



OPEN

# Optimizing strategy for the discovery of compositionally-biased or low-complexity regions in proteins

Paul M. Harrison

Proteins can contain tracts dominated by a subset of amino acids and that have a functional significance. These are often termed 'low-complexity regions' (LCRs) or 'compositionally-biased regions' (CBRs). However, a wide spectrum of compositional bias is possible, and program parameters used to annotate these regions are often arbitrarily chosen. Also, investigators are sometimes interested in longer regions, or sometimes very short ones. Here, two programs for annotating LCRs/CBRs, namely SEG and fLPS, are investigated in detail across the whole expanse of their parameter spaces. In doing so, boundary behaviours are resolved that are used to derive an optimized systematic strategy for annotating LCRs/CBRs. Sets of parameters that progressively annotate or 'cover' more of protein sequence space and are optimized for a given target length have been derived. This progressive annotation can be applied to discern the biological relevance of CBRs, e.g., in parsing domains for experimental constructs and in generating hypotheses. It is also useful for picking out candidate regions of interest of a given target length and bias signature, and for assessing the parameter dependence of annotations. This latter application is demonstrated for a set of human intrinsically-disordered proteins associated with cancer.

Despite being composed of an alphabet of twenty diverse amino acids, proteins can often demonstrate a compositional bias (CB) for a small subset of this residue alphabet. For example, the sequence PPQPPSPPPSPPPPPQPPP is biased for the single residue P (proline), and the sequence ARGGRARGARRRGAAAGAGGRAGSAG is biased for A (alanine), R (arginine) and G (glycine). Simpler, more repetitive regions that tend to be shorter are often termed 'low-complexity' regions (LCRs). However, longer compositionally-biased regions (CBRs) can have quite a mild compositional skew, and LCRs can be considered a subset of CBRs. Many CBRs are composed of tandem repeats of sequence units several residues long. CBRs can be found in intrinsically disordered proteins, fibrous proteins, cell-structural proteins, functional amyloids and prions, or globular domains with specific functional roles such as metal binding<sup>1</sup>. They are also a significant component of the 'dark proteome' that has been chronically un- or understudied<sup>2</sup>.

Discovery of LCRs/CBRs in proteins has been actively researched, with several programs being developed. These include SIMPLE<sup>3,4</sup>, SEG<sup>5,6</sup>, Oj.py<sup>7</sup>, ScanCom<sup>8</sup>, CARD<sup>9</sup>, BIAS<sup>10</sup>, SARP<sup>11</sup>, LCD-Composer<sup>12,13</sup> and LPS/fLPS<sup>1,14-16</sup>. Lee, et al. developed a method for picking out low-complexity regions using image processing of dot plots<sup>17</sup>. Furthermore, servers pooling results of multiple methods have been produced, including LCT, LCRExplorer and PLaToLoCo<sup>18-20</sup>. SEG labels LCRs by scanning sequences with a fixed window length and applying thresholds for sequence entropy<sup>5</sup>. This algorithm has long been a component of the BLAST sequence alignment suite, wherein it can filter for false positive sequence matches arising because of simple, low-complexity sequence composition<sup>21</sup>. The program fLPS uses binomial probability to pick out low-probability sequence tracts<sup>14,15</sup>. It has been applied to studying the evolution of prions and prion-like regions<sup>22-26</sup>, and to characterizing the 'dark proteome'<sup>2</sup>, and by many other investigators to aid in characterization of protein sub-domains<sup>15</sup>.

SEG and fLPS are particularly useful for large-scale automated analysis of CBRs, since they do not require specification of residue type lists. They can characterize regions made from multiple-residue bias; they can also delineate milder biased regions. The benefits of SEG include that background amino-acid frequencies do not

Department of Biology, McGill University, Montreal, QC, Canada. email: paul.harrison@mcgill.ca

need to be considered, and its rapid calculations<sup>5</sup>. fLPS is a faster algorithm that works by detecting specific amino-acid biases<sup>14,15</sup>. Single-residue and multiple-residue CBRs are calculated explicitly, several options for background amino-acid frequencies are offered, and its two-window system is designed to capture a diversity of region lengths.

What is a low-complexity region and what is not? When is a protein domain compositionally biased? The answer to these questions is not simple. Many different thresholds are possible for labelling these regions that will ‘cover’ smaller or larger amounts of protein sequences. Also, investigators might be interested in very short regions, or sometimes longer ones. Although, algorithms for their discovery have been published with ‘recommended’ parameter sets for CBR annotation, there has been no systematic, thorough examination of what parameters are suitable, and how parameter choice relates to region length. Here, I examine the performance of SEG and fLPS across the whole expanse of their parameter spaces. In doing so, boundary behaviours are discovered that are used to derive an optimized strategy for annotation of LCRs or CBRs of given target lengths.

