

Comparison of genome sequences via projection extractor upon virtual mixer

Hongjie Yu ^{Corresp., 1}, Yuan-Ting Zhang ¹, Wei Fang ¹

¹ Department of Mathematics, School of Information and Network Engineering, Anhui Science and Technology University, Bengbu, Anhui, People Republic of China

Corresponding Author: Hongjie Yu
Email address: yhj70@mail.ustc.edu.cn

To compare multiple genome sequences, we transform each primary genome sequence into corresponding k -mer-based vectors. According to the principle of independent component analysis (ICA), the operation can be regarded as mixing multiple source genomic signals via several sensors, through which we can obtain the mixed vectors with equal-length from the corresponding genome sequences with different length. However, this mixing operation is performed by counting all the k -mer-based frequencies, instead of using real hardware of sensors. Thus, we name this preprocessing operation as virtual mixer (VM). Using ICA-based transformation, we projected all the vectors upon their independent components to capture the coefficients-based feature vector through the projection extractor (PE), which has been proved to have a property of distance preserving. Then, we used the proposed VMPE model upon three representative real datasets of genome sequence to test the efficiency for the model. The contrastive analysis results indicate that the proposed VMPE model performs well in similarity analysis.

1 Comparison of genome sequences via projection 2 extractor upon virtual mixer

3
4 Hong-Jie Yu, Yuan-Ting Zhang and Wei Fang

5 Department of Mathematics, School of Information and Network Engineering, Anhui Science
6 and Technology University, China

7
8 Corresponding Author:

9 Hong-Jie Yu

10 1501 Huangshan Avenue, Bengbu, Anhui Province, 233000, P.R. China

11 Email address: yhj70@mail.ustc.edu.cn

12 13 ABSTRACT

14 To compare multiple genome sequences, we transform each primary genome sequence into
15 corresponding k -mer-based vectors. According to the principle of independent component analysis
16 (ICA), the operation can be regarded as mixing multiple source genomic signals via several sensors,
17 through which we can obtain the mixed vectors with equal-length from the corresponding genome
18 sequences with different length. However, this mixing operation is performed by counting all the k -mer-
19 based frequencies, instead of using real hardware of sensors. Thus, we name this preprocessing
20 operation as virtual mixer (VM). Using ICA-based transformation, we projected all the vectors upon
21 their independent components to capture the coefficients-based feature vector through the projection
22 extractor (PE), which has been proved to have a property of distance preserving. Then, we used the
23 proposed VMPE model upon three representative real datasets of genome sequence to test the efficiency
24 for the model. The contrastive analysis results indicate that the proposed VMPE model performs well in
25 similarity analysis.

26 27 INTRODUCTION

28 In the field of bioinformatics, sequence comparison aims to discover relationship of similarity
29 among various biological sequences. Sequence comparison of several genomes at the nucleotide
30 sequence level can be accomplished by multiple sequence alignment (MSA). Recently, a web-
31 based alignment services have been proposed (*Nguyen, 2012*). Generally, based on an
32 appropriate base-substitution model, the MSA are used to compute similarity scores. However,
33 since species diverge widely over time, both genomic rearrangements and insertions/deletions
34 make MSA a little difficult in genome comparison.

35 In fact, at present, the genomics necessarily requires approaches within the non-coding area of
36 genomic sequence. Obviously, it can be used to compare the whole genomes (including coding
37 and non-coding regions). It is an urgent need to develop a method which can be free of a specific
38 gene set, and can also analyze nongenic regions. Later, alignment-free approaches are
39 successively proposed to achieve this target (*Deng et al., 2011, Dong et al., 2018, Gao & Qi,*
40 *2007, Vinga & Almeida, 2003, Yu et al., 2010*). The representative method is involved in

41 frequency-based method called feature frequency profiles (FFP) (*Jun et al., 2010, Sims et al.,*
42 *2009, Sims et al., 2009, Sims & Kim, 2011*), which can be used as the comparison of whole
43 genomes or genomic segments that may not be closely related and have latent remarkable
44 rearrangement or have not shared a common set of genes, e.g. regulatory, intronic or nongenic
45 regions.

46 As a classical alignment-free method, Blaisdell first introduced k -mer for biological sequence
47 comparison (*Blaisdell, 1986*). The relative approaches can be found in the reference (*Luczak et*
48 *al., 2017*). Vinga reviewed more developments about alignment-free comparison (*Vinga &*
49 *Almeida, 2003*), which described many developed approaches to mining data and comparing
50 multiple sequences. As a variation of approach for text comparison, Sims et al. (*Sims et al., 2009*)
51 insist that, for every two different texts, the ‘distance’ between these two word frequency profiles
52 can be regarded as a means of the dissimilarity between corresponding two texts. However,
53 because there are no ‘words’ within a long string of base-pairs that constitute genome sequences,
54 the differences among relative k -mer frequencies can be used to calculate distance values. For a
55 given length sequence, the frequency information for all of the possible features (k -mers) is
56 assembled into a FFP, where either the resolution or length of the features is the most important
57 parameter.

58 The information theory (IT) has been widely used in the area of computational biology. Based
59 on the use of digital signal processing (DSP) theory and algorithms, using genomic signal
60 processing (GSP), one can analyze DNA or protein sequences. GSP methods convert DNA data
61 to numerical values, thus it would offer the opportunity of applying existing DSP methods for
62 genomic data. Examples of the use of GSP methods include performing cluster analysis of
63 biological sequences and the K-means algorithm, proposing visualization method to inspect and
64 analyze possible hidden behaviors (*Mendizabal-Ruiz et al., 2018*).

65 At the nucleotide, codon and amino acid levels, the researcher has derived optimal symbolic-
66 to-digital mappings of the linear, nucleic acid strands into real or complex genomic signals
67 (*Cristea, 2002*). The proposed approach converts the sequences of polypeptides or nucleotides
68 into digital genomic signals, and provides the possibility for using a large variety of signal
69 processing approaches to their handling and analysis. It is also shown that using this
70 representation one can better extracted some essential features of the nucleotide sequences.

71 Alignment-free sequence comparison and analysis greatly benefited from concepts which were
72 derived from IT, such as mutual information and entropy. Within the review (*Vinga, 2014*), the
73 author investigated many aspects of IT applications, such as resolution-free metrics based on
74 iterative maps, block-entropy estimation, prediction of transcription factor binding sites,
75 comprising the classification of motifs and sequence characterization based on linguistic
76 complexity and entropic profiles. As a function of the genomic location, the Entropic Profiler
77 (EP) captures the essential region with respect to the whole genome (*Comin & Antonello, 2013,*
78 *Fernandes et al., 2009*).

79 A method (*Stuart et al., 2002*) have been developed to produce comprehensive gene and to
80 reconstruct species phylogenies from the unaligned whole genome data, where one can use the

81 singular value decomposition (SVD) approach to character string frequencies analysis. Within
82 the SVD-based dimension-reduced space, a quantitative comparison for the relative orientations
83 of protein vectors provided straightforward and accurate estimations of sequence similarity,
84 which can in turn be used to construct comprehensive gene trees. Later (*Stuart & Berry, 2004*),
85 using this SVD-based alignment-free method, the authors compared the predicted protein
86 complement of 9 whole eukaryotic genomes ranging from yeast to man. They got simultaneous
87 identification and definition of a large number of well conserved motifs and gene families.
88 Meanwhile, a species tree was constructed, which supports one of two conflicting hypotheses for
89 metazoan relationships. Based on SVD, the analysis of the entire protein complement of 9 whole
90 eukaryotic genomes suggests that highly conserved motifs and gene families can be identified,
91 and one can effectively compare within a single coherent definition space, where the extraction
92 of gene and species trees can be easily implemented. The analysis can provide a basis for these
93 definitions, when there is no explicit definition of orthologous or homologous sites.

94 *Comon* developed *Independent Component Analysis* (ICA) to find a linear representation of
95 non-gaussian data so that all the components are statistically independent with each other, or as
96 independent as possible (*Comon, 1994*). Using such a representation, it seems to capture the
97 essential structure of the data in several applications, including signal separation and feature
98 extraction.

99 As an no-alignment method, the composition vector (CV) method (*Chan et al., 2012*) has
100 been extensively studied recently. The abilities for ICA in feature extraction (*Huang & Zheng,*
101 *2006, Yeredor, 2002*), plus the successful use of composition vector or *k*-mer method in the
102 comparison of sequences inspire us to combine them for improving the performance on
103 similarity analysis.

104 In this study, we propose an ICA-based model for similarity analysis of genome sequences,
105 i.e., virtual mixer & projection extractor model (VMPE), where we extracted the projections as
106 the features for genome through optimizing the model. We also cluster the genome upon three
107 real datasets, and visually inspect features of genome and analyze their similarities. Our results
108 indicate the feasibility of applying the proposed method to compare genomes.

109

110 **MATERIALS & METHODS**

111 **Model for the feature extraction from genome sequences**

112 Through transforming the sequences of polypeptides or nucleotides into digital genomic signals,
113 several methods offer the possibility to employ a large variety of signal processing approaches to
114 deal with and analyze the sequences (*Cristea, 2002*). In this section, we propose a novel model to
115 extract information from genome sequences, and investigate the properties of the model.

