

REDUCED AMINO ACID ALPHABET-BASED ENCODING AND ITS IMPACT ON MODELING INFLUENZA ANTIGENIC EVOLUTIONForghani M.^a,Firstkov A. L.^a,AlyanNezhadi M. M.^b,Danilenko D. M.^c,Komissarov A. B.^c

^aN.N. Krasovskii Institute of Mathematics and Mechanics of the Ural Branch of the Russian Academy of Sciences (IMM UB RAS), Yekaterinburg, Russia

^bMazandaran University of Science and Technology, Iran.

^c Smorodintsev Research Institute of Influenza, Saint Petersburg, Russia

КОДИРОВАНИЕ С ПОМОЩЬЮ СОКРАЩЁННОГО АМИНОКИСЛОТНОГО АЛФАВИТА И ЕГО ВЛИЯНИЕ НА МОДЕЛИРОВАНИЕ АНТИГЕННОЙ ЭВОЛЮЦИИ ГРИППАФоргани М.¹,Фирстков А. Л.¹,Аляннеджади М. М.²,Даниленко Д. М.³,Комиссаров А. Б.³

¹ Федеральное государственное бюджетное учреждение науки Институт математики и механики им. Н. Н. Красовского Уральского отделения Российской академии наук (ИММ УрО РАН), Екатеринбург, Россия

² Университет науки и технологии Мазандарана, Иран.

³ ФГБУ Научно-исследовательский институт гриппа имени А.А. Смородинцева, Минздрава России.

ABSTRACT. Currently, vaccination is one of the most efficient ways to control and prevent influenza infection. Vaccine production largely relies on the results of laboratory assays, including hemagglutination inhibition and microneutralization assays, which are time-consuming and laborious. Viruses can escape from the immune response that results in the need to revise and update vaccines biannually. The hemagglutination inhibition assay can measure how effectively antibodies against a reference strain bind and block an antigen of the test strain. Various computer-aided models have been developed to optimize candidate vaccine strain selection. A general problem in modeling of antigenic evolution is the representation of genetic sequences for input into the research model. Our motivation stems from the well-known problem of encoding genetic information for modeling antigenic evolution. This paper introduces a two-fold encoding approach based on reduced amino acid alphabet and amino acid index databases called AAindex. We propose to apply a simplified amino acid alphabet in modeling of antigenic evolution. A simplified alphabet, also called a sub-alphabet or reduced amino acid alphabet, implies to use the 20 amino acids being clustered and divided into amino acid groups. The proposed encoding allows to redefine mutations termed for amino acid groups located in reduced alphabets. We investigated 40 reduced amino acid sets and their performance in modeling antigenic evolution. The experimental results indicate that the proposed reduced amino acid alphabets can achieve the performance of the standard alphabet in its accuracy. Moreover, these alphabets provide deeper insight into various aspects of the relationship between mutation and antigenic variation. By checking identified high-impact sites in the Influenza Research Database, we found that not only antigenic sites have a significant influence on antigenicity, but also other amino acids located in close proximity. The results indicate that all selected non-antigenic sites are related to immune responses. According to the Influenza Research Database, these have been experimentally determined to be T-cell epitopes, B-cell epitopes, and MHC-binding epitopes of different classes. This highlighted a caveat: while simulating antigenic

evolution, the model should consider not only the genetic information on antigenic sites, but also that of neighboring positions, as they may indirectly impact antigenicity. Additionally, our findings indicate that structural and charge characteristics are the most beneficial in modeling antigenic evolution, which is in agreement with previous studies.

Keywords: AAindex, antigenic evolution, hemagglutinin, influenza, modeling, reduced amino acid alphabet

РЕЗЮМЕ. В настоящее время, вакцинация является одним из наиболее эффективных способов контроля и профилактики гриппозной инфекции. Производство вакцин в основном зависит от результатов лабораторных анализов, включая анализ реакции торможения гемагглютинации и микронеутрализации, которые требуют много времени и труда. Вирусы могут избегать иммунного ответа, что приводит к необходимости пересмотра и обновления вакцин два раза в год. Анализ реакции торможения гемагглютинации позволяет измерить, насколько эффективно антитела против эталонного штамма связывают и блокируют антиген испытуемого штамма. Для оптимизации выбора вакцинного штамма-кандидата были разработаны различные компьютерные модели. Одна из общих проблем в моделировании антигенной эволюции является представление генетических последовательностей для ввода в исследовательскую модель. Наша мотивация связана с хорошо известной проблемой кодирования генетической информации для моделирования антигенной эволюции. В данной работе представлен двухэтапный подход к кодированию, основанный на сокращенных аминокислотных алфавитах и базах данных аминокислотных индексов под названием AAindex. Мы предлагаем использовать упрощенные аминокислотные алфавиты для моделирования антигенной эволюции.

Упрощённый алфавит, также называемый субалфавитом или сокращённым аминокислотным алфавитом, это алфавит, в котором 20 аминокислот разделены на группы. Предложенное кодирование позволяет переопределить мутации в терминах групп аминокислот, расположенных в сокращённом алфавите. Мы исследовали 40 сокращённых алфавитов и их эффективность при моделировании антигенной эволюции. Результаты экспериментов показывают, что предложенные сокращённые аминокислотные алфавиты могут достичь показателей стандартного алфавита по точности. Более того, эти алфавиты позволяют лучше понять взаимосвязь между мутациями и антигенными изменениями с различных точек зрения. Проверив полученные высокоэффективные сайты в исследовательской базе данных гриппа (Influenza Research Database), мы обнаружили, что не только антигенные сайты оказывают значительное влияние на антигенность, но и их соседние аминокислоты. Результаты показывают, что все выбранные неантигенные участки связаны с иммунным ответом. Согласно исследовательской базе данных гриппа, экспериментально установлено, что это эпитопы Т-клеток, эпитопы В-клеток и МНС-связывающие эпитопы различных классов. Это подчёркивает значимость того, что: при моделировании антигенной эволюции модель должна учитывать не только генетическую информацию антигенных участков, но и генетическую информацию соседних позиций, поскольку они могут косвенно влиять на антигенность. Кроме того, наши результаты показывают, что, в соответствии с предыдущими исследованиями, структурные и зарядовые характеристики аминокислот являются наиболее значимыми при моделировании антигенной эволюции.

Ключевые слова: AAindex, антигенная эволюция, гемагглютинин, грипп, моделирование, сокращённый аминокислотный алфавит

1 **INTRODUCTION.** Influenza is a contagious respiratory infection that affects 5%-
2 15% of the population worldwide annually, resulting in 3-5 million cases of severe
3 illness and 250,000 to 500,000 deaths [36]. Influenza epidemics influence public
4 health and involve severe economic consequences, making it the subject of various
5 economic studies [4]. The World Health Organization (WHO) continuously
6 monitors viral pathogens, especially those that can become epidemics or pandemics,
7 and decides on strategies to combat them. Given the special status of influenza, the
8 WHO created the Global Influenza Surveillance and Response System, the primary
9 function of which is to monitor the evolution of the influenza virus and to provide
10 recommendations for the annual vaccine's composition for the Northern and
11 Southern Hemispheres.

12 Influenza viruses are part of the *Orthomyxoviridae* family.
13 According to antigenic characteristics of their nuclear proteins, they are grouped into
14 four types: IVA (A); IVB (B); IVC (C); and IVD (D). Among them, types A and B
15 are associated with influenza outbreaks. Type C appears to evolve slowly and leads
16 to less severe and less significant health consequences. Type D is an influenza C-
17 like virus that is observed in non-human hosts, e.g., cattle and swine [30]. Type A is
18 further classified according to the combination of hemagglutinin (HA) and
19 neuraminidase (NA), the two main surface antigens of influenza that play a key role
20 in infectivity and immune responses. HA has 18 variants (H1-H18), while the NA
21 protein can be one of 11 variants (N1-N11). Hence, the virus can theoretically be
22 any of 198 different subtypes; this provides an ability to infect a broad spectrum of
23 various hosts [37]. Despite this diversity, humans are infected with only a limited
24 number of influenza A subtypes (i.e., H1N1, H2N2, H3N2), with H1N1 and H3N2
25 being currently relevant. Thus, we consider them in this paper. Other zoonotic
26 subtypes represent only sporadic infections and are out of the scope of this study.

27 Influenza A viruses are capable of enormous genetic variation, both through
28 continuous, gradual mutation and by reassortment of gene segments between

29 viruses, resulting in emerging novel antigenic variants. Epidemics are the result of
30 gradual evolutionary changes called antigenic drift, which leads to the generation of
31 new strains from existing ones through mutation. In addition to antigenic drift, the
32 influenza virus can be altered by antigenic shift. It is an abrupt significant change in
33 influenza viruses resulting in the emergence of new HA and/or NA. It is the process
34 by which at least two subtypes combine into a new subtype that has a mixture of
35 surface antigens of two or more strains [35].

36 The only effective method to control influenza is vaccination, eliciting
37 protective neutralizing antibodies and memory T-cell responses. Since HA antigen
38 abundance on the viral surface is approximately four-fold greater than NA [31], it is
39 the primary component in vaccine compositions. This is the reason why we consider
40 only HA protein sequences in this paper.

41 The influenza vaccine requires an update if the vaccine composition strains
42 are antigenically distinct from currently circulating viruses. A gold-standard and
43 widespread laboratory procedure called hemagglutination inhibition (HI) assay is
44 used to assess the measure of antigenic similarity between strains. The HI assay can
45 measure how effectively antibodies against a reference strain bind and block an
46 antigen of the test strain. High HI titers indicate a high degree of antigenic similarity
47 between strains [16]. The main conclusion of HI assay analysis is determining
48 antigenic distance (i.e., similarity between reference and test antigens), which
49 further can be presented in terms of a binary variable called antigenic variant.
50 Currently, there are two widely used definitions of antigenic distance [18, 27]:

$$d_1(i, j) = \log_2\left(\frac{M}{H_{i,j}}\right) \quad (1)$$

$$d_2(i, j) = \sqrt{\frac{H_{i,i} \times H_{j,j}}{H_{i,j} \times H_{j,i}}} \quad (2)$$

51 where $H_{i,j}$ is the obtained HI titer for antiserum of (reference) strain j against the
52 antigen of (test) strain i , and M is the maximum titer observed for antiserum j against
53 any antigen in the HI table. The antigenic variant is determined by applying the
54 threshold to the obtained antigenic distance. The pair of test and reference viruses
55 whose antigenic distance meets the threshold are designated as antigenic variants;
56 otherwise, they are only antigenically similar.