## Methods

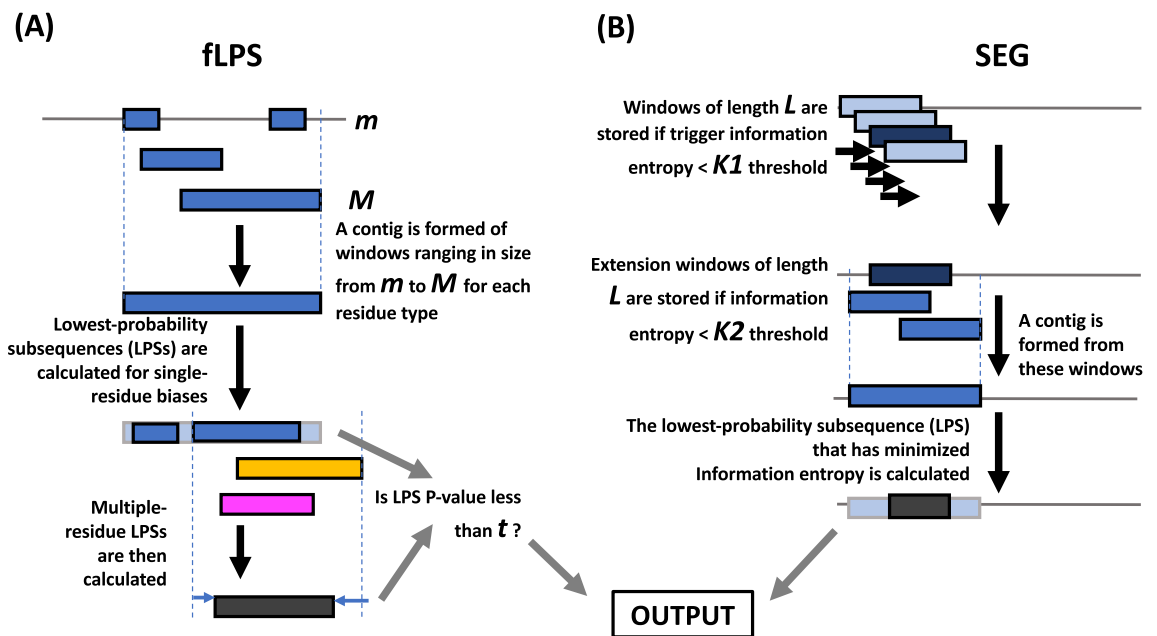
### Data

Two sequence sets were downloaded from UniProt (uniprot.org)<sup>27</sup> in June 2022. These are: (i) *Saccharomyces cerevisiae* strain S288C proteome (number UP000002311, 5879 sequences); (ii) UniRef50 representative protein set. The latter was reduced to a 0.1% random sample (52,523 sequences, i.e., every 1000th sequence). A data set of water-soluble non-membrane protein domain sequences <40% identical to each other was generated from the ASTRAL sequence data available at SCOPe (scop.berkeley.edu)<sup>28</sup>. Atom record sequences were used to avoid including intrinsically disordered regions that have no electron density in crystallographic data.

A data set of 137 human intrinsically-disordered proteins (IDPs) associated by database curators with cancer was downloaded from the Uniprot database (uniprot.org) in October 2023.

### Running the SEG and fLPS programs

The two programs SEG and fLPS were investigated across a thorough sample of their parameter spaces. Both programs use three main parameters for CBR discovery, but extract fundamentally different information from the sequences. Figure 1 details both algorithms. The fLPS program (Fig. 1A) works through a process of binomial probability minimization<sup>14</sup>. A binomial P-value can be calculated for the amino-acid biases of any sequence tract. SEG uses sequence entropy to search for CBRs/LCRs (Fig. 1B).



**Figure 1.** Schematics of the algorithms. **(A)** fLPS algorithm. The three main parameters are maximum and minimum window lengths  $M$  and  $m$ , and an output P-value threshold  $t$ . Window lengths down from  $M$  to  $m$  are searched for single-residue biases with P-values less than a fixed high threshold ( $=0.001$  here). These are then used to form contigs out of which lowest-probability subsequences (LPSs) are calculated (multiple LPSs from the same contig are possible). Then, multiple-residue biased regions are tested for, which includes trimming or extending to obtain the multiple-residue LPS. Finally, the output is filtered with  $t$ , the threshold P-value. **(B)** SEG algorithm. Regions with low sequence entropy are more ‘ordered’ since they are dominated by a few of the possible amino-acid residue types. SEG works by scanning along sequences for windows of length  $L$  that have sequence entropy  $\leq K1$ , a trigger threshold. Then, these ‘trigger windows’ are extended with further windows whose sequence entropy is  $\leq K2$ , the extension threshold, to form a contig. SEG LPSs are then calculated from these contigs using recursion.

The parameter spaces searched are tabulated (Table 1). For fLPS, default background ‘domains’ amino-acid frequencies were used. These were updated using the downloaded ASTRAL sequences (however, the frequencies have changed minimally indicating likely convergence). For both programs, there are recommended parameter sets to discover shorter low-complexity regions (Table 1). Also, there are SEG parameter sets to label longer biased regions, such as those made of longer tandem repeats, whose repetitiveness is only captured by longer L windows. The default fLPS settings are intended as a ‘catch-all’ set of very loose parameters that informs users of all CBRs in their sequences, even very mild short tracts. After a default run, the user is expected to home in on biases of interest with smaller  $t$  P-value thresholds. For example,  $t = 1e-10$  was used to study prion-like proteins, since  $1e-10$  was the highest P-value observed for known prion-forming protein tracts (Table 1).

### Metrics for assessing region discovery

To assess and compare parameter set performance, three *metrics* were derived: coverage (*Cov*), median (*Med*), and interquartile range (*IQR*) (Fig. 2).

Coverage (*Cov*) is the proportion of a protein sequence set that is annotated. More liberal parameters lead to greater coverage, and discern more, more mildly-biased regions. *Cov* is calculated taking account of any annotated region overlap. So, residues that appear in multiple annotated regions are counted only once. Median region length (*Med*) and interquartile range (*IQR*) are also calculated from the distribution of region lengths (Fig. 2). *IQR* is an indicator of region length diversity, with smaller *IQR* values for more limited variance in region lengths.

These metrics behave consistently for the two programs and two data sets, Uniref50 and the yeast proteome, with distinct behaviour only for the third data set, the ASTRAL domains, which is to be expected, since it comprises structured regions only (Suppl. Table 1). *IQR/Med* and *Cov* are correlated in all cases, implying that increased coverage comes with an increased range of region sizes. For the UniRef50 data random samples of a third of the size of the sample studied yield highly correlated values for *Cov*, *Med* and *IQR* for both programs (Pearson  $R^2 > 0.99$ ), indicating sufficient sample size. The metrics are also highly correlated between the ASTRAL structural domain set and the UniRef50 sample and yeast proteome (Pearson  $R^2 > 0.94$ ).

### Deriving curves indicating optimal strategies for CBR/LCR annotation

The behaviour of parameter sets was probed using plots of *IQR/Med* versus *Med*, for intervals of *Cov*. The following *Cov* intervals were used since they have approximately equal numbers of points: 0.015–0.025 (~2%), 0.04–0.06 (~5%), 0.08–0.12 (~12%), 0.2–0.3 (~25%), 0.35–0.45 (~40%). Upper and lower boundaries for the point distributions were derived for each plot, with logarithmic equations almost universally best fitting (Fig. 2). Boundary points were defined as extreme relative to all the points above or below them within a margin added around the point along the *Med* axis, with different margins (in the range 3–9) being tried, with 3 discovered as optimal. Because of the characteristic banding on these plots, the margin was skewed to lower values for the lower boundary and higher values for the higher (e.g.,  $x - 5$  to  $x + 1$  for the point  $x$  for the lower boundary). From the average of these two boundary equations, ‘middle curves’ were calculated (Fig. 2). The points nearest these boundary and middle curves (with a tolerance of  $\pm 0.05$  *IQR/Med*) were then analyzed for relationships between metrics and parameters. These relationships were also derived using appropriate line fitting, with characteristic power-law or straight-line relationships between parameters being discovered (discussed in detail in “[Results and discussion](#)”). These relationships are robust to missing points that are partitioned to different plots dependent on the *Cov* intervals examined. These were analysed collectively to derive an optimal annotation protocol for LCRs/CBRs of a given target length.