116 In the proposed approach, the procedure is composed of four stages:

- 117 a) Transforming each primary genome sequence into corresponding *k*-mer-based vectors via
118 our designed 'virtual mixer' (VM).
- 119 b) Projecting all these 'mixed' vectors upon independent component coordinate system to
120 extract features through the hierarchy 'projective extractor' (PE);

121 c) Optimizing the number of segments s^* for the best segmentation scheme; and d) Applying
 122 the final dimension-reduced feature vectors obtained by our hierarchy VMPE model on the
 123 similarities analysis among genome sequences. The application on real dataset demonstrated
 124 the validity of our proposed approach.

125 **Preprocessing of sequences via K-mer-based virtual mixer (VM)**

126 Hao Bailin's laboratory has proposed k-mer-based composition vector (CV) approach, where
 127 background 'noise' can be subtracted via a Markov chain estimator. Based on this no-alignment
 128 approach, the author obtained valuable results for both genome and protein sequences (*Gao & Qi,*
 129 *2007, Qi et al., 2004*).

130 **Description for k-mer approach**

131 As for the k-mer approach, a description of the details can be described as follows. Let s be the
 132 primary genome sequence with length L , $s = N_1 N_2 \cdots N_L$, where $N_l \in \{A, T, G, C\}$,
 133 $l = 1, 2, \dots, L$.

134 Generally, a k-mer mode is a series of k consecutive characters within a sequence. The usual
 135 handling way to count k-mers within a sequence with length L is to use a sliding window with
 136 length k , shifting one frame base each time from position 1 to $L-k+1$, until all the entire genomes
 137 have been scanned. Thus, the k-mer-based feature vector can be denoted as:

$$138 \bar{m}^{(i)} = (m_{i1}, m_{i2}, \dots, m_{i,4^k}), \quad (1)$$

139 where m_{ij} denotes the frequency of the counting number for each corresponding pattern with k -
 140 consecutive characters, and i indicates the label of genome sequence. Meanwhile, all the $\bar{m}^{(i)}$
 141 have been normalized by its own length, i.e. $L-k+1$, which can keep the obtained feature vector
 142 be freed from the negative influence of different lengths for each genome sequence.

143 The number of the order for k-mer mode is just denoted as k . Since the parameter k has a great
 144 influence on the results of sequence analysis, it is a key to pick out an appropriate k . Some works
 145 have investigated the selection of k . For example, Wu et al. (*Wu et al., 2005*) proposed an
 146 optimal word size to calculate dissimilarity, where the optimal word size can be determined by
 147 length of sequence considered. Another solution was investigated by Sims et al. (*Sims et al.,*
 148 *2009, Sims et al., 2009*), who gave the lower limit and upper limit for the range of optimal length.
 149 The lower limit can be approximately determined by the average length of multiple sequences,
 150 named as \bar{L} . Then one can let k be the integer part of $\log_4(\bar{L})$.

151 **Observing the mixed genomic signals from sequence**

152 As proposed by Pierre Comon (*Comon, 1994*), Independent Component Analysis (ICA) can
 153 lessen the effects of noise or artifacts within the data as it focuses on separating a mixture of
 154 signals into their each different sources, respectively. ICA models observations as linear
 155 combinations of several certain components, or variables, which are selected as statistically
 156 independent as possible, i.e. the different components represent different non-overlapping
 157 information.

158 Thus, we can regard the feature vector $\bar{m}^{(i)}$ in Eq. (1) as the mixture from several independent
 159 components (ICs). However, in this study, the course of mixture need not perform through real

160 hardware, i.e. just via virtual k-mer-based counter.

161 Assuming further that there are several virtual ‘sensors’ or virtual ‘receivers’, these sensors
162 have different ‘inherent frequencies’ so that each ‘sensor’ records a mixture of the original
163 source genomic signals with slightly different patterns.

164 The observed data m_{ij} are obtained through “sensor”, i.e., virtual mixer (VM) that gives the
165 mixing weights among several latent independent components. In the left side hand of Eq. (1),
166 each observed vector $\bar{m}^{(i)}$ can be expressed as a linear combination of n latent components $\bar{c}^{(1)}$,
167 $\bar{c}^{(2)}$, ..., and $\bar{c}^{(n)}$, where i just denotes the label for the genome sequence. All that we directly
168 observed are the n genomic signals, $\bar{m}^{(1)}$, $\bar{m}^{(2)}$, ..., and $\bar{m}^{(n)}$. So the combination can be
169 represented as follows:

$$170 \bar{m}^{(i)} = f_{i1} \cdot \bar{c}^{(1)} + f_{i2} \cdot \bar{c}^{(2)} + \dots + f_{in} \cdot \bar{c}^{(n)}, \quad (2)$$

171 where f_{ij} are some coefficients that define the representation, $i = 1, 2, \dots, n$, and i is just the label
172 for the original sequence or species. Also, $\bar{c}^{(1)}$, $\bar{c}^{(2)}$, ..., and $\bar{c}^{(n)}$ are as independent as possible
173 with each other.

174 How can we estimate the coefficients f_{ij} in Eq. (2)? Collecting all these coefficients f_{ij} into a
175 matrix F , we want to use some general statistical properties to find the matrix F through which
176 we can represent the observed genomic signals by several latent independent signals. Then, the
177 question can just boil down to the very ICA-based problem which is started with how to find a
178 good representation of multivariate data.

179 **Extracting features from mixed genomic signals**

180 However, a general solution to the problem can be found by using independent component
181 analysis (ICA) upon the observed data $\bar{m}^{(1)}$, $\bar{m}^{(2)}$, ..., $\bar{m}^{(n)}$, which are likewise collected into the
182 k -mer-based matrix K .

183 The mixed k -mer-based matrix K comprises n counting vectors as follows:

$$184 \begin{pmatrix} \bar{m}^{(1)} \\ \bar{m}^{(2)} \\ \vdots \\ \bar{m}^{(n)} \end{pmatrix} = \begin{pmatrix} m_{11} & m_{12} & \cdots & m_{1,4^k} \\ m_{21} & m_{22} & \vdots & m_{2,4^k} \\ \vdots & \vdots & \ddots & \vdots \\ m_{n1} & m_{n2} & \cdots & m_{n,4^k} \end{pmatrix} = K, \quad (3)$$

185 where $i = 1, 2, \dots, n$, and $\bar{m}^{(i)}$ is a 4^k dimensional vector comprising of all the mixed elements
186 transformed from the corresponding sequence $s^{(i)}$ via the virtual mixer (VM)

187 Thus, we can estimate the coefficient matrix F for the genomic signals from original genome
188 sequences (The elements f_{ij} in coefficient matrix F can be estimated by the developed algorithm,
189 e.g. FastICA). Thus, Eq. (1) can be regarded as a statistical “latent variables” model. These latent
190 components are assumed *unknown*, because we cannot know the values of c_{ij} without knowing
191 all the properties of the virtual extractor system, which can be difficult in general. The projection
192 coefficients f_{ij} are *unknown as well*, which we cannot record directly. The problem can then be

193 rephrased as that of how to determine the coefficients f_{ij} within matrix F .

194 Moreover, it is usually more convenient to introduce a vector-matrix notation instead of all the
 195 n equations in Eq. (2). Let us denote the three series of variables by the corresponding row
 196 vectors, respectively. Thus, according to the block-matrix appearance, the relationship of K and
 197 F can be written as:

$$198 \begin{pmatrix} \bar{m}^{(1)} \\ \bar{m}^{(2)} \\ \vdots \\ \bar{m}^{(n)} \end{pmatrix} = \begin{pmatrix} \bar{f}^{(1)} \\ \bar{f}^{(2)} \\ \vdots \\ \bar{f}^{(n)} \end{pmatrix} \times \begin{pmatrix} \bar{c}^{(1)} \\ \bar{c}^{(2)} \\ \vdots \\ \bar{c}^{(n)} \end{pmatrix}, \quad (4)$$

199 where i is just the label for different genome sequence, $i = 1, 2, \dots, n$, and $\bar{c}^{(i)}$ are n independent
 200 4^k dimensional vectors which are obtained by using ICA upon all the n observed 4^k
 201 dimensional row vectors in matrix K .

202 Virtual mixer & projection extractor (VM-PE)

203 VM-PE model for capturing genomic signals

204 Below are the descriptions for the VM-PE model, which is based on k -mer and ICA.

205 Considering two transformations:

206 (a) $S \xrightarrow{\tau_1} K$

207 The matrix S denotes all the n involved primary genome sequences with different length, while
 208 the matrix $K \in \mathcal{R}^{n \times 4^k}$ stands for the corresponding matrices which were mapped from all these n
 209 primary sequences. (Please see the first box shown in Fig. 1).