57 The HI assay is a labor-intensive and time-consuming procedure, while
58 vaccine development is under time pressure. Over the past decade, various
59 computer-aided approaches have been developed to speed up the process of strain
60 selection and to increase the quality of vaccine production. Klingten *et al.* [16] has
61 provided a comprehensive review of antigenic evolution prediction associated with
62 vaccine production. They classified the approaches into phylogenetic and population
63 genetics-based methods, statistical methods, epidemiological models, and other
64 methods based on information and graph theories. The approaches employ different
65 data types serving as model inputs, e.g., viral sequence, HI assay data, protein
66 structure, physicochemical properties, etc. A critical step in antigenic variant
67 modeling is describing the biological significance of a mutation between test and
68 reference viruses and linking it to antigenicity.

69 Unfortunately, the exact roles and how they affect biological properties within
70 evolution are not yet fully understood for many such changes. Generally, it is known
71 that evolution is influenced by several biological properties, especially the volume
72 and hydrophobicity of amino acids [32]. Studies on amino acid property changes
73 provide fundamental information about the evolution of specific proteins. Earlier
74 studies indicated that HA antigen is positively charged, while on the contrary, the
75 glycan receptors of the host cell are negatively charged. Thus, changes in
76 electrostatic charge due to mutation can play a significant role in receptor specificity,
77 enhancing or diminishing the receptor binding affinity and avidity [2, 17]. Moreover,
78 Huang *et al.* [14] recently showed that charged amino acid mutations impact

79 influenza virus evolution and are beneficial in vaccine research. Accordingly,
80 mutation can be considered a multidimensional event, wherein each dimension
81 represents an amino acid attribute.

82 Several techniques reflect the biological characteristics of mutation in
83 numerical domains, among which application of the AAindex database [15] is the
84 most popular. The AAindex database is a comprehensive collection of biological,
85 physical, and chemical amino acid properties collected from scientific papers and
86 accessed through *www.genome.jp*. The database mainly consists of three sections:
87 AAindex1; AAindex2; and AAindex3. AAindex1 includes various amino acid
88 indices, each of which can be represented as a numerical vector of 20 numbers
89 representing 20 standard amino acids. AAindex2 contains different amino acid
90 mutation matrices, while AAindex3 consists of statistical protein contact potentials.
91 The AAindex database (ver.9.2) currently covers 566, 94, and 47 records for
92 AAindex1, AAindex2, and AAindex3, respectively.

93 As mentioned, the AAindex database has been employed for encoding protein
94 sequence in various studies. Here, we mention some of the more relevant studies in
95 which the AAindex database was used for exploring genetic and antigenic evolution.
96 Yao *et al.* [39] proposed an algorithm called joint random forest regression to predict
97 antigenic variants. They compared 95 amino acid matrices, including AAindex2, to
98 assess the relationship between genetic and antigenic evolution by amino acid
99 attributes at different protein sites. Their results indicated that structural features are
100 more significant to the antigenicity of the influenza virus. Wang *et al* [34] suggested
101 an approach based on matrix completion for predicting antigenic evolution. They
102 studied the impact of 65 amino acid substitution matrices taken from the AAindex
103 database to predict antigenic evolution. Their results suggested that the “homologous
104 structure derived matrix (called HSDM) for alignment of distantly related
105 sequences” outperformed others in terms of RMSE.

106 Moreover, Qui *et al.* [24] developed a structure-based antigenicity scoring
107 model. Their model engages antigenically dominant positions according to structural
108 context, including correlation with local amino acid attribute changes, to analyze
109 antigenicity. They demonstrated that incorporating the structural context of protein
110 can enhance antigenic evolution prediction. Additionally, Forghani and
111 Khachay [10] carried out a principal component analysis on AAindex1 and
112 introduced 11 indices that explained 91% of the total variation in the database. The
113 new indices are further used to encode HA protein sequence and create an input
114 tensor fed into a convolutional neural network. Their model achieves a mean
115 absolute error of 0.935 antigenic units for yearly, non-anticipating prediction of
116 antigenic distance for subtype H1N1 (2001-2009). Cui *et al.* [6] suggested modeling
117 influenza virus antigenicity by selecting the most significant sites, clustering the
118 AAindex1 based on mutual information, and encoding the sites by the representative
119 from clusters to form the feature vector. The feature vector is further given to a
120 classifier to discretize antigenic variant classes. Recently, we performed a
121 preliminary analysis to study the impact of amino acid encoding on modeling the
122 antigenic evolution of the influenza virus [11]. Apart from Cui *et al.*'s work, our
123 work introduces an early-stage mutation encoding by applying reduced amino acid
124 alphabets.

125 The current paper addresses one of the fundamental challenges in
126 bioinformatics: deciding how to represent input genetic information for modeling
127 more efficiently and meaningfully. In response to this problem, we employed
128 simplified amino acid alphabets. A simplified alphabet, also called a sub-alphabet or
129 reduced amino acid alphabet (RAAA), is an alphabet in which the 20 amino acids
130 are clustered and divided into amino acids groups. RAAA construction is a problem
131 that belongs to the set partitioning problem, which is out of this paper's scope.
132 Previous studies have shown that RAAAs have been successfully applied in various
133 domains, including: protein annotation and description; protein functionality

134 prediction [21, 41]; protein folding assessment; sequence classification [19];
135 consensus sequence search; and genetic pattern identification [5].

136 A reduced amino acid set simplifies protein system complexity, providing a
137 better insight into structural similarities across protein sequences [42]. We
138 considered different definitions of similarity via RAAAs to reconstruct the
139 relationship between genotype and phenotype. A RAAA represents genetic
140 information on a coarse scale, which may highlight attributes that drive antigenic
141 evolution of the influenza virus.

142 In our approach, encoding is conducted in two steps. In addition to the
143 standard amino acid alphabet, the first step employs a RAAA to represent the
144 mutation in different structural, biological, and physicochemical contexts. Further,
145 the second step encodes the alphabetical information of the encoded genetic
146 sequence into a numerical one, which enables its use in various types of
147 mathematical modeling. Preliminary results indicate that some RAAA-based models
148 outperform models based on the standard amino acid alphabet in terms of accuracy.

149 In this paper, we take a step forward and perform a comprehensive analysis to
150 further refine result accuracy. The contributions of this paper are three-fold:

- 151 1. We propose a novel encoding method using reduced amino acid alphabets,
152 which helps to clarify the genetic/antigenic relationship.
- 153 2. Relative to similar previous studies [6, 11], we improve the approach at several
154 levels:
 - 155 2.1. Increasing the resolution of thresholds.
 - 156 2.2. Clustering by several methods and comparing their results to find the
157 optimal number of clusters.
 - 158 2.3. Selecting the closest index to the center of a cluster as its representative.
 - 159 2.4. Applying five well-known classification algorithms.

160 2.5. Optimizing of classifier hyperparameters through a comprehensive grid
161 search.

162 3. Relying on experimental results, we found that incorporating structural and
163 charge properties can enhance modeling quality, which is in agreement with
164 previous studies.

165 The rest of the paper is organized as follows. Section II describes the general
166 computational pipeline, data preparation, and all necessary algorithms for primary
167 and secondary encoding. Experimental setup and its outcomes are presented in
168 Section III. This section also covers interpretation and discussion of the obtained
169 results. Finally, the conclusion is given in Section IV.

170

171 **II. METHODOLOGY**

172 As mentioned earlier, our experimental design was inspired by a published
173 methodology [6]. However, we propose some modifications and enhancements to
174 improve modeling quality. Our approach is mainly divided into five steps: encoding
175 genetic sequences by RAAA; selecting the most relevant sites; clustering the
176 AAindex1 data set based on selected sites; encoding the selected sites by a
177 representative from each cluster; and modeling antigenic variants by a classifier. The
178 general schema of our pipeline is shown in Figure 1.

179 ***2.1. Data Preparation***

180 Our approach relies on three database types, each of which requires specific
181 preparation in order to be used in the computational pipeline.

182 ***2.1.1. Simplified Amino Acid Alphabets***

183 Apart from the standard amino acid alphabet with 20 letters, there are various
184 RAAAs, in which the number of letters is less than 20. Typically, an RAAA is
185 obtained by grouping the 20 amino acids. There are several strategies to perform

186 this, some of which have been described [29]. For example, the set of 20 amino acids
187 can be divided into three groups based on Van Der Waals volume by setting three
188 ranges (0-2.78, 2.95-4.00, 4.43-8.08), resulting in three partitions: GASCTPD;
189 NVEQIL; and MHKFRYW. This permits new interpretation of the mutation from a
190 different point of view, such as change in hydrophobicity. In total, 40 published
191 RAAAs were collected [8, 29, 38] and are presented in Table 1.

192 **2.1.2. HI Assay Database**

193 Typically, an HI assay database record includes three fields of information:
194 test virus identifier; reference antiserum identifier; and HI titer. Sometimes
195 additional metadata, such as experiment date, may be appended. HI assay results can
196 be presented in four forms: raw HI titer; standardized HI titer; antigenic distance;
197 and antigenic variant. At this point, we only used the antigenic variant obtained via
198 the antigenic distance threshold. We employed Eq. (1) with threshold 4 for
199 calculating the antigenic variant. Duplicated entries were averaged in terms of titer.
200 Therefore, each test/reference virus combination is unique within the database. Here,
201 we considered two subtypes in the influenza vaccine, H1N1 and H3N2. The HI assay
202 database was taken from references [13, 34]. The final obtained database had 7,449
203 H3N2 and 3,747 H1N1 entries. There were 506 viruses for the H1N1 subtype (506
204 test against 44 references) and 772 for H3N2 (666 test against 130 references).

205 **2.1.3. AAindex1 Database**

206 The latest version of AAindex1, ver. 9.2, consists of 566 entries. A typical
207 database entry includes a vector of 20 numbers, each of which is assigned to a
208 standard amino acid. Since the range of numbers in vectors varies within the
209 database, we individually scaled each vector into the unit interval $[0,1]$. After
210 removing vectors with missing values, 553 remaining entries were used for analysis.

211 **2.2. Encoding of HA Sequences**

212 Here, we use RAAAs to take into account the impact of the mutation on
213 antigenic evolution from different physicochemical (amino acid) perspectives. The

214 first step of encoding the sequence by RAAA is selecting an arbitrary amino acid
 215 from each group in the alphabet as a group representative. Further, we replace all
 216 members of the group with its representative in the protein sequence. This step does
 217 not influence data if the standard amino acid alphabet is chosen since this alphabet
 218 has 20 groups, not less.