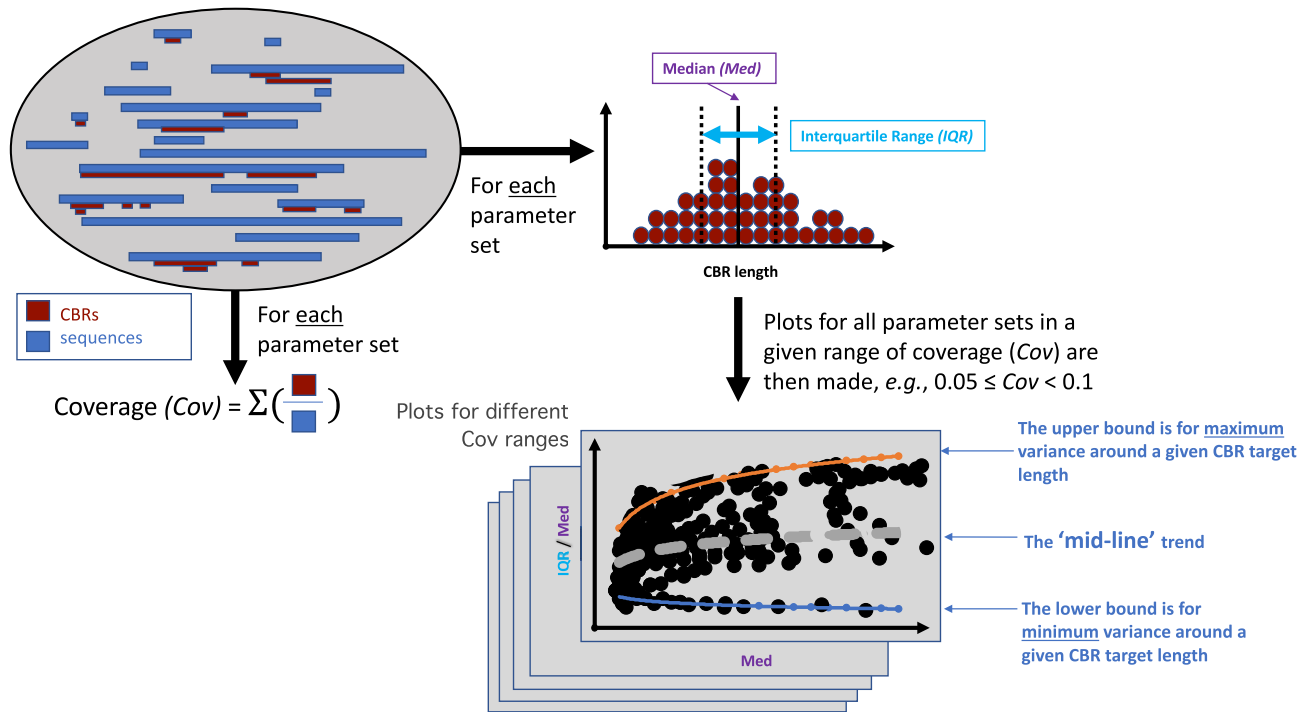
### Further fLPS parameters

There are extra fLPS parameters for: (i) expected amino-acid frequencies ( $-c$  option), (ii) initial search granularity ( $-z$  option)<sup>15</sup>. The  $-c$  option can be: ‘equal’ (=0.05 for each amino acid), ‘domains’ (frequencies from ASTRAL domains<sup>28</sup>), or ‘user’ (from input sequences). *Cov* shows a clear trend for  $-c$ , with all other parameters

	Recommended parameter sets	Parameter space tested
fLPS*	<p>For shorter low-complexity regions:  <math>m = 5, M = 25, t = 1e-05</math>  <math>m = 5, M = 25, t = 1e-06</math></p> <p>To find most shorter low-complexity tracts, and longer compositionally biased regions (default), i.e., ‘catch-all’ default parameters:  <math>m = 15, M = 500, t = 1e-03</math></p> <p>Parameters used for identifying prion-like compositional biases or biased ‘dark matter’: †  <math>m = 15, M = 500, t = 1e-10</math>  <math>m = 15, M = 100, t = 1e-10</math></p>	$5 \leq m \leq 100,$ $5 \leq M \leq 1000,$ $1e-03 \geq t \geq 1e-12$
SEG**	<p>For shorter low-complexity regions (default):  <math>L = 12, K1 = 2.2, K2 = 2.5</math></p> <p>For longer compositionally biased regions, such as those made of longer tandem repeats:  <math>L = 25, K1 = 3.0, K2 = 3.3</math>  <math>L = 45, K1 = 3.4, K2 = 3.75</math></p>	$6 \leq L \leq 250,$ $0.2 \leq K1 \leq 4.2,$ $0.2 \leq K2 \leq 4.2$

**Table 1.** Parameter values analyzed for SEG and fLPS. \*The minimum window  $m \leq$  maximum window  $M$ .

\*\*The trigger information entropy threshold  $K1$  is  $\leq K2$  for all runs. The maximum value for  $K1$  and  $K2$  in amino-acid sequences is  $\log_2(20) = 4.3$ . The recommended parameter sets are taken from the original reference for SEG<sup>5</sup>. †Refs.<sup>2,22,23,25,26</sup>.



**Figure 2.** Analyzing the algorithm parameter spaces. A set of protein sequences is analyzed for each parameter set to extract three metrics: (i) the coverage ( $Cov$ ), which is the proportion of the protein data set annotated by the algorithm; (ii) the median length of annotated regions ( $Med$ ); (iii) the interquartile range of the distribution of regions lengths ( $IQR$ ). Plots of  $IQR/Med$  versus  $Med$  are derived for intervals of  $Cov$ , and upper and lower boundary curves are fitted, then average 'mid-line' trends calculated. The parameter sets that yield the lower and mid-line bounds are extracted, as described in *Methods*.

set equal: >99% of the time, 'user' yields greatest coverage, then 'domains', then 'equal'. Option  $-z$  can be: 'fast' (default initial upper P-value =  $1e-03$ ), 'medium' (=0.01), and 'thorough' (=0.01). With higher values, biased regions made from longer lists of amino-acid types are detected. Indeed, with other parameters set equal, the mean number of residue types defining CBRs increases from 2.7 (fast) to 3.5 (medium) to 4.3 (thorough), with  $Cov$  trending similarly (83% of the time thorough > medium > fast).

## Results and discussion

### How do 'recommended' parameter sets perform?

What is a low-complexity region (LCR)? What is a compositionally biased region (CBR)? These questions are typically answered by applying the recommended or default parameter sets of programs that annotate them. For example, often LCRs are defined through default application of SEG simply because researchers have always tended to define them that way. However, these default parameters have been quite arbitrarily chosen. Indeed, LCRs and CBRs exist on a spectrum of compositional bias, with LCRs generally shorter and more repetitive, but some cases may also be long (Fig. 3). The most extreme LCRs are, of course, homopeptides<sup>29</sup>.

