

1

Research article

2

3 ***k*-mer Distributions of Aminoacid Sequences are Optimised Across**
4 **the Proteome**

5

A.A. Morozov^{1*}

6

¹Cell Ultrastructure Department, Limnological Institute SB RAS

7

Address: 3, Ulan-Batorskaya, Irkutsk, 664033, Russia, P.O. box 278

8

* corresponding author. e-mail: morozov@lin.irk.ru

9

10 **Running title:** *Morozov A / k-mer Distributions in Proteins*

11

12

13 *Total word count:* 1589

14 *Figures:* 2

15 *Tables:* 0

16 *Supplementary materials:* None

17

18 **Abstract**

19 *k*-mer based methods are widely utilized for the analysis of nucleotide sequences and were
20 successfully applied to proteins in several works. However, the reasons for the species-specificity of
21 aminoacid *k*-mer distributions are unknown. In this work I show that performance of these methods
22 is not only due to orthology between *k*-mers in different proteomes, which implies the existence of
23 some factors optimizing *k*-mer distributions of proteins in a species-specific manner. Whatever
24 these factors could be, they are affecting most if not all proteins and are more pronounced in
25 structurally organized regions.

26

27 **KEYWORDS:** bayesian classifiers; composition bias; *k*-mer composition

28

29 **Introduction**

30 k-mer based methods are widely used in metagenomic studies because of their relatively low
31 computational cost compared to aligning reads to reference database. The exact algorithms vary
32 between implementations [1-3], but the idea is that *k*-mer spectra (or distributions) of
33 phylogenetically close taxa are more similar to each other than they are to those of more distant
34 groups. There is a plenty of empirical data to support this notion. The above-mentioned
35 metagenomic approaches perform rather well on both simulated and real datasets, and *k*-mer based
36 distance metrics have been used to reconstruct large-scale phylogenomic trees which were
37 consistent with trees produced by more orthodox methods [4].

38 Most of the *k*-mer-related work in bioinformatics was performed on nucleotide sequences,
39 but there is nothing inherently DNA-specific in this kind of analysis. There are works that have
40 translated *k*-mer based methods, initially designed for DNA, to proteomics. Using a distance metric
41 based on relative frequencies of *k*-mers, [5] have reconstructed a phylogenetic tree of 109 different
42 organisms from all major taxa. The topology of this tree does not contradict results produced by
43 other methods. In more recent work [6], a tree of approx. 900 bacteria with some eukaryotic
44 outgroups was built using a different distance metric, again pretty consistent with the consensus on
45 bacterial evolution. A recent metagenomic classifier named Kaiju [1] leverages protein
46 conservativity to classify sequences that don't have any close relatives in the reference database.
47 Thus, there is no question of whether *k*-mer distribution in aminoacid sequences is species-specific
48 or whether the divergence of these distributions correlates with evolutionary distances. However,
49 there is no answer to *why* it does.

50 The most common explanation relies on the orthology between *k*-mers in query sequence
51 and database. When the classifier is concerned with orthologous sequences, as eg in case of
52 classifying SSU RNA reads via RDP classifier [7], with sufficient value of *k* the chance of identical
53 *k*-mers appearing in non-homologous parts of sequences by random coincidence is negligible.
54 Somewhat similarly, protein-level metagenomic classification in Kaiju relies on finding MEMs

55 (maximum exact matches) and extending them to inexact shared k -mers. While not stated explicitly,
56 the phylogenetic importance of shared subsequences is also based on the orthology assumption.
57 However, performance of k -mer-based classifiers and distance metrics on divergent bacterial
58 proteomes with relatively few shared genes suggests there may be more to k -mer distribution than
59 MEMs. In this work I show that this specificity holds even in the complete absence of the
60 orthology.

