

## Распознавание рода коронавируса на основе прототипных штаммов

Чалей М.Б.<sup>\*1</sup>, Кутыркин В.А.<sup>\*\*2</sup>

<sup>1</sup>Институт математических проблем биологии – филиал ИПМ им. М.В. Келдыша РАН, Пущино, Московская область, Россия

<sup>2</sup>Московский государственный технический университет им. Н.Э. Баумана, Москва, Россия

**Аннотация.** Предложен вариантный подход к распознаванию рода коронавируса на основе распределения частот кодонов в ORF1ab и генах структурных белков (S, M и N). Метод основан на модифицированной статистике, ранее показавшей свою эффективность при распознавании видов флавивирусов. Вариантный подход использует для распознавания рода коронавирусов как различные комбинации нескольких генов коронавируса, так и отдельные гены. Род коронавируса окончательно определяется по результатам анализа всех рассматриваемых вариантов. Предлагаемый метод разработан с помощью обучающей выборки геномов прототипных штаммов коронавирусов родов *Alphacoronavirus*, *Betacoronavirus*, *Deltacoronavirus* и *Gammacoronavirus*. Использование вариантного подхода для распознавания рода коронавируса показало его высокую достоверность на уровне 95 %. Среди всех вариантов совместного анализа наибольшую надежность (98 %) показало использование распознавания рода коронавирусов на основе частот кодонов ORF1ab. Вариантный подход выявил явление мозаичности в геномах коронавирусов при распознавании рода, когда результаты распознавания рода по отдельным генам расходились с окончательным определением рода коронавируса. Такое явление, по-видимому, отражает гомологичные рекомбинации генов между коронавирусами различных видов и пластичность генома коронавирусов в эволюционных процессах.

**Ключевые слова:** геном коронавируса, ORF1ab, S-ген, M-ген, N-ген, статистический анализ, вариантный подход к распознаванию рода коронавируса.

### ВВЕДЕНИЕ

Коронавирусы (CoV, семейство *Coronaviridae*) представляют собой обширное семейство РНК-содержащих вирусов, инфицирующих млекопитающих, в том числе, человека, птиц и земноводных. Впервые коронавирус был выделен в 1931 г. как вирус, вызывающий инфекционный бронхит у цыплят (*Infectious bronchitis virus* – IBV), который в настоящее время носит название коронавируса птиц (*Avian coronavirus* – ACoV) [1]. Впоследствии оказалось, что коронавирусы вызывают заболевания различной тяжести у многих других видов птиц, как домашних, так и диких. Коронавирусы также вызывают гастроэнтериты и поражают респираторный тракт у домашних свиней, крупного рогатого скота, у кошек и собак, верблюдов, землероек, крыс и мышевидных грызунов, ежей, летучих мышей, китов, и др. [2]. В настоящее время известны семь коронавирусов, вызывающих заболевания у человека. Обычно их разделяют на две группы: HCoV 229E, HCoV NL63, HCoV HKU1 и HCoV OC43, при которых заболевание протекает в основном по типу ОРВИ; и особо опасные, такие как

\*maramaria@yandex.ru

\*\*vkutyркиn@yandex.ru

SARS-CoV-1 (*severe acute respiratory syndrome-related coronavirus*), MERS-CoV (*middle east respiratory syndrome coronavirus*) и SARS-CoV-2 (*severe acute respiratory syndrome 2 coronavirus*), вызвавший пандемию COVID-19 [3]. Более подробную информацию об этих коронавирусах и их природных источниках можно найти, например, в [4, 5].

Таксономическая классификация коронавирусов сложилась к 2011 г., когда в каталоге Международного комитета по таксономии вирусов (International Committee on Taxonomy of Viruses – ICTV) род *Coronavirus* перешел в статус подсемейства *Coronavirinae*, включающего четыре рода: *Alphacoronavirus*, *Betacoronavirus*, *Deltacoronavirus* и *Gammacoronavirus* [6]. Таксономическая структура коронавирусов была существенно пересмотрена в 2018 г. в результате введения понятие подрода и математической оценки некоторых рангов (подрод, род, семейство) по филогенетическому дереву, построенному методом наибольшего правдоподобия после множественного выравнивания полноразмерных геномов всех известных представителей семейства. Исходя из суммы различий по совокупности признаков между двумя узлами филогенетического дерева, измеренной вдоль соединяющих их ребер, были установлены соответствующие пороги: для подрода – 0.186, для рода – 0.789 и для подсемейства – 1.583. С расширением множества известных коронавирусов, эти значения могут постепенно меняться. Кроме того, обнаруженный в 2018 г. методом метагеномного анализа новый коронавирус [7], получивший название летовирус узкоротых квакш (*Microhyla fissipes*) 1-го типа был выделен ICTV в отдельное подсемейство *Letovirinae*, в то время как все ранее известные коронавирусы составили подсемейство *Orthocoronavirinae* [8].

В настоящей работе рассматриваются коронавирусы только подсемейства *Orthocoronavirinae*. Установлено, что это подсемейство связано с крылатыми животными, являющихся резервуаром коронавирусов: рукокрылые – для родов *Alphacoronavirus* ( $\alpha$ -CoV) и *Betacoronavirus* ( $\beta$ -CoV), птицы – для родов *Deltacoronavirus* ( $\delta$ -CoV) и *Gammacoronavirus* ( $\gamma$ -CoV) [4].

Свое название коронавирусы получили после того, как с помощью электронной микроскопии была выявлена их форма, напоминающая солнечную корону [9]. Такая форма вызвана тем, что вирион коронавируса обрамляют булавовидные шипы поверхностного гликопротеина S (spike protein, или S-белка). S-белок отвечает за проникновение вируса в клетку, связываясь с её определенными трансмембранными рецепторами. В липидную оболочку вириона закорены еще два структурных белка: каналообразующий оболочечный белок E (envelope protein) и трансмембранный белок M (membrane protein). У коронавируса есть нуклеокапсид, внутри которого находятся вирионная геномная (+)РНК (РНК позитивной полярности, на которой сразу может идти синтез белка на рибосомах инфицированной клетки), и связанный с нуклеокапсидом N-белок (nucleocapsid protein) [10, 11].

Геномная РНК вируса выступает в качестве мРНК для синтеза двух протяженных полипротеинов pp1a и pp1ab длиной около 4000 и 7000 аминокислотных остатков соответственно. Полипротеин pp1ab образуется в результате игнорирования рибосомой стоп-сигнала из-за шпильки, смещающей рамку считывания, как правило на –1 нуклеотид. Рамки считывания таких полипротеинов в геноме соответственно обозначаются как ORF1a и ORF1ab. Далее из полипротеинов pp1a и pp1ab высвобождаются две протеазы: главная протеаза Mpro (major protease) и папаин-подобная протеаза PLpro (papain-like protease), которые затем расщепляют сплошной полипротеин pp1a на 11 и полипротеин pp1ab на 16 отдельных неструктурных белков (nsp – non-structural protein). Эти неструктурные белки выполняют различные важные функции в жизненном цикле вируса [2, 11].

Например, nsp12 – РНК-зависимая РНК-полимераза RdRp (RNA-dependent RNA-polymerase) синтезирует (–)РНК, комплементарную геномной (+)РНК, которая, в свою очередь, выступает в качестве матрицы для синтеза геномных РНК для дочерних

вирионов. Помимо этого, RdRp на матрице вирионной геномной (+)РНК синтезирует серию субгеномных (-)РНК, далее используемых этой же полимеразой для синтеза субгеномных матричных (+)РНК, с которых считываются некоторые неструктурные и структурные белки, в частности, S-, M- и N-белок. Коронавирусы имеют самый длинный геном из всех геномов РНК-вирусов, который составляет порядка 30000 нукл., из которых около 2/3 занимает ORF1ab.

Выравнивание последовательности генома коронавирусов показывает гомологию 58 % в области, кодирующей неструктурные белки, 43 % – в области, кодирующей структурные белки и 54 % – на уровне всего генома, предполагая, что неструктурные белки более консервативны, а структурные белки более разнообразны и поддерживают адаптацию вируса к новым хозяевам [11].

Естественная изменчивость генома коронавирусов обеспечивается спонтанными мутациями в ходе его репликации и гомологичной рекомбинацией с другими геномами коронавирусов. Коронавирусы накапливают меньше мутаций, чем большинство РНК-вирусов, потому что они кодируют фермент nsp14, исправляющий ошибки репликации [12]. Наиболее вариабельным у коронавирусов является ген S-белка, отвечающий за их межвидовую передачу, вирулентность и контагиозность. Скорость его мутаций у разных видов коронавирусов различается. Например, у коронавирусов IBV, вызывающих инфекции респираторного тракта у птиц, скорость мутаций в S-гене на порядок выше ( $0.3-0.6 \times 10^{-2}$  замен в год/сайт [13]), чем у коронавирусов HCoV 229E, вызывающих относительно легко протекающие респираторные заболевания человека ( $3 \times 10^{-4}$  замен в год/сайт [14]). При такой скорости мутирования вид коронавируса HCoV NL63, провоцирующий развитие кашля у детей, по оценкам, разошелся со своим предком HCoV 229E в 11 веке [14]. Однако, следует понимать, что мутации охватывают гораздо больше генов, чем только ген белка S.

