

# Abundance and Distributions of Eukaryote Protein Simple Sequences

Kim Lan Sim and Trevor P. Creamer \*

Center for Structural Biology

Department of Molecular and Cellular Biochemistry

University of Kentucky

800 Rose Street

Lexington, KY 40536-0298

U.S.A.

\* Corresponding author: Phone: (859) 323-6037

Fax: (859) 323-1037

Email: tprea0@uky.edu

Running title: Eukaryote protein simple sequences.

## Abbreviations

AF, *A. fulgidus*; AgT, *A. tumefaciens* C58; AP, *A. pernix* K1; AT, *A. thaliana*; BH, *B. halodurans*; BM, *B. melintensis* 16M chr1; BS, *B. subtilis*; CA, *C. acetobutylicum* ATCC824; CE, *C. elegans*; DM, *D. melanogaster*; DR, *D. radiodurans* chr1; EC, *E. coli* K-12; HI, *H. influenzae*; HP, *H. pylori* 26695; HS, *Halobacterium* sp. NRC-1; MG, *M. genitalium*; MJ, *M. janaschii*; MP, *M. pneumoniae*; MT, *M. thermoautotrophicum*; Nos, *Nostoc* sp. PCC7120; PA, *P. abyssi*; PAe, *P. aerophilum*; PH, *P. horikoshii*; SC, *S. cerevisiae*; SS, *Synechocytis* sp. PCC6803; SSol, *S. solfataricus*; ST, *S. tokodaii*; TA, *T. acidophilum*; TV, *T. volcanium*; VC, *V. cholerae* chr1.

## Summary

Protein simple sequences are a subclass of low-complexity regions of sequence that are highly enriched in one or a few residue types. Such sequences are common in transcription regulatory proteins, structural proteins, proteins involved in nucleic acid interactions, and in mediating protein-protein interactions. Simple sequences of ten or more residues, containing  $\geq 50\%$  of a single residue type are surveyed in this work. Both eukaryote and prokaryote proteomes are investigated, with emphasis on the eukaryotes. Very large numbers of such sequences are found in all organisms surveyed. It is found that eukaryotes possess far more simple sequences per protein than prokaryotes. Prokaryotes display a linear relationship between number of proteins containing simple sequences and proteome size, whereas it is not clear that such a relationship holds for eukaryotes. Strikingly, it is found that each eukaryote possesses its own unique distribution of simple sequences. Within those distributions it is found that simple sequences enriched in certain residue types are clearly favored, whereas others are just as clearly discriminated against. The preferences observed are not correlated with residue occurrence. An analysis of classes of proteins of known function suggests that simple sequence occurrence and distribution may be related to protein function. Based upon this analysis, the large number of simple sequences found above that which would be expected from a simple statistical model, plus the known functional importance of numerous such sequences, it is postulated that eukaryotes have evolved to not only tolerate large numbers of simple sequences, but also to require them.

## Introduction

Protein simple sequences are stretches of sequence highly enriched in one or a few residue types. These sequences form a major subclass of low-complexity sequences (1). Such sequences are common in transcription regulatory proteins, where they are often enriched in glutamine, proline or charged residues, and tend to be highly conserved (2-5). Glutamine-enriched sequences are thought to be the most common simple sequences (5, 6), and have been associated with a number of human neurological disorders such as Huntington's disease (7-10). Proline-rich sequences are known to have important roles as structural elements and in mediating protein-protein interactions (11, 12). Sequences enriched in charged residues have been associated with DNA and RNA processing, chromatin structure, ion binding and protein-protein interactions (10, 13). Various simple sequences have been implicated as protein domain linkers (14), or as markers for disordered proteins (15, 16). Clearly there are numerous instances where such sequences play important functional roles. In addition, Kashi *et al.* (4) have noted that DNA simple sequences are a potential source of genetic variation. Some of these DNA sequences fall within coding regions, leading to variation at the protein level. The recent explosion in available genomic, and consequently proteomic data, has provided the opportunity to examine the occurrence and distribution of protein simple sequences at a level of detail not previously possible. Here we present a survey of the occurrence and distribution of protein simple sequences highly enriched in a single residue type in the proteomes of four eukaryotes whose genomes have been fully sequenced. The occurrence of eukaryote simple sequences is compared to the occurrence of such sequences in the proteomes of twenty-six prokaryotes.

Some previous studies of protein simple sequences have employed somewhat limited protein databases and have not necessarily compared organisms (5, 6, 17, 18). Other surveys have

considered whole proteomes, but often remove sequences considered redundant (19, 20). There are a number of surveys where simple sequences enriched in a particular residue type or associated with a particular function have been examined (2, 3, 9, 14). Some recent studies have focused on comparisons between organisms (10, 21-23), but have mostly considered only homopolymeric sequences. Our current study differs from prior work in that we employ only intact proteomes from fully-sequenced genomes, including sequences annotated as hypothetical proteins. We focus solely on non-overlapping simple sequences, of ten or more residues in length, highly enriched in a single residue type ( $\geq 50\%$  composition). This approach provides a non-biased view of the distribution of this set of protein simple sequences, as well as allowing for ready comparison of their occurrence in the organisms examined. The eukaryotes surveyed, namely a yeast, worm, fruitfly and a plant, comprise a diverse sample of members of the eukaryote kingdom. We have chosen not to include the human proteome given the current uncertain state of its completion. In addition, for comparison we have surveyed twenty-six prokaryotes, including twelve archaea, two cyanobacteria, and six gram-negative and six gram-positive bacteria.

We find that highly-enriched simple sequences are remarkably common in all of the organisms examined. Eukaryotes are found to possess more simple sequences per protein than prokaryotes, in keeping with the findings of other groups (19, 21, 23). The occurrence of prokaryote proteins containing simple sequences is linearly correlated with proteome size. Given the limited number of organisms examined, it is not clear that this is the case for the eukaryotes. Perhaps most notably, each organism examined possesses its own unique distribution of simple sequences. We find that simple sequences display surprising length dependencies, with some residues preferentially populating long simple sequences regions, while others clearly prefer short

simple sequences. There is no discernible correlation with residue occurrence. For example, leucine-enriched sequences appear to be discriminated against, despite leucine being the most common residue in most organisms. Some observed length dependencies can be explained in structural and functional terms, although many remain enigmatic. We have also found that simple sequence distributions vary according to functional groupings. For example, leucine-rich regions, despite being discriminated against in the overall distributions, are among the most common simple sequences found in membrane-associated proteins. It is clear from the sheer number found that all organisms examined, particularly eukaryotes, tolerate, and perhaps even require, large numbers of protein simple sequences. The data presented here will provide the basis for future studies of these ubiquitous and potentially extremely important sequences.

## Experimental Procedures

Complete proteomes from the fully-sequenced genomes of four eukaryotes and twenty-six prokaryotes were used in our studies (Table 1). Sequences were obtained as FASTA format files from the European Bioinformatics Institute (<http://www.ebi.ac.uk/genomes/>). We employ the entire proteome for each organism, including all proteins marked “hypothetical”, “putative” or “probable”, as well as all proteins that have no annotation. The one exception to this is the proteome of *A. thaliana* (AT), in which 782 of the protein sequences were found to be incomplete (3% of the proteome). We therefore used only the 26,496 complete sequences in the AT proteome.

We arbitrarily define simple sequences as stretches of sequence that:

- 1) are at least ten residues in length,

- 2) are composed of  $\geq 50\%$  of a single type of residue,
- 3) begin and end with the residue of interest,
- 4) and do not possess gaps (runs without residue of interest) of more than five residues in length.

We represent a protein sequence of length  $L$  as a string,  $a_1a_2a_3a_4\dots a_L$ , where  $a_i$  is the residue at position  $i$ . When searching for a simple sequence enriched in a certain residue type, the numerical positions in the protein string for that residue are first generated as a string of  $i$  values. Putative simple sequences are extracted based on the positions of the  $i$  values, given that gaps of six or more residues in length are not allowed within a simple sequence. Putative simple sequences of many lengths are identified, with all  $i$  values corresponding to the residue of interest being output. Since only the residue of interest is selected, the process automatically generates only sequences that begin and end with the residue of interest. Subsequent filtering removes sequences that are less than ten residues long. Remaining sequences are tested to satisfy the  $\geq 50\%$  threshold for the residue of interest. Sequences that do not satisfy the criteria are further analyzed in order to determine whether shorter simple sequences satisfying our criteria are within them. The entire process results in the identification of all non-overlapping simple sequences within the proteomes which satisfy all four of the above criteria. The computer programs used to identify simple sequences were written in Python/C++ and executed on a SGI workstation.