219 **2.3. Selection of High Impact Sites**

220 The model's input is a feature vector produced from encoded relevant sites in
 221 the genetic sequence. The model utilizes these sites for reconstructing the
 222 relationship between genetic and antigenic evolution in the feature space. Therefore,
 223 it is necessary to measure the relevance of site mutations according to the antigenic
 224 variation. Cui *et al.* [6] proposed measurement by introducing the below score for
 225 the site's antigenic significance:

$$S_i = |\Phi_i| \times E_i \quad (3)$$

226 where i is the index of the site in the sequence, S_i is the significance score, and E_i is
 227 Shannon's entropy of site i in the whole database as computed by the following
 228 formula:

$$E_i = - \sum_{j=1}^{20} P_{i,j} \log P_{i,j} \quad (4)$$

229 where $P_{i,j}$ is the probability of amino acid j occurrence at position i . Φ_i is a coefficient
 230 expressed with the following formula:

$$\Phi_i = \frac{(N_{11} \times N_{00} - N_{10} \times N_{01})}{\sqrt{N_{X1} \times N_{X0} \times N_{1Y} \times N_{0Y}}} \quad (5)$$

231 where N_{mn} ($m, n \in \{0, 1\}$) is the number of HI entries with $X=m$ and $Y=n$. The variable
 232 X represents the occurrence of mutation at site i (0 or 1 for conserved or mutated
 233 cases, respectively). The variable Y expresses the antigenic relationship between the
 234 test-reference pair of viruses in HI entries. If the test and reference are antigenically

235 similar, the Y variable value is zero. Otherwise, they are variants, and it takes the
236 value of one. $N_{X,n}$ denotes the number of entries with $Y=n$, whereas X can take any
237 value from $\{0,1\}$. Similarly, $N_{m,Y}$ represents the number of entries with $X=m$, while
238 Y has a value from $\{0,1\}$. Note that all variables in Eq. (5) are calculated only for
239 site i . In the case of a conserved site, the significance score is set to zero.

240 The application of Eq. (3) can be extended to sequences encoded by RAAAs.
241 Encoding genetic sequences by such an alphabet notably changes the entropy and Φ
242 values and, accordingly, the significance score. The significance score for all sites
243 obtained by applying a RAAA is further scaled into the unit interval $[0,1]$. This
244 allows us to compare the significance of a specific site considering different
245 alphabets. The final high-impact sites are determined by setting a threshold on the
246 results of the scaled significance score. The threshold value is selected from the set
247 $\{0.2,0.3,0.4,0.5,0.6,0.7,0.8\}$. It's worth noting that a site is selected if its scaled
248 significance score is higher than the target threshold. Obviously, decreasing the
249 threshold leads to an increase in the number of selected (high-impact) sites.

250 **2.4. Clustering the AAindex1 Database**

251 The AAindex1 database is used to perform the second stage of encoding. We
252 select some entries from AAindex1 (called representatives) that are further used to
253 encode the genetic information of obtained selected sites in the previous step. It is
254 known that there is a high correlation between AAindex1 entries. Therefore, we
255 cluster them and choose a representative from each cluster to diminish the number
256 and correlation of final features. Clustering should be performed so that the objects
257 of a cluster have almost the same encoding impact on antigenicity modeling.

258 **2.4.1. Computing Mutual Information**

259 To cluster the AAindex1 database, we create a feature vector for each entry
260 by a similar scenario as described [6] with a modification for RAAAs. In the
261 suggested method, the feature vector characterizes the AAindex1 entry by mutual

262 information (MI). The MI value expresses not only the significance of genetic
263 information but also the impact of encoding for a selected site individually. Note that
264 the size of the feature vector for clustering AAindex1 is the same as the size of
265 selected sites. Indeed, each element of the feature vector is the measure of mutual
266 dependency between the changes in a selected site, encoded by an AAindex1 entry,
267 and antigenic variants within the HI database.

268 The number of amino acids in the RAAA is less than in the standard alphabet.
269 Thus, a question arises on how encoding is carried out using AAindex1 regarding a
270 RAAA. In order to solve this issue, we define a new database, called pseudo-
271 AAindex1, derived from the original AAindex1 database. The procedure of
272 generating pseudo-AAindex1 is described in Figure 2.

273 As previously stated, each amino acid group has a representative, which
274 replaces all amino acids of the group in protein sequences. In order to assign a value
275 to the representative, we compute the average of AAindex1 values for the amino
276 acids within the group. This allows each amino acid to participate and have its own
277 effect through the representative. Thus, a pseudo-AAindex1 is created for each
278 RAAA, making it possible to calculate the mutual information in RAAA encoding.
279 For simplicity, we hereafter refer to both the original AAindex1 and the pseudo-
280 AAindex1 simply as AAindex1.

281 ***2.4.2. Determining the of optimal number of cluster***

282 When considering an alphabet, we create a feature vector for each AAindex1
283 entry, the size of which depends on the number of determined high-impact sites.
284 Before clustering AAindex1, it is necessary to determine the optimal number of
285 clusters. Indeed, this number affects the final feature vector, which is used for
286 antigenic variant modeling. For this purpose, we conduct a comprehensive search
287 for the optimal number by employing three algorithms: K-means; agglomerative
288 clustering with different linkage criteria; and spectral clustering.

289 First, we determine the number of unique feature vectors. Clustering is not
290 required if the number is less than a threshold (e.g., five). When the number of
291 unique vectors is more than the threshold, we cluster the set of vectors, while the
292 number of clusters starts from two and increases up to ten.

293 Generally, six clustering variants are applied, including: K-means;
294 agglomerative clustering with four different criteria (ward, average,
295 complete/maximum, and single/minimum); and spectral clustering. The obtained
296 clustering from each algorithm is individually evaluated by four scores, including
297 Silhouette, Calinski-Harabasz, Davies-Bouldin, and the sum of squared distances of
298 objects to their closest cluster. Further, the results are plotted and manually checked
299 to decide the optimal number of clusters for AAindex1 associated with an alphabet.

300 ***2.4.3. Clustering***

301 Generally speaking, the aim of clustering is to decrease correlation between
302 AAindex1 entries. This also leads to diminishing the number of features, which are
303 used in the final classification. To cluster the AAindex1 database, we apply the
304 associated clustering algorithm by which the optimal number of clusters was
305 determined from the previous subsection. Next, we select a representative from each
306 cluster. The representative of a cluster is the closest object to its center. The
307 representative is further employed to encode the information of high-impact sites for
308 the classification.

309 ***2.5. Classification of Antigenic Variants***

310 We use the obtained cluster representatives to apply the secondary encoding.
311 This is carried out by replacing an amino acid group representative in the selected
312 sites with its numerical value from the cluster's representative. Then, we
313 individually calculate the differences between the test and the reference strains for
314 each HI assay database entry by subtracting their encoded selected sites (or feature
315 vectors). If we denote the number of high-impact sites and number of clusters

316 representatives with N and M , respectively, then the final feature vector has the size
317 of $N \times M$.

318 Before performing the final classification, the last step is to determine the best
319 classifier. To decrease the effect of the classifier on the results, we consider five
320 different classifiers, including random forest, multilayer perceptron, logistic
321 regression, support vector, and Gaussian naïve Bayes. Each classifier has its own
322 parameters optimized through grid search (parameter list in Table II).

323 Grid search is carried out by cross-validation with different parameter
324 combinations. Note that the Gaussian naïve Bayes classifier has no parameters for
325 grid search, but it assumes that features are independent. Thus, we perform principal
326 component analysis on the feature matrix to decrease the dependency. A threshold
327 on the percentage of variance explained by the selected components was set as a
328 parameter for Gaussian naïve Bayes.

329 By comparing grid search results, we were able to choose the best classifier
330 with high performance in terms of accuracy. Note that the selection of optimal
331 classifier depends on the results of three procedures:

- 332● Encoding by the alphabet (primary encoding)
- 333● Selection of high-impact sites
- 334● Clustering the AAindex1 database and choosing representatives for secondary
335 encoding

336 Among these procedures, the first has the most decisive influence on
337 classification results. In fact, it changes the amino acid space globally, resulting in
338 different representations of genetic variation, as well as different relationships
339 between genotype and phenotype.

340

341 **III. RESULTS & DISCUSSION**

342 Considering all parameters, we ran 224,147 fits (41 alphabets \times 7
343 thresholds \times 781 5-fold cross-validations) for each subtype (H3N2 and H1N1) in the

344 experimental data to obtain the best classifiers. Knowing the best classifier for each
345 triple-combination case (subtype, alphabet, threshold), we performed a 10-fold
346 cross-validation by applying its best classifier. A comprehensive report of the results,
347 including the evaluation criteria, is publicly available at:
348 github.com/viroinformatics/Simplified_Alphabets.

349 The maximum accuracy achieved by each threshold is presented in Table III.
350 Since the length of the feature vector is decreased by increasing the threshold, this
351 also leads to accuracy reduction. From Table III, it is observed that threshold 0.4
352 seems to be a good choice for modeling the antigenic variants. Compared with
353 previous studies [10, 11], our results indicate a high degree of accuracy, especially
354 for H3N2, which suggests potential application in the field of public health.

355 As expected, some RAAAs achieved the same accuracy as the standard amino
356 acid alphabet. Table IV presents the alphabets with the highest performance for
357 different thresholds and subtypes. In the case of subtype H1N1 with thresholds 0.3
358 and 0.5, there are alphabets, the accuracy of which are slightly less (about 0.01) than
359 the standard and Prlic-SDM12-2000 alphabets, but are not added to the table. Since
360 prediction accuracy significantly drops from threshold 0.5 to threshold 0.8 (Table
361 III), we did not consider their results in Table IV. Interestingly, the Risler-88 and Li-
362 2003 alphabets are observed in the list of each subtype.

363 Moreover, the Cannata-2002 alphabet seems to be more informative for
364 subtype H3N2 rather than for H1N1. In some cases, e.g., subtype H1N1 with
365 threshold 0.4 and subtype H3N2 with threshold 0.3, the feature vector obtained from
366 RAAAs is shorter in length than that obtained from the standard alphabet, while their
367 accuracy is the same. This indicates that the amino acid space represented by the
368 standard alphabet has redundant dimensions to express genetic variation of antigenic
369 variants. Next, we briefly discuss each of the alphabets from Table IV.

370 Stephenson & Freeland analyzed 34 different RAAAs [29] and classified
371 them into five classes based on how grouping was carried out. The classes are
372 chemistry, sequence alignment, structural alignment, contact potential, and protein

373 blocks. Of the alphabets in Table IV, four are based on sequence alignment methods,
374 whereas two rely on structural alignment. Only one alphabet (Zou-2009) was created
375 by protein blocks. The complete classification of RAAAs presented in Table I is
376 based on published work [29].