Наибольшую пластичность геному коронавируса создает возможность обширной гомологичной рекомбинации с геномами других коронавирусов, циркулирующих в одних и тех же хозяевах. По этой причине геном коронавирусов в природных резервуарах всегда является мозаичным. В основе такой рекомбинации лежит сложный жизненный цикл коронавирусов, включающий в себя образование набора субгеномных РНК. Точный механизм генетической рекомбинации у коронавирусов остается неясным: рекомбинационные сайты кажутся расположенными по геному случайным образом [15, 16].

В связи с продолжающейся пандемией COVID-19 полногеномное секвенирование изолятов вируса, хотя и рассматривается как наилучший способ генетического анализа по распространению новых патогенных штаммов мутирующего вируса SARS-CoV-2 [17], однако для отслеживания мутантных штаммов SARS-CoV-2 и других видов коронавирусов возможен более быстрый и экономичный подход таргетного секвенирования [18]. При таком подходе секвенированию подвергается только один геномный фрагмент коронавируса, например, фрагмент гена шиповидного S-белка, учитывая его способность к связыванию с клеточными рецепторами при инфицировании.

Вот уже два с лишним года не затухает пандемия, вызванная коронавирусом SARS-CoV-2, переживая все новые волны распространения мутантных штаммов. Как никогда остро встал вопрос оценки рисков возникновения возможных новых эпидемических источников, инициированных коронавирусами. Чем разнообразнее и точнее инструментарий идентификации и анализа коронавирусов, тем скорее и успешнее может быть выполнена такая оценка. Исходя из предыдущего опыта создания статистического метода для быстрого распознавания вида флавивируса (среди которых вирус желтой лихорадки, лихорадки денге, лихорадки Западного Нила и др.) по известной последовательности его генома [19], в настоящей работе предлагается надежный метод

распознавания рода коронавируса. Этот метод получил название вариантного подхода к распознаванию рода коронавируса. Вариантный подход использует для распознавания как комбинации нескольких структурных генов коронавируса, так и отдельные гены, секвенирование которых займет меньше времени, чем определение последовательности всего генома.

### Вариантный подход к распознаванию рода коронавируса

По результатам исследования структуры кодирования ORF1ab, S, M и N генов коронавирусов, было показано, что в последовательностях этих генов наблюдается скрытая триплетная профильная периодичность [20]. Такая периодичность характеризуется 3-профильной матрицей, в которой в каждом из тех столбцов стоит распределение букв четырехбуквенного алфавита ДНК. Для каждой структурной единицы из одного рода создавалась своя усредненная 3-профильная матрица. Было отмечено, что для всех родов коронавирусов ( $\alpha$ -CoV,  $\beta$ -CoV,  $\delta$ -CoV и  $\gamma$ -CoV) в гене мембранного белка M скрытая 3-профильная периодичность слабо выражена. Поэтому для этого гена во всех родах коронавирусов скрытые 3-профильные матрицы не вычислялись. Для каждой из трех остальных структурных единиц сравнение усредненных 3-профильных матриц между родами не выявило отличий, которые можно было бы использовать для создания методов или процедур идентификации рода коронавируса.

Ранее, при создании метода распознавания видов флавивирусов использовалась статистика, основанная на сравнении распределений кодонов в гене полипротеина флавивируса [19]. Если геном флавивируса обычно представлен непрерывной цепочкой кодирующих последовательностей вирусных генов, то в геноме коронавируса более сложная структурная организация генов. Кроме того, для создания метода идентификации из генома коронавируса были отобраны отдельные, особо значимые гены: ген ORF1ab; ген нуклеопротеина N – белка, формирующего комплекс с РНК-геномом вируса; ген мембранного M-белка; ген S-белка, образующего шипы на поверхности вирусного вириона. Эти гены не содержат интронов и из них можно сформировать псевдо-ген полипротеина, аналогичный гену полипротеина флавивируса. Такой псевдо-ген состоит из существенно отличающихся по длине структурных единиц (генов), которые вносят неравноценный вклад в общее распределение частот кодонов. Поэтому при создании статистики для распознавания рода коронавируса учитывалась длина каждой из структурных единиц псевдо-гена.

Целью настоящей работы было создание достаточно простого, быстрого и надежного метода распознавания рода коронавируса. При разработке метода мы опирались на опыт предыдущей работы по эффективному распознаванию вида флавивируса на основе распределения частот кодонов в кодирующей части генома, транскрибируемой единым полипротеином. По аналогии с геном полипротеина флавивируса можно рассмотреть псевдо-ген, составленный из отдельных генов структурных белков коронавируса. В качестве таких генов в работе были выбраны ген полипротеина pp1ab – ORF1ab, кодирующий 16 неструктурных белков, и три гена структурных белков: S-ген, M-ген и N-ген. Совместное распределение частот кодонов этих генов, учитывающее длину каждого в псевдогене, использовалось в первом варианте распознавания рода коронавируса. Для этого первого варианта  $V_1$  использовалось обозначение S+N+M+ORF1ab. Поскольку каждый из генов, входящих в вариант  $V_1$ , имеет свою степень пластичности и консервативности, для распознавания рода коронавируса дополнительно были исследованы еще пять других вариантов сочетания этих генов:  $V_2$ ) S+N+M,  $V_3$ ) S+N,  $V_4$ ) ORF1ab,  $V_5$ ) S и  $V_6$ ) N.

Результат распознавания коронавируса в вариантном подходе представляется в виде вариантной строки  $v = (v_1, v_2, v_3, v_4, v_5, v_6)$ , где  $v_i$  – результат распознавания рода коронавируса с помощью варианта  $V_i$  для  $i = \overline{1,6}$ , т.е.  $v_i \in \{\alpha, \beta, \delta, \gamma\} - CoV$ .

Однозначный результат вариантного подхода для анализируемого генома коронавируса определяется в случае четырех одинаковых компонент в вариантной строке. Такой подход позволяет сравнить надежность всех рассматриваемых вариантов распознавания и степень консервативности рассматриваемых генов. Мозаичность вариантной строки, т.е. наличие в ней неодинаковых компонент, может иметь различное биологическое объяснение, в частности, с помощью предположения о возможной гомологичной рекомбинации между генами коронавирусами различных видов [21–23].

### МАТЕРИАЛЫ

Для создания в каждом роде коронавирусов усредненного распределения кодонов использовались соответствующие гены прототипных штаммов рода [2]. Данные о кодах доступа GenBank для набора последовательностей геномов прототипных штаммов, из которых выделены анализируемые в работе гены ORF1ab, S, M и N, представлены в таблице 1.

В настоящее время род *Alphacoronavirus* подразделяется на 12 подродов (22 прототипных штамма), род *Betacoronavirus* – на пять подродов (28 прототипных штаммов), род *Deltacoronavirus* – на четыре подрода (10 прототипных штаммов) и род *Gammacoronavirus* – на два подрода (7 прототипных штаммов). В настоящей работе для удобства анализа и описания результатов вводится отдельная нумерация прототипных штаммов каждого рода (см. табл. 1).

**Таблица 1.** Коды доступа прототипных (п/т) штаммов коронавирусов в GenBank и их нумерация, принятая в работе

Род (Genus)	Подрод (Subgenus)	Вид (Species)	GenBank ID п/т штамма	№ п/т штамма рода
<i>Alphacoronavirus</i>	<i>Colacovirus</i>	Коронавирус летучих мышей CDPHE15 (Bat coronavirus CDPHE15)	NC_022103	1
	<i>Decacovirus</i>	Альфакоронавирус больших подковоносов HuB2013 ( <i>Rhinolophus ferrumequinum</i> alphacoronavirus HuB2013)	NC_028814	2
		Коронавирус летучих мышей HKU10 (Bat coronavirus HKU10)	NC_018871	3
	<i>Duvinacovirus</i>	Коронавирус человека 229E (Human coronavirus 229E)	NC_002645	4
	<i>Luchacovirus</i>	Коронавирус крыс Лунцюань Rn (Lucheng Rn rat coronavirus)	NC_032730	5
	<i>Minacovirus</i>	Коронавирус норки 1-го типа (Mink coronavirus 1)	NC_023760	6
		Коронавирус хорьков (Ferret coronavirus)	KX512809 KX512810	7 8
	<i>Minunacovirus</i>	Коронавирус длиннокрылов 1-го типа ( <i>Miniopterus bat coronavirus 1</i> )	EU420138	9
		Коронавирус длиннокрылов HKU8 ( <i>Miniopterus bat coronavirus HKU8</i> )	NC_010438	10
	<i>Miotacovirus</i>	Альфакоронавирус азиатских рыбоядных ночниц Sax-2011 ( <i>Myotis ricketti</i> alphacoronavirus Sax2011)	NC_28811	11