210 (b) $K \xrightarrow{\tau_2} F$

211 The matrix F comprises of all the n extracted feature vectors:

$$212 \bar{f}^{(i)} = (f_{i1}, f_{i2}, \dots, f_{in}), \quad (5)$$

213 which are transformed from k -mer-based mixer matrix K via ICA-based projection extractor (PE)
 214 (Please see the second box shown in Fig. 1).

215 Likewise, if we collect all the feature vectors into a matrix F , then Eq. (5) can be rewritten as:

$$216 \begin{pmatrix} \bar{f}^{(1)} \\ \bar{f}^{(2)} \\ \vdots \\ \bar{f}^{(n)} \end{pmatrix} = \begin{pmatrix} f_{11} & f_{12} & \cdots & f_{1n} \\ f_{21} & f_{22} & \vdots & f_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ f_{n1} & f_{n2} & \cdots & f_{nn} \end{pmatrix}^A = F, \quad (6)$$

217 Thus, the relationship between the matrices K and F shown in Fig. 1 can be described as:

$$218 K = F * C. \quad (7)$$

219 That is,

$$\begin{matrix} 220 \\ 221 \\ 222 \\ 223 \\ 224 \end{matrix}
 \begin{pmatrix} m_{11} & m_{12} & \cdots & m_{1,4^k} \\ m_{21} & m_{22} & \cdots & m_{2,4^k} \\ \vdots & \vdots & \ddots & \vdots \\ m_{n1} & m_{n2} & \cdots & m_{n,4^k} \end{pmatrix} = \begin{pmatrix} f_{11} & f_{12} & \cdots & f_{1n} \\ f_{21} & f_{22} & \cdots & f_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ f_{n1} & f_{n2} & \cdots & f_{nn} \end{pmatrix} \times \begin{pmatrix} c_{11} & c_{12} & \cdots & c_{1,4^k} \\ c_{21} & c_{22} & \cdots & c_{2,4^k} \\ \vdots & \vdots & \ddots & \vdots \\ c_{n1} & c_{n2} & \cdots & c_{n,4^k} \end{pmatrix} \quad (8)$$

Obviously, the matrix F belongs to the space $\mathcal{R}^{n \times n}$, where the dimension of all the row vectors is uniform, i.e. n . Moreover, compared with matrix K , it can be seen from Eqs. (7) or (8) that the dimension for feature matrix F is greatly reduced. Thus, the final feature matrix F deserves our special attention.

As shown in Fig. 1, we design a k -mer-based virtual mixer (VM) and ICA-based projection extractor (PE), through which we can extract the dimension-reduced essential features from the genome sequences.

228

229 **Ensemble transformation**

In fact, all the n obtained feature vectors $\bar{f}^{(i)}$, $i = 1, 2, \dots, n$, can be regarded as the projections, which are captured by projecting the k -mer-based ‘observed variables’ $\bar{m}^{(i)}$ upon the n latent independent 4^k dimensional vectors $\bar{c}^{(j)}$ in matrix C , where f_{ij} , $i, j = 1, 2, \dots, n$, are some real coefficients. Meanwhile, as the description of Eq. (2), all $\bar{c}^{(i)}$ are statistically mutually independent.

In summary, as a whole, the compound transformation can be described as follows:

$$\begin{matrix} 230 \\ 231 \\ 232 \\ 233 \\ 234 \end{matrix}
 \tau_2 \circ \tau_1 : s^{(i)} \mapsto \bar{f}^{(i)} = (f_{i1}, f_{i2}, \dots, f_{in}), \quad (9)$$

through which one can freely extract the features of the multiple genome sequences. Also, the ensemble transformation can be depicted as:

$$\begin{matrix} 235 \\ 236 \\ 237 \\ 238 \end{matrix}
 \text{Ker } \varphi : S^{n \times 1} \xrightarrow{\tau} \mathcal{R}^{\frac{n(n-1)}{2}} \quad (10)$$

where $S^{n \times 1}$ stands for the *original string sequence space* comprising of n genome sequences with different length, while $\mathcal{R}^{\frac{n(n-1)}{2}}$ denotes the *objective distance space* transformed from the *original space* via the proposed VM-PE scheme shown in Fig. 1.

243 **Distance-preserving properties of transformation**

As can be seen from the following proposition, the above-mentioned compound transformation embodies the essential property of the genome sequences. Thus, the φ can be interpreted as a kernel operator. Some properties for the proposed VM-PE algorithm can be depicted as:

Definition 1: Within the *original string sequence space*, the distance between every two different sequences $D(s^{(i)}, s^{(j)})$, is defined as:

$$\begin{matrix} 244 \\ 245 \\ 246 \\ 247 \\ 248 \\ 249 \end{matrix}
 D(s^{(i)}, s^{(j)}) \stackrel{\text{def}}{=} \text{corr}(\bar{m}^{(i)}, \bar{m}^{(j)}) \quad (11)$$

where $\bar{m}^{(i)}$ is the k -mer-based mixed vector for sequence $s^{(i)}$, and so is $\bar{m}^{(j)}$, $i, j = 1, 2, \dots, n$. Here,

251

252 the function $corr(\cdot, \cdot)$ denotes the correlation degree between two vectors, i.e., $\bar{\mathbf{m}}^{(i)}$ and $\bar{\mathbf{m}}^{(j)}$.

253 Usually, it can be defined by the values of correlation coefficients from a pair of vectors.

254 **Definition 2:** Let \mathbf{R}^d be a real normed space with dimensions d , and let $\varphi: \mathcal{S}^{n \times 1} \mapsto \mathbf{R}^{1 \times d}$ be a

255 function from $\mathcal{S}^{n \times 1}$ to $\mathbf{R}^{1 \times d}$, where $d = \frac{n \cdot (n-1)}{2}$. For any element within the space $\mathcal{S}^{n \times 1}$, such as

256 $\mathbf{s}^{(i)}$ and $\mathbf{s}^{(j)}$, if $D(\mathbf{s}^{(i)}, \mathbf{s}^{(j)}) = \delta$ implies $D(\varphi(\mathbf{s}^{(i)}), \varphi(\mathbf{s}^{(j)})) = \delta$, the function φ can be called δ -

257 distance preserving.

258 **Theorem 1:** $\mathcal{T}_2 \circ \mathcal{T}_1: \mathbf{s}^{(i)} \mapsto \bar{\mathbf{f}}^{(i)} = (f_{i1}, f_{i2}, \dots, f_{in})$ is a distance-preserving transformation.

259 **Proof:** Since $\bar{\mathbf{m}}^{(i)}$ and $\bar{\mathbf{m}}^{(j)}$ are the k -mer-based mixed vectors of two corresponding sequences $\mathbf{s}^{(i)}$
 260 and $\mathbf{s}^{(j)}$, respectively, $i, j = 1, 2, \dots, n$.

261 Let $\mathcal{T}_2 \circ \mathcal{T}_1$ be the compound transformation from the original string sequence space to
 262 objective distance space.

263 Then, the following equations:

$$264 \varphi(\mathbf{s}^{(i)}) = (\mathcal{T}_2 \circ \mathcal{T}_1)(\mathbf{s}^{(i)}) = \mathcal{T}_2[\mathcal{T}_1(\mathbf{s}^{(i)})] = \mathcal{T}_2(\bar{\mathbf{m}}^{(i)}) = \bar{\mathbf{f}}^{(i)} \quad (12)$$

265 will hold.

266 According to Eq. (2), we can rewrite them in summation form as follows:

$$267 \bar{\mathbf{m}}^{(i)} = \sum_{p=1}^n f_{ip} \cdot \bar{\mathbf{c}}^{(p)}, \quad (13)$$

268 while

$$269 \bar{\mathbf{m}}^{(j)} = \sum_{q=1}^n f_{jq} \cdot \bar{\mathbf{c}}^{(q)}. \quad (14)$$

270 Thus, there have

$$271 D(\mathbf{s}^{(i)}, \mathbf{s}^{(j)}) = D(\bar{\mathbf{m}}^{(i)}, \bar{\mathbf{m}}^{(j)}) = corr(\bar{\mathbf{m}}^{(i)}, \bar{\mathbf{m}}^{(j)}) = corr\left(\sum_{p=1}^n f_{ip} \cdot \bar{\mathbf{c}}^{(p)}, \sum_{q=1}^n f_{jq} \cdot \bar{\mathbf{c}}^{(q)}\right)$$

$$272 = \sum_{p=1}^n \sum_{q=1}^n (corr(f_{ip} \cdot \bar{\mathbf{c}}^{(p)}, f_{jq} \cdot \bar{\mathbf{c}}^{(q)}))$$

$$273 = \sum_{\substack{p=1 \\ p=q}}^n (corr(f_{ip} \cdot \bar{\mathbf{c}}^{(p)}, f_{jq} \cdot \bar{\mathbf{c}}^{(q)})) + \sum_{\substack{p=1 \\ p \neq q}}^n \sum_{q=1}^n (corr(f_{ip} \cdot \bar{\mathbf{c}}^{(p)}, f_{jq} \cdot \bar{\mathbf{c}}^{(q)})) \quad (15)$$

274 Since, every two independent component $\bar{\mathbf{c}}^{(p)}$ and $\bar{\mathbf{c}}^{(q)}$, $p \neq q$, are statistically mutually
 275 independent, then, according to the calculation for the correlation coefficients, it can be obtained
 276 that the second part within the right-hand side of Eq. (15) is just zero.