377 Similarity between amino acids can be defined from various viewpoints, e.g.,
378 hydrophobic residues (I, V) and aromatic residues (F, W, Y). The main idea of
379 constructing a RAAA is to use amino acid properties to define similarity, with
380 placing of similar amino acids in a group. For example, the RAAA presented by Li
381 *et al.* [20] was obtained from amino acids substitutions by scoring similarities that
382 may be beneficial in recognition of protein folds. Their results imply that at least ten
383 amino acid types are required to characterize protein complexity.

384 Cannata *et al.* [5] presented a method to produce RAAAs by scoring different
385 amino acid compositions using a branch and bound algorithm and substitution
386 matrix. Their alphabet belongs to the 'alignment-based methods' class of sequences.
387 Furthermore, Lenckowski *et al.* [19] suggested an alphabet generated using a genetic
388 algorithm and strategy based on global sequence alignment. Their results indicate
389 that the proposed alphabet outperformed the standard amino acid set and other
390 RAAAs in the sequence classification task. Andersen and Brunak's RAAA [1]
391 includes 13 letters; it is also constructed based on sequence alignment. In contrast,
392 RAAAs proposed by Prlic *et al.* [23] and Risler *et al.* [25] are both derived by
393 substitution frequency of structural alignments.

394 Zou *et al.* [42] applied reduced amino acid alphabets to predict defensin
395 family and subfamily. They clustered amino acids by the protein blocks (PBs)
396 method [7, 9], in which the distribution of amino acids in PBs was used to generate
397 clusters of equivalent amino acids with respect to local structure. Indeed, this kind
398 of alphabet can be considered a structural alphabet. Their results indicate that use of
399 such alphabets can improve prediction accuracy with defensin family and subfamily.
400 Surprisingly, no alphabet based on attributes of individual amino acids attained a

401 high level of performance. Taken together, the high-performing RAAAs emphasize
402 the role of structural features in antigenic evolution modeling.

403 By checking the high-impact sites in the Influenza Research Database
404 (IRD) [40], we found that not only antigenic sites have a significant influence on
405 antigenicity, but also other amino acids located in close proximity. The results
406 indicate that all selected non-antigenic sites are related to immune responses.
407 According to IRD, these have been experimentally determined to be T-cell epitopes,
408 B-cell epitopes, and MHC-binding epitopes of different classes. This highlighted a
409 caveat: In modeling of antigenic evolution, the model should consider not only the
410 genetic information of antigenic sites, but also that of neighboring positions, as they
411 may indirectly impact antigenicity. Note that feature vector construction relies on
412 high-impact sites, but the evolutionary history showed that even one amino acid
413 substitution can change the antigenic cluster of the influenza virus [28]. Such a
414 substitution may present a low impact through the mutual information score. We
415 believe a desirable model must take into account the effects of both high and low
416 impact sites. The visualizations of selected high-impact sites for H1N1 (threshold
417 0.3), and H3N2 (threshold 0.4), are presented in Figure 3. These are cases with high
418 accuracy and shorter feature vector length.

419 Various AAindex1 entries were designated as representatives during all
420 experiments with optimized classifiers. The top ten entries and their frequencies are
421 listed in Table V. The complete list of AAindex1 entries and their frequencies is
422 available (github.com/viroinformatics/Simplified_Alphabets).

423 It can be seen that the majority of AAindex1 attributes used in model
424 construction are associated with charge properties. This emphasizes that antigenicity
425 notably depends on protein conformation, which cannot be fully reflected in a one-
426 dimensional representation of protein as a sequence. However, the model can capture
427 some attributes information by encoding the genetic sequence using
428 physicochemical properties presented in the AAindex1 database.

429 Table V also indicates that nine out of the ten most frequently AAindex1
430 entries are common in both subtypes. The last AAindex1 entry in the list of each
431 subtype is different. To better understand the characteristics of entries in Table V, we
432 computed the Pearson correlation coefficient and visualized it in Figure 4. It is
433 observed that the majority of entries in Table V are not correlated, with two
434 exceptions: FAUJ880111/KLEP840101 and FAUJ880105/CHAM830103.

435 Although the uncommon entries between H1N1 and H3N2 are different, it can
436 be seen that they are correlated. In addition to the correlation matrix, we used
437 principal component analysis (PCA) to identify the main components of 11 distinct
438 AAindex1 entries in Table V and their expression in terms of explained variance.
439 Figure 5 indicates that the first six components describe more than 90% of the
440 explained variance. Seven and nine components represent 95% and 99% of the
441 explained variance, respectively.

442 We also considered the performance of classifiers for antigenic variant
443 modeling. Among five classifiers, random forest and multilayer perceptron
444 outperformed others, in terms of accuracy, for both the H1N1 and H3N2 subtypes.
445 The Gaussian naive Bayes classifier gave the worst results, so it may not be suitable
446 for this kind of modeling.

447 In summary, the outstanding ability of our approach is based on redefining the
448 mutation by RAAA and amino acid attributes used for encoding through a two-fold
449 procedure. The primary encoding plays the main role with high priority, whereas
450 secondary encoding has a supplementary role. From one point of view, the primary
451 encoding determines the high-impact sites, while the secondary encoding gives the
452 numerical interpretation to the genetic information of selected sites. From another
453 point of view, the primary encoding interprets the mutation, and the secondary
454 encoding reconstructs the specific relationship between genetic and antigenic
455 differences (for the test and reference strains).

456 The proposed two-fold encoding approach revealed some aspects of
457 mutations related to the antigenicity. Our findings indicate that encoding associated

458 with structural or charge properties of the protein dramatically impacts the
459 performance of the antigenic model. This is in agreement with recent studies done
460 by other researchers [14, 39]. In addition, RAAA encoding can lead to a smaller
461 feature space dimension, while performance is maintained or improved. So far, this
462 approach was applied only to seasonal human influenza strains. However, there are
463 no theoretical limitations that would prevent further testing as a universal
464 computational model for predicting antigenicity in other influenza subtypes, such as
465 zoonotic H5, H7, H9, or other relevant influenza A subtypes that cause sporadic
466 human infection.

467 **III. Conclusion**

468 Determining the degree of antigenic similarity between influenza virus strains
469 is crucial in choosing candidate vaccine strains and subsequent timely vaccine
470 production. Currently, the degree is measured via HI assay, a widespread laboratory
471 procedure. Although HI assay is the gold standard method, it suffers from several
472 shortcomings. Therefore, it has been suggested to employ computer-aided models as
473 auxiliary tools to assess preliminary information about viral antigenicity prior to HI
474 assay.

475 A notable problem in modeling antigenic evolution is the representation of
476 genetic information to better express the relationship between genetic and antigenic
477 variations. This paper proposes a two-fold encoding approach to genetic information
478 using both a reduced amino acid alphabet (RAAA) and an amino acid index
479 database. By applying a RAAA, we redefine the mutation as changes between amino
480 acid groups of the alphabet, while the output sequence of the primary encoding is
481 still alphabetical. The secondary encoding uses representatives from the AAindex1
482 database to convert the alphabetical sequence of the primary encoding into the
483 numerical. The experimental results indicate that models built using RAAA
484 encoding are able to achieve the same accuracy as models using the standard amino
485 acid alphabet. The RAAA-based approach, however, features reduced computational
486 complexity and associated cost.

487 Moreover, the suggested encoding can reveal the amino acid attributes which
488 drive antigenic evolution. In agreement with previous studies, we find that structural
489 and charge characteristics are the most beneficial in modeling antigenic evolution.
490 Although the results obtained by our approach are desirable and promising, they are
491 achieved by taking into account only high-impact sites. It is known that even one
492 substitution can change the antigenic cluster, so we believe that further incorporating
493 the role of low-impact sites into the model may enhance its accuracy and prediction
494 potential; this will be addressed in future studies. Additionally, the model can be
495 improved by: introducing new reduced amino acid alphabets; employing more
496 significant and descriptive criteria for selecting key sites; and incorporating
497 neighboring amino acid effects into the model.

498 Computational approaches for predicting antigenic properties from genetic
499 sequence are also quite relevant for highly virulent influenza viruses. Laboratory
500 testing of these pathogens requires high biosafety certification levels, and such
501 analysis is not only time-consuming and labor-intensive, but also costly. Unlike
502 current laboratory approaches, computational prediction of antigenic properties from
503 viral sequence has the potential to enable rapid, large-scale antigenic
504 characterization of influenza viruses. It is worth mentioning that application of our
505 approach is not limited to modeling of antigenic evolution. It can be used in
506 modeling any phenotype that is based on protein sequence, such as interactions with
507 monoclonal antibodies.

508

509 **IV. Funding**

510 This study was funded by the Russian Foundation for Basic Research (RFBR),
511 project number 19-31-60025.

512

513 **V. Acknowledgments**

514 Our work was performed using the 'Uran' supercomputer (IMM UB RAS).
515 The authors would like to thank Edward S. Ramsay for his valuable improvements
516 on the manuscript.

TABLES

Table I. The list of alternative and standard amino acid alphabets employed to encode the protein sequences in our experiments. The alphabets are borrowed from

Таблица I. Список исследованных альтернативных и стандартных аминокислотных алфавитов, использованных для кодирования белковых последовательностей. Алфавиты заимствованы из [8, 29, 38]. Цвет RAAA определяет метод его получения. Классификация алфавитов заимствована из [29].