Род (Genus)	Подрод (Subgenus)	Вид (Species)	GenBank ID п/т штамма	№ п/т штамма рода		
	<i>Nyctacovirus</i>	Альфакоронавирус китайских вечерниц SC2013 ( <i>Nyctacus velutinus</i> alhacoronavirus SC2013)	NC_028833	12		
	<i>Pedacovirus</i>	Вирус эпизоотической диареи свиней (Porcine epidemic diarrhea virus)	KT323979	13		
		Коронавирус домашних гладконосов 512 ( <i>Scotophilus bat coronavirus</i> 512)	NC_009657	14		
	<i>Rhinacovirus</i>	Коронавирус подковоносов HKU2 ( <i>Rhinolopus bat coronavirus</i> HKU2)	NC_009988	15		
	<i>Setracovirus</i>	Коронавирус человека NL63 ( <i>Human coronavirus</i> NL63)	AY567487	16		
		NL63-подобный коронавирус BtKYNL63-9b (NL63-related bat coronavirus BtKYNL63-9b)	KY073745	17		
	<i>Tegacovirus</i>	Альфакоронавирус 1-го типа ( <i>Alphacoronavirus</i> 1)	KP981644	18		
			FJ938051	19		
			AY994055	20		
			KR270796	21		
			NC_038861	22		
<i>Betacoronavirus</i>	<i>Embecovirus</i>	Бетакоронавирус 1-го типа ( <i>Betacoronavirus</i> 1)	KF294357	1		
			BCU00735	2		
			KX432213	3		
			EF446615	4		
			AY391777	5		
			NC_017083	6		
			MF083115	7		
	<i>Embecovirus</i>	Коронавирус крыс Китая HKU24 ( <i>China Rattus coronavirus</i> HKU24)	NC_026011	8		
			Коронавирус мышей ( <i>Murine coronavirus</i> )	AC_000192	9	
				KF294371	10	
	<i>Embecovirus</i>	Коронавирус человека HKU1 ( <i>Human coronavirus</i> HKU1)	NC_012936	11		
			NC_006577	12		
			<i>Hibecovirus</i>	Бетакоронавирус листоносов Пратта Zhejiang2013 ( <i>Bat Hrbetacoronavirus</i> Zhejiang2013)	NC_025217	13
	<i>Merbecovirus</i>	Коронавирус Ближневосточного респираторного синдрома (Middle East respiratory syndrome-related coronavirus)			KF917527	14
					JX869059	15
			MG596803	16		
	<i>Merbecovirus</i>	Коронавирус ежей 1-го типа ( <i>Hedgehog coronavirus</i> 1)	МК679660	17		
			Коронавирус косолапых кожанов HKU4 ( <i>Tylonycteris bat coronavirus</i> HKU4)	NC_009019	18	
				Коронавирус нетопырей HKU5 ( <i>Pipustrellus bat coronavirus</i> HKU5)	NC_009020	19
	<i>Nobecovirus</i>	Коронавирус ночных крыланов GCCDC1 ( <i>Rousettus bat coronavirus</i> GCCDC1)	NC_030886	20		
			Коронавирус ночных крыланов HKU9 ( <i>Rousettus bat coronavirus</i> HKU9)	NC_009021	21	
	<i>Sarbecovirus</i>		MG772933	22		

Род (Genus)	Подрод (Subgenus)	Вид (Species)	GenBank ID п/т штамма	№ п/т штамма рода		
		Коронавирус китайских подковоносов ( <i>Rhinolophus sinicus coronavirus</i> )	MG772934	23		
		Коронавирус тяжелого острого респираторного синдрома (Severe acute respiratory syndrome-related coronavirus)	AY278489	24		
		Коронавирус тяжелого острого респираторного синдрома 2-го типа (Severe acute respiratory syndrome-related coronavirus 2)	NC_045512	26		
			MT121216	27		
			MN996532	28		
		<i>Deltacoronavirus</i>	<i>Andecovirus</i>	Коронавирус связей HKU20 (Wigeon coronavirus KCU20)	NC_016995	1
			<i>Buldecovirus</i>	Дельтакоронавирус свиней (Porcine deltacoronavirus)	JQ065042	2
KJ569769	3					
NC_016992	4					
Коронавирус белоглазок HKU16 (White eye coronavirus HKU16)	NC_016991			5		
Коронавирус бюльбюлей HKU11 (Bulbul coronavirus HKU11)	FJ376620			6		
Коронавирус муний HKU13 (Munia coronavirus HKU13)	NC_011550		7			
	NC_016993		8			
<i>Herdecovirus</i>	Коронавирус квакв HKU19 (Night heron coronavirus HKU19)		NC_016994	9		
<i>Moordecovirus</i>	Коронавирус камышниц HKU21 (Common morgen coronavirus HKU21)		NC_016996	10		
<i>Gammacoronavirus</i>	<i>Cegacovirus</i>	Коронавирус китообразных (Cetacean coronavirus)	EU111742	1		
			KF793826	2		
	<i>Igacovirus</i>	Коронавирус птиц (Avian coronavirus)	KF696629	3		
			GQ504724	4		
			NC_010800	5		
			AY641576	6		
			MK423877	7		

Для каждого из четырех генов ORF1ab, S, M и N последовательности генома прототипного штамма коронавируса рассматривалось распределение частот кодонов. Количество анализируемых геномов прототипных штаммов каждого рода составляло: 22, 28, 10 и 7 для родов  $\alpha$ -CoV,  $\beta$ -CoV,  $\delta$ -CoV и  $\gamma$ -CoV, соответственно (см. табл. 1).

Помимо геномов прототипных штаммов в работе использовались геномы коронавирусов четырех родов, полученные из базы данных нуклеотидных последовательностей GenBank выпуска 237 [24]. Количество всех анализируемых последовательностей геномов (вместе с прототипными штаммами) для родов  $\alpha$ -CoV,  $\beta$ -CoV,  $\delta$ -CoV и  $\gamma$ -CoV составляло: 924, 1954, 19 и 345, соответственно.

## МЕТОДЫ

Введем количественные характеристики, необходимые для предлагаемого в работе метода распознавания рода коронавируса, к которому относится анализируемая последовательность генома. Средняя длина по всем анализируемым в работе генам (их кодирующим районам – CDS) для каждого рода коронавирусов показана в таблице 2.

Таблица 2. Средняя длина CDS анализируемых генов коронавирусов в каждом роде

ген	$\alpha$ -CoV, нукл	$\beta$ -CoV, нукл	$\delta$ -CoV, нукл	$\gamma$ -CoV, нукл
ORF1ab	$L_{ORF}^{\alpha} = 20284$	$L_{ORF}^{\beta} = 21275$	$L_{ORF}^{\delta} = 18794$	$L_{ORF}^{\gamma} = 19861$
S	$L_S^{\alpha} = 4105$	$L_S^{\beta} = 3926$	$L_S^{\delta} = 3590$	$L_S^{\gamma} = 3591$
N	$L_N^{\alpha} = 1270$	$L_N^{\beta} = 1267$	$L_N^{\delta} = 1044$	$L_N^{\gamma} = 1229$
M	$L_M^{\alpha} = 692$	$L_M^{\beta} = 669$	$L_M^{\delta} = 679$	$L_M^{\gamma} = 678$
S+N+M+ORF1ab	$L^{\alpha} = 26351$	$L^{\beta} = 27137$	$L^{\delta} = 24107$	$L^{\gamma} = 25359$

В качестве обучающей выборки для распознавания рода использовались известные прототипные штаммы коронавирусов, коды доступа которых в GenBank исходно были даны в работе [2] и представлены в таблице 1. Для рода *Alphacoronavirus* количественные характеристики имеют вид:

$n = \overline{1,22}$  – номер прототипного штамма для рода  $\alpha$ -CoV;

$P_{ORF}^{\alpha}(n)$ ,  $P_S^{\alpha}(n)$ ,  $P_M^{\alpha}(n)$ ,  $P_N^{\alpha}(n)$  – распределение частот кодонов в CDS генов ORF1ab, S, M, N для прототипа с номером  $n$ ;

$l_{ORF}^{\alpha}(n)$ ,  $l_S^{\alpha}(n)$ ,  $l_M^{\alpha}(n)$ ,  $l_N^{\alpha}(n)$  – длины CDS генов ORF1ab, S, M, N в прототипе с номером  $n$ ;

$l^{\alpha}(n) = l_{ORF}^{\alpha}(n) + l_S^{\alpha}(n) + l_M^{\alpha}(n) + l_N^{\alpha}(n)$  – сумма длин соответствующих генов в прототипе с номером  $n$ .