We use the Poisson distribution (9, 24) to model the probability of random occurrence of simple sequences containing a given residue type in the eukaryote proteomes. This is given by

$$f(n) = \frac{e^{-m}m^n}{n!}, \quad [1]$$

where  $f(n)$  is the probability of an event happening  $n$  times. In our studies  $l$  is the length of the simple sequence,  $n$  is the threshold value, and  $m$  is derived from

$$m = \frac{l \times (\% \text{occurrence of residue})}{100}. \quad [2]$$

The expected number of simple sequences of length  $l$  in a proteome is then

$$SS_{expect} = f(n) \times T_l, \quad [3]$$

where  $T_l$  is the total of number of sequence windows of length  $l$  in the proteome.

The difference between the actual number of simple sequences,  $SS_{Tot}$ , of length  $l$  found, and the number expected from the Poisson distribution is then

$$\Delta = SS_{Tot} - SS_{expect}. \quad [4]$$

For simple sequences longer than about twenty-five residues,  $SS_{expect}$  is essentially zero, in which case  $\Delta$  is equal to the number of simple sequences found. Finally, in order to compare the occurrence of simple sequences among organisms, we define  $\Delta_R$  as follows:

$$\Delta_R = \frac{\Delta}{\text{Number of proteins in proteome}}. \quad [5]$$

## Results and Discussion

### ***Simple sequence definition***

Our criteria for identifying protein simple sequences ensures that we find sequences that would satisfy any definition of simple sequences, such as the low-complexity measures of



Wootton and co-workers (1) or the definition employed by Golding (20). We chose to use this definition since it is relatively straightforward to apply, and the sequences identified are unambiguous in nature. The allowable gap used (five or fewer residues) was chosen because this is the largest gap possible in a ten residue sequence, the shortest considered, while still satisfying our  $\geq 50\%$  threshold requirement. The  $\geq 50\%$  threshold ensures that even the shorter sequences identified are relatively unlikely to have occurred as a result of randomness in protein sequences. As will be demonstrated below, for many residues the  $\Delta$  values obtained tend to be large and positive, indicating that we did indeed identify many more sequences than would be expected were sequences random. If the threshold is decreased to  $\geq 30\%$ , we find significantly more simple sequences at all lengths, however many of these, particularly short sequences, are accounted for by the number expected using the Poisson distribution model (data not shown). If the threshold is increased to  $\geq 70\%$  we find relatively few sequences (data not shown).

### ***Inclusion of Potentially Incorrect Protein Sequences***

We have chosen to include all complete protein sequences in the proteomes that we have examined. This includes those marked “hypothetical”, “putative” or “probable”, and those proteins that have not as yet been annotated. Redundant sequences have also been included. This choice was made so as to be able to perform a more complete analysis of the proteomes, leading to a “unbiased” view. It is possible that some of the simple sequences found come from sequences that are not expressed as proteins. Bork and Copley (25) have pointed out that the identification of genes in sequenced genomes is difficult. It is particularly difficult for eukaryote genes where the identification of exons is error-prone. Ideally the analyses presented below should be repeated leaving out those proteins marked hypothetical or not annotated. This is however extremely difficult due to the wide variety of annotations used to denote such putative

protein sequences. We have thus chosen to present the analyses of the complete proteomes with the caveat that some of the results may be slightly skewed by the presence of incorrect protein sequences.

### ***Abundance of Protein Simple Sequences***

All of the organisms surveyed possess a remarkable number of simple sequences in their proteomes (Table 1). The number found ranges from 251 in the small proteome of MG (480 proteins), up to 27,542 in the proteome of AT (26,496 protein sequences surveyed). Furthermore, a remarkable fraction of proteins in each proteome possess at least one simple sequence. Figure 1a is a plot of the number of proteins possessing one or more simple sequences,  $Prot_{SS}$ , against the number of proteins in each proteome. At first glance one might deduce that there is a linear relationship between the number of simple sequence-containing proteins and the total number of proteins. The line of best fit drawn in Figure 1a has a correlation coefficient of 0.99. However, the eukaryotes possess significantly larger proteomes than the prokaryotes, and consequently far more simple sequences. In effect, the fit to the data is reduced to a fit to five points – the four eukaryotes plus the prokaryotes essentially as a single point.

If one considers just the four eukaryotes surveyed, a line of best fit through the data in Figure 1a would yield a correlation coefficient of 0.99. Note however, this is just a four point fit, and that it may well be that there is not a linear relationship between eukaryote proteome size and  $Prot_{SS}$ . Clearly the complete proteomes of more eukaryotes need to be examined, once they become available, in order to gain a better understanding of this relationship. What can be concluded from this figure, and the data in Table 1, is that a remarkable number of the proteins in the eukaryote proteomes surveyed possess at least one simple sequence as defined in this work.

The individual amounts are 53% of the proteins in SC, 51% in CE, 59% in DM and 55% in AT. Why DM would possess a significantly higher fraction of proteins with at least one simple sequence is not clear. Karlin *et al.* (10), in a recent survey of homopolymeric runs in proteins  $\geq 200$  residues in size, found that DM possessed far more than other eukaryotes. They also found that human proteins possessed more of these runs than proteins from CE, despite there being more CE proteins surveyed. Such data suggest that the human proteome may also possess a larger fraction of proteins containing simple sequences than the average observed in this work.

Figure 1b is a plot of  $\text{Prot}_{\text{SS}}$  against the number of proteins in each proteome for the twenty-six prokaryotes surveyed. There is a clear linear correlation, with the line of best fit having a correlation coefficient of 0.92. Two prokaryotes, the archaea HS and the bacteria DR, appear to be outliers. Excluding these from the fit results in a correlation coefficient of 0.96. The strong linear correlation observed for the prokaryotes might suggest that these simple sequences have arisen via random events, leading to random distributions that depend only upon the number of proteins in each proteome. As will be demonstrated below however, our data suggest the opposite, that the occurrence and distributions of simple sequences is not random in nature and that many of these sequences may possess biological significance.

Figure 1c, a bar plot of the ratio of number of simple sequences found,  $\text{SS}_{\text{Tot}}$ , to  $\text{Prot}_{\text{SS}}$  for each organism surveyed, illustrates the difference in occurrence of protein simple sequences in prokaryotes and eukaryotes. Prokaryotes have far fewer simple sequences per protein than do the eukaryotes. In all cases, the prokaryotes have fewer simple sequences than the total number of proteins in their proteomes, whereas the eukaryotes possess more (Table 1). The prokaryotes average 1.40 simple sequences per protein possessing at least one simple sequence (the dashed line on Figure 1c). Once again, HS and DR are clear outliers among the prokaryotes, possessing

$SS_{Tot}/Prot_{SS}$  ratios of 1.68 and 1.73 respectively, both values greater than two standard deviations from the mean for prokaryotes. The eukaryotes have ratios that range from 1.88 in AT, through 2.09 in CE and 2.18 in SC, up to 3.09 simple sequences per protein possessing at least one in DM. Eukaryotes clearly not only tolerate a significantly higher occurrence of these sequences than do the prokaryotes, they are also more likely to possess multiple simple sequences in each protein.

The ratio  $SS_{Tot}/Prot_{SS}$  is of course dependent upon our definition of protein simple sequences. One can imagine that increasing the size of the allowable gap (currently set at five or less residues) will result in some of the simple sequences merging, resulting in fewer overall, but an increase in the number of longer sequences. The result will be lower values of  $SS_{Tot}/Prot_{SS}$  for each proteome.

A number of groups have examined the occurrence of homopolymeric runs of sequence and noted that eukaryotes possess more per protein than prokaryotes (19, 21, 23). Nishizawa *et al.* (23) note that “modern” tissue-specific proteins have a higher tendency to possess homopolymeric stretches of up to twenty residues in length as compared to ancient proteins. They go on to postulate that this repetitiveness enhances the chance for intermolecular interactions. This hypothesis is supported by observations that simple sequences enriched in glutamine, proline or charged residues are often found in protein interaction domains of transcription regulatory proteins (2-5), and that proline-rich sequences are common protein-protein interaction domains (11, 12). It seems likely then that eukaryotes, in particular the multicellular organisms, have evolved to require numerous protein simple sequences for functional purposes.

It is not clear why HS and DR would be outliers among the prokaryotes in Figure 1. HS is an extreme halophile (26), the only one in the set of organisms surveyed. It is tempting to postulate that HS might possess a higher proportion of simple sequences as a result of evolving to survive in such an unusual environment. Ng *et al.* (26) pointed out that 36% of the putative proteins in the HS proteome were unrelated to any previously reported at that time, and that these may well provide the mechanisms by which HS can survive extreme salt concentrations. However, the HS proteome has not been analyzed in sufficient detail to know whether those proteins are particularly enriched in simple sequences, so we cannot draw any conclusions at this point.