Also, different parameters can yield widely differing results. Some 'recommended' or default parameter sets tend to annotate longer regions, others shorter ones (higher or lower median ( $Med$ ) values in Table 2). Coverage ( $Cov$ ) values vary widely, with default SEG annotating ~9% of proteins, but with alternative SEG parameters for longer regions covering >3 times as many residues. fLPS parameter sets demonstrate a corresponding range of  $Med$  and  $Cov$ . Default 'catch-all' fLPS parameters yield high coverage (>60%), since they are designed to comprehensively capture regions with a compositional perturbation; for these parameters, any remaining unannotated sequence regions can be considered 'high-complexity'. In general, annotations made by fLPS have a greater diversity of lengths (wider  $IQR$ ), than those made by SEG.

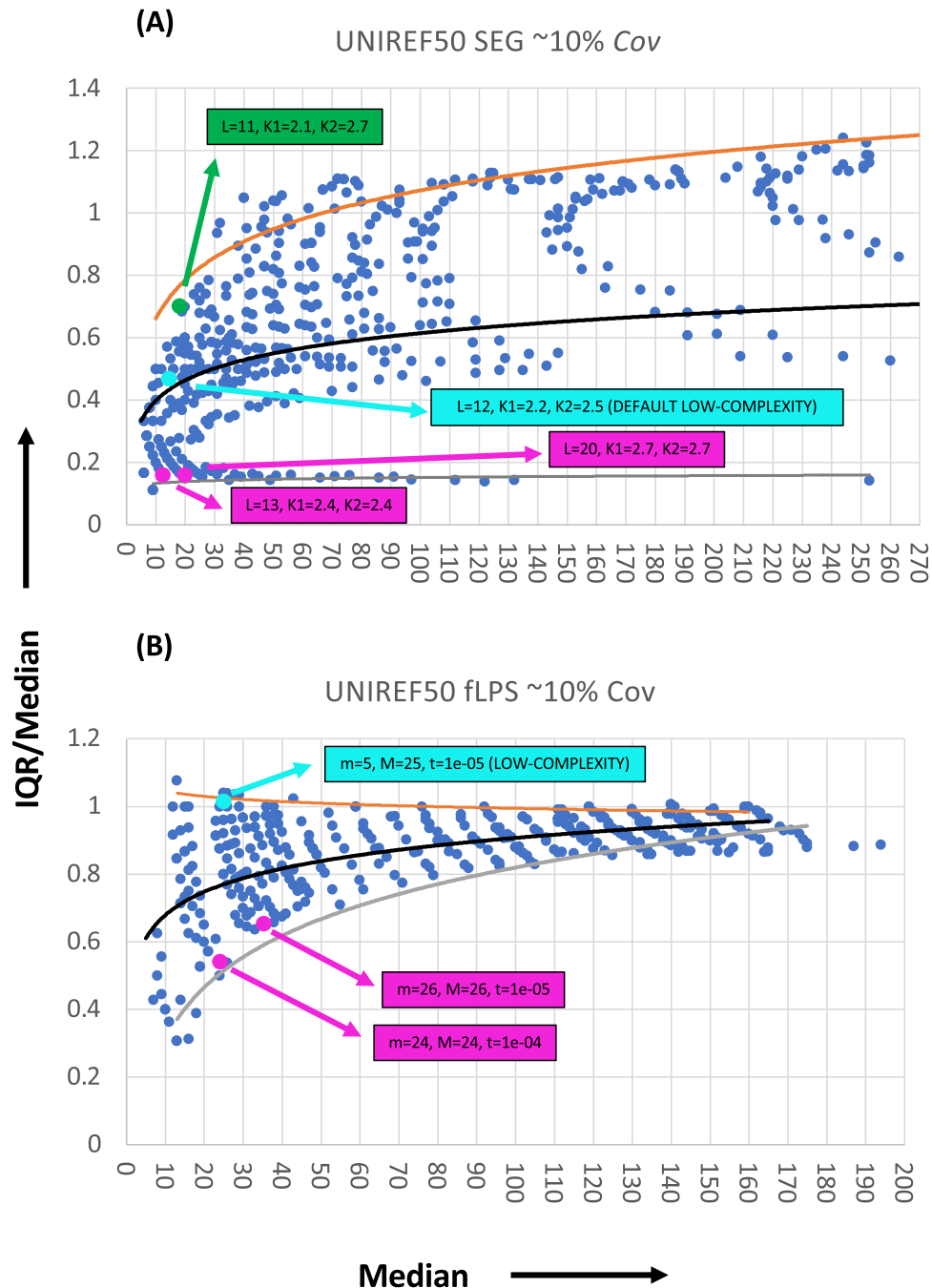
The annotations made by Lee et al. are derived from an image-processing algorithm applied to dot plots<sup>17</sup>. These annotations for budding yeast are short with very low length diversity ( $Med = 13$ ;  $IQR/Median = 0.385$ ), and have very low coverage = 2.3%.

Low-complexity or compositionally-biased sequence in structured protein domains is clearly rarer and less diverse lengthwise, regardless of the parameters chosen (Table 1). The ASTRAL set stands out as always having lower  $IQR/Med$  values, and having much lower  $Cov$  values generally. Thus, sequence complexity is higher at every resolution in the structured parts of proteins.

### The recommended parameters located in parameter space

How do these default or recommended parameter sets compare to the rest of their parameter spaces? How are these parameter sets *special*? To gain answers to these questions, the plots of  $IQR/Med$  versus  $Med$  containing the