Name of alphabet Название алфавита	Groups группы
Standard Стандартный	A, C, D, E, Q, F, Y, G, H, I, V, K, R, L, M, N, P, S, T, W
Hydrophobic Гидрофобный	RKEDQN, GASTPHY, CVLIMFW
V Объем Ван-дер-Ваальса	GASCTPD, NVEQIL, MHKFRYW
Polarity Полярность	LIFWCMVY, PATGS, HQRKNE
Polarizability Поляризуемость	GASDT, CPNVEQIL, KMHFRYW
Mahler 1966	DE, KRH, QN, ST, P, CM, WYF, GALIV
Lehninger 1970	DE, KRH, NQSTGQY, PAWFMLIV
Dickerson 1983	DENKRQH, STGPACWY, FMLIV
Taylor 1986	DE, N, KRH, Q, T, SGAC, P, YWF, M, LIV
Weathers 2004	DENRH, KQST, GPACM, WYFLIV
SE-B(14)	A, C, D, EQ, FY, G, H, IV, KR, LM, N, P, ST, W
SE-B(8)	AST, C, DHN, EKQR, FWY, G, ILMV, P
Risler 1988	D, E, N, KRQ, S, T, G, P, H, A, C, W, YF, ML, IV
Riddle 1997	DE, NKRQS, THA, GP, CWYFMLIV
Mirny 1999	DE, KR, NQST, GP, HWYF, ACMLIV
Prlic SDM12 2000	D, N, EKR, QST, G, P, H, A, C, W, YF, MLIV
Prlic SDM17 2000	D, EK, N, R, Q, S, T, G, P, H, A, C, W, Y, F, M, LIV
Melo 2005	DENKRQSTP, GA, H, C, WYFMLIV
Robson 1976	DKR, EA, GP, STNQ, H, C, WY, FMLIV
Solis(G) 2000	D, N, S, T, G, P, H, C, Y, EKRQAWFMLIV
Solis(D) 2000	DNS, EKRQ, TH, GP, AM, C, W, F, YL, IV
Rogov 2001	DNSTA, EKRQ, G, P, H, C, W, M, YFLIV
Etchebest 2007	DN, EKRQ, SH, TC, G, P, WYF, AML, IV
Solis GBMR4 2009	DENKRQSTA, G, P, HCWYFMLIV
Zuo 2009	DN, E, KRQ, SH, T, G, P, A, C, WYF, M, L, IV
Dayhoff 1978	DENQ, KRH, STGPA, C, WYF, MLIV
Murphy 2000	DENQ, KR, ST, G, P, H, A, C, WYF, MLIV
Cannata 2002	D, E, N, KR, Q, ST, G, P, H, A, C, W, Y, F, ML, IV
Fan 2003	DEQ, KR, STA, G, P, NH, C, WYF, ML, IV
Li 2003	DE, KRQ, ST, G, P, NH, AC, WYF, ML, IV
Edgar Se-B 2003	DN, EQ, KR, STA, G, P, HW, C, YF, MLIV
Edgar Se-V 2003	DEN, KRQ, STA, G, P, H, C, W, YF, MLIV

RAAA ENCODING & ANTIGENIC EVOLUTION

Kosiol 2003	DENKRQSTGPHA, C, W, YF, MLIV
Anderson 2004	D, E, KRQ, NS, T, G, P, H, A, C, WYF, ML, IV
Lenckowski 2007	DSHFM, ERQL, KPAC, NTWY, GIV
Crippen 1990	ENRSGHY, DKQTPW, AV, CFMLI
Maiorov 1992	DENQ, KR, G, P, AV, STHWY, CFMLI
Thomas 1996	DE, KR, QSTNGPH, C, AWYFMLIV
Wang 1999	DE, NKRQS, GP, THA, CWYFMLIV
Ceiplak 2001	DENRQSTG, K, HA, CWYMV, FLI
Liu 2002	DE, KR, NQSTGPHY, ACW, FMLIV

Amino acid physicochemical attributes Физико-химические свойства аминокислот
Substitution frequency -- Structural alignment Частота замещения -- структурное выравнивание
Spatial frequency -- Protein blocks Пространственная частота -- Белковые блоки
Substitution frequency -- Sequence alignment Частота замены -- Выравнивание последовательности
Spatial frequency -- Contact potential Пространственная частота -- Контактный потенциал

Table II. Parameters used in the grid search. Parameter names are based on the machine learning package Scikit-learn [22].

Таблица 2. Параметры, используемые при поиске по сетке. Имена параметров основаны на пакете машинного обучения Scikit-learn [22].

Method Метод	Parameters Параметры	Total cases in grid search Общее число случаев при поиске по сетке
Random forest алгоритм случайного леса	Criterion, n_estimator, min_samples_split	Fitting 5-fold cross-validation for each of 40 candidates, totaling 200 fits фитинг 5-кратной перекрестной проверки для каждого из 40 кандидатов, всего 200 фитингов
Logistic regression Логистическая регрессия	Solver, penalty, max_iter	Fitting 5 folds for each of 66 candidates, total 330 fits фитинг 5 кратностей для каждого из 66 кандидатов, всего 330 фитингов
Multilayer perceptron Многослойный перцептрон	Solver, learning_rate, activation, max_iter, learning_rate_init, hidden_layer_sizes	Fitting 5 folds for each of 648 candidates, total 3240 fits Подгонка 5 кратностей для каждого из 648 кандидатов, всего 3240 фитингов
SVM	Kernel, gamma, C, degree	Fitting 5 folds for each of 24 candidates, total 120fits фитинг 5 кратностей для каждого из 24 кандидатов, всего 120 фитингов
Gaussian naive bayes Гауссовский наивный байесовский классификатор	PCA_threshold	Fitting 5 fold for each of 3 candidates, total 15 fits фитинг 5 кратностей для каждого из 3 кандидатов, всего 15 фитингов

Table III. The maximum accuracy was obtained by 10-fold cross-validation using different thresholds for scaled significance score. Note that increasing the threshold may lead to shorter feature vector length and consequent reduced accuracy.

Таблица 3. Максимальная точность была получена путем 10-кратной перекрестной проверки с использованием различных пороговых значений для масштабируемой оценки значимости. Увеличение порога может привести к уменьшению длины вектора признаков и, как следствие, к снижению точности.

Threshold							
Порог	0.2	0.3	0.4	0.5	0.6	0.7	0.8
Subtype							
Подтип							
H1N1	0.88	0.88	0.87	0.84	0.77	0.77	0.77
H3N2	0.92	0.92	0.92	0.9	0.88	0.85	0.83

Table IV. High performance alphabets in our experiments by virus type and site selection threshold. The amino acid groups for each alphabet are presented in Table I.

Таблица 4. Высокоэффективные алфавиты в экспериментах по типу вируса и порогу выбора сайта. Аминокислотные группы для каждого алфавита представлены в таблице I.

Threshold Порог	H1N1	H3N2
0.2	Standard, Risler-88, Li-2003, Anderson-2004 Стандартный, Risler-88, Li-2003, Anderson-2004	Standard, Risler-88, Cannata-2002, Zuo-2009 Стандартный, Risler-88, Cannata-2002, Zuo-2009
0.3	Standard* Стандартный *	Standard, Cannata-2002 Стандартный, Cannata-2002
0.4	Standard, Lenckowski-2007 Стандартный, Lenckowski-2007	Standard, Cannata-2002 Стандартный, Cannata-2002
0.5	Prlic-SDM12-2000*	Risler-88, Li-2003, Zuo-2009

* There are other alphabets, not listed, whose accuracy was slightly less than the alphabets shown in the table.

* Существуют и другие алфавиты, не указанные в списке, точность которых была несколько меньше, чем у алфавитов, приведенных в таблице.

Table V. List of the top ten most frequent AAindex1 entries in the experiments with optimized classifiers.

Таблица V. Список первых десяти наиболее частых записей AAindex1 в экспериментах с оптимизированными классификаторами.

ID Идентификатор	Description Описание	Freq. Частота
H1N1		
ANDN920101	alpha-CH chemical shifts химические сдвиги альфа-СН (Andersen <i>et al.</i> , 1992)	147
CHAM830104	The number of atoms in the side chain labelled 2+1 Количество атомов в боковой цепи, помеченной 2+1 (Charton-	106
KLEP840101	Net charge Результирующий заряд (Klein <i>et</i>	97
CHAM830103	The number of atoms in the side chain labelled 1+1 Количество атомов в боковой цепи, помеченной 1+1 (Charton-Charton, 1983)	92
FAUJ880111	Positive charge Положительный заряд (Fauchere <i>et al.</i> , 1988)	83
CHAM830107	A parameter of charge transfer capability Параметр способности переноса заряда (Charton-Charton, 1983)	83
VENT840101	Bitterness Горечь (Venanzi, 1984)	74
FAUJ880112	Negative charge Отрицательный заряд (Fauchere <i>et al.</i> , 1988)	70
FAUJ880105	STERIMOL minimum width of the side chain минимальная ширина боковой цепи (Fauchere <i>et al.</i> , 1988)	47
CHAM830105	The number of atoms in the side chain labelled 3+1 Количество атомов в боковой цепи, помеченной как 3+1 (Charton-Charton, 1983)	38
H3N2		
VENT840101	Bitterness Горечь (Venanzi, 1984)	119
CHAM830103	The number of atoms in the side chain labelled 1+1 Количество атомов в боковой цепи, помеченной 1+1 (Charton-Charton, 1983)	117
FAUJ880111	Positive charge Положительный заряд (Fauchere <i>et al.</i> , 1988)	101
ANDN920101	alpha-CH chemical shifts химические сдвиги альфа-СН (Andersen <i>et al.</i> , 1992)	101
KLEP840101	Net charge Результирующий заряд (Klein <i>et al.</i> , 1984)	88
FAUJ880112	Negative charge Отрицательный заряд (Fauchere <i>et al.</i> , 1988)	87
CHAM830107	A parameter of charge transfer capability Параметр способности переноса заряда (Charton-Charton, 1983)	66
CHAM830104	The number of atoms in the side chain labelled 2+1 Количество атомов в боковой цепи, помеченной 2+1 (Charton-Charton, 1983)	60

FAUJ880105	STERIMOL minimum width of the side chain STERIMOL минимальная ширина боковой цепи (Fauchere <i>et al.</i> , 1988)	59
FAUJ880109	Number of hydrogen bond donors Количество доноров водородных связей (Fauchere <i>et al.</i> , 1988)	58

FIGURES

Figure 1. General scheme of the computational pipeline. It consists of five parts: encoding HA sequences by a reduced amino acid alphabet; selecting significant sites; clustering the AAindex1 database using mutual information of selected sites; encoding the sites by a representative from each cluster; and finally training the classifier.

Рис. 1. Общая схема вычислительного конвейера. Он состоит из пяти частей: кодирование последовательностей HA сокращенным аминокислотным алфавитом; выбор значимых участков; кластеризация базы данных AAindex1 с использованием взаимной информации выбранных сайтов; кодирование сайтов представителем от каждого кластера; и, наконец, обучение классификатора.