В рамках указанных выше обозначений для усредненного взвешенного распределения частот кодонов  $P^{\alpha}$  рода  $\alpha$ -CoV используется формула:

$$P^{\alpha} = \frac{1}{22} \sum_{n=1}^{22} \left( \frac{l_{ORF}^{\alpha}(n)}{l^{\alpha}(n)} P_{ORF}^{\alpha}(n) + \frac{l_S^{\alpha}(n)}{l^{\alpha}(n)} P_S^{\alpha}(n) + \frac{l_M^{\alpha}(n)}{l^{\alpha}(n)} P_M^{\alpha}(n) + \frac{l_N^{\alpha}(n)}{l^{\alpha}(n)} P_N^{\alpha}(n) \right). \quad (1)$$

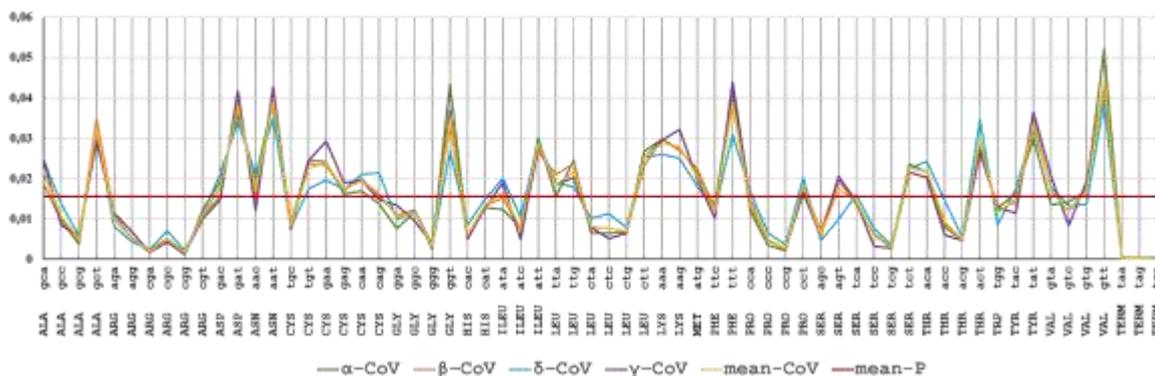


Рис. 1. Усредненные взвешенные распределения  $P^x$  ( $x \in \{\alpha, \beta, \delta, \gamma\}$ ) частот кодонов для первого варианта объединения генов S+N+M+ORF1ab из геномов прототипных штаммов коронавирусов для родов  $\alpha$ -CoV,  $\beta$ -CoV,  $\delta$ -CoV и  $\gamma$ -CoV (см. табл. 1). Среднее арифметическое распределение частот кодонов для всех родов показано желтой линией (mean-CoV). Красной линией (mean-P) показано среднее значение частот кодонов в таком среднем арифметическом распределении. Триплеты кодонов указаны в паре с названием кодируемой аминокислоты (нижняя линия подписей по горизонтальной оси).

Аналогично формуле (1) рассчитываются усредненные взвешенные распределения частот кодонов  $P^{\beta}$ ,  $P^{\delta}$  и  $P^{\gamma}$  для родов  $\beta$ -CoV,  $\delta$ -CoV и  $\gamma$ -CoV, соответственно. На

рисунке 1 представлены графики усредненных взвешенных распределений частот кодонов  $P$  для рассматриваемых родов коронавируса.

Результаты численных экспериментов показали, что распознавание рода коронавируса существенно улучшается, если из распределений кодонов исключить самые низкие частоты. Такими частотами оказались частоты всех трех стоп-кодонов (TERM) и двух (сga и сgg) из шести синонимичных кодонов аргинина (ARG), что учитывалось при вычислении усредненных распределений частот кодонов по формуле (1). В частности, для рода  $\alpha$ -CoV формат усредненного взвешенного распределения  $P^\alpha$  имеет вид:  $P^\alpha = (P_1^\alpha, P_2^\alpha, \dots, P_{59}^\alpha)$ .

Если рассматривается геном коронавируса неизвестного рода, то для его распределений частот кодонов (после исключения указанных выше кодонов) в генах ORF1ab, S, M и N используются обозначения:  $p_{\text{ORF}}$ ,  $p_S$ ,  $p_M$  и  $p_N$ . Следовательно, аналогично формуле (1), взвешенное распределение  $p$  частот кодонов в этом геноме вычисляется согласно формуле:

$$p = \frac{l_{\text{ORF}}}{l} p_{\text{ORF}} + \frac{l_S}{l} p_S + \frac{l_M}{l} p_M + \frac{l_N}{l} p_N, \quad (2)$$

где  $l_{\text{ORF}}$ ,  $l_S$ ,  $l_M$ ,  $l_N$  – длины соответствующих генов в рассматриваемом геноме коронавируса и  $l = l_{\text{ORF}} + l_S + l_M + l_N$ . Кроме того, формат распределения  $p$  имеет вид:  $p = (p_1, p_2, \dots, p_{59})$ .

Пусть  $x \in \{\alpha, \beta, \delta, \gamma\}$ , т.е.  $x$  – символ, обозначающий род коронавируса. Используя это обозначение, отклонение  $D(P^x, p)$  взвешенного распределения  $p$  частот кодонов в анализируемом геноме отдельного прототипного штамма или генома коронавируса неизвестного рода от усредненного взвешенного распределения  $P^x$  частот кодонов рода  $x$ -CoV вычисляется по формуле:

$$D(P^x, p) = \frac{1}{7} \sum_{i=1}^{59} \frac{|P_i^x - p_i|}{P_i^x}. \quad (3)$$

Формулы (1)–(3) введены для случая, когда совместно рассматривают четыре гена ORF1ab, S, M и N, для чего вводится обозначение:  $V_1$ ) S+N+M+ORF1ab. Аналогичные формулы также используются для вариантов объединения генов:  $V_2$ ) S+N+M,  $V_3$ ) S+N и для отдельных генов  $V_4$ ) ORF1ab,  $V_5$ ) S и  $V_6$ ) N. Отметим, что при анализе отдельного гена N из 64 кодонов исключаются пять кодонов, среди которых три кодона терминации (TERM) и два кодона (tgt и tgc) цистеина (CYS), имеющие наименьшую частоту встречаемости только для N-генов коронавируса. Среди отклонений  $D(P^\alpha, p)$ ,  $D(P^\beta, p)$ ,  $D(P^\delta, p)$  и  $D(P^\gamma, p)$  выбирается минимальное, которое предлагается как результат распознавания рода анализируемого генома коронавируса с помощью выбранного варианта объединения генов.

## РЕЗУЛЬТАТЫ И ОБСУЖДЕНИЕ

Вначале рассмотрим результаты применения предлагаемого вариантного подхода к распознаванию рода коронавируса на основе обучающей выборки из прототипных штаммов.

### Распознавание рода коронавируса для прототипных штаммов

Результаты применения статистики (3) к обучающим выборкам для различных вариантов объединения генов прототипных штаммов показаны в таблице 3а и таблице 3б.

**Таблица 3а.** Количественные результаты правильного распознавания рода прототипных (п/т) штаммов коронавирусов с помощью различных вариантов объединения генов

Варианты объединения генов		$\alpha$ -CoV Всего 22 п/т штамма	$\beta$ -CoV Всего 28 п/т штаммов	$\delta$ -CoV Всего 10 п/т штаммов	$\gamma$ -CoV Всего 7 п/т штаммов
V <sub>1</sub>	S+N+M+ORF1ab	20	25	9	7
V <sub>2</sub>	S+N+M	22	24	9	7
V <sub>3</sub>	S+N	22	25	9	7
V <sub>4</sub>	ORF1ab	20	25	9	7
V <sub>5</sub>	S	22	23	9	7
V <sub>6</sub>	N	19	28	10	7

Из таблицы 3а следует, что применение предложенного метода распознавания рода коронавируса по обучающей выборке показало результат в рамках допустимой статистической погрешности. В среднем, надежность распознавания по всем шести вариантам составила 93 %.

Конкретный вид ошибок в распознавании каждого из четырех родов коронавирусов показан в таблице 3б. Для каждого рода указывается номер прототипного штамма по таблице 1 и ошибочная принадлежность к другому роду отмечается прописной буквой  $\alpha$ ,  $\beta$ ,  $\delta$ ,  $\gamma$  в соответствии с названием рода *Alphacoronavirus* ( $\alpha$ -CoV), *Betacoronavirus* ( $\beta$ -CoV), *Deltacoronavirus* ( $\delta$ -CoV), *Gammacoronavirus* ( $\gamma$ -CoV). Так, например, пятый прототипный штамм рода *Alphacoronavirus* (см. табл. 1) распознан усредненным взвешенным распределением кодонов первого варианта объединения генов V<sub>1</sub>, как относящийся к роду *Betacoronavirus*, что в соответствующей строке обозначено 5 $\beta$ .