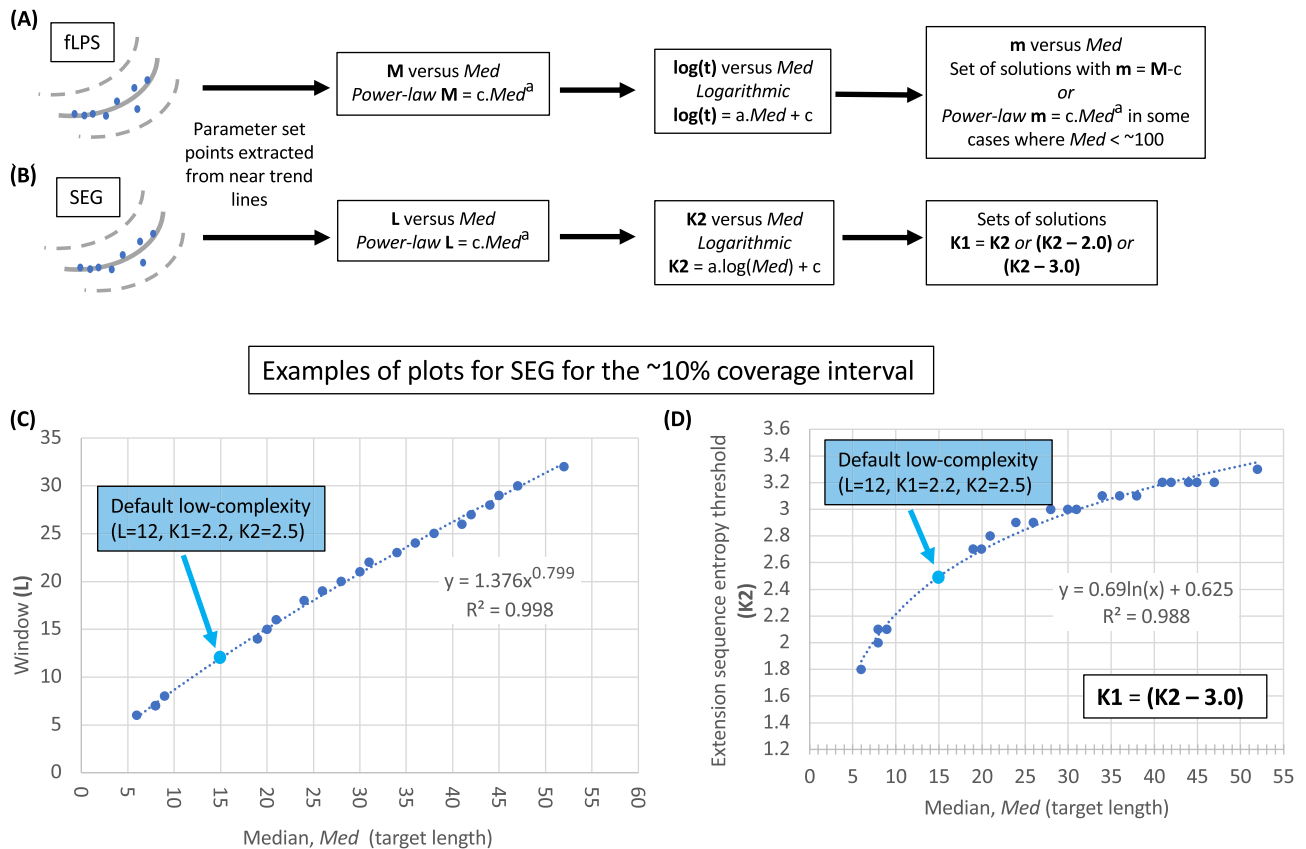


**Figure 4.** Recommended parameters located within parameter space. Plots of  $IQR/Med$  versus  $Med$  for the  $Cov$  interval 0.08–0.12 (~10% coverage) for: (A) default SEG parameters, and (B) fLPS parameters that have been recommended for annotating shorter, low-complexity regions. The Uniref50 protein sequence data were used. The recommended parameters are labelled in light blue. Examples of parameter sets that are below these and close to the lower boundary are labelled in magenta, while those above and close to the upper boundary are in green.

#### Discovery of trends in program parameters

Depending on their specific proteins of interest and their biological contexts, researchers may have an interest in regions of a particular *target length*, e.g., very short regions that may be motifs, or much longer ones, such as the prion-like regions<sup>22,23,30</sup>. The  $Med$  value of LCRs/CBRs calculated here can be considered such a target length.

What program parameters are producing the trend lines in Fig. 4, and how is this related to a specific target length? To answer this, the points near these lines were extracted and the relationship between  $Med$  and program parameters discerned (as described in *Methods*) (Fig. 5).



**Figure 5.** Discovering trends in program parameters. Points near the trendlines in the *IQR/Med* versus *Med* plots are extracted and examined. For both programs (A) fLPS and (B) SEG, there are power-law relationships between window size and *Med*, and logarithmic relationships with *Med* for a second parameter that measures the degree of bias. Different sets of solutions for the third parameter (*m* for fLPS and *K1* for SEG), then arise in the data. In (C) and (D), there are examples of plots of *Med* versus *L* and *Med* versus *K2* for a coverage interval of ~10% (*Cov* 0.08–0.12).

For both programs there are standard relationships between this *Med* target length and chosen parameters for the mid-line trends. All correlations for these trends are significant ( $P < 0.0001$ ). For example, there are power-law relationships between window lengths (*M* for fLPS and *L* for SEG) and *Med* target lengths, with Pearson  $R^2 > 0.92$  for fLPS and  $> 0.99$  for SEG for all *Cov* ranges studied (Fig. 5A,B). An example of this is shown for SEG and the ~10% coverage interval (Fig. 5C). This power-law may thus be an inherent property of such window-based algorithms. Also, logarithmic relationships are standard for the second parameter that determines bias levels (namely the thresholds *K2* for SEG with  $R^2$  values  $> 0.95$  and *t* for fLPS with  $R^2$  values  $> 0.92$ ). An example is shown in Fig. 5D.

The lower-bound trends yield the narrowest focus around a particular target length. The solution of these trends is very simple. For fLPS,  $m = M$  and *M* and  $\log(t)$  are proportional to *Med* and to each other ( $R^2$  values  $> 0.89$ ). For SEG,  $K1 = K2$ , and *K2* and the window length *L* are proportional to *Med*, and to each other ( $R^2$  values  $> 0.92$ ).

However, parameter sets extracted from near the upper bound of these plots for SEG demonstrate that it is not meaningful to consider these solutions. They only exhibit a logarithmic correlation between *Med* and *K2* (e.g.,  $R^2 = 0.8$  for the ~10% coverage interval) with both *L* and *K1* being un-correlated with *Med*, and there no subsets with regular patterning of *K1*–*K2*, as are observed for the other trends (Fig. 5B). This implies that the value of the *K1* trigger threshold is arbitrary, and the annotations thus have no biological meaning, since they are an arbitrary subset of what is possible. This highlights a key feature of the SEG algorithm to be aware of, in that the LCRs/CBRs it annotates are all based on a core region that can have lower sequence entropy than the rest of the LCR. For the fLPS algorithm, these upper bounds are meaningful, with maximum IQR given by *m* set equal to the lowest value 5 and *M* and *t* correlated with target length (*Med*) (e.g., *Med* vs. *M*,  $R^2 = 0.98$ ; *Med* vs.  $\log(t)$ ,  $R^2 = 0.97$  for the 10% *Cov* interval). However, in general, for simplicity, we decided not to use upper-bound trends in these analyses.