277 And hence, the whole Eq. (15) will be simplified into:

$$278 D(\mathbf{s}^{(i)}, \mathbf{s}^{(j)}) = \sum_{p=1}^n corr(f_{ip} \bar{\mathbf{c}}^{(p)}, f_{jp} \bar{\mathbf{c}}^{(p)}) = corr(\bar{\mathbf{f}}^{(i)}, \bar{\mathbf{f}}^{(j)}), \quad (16)$$

279 where $i, j = 1, 2, \dots, n$.

280 By Eq. (12) and *Definition 1*, we have:

$$281 D(\varphi(s^{(i)}), \varphi(s^{(j)})) = D(\bar{f}^{(i)}, \bar{f}^{(j)}) = \text{corr}(\bar{f}^{(i)}, \bar{f}^{(j)}). \quad (17)$$

282 Comparing Eq. (16) with Eq. (17), by *Definition 2*, the following equation holds:

$$283 D(s^{(i)}, s^{(j)}) = D(\varphi(s^{(i)}), \varphi(s^{(j)})), \quad i, j = 1, 2, \dots, n$$

284 Therefore, for two given genome sequences, $s^{(i)}$ and $s^{(j)}$, it can be seen that the above-
285 mentioned ensemble transformation τ indeed has a property of δ -distance preserving. Here,
286 superscript i or j denotes the label of sequences, and the sequences are usually non-equal length.

287 *QED*

288

289 According to *Theorem 1* and *Definition 2*, we can calculate all the row vectors for matrix F ,

290 such as $\bar{f}^{(i)} = (f_{i1}, f_{i2}, \dots, f_{in})$, $i = 1, 2, \dots, n$, where i denotes the label of n primary genome or

291 proteome sequences. Then we can get n corresponding combinational coefficients vectors with

292 the dimension n , which can be regarded as the features extracted from the original genome

293 sequence via the proposed VM-PE algorithm, whose steps can be summarized as follows.

294

Input: multiple genome sequences with different length: $s^{(1)}, s^{(2)}, \dots, s^{(n)}$

begin

for $i=1$ **to** n **do**

Through virtual mixer (VE), transform each genome sequence $s^{(i)}$ into a k -mer-based 4^k
dimensional vector $\bar{m}^{(i)}$, which comprises n by 4^k observed matrix K

end for

Using FastICA-based projection extractor (PE), factorize matrix K into independent
component matrix C which is left multiplied by projection feature matrix F

for $i=1$ **to** $n-1$ **do**

for $j=i+1$ **to** n **do**

Calculate pairwise distances using F by $D(s^{(i)}, s^{(j)}) = \|\bar{f}^{(i)} - \bar{f}^{(j)}\|_2$

end for

end for

for $i=1$ **to** n **do**

The feature vector for the i -th genome sequence \leftarrow The i -th row vector of feature matrix F

end for

Draw the dendrogram using the pairwise distances matrix

end

295

296 RESULTS

297 APPLICATION FOR THE REAL GENOME DATASET

298 We apply the proposed VMPE model upon the real genome dataset.

299 Data preparation

300 Table 1 shows the concise information for these 20 species in the GenBank. The first dataset
301 includes the mitochondrial genome sequences from 20 eutherian species, which has also been
302 investigated by several works (*Dai et al., 2011, Deng et al., 2011, Huang & Wang, 2011, Huang*
303 *et al., 2011, Yang et al., 2012, Yu & Huang, 2012*).

304 Determining k^* for virtual mixer (VM)

305 Given multiple genome sequences with average length about 17,000 bps, the k -mer count
306 vector $\bar{m}^{(i)}$ becomes too sparse for $k \geq 8$ (*Yu, 2013*). Thus, based on the higher order k -mer, the
307 comparisons among multiple long sequences may not capture the essential feature of sequences.
308 Therefore, we need to consider how to determine a rational order k^* for k -mers according to the
309 approximate formula investigated in Ref. (*Sims et al., 2009*). Therefore, for the above-mentioned
310 dataset, the integer k^* can be calculated as:

$$311 k^* \approx \log_4(17000), \quad (18)$$

312 i.e. $k^*=7$.

313 Therefore, for the given genome sequence $s^{(i)}$, $k^*=7$ is just the optimal order of k -mer for
314 VMPE model at the stage of virtually mixing, where all the original sequences can be
315 transformed into the corresponding mixed 4^7 -dimensional vectors. In other words, through
316 virtual mixer, all the obtained mixed 7-mer vector $\bar{m}^{(i)}$ should uniformly be located at in
317 dimensional space of 16384, $i=1, 2, \dots, n$.

318 Comparison of VMPE model with other works

319 We apply the VM-PE model upon the above-mentioned data set.

320 In order to assure our results, we list the observations for pairwise distance between Human
321 and every one of the rest 19 species. As shown in Table 2, Columns 2~4 are the observations
322 which are extracted from the pairwise distances matrices calculated by different approaches, i.e.
323 Kolmogorov complexity (*Li et al., 2001*), OPT-SVD (*Yu & Huang, 2012*) and our proposed
324 VMPE model, respectively.

325 In general, the correlation degree between every two different results from each case is an
326 effective way for comparing every two different approach. The higher correlation degree with
327 the results of the existing approach and a newly developed one means that the latter has the same
328 efficiency as the former.

329 For Kolmogorov complexity (*Li et al., 2001*), the calculated correlation coefficient value,
330 between Column 2 and Column 4, is 0.95. Likewise, we can compute correlation coefficient
331 value between Column 3 and Column 4, which is 0.8312. These phenomena indicate that our
332 proposed VMPE model achieves the same effect as those representative approaches to the
333 similarity analysis of genome sequences.

334 Phylogeny analysis of genomes via VMPE model

335 Generally, the steps for phylogenetic analysis among multiple genome sequences can be
336 described as follows:

- 337 (1) Firstly, we can calculate all the corresponding projection feature vectors (FV) with
338 different segmentation schemes for each genome sequence via VMPE;

- 339 (2) Secondly, selecting ‘euclidean’ distance metric, we calculate pair-wise distance matrix;
340 (3) Finally, we investigate the effect of dendrogram.

341 Fig. 2 illustrates the dendrogram derived from our proposed VMPE model, where these 20
342 species are separated clearly:

- 343 (a) Outgroup (Platypus, (Opossum, Wallaroo)) is far away from other clusters;
344 (b) Seven Primates are grouped together;
345 (c) Two Rodents (Mouse and Rat) stand at the same branch;
346 (d) The rest ones are clustered into a close group.

347 Meanwhile, these results are in agreement with both the evolutionary facts and the conclusions
348 in (*Li et al., 2001, Yu & Huang, 2012*). Thus, it is shown that the proposed approach is effective
349 in multiple genome sequences comparison. However, the result suggests the hypothesis of
350 (Primates, (Rodents, Ferungulates)), which can also be found in Ref. (*Cao et al., 1998*), where
351 there is a controversy on the hypothesis.

352 As a contrast, using MEGA software, we rebuild the Neighboring-Joining (NJ) tree through
353 alignment-based method. As for the tree shown in Fig. 3, when compared to Fig. 2, it can be
354 found that the proposed VMPE model yields similar results to that from the traditional approach.

355 **Application upon another two genome datasets**

356 To further verify the efficiency for our proposed VMPE model, we select another two genome
357 datasets: 1) with 34 mammalian sequences, which has been also investigated by many works
358 (*Huang et al., 2011, Yu et al., 2010, Yu, 2013*); 2) with 16 longer sequences from a subfamily of
359 Archaea, which has also been investigated by (*Qi et al., 2004*).

360 **Larger dataset with 34 mammalian genome sequences**

361 Likewise, according to procedure depicted, we obtain 34 feature vectors extracted from
362 corresponding primary genome sequences via VMPE model. Moreover, these feature vectors
363 have uniform dimension reduced to 34. Then, using these 34 vectors, we can calculate the
364 pairwise distance matrix, through which we can construct the dendrogram for these 34 sequences.
365 Fig. 4 illustrates the dendrogram based on VMPE model, while Fig. 5 shows the Neighboring-
366 Joining (NJ) tree via alignment-based approach using MEGA software. Compare with Fig. 5, Fig.
367 4 illustrates that our proposed model produce a reasonable result similar to that from the
368 traditional approach. Meanwhile, the clustering results of our model are also consistent with
369 those in the representative published works (*Huang et al., 2011, Yu et al., 2010, Yu, 2013*).

370 **Longer dataset with 16 Archaea sequences**

371 Even for a subfamily, where all the sequences are closer to each other, our proposed VMPE
372 model still distinguishes them clear. The concise information the third dataset is listed in Table 3,
373 while Fig. 6 shows the dendrogram.

