div align="center">LITERATURE CITED</div>

135

136 Çinlar, E., 2011 Conditioning, pp. 139-170 in *Probability and Stochastics*, edited by S. Axler and K. A.

137 Ribet. Springer, New York.

138 Komjáth, P., and V. Totik, 2006 Countability, pp. 9-12 in *Problems and Theorems in Classical Set*

139 *Theory*, edited by P. Winkler. Springer, New York.

140 Vestrup, E. M., 2003a Set Systems, pp. 1-34 in *The theory of measures and integration*. John Wiley &

141 Sons, Inc., Hoboken.

142 Vestrup, E. M., 2003b Lebesgue Measure, pp. 113-162 in *The theory of measures and integration*. John

143 Wiley & Sons, Inc., Hoboken.

144 Vestrup, E. M., 2003c Measures, pp. 35-61 in *The theory of measures and integration*. John Wiley &

145 Sons, Inc., Hoboken.

146 Vestrup, E. M., 2003d Measurable Functions, pp. 163-207 in *The theory of measures and integration*.

147 John Wiley & Sons, Inc., Hoboken.

148 Vestrup, E. M., 2003e The Radon-Nikodym Theorem, pp. 367-436 in *The theory of measures and*

149 *integration*. John Wiley & Sons, Inc., Hoboken.