### Software

Figure 6 shows how program parameters scale as the implied coverage of annotations increases. Parameters for shorter and longer LCRs/CBRs clearly behave very differently. Also, we can see examples of parameter settings

SEG		~ COVERAGE	fLPS	
L, K1, K2			m, M, t	
SHORT	LONG		SHORT	LONG
(15 residues)	(150 residues)		(15 residues)	(150 residues)
12, 1.75, 2.05	87, 2.98, 3.28	2%	8, 10, 5.1e-07	N/A
12, 2.02, 2.32	72, 3.17, 3.57	5%	10, 14, 1.7e-05	40, 44, 9.9e-19
12, 2.19, 2.49	75, 3.37, 3.67	10%	7, 17, 3.5e-05	49, 59, 1.3e-12
12, 2.56, 2.86	69, 3.49, 3.79	25%	N/A	47, 97, 2.6e-08
13, 2.69, 2.89	85, 3.64, 3.84	40%	N/A	46, 126, 3.1e-07

**Figure 6.** Optimal parameters for a given target length. Considering the median *Med* to be a target length, parameters can be selected to aim at this target. These are listed for a ‘short’ target length (15 residues) and a longer one (150 residues). They are taken from the ‘mid-line’ analysis, which is termed a ‘diverse’ focus in the distributed software *fLPSparameters* and *SEGparameters*. The approximate coverage expected for each parameter set is listed in the middle.

that yield short median lengths and higher coverage, or vice versa. The default SEG parameters re-emerge in the table (blue highlights); equivalent fLPS parameters are also highlighted.

A pair of programs **fLPSparameters** and **SEGparameters** to choose parameters to perform such analysis are available at: <https://github.com/pmharrison/parameters>. They allow for either a ‘diverse’ focus (the mid-line trend) or a ‘narrow’ focus (the lower boundary trend), *i.e.*, with the least possible diversity of region lengths. The discussion below just uses the default ‘diverse’ focus, for greater simplicity.

### An optimized strategy: progressive parsing of CBRs/LCRs annotated across multiple target lengths

So, what is an optimized strategy for annotating regions of a given target length? The best answer is to examine the results from all these progressive program runs (Fig. 6), and assess the biological relevance at each stage. Such an approach may be productive for large-scale bioinformatical analyses involving cross-referencing with other information about function. It may also be useful in directing the parsing of domains to fashion experimental constructs. Thus, what is a meaningfully defined LCR/CBR is determined by such progressive analysis within the relevant biological context for a specific protein under experimental study.

An example of applying these scaled parameter sets is shown (Fig. 7). Human prion protein PrP (UniProt accession P04156) is dissected with parameters for ‘short’ and ‘long’ target lengths (15 and 150 residues). PrP underlies mammalian prion diseases, through amyloid formation, and functions in copper metabolism and circadian control<sup>31</sup>. For fLPS, the ‘long’ parameters annotate the protein’s repetitive copper-binding tract, which converges to a maximum length for parameters with estimated coverage > ~ 10%. The same tract is also found by ‘long’ SEG parameters, but it is lengthened to include the A/G-rich tract that is transmembrane in some PrP isoforms, and is implicated in conversion to amyloid<sup>32,33</sup>. This region is annotated separately by fLPS for the ‘short’ parameters, along with two other tracts that may be biologically significant. The ‘short’ SEG analysis evidences a slow, gradual filling-in of the whole sequence as estimated coverage rises (Fig. 7).

Further examples of this progressive CBR parsing show that some CBRs are detected regardless of target length, but some are only detected with short target length (Suppl. Fig. 2). For a human collagen (Suppl. Fig. 2B), short proline-rich tracts appear at target length = 15, which then expand into longer regions as estimated coverage is increased. In *Saccharomyces cerevisiae* MRN1 RNA-binding protein (Suppl. Fig. 2C), a putative intrinsically-disordered region up to residue 195 (predicted by AlphaFold, in UniProt entry Q08925) is parsed into distinct sub-regions. Arguably, in general, such results based on parsing CBRs according to compositional biases, rather than sequence entropy, are more meaningful biologically, since sequence entropy per se is less likely to be under selection than, say, a bias for glutamine or glutamate residues linked to specific functional roles.

### Application to searching for CBRs of a given target length

Parameter choice focused on CBR target length can also be used to pick out a data set of CBRs with a similar bias. This is illustrated for the M domain of Sup35 protein from *S. cerevisiae*, a domain which mediates pH sensing during reversible condensate formation in response to stress<sup>34</sup>. A ~ 90-residue {KE}-rich CBR that corresponds to the M-domain was discovered using intermediate target lengths and > 5% estimated coverage (Suppl. Fig. 3A,B);



## Human Prion Protein PrP

fLPS 'short'

(15 residues target length)

≥2% coverage [m=8, M=10, t=5.1e-07]

Add in at ≥5% coverage [m=10, M=14, t=1.7e-05]

MANLGCWMLVLFVATWSDLGLCKKRPKPKGGWNTGGSRYPGQGSPGGNRYPPQGGGGWGQPHGGG  
 WGQPHGGGGWGQPHGGGGWGQPHGGGGWGQGGGTHSQWNKPSKPKTNMKHMAGAAAAGAVVGGGLGGY  
 MLGSAMSRPIIHFGSDYEDRYRENMHRYPNQVYYRPMDEYSNQNNFVHDCVNIITIKQHTVTTT  
 TKGENFTETDVKMMERVVEQMCITQYERESQAYYQRGSSMVLFSPPVILLISFLIFLIVG

fLPS 'long'

(150 residues target length)

≥2% coverage [m=40, M=44, t=9.9e-19]

Add in at ≥10% coverage [m=49, M=59, t=1.3e-12]

MANLGCWMLVLFVATWSDLGLCKKRPKGGWNTGGSRYPGQGSPGGNRYPPQGGGGWGQPHGGG  
 WGQPHGGGGWGQPHGGGGWGQPHGGGGWGQGGGTHSQWNKPSKPKTNMKHMAGAAAAGAVVGGGLGGY  
 MLGSAMSRPIIHFGSDYEDRYRENMHRYPNQVYYRPMDEYSNQNNFVHDCVNIITIKQHTVTTT  
 TKGENFTETDVKMMERVVEQMCITQYERESQAYYQRGSSMVLFSPPVILLISFLIFLIVG

SEG 'short'

(15 residues target length)

≥2% coverage [L=12, K1=1.75, K2=2.05]

Adds in at ≥5% coverage [L=12, K1=2.02, K2=2.32]

MANLGCWMLVLFVATWSDLGLCKKRPKGGWNTGGSRYPGQGSPGGNRYPPQGGGGWGQPHGGG  
 WGQPHGGGGWGQPHGGGGWGQPHGGGGWGQGGGTHSQWNKPSKPKTNMKHMAGAAAAGAVVGGGLGGY  
 MLGSAMSRPIIHFGSDYEDRYRENMHRYPNQVYYRPMDEYSNQNNFVHDCVNIITIKQHTVTTT  
 TKGENFTETDVKMMERVVEQMCITQYERESQAYYQRGSSMVLFSPPVILLISFLIFLIVG