374

375 **DISCUSSION**

376 In order to infer evolutionary relationship among organism, we can use alignment-free approaches
377 for estimating the evolutionary distances among multiple DNA/protein sequences and
378 constructing phylogenetic trees. Most methods are based on *k*-mer patterns or transformation

379 from original sequences into numerical vectors. In general, such approaches are less accurate
380 than those traditional phylogeny ones that are based on multiple sequence alignments (MSA). In
381 k -mer based transformation, similarity analysis among multiple genome sequences may easily
382 cause dimension curse, because it is a typical high dimension and small sample problem.

383 In this paper, our goal is to introduce a new methodology for the comparative genomics
384 research community. As an alignment-free approach, our proposed VMPE model does not
385 require any Human intervention and any evolution assumption. It is mathematically well-
386 founded according to ICA-based (independent component analysis) distance preserving
387 transformation which was strictly proved in theory. The latent application of our model is in the
388 field of feature extraction with dimension reduced greatly. It works well when we use hierarchy
389 ICA-based extractor to capture essential information to compare genome sequence within lower
390 dimension space. Motivated by this approach, we proposed to use projected coordinates vector
391 instead of the traditionally used k -mer based vectors directly to calculate the genetic distances
392 among multiple sequences and to reconstruct phylogenetic trees.

393 Although a possible criticism for our approach is that it depends on the existing ICA technique,
394 it is worth our while to stress that the dependence upon ICA-based approach is just to obtain the
395 projected coordinates from multiple genome sequences, so that we can efficiently characterize
396 the sequences within lower dimensional space. In fact, through projection, the proposed
397 approach has a property of distances preserving, which achieves the reasonable results for the
398 comparison of multiple genome sequences. Fig. 2 and Fig. 3 demonstrate that the dendrogram
399 from our approach accords with the evolution tree obtained by traditional alignment-based. So do
400 the results of Fig. 4 and Fig. 5.

401 However, the proposed model depends on an assumption that both local similarity and global
402 similarity should be considered simultaneously for each genome sequence. Our preliminary
403 experimental results have demonstrated the validity of the assumption. Moreover, our results
404 have demonstrated that if the mixed vectors are projected in their independent-components-based
405 coordinate system to extract the corresponding feature vectors within dimension-reduced space,
406 the obtained lower-dimensional vectors are of great value to be applied into the fields of
407 clustering and classification.

408

409 CONCLUSIONS

410 As an alignment free approach, our proposed VMPE model does not require any Human
411 intervention and any evolution assumption. It is mathematically well-founded according to ICA-
412 based distance preserving transformation which was strictly proved in theory. The latent
413 application of our model is in the field of feature extraction with dimension reduced greatly. It
414 works well when we use hierarchy ICA-based extractor to capture essential information to
415 compare genome sequence within lower dimension space. Although a possible criticism for our
416 approach is that it depends on alignment-based results, what is worth while stressing is that the
417 dependence upon alignment is just to obtain the optimized segmentation scheme. In fact, through

418 optimization, the optimal segment number is about 3, which achieves the best performance on
419 the comparison of genome.

420 However, the proposed model depends on an assumption that both local similarity and global
421 similarity should be considered simultaneously for each genome sequence. Our preliminary
422 experimental results have demonstrated the validity of the assumption. Moreover, our results
423 have demonstrated that if the mixed vectors are projected in their independent-components-based
424 coordinate system to extract the corresponding feature vectors within dimension-reduced space,
425 the obtained lower-dimensional vectors are of great value to be applied into the fields of
426 clustering and classification. In the future, we are planning to design a new output-controllable
427 model to further improve the performance on similarity analysis or to extend its applications into
428 other fields.

429

430 ACKNOWLEDGEMENTS

431

432 REFERENCES

- 433 **Blaisdell BE. 1986.** A measure of the similarity of sets of sequences not requiring sequence alignment.
434 *Proceedings of the National Academy of Sciences* **83(14)**:5155-5159 DOI 10.1073/pnas.83.14.5155.
- 435 **Cao Y, Janke A, Waddell PJ, Westerman M, Takenaka O, Murata S, Okada N, Pääbo S, Hasegawa M.**
436 **1998.** Conflict among individual mitochondrial proteins in resolving the phylogeny of eutherian orders.
437 *Journal of Molecular Evolution* **47(3)**:307-322 DOI 10.1007/PL00006389.
- 438 **Chan RH, Chan TH, Yeung HM, Wang RW. 2012.** Composition vector method based on maximum entropy
439 principle for sequence comparison. *IEEE/ACM Trans Comput Biol Bioinform* **9(1)**:79-87 DOI
440 10.1109/TCBB.2011.45.
- 441 **Comon P. 1994.** Independent component analysis, a new concept? *Signal Processing* **36(3)**:287-314 DOI
442 10.1016/0165-1684(94)90029-9.
- 443 **Cristea PD. 2002.** Conversion of nucleotide sequences into genomic signals. *Journal of Cellular and*
444 *Molecular Medicine* **6(2)**:25 DOI 10.1111/j.1582-4934.2002.tb00196.x.
- 445 **Dai Q, Liu X, Yuhua Yao, Zhao F. 2011.** Numerical characteristics of word frequencies and their application
446 to dissimilarity measure for sequence comparison. *Journal of Theoretical Biology* **276(1)**:174-180 DOI
447 10.1016/j.jtbi.2011.02.005.
- 448 **Deng M, Yu C, Liang Q, He RL, Yau SS-T. 2011.** A novel method of characterizing genetic sequences:
449 genome space with biological distance and applications. *PLoS ONE* **6(3)**:e17293 DOI
450 10.1371/journal.pone.0017293.g001.
- 451 **Dong R, Zhu Z, Yin C, He RL, Yau SS. 2018.** A new method to cluster genomes based on cumulative
452 Fourier power spectrum. *Gene* **673**:239-250 DOI 10.1016/j.gene.2018.06.042.
- 453 **Fernandes F, Freitas AT, Almeida JS, Vinga S. 2009.** Entropic profiler - detection of conservation in
454 genomes using information theory. *BMC Res Notes* **2**:72 DOI 10.1186/1756-0500-2-72.
- 455 **Gao L, Qi J. 2007.** Whole genome molecular phylogeny of large dsDNA viruses using composition vector
456 method. *BMC Evol Biol* **7**:41 DOI 10.1186/1471-2148-7-41.
- 457 **Huang DS, Zheng CH. 2006.** Independent component analysis-based penalized discriminant method for
458 tumor classification using gene expression data. *Bioinformatics* **22(15)**:1855-1862 DOI
459 10.1093/bioinformatics/btl190.
- 460 **Huang G, Zhou H, Li Y, Xu L. 2011.** Alignment-free comparison of genome sequences by a new numerical
461 characterization. *Journal of Theoretical Biology* **281(1)**:107-112 DOI 10.1016/j.jtbi.2011.04.003.