DR has been nicknamed “Conan the bacterium” for its amazing ability to resist very high doses of ionizing radiation and UV-irradiation (27), and is the only organism surveyed to possess these remarkable traits. It has been speculated that the radiation resistance of DR is due to its unique polyploid nature, and the abundant DNA repeat elements in its genome. These DNA repeats may function to regulate DNA degradation after damage to this organism. The high number of protein simple sequences identified in this species may be attributed to such repeats, although not to the polyploid nature of DR. This organism possesses more simple sequences *per protein* than other prokaryotes (Table 1). Simply possessing multiple copies of each gene would not raise the number of simple sequences per protein. The protein simple sequences may have arisen over time as a result of errors made by the DNA repair apparatus of DR while “rebuilding” its genome from multiple gene copies after exposure to extreme conditions such as radiation. On the other hand, some of these simple sequences may play an active role in the survival mechanisms developed by DR. Further functional analysis of the DR proteome is required in order to better understand why this organisms possesses so many protein simple sequences.

For reasons of clarity and focus, the remainder of this manuscript will focus on the occurrence and distributions of protein simple sequences in eukaryotes.

### ***Overall Length Distributions***

Figure 2, a log-log plot of the number of simple sequences found against their length, is a clear illustration of the remarkable simple sequence length distributions observed in the four eukaryotes examined. Prokaryotes display similar length distributions, although generally the longest prokaryotic simple sequences are shorter than the longest eukaryotic sequences (data not shown). At the shorter simple sequence lengths a periodicity in the data can be seen, with there being fewer occurrences where the length is an odd number as compared to adjacent even numbered lengths. This is a consequence of the algorithm used to identify the simple sequences. As an example, given the threshold of  $\geq 50\%$ , a simple sequence eleven residues long must possess at least six residues of a given type. This amounts to a minimum of 55% enrichment, whereas a twelve residue simple sequence can also possess six residues, leading to a minimum of 50% enrichment. This periodicity tends to be damped out at long simple sequence lengths.

Not surprisingly there is a sharp decrease in the number of simple sequences found with increasing length. The shorter simple sequences are extremely common. Of course such observations are in part a result of our definition of a protein simple sequence. Lowering or raising the threshold of  $\geq 50\%$  enrichment will change these numbers, as will changing the gap between putative simple sequences. Nonetheless, protein simple sequences highly enriched in a single residue type are remarkably common.

The longest simple sequence was found in AT, is 410 residues long and is enriched in glycine. AT is not alone in possessing remarkably long simple sequences. The longest in SC is

246 residues long and is enriched in serine. The longest in CE is threonine-rich and is 291 residues long, while DM possesses a 322 residue long glycine-rich sequence. Notably, all four of these simple sequences occur in proteins that have been annotated as being hypothetical. The majority of the simple sequences found are of course much shorter than these, the vast majority being sixty or fewer residues in length (~99.5%; Figure 2).

Figure 3 is a bar plot of the ratio of the number of simple sequences found to the number of proteins in the proteome as a function of length for the four eukaryotes. Division by the proteome size allows for direct comparison of the organisms. The data are split into three length scales; ten to twenty residues (Figure 3, panel a), twenty to forty (panel b) and forty to sixty (panel c). The periodicity observed in Figure 2 is obvious in Figure 3a, and can be seen to have subsided in Figure 3b. It is clear from Figure 3 that DM averages more simple sequences per protein at all lengths than the other organisms, despite AT possessing more in total, and CE possessing a similar number (Table 1). In fact, DM possesses more than twice as many simple sequences of lengths  $\geq 20$  residues per protein than any of the other three eukaryotes examined. Clearly DM has evolved to tolerate large numbers of simple sequences. What is not entirely clear is whether this observation is linked to the functional requirements of DM. Nishizawa *et al.* (23) have pointed out that neural and immune system-specific proteins have a higher propensity to possess short runs of sequence consisting entirely of one residue type. One could reasonably expect that this would extend to the highly enriched simple sequences found in this survey. If so, it may not in fact be surprising that DM possesses such an abundance of these simple sequences as compared to the other eukaryotes examined. An analysis of the functions of the proteins in DM possessing simple sequences will shed light on this, as will surveys of the proteomes of other eukaryotes as they become available.

What is perhaps most surprising in Figure 3 is that SC possesses the second largest fraction of proteins in its proteome averaging a single simple sequence at almost all lengths up to sixty residues. AT generally possesses the fewest. Huntley and Golding (19) had previously noted that SC possesses a high proportion of protein simple sequences, although they could not explain why. The common theme from Figure 3 is that many of the proteins in the proteomes of the four of the eukaryotes examined possess simple sequence regions. In fact, it has previously been observed that protein simple sequences are the most commonly shared sequence pattern among the eukaryotes (19). Huntley and Golding (19) have suggested that protein simple sequences are the equivalent of “junk DNA”, serving little purpose. However, given the number of simple sequences found, coupled with the known functions of some of them (2, 3, 10-12), it is tempting to postulate that eukaryotes tolerate and even require large numbers of simple sequences for functional reasons.

### ***Residue Length Dependencies***

From Figures 2 and 3 it would appear that the four eukaryotes examined have similar protein simple sequence distributions, albeit with differences in relative abundance. Striking differences between the organisms are revealed when simple sequence distributions are considered at the level of individual residue types. Figure 4 shows the ratio of the number of simple sequences found above that expected from the Poisson distribution to the number of proteins in each organisms proteome,  $\Delta_R$ , plotted against simple sequence length, for each residue. The sequence lengths are binned into ranges: 10-20 (Figure 4a), 21-40 (Figure 4b) and 41 and more (Figure 4c) residues. Data for cysteine, methionine and tryptophan are omitted since we found very few simple sequences containing these rare residues. The ratio  $\Delta_R$  is a measure of how common simple sequences are, above the Poisson distribution predictions, per protein in



each eukaryote proteome. This ratio allows for easy comparison of the organisms. A higher  $\Delta_R$  indicates that simple sequences of a given length in an organism are more common in comparison to the other organisms, even though the actual number found might be the same or even lower. A negative value of  $\Delta_R$  indicates that those simple sequences are found less often than predicted from the Poisson distribution. Such sequences are presumably discriminated against for various reasons.

### Features Common to All Eukaryotes Examined

Before considering differences between the distributions of simple sequences for each organism, there are some features common to all four eukaryotes worth looking at in Figure 4. Perhaps the most obvious common features are the negative values of  $\Delta_R$  at short lengths (ten to twenty residues) observed for the small apolar residues isoleucine, leucine and valine (Figure 4a). These negative values indicate that there are fewer such simple sequences than might be expected from the Poisson distribution. The most striking observation is that for leucine, the most common residue, where we find hundreds fewer leucine-rich sequences at short lengths than expected. This is particularly apparent for SC and AT, which have  $\Delta$  values of -585 and -1523 respectively. CE and DM also have large negative  $\Delta$  values, which are -497 and -397 respectively. For simple sequences of 21 to 40 residues (Figure 4b), the  $\Delta_R$  values for leucine, isoleucine and valine become positive, but are small. For even longer lengths the  $\Delta_R$  values are zero or very small. We appear to be observing a discrimination against simple sequences highly enriched in these small apolars. This has previously been observed by Green and Wang (6), Katti *et al.* (5) and Karlin *et al.* (10), who all found very few occurrences of runs of these residues longer than ten residues in length. In these studies a run of residues was defined as consisting solely of a single type of

residue, except for in the study of Katti *et al.* (5) where a 10% mismatch was allowed for sequence runs greater than twenty residues in length. We may be observing a biophysical effect here. Sequences of ten to twenty residues in length with  $\geq 50\%$  leucine, isoleucine or valine will be highly hydrophobic and may pose an aggregation risk for proteins that contain them. Hence they are evolutionarily discriminated against. Moderate lengths, twenty-one to around thirty residues become more likely since such sequences could act as membrane-spanning regions, as suggested by Schwartz *et al.* (28).

We should note that the actual number of leucine-, isoleucine- and valine-enriched simple sequences found can be quite large. For example, in DM we find 1841, 96 and 223 simple sequences of length ten enriched in each of these residues respectively. Due to the relative abundance of these residues however, the Poisson distribution predictions are also large (1445, 95 and 218 respectively), leading to small or negative values of  $\Delta$  and  $\Delta_R$ .