Add in at ≥10% coverage [L=12, K1=2.19, K2=2.49]

Add in at ≥25% coverage [L=12, K1=2.56, K2=2.86]

Add in at ≥40% coverage [L=13, K1=2.69, K2=2.89]

SEG 'long'

(150 residues target length)

≥2% coverage [L=87, K1=2.98, K2=3.28]

Adds in at ≥10% coverage [L=72, K1=3.27, K2=3.57]

MANLGCWMLVLFVATWSDLGLCKKRPKGGWNTGGSRYPGQGSPGGNRYPPQGGGGWGQPHGGG  
 WGQPHGGGGWGQPHGGGGWGQPHGGGGWGQGGGTHSQWNKPSKPKTNMKHMAGAAAAGAVVGGGLGGY  
 MLGSAMSRPIIHFGSDYEDRYRENMHRYPNQVYYRPMDEYSNQNNFVHDCVNIITIKQHTVTTT  
 TKGENFTETDVKMMERVVEQMCITQYERESQAYYQRGSSMVLFSPPVILLISFLIFLIVG

**Figure 7.** How fLPS and SEG can parse a protein progressively using short and long target lengths. The examples of 'short' (15-residue) and 'long' (150-residue) target lengths from Fig. 6 are employed on the human prion protein PrP as an example, using a 'diverse' focus. The sequence parts that add in at each level are coloured progressively. The definitions of LCRs from Lee et al.<sup>17</sup> are underlined.

the yeast proteome was then scanned for {KE}/{EK}-rich CBRs of target length = 90 with estimated coverage of 5% (Suppl. Fig. 3C). Significant Gene Ontology category enrichments were observed for these CBRs that are linked to rRNA and ribosomal processing (Suppl. Fig. 3C). Interestingly the Lee et al. annotations do not contain

this biased M domain, nor most of the repetitive prion-forming domain (Suppl. Fig. 3A), and only cover ~ 36% of the extent of the {KE}/{EK}-rich regions analyzed here.

### Application to checking parameter dependence of annotations in a large data set

Examining CBR prevalences across a wide range of target lengths can be used to pick out CBR types that are prevalent regardless of parameter choice, and to home in on regions that are only detectable with shorter or longer target lengths, or with lower or higher estimated coverage. This is illustrated for a data set of 137 cancer-associated human intrinsically disordered proteins (IDPs). The top five bias signatures for CBRs are listed for sets of fLPS parameters for target lengths ranging from 10 to 200, and estimated coverage between 2 and 40%.

We can see that some bias signatures are prevalent regardless of parameter choice, e.g., {P}, whereas others are only detected at lower target lengths, e.g., {R} (Suppl. Fig. 4A). {K}-rich regions are only numerous at higher coverage values, i.e., many of them are more mildly-biased. The {R}-rich regions detected are linked to Gene Ontology categories such as protein kinase activity and adenylyl nucleotide binding, with {P}-rich regions being associated with enzyme binding and  $\beta$ -catenin binding (Suppl. Fig. 4B). The {R}-rich regions are only detected as prominent, if lower target lengths and medium to high coverage levels are applied.

### Conclusions

An optimized strategy for discovering LCRs/CBRs of a given target length has been derived. Such an approach is suitable for large-scale bioinformatical analyses, for fishing out similar regions of similar length, and in guiding experimental hypotheses about functionally significant protein tracts. We saw how the CBR annotation problem could generally be simplified by choosing specific boundary or mid-line trends. Also, clear highly correlated power-law and logarithmic relationships between target lengths and program parameters were discovered, that were indicated to be general features of window-based algorithms such as fLPS and SEG. This sort of analysis could be combined with application of other tools that can further dissect the character of the CBRs that are discovered, e.g., the LCR server for visualizing repetitiveness<sup>20</sup>, or LCD-Composer, which can assess the dispersion of residue types within a CBR<sup>13</sup>.

The results here are of utility for the development of further improved algorithms for characterization of LCRs and CBRs, and for informing the combination of different algorithms to provide insights into biologically relevant features.

### Data availability

The data sets analyzed here are available from the public databases UniProt (<https://ftp.uniprot.org/pub/databases/uniprot/uniref/uniref50/> and <https://www.uniprot.org/proteomes/UP000002311>) and SCOPe (<http://scop.berkeley.edu/astral/>). Results of the research are available in code on GitHub at: <https://github.com/pmharrison/parameters>.

Received: 25 August 2023; Accepted: 28 December 2023

Published online: 05 January 2024

### References

- Harrison, P. M. Exhaustive assignment of compositional bias reveals universally prevalent biased regions: Analysis of functional associations in human and *Drosophila*. *BMC Bioinform.* **7**, 441. <https://doi.org/10.1186/1471-2105-7-441> (2006).
- Harrison, P. M. Compositionally biased dark matter in the protein universe. *Proteomics* **18**, e1800069. <https://doi.org/10.1002/pmic.201800069> (2018).
- Hancock, J. M. & Armstrong, J. S. SIMPLE34: An improved and enhanced implementation for VAX and Sun computers of the SIMPLE algorithm for analysis of clustered repetitive motifs in nucleotide sequences. *Comput. Appl. Biosci.* **10**, 67–70. <https://doi.org/10.1093/bioinformatics/10.1.67> (1994).
- Alba, M. M., Laskowski, R. A. & Hancock, J. M. Detecting cryptically simple protein sequences using the SIMPLE algorithm. *Bioinformatics* **18**, 672–678. <https://doi.org/10.1093/bioinformatics/18.5.672> (2002).
- Wootton, J. C. & Federhen, S. Analysis of compositionally biased regions in sequence databases. *Methods Enzymol.* **266**, 554–571 (1996).
- Promponas, V. J. *et al.* CAST: An iterative algorithm for the complexity analysis of sequence tracts Complexity analysis of sequence tracts. *Bioinformatics* **16**, 915–922 (2000).
- Wise, M. J. 0j.py: A software tool for low complexity proteins and protein domains. *Bioinformatics* **17**(Suppl 1), S288–295. [https://doi.org/10.1093/bioinformatics/17.suppl\\_1.s288](https://doi.org/10.1093/bioinformatics/17.suppl_1.s288) (2001).
- Nandi, T. *et al.* A novel complexity measure for comparative analysis of protein sequences from complete genomes. *J. Biomol. Struct. Dyn.* **20**, 657–668. <https://doi.org/10.1080/07391102.2003.10506882> (2003).
- Shin, S. W. & Kim, S. M. A new algorithm for detecting low-complexity regions in protein sequences. *Bioinformatics* **21**, 160–170. <https://doi.org/10.1093/bioinformatics/bth497> (2005).
- Kuznetsov, I. B. & Hwang, S. A novel sensitive method for the detection of user-defined compositional bias in biological sequences. *Bioinformatics* **22**, 1055–1063. <https://doi.org/10.1093/bioinformatics/btl049> (2006).
- Antonets, K. S. & Nizhnikov, A. A. SARP: A novel algorithm to assess compositional biases in protein sequences. *Evol. Bioinform. Online* **9**, 263–273. <https://doi.org/10.4137/EBO.S12299> (2013).
- Cascarina, S. M., King, D. C., Osborne Nishimura, E. & Ross, E. D. LCD-Composer: An intuitive, composition-centric method enabling the identification and detailed functional mapping of low-complexity domains. *NAR Genom. Bioinform.* **3**, lqab048. <https://doi.org/10.1093/nargab/lqab048> (2021).
- Cascarina, S. M. & Ross, E. D. The LCD-Composer webserver: High-specificity identification and functional analysis of low-complexity domains in proteins. *Bioinformatics* **38**, 5446–5448. <https://doi.org/10.1093/bioinformatics/btac699> (2022).
- Harrison, P. M. fLPS: Fast discovery of compositional biases for the protein universe. *BMC Bioinform.* **18**, 476. <https://doi.org/10.1186/s12859-017-1906-3> (2017).
- Harrison, P. M. fLPS 2.0: Rapid annotation of compositionally-biased regions in biological sequences. *PeerJ* **9**, e12363. <https://doi.org/10.7717/peerj.12363> (2021).