**Таблица 3б.** Результаты ошибочного распознавания рода у прототипных штаммов четырех родов коронавирусов ( $\alpha$ -CoV,  $\beta$ -CoV,  $\delta$ -CoV и  $\gamma$ -CoV)

Варианты объединения генов	$\alpha$ -CoV				$\beta$ -CoV				$\delta$ -CoV	$\gamma$ -CoV	
	V <sub>1</sub>	5 $\beta$			18 $\beta$	12 $\gamma$		19 $\delta$	20 $\delta$		10 $\gamma$
V <sub>2</sub>					12 $\alpha$	13 $\delta$		20 $\delta$	21 $\alpha$	10 $\beta$	-
V <sub>3</sub>					12 $\alpha$	13 $\delta$			21 $\alpha$	10 $\alpha$	-
V <sub>4</sub>	5 $\beta$			18 $\beta$	12 $\gamma$		19 $\delta$	20 $\delta$		10 $\beta$	-
V <sub>5</sub>					12 $\alpha$	13 $\delta$	19 $\delta$	20 $\alpha$	21 $\alpha$	10 $\alpha$	-
V <sub>6</sub>		11 $\beta$	12 $\beta$	17 $\beta$						-	-

### Распознавание рода среди последовательностей геномов коронавирусов в GenBank

В силу достаточно высокой надежности распознавания рода коронавирусов, полученной по обучающей выборке из прототипных штаммов, применим статистику (3) для распознавания рода последовательностей вирусных геномов из базы GenBank,

используя те же усредненные взвешенные распределения частот кодонов, полученные для прототипов каждого рода коронавирусов. Результаты такого распознавания рода коронавирусов для шести вариантов объединения генов показаны в таблице 4. В среднем, надежность распознавания по всем шести вариантам составила 95.6%, т.е. полученный результат находится в рамках стандартной статистической погрешности.

**Таблица 4.** Количественные результаты правильного распознавания рода последовательностей геномов коронавирусов (вместе с геномами рассматриваемых прототипных штаммов) из базы GenBank с помощью различных вариантов объединения генов

Варианты объединения генов		$\alpha$ -CoV Всего геномов: 924	$\beta$ -CoV Всего геномов: 1954	$\delta$ -CoV Всего геномов: 19	$\gamma$ -CoV Всего геномов: 345
V <sub>1</sub>	S+N+M+ORF1ab	911	1883	16	342
V <sub>2</sub>	S+N+M	903	1887	17	342
V <sub>3</sub>	S+N	831	1898	17	342
V <sub>4</sub>	ORF1ab	910	1906	16	343
V <sub>5</sub>	S	898	1898	17	192
V <sub>6</sub>	N	743	1924	19	343

### Вариантный подход к распознаванию рода коронавирусов

Рассмотрим вариантный подход к распознаванию рода коронавируса, основанного на анализе результатов распознавания рода с помощью нескольких вариантов объединения генов, которые перечислены в таблице 3а и таблице 4. В вариантном подходе результат распознавания рода для последовательности вирусного генома выглядит в виде вариантной строки  $\nu = (v_1, v_2, v_3, v_4, v_5, v_6)$ , где  $v_i \in \{\alpha, \beta, \delta, \gamma\}$  – есть результат распознавания рода (*Alphacoronavirus*, *Betacoronavirus*, *Deltacoronavirus*, *Gammacoronavirus*) коронавируса с использованием  $i$ -го ( $i = \overline{1, 6}$ ) варианта объединения генов.

В подавляющем большинстве случаев компоненты строки вариантного распознавания представлены одной буквой. Однако, встречаются случаи, когда в вариантной строке распознавания встречаются разные компоненты, то есть вариантный подход проявляет некоторую мозаичность в распознавании рода по геному коронавируса.

В таблице 5 представлены количественные результаты распознавания рода с помощью вариантного подхода для анализируемых последовательностей коронавирусных геномов из GenBank. Сначала для каждого рода указано число правильных распознаваний по всем шести вариантам объединения генов. Явление мозаичности представлено в следующих двух строках таблицы, где указаны числа правильных распознаваний рода по пяти и четырем компонентам вариантной строки распознавания.

Для вариантного подхода можно рассмотреть следующее правило распознавания рода коронавируса по анализируемой геномной последовательности. Если среди шести компонент вариантной строки не менее четырех указывают на один и тот же род коронавируса, то этот род и признается в качестве результата распознавания. Последняя строка в таблице 5 суммирует количественные результаты распознавания четырех родов коронавирусов для последовательностей геномов из GenBank с использованием такого правила вариантного подхода.

**Таблица 5.** Результаты распознавания рода для отдельных последовательностей геномов коронавируса из базы GenBank с помощью правила вариантного подхода

Число одинаковых компонент в вариантной строке, корректно указывающих род коронавируса	$\alpha$ -CoV Всего геномов 924	$\beta$ -CoV Всего геномов 1954	$\delta$ -CoV Всего геномов 19	$\gamma$ -CoV Всего геномов 345
6	687	1867	16	192
5	168	11	-	149
4	50	9	1	-
	распознано: 905	распознано: 1887	распознано: 17	распознано: 341

Как видно из таблицы 5, применение правила вариантного подхода гарантирует надежность распознавания рода коронавируса в 97 %. Следует отметить, что, согласно таблице 4, распознавание рода коронавируса только на основе гена ORF1ab обеспечивает надежность в 97.9 %. Однако, вариантный подход может способствовать выявлению структурных особенностей в генах коронавирусов.

### Явление мозаичности генов в геномах коронавирусов при использовании вариантного подхода

В вариантном подходе при распознавании рода коронавируса отдельного генома число одинаковых компонент в вариантной строке может варьироваться от двух до шести. Однако, в случаях, когда не все компоненты одинаковы, пять одинаковых компонент встречаются наиболее часто. В качестве примера таких случаев рассмотрим результаты применения вариантного подхода к распознаванию рода *Gammacoronavirus* в специальной выборке 149 геномов коронавирусов этого рода из базы GenBank (см. табл. 6). Полная выборка анализируемых в работе геномов коронавирусов рода *Gammacoronavirus* из базы GenBank составляет 345 последовательностей (см. табл. 5).

Таблица 6 представляет варианты строки, полученные при распознавании рода коронавируса для 149 последовательностей геномов из GenBank в результате применения вариантного подхода. Из таблицы 6 следует, что среди всех 149 геномов варианты строки имеют пять одинаковых компонент (с первой по четвертую и шестую), правильно определяющих род *Gammacoronavirus*. Ошибка в определении рода во всех этих вариантных строках находится в пятой компоненте, соответствующей S-гену спайк-белка. При этом пятая компонента указывает только на род *Alphacoronavirus* (15.4 %) или род *Betacoronavirus* (84.6 %). Аналогичное явление мозаичности можно наблюдать при распознавании рода *Alphacoronavirus* (см. таблицу 7).

Пример явления мозаичности в геномах коронавирусов из рода *Alphacoronavirus*, связанный, главным образом, с N-геном белка нуклеокапсида, представлен в таблице 7. Несмотря на то, что для родов *Betacoronavirus*, *Deltacoronavirus* и *Gammacoronavirus* достоверность распознавания на основе только одного N-гена белка нуклеокапсида имеет, практически, 100%-ю надежность (см. табл. 4), при распознавании рода *Alphacoronavirus* выявляется 131 вариантная строка, где шестая компонента, соответствующая N-гену, ошибочно указывает на род *Betacoronavirus*. Кроме того, в 39 вариантных строках содержится такое же указание на *Betacoronavirus* только в двух компонентах: в шестой и в третьей, отвечающей за вариант объединения S+N (см. табл. 7). Таким образом, ген нуклеокапсида не позволяет однозначно идентифицировать роды  $\alpha$ -CoV и  $\beta$ -CoV.

**Таблица 6.** Явление мозаичности в распознавании коронавирусных геномов рода *Gammacoronavirus* с помощью вариантного подхода (см. табл. 3а, табл. 3б)