It is notable that we find positive values of  $\Delta_R$  for phenylalanine and tyrosine at short, moderate, and even long lengths (Figure 4). The tyrosine-rich sequences are particularly surprising given that this is one of the rarer residues. One might expect that sequences enriched in such large hydrophobic residues might be disfavored, and yet this does not appear to be the case. It is not clear as to why such sequences would be tolerated.

Careful inspection of Figure 4 reveals that sequences highly enriched in serine, glutamate, lysine and alanine appear to be favored by all four of the eukaryotes examined at short lengths (Figure 4a). At moderate lengths alanine-rich sequences become less common (Figure 4b), while at long lengths glycine-rich sequences seem to be favored. Similar distributions of runs of sequence containing these residues were observed by Green and Wang (6) and Katti *et al.* (5),

although these authors did not normalize their data for residue occurrence, nor for what might be expected were sequences random. It is not entirely apparent why sequences highly-enriched in serine are tolerated, or even required, by eukaryotes, although there are certainly examples of important protein domains enriched in this residue. One such example is the C-terminal domain of RNA polymerase II, which is functionally essential and consists of the heptad YSPTSPS repeated between 26 and 52 times in various organisms (29). Interestingly, this serine-enriched region (~ 43 % serine) is known to interact with proline-rich regions (12), as well as a family of serine/arginine-rich proteins (30). Wootton and Drummond (14) have suggested that sequences enriched in serine may act as flexible linkers between protein domains in much the same way as postulated for glycine-rich sequences.

Sequences enriched in charged residues, such as the lysine- and glutamate-rich sequences seen to be favored by the eukaryotes (Figure 4), have been associated with DNA and RNA processing, chromatin structure, ion binding and protein-protein interactions (13). The involvement of such simple sequences in a wide variety of functional roles might therefore explain their relative abundance. Alanine is known to be the most energetically favorable residue in  $\alpha$ -helices (31, 32). One might therefore expect that sequences that are 50% or greater alanine in composition will have a tendency to be  $\alpha$ -helical, although this will of course be modulated by the nature of the other residues in the sequence, as well as by the tertiary structure of the proteins they are part of. The preference for short alanine-rich sequences that we have observed might then be related to secondary structure requirements. The long glycine-rich sequences found are probably tolerated for the opposite reason. That is, these most likely represent flexible linkers between protein domains.

One of the more surprising observations from Figure 4 is that of simple sequences highly enriched in histidine. Although there are not many of these at all lengths, the number we find above that expected is significant. Some of these are quite long. For example, the four longest histidine-rich sequences in DM are 46, 51, 54 and 56 residues long. In CE the four longest are 50, 51, 84 and 251 residues long, although the longest of these is in a protein annotated as being hypothetical and could in fact be an indication that this is not an expressed protein. Histidine is one of the most rare residues, comprising just 2.2 to 2.7% of all residues in the four proteomes. By comparison, methionine has a similar level of occurrence and yet we find almost no simple sequences enriched in this residue above the Poisson distribution predictions. Similarly, we find very few tryptophan- and cysteine-enriched sequences. One might postulate that histidine-rich sequences have some kind of ion-binding function, although this has not been demonstrated.

### Distribution Differences Among the Eukaryotes Examined

It is immediately apparent from Figure 4 that DM has a markedly different distribution of simple sequences when compared to the other three eukaryotes. The data for DM demonstrates a preference for simple sequences of all lengths enriched in alanine, glutamine, glycine and serine. At short to moderate lengths, ten to forty residues (Figure 4, panels a and b), DM also shows some preference for asparagine-, proline-, threonine-, and perhaps most surprisingly, histidine-enriched sequences. To a lesser extent, there may also be a preference for aspartate- and arginine-rich sequences. These observed preferences are in large part responsible for the “unusually” high  $SS_{Tot}/Prot_{SS}$  ratio observed for DM (Figure 1c). The large numbers of glutamine-, and to some extent asparagine-rich sequences in DM were also observed by Michelitsch and Weissman (9), who suggested that many of these may act as protein-protein interaction domains. It is not clear

why DM would tolerate, and perhaps even require, large numbers of alanine-, glycine- and serine-rich sequences.

Although DM is clearly different to the other three eukaryotes, it would be a mistake to assume that there are no significant differences between the distributions observed for the other organisms. SC has preferences for asparagine- and aspartate-enriched sequences at all lengths, along with a striking preference for moderate length to long serine-rich sequences. Furthermore, SC disfavors leucine- and isoleucine-rich sequences more than the other eukaryotes, and is somewhat less tolerant of arginine-, glycine- and proline-rich sequences. The reasons behind each of these preferences are not always clear. For example, the reasons for the large preference for moderate length to long serine-rich sequences are unknown. Wootton and Drummond (14) have suggested that sequences rich in serine form flexible linkers between protein domains. If this is true, then the preference for serine-rich sequences in SC may be linked to the observed lower tolerance for glycine-rich sequences (Figure 4). SC may have evolved to use serine-rich sequences as linkers instead of the glycine-rich sequences that the other eukaryotes seem to prefer. Another potential role for serine-rich regions is discussed below. Michelitsch and Weissman (9) have previously observed large numbers of asparagine-rich sequences in SC, as well as in other eukaryotes. These authors postulate that such regions act as modulators of protein-protein interactions. Why SC would require a larger fraction of asparagine-rich sequences for such interactions as compared to the eukaryotes is unclear. The lower tolerance for proline-rich sequences is probably due to the unicellular nature of SC. It has no need for the proline-rich extracellular structural proteins that the multicellular eukaryotes require. The reasons for the discrimination against leucine- and isoleucine-enriched sequences, and the lower tolerance for arginine-rich sequences remain enigmatic.

The worm CE also possesses its own unique distribution of protein simple sequences. From Figure 4 it can be seen that CE has some preference for short phenylalanine-rich sequences, and for long glutamine- and serine-enriched sequences. CE also appear to be less tolerant of asparagine-rich sequences than SC, DM, and perhaps AT, and is less tolerant than DM and AT of long proline-rich sequences. AT has little tolerance for threonine-rich sequences and a lowered tolerance for glutamine-rich sequences (Figure 4). AT does not appear to have a heightened preference for any particular simple sequences at any length scale compared to the other eukaryotes.

It is clear that each of the four eukaryotes examined possesses its own unique distribution of simple sequences (Figure 4). Based upon the analysis of homopolymeric runs performed by Karlin *et al.* (10), and an analysis by Kreil and Kreil (33) of asparagine-rich sequences, it seems clear that the human proteome will also display a unique simple sequence distribution. Some of the differences observed for the four eukaryotes examined arise for understandable reasons. For example, SC would not be expected to possess as many proline-rich sequences as the other eukaryotes examined since it does not have the same requirements for proline-rich structural proteins. However, as noted repeatedly above, the reasons for many of the various simple sequence preferences observed are not known. Some differences might well arise as a result of an organism using particular residues for the same purposes as other organisms use a different set of residues. For example, as suggested, SC might employ serine-rich regions as flexible linkers where CE, DM and AT use glycine-rich sequences. A detailed analysis of the conservation of simple sequence regions will aid in resolving such issues. Huntley and Golding (19) have noted that simple sequences are the most commonly shared feature between proteins, but that the identity of the residues within the sequences can vary between organisms.

## Functional Analysis of Protein Simple Sequence Occurrence

Our survey of the eukaryote (and prokaryote) proteomes has resulted in the identification of an enormous number of protein simple sequences, far more than would be expected were sequences random in nature. We have postulated that many of these sequences play some kind of functional role. This postulate is supported by a limited amount of experimental and bioinformatic evidence (3, 5, 9-12, 29, 30, 34). In order to further examine this issue we have examined the distribution of simple sequences in proteins of known function. Specifically, we have collected the sequences of all proteins from each of the four eukaryotes that are annotated in the SWISS-PROT database (35, 36) as being involved in a protein class (e.g. membrane proteins) or set of processes (e.g. transcription). The occurrence and distributions of simple sequences in these proteins were then analyzed employing the approaches used on intact proteomes above. The results are shown in Table 2. Note that the data shown is highly dependent upon the completeness and accuracy of the annotations in SWISS-PROT, as well as how well-studied the particular classes of proteins are in each organism. As a result of these limitations we have found comparatively few protein sequences in most cases. In addition, some proteins may appear in more than one classification in Table 2. Thus, it is difficult to make direct comparisons between classes as to the number of simple sequences found, as well as between organisms. However, it is feasible to consider the most common types of simple sequence found (Table 2).