16. Harrison, P. M. & Gerstein, M. A method to assess compositional bias in biological sequences and its application to prion-like glutamine/asparagine-rich domains in eukaryotic proteomes. *Genome Biol.* **4**, R40. <https://doi.org/10.1186/gb-2003-4-6-r40> (2003).
17. Lee, B., Jaber-Lashkari, N. & Calo, E. A unified view of low complexity regions (LCRs) across species. *Elife* **11**, 77058. <https://doi.org/10.7554/eLife.77058> (2022).
18. Kirmitzoglou, I. & Promponas, V. J. LCR-eXXXplorer: A web platform to search, visualize and share data for low complexity regions in protein sequences. *Bioinformatics* **31**, 2208–2210. <https://doi.org/10.1093/bioinformatics/btv115> (2015).
19. Jarnot, P. *et al.* PlaToLoCo: The first web meta-server for visualization and annotation of low complexity regions in proteins. *Nucleic Acids Res.* **48**, W77–W84. <https://doi.org/10.1093/nar/gkaa339> (2020).
20. Mier, P. & Andrade-Navarro, M. A. Assessing the low complexity of protein sequences via the low complexity triangle. *PLoS ONE* **15**, e0239154. <https://doi.org/10.1371/journal.pone.0239154> (2020).
21. Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: A new generation of protein database search programs. *Nucleic Acids Res.* **25**, 3389–3402 (1997).
22. An, L., Fitzpatrick, D. & Harrison, P. M. Emergence and evolution of yeast prion and prion-like proteins. *BMC Evol. Biol.* **16**, 24. <https://doi.org/10.1186/s12862-016-0594-3> (2016).
23. An, L. & Harrison, P. M. The evolutionary scope and neurological disease linkage of yeast-prion-like proteins in humans. *Biol. Direct* **11**, 32. <https://doi.org/10.1186/s13062-016-0134-5> (2016).
24. Su, T. Y. & Harrison, P. M. Conservation of prion-like composition and sequence in prion-formers and prion-like proteins of *Saccharomyces cerevisiae*. *Front. Mol. Biosci.* **6**, 54. <https://doi.org/10.3389/fmolb.2019.00054> (2019).
25. Su, W. C. & Harrison, P. M. Deep conservation of prion-like composition in the eukaryotic prion-former Pub1/Tia1 family and its relatives. *PeerJ* **8**, e9023. <https://doi.org/10.7717/peerj.9023> (2020).
26. Luo, J. & Harrison, P. M. Evolution of sequence traits of prion-like proteins linked to amyotrophic lateral sclerosis (ALS). *PeerJ* **10**, e14417. <https://doi.org/10.7717/peerj.14417> (2022).
27. UniProt, C. UniProt: The universal protein knowledgebase in 2023. *Nucleic Acids Res.* **51**, D523–D531. <https://doi.org/10.1093/nar/gkac1052> (2023).
28. Fox, N. K., Brenner, S. E. & Chandonia, J. M. SCOPe: Structural Classification of Proteins—extended, integrating SCOP and ASTRAL data and classification of new structures. *Nucleic Acids Res.* **42**, D304–309. <https://doi.org/10.1093/nar/gkt1240> (2014).
29. Wang, Y. & Harrison, P. M. Homopeptide and homocodon levels across fungi are coupled to GC/AT-bias and intrinsic disorder, with unique behaviours for some amino acids. *Sci. Rep.* **11**, 10025. <https://doi.org/10.1038/s41598-021-89650-1> (2021).
30. Harbi, D. & Harrison, P. M. Classifying prion and prion-like phenomena. *Prion* **8**, 161–165 (2014).
31. Castle, A. R. & Gill, A. C. Physiological functions of the cellular prion protein. *Front. Mol. Biosci.* **4**, 19. <https://doi.org/10.3389/fmolb.2017.00019> (2017).
32. Harrison, P. M., Bamborough, P., Daggett, V., Prusiner, S. B. & Cohen, F. E. The prion folding problem. *Curr. Opin. Struct. Biol.* **7**, 53–59. [https://doi.org/10.1016/s0959-440x\(97\)80007-3](https://doi.org/10.1016/s0959-440x(97)80007-3) (1997).
33. Hegde, R. S. *et al.* A transmembrane form of the prion protein in neurodegenerative disease. *Science* **279**, 827–834. <https://doi.org/10.1126/science.279.5352.827> (1998).
34. Franzmann, T. M. & Alberti, S. Protein phase separation as a stress survival strategy. *Cold Spring Harb. Perspect. Biol.* **11**, a034058. <https://doi.org/10.1101/cshperspect.a034058> (2019).

## Acknowledgements

This work was performed on computing resources funded by the Natural Sciences and Engineering Research Council of Canada.

## Author contributions

P.H. performed the research and wrote the manuscript.

## Competing interests

The author declares no competing interests.

## Additional information

**Supplementary Information** The online version contains supplementary material available at <https://doi.org/10.1038/s41598-023-50991-8>.

**Correspondence** and requests for materials should be addressed to P.M.H.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

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



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

© The Author(s) 2024