61 **Results and Discussion**

62 Performance of the naïve bayesian classifier on CEGMA dataset is shown at fig.1. In
63 practically all cases this classifier performs better than random, and with optimal k of 5-7 more than
64 50% of sequences are assigned correctly. There is no possible orthology between sequences from
65 the same species' training and test sets. In fact, there is a risk that a protein from test set has an
66 ortholog in the *wrong* species' training set. k -mer distribution specificity persists even despite the
67 lack of orthology, which suggests that it is formed by species-specific factors on the proteomic
68 scale, rather than solely by the requirements of a particular protein family. Expanding the dataset to
69 the entire proteomes leads to precision skyrocketing to almost 100%. Although some part of the
70 precision increase can be explained by the presence of recently duplicated paralogs and isoforms, it
71 still suggests that most, if not all, proteins are affected by these factors.

72 To study the effect of these factors on a finer scale, we have built k -mer distributions for
73 protein features from the complete proteomes of the same species according to UNIPROT
74 annotations. Distances between the k -mer distribution of the feature in a particular species and the
75 summary distribution for this feature across the entire dataset were calculated. The higher this
76 distance, the more different these features in one organism are (on average) from their counterparts
77 from other species, which allows to use them as a proxy for the species-specificity of k -mer
78 distribution in protein fragments. As only structural features and entire domains have both average
79 length and feature counts sufficient for a reliable estimation of k -mer distribution, various binding
80 sites and signal peptides are omitted. Box-plots of these distances among different species are

81 shown at fig. 2.

82 For all structurally organised elements (*ie* helices and beta-strands) k -mer distributions are
83 more species-specific than they are for protein sequences as a whole (fig. 2), which means that
84 pressure for k -mer adaptation is greater in this regions. The same is true for functional domains,
85 whose k -mer distributions are optimised above protein-average level. This is strikingly similar to
86 codon usage adaptation on DNA level, where the use of different codons is regulating kinetics of
87 translation and folding. In particular, quickly translating high-frequency codons are common in
88 alpha helices, while rare, slower ones are more likely to be found in random coils [11]. Several
89 mechanisms can be proposed to explain this specificity on protein level. It's possible that the
90 evolutionary advantage or disadvantage of particular k -mers is related to protein creation specifics,
91 *eg* quicker and more efficient folding of optimal aminoacid sequence. Different aminoacid
92 composition can also be invoked as one of the explanations, although different frequencies of k -
93 mers with similar aminoacid composition prevent it from being considered the sole source of k -mer
94 distribution. Some of the specificity can be the effect of translating DNA with a specific distribution
95 of $3k$ -mers, which in turn is created by a range of DNA-specific factors such as GC-content, codon
96 usage, presence of specific sites like splicing regulators and so on. If the analogy with codon usage
97 bias is anything to go by, though, we should presume that there isn't a single source of selective
98 pressure on k -mer composition. All the factors described above probably apply to some degree, as
99 well as many others.

100 **Material and methods**

101 CEGMA dataset of highly conserved genes from six model eukaryotic species (*A. thaliana*,
102 *C. elegans*, *D. melanogaster*, *H. sapiens*, *S. cerevisiae*, *S. pombe*) was used. These are genes from
103 459 distinct orthogroups, each of which is represented by no more than one sequence from every
104 species, for a total of 456-458 proteins per species [8].

105 50 randomly selected proteins from each species were used as a test set, and naïve Bayesian
106 classifier (similar to multinomial classifier in [9]) was trained on the remaining ones. Test set

107 sequences were assigned to the proteomes using this classifier using for values of k between 3 and
108 10. Similar procedure was performed on the complete proteomes of these species, using 10% of
109 proteins randomly sampled as a testing set. All distances between k -mer distributions were
110 calculated using FFP distance metric [10].

111 **Authors' contributions**

112 AM has conceived the analysis, performed it and written the paper.

113 **Competing interests**

114 The author has declared no competing interest

115 **Acknowledgements**

116 This work was supported by the Federal Agency of Scientific Organisations (Russian
117 Federation) project #0345-2016-0031. Author is grateful to Dr. Y.P. Galachyants and A.N. Gurkov
118 for the valuable discussion.