Immediately noticeable in Table 2 is the abundance of serine-rich sequences in almost all classes of proteins examined. Serine-rich sequences are the most common in all four organisms, particularly at the most abundant short length scales (Figure 4), so perhaps this finding is not surprising. The role of serine-rich sequences is not however clear. As noted above, it has been proposed that such sequences can act as flexible linkers between protein domains (14) or as

protein interaction domains (29, 30). Proteins containing regions enriched in both serine and arginine have also been shown to be involved in mRNA splicing control (37). Serine-rich regions may also function as some form of phosphorylation switch, much as the C-terminal domain of RNA polymerase II operates (29).

Considering now each class of protein in Table 2, it can be seen that the most common simple sequences in the limited set of cell cycle proteins identified are serine-rich. Very few cell cycle proteins were found, except in the case of SC, which is perhaps the model system for studying these processes. The 102 SC cell cycle proteins found possess a total of 177 simple sequences, over a third of which (sixty-four) are serine-rich. This is a clear enrichment of such sequences as compared to the overall distribution of simple sequences in SC (Figure 4). Potential roles for these sequences are as discussed above.

We found relatively few proteins with the keyword “metabolism” in their annotations (Table 2). With the preceding as a caveat, it is notable that there are fewer simple sequences per protein in metabolism-related proteins (significantly less than one per protein) than the average over intact proteomes (slightly more than one per protein; Table 1). This would suggest that simple sequences are either generally not required in metabolism-related proteins, or that they are discriminated against in comparison to other protein classes. However, as already noted, few proteins were identified in this class and we could simply be observing the vagaries of poor statistics.

Using “signal” as a keyword, we have identified a significant number of proteins in all four eukaryotes (Table 2). These possess a significant number of simple sequences, the most common of which are enriched in serine, threonine, proline and, perhaps surprisingly, leucine.



Given that signal transduction processes involve significant numbers of phosphorylation, and dephosphorylation, events, it is perhaps not so remarkable that serine- and threonine-rich sequences are common in signaling proteins. There are also a number of small protein interaction domains common in signaling processes (e.g. SH3 domains) that bind to proline-rich sequences (11), leading to an enrichment in such sequences in this class. Thus, the occurrence of serine-, threonine- and proline-rich sequences in this class of proteins would appear to be biologically significant. The occurrence of a significant number of leucine-rich sequences is at first puzzling, particularly given that such sequences are found at levels lower than would be predicted using our Poisson-distribution model (Figure 4). However, it is possible that a reasonable number of the proteins in this class possess membrane-spanning segments which, as will be discussed below, can be leucine-rich.

A significant number of transcription-related proteins were also identified (Table 2). Remarkably, the transcription-related proteins in SC and DM possess enormous numbers of simple sequences (677 in 274 proteins and 838 in 177 proteins respectively). Although the same level of enrichment is not seen in CE and AT, it is tempting to postulate that large numbers of simple sequences indicate important functional roles in transcription processes. Indeed, such proteins are known to often possess glutamine-rich sequences (3), so it is not surprising that such sequences are common in DM transcription-related proteins. We also find large numbers of serine-rich sequences (Table 2). Perhaps the best-known example of a serine-rich region acting as a phosphorylation switch region is at the heart of transcription, being found in RNA polymerase II (29). Although not enriched in serine enough to be found in our surveys, this region is known to interact with a variety of transcription factors when not phosphorylated. These interactions, and consequently transcription, are interrupted when serines become phosphorylated. It is

possible that there are similar serine-rich switch/interaction regions in other transcription-related proteins.

A reasonable number of transport-related proteins were also found (Table 2). These possess approximately the same numbers of simple sequences as would be expected from the overall average values for the four eukaryotes (Table 1). Leucine-, alanine- and serine-rich sequences are the most common. A significant number of transport-related proteins will be associated with membranes given that transport of molecules through membranes is a common and vital set of processes. The large numbers of leucine-rich, and perhaps alanine-rich, regions are then most likely indicative a membrane-spanning regions, as suggested by Schwartz *et al.* (28).

Finally, we have identified numerous membrane-associated proteins, many of which contain simple sequences (Table 2). Presumably for the reasons noted above, large numbers of leucine-rich sequences are found in this class of proteins. In fact, many of these leucine-rich regions are annotated as being membrane-spanning in the SWISS-PROT files for these proteins. Why serine-rich regions would be so abundant is not clear. Some of these are probably found in signaling proteins associated with the membrane (see above), while others may be acting as flexible linkers separating soluble domains from integral membrane domains. Wootton and Drummond (14) have hypothesized that serine-rich regions act as flexible linkers. Notably, glycine-rich regions, also thought to act as linkers, are common in DM membrane-related proteins. Perhaps serine-rich regions are substituted for glycine-rich in the other organisms (Table 2).

## Simple Sequence Structure

It would of course be useful to know the types of structures adopted by protein simple sequences. Unfortunately little is known about the structural properties of such sequences. Saqi (17) and more recently Huntley and Golding (38) have looked for all occurrences of simple sequences in protein structures in the Protein Data Bank (PDB) (39). Very few were found. Huntley and Golding (38) point out that simple sequences are under-represented in the PDB and hypothesize that this indicates that such regions are intrinsically disordered. Intrinsically disordered regions of proteins are a barrier to structure determination and are consequently routinely deleted from proteins by structural biologists. That simple sequences, particularly relatively long sequence, are disordered is supported by the work of Dunker and co-workers (15, 16, 40), who use low-complexity sequences as identifiers of intrinsically disordered proteins. There are indications however that not all protein simple sequences are unstructured. For example, leucine-rich membrane-spanning sequences will be highly structured, most likely  $\alpha$ -helices, in the membrane. Proline-rich regions are believed, and in many cases have been shown, to adopt the left-handed polyproline II helical conformation (11). It would be a mistake to assume that all simple sequences are unstructured. This is an area that clearly requires further investigation.

## **Conclusions**

We have presented here a survey of protein simple sequences highly enriched in a single residue type ( $\geq 50\%$ ) in the proteomes of four eukaryotes. For comparison we have also surveyed the proteomes of twenty-six prokaryotes. Strikingly large number of simple sequences are found in all of the organisms surveyed (Table 1). We find that eukaryotes possess, on average, one or

more such simple sequence per protein, whereas prokaryotes average less than one simple sequence for each protein in their proteomes. Furthermore, proteins in eukaryotes that possess at least one simple sequence average between just under two up to slightly more than three simple sequences per protein. These findings are consistent with the work of others (19, 21, 23). The number of simple sequences in the proteomes of prokaryotes is strongly correlated with the number of proteins in their proteomes (Figure 1b). Given that we have only surveyed four proteomes, it is not clear that a linear relationship will be applicable to eukaryotes.

Among the eukaryotes we find that DM possesses more simple sequences per protein than any of the other three eukaryotes (Table 1). This is true for all simple sequence lengths (Figure 3). By comparison, SC, CE and AT possess similar number of simple sequences per protein at most lengths, with SC perhaps showing some preference for long simple sequences (Figure 3). In the distributions for the intact proteomes, we find that simple sequences enriched in certain residues, for example alanine, glutamine, glutamate, glycine and serine, appear to be favored, whereas other residues, specifically leucine, isoleucine and valine, are discriminated against. These preferences do not correlate with residue occurrence. Some of these observed preferences can be rationalized in terms of structure and/or function, while others remain enigmatic.

The most notable finding of these surveys is that each of the eukaryotes possesses its own unique distribution of protein simple sequences. We find that each organism apparently has preferences for simple sequences enriched in certain residues, while at times disfavoring simple sequences enriched in other residues. It is not clear why these eukaryotes have evolved to have differing simple sequence distributions. However, given the sheer number of such sequences found, plus the known functional importance of those simple sequences that have been studied in detail, it is tempting to postulate that not only have eukaryotes evolved to tolerate large numbers

of simple sequences, but also that they require many of these. A simple analysis of simple sequences in classes of proteins indicates that some classes may favor simple sequences enriched in certain residues (Table 2).

The data presented here raise questions that can only be answered by further study and analysis. For example, is there an association between type of simple sequence and function ? The data in Table 2 are suggestive, but by no means conclusive. Do different organisms use different types of simple sequence for the same function ? The fact that each organism possesses a unique distribution implies that this may be the case, but we have no direct evidence. What are the structural properties of such sequences ? Little structural data is currently available, although it is clear that it would be incorrect to assume that all simple sequences will be disordered. Answers to questions such as these will shed light on the abundance and distributions of simple sequences highlighted here.

## **Acknowledgements**

K.L.S. was supported in part by a Post-Doctoral Research Grant in Computational Genomics from the Pharmacia Corporation. This work was supported in part by a grant from the National Science Foundation to T.P.C. (MCB-00110720). We wish to thank Brian Chellgren, Marnie Campbell, David Rodgers, Rajeev Aurora and George Rose for helpful discussions.