GenBank ID		$v_1 v_2 v_3 v_4 v_5 v_6$	GenBank ID		$v_1 v_2 v_3 v_4 v_5 v_6$	GenBank ID		$v_1 v_2 v_3 v_4 v_5 v_6$
1	MK071267	γγγγβγ	51	KX348115	γγγγβγ	101	KP118892	γγγγβγ
2	KX258195	γγγγβγ	52	KP118882	γγγγβγ	102	KP118888	γγγγβγ
3	KY933089	γγγγβγ	53	KP118881	γγγγβγ	103	KX252787	γγγγαγ
4	KY626044	γγγγβγ	54	KX364294	γγγγαγ	104	KX236001	γγγγβγ
5	KY626045	γγγγαγ	55	JX195175	γγγγβγ	105	EU637854	γγγγβγ
6	AY338732	γγγγβγ	56	KJ425486	γγγγαγ	106	KX364291	γγγγβγ
7	AY319651	γγγγβγ	57	KX236008	γγγγβγ	107	KP118893	γγγγβγ
8	AY646283	γγγγβγ	58	KP790143	γγγγβγ	108	KJ425510	γγγγβγ
9	LN610099	γγγγβγ	59	KX236015	γγγγβγ	109	KP118884	γγγγβγ
10	KX266757	γγγγβγ	60	KX077987	γγγγαγ	110	KP118885	γγγγβγ
11	MH779860	γγγγβγ	61	JX195178	γγγγαγ	111	KX236003	γγγγβγ
12	MH779856	γγγγβγ	62	JX195177	γγγγαγ	112	KX236002	γγγγαγ
13	MH779857	γγγγβγ	63	KX302874	γγγγαγ	113	KX236013	γγγγβγ
14	MH779858	γγγγβγ	64	KX640829	γγγγβγ	114	KX219798	γγγγαγ
15	MH779859	γγγγβγ	65	KP343691	γγγγβγ	115	KX372250	γγγγβγ
16	KU556805	γγγγβγ	66	KX400753	γγγγβγ	116	KX219799	γγγγβγ
17	KU556806	γγγγβγ	67	MK574042	γγγγβγ	117	KX219796	γγγγβγ
18	KU556807	γγγγβγ	68	MK574043	γγγγβγ	118	JQ088078	γγγγβγ
19	KU556804	γγγγβγ	69	KX219793	γγγγβγ	119	KY776701	γγγγβγ
20	KF663559	γγγγβγ	70	KX247127	γγγγαγ	120	KY407557	γγγγβγ
21	KP662631	γγγγβγ	71	KX236014	γγγγαγ	121	MH924835	γγγγαγ
22	KT736032	γγγγβγ	72	KP036503	γγγγβγ	122	KY776700	γγγγβγ
23	MN566147	γγγγβγ	73	KP118883	γγγγβγ	123	MK032180	γγγγβγ
24	MN599049	γγγγβγ	74	MG448607	γγγγβγ	124	MK581203	γγγγβγ
25	MK878536	γγγγβγ	75	KX247128	γγγγβγ	125	MK581206	γγγγβγ
26	MN512434	γγγγβγ	76	KP118890	γγγγβγ	126	MK581205	γγγγβγ
27	MN512435	γγγγβγ	77	KX252774	γγγγβγ	127	MK581200	γγγγαγ
28	MN512436	γγγγβγ	78	KP036504	γγγγβγ	128	MK581201	γγγγαγ
29	MN512437	γγγγβγ	79	KP118880	γγγγβγ	129	MK581202	γγγγβγ
30	MN512438	γγγγβγ	80	KX252776	γγγγβγ	130	MK581207	γγγγβγ
31	KY805846	γγγγαγ	81	KP118887	γγγγβγ	131	MK581208	γγγγβγ
32	MG233398	γγγγαγ	82	KX252773	γγγγβγ	132	GQ504722	γγγγβγ
33	KY273667	γγγγαγ	83	JF274479	γγγγβγ	133	GQ504723	γγγγβγ
34	KR608272	γγγγβγ	84	KP868572	γγγγβγ	134	FJ888351	γγγγαγ
35	KY588135	γγγγαγ	85	KP118891	γγγγβγ	135	MK217372	γγγγβγ
36	KT852992	γγγγβγ	86	KX434788	γγγγβγ	136	MK937829	γγγγαγ
37	KF931628	γγγγβγ	87	KX302870	γγγγβγ	137	MK217373	γγγγβγ
38	KF460437	γγγγβγ	88	KP790144	γγγγβγ	138	MK217374	γγγγβγ
39	GU393331	γγγγβγ	89	KP790146	γγγγβγ	139	MK217375	γγγγβγ
40	AY514485	γγγγβγ	90	KP036502	γγγγβγ	140	KX185056	γγγγβγ
41	GU393332	γγγγβγ	91	KP790145	γγγγβγ	141	KX185059	γγγγβγ
42	GU393334	γγγγαγ	92	KJ425508	γγγγαγ	142	KU900739	γγγγβγ
43	GU393338	γγγγβγ	93	KX302866	γγγγβγ	143	KU900740	γγγγβγ
44	KF377577	γγγγβγ	94	KJ425509	γγγγβγ	144	KU900743	γγγγβγ
45	GQ504720	γγγγβγ	95	KP036505	γγγγβγ	145	MN128086	γγγγβγ
46	GQ504721	γγγγβγ	96	KX364297	γγγγβγ	146	MN128088	γγγγβγ
47	EU418975	γγγγβγ	97	MK937833	γγγγαγ	147	MN128087	γγγγβγ
48	EU418976	γγγγβγ	98	KX434790	γγγγβγ	148	KC008600	γγγγαγ
49	FJ904714	γγγγβγ	99	KX302861	γγγγβγ	149	GQ427174	γγγγβγ
50	FJ904715	γγγγβγ	100	KX302865	γγγγβγ			

Таблица 7. Явление мозаичности в распознавании коронавируса геномов рода *Alphacoronavirus* с помощью вариантного подхода (см. табл. 3а, табл. 3б)