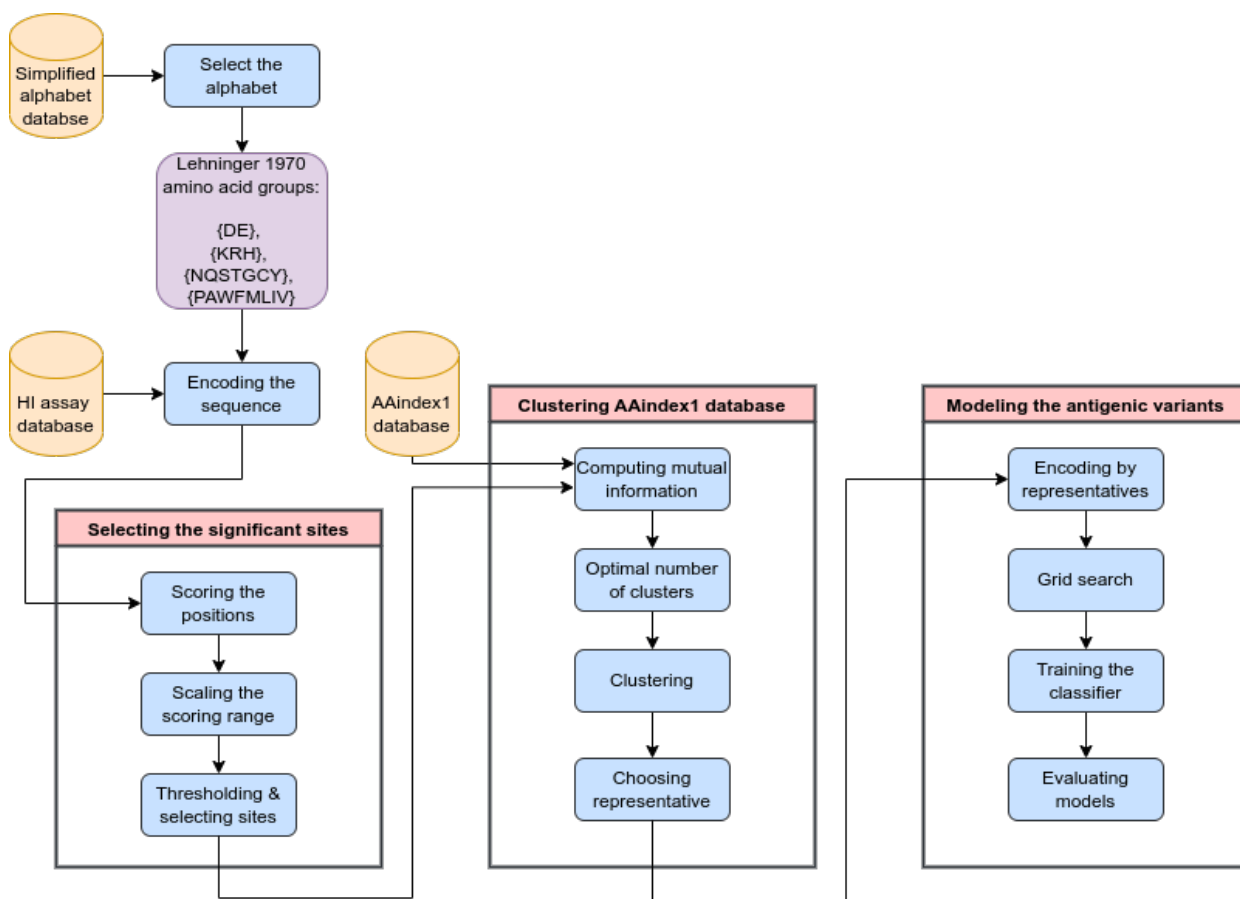


Figure 2. Generation of the pseudo-AAindex1 database from the hydrophobicity index. The pseudo database is created based on the selected RAAA. Note that the value assigned to each group is the average of the group's amino acid values in the scaled AAindex1 vector.

Рисунок 2. Создание базы данных псевдо-AAindex1 из индекса гидрофобности. Псевдобаза данных создается на основе выбранного RAAA. Обратите внимание, что значение, присвоенное каждой группе, представляет собой среднее значение величин аминокислот группы в масштабированном векторе AAindex1.

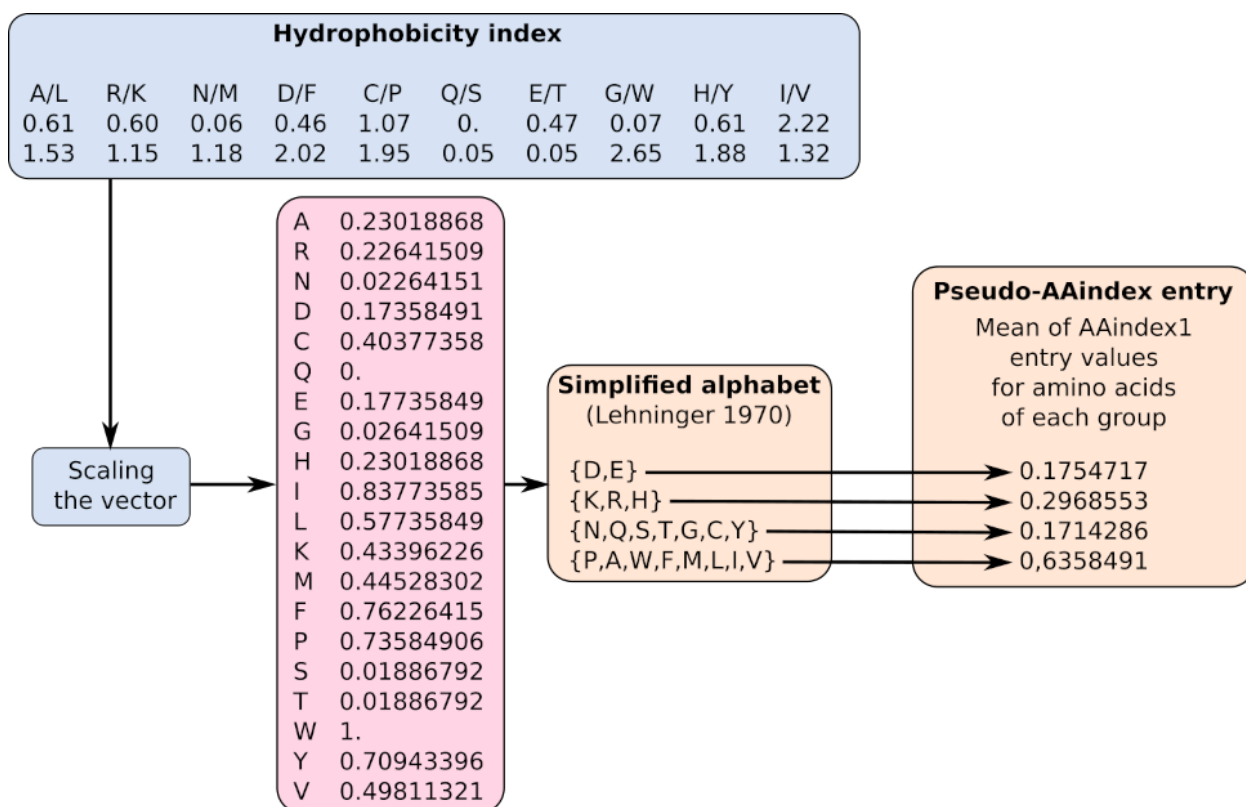


Figure 3. Visualization of high-impact sites on the surface of hemagglutinin protein by PyMOL [26]. Top – H1 protein (PDB ID: 1RUY [3, 12]). Bottom – H3 protein (PDB ID: 5THF [3, 33]). Note that the highlighted sites include not only the antigenic sites but also those experimentally determined as T-cell epitopes, B-cell epitopes, as well as MHC-binding epitopes of different classes.

Рисунок 3. Визуализация участков сильного воздействия на поверхности белка гемагглютинаина с помощью PyMOL [26]. Вверху – белок H1 (PDB ID: 1RUY [3, 12]). Внизу – белок H3 (PDB ID: 5THF [3, 33]). Обратите внимание, что выделенные сайты включают не только антигенные сайты, но и те, которые экспериментально определены как Т-клеточные эпитопы, В-клеточные эпитопы, а также МНС-связывающие эпитопы разных классов.

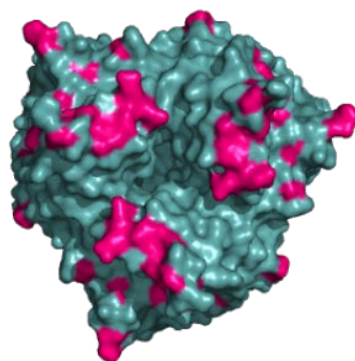
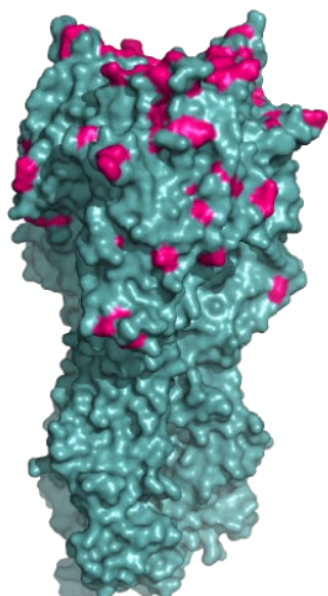
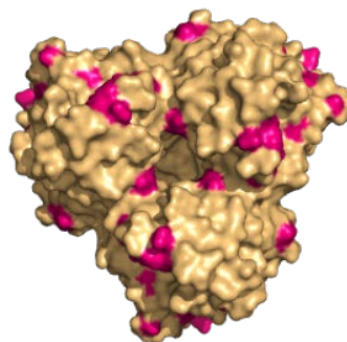
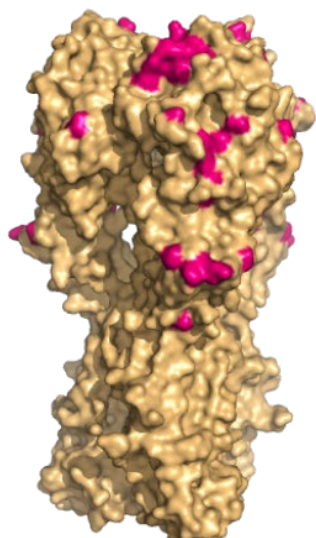


Figure 4. Correlation matrix of 11 unique AAindex1 entries from Table V. Note that the majority of indices have low correlation.

Рисунок 4. Матрица корреляции 11 уникальных записей AAindex1 из таблицы V. Обратите внимание, что большинство индексов имеют низкую корреляцию.