- 462 **Huang Y, Wang T. 2011.** Phylogenetic analysis of DNA sequences with a novel characteristic vector. *Journal*
463 *of Mathematical Chemistry* **49(8)**:1479-1492 DOI 10.1007/s10910-011-9811-x.
- 464 **Huang Y, Yang L, Wang T. 2011.** Phylogenetic analysis of DNA sequences based on the generalized pseudo-
465 amino acid composition. *Journal of Theoretical Biology* **269(1)**:217-223 DOI 10.1016/j.jtbi.2010.10.027.
- 466 **Jun S-R, Sims GE, Wu GA, Kim S-H. 2010.** Whole-proteome phylogeny of prokaryotes by feature
467 frequency profiles: an alignment-free method with optimal feature resolution. *Proceedings of the National*
468 *Academy of Sciences* **107(1)**:6 DOI 10.1073/pnas.0913033107.
- 469 **Li M, Badger JH, Chen X, Kwong S, Kearney P, Zhang H. 2001.** An information-based sequence distance
470 and its application to whole mitochondrial genome phylogeny. *Bioinformatics* **17(2)**:149-154 DOI
471 10.1093/bioinformatics/17.2.149.
- 472 **Luczak BB, James BT, Girgis HZ. 2017.** A survey and evaluations of histogram-based statistics in
473 alignment-free sequence comparison. *Brief Bioinform* (1):16 DOI 10.1093/bib/bbx161.
- 474 **Mendizabal-Ruiz G, Roman-Godinez I, Torres-Ramos S, Salido-Ruiz RA, Velez-Perez H, Morales JA.**
475 **2018.** Genomic signal processing for DNA sequence clustering. *PeerJ* **6(3)**:e4264 DOI 10.7717/peerj.4264.
- 476 **Nguyen KD. 2012.** On the edge of web-based multiple sequence alignment services. *Tsing Hua Science and*
477 *Technology* **17(6)**:629-637 DOI 10.1109/TST.2012.6374364.
- 478 **Qi J, Wang B, Hao B-I. 2004.** Whole proteome prokaryote phylogeny without sequence alignment a: k-string
479 composition approach. *Journal of Molecular Evolution* **58**:1-11 DOI 10.1007/s00239-003-2493-7.
- 480 **Sims GE, Jun S-R, Wu GA, Kim S-H. 2009.** Alignment-free genome comparison with feature frequency
481 profiles (FFP) and optimal resolutions. *Proceedings of the National Academy of Sciences* **106(8)**: 2677-
482 2682 DOI 10.1073/pnas.0813249106.
- 483 **Sims GE, Jun SR, Wu GA, Kim SH. 2009.** Whole-genome phylogeny of mammals: Evolutionary
484 information in genic and nongenic regions. *Proceedings of the National Academy of Sciences*
485 **106(40)**:17077-17082 DOI 10.1073/pnas.0909377106.
- 486 **Sims GE, Kim S-H. 2011.** Whole-genome phylogeny of Escherichia coli/Shigella group by feature frequency
487 profiles (FFPs). *Proceedings of the National Academy of Sciences* **108 (20)**:8329-8334 DOI 10.1073/pnas.
- 488 **Stuart GW, Berry MW. 2004.** An SVD-based comparison of nine whole eukaryotic genomes supports a
489 coelomate rather than ecdysozoan lineage. *BMC Bioinformatics* **5**:1-13 DOI 10.1186/1471-2105-5-204.
- 490 **Stuart GW, Moffett K, Leader JJ. 2002.** A comprehensive vertebrate phylogeny using vector representations
491 of protein sequences from whole genomes. *Molecular Biology Evolution* **19(4)**:554-562 DOI .
- 492 **Vinga S. 2014.** Information theory applications for biological sequence analysis. *Brief Bioinform* **15(3)**:376-
493 389 DOI 10.1093/bib/bbt068.
- 494 **Vinga S, Almeida J. 2003.** Alignment-free sequence comparison--a review. *Bioinformatics* **19(4)**:513-523
495 DOI 10.1093/bioinformatics/btg005.
- 496 **Wu T-J, Huang Y-H, Li L-A. 2005.** Optimal word sizes for dissimilarity measures and estimation of the
497 degree of dissimilarity between DNA sequences. *Bioinformatics* **21(22)**:4125-4132 DOI
498 10.1093/bioinformatics/bti658.
- 499 **Yang L, Zhang X, Zhu H. 2012.** Alignment free comparison: similarity distribution between the DNA
500 primary sequences based on the shortest absent word. *Journal of Theoretical Biology* **295**:125-131 DOI
501 10.1016/j.jtbi.2011.11.021.
- 502 **Yeredor A. 2002.** Non-orthogonal joint diagonalization in the least-squares sense with application in blind
503 source separation. *IEEE Transactions on Signal Processing* **50(7)**:1545-1553 DOI
504 10.1109/TSP.2002.1011195.
- 505 **Yu C, Liang Q, Yin C, He RL, Yau SST. 2010.** A novel construction of genome space with biological
506 geometry. *DNA Research* **17(3)**:155-168 DOI 10.1093/dnares/dsq008.

- 507 **Yu H-J. 2013.** Segmented K-mer and its application on similarity analysis of mitochondrial genome sequences.
508 *Gene* **518(2)**:419-424 DOI 10.1016/j.gene.2012.12.079.
- 509 **Yu H-J, Huang D-S. 2012.** Novel graphical representation of genome sequence and its applications in
510 similarity analysis. *Physica A: Statistical Mechanics and its Applications* **391(23)**:6128-6136 DOI
511 10.1016/j.physa.2012.07.020.
- 512
- 513

- 514 **Fig. 1.** Scheme for the monolayer ICA upon k -mers of corresponding multiple sequences.
- 515 **Fig. 2.** The dendrogram based on segmented 7-mer using 20 mitochondrial genome sequences of eutherian
516 species.
- 517 **Fig. 3.** The Neighboring-Joining (NJ) tree of the eutherians constructed from the mitochondrial genome
518 sequences of the 20 species via traditional alignment-based approach using MEGA software.
- 519 **Fig. 4.** The dendrogram based on VMPE model using 34 mitochondrial genome sequences from mammalian
520 species.
- 521 **Fig. 5.** The Neighboring-Joining (NJ) tree of the mammals constructed from the mitochondrial genome
522 sequences of the 34 species via traditional alignment-based approach using MEGA software.
- 523 **Fig. 6.** The dendrogram based on VMPE model using 16 Archaea sequences.
- 524
- 525 **Table 1.** Summary information for the 20 eutherians species.
- 526 **Table 2.** Comparison of performance values with other two published works for the 20 eutherians species via
527 the pairwise distances among Human and the rest 19 ones.
- 528 **Table 3.** 16 longer sequences of Archaea with names, abbreviations, and NCBI accession numbers.

Table 1 (on next page)

Summary information for the 20 eutherians species

1 **Table 1.** Summary information for the 20 eutherians species

2

Accession no.	Species	Length
V00662	Human	16,569
D38116	Pigmy chimpanzee	16,563
D38113	Common chimpanzee	16,554
D38114	Gorilla	16,364
D38115	Bornean orangutan	16,389
X99256	Gibbon	16,472
Y18001	Baboon	16,521
X79547	Horse	16,660
Y07726	White rhinoceros	16,832
X63726	Harbor seal	16,826
X72004	Gray seal	16,797
U20753	Cat	17,009
X61145	Fine Whale	16,398
X72204	Blue Whale	16,402
V00654	Cow	16,338
X14848	Norway rat	16,300
V00711	Mouse	16,295
Z29573	Opossum	17,084
Y10524	Wallaroo	16,896
X83427	Platypus	17,019

3

Table 2 (on next page)

Comparison of performance values with other two published works for the 20 eutherians species via the pairwise distances among Human and the rest 19 ones

1 **Table 2.** Comparison of performance values with other two published works for the 20 eutherians species via
 2 the pairwise distances among Human and the rest 19 ones

3

Human Vs. rest species	Kolmogorov complexity (<i>Li, Badger, Chen, Kwong, Kearney & Zhang, 2001</i>)	OPT-SVD (<i>Yu & Huang, 2012</i>)	VMPE ($\times 10^{-2}$)
P.Chim.	0.654234	0.059285712	0.414459
C.Chim.	0.657387	0.041232607	0.413945
Gorilla	0.732325	0.041462402	0.442228
Orang.	0.847139	0.068655059	0.486724
Gibbon	0.880203	0.065596946	0.496526
Baboon	0.841775	0.070727095	0.533934
Horse	0.971558	0.144379093	0.555873
W.Rhin.	0.973694	0.168378304	0.557738
H.Seal	0.974737	0.211514019	0.578971
G.Seal	0.97576	0.207692966	0.579792
Cat	0.977328	0.276304513	0.577581
F.Whale	0.980493	0.209182609	0.562304
B.Whale	0.976034	0.197653729	0.557521
Cow	0.97362	0.283610943	0.565623
Rat	0.981715	0.276197913	0.569249
Mouse	0.9804	0.372617581	0.585258
Oposs.	0.986243	0.54735028	0.626082
Walla.	0.985926	0.258992019	0.569803
Platypus	0.988041	0.445501566	0.612848

4

5

6

7 **Li M, Badger JH, Chen X, Kwong S, Kearney P, Zhang H. 2001.** An information-based sequence
 8 distance and its application to whole mitochondrial genome phylogeny, *Bioinformatics*, **17(2)**:149-
 9 154 DOI 10.1093/bioinformatics/17.2.149.

10 **Yu H-J, Huang D-S. 2012.** Novel graphical representation of genome sequence and its applications in
 11 similarity analysis, *Physica A: Statistical Mechanics and its Applications*, **391(23)**:6128-6136 DOI
 12 10.1016/j.physa.2012.07.020.

13

Table 3 (on next page)

16 longer sequences of Archaea with names, abbreviations, and NCBI accession numbers

1 **Table 3.** 16 longer sequences of Archaea with names, abbreviations, and NCBI accession numbers

2

Species/strain	Abbrev.	Accession No.	Length (nt)
Pyrobaculum aerophilum	01Pyrae	NC_003364	2222430
Aeropyrum pernix K1	02Aerpe	NC_000854	1669696
Sulfolobus solfataricus	03Sulso	NC_002754	2992245
Sulfolobus tokodaii	04Sulto	NC_003106	2694756
Methanobacterium thermoautotrophicus	05Metth	NC_000916	1751377
Methanococcus jannaschii	06Metja	NC_000909	1664970
Methanosarcina acetivorans strain C2A	07Metac	NC_003552	5751492
Methanosarcina mazei Goel	08Metma	NC_003901	4096345
Halobacterium sp. NRC-1	09Halsp	NC_002607	2014239
Thermoplasma acidophilum	10Theac	NC_002578	1564906
Thermoplasma volcanium	11Thevo	NC_002689	1584804
Pyrococcus abyssi	12Pyrab	NC_000868	1765118
Pyrococcus furiosus	13Pyrfu	NC_003413	1908256
Pyrococcus horikoshii	14Pyrho	NC_000961	1738505
Archaeoglobus fulgidus	15Arcfu	NC_000917	2178400
Methanopyrus kandleri AV19	16Metka	NC_003551	1694969