## References

1. Wootton, J. C., Federhen, S. (1996) Analysis of compositionally biased regions in sequence databases. *Meth. Enzymol.* 266, 554-571
2. Brendel, V., Karlin, S. (1989) Association of charge clusters with functional domains of cellular transcription factors. *Proc. Natl. Acad. Sci., USA* 86, 5698-5702
3. Gerber, H. P., Seipel, K., Georgiev, O., Hofferer, M., Hug, M., Rusconi, S., Schaffner, W. (1994) Transcriptional activation modulated by homopolymeric glutamine and proline stretches. *Science* 263, 808-811
4. Kashi, Y., King, D., Soller, M. (1997) Simple sequence repeats as a source of quantitative genetic variation. *Trends Gen.* 13, 74-78
5. Katti, M. V., Sami-Subbu, R., Ranjekar, P. K., Gupta, V. S. (2000) Amino acid repeat patterns in protein sequences: their diversity and structural-functional implications. *Protein Sci.* 9, 1203-1209
6. Green, H., Wang, N. (1994) Codon reiteration and the evolution of proteins. *Proc. Natl. Acad. Sci., USA* 91, 4298-4302
7. Cummings, C. J., Zoghbi, H. Y. (2000) Trinucleotide repeats: mechanisms and pathophysiology. *Annu. Rev. Genom. Hum. Gen.* 1, 281-328
8. Karlin, S., Burge, C. (1996) Trinucleotide repeats and long homopeptides in genes and proteins associated with nervous system disease and development. *Proc. Natl. Acad. Sci., USA* 93, 1560-1565

9. Michelitsch, M. D., Weissman, J. S. (2000) A census of glutamine/asparagine-rich regions: implications for their conserved function and the prediction of novel prions. *Proc. Natl. Acad. Sci., USA* 97, 11910-11915
10. Karlin, S., Brocchieri, L., Bergman, A., Mrazek, J., Gentles, A. J. (2002) Amino acid runs in eukaryotic proteomes and disease associations. *Proc. Natl. Acad. Sci., USA* 99, 333-338
11. Kay, B. K., Williamson, M. P., Sudol, M. (2000) The importance of being proline: The interaction of proline-rich motifs in signaling proteins with their cognate domains. *FASEB J.* 14, 231-241
12. Williamson, M. P. (1994) The structure and function of proline-rich regions in proteins. *Biochem. J.* 297, 249-260
13. Karlin, S. (1995) Statistical significance of sequence patterns in proteins. *Curr. Opin. Struct. Biol.* 5, 360-371
14. Wootton, J. C., Drummond, M. H. (1989) The Q-linker: a class of interdomain sequences found in bacterial multidomain regulatory proteins. *Protein Eng.* 2, 535-543
15. Romero, P., Obradovic, Z., Li, X., Garner, E. C., Brown, C. J., Dunker, A. K. (2001) Sequence complexity of disordered protein. *Proteins* 42, 38-48
16. Dunker, A. K., Obradovic, Z., Romero, P., Garner, E. C., Brown, C. J. (2000) Intrinsic disorder in complete genomes. *Genome Inform. Ser. Workshop Genome Inform.* 11, 161-171
17. Saqi, M. (1995) An analysis of structural instances of low complexity segments. *Protein Eng.* 8, 1069-1073

18. Meyer, E. F., Tollet Jr., W. J. (2001) WWWWhy does nature stutter ? A survey of strands of repeated amino acids. *Acta Cryst. D Biol. Cryst.* D57, 181-186
19. Huntley, M., Golding, G. B. (2000) Evolution of simple sequence in proteins. *J. Mol. Evol.* 51, 131-140
20. Golding, G. B. (1999) Simple sequence is abundant in eukaryotic proteins. *Protein Sci.* 8, 1358-1361
21. Marcotte, E. M., Pellegrini, M., Yeates, T. O., Eisenberg, D. (1999) A census of protein repeats. *J. Mol. Biol.* 293, 151-160
22. Katti, M. V., Ranjekar, P. K., Gupta, V. S. (2001) Differential distribution of simple sequence repeats in eukaryotic genome sequences. *Mol. Biol. Evol.* 18, 1161-1167
23. Nishizawa, K., Nishizawa, M., Kim, K. S. (1999) Tendency for local repetitiveness in amino acid usages in modern proteins. *J. Mol. Biol.* 294, 937-953
24. Soper, H. E. (1914) Tables of Poisson's exponential limit. *Biometrika* 10, 25-35
25. Bork, P., Copley, R. (2001) Filling in the gaps. *Nature* 409, 818-820
26. Ng, W. V., Kennedy, S. P., Mahairas, G. G., Berquist, B., Pan, M., Shukla, H. D., Lasky, S. R., Baliga, N. S., Thorsson, V., Sbrogna, J., Swartzell, S., Weir, D., Hall, J., Dahl, T. A., Welti, R., Goo, Y. A., Leithauser, B., Keller, K., Cruz, R., Danson, M. J., Hough, D. W., Maddocks, D. G., Jablonski, P. E., Krebs, M. P., Agevine, C. M., Dale, H., Isenbarger, T. A., Peck, R. F., Pohlschroder, M., Spudich, J. L., Jung, K. W., Alam, M., Freitas, T., Hou, S., Daniels, C. J., Dennis, P. P., Omer, A. D., Ebhardt, H., Lowe, T. M., Liang, P., Riley, M., Hood, L., DasSarma, S. (2000) Genome sequence of *Halobacterium* species NRC-1. *Proc. Natl. Acad. Sci., USA* 97, 12176-12181



27. White, O., Eisen, J. A., Heidelberg, J. F., Hickey, E. K., Peterson, J. D., Dodson, R. J., Haft, D. H., Gwinn, M. L., Nelson, W. C., Richardson, D. L., Moffat, K. S., Qin, H., Jiang, L., Pamphile, W., Crosby, M., Shen, M., Vamathevan, J. J., Lam, P., McDonald, L., Utterback, T., Zalewski, C., Makarova, K. S., Aravind, L., Daly, M. J., Minton, K. W., Fleischmann, R. D., Ketchum, K. A., Nelson, K. E., Salzberg, S., Smith, H. O., Venter, J. C., Fraser, C. M. (1999) Genome sequence of the radioresistant bacterium *Deinococcus radiodurans* R1. *Science* 286, 1571-1577
28. Schwartz, R., Istrail, S., King, J. (2001) Frequencies of amino acid strings in globular protein sequences indicate suppression of blocks of consecutive hydrophobic residues. *Protein Sci.* 10, 1023-1031
29. Corden, J. L. (1990) Tails of RNA polymerase II. *Trends Biochem. Sci.* 15, 383-387
30. Yuryev, A., Patturajan, M., Litingtung, Y., Joshi, R. V., Gentile, C., Gebara, M., Corden, J. L. (1996) The C-terminal domain of the largest subunit of RNA polymerase II interacts with a novel set of serine/arginine-rich proteins. *Proc. Natl. Acad. Sci., USA* 93, 6975-6980
31. Aurora, R., Creamer, T. P., Srinivasan, R., Rose, G. D. (1997) Local interactions in protein folding: Lessons from the  $\alpha$ -helix. *J. Biol. Chem.* 272, 1413-1416
32. Chakrabartty, A., Kortemme, T., Baldwin, R. L. (1994) Helix propensities of the amino acids measured in alanine-based peptides without helix-stabilizing side-chain interactions. *Protein Sci.* 3, 843-852
33. Kreil, D. P., Kreil, G. (2000) Asparagine repeats are rare in mammalian proteins. *Trends Biochem. Sci.* 25, 270-271