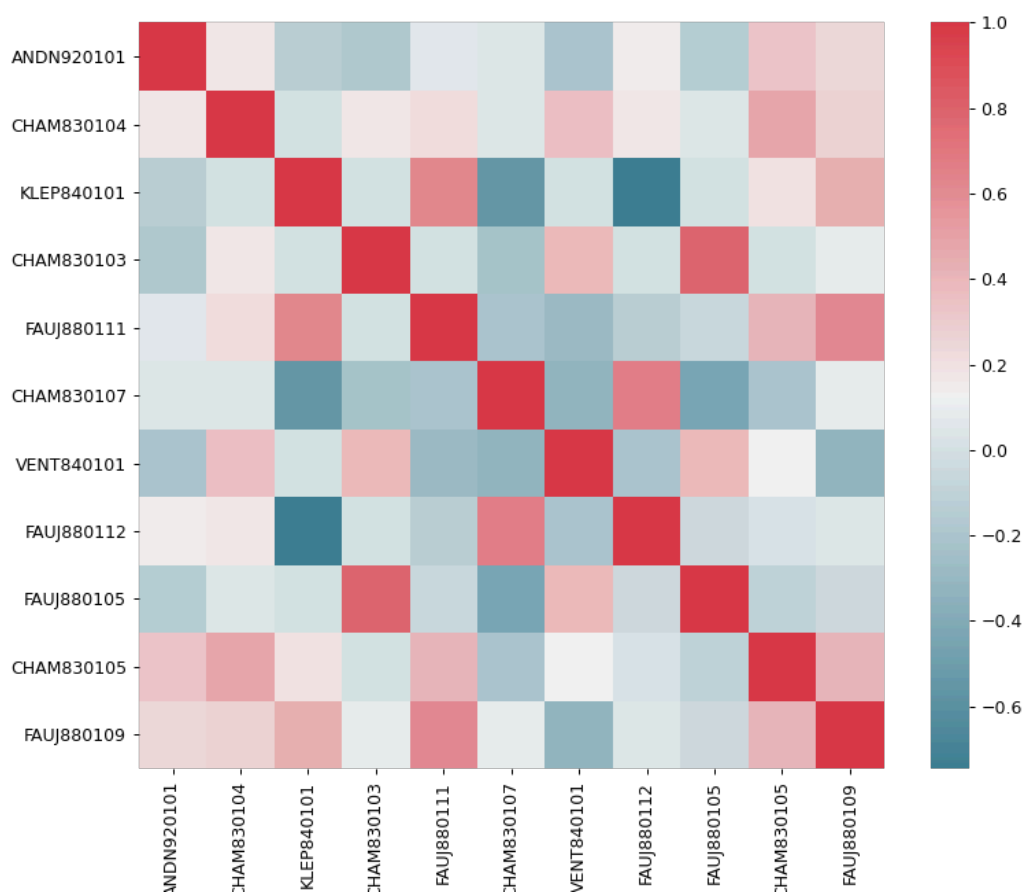
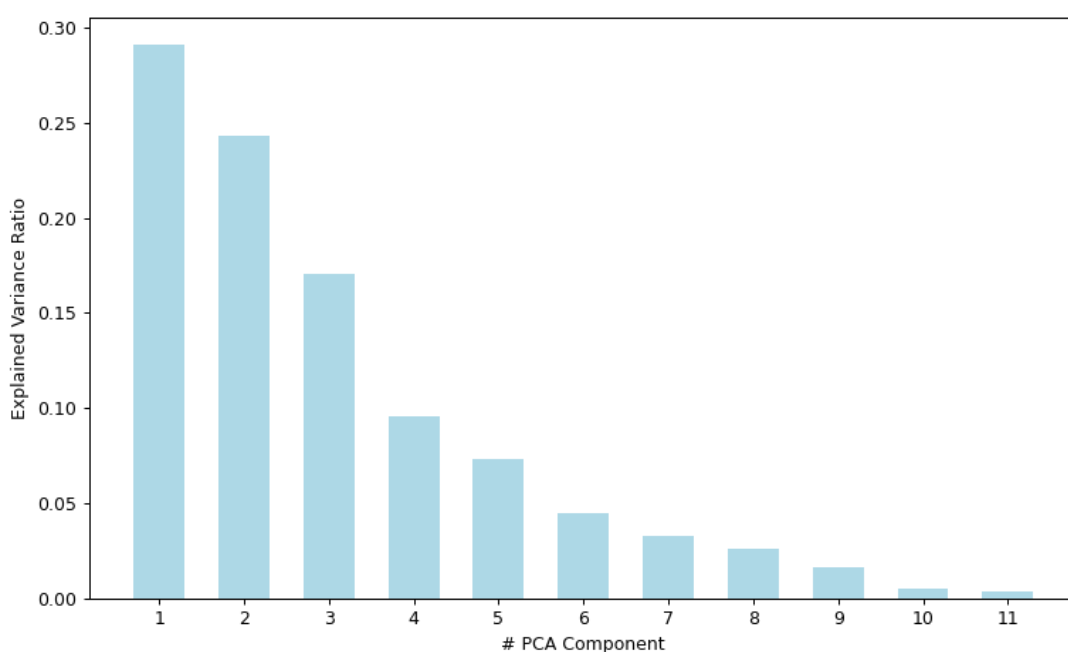


Figure 5. Explained variance ratios for PCA analysis components. Analysis was performed for 11 unique AAindex1 indices from Table V. The result shows that the first six components represent more than 90% of the identified variance.

Рисунок 5. Выявленные коэффициенты дисперсии для компонентов анализа PCA. Анализ был проведен для 11 уникальных индексов AAindex1 из таблицы V. Результат показывает, что первые шесть компонентов представляют более 90% выявленной дисперсии.



TITLE PAGE_METADATA

CORRESPONDING AUTHOR: Forghani Majid – Ph.D. in Physics and Mathematics, researcher (N.N. Krasovskii Institute of Mathematics and Mechanics), senior Researcher (Ural Federal University), 16 S.Kovalevskaya Str., Yekaterinburg, 620108 Russia, +7 (906) 811 9201, forghani@imm.uran.ru

Форгани Маджид – кандидат физико-математических наук, научный сотрудник отдела математического программирования, Федеральное государственное бюджетное учреждение науки Институт математики и механики им. Н. Н. Красовского Уральского отделения Российской академии наук (ИММ УрО РАН), старший научный сотрудник Уральского Федерального Университета (УрФУ)

620108, Россия, г. Екатеринбург, ул. Софьи Ковалевской, д. 16, +7 (906) 811 9201, forghani@imm.uran.ru

Фирстков А.Л – математик первой категории отдела математического программирования, Федеральное государственное бюджетное учреждение науки Институт математики и механики им. Н. Н. Красовского Уральского отделения Российской академии наук (ИММ УрО РАН)

Даниленко Д.М. – кандидат биологических наук, заместитель директора по научной работе, заведующий отделом этиологии и эпидемиологии (ФГБУ Научно-исследовательский институт гриппа имени А.А. Смородинцева, Минздрава России)

Комиссаров А.Б. – заведующий лабораторией молекулярной вирусологии (ФГБУ Научно-исследовательский институт гриппа имени А.А. Смородинцева, Минздрава России)

Running head/Краткое название: Вложение САА и антигенная эволюция/RAAA encoding & antigenic evolution

Keywords: AAindex, antigenic evolution, hemagglutinin, influenza, modeling, reduced amino acid alphabet

Ключевые слова: AAindex, антигенная эволюция, гемагглютинин, грипп, моделирование, сокращённый аминокислотный алфавит

Original article

20 pages, 5 table, 5 figure

31.05.2022

REFERENCES

Порядковый номер ссылки	Авторы, название публикации и источника, где она опубликована, выходные данные	ФИО, название публикации и источника на английском	Полный интернет-адрес (URL) цитируемой статьи и/или
1	Andersen, C. A., & Brunak, S., Representation of protein-sequence information by amino acid subalphabets. AI magazine, 2004, vol. 25, no. 1, pp. 97–97	---	https://doi.org/10.1609/aimag.v25i1.1750 [10.1609/aimag.v25i1.1750]
2	Arinaminpathy, N., & Grenfell, B. Dynamics of glycoprotein charge in the evolutionary history of human influenza. PloS one, 2010, vol. 5, no. 12, pp. e15674.	---	https://doi.org/10.1371/journal.pone.0015674 [10.1371/journal.pone.0015674]

3	Berman H. M. et al. The protein data bank. Nucleic acids research. Vol. 28, no. 1, pp. 235–242, 2000.	---	https://doi.org/10.1093/nar/28.1.235 [10.1093/nar/28.1.235]
4	Burns, A., Van der Mensbrugghe, D., & Timmer, H. Evaluating the economic consequences of avian influenza, in Plastics, World Bank Washington, DC, 2006.	---	https://www.academia.edu/download/72419367/474170WP0Evalu101PUBLIC10Box334133B.pdf
5	Cannata, N., Toppo, S., Romualdi, C., & Valle, G. Simplifying amino acid alphabets by means of a branch and bound algorithm and substitution matrices. Bioinformatics, vol. 18, no. 8, pp. 1102–1108, 2002.	---	https://doi.org/10.1093/bioinformatics/18.8.1102 [10.1093/bioinformatics/18.8.1102]

6	Cui, H., Wei, X., Huang, Y., Hu, B., Fang, Y., & Wang, J. Using multiple linear regression and physicochemical changes of amino acid mutations to predict antigenic variants of influenza A/H3N2 viruses. <i>Bio-medical materials and engineering</i> , vol. 24, no. 6, pp. 3729–3735, 2014.	---	https://doi.org/10.3233/BME-141201 [10.3233/BME-141201]
7	de Brevern, A. G. New assessment of a structural alphabet. <i>In silico biology</i> , vol. 5, no. 3, pp. 283–289, 2005.	---	https://content.iospress.com/articles/in-silico-biology/isb00186
8	Edgar, R. C. Local homology recognition and distance measures in linear time using compressed amino acid alphabets. <i>Nucleic acids research</i> , vol. 32, no. 1, pp. 380–385, 2004.	---	https://doi.org/10.1093/nar/gkh180 [10.1093/nar/gkh180]

9	Etchebest, C., Benros, C., Bornot, A., Camproux, A. C., & De Brevern, A. G. A reduced amino acid alphabet for understanding and designing protein adaptation to mutation. <i>European Biophysics Journal</i> , vol. 36, no. 8, pp. 1059–1069, 2007.	---	https://doi.org/10.1007/s00249-007-0188-5 [10.1007/s00249-007-0188-5]
10	Forghani, M., & Khachay, M. Convolutional Neural Network Based Approach to in Silico Non-Anticipating Prediction of Antigenic Distance for Influenza Virus. <i>Viruses</i> , vol. 12, no. 9, pp. 1019, 2020.	---	https://doi.org/10.3390/v12091019 [10.3390/v12091019]
11	Forghani, M., Khachay, M., & AlyanNezhadi, M. M. The Impact of Amino Acid Encoding on the Prediction of Antigenic Variants. In 2020 6th Iranian Conference on Signal Processing and Intelligent Systems (ICSPIS), pp. 1–5, 2020.	---	https://doi.org/10.1109/ICSPIS51611.2020.9349560 [10.1109/ICSPIS51611.2020.9349560]