3

Figure 1

Scheme for the monolayer ICA upon k-mers of corresponding multiple sequences

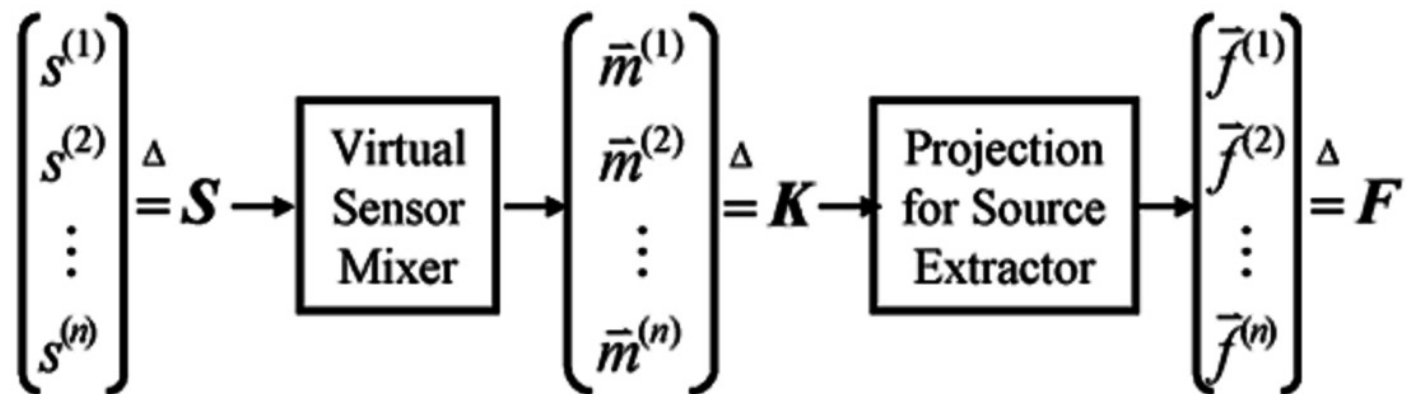


Figure 2

The dendrogram based on segmented 7-mer using 20 mitochondrial genome sequences of eutherian 366 species

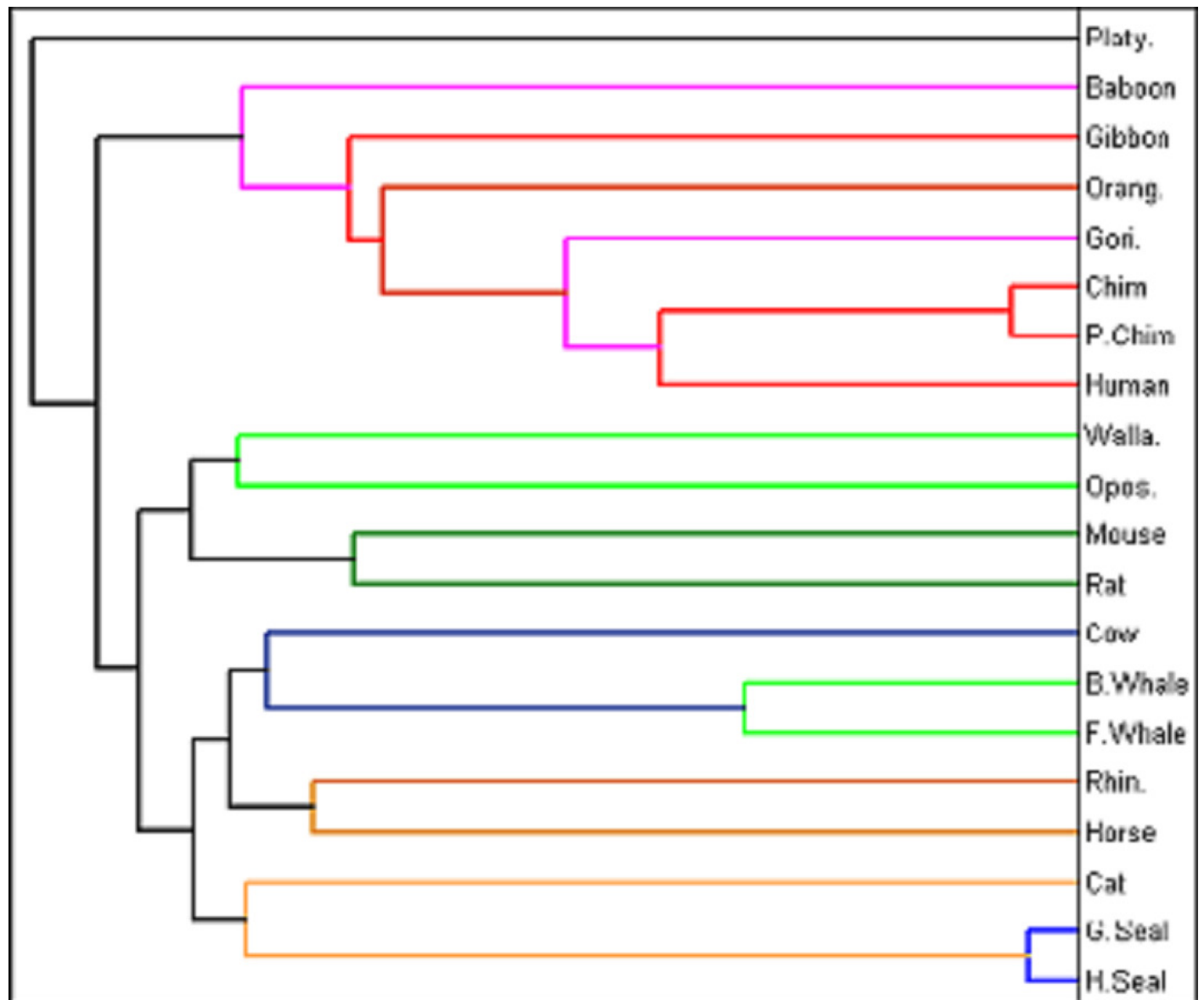


Figure 3

The Neighboring-Joining (NJ) tree of the eutherians constructed from the mitochondrial genome 370 sequences of the 20 species via traditional alignment-based approach using MEGA software

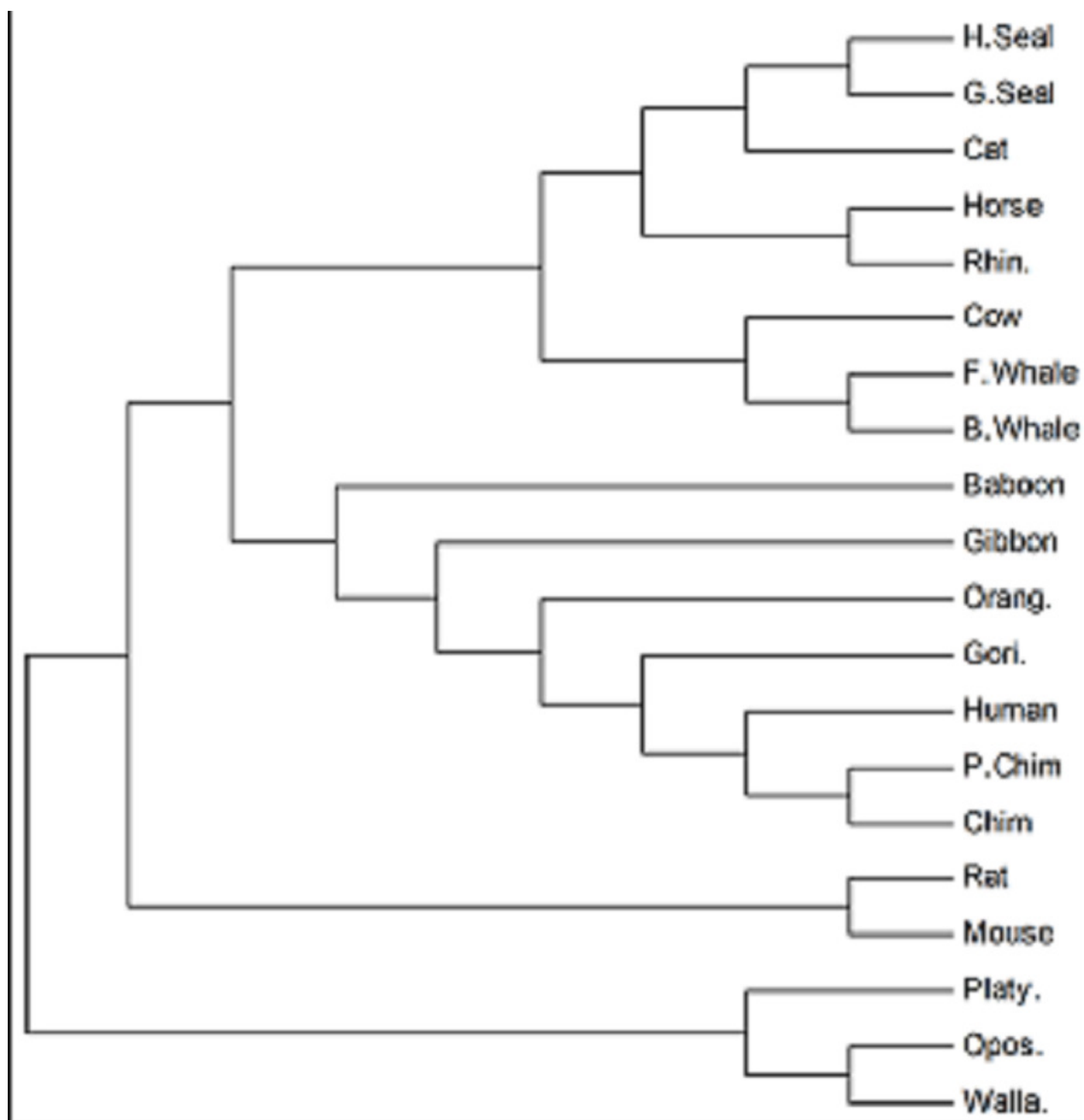


Figure 4

The dendrogram based on VMPE model using 34 mitochondrial genome sequences from mammalian 400 species