GenBank ID		$v_1 v_2 v_3 v_4 v_5 v_6$	GenBank ID		$v_1 v_2 v_3 v_4 v_5 v_6$	GenBank ID		$v_1 v_2 v_3 v_4 v_5 v_6$
1	MN938448	$\alpha \alpha \alpha \alpha \alpha \beta$	58	LT900499	$\alpha \alpha \alpha \alpha \alpha \beta$	115	KM609213	$\alpha \alpha \alpha \alpha \alpha \beta$
2	MN938450	$\alpha \alpha \alpha \alpha \alpha \beta$	59	LT900500	$\alpha \alpha \alpha \alpha \alpha \beta$	116	MH593900	$\alpha \alpha \alpha \alpha \alpha \beta$
3	MN938449	$\alpha \alpha \alpha \alpha \alpha \beta$	60	LT900501	$\alpha \alpha \alpha \alpha \alpha \beta$	117	KJ196348	$\alpha \alpha \alpha \alpha \alpha \beta$
4	MK472068	$\alpha \alpha \alpha \alpha \alpha \beta$	61	LT900502	$\alpha \alpha \alpha \alpha \alpha \beta$	118	KU297956	$\alpha \alpha \alpha \alpha \alpha \beta$
5	MK472069	$\alpha \alpha \alpha \alpha \alpha \beta$	62	KR003452	$\alpha \alpha \alpha \alpha \alpha \beta$	119	KP403954	$\alpha \alpha \alpha \alpha \alpha \beta$
6	MK472070	$\alpha \alpha \alpha \alpha \alpha \beta$	63	JN825712	$\alpha \alpha \alpha \alpha \alpha \beta$	120	KJ645682	$\alpha \alpha \alpha \alpha \alpha \beta$
7	MK472071	$\alpha \alpha \alpha \alpha \alpha \beta$	64	KX016034	$\alpha \alpha \alpha \alpha \alpha \beta$	121	MK138516	$\alpha \alpha \alpha \alpha \alpha \beta$
8	MH687935	$\alpha \alpha \alpha \alpha \alpha \beta$	65	KR809885	$\alpha \alpha \alpha \alpha \alpha \beta$	122	KX883635	$\alpha \alpha \alpha \alpha \alpha \beta$
9	MH687936	$\alpha \alpha \alpha \alpha \alpha \beta$	66	KT199103	$\alpha \alpha \alpha \alpha \alpha \beta$	123	KR818832	$\alpha \alpha \alpha \alpha \alpha \beta$
10	MH687939	$\alpha \alpha \alpha \alpha \alpha \beta$	67	MF346935	$\alpha \alpha \alpha \alpha \alpha \beta$	124	MK250953	$\alpha \alpha \alpha \alpha \alpha \beta$
11	MH687940	$\alpha \alpha \alpha \alpha \alpha \beta$	68	MK606369	$\alpha \alpha \alpha \alpha \alpha \beta$	125	MK720945	$\alpha \alpha \alpha \alpha \alpha \beta$
12	MH687949	$\alpha \alpha \alpha \alpha \alpha \beta$	69	MG837012	$\alpha \alpha \alpha \alpha \alpha \beta$	126	MK720946	$\alpha \alpha \alpha \alpha \alpha \beta$
13	MH687950	$\alpha \alpha \alpha \alpha \alpha \beta$	70	MN114121	$\alpha \alpha \alpha \alpha \alpha \beta$	127	MN611517	$\alpha \alpha \alpha \alpha \alpha \beta$
14	MH687951	$\alpha \alpha \alpha \alpha \alpha \beta$	71	MK644601	$\alpha \alpha \alpha \alpha \alpha \beta$	128	MN611521	$\alpha \alpha \alpha \alpha \alpha \beta$
15	MH687959	$\alpha \alpha \alpha \alpha \alpha \beta$	72	MH726408	$\alpha \alpha \alpha \alpha \alpha \beta$	129	MK720944	$\alpha \alpha \alpha \alpha \alpha \beta$
16	MH687962	$\alpha \alpha \alpha \alpha \alpha \beta$	73	MH726374	$\alpha \alpha \alpha \alpha \alpha \beta$	130	NC_028811	$\alpha \alpha \alpha \alpha \alpha \beta$
17	KJ473799	$\alpha \alpha \alpha \alpha \alpha \beta$	74	MH726393	$\alpha \alpha \alpha \alpha \alpha \beta$	131	NC_028833	$\alpha \alpha \alpha \alpha \alpha \beta$
18	KJ473800	$\alpha \alpha \alpha \alpha \alpha \beta$	75	MH726368	$\alpha \alpha \alpha \alpha \alpha \beta$	132	LT906582	$\alpha \alpha \beta \alpha \alpha \beta$
19	KJ473798	$\alpha \alpha \alpha \alpha \alpha \beta$	76	MH726372	$\alpha \alpha \alpha \alpha \alpha \beta$	133	KX981440	$\alpha \alpha \beta \alpha \alpha \beta$
20	KJ473806	$\alpha \alpha \alpha \alpha \alpha \beta$	77	MH726369	$\alpha \alpha \alpha \alpha \alpha \beta$	134	MF375374	$\alpha \alpha \beta \alpha \alpha \beta$
21	KJ473809	$\alpha \alpha \alpha \alpha \alpha \beta$	78	MH726394	$\alpha \alpha \alpha \alpha \alpha \beta$	135	KM887144	$\alpha \alpha \beta \alpha \alpha \beta$
22	MK211373	$\alpha \alpha \alpha \alpha \alpha \beta$	79	MH726395	$\alpha \alpha \alpha \alpha \alpha \beta$	136	LT905450	$\alpha \alpha \beta \alpha \alpha \beta$
23	KY799179	$\alpha \alpha \alpha \alpha \alpha \beta$	80	MH726381	$\alpha \alpha \alpha \alpha \alpha \beta$	137	LT905451	$\alpha \alpha \beta \alpha \alpha \beta$
24	KY073746	$\alpha \alpha \alpha \alpha \alpha \beta$	81	MH726382	$\alpha \alpha \alpha \alpha \alpha \beta$	138	LT906620	$\alpha \alpha \beta \alpha \alpha \beta$
25	KY073745	$\alpha \alpha \alpha \alpha \alpha \beta$	82	MH726405	$\alpha \alpha \alpha \alpha \alpha \beta$	139	MH726366	$\alpha \alpha \beta \alpha \alpha \beta$
26	LM645058	$\alpha \alpha \alpha \alpha \alpha \beta$	83	MH726383	$\alpha \alpha \alpha \alpha \alpha \beta$	140	MH726367	$\alpha \alpha \beta \alpha \alpha \beta$
27	LT898408	$\alpha \alpha \alpha \alpha \alpha \beta$	84	MK690502	$\alpha \alpha \alpha \alpha \alpha \beta$	141	MH726380	$\alpha \alpha \beta \alpha \alpha \beta$
28	LT898409	$\alpha \alpha \alpha \alpha \alpha \beta$	85	KX534206	$\alpha \alpha \alpha \alpha \alpha \beta$	142	KY070587	$\alpha \alpha \beta \alpha \alpha \beta$
29	LT898410	$\alpha \alpha \alpha \alpha \alpha \beta$	86	MK482396	$\alpha \alpha \alpha \alpha \alpha \beta$	143	MK841494	$\alpha \alpha \beta \alpha \alpha \beta$
30	LT898411	$\alpha \alpha \alpha \alpha \alpha \beta$	87	MK482397	$\alpha \alpha \alpha \alpha \alpha \beta$	144	LT897799	$\alpha \alpha \beta \alpha \alpha \beta$
31	LT898412	$\alpha \alpha \alpha \alpha \alpha \beta$	88	LT898447	$\alpha \alpha \alpha \alpha \alpha \beta$	145	KY007140	$\alpha \alpha \beta \alpha \alpha \beta$
32	LT898413	$\alpha \alpha \alpha \alpha \alpha \beta$	89	KX812523	$\alpha \alpha \alpha \alpha \alpha \beta$	146	GU937797	$\alpha \alpha \beta \alpha \alpha \beta$
33	LT898414	$\alpha \alpha \alpha \alpha \alpha \beta$	90	KX812524	$\alpha \alpha \alpha \alpha \alpha \beta$	147	GU937797	$\alpha \alpha \beta \alpha \alpha \beta$
34	LT898415	$\alpha \alpha \alpha \alpha \alpha \beta$	91	KU558701	$\alpha \alpha \alpha \alpha \alpha \beta$	148	MK702008	$\alpha \alpha \beta \alpha \alpha \beta$
35	LT898416	$\alpha \alpha \alpha \alpha \alpha \beta$	92	LM645057	$\alpha \alpha \alpha \alpha \alpha \beta$	149	MK409659	$\alpha \alpha \beta \alpha \alpha \beta$
36	LT898417	$\alpha \alpha \alpha \alpha \alpha \beta$	93	MF807951	$\alpha \alpha \alpha \alpha \alpha \beta$	150	MK409657	$\alpha \alpha \beta \alpha \alpha \beta$
37	LT898418	$\alpha \alpha \alpha \alpha \alpha \beta$	94	KR153325	$\alpha \alpha \alpha \alpha \alpha \beta$	151	MK409658	$\alpha \alpha \beta \alpha \alpha \beta$
38	LT898421	$\alpha \alpha \alpha \alpha \alpha \beta$	95	KR153326	$\alpha \alpha \alpha \alpha \alpha \beta$	152	KU664503	$\alpha \alpha \beta \alpha \alpha \beta$
39	LT898423	$\alpha \alpha \alpha \alpha \alpha \beta$	96	KY793536	$\alpha \alpha \alpha \alpha \alpha \beta$	153	KX839246	$\alpha \alpha \beta \alpha \alpha \beta$
40	LT898425	$\alpha \alpha \alpha \alpha \alpha \beta$	97	KR095279	$\alpha \alpha \alpha \alpha \alpha \beta$	154	KY486713	$\alpha \alpha \beta \alpha \alpha \beta$
41	LT898426	$\alpha \alpha \alpha \alpha \alpha \beta$	98	KF760557	$\alpha \alpha \alpha \alpha \alpha \beta$	155	KY486714	$\alpha \alpha \beta \alpha \alpha \beta$
42	LT898427	$\alpha \alpha \alpha \alpha \alpha \beta$	99	MH061339	$\alpha \alpha \alpha \alpha \alpha \beta$	156	KX839247	$\alpha \alpha \beta \alpha \alpha \beta$
43	LT898430	$\alpha \alpha \alpha \alpha \alpha \beta$	100	MH061336	$\alpha \alpha \alpha \alpha \alpha \beta$	157	KX839248	$\alpha \alpha \beta \alpha \alpha \beta$
44	LT898431	$\alpha \alpha \alpha \alpha \alpha \beta$	101	MH061343	$\alpha \alpha \alpha \alpha \alpha \beta$	158	KX839249	$\alpha \alpha \beta \alpha \alpha \beta$
45	LT898432	$\alpha \alpha \alpha \alpha \alpha \beta$	102	MH061337	$\alpha \alpha \alpha \alpha \alpha \beta$	159	KX839250	$\alpha \alpha \beta \alpha \alpha \beta$
46	LT898433	$\alpha \alpha \alpha \alpha \alpha \beta$	103	MH061338	$\alpha \alpha \alpha \alpha \alpha \beta$	160	KX839251	$\alpha \alpha \beta \alpha \alpha \beta$
47	LT898435	$\alpha \alpha \alpha \alpha \alpha \beta$	104	KC196276	$\alpha \alpha \alpha \alpha \alpha \beta$	161	KP890336	$\alpha \alpha \beta \alpha \alpha \beta$
48	LT898436	$\alpha \alpha \alpha \alpha \alpha \beta$	105	KR011756	$\alpha \alpha \alpha \alpha \alpha \beta$	162	MH061341	$\alpha \alpha \beta \alpha \alpha \beta$
49	LT898438	$\alpha \alpha \alpha \alpha \alpha \beta$	106	MK138353	$\alpha \alpha \alpha \alpha \alpha \beta$	163	MH061340	$\alpha \alpha \beta \alpha \alpha \beta$
50	LT898439	$\alpha \alpha \alpha \alpha \alpha \beta$	107	JX112709	$\alpha \alpha \alpha \alpha \alpha \beta$	164	AF353511	$\alpha \alpha \beta \alpha \alpha \beta$
51	LT898440	$\alpha \alpha \alpha \alpha \alpha \beta$	108	MH708243	$\alpha \alpha \alpha \alpha \alpha \beta$	165	EF185992	$\alpha \alpha \beta \alpha \alpha \beta$
52	LT898441	$\alpha \alpha \alpha \alpha \alpha \beta$	109	MK862249	$\alpha \alpha \alpha \alpha \alpha \beta$	166	MF782687	$\alpha \alpha \beta \alpha \alpha \beta$
53	LT898443	$\alpha \alpha \alpha \alpha \alpha \beta$	110	KT941120	$\alpha \alpha \alpha \alpha \alpha \beta$	167	KM609208	$\alpha \alpha \beta \alpha \alpha \beta$
54	LT898444	$\alpha \alpha \alpha \alpha \alpha \beta$	111	MF782686	$\alpha \alpha \alpha \alpha \alpha \beta$	168	KU252649	$\alpha \alpha \beta \alpha \alpha \beta$
55	LT898445	$\alpha \alpha \alpha \alpha \alpha \beta$	112	KY111278	$\alpha \alpha \alpha \alpha \alpha \beta$	169	KX550281	$\alpha \alpha \beta \alpha \alpha \beta$

GenBank ID	$v_1 v_2 v_3 v_4 v_5 v_6$	GenBank ID	$v_1 v_2 v_3 v_4 v_5 v_6$	GenBank ID	$v_1 v_2 v_3 v_4 v_5 v_6$
56	LT898446	$\alpha \alpha \alpha \alpha \alpha \beta$	113	KM609211	$\alpha \alpha \alpha \alpha \alpha \beta$
57	LT900498	$\alpha \alpha \alpha \alpha \alpha \beta$	114	KM609212	$\alpha \alpha \alpha \alpha \alpha \beta$

Биологическую значимость наблюдаемого явления мозаичности можно будет оценить только при проведении дальнейших исследований. По-видимому, это явление определяется высокой генетической изменчивостью коронавирусов и их способностью к рекомбинации в процессе межпопуляционных взаимоотношений вирусов и их природных резервуаров.