12	Gamblin, S. J., et al. The structure and receptor binding properties of the 1918 influenza hemagglutinin. <i>Science</i> , vol. 303, no. 5665, pp. 1838–1842, 2004.	---	https://doi.org/10.1126/science.1093155 [10.1126/science.1093155]
13	Gregory, V., et al. Human former seasonal Influenza A (H1N1) haemagglutination inhibition data 1977-2009 from the WHO Collaborating Centre for Reference and Research on Influenza, London, UK. University of Glasgow, 2016.	---	http://dx.doi.org/10.5525/gla.researchdata.289 [10.5525/gla.researchdata.289]
14	Huang, Z. Z., Yu, L., Huang, P., Liang, L. J., & Guo, Q. Charged amino acid variability related to N-glyco-sylation and epitopes in A/H3N2 influenza: Hem-agglutinin and neuraminidase. <i>PloS one</i> , vol. 12, no. 7, pp. e0178231, 2017.	---	https://doi.org/10.1371/journal.pone.0178231 [10.1371/journal.pone.0178231]

15	Kawashima, S., Pokarowski, P., Pokarowska, M., Kolinski, A., Katayama, T., & Kanehisa, M. AAindex: amino acid index database, progress report 2008. Nucleic acids research, vol. 36, no. suppl 1, pp. D202–D205, 2007.	---	https://doi.org/10.1093/nar/gkm998 [10.1093/nar/gkm998]
16	Klingen, T. R., Reimering, S., Guzmán, C. A., & McHardy, A. C. In silico vaccine strain prediction for human influenza viruses. Trends in microbiology, vol. 26, no. 2, pp. 119–131, 2018.	---	https://doi.org/10.1016/j.tim.2017.09.001 [10.1016/j.tim.2017.09.001]
17	Kobayashi, Y., & Suzuki, Y. Compensatory evolution of net-charge in influenza A virus hemagglutinin. PloS one, vol. 7, no. 7, pp. E40422, 2012.	---	https://doi.org/10.1371/journal.pone.0040422 [10.1371/journal.pone.0040422]

18	Lee, M. S., & Chen, J. S. E. Predicting antigenic variants of influenza A/H3N2 viruses. <i>Emerging infectious diseases</i> , vol. 10, no. 8, pp. 1385, 2004.	---	https://doi.org/10.3201%2Faid1008.040107 [10.3201%2Faid1008.040107]
19	Lenckowski, J., & Walczak, K. Simplifying amino acid alphabets using a genetic algorithm and sequence alignment. In <i>European Conference on Evolutionary Computation, Machine Learning and Data Mining in Bioinformatics</i> , pp. 122–131, 2007.	---	https://doi.org/10.1007/978-3-540-71783-6_12 [10.1007/978-3-540-71783-6_12]
20	Li, T., Fan, K., Wang, J., & Wang, W. Reduction of protein sequence complexity by residue grouping. <i>Protein Engineering</i> , vol. 16, no. 5, pp. 323–330, 2003.	---	https://doi.org/10.1093/protein/gzg044 [10.1093/protein/gzg044]

21	Nanni, L., & Lumini, A. A genetic approach for building different alphabets for peptide and protein classification. BMC bioinformatics, vol. 9, no. 1, pp. 1–10, 2008.	---	https://doi.org/10.1186/1471-2105-9-45 [10.1186/1471-2105-9-45]
22	Pedregosa, F., et al. Scikit-learn: Machine learning in Python. Journal of Machine Learning Research, vol. 12, pp. 2825–2830, 2011.	---	https://www.jmlr.org/papers/volume12/pedregosa11a/pedregosa11a.pdf?ref=https://githubhelp.com
23	Prlić, A., Domingues, F. S., & Sippl, M. J. Structure-derived substitution matrices for alignment of distantly related sequences. Protein Engineering, vol. 13, no. 8, pp. 545–550, 2000.	---	https://doi.org/10.1093/protein/13.8.545 [10.1093/protein/13.8.545]

24	Qiu, J., Qiu, T., Yang, Y., Wu, D., & Cao, Z. Incorporating structure context of HA protein to improve antigenicity calculation for influenza virus A/H3N2. <i>Scientific reports</i> , vol. 6, no. 1, pp. 1–9, 2016.	---	https://doi.org/10.1038/srep31156 [10.1038/srep31156]
25	Risler, J. L., Delorme, M. O., Delacroix, H., & Henaut, A. Amino acid substitutions in structurally related proteins a pattern recognition approach: Determination of a new and efficient scoring matrix. <i>Journal of molecular biology</i> , vol. 204, no. 4, pp. 1019–1029, 1988.	---	https://doi.org/10.1016/0022-2836(88)90058-7 [10.1016/0022-2836(88)90058-7]
26	Schrödinger, L. L. C. The PyMOL molecular graphics system, version 1.8, 2015.	---	https://pymol.org/2/

27	Smith, D. J., Forrest, S., Ackley, D. H., & Perelson, A. S. Variable efficacy of repeated annual influenza vaccination. <i>Proceedings of the National Academy of Sciences</i> , vol. 96, no. 24, pp. 14001–14006, 1999.	---	https://doi.org/10.1073/pnas.96.24.14001 [10.1073/pnas.96.24.14001]
28	Smith, D. J., Lapedes, A. S., De Jong, J. C., Bestebroer, T. M., Rimmelzwaan, G. F., Osterhaus, A. D., & Fouchier, R. A. Mapping the antigenic and genetic evolution of influenza virus. <i>Science</i> , vol. 305, no. 5682, pp. 371–376, 2004.	---	https://doi.org/10.1126/science.1097211 [10.1126/science.1097211]
29	Stephenson, J. D., & Freeland, S. J. Unearthing the root of amino acid similarity. <i>Journal of molecular evolution</i> , vol. 77, no. 4, pp. 159–169, 2013.	---	https://doi.org/10.1007/s00239-013-9565-0 [10.1007/s00239-013-9565-0]

30	Su, S., Fu, X., Li, G., Kerlin, F., & Veit, M. Novel Influenza D virus: Epidemiology, pathology, evolution and biological characteristics. <i>Virulence</i> , vol. 8, no. 8, pp. 1580–1591, 2017.	---	https://doi.org/10.1080/21505594.2017.1365216 [10.1080/21505594.2017.1365216]
31	Sylte, M. J., & Suarez, D. L. Influenza neuraminidase as a vaccine antigen. <i>Vaccines for Pandemic Influenza</i> , pp. 227–241, 2009.	---	https://doi.org/10.1007/978-3-540-92165-3_12 [10.1007/978-3-540-92165-3_12]
32	Tomii, K., & Kanehisa, M. Analysis of amino acid indices and mutation matrices for sequence comparison and structure prediction of proteins. <i>Protein Engineering, Design and Selection</i> , vol. 9, no. 1, pp. 27–36, 1996.	---	https://doi.org/10.1093/protein/9.1.27 [10.1093/protein/9.1.27]

33	Tzarum N. et al. The 150-loop restricts the host specificity of human H10N8 influenza virus. Cell reports. vol. 19, no. 2, pp. 235–245, 2017.	---	https://doi.org/10.1016/j.celrep.2017.03.054 [10.1016/j.celrep.2017.03.054]
34	Wang, P., Zhu, W., Liao, B., Cai, L., Peng, L., & Yang, J. Predicting influenza antigenicity by matrix completion with antigen and antiserum similarity. Frontiers in microbiology, vol. 9, pp. 2500, 2018.	---	https://doi.org/10.3389/fmicb.2018.02500 [10.3389/fmicb.2018.02500]
35	Wikramaratna, P. S., Sandeman, M., Recker, M., & Gupta, S. The antigenic evolution of influenza: drift or thrift?. Philosophical Transactions of the Royal Society B: Biological Sciences, vol. 368, no. 1614, pp. 20120200, 2013.	---	https://doi.org/10.1098/rstb.2012.0200 [10.1098/rstb.2012.0200]

36	World Health Organization, Influenza fact sheet: Overview, Weekly Epidemiological Record= Relevé épidémiologique hebdomadaire, vol. 78, no. 11, pp. 77–80, 2003.	---	https://apps.who.int/iris/handle/10665/232113
37	Yang, H., Carney, P. J., Chang, J. C., Guo, Z., Villanueva, J. M., & Stevens, J. Structure and receptor binding preferences of recombinant human A (H3N2) virus hemagglutinins. Virology, vol. 477, pp. 18–31, 2015.	---	https://doi.org/10.1016/j.virol.2014.12.024 [10.1016/j.virol.2014.12.024]
38	Yang, X. Y., Shi, X. H., Meng, X., Li, X. L., Lin, K., Qian, Z. L., ... & Cai, Y. D. Classification of transcription factors using protein primary structure. Protein and peptide letters, vol. 17, no. 7, pp. 899–908, 2010.	---	https://doi.org/10.2174/092986610791306670 [10.2174/092986610791306670]

39	Yao, Y., Li, X., Liao, B., Huang, L., He, P., Wang, F., ... & Yang, J. Predicting influenza antigenicity from Hemagglutinin sequence data based on a joint random forest method. <i>Scientific reports</i> , vol. 7, no. 1, pp. 1–10, 2017.	---	https://doi.org/10.1038/s41598-017-01699-z [10.1038/s41598-017-01699-z]
40	Zhang, Y., Aebermann, B. D., Anderson, T. K., Burke, D. F., Dauphin, G., Gu, Z., ... & Scheuermann, R. H. Influenza Research Database: An integrated bioinformatics resource for influenza virus research. <i>Nucleic acids research</i> , vol. 45, no. D1, pp. D466–D474, 2017.	---	https://doi.org/10.1093/nar/gkw857 [10.1093/nar/gkw857]
41	Zhang, Z. H., Wang, Z. H., Zhang, Z. R., & Wang, Y. X. A novel method for apoptosis protein subcellular localization prediction combining encoding based on grouped weight and support vector machine. <i>FEBS letters</i> , vol. 580, no. 26, pp. 6169–6174, 2006.	---	https://doi.org/10.1016/j.febslet.2006.10.017 [10.1016/j.febslet.2006.10.017]

42	Zuo, Y. C., & Li, Q. Z. Using reduced amino acid composition to predict defensin family and subfamily: Integrating similarity measure and structural alphabet. <i>Peptides</i> , vol. 30, no. 10, pp. 1788–1793, 2009.	---	https://doi.org/10.1016/j.peptides.2009.06.032 [10.1016/j.peptides.2009.06.032]
----	--	-----	---