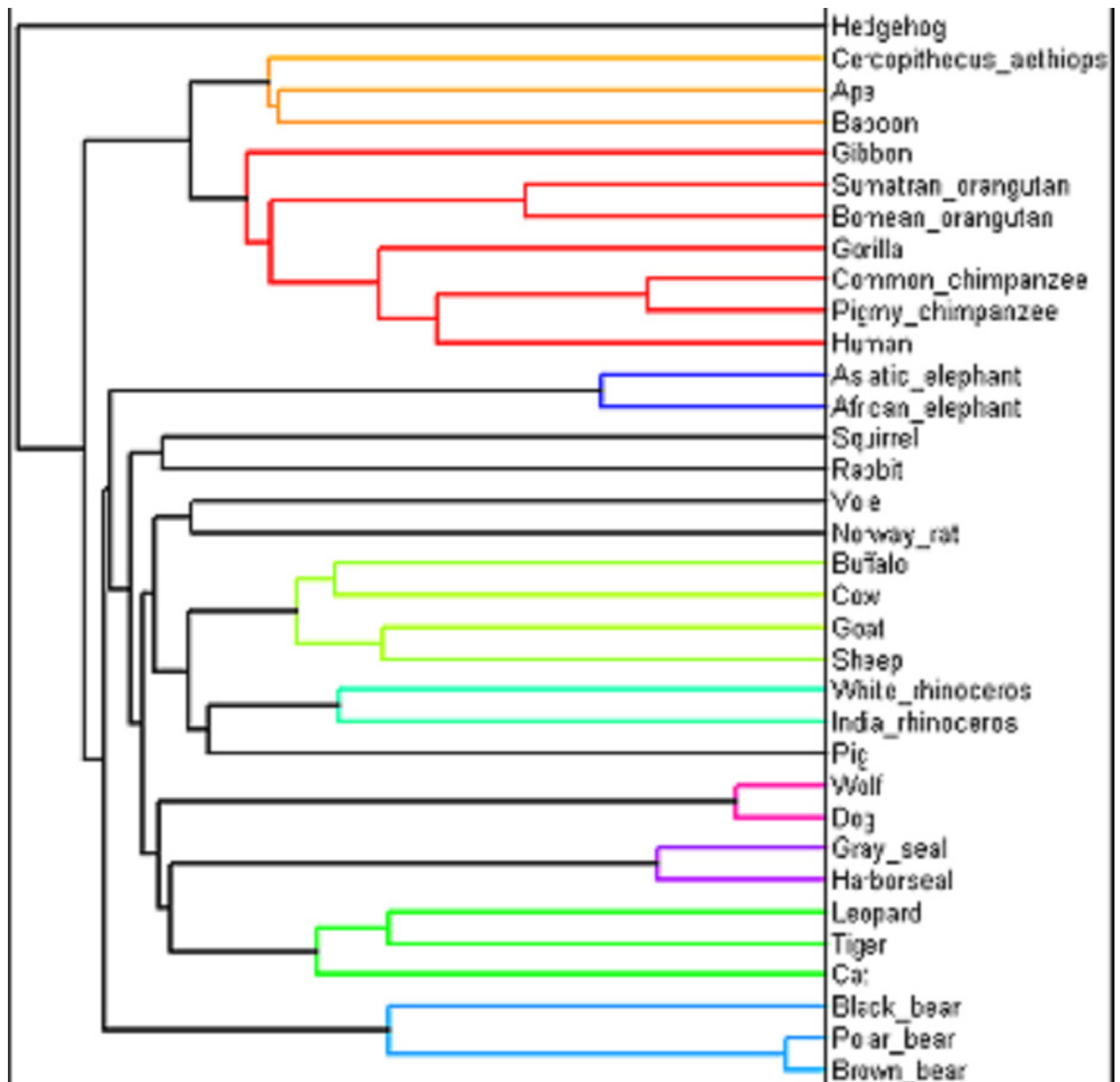


Figure 5

The Neighboring-Joining (NJ) tree of the mammals constructed from the mitochondrial genome 405 sequences of the 34 species via traditional alignment-based approach using MEGA software

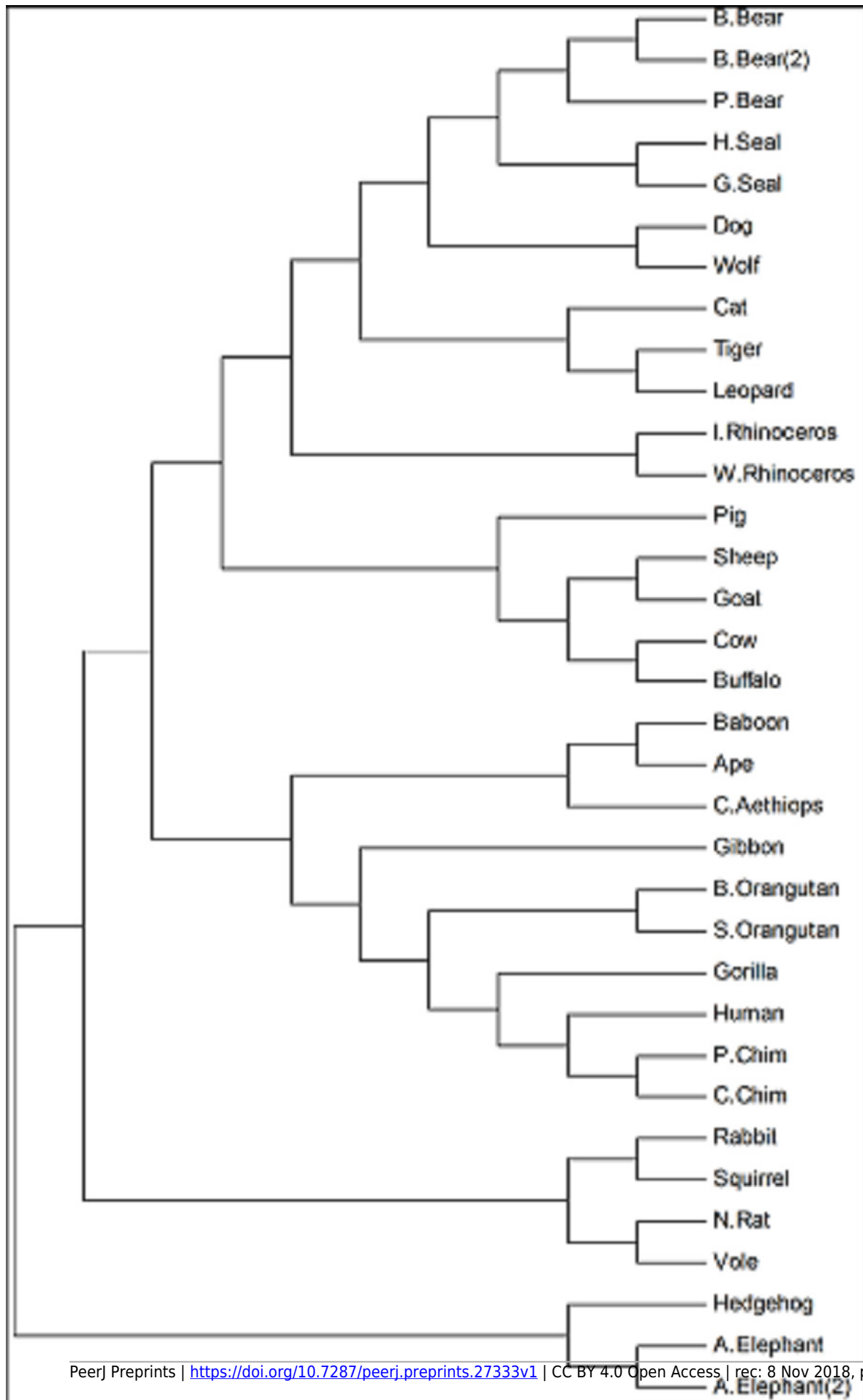


Figure 6

The dendrogram based on VMPE model using 16 Archaea sequences