## ЗАКЛЮЧЕНИЕ

В работе предложены методы распознавания рода коронавируса на основе модифицированной статистики, ранее показавшей свою эффективность при распознавании видов флавивирусов. В основе этой статистики лежит использование распределений кодонов в ORF1ab, кодирующей неструктурные белки, и в генах структурных белков (S, M и N) коронавируса. Предлагаемые методы разработаны с помощью обучающей выборки геномов прототипных штаммов.

В отличие от однозначной трансляции единого полипротеина в геноме флавивируса, механизм трансляции неструктурных и структурных белков коронавирусов многостадийный, поэтому учитывались как длина генов, так и различные варианты их совместного анализа. Рассматривались шесть вариантов совместного анализа генов:  $V_1$ ) S+N+M+ORF1ab,  $V_2$ ) S+N+M,  $V_3$ ) S+N,  $V_4$ ) ORF1ab,  $V_5$ ) S и  $V_6$ ) N.

Использование предложенных статистик для распознавания рода коронавируса показало их высокую достоверность на уровне 95 %. Среди всех вариантов совместного анализа наибольшую надежность (98 %) показало использование распознавания рода коронавирусов на основе частот кодонов ORF1ab.

В работе был рассмотрен вариантный подход, в котором результат распознавания рода коронавирусов представлялся в виде строки из шести компонент. В этой строке каждая компонента представляла результат распознавания рода анализируемого генома коронавируса на основе одного из шести вариантов совместного анализа четырех генов коронавируса. При таком подходе однозначный результат распознавания признавался при наличии не менее четырех одинаковых компонент в вариантной строке. Однако надежность такого подхода оказалась немного ниже (97 %) надежности только на основе гена ORF1ab. С другой стороны, вариантный подход выявил явление мозаичности при распознавании рода коронавируса. Так, например, среди распознанных 341 геномов (из всех 345) рода *Gammacoronavirus* в пятой компоненте вариантной строки, отвечающей за S-ген, «распознавался» род *Alphacoronavirus* (22) или *Betacoronavirus* (127). Также среди 905 распознанных геномов (из всех 924) рода *Alphacoronavirus* у 170 геномов был «распознан» род *Betacoronavirus* в шестой компоненте, соответствующей N-гену белка нуклеокапсида. Можно предположить, что явление мозаичности отражает пластичность генома коронавирусов и, своего рода, степень готовности к освоению не только новых ареалов существования, но и к поиску новых хозяев.

## СПИСОК ЛИТЕРАТУРЫ

1. Schalk A.F., Hawn M.C. An apparently new respiratory disease of baby chicks. *J. Am. Vet. Med. Assoc.* 1931. V. 78. P. 413–423.
2. Щелканов М.Ю., Попова А.Ю., Дедков В.Г., Акимкин В.Г., Малеев В.В. История изучения и современная классификация коронавирусов (*Nidovirales*:

- Coronaviridae*). *Инфекция и иммунитет*. 2020. Т. 10. № 2. С. 221–246. doi: [10.15789/2220-7619-HOI-1412](https://doi.org/10.15789/2220-7619-HOI-1412)
3. Хайтович А.Б. Коронавирусы (таксономия, структура вируса). *Крымский журнал экспериментальной и клинической медицины*. 2020. Т. 10. № 3. С. 69–81. doi: [10.37279/2224-6444-2020-10-3-69-81](https://doi.org/10.37279/2224-6444-2020-10-3-69-81)
  4. Львов Д.К., Альховский С.В. Истоки пандемии COVID-19: экология и генетика коронавирусов (*Betacoronavirus: Coronaviridae*) SARS-CoV, SARS-CoV-2 (подрод *Sarbecovirus*), MERS-CoV (подрод *Merbecovirus*). *Вопросы вирусологии*. 2020. Т. 65. № 2. С. 62–70. doi: [10.36233/0507-4088-2020-65-2-62-70](https://doi.org/10.36233/0507-4088-2020-65-2-62-70)
  5. Львов Д.К., Альховский С.В., Колобухина Л.В., Бурцева Е.И. Этиология эпидемической вспышки COVID-19 в г. Ухань (провинция Хубэй, Китайская Народная Республика), ассоциированной с вирусом 2019-nCoV (*Nidovirales, Coronaviridae, Coronavirinae, Betacoronavirus*, подрод *Sarbecovirus*): уроки эпидемии SARS-CoV. *Вопросы вирусологии*. 2020. Т. 65. № 1. С. 6–15. doi: [10.36233/0507-4088-2020-65-1-6-15](https://doi.org/10.36233/0507-4088-2020-65-1-6-15)
  6. *Virus Taxonomy. Classification and Nomenclature of Viruses. Ninth Report of the International Committee on Taxonomy of Viruses*. Eds. King A.M.Q., Adams M.J., Carstens E.B., Lefkowitz E.J. Elsevier Academic Press, 2011. 1338 p.
  7. Bukhari K., Mulley G., Gulyaeva A.A., Zhao L., Shu G., Jiang J., Neuman B.W. Description and initial characterization of metatranscriptomic nidovirus-like genomes from the proposed new family *Abyssoviridae*, and from a sister group to the *Coronavirinae*, the proposed genus *Alphaletovirus*. *Virology*. 2018. V. 524. P. 160–171. doi: 10.1016/j.virol.2018.08.010
  8. Ziebuhr J., Baric R.S., Baker S., de Groot R.J., Drosten C., Gulyaeva A., Haagmans B.L., Neuman B.W., Perlman S., Poon L.L.M., Sola I., Gorbalenya A.E. *Reorganization of the family Coronaviridae into two families, Coronaviridae (including the current subfamily Coronavirinae and the new subfamily Letovirinae) and the new family Tobaniviridae (accommodating the current subfamily Torovirinae and three other subfamilies), revision of the genus rank structure and introduction of a new subgenus rank: Proposal 2017.013S (08.08.2018) for International Committee on Taxonomy of Viruses*.
  9. Neuman B.W., Adair B.D., Yoshioka C., Quispe J.D., Kuhn G.O.P., Milligan R.A., Yeager M., Buchmeier M.J. Supramolecular architecture of severe acute respiratory syndrome coronavirus revealed by electron cryomicroscopy. *J. Virol.* 2006. V. 80. No. 16. P. 7918–7928. doi: [10.1128/JVI.00645-06](https://doi.org/10.1128/JVI.00645-06)
  10. *Virology: principles and applications*. Eds. J. Carter, V. Saunders. Chichester, England: John Wiley & Sons Ltd, 2007. 358 p. doi: [10.1093/tropej/fmn001](https://doi.org/10.1093/tropej/fmn001)
  11. Chen Y., Liu Q., Guo D. Emerging coronaviruses: Genome structure, replication, and pathogenesis. *J. Med. Virol.* 2020. V. 92. P. 418–423. doi: [10.1002/jmv.25681](https://doi.org/10.1002/jmv.25681)
  12. Ma Y., Wu L., Shaw N., Gao Y., Wang J., Sun Y., Lou Z., Yan L., Zhang R., Rao Z. Structural basis and functional analysis of the SARS coronavirus nspl4-nspl0 complex. *PNAS*. 2015. V. 112. No. 30. P. 9436–9441. doi: [10.1073/pnas.1508686112](https://doi.org/10.1073/pnas.1508686112)
  13. Cavanagh D., Mawditt K., Adzharet A., Gough R.E., Picault J.P., Naylor C.J., Haydon D., Shaw K., Britton P. Does IBV change slowly despite the capacity of the spike protein to vary greatly? *Adv. Exp. Med. Biol.* 1998. V. 440. P. 729–734. doi: [10.1007/978-1-4615-5331-1\\_94](https://doi.org/10.1007/978-1-4615-5331-1_94)
  14. Pyrc K., Dijkman R., Deng L., Jebbink M.F., Ross H.A., Berkhout B., der Hoek L. Mosaic structure of human coronavirus NL63, one