34. Tonjum, T., Caugant, D. A., Dunham, S. A., Koomey, M. (1998) Structure and function of repetitive sequence elements associated with a highly polymorphic domain of *Neisseria meningitidis* PiLQ protein. *Mol. Microbiol.* 29, 111-124
35. Bairoch, A., Apweiler, R. (1997) The SWISS-PROT protein sequence database: Its relevance to human molecular medical research. *J. Mol. Med.* 75, 312-316
36. Bairoch, A., Boeckmann, B. (1991) The SWISS-PROT protein sequence data bank. *Nuc. Acids Res.* 19, 2247-2249
37. Manley, J. L., Tacke, R. (1996) SR proteins and splicing control. *Genes Dev.* 10, 1569-1579
38. Huntley, M. A., Golding, G. B. (2002) Simple sequences are rare in the Protein Data Bank. *Proteins* 48, 134-140
39. Bernstein, F. C., Koetzle, T. F., Williams, G. J. B., Meyer Jr., E. F., Brice, M. D., Rodgers, J. R., Kennard, O., Shimanouchi, T., Tasumi, M. (1977) The Protein Data Bank: A computer-based archival file for macromolecular structures. *J. Mol. Biol.* 112, 535-542
40. Dunker, A. K., Lawson, J. D., Brown, C. J., Williams, R. M., Romero, P., Oh, J. S., Oldfield, C. J., Campen, A. M., Ratliff, C. M., Hipps, K. W., Ausio, J., Nissen, M. S., Reeves, R., Kang, C., Kissinger, C. R., Bailey, R. W., Griswold, M. D., Chiu, W., Garner, E. C., Obradovic, Z. (2001) Intrinsically disordered protein. *J. Mol. Graphics Model.* 19, 26-59

## Figure Legends

**Figure 1:** **Panel a** shows the number of proteins possessing at least one simple sequence plotted against the total number of proteins in the proteome for all organisms. **Panel b** shows the same data for just the prokaryotes. Line of best fit was calculated excluding HS and DR. **Panel c** is a bar plot of the ratio of simple sequences to the number of proteins possessing simple sequences for each organism. The dashed line denotes the average value for prokaryotes (1.40). The ratios for the eukaryotes and the two outlying prokaryotes (HS and DR) are provided.

**Figure 2:** A log-log plot of the total number of simple sequences in each of the eukaryotes plotted against simple sequence length.

**Figure 3:** Ratio of total number of simple sequences to the number of proteins in each proteome for various simple sequence lengths. **Panel a** is data for simple sequences of ten to twenty residues in length, **panel b** for twenty to forty, and **panel c** forty to sixty.

**Figure 4:** Ratio of the number of simple sequences found above that expected from the Poisson distribution model,  $\Delta_R$ , to the number of proteins in each eukaryote proteome, plotted for each residue type. Data for cysteine, methionine and tryptophan are excluded due to the low numbers of simple sequences found for these residues. **Panel a** shows  $\Delta_R$  for simple sequences ten to twenty residues in length, **panel b** twenty-one to forty residues, and **panel c** forty-one or more residues.

**Table 1:** Organisms surveyed for protein simple sequences, the number of proteins in each proteome, total number of simple sequences found ( $SS_{Tot}$ ) and the number of proteins containing at least one simple sequence ( $Prot_{SS}$ ).

Organism	Two-Letter Code	Type	Number of Proteins in Proteome	Total Number of Simple Sequences ( $SS_{Tot}$ )	Proteins with Simple Sequences ( $Prot_{SS}$ )	$SS_{Tot}/Prot_{SS}$
<i>S. cerevisiae</i>	SC	Eukaryote	6,203	7,177	3,293	2.18
<i>C. elegans</i>	CE		21,962	23,295	11,125	2.09
<i>D. melanogaster</i>	DM		13,608	24,725	7,989	3.09
<i>A. thaliana</i> *	AT		26,496	27,542	14,637	1.88
<i>Synechocytis</i> sp. PCC6803	SS	Cyanobacteria	3,169	1,493	1,034	1.44
<i>Nostoc</i> sp. PCC7120	Nos		5,368	2,497	1,762	1.42
<i>E. coli</i> K-12	EC	Gram-negative bacteria	4,289	2,064	1,636	1.26
<i>H. influenzae</i>	HI		1,709	614	476	1.29

<i>V. cholerae</i> chr1	<b>VC</b>		2,736	1,219	881	1.38
<i>H. pylori</i> 26695	<b>HP</b>		1,566	699	500	1.40
<i>B. melintensis</i> 16M chr1	<b>BM</b>		2,059	1,067	756	1.41
<i>A. tumefaciens</i> C58	<b>AgT</b>		2,722	1,610	1,078	1.49
<i>B. subtilis</i>	<b>BS</b>	Gram-positive bacteria	4,367	1,723	1,270	1.36
<i>B. halodurans</i>	<b>BH</b>		4,066	1,597	1,182	1.35
<i>M. pneumoniae</i>	<b>MP</b>		677	360	242	1.49
<i>M. genitalium</i>	<b>MG</b>		480	251	170	1.48
<i>D. radiodurans</i> chr1	<b>DR</b>		2,579	2,274	1,311	1.73
<i>C. acetobutylicum</i> ATCC824	<b>CA</b>		3,672	1,550	1,132	1.37
<i>A. fulgidus</i>	<b>AF</b>	Archaea	2,421	990	747	1.32
<i>A. pernix</i> K1	<b>AP</b>		2,694	1,827	1,176	1.55
<i>M. thermoautotrophicum</i>	<b>MT</b>		1,869	691	531	1.30

<i>M. janaschii</i>	<b>MJ</b>		1,715	875	641	1.36
<i>P. abyssi</i>	<b>PA</b>		1,765	847	613	1.38
<i>P. horikoshii</i>	<b>PH</b>		2,064	973	730	1.33
<i>Halobacterium</i> sp. NRC-1	<b>HS</b>		2,058	1,785	1,060	1.68
<i>T. acidophilum</i>	<b>TA</b>		1,478	511	395	1.29
<i>T. volcanium</i>	<b>TV</b>		1,478	452	376	1.20
<i>P. aerophilum</i>	<b>PAe</b>		2,605	1,220	866	1.41
<i>S. tokodaii</i>	<b>ST</b>		2,826	1,203	866	1.39
<i>S. solfataricus</i>	<b>SSol</b>		2,994	1,335	962	1.39

\* Some AT protein sequences were incomplete and were not included in the analysis. The number of proteins listed for AT corresponds to the number used.

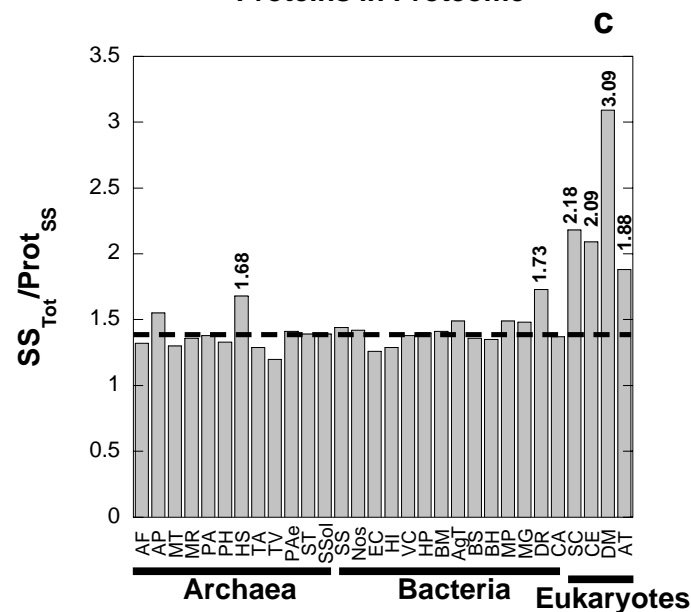
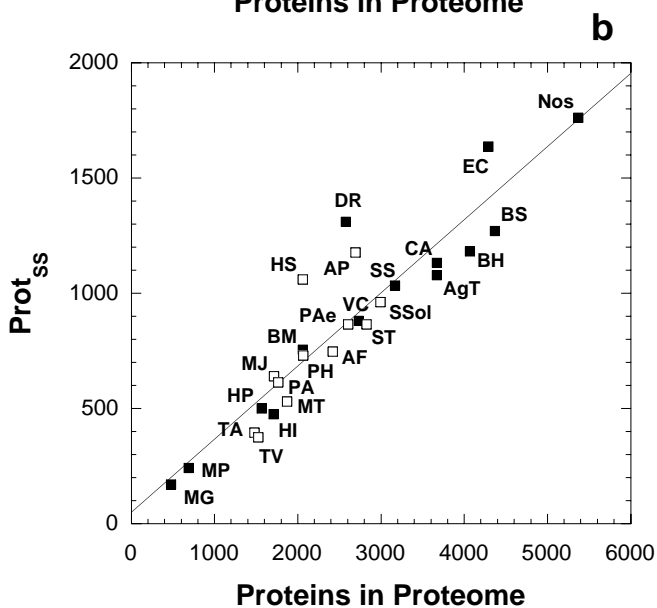
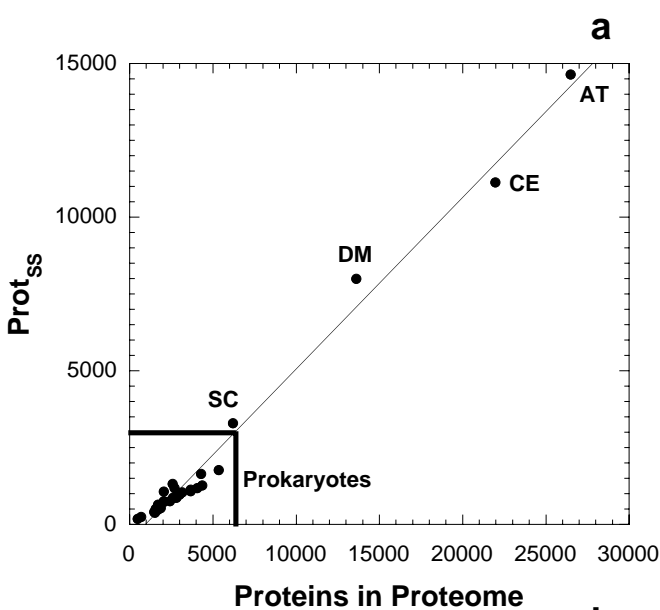
**Table 2:** Simple sequence distribution among proteins grouped according to class or process.