119 **References**

- 120 1. Menzel P, Ng KL, Krogh A. Fast and sensitive taxonomic
121 classification for metagenomics with Kaiju. *Nature Communications* 2016; 7:11257
- 122 2. Wood DE, Saltzberg S. Kraken: ultrafast metagenomic sequence
123 classification using exact alignments. *BMC Genome Biology* 2014; 15:R46
- 124 3. Ounit R, Wanamaker S, Close TJ, Leonardi S. CLARK: fast and
125 accurate classification of metagenomic and genomic sequences using discriminative
126 k -mers. *BMC Genomics* 2015; 16:236
- 127 4. Sims GE, Kim SH. Whole-genome phylogeny of
128 *Escherichia/Shigella* group by feature frequency profiles (FFPs). *PNAS* 2011;
129 108(20):8329-8334
- 130 5. Qi, Wang, Hao. Whole Proteome Prokaryote Phylogeny Without
131 Sequence Alignment: a K-String Composition Approach. *Journal of Molecular*
132 *Evolution* 2004; 58:1-11

- 133 6. Jun SR, Sims GE, Wu GA, Kim SH. 2010 Whole-proteome
134 phylogeny of prokaryotes by feature frequency profiles: An alignment-free method
135 with optimal feature resolution. *PNAS* 2010; 107(1):133-138
- 136 7. Wang Q, Garrity GM, Tiedje JM, Cole JR. Naive Bayesian Classifier
137 for Rapid Assignment of rRNA Sequences into the New Bacterial Taxonomy.
138 *Applied and Environmental Microbiology* 2007; 73(16):5261-5267
- 139 8. Parra G, Bradnam K, Korf I. CEGMA: a pipeline to accurately
140 annotate core genes in eukaryotic genomes. *Bioinformatics* 2007; 23(17):1061-1067
- 141 9. Liu KL, Wong TT. Naive Bayesian Classifiers with Multinomial
142 Models for rRNA Taxonomic Assignment. *IEEE/ACM Transactions on*
143 *computational biology and bioinformatics* 2013; 10(5): 1334-1339
- 144 10. Sims GE, Jun SR, Wu GA, Kim SH. Alignment-free genome
145 comparison with feature frequency profiles (FFP) and optimal resolutions. *PNAS*
146 2009; 106(8):2677-2682
- 147 11. Angov E. 2011 Codon usage: Nature's roadmap to expression and
148 folding of proteins. *Biotechnology Journal* 2011; 6:650-659

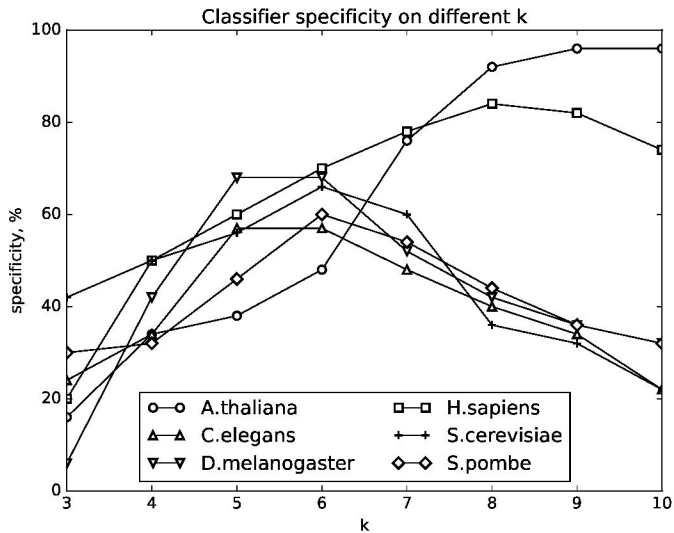
149

150 **Figure legends**

151 **Figure 1 Specificity of naive Bayesian classifier on CEGMA dataset under different**
152 **values of k .**

153 **Figure 2 Species-specificity of k -mer distribution on different features across six**
154 **proteomes. "Chain" feature represents protein sequence as a whole.**

155



Feature species-specificity