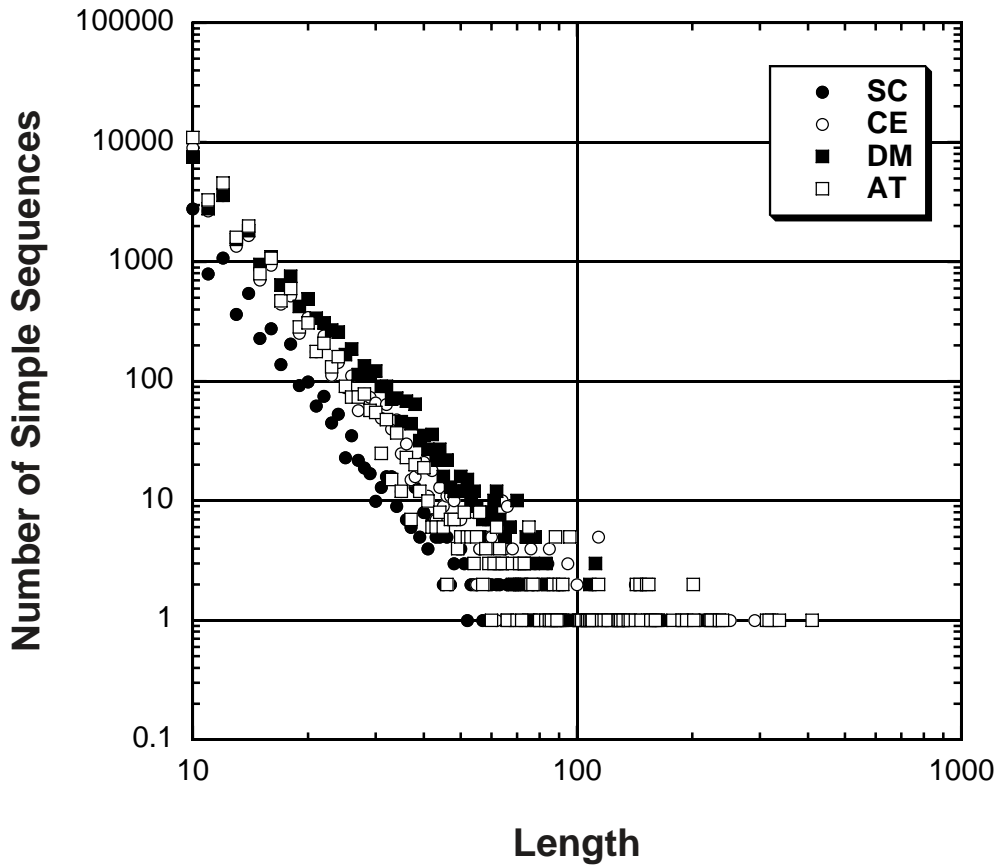
	<i>S. cerevisiae</i> (SC)		<i>C. elegans</i> (CE)		<i>D. Melanogaster</i> (DM)		<i>A. thaliana</i> (AT)	
<b>Keyword *</b>	<b>Proteins</b>	<b>SS<sub>found</sub></b>	<b>Proteins</b>	<b>SS<sub>found</sub></b>	<b>Proteins</b>	<b>SS<sub>found</sub></b>	<b>Proteins</b>	<b>SS<sub>found</sub></b>
<b>Cell cycle</b>	102	177	11	9	17	37	11	10
<b>Metabolism</b>	75	59	21	11	16	5	27	16
<b>Signal</b>	229	455	189	202	251	399	201	176
<b>Transcription</b>	274	677	87	101	177	838	105	130
<b>Transport</b>	440	449	148	85	102	133	160	118
<b>Membrane</b>	1,004	1,192	399	388	426	642	311	298
<b>Most Common Simple Sequence Type (Number Found)</b>								
	<i>S. cerevisiae</i> (SC)		<i>C. elegans</i> (CE)		<i>D. melanogaster</i> (DM)		<i>A. thaliana</i> (AT)	

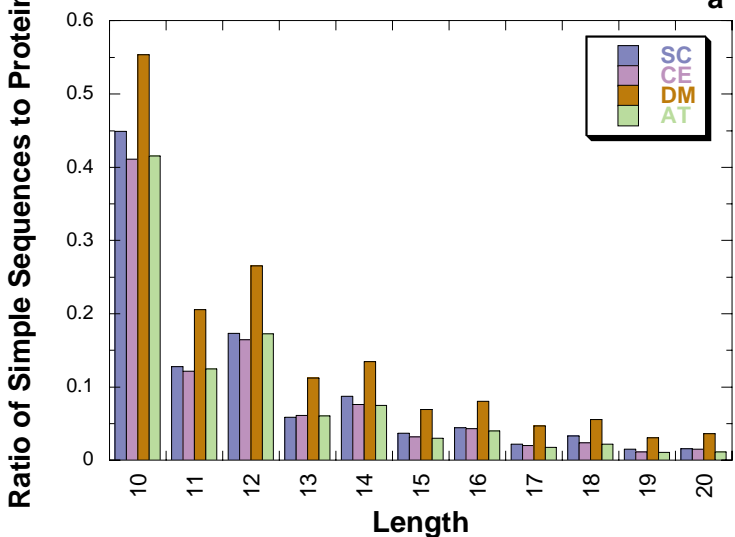
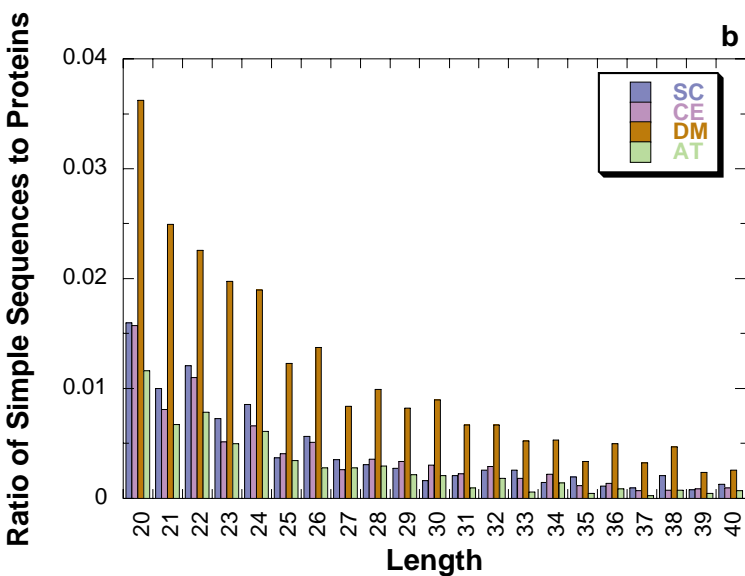
<b>Cell cycle</b>	S (64)	A (3)	S (11)	S (3)
<b>Metabolism</b>	A (13)	G (4), L (4)	V (2)	S (6)
<b>Signal</b>	S (173), T (133)	P (34), S (24), L (24), E (24)	L (63), S (55), Q (48), A (40)	P (61), L (42)
<b>Transcription</b>	S (138), N (114)	S (42)	S (194), Q (161), A (136), G (99)	S (35), A (21)
<b>Transport</b>	S (88), L (75)	A (16), L (16)	L (22), A (21)	A (32), S (20)
<b>Membrane</b>	S (283), L (235), T (111)	L (77), S(47)	L (178), S (77), G (71), A (63)	L (85), A (45)

\* Keyword used to search SWISS-PROT database for related proteins.







**a****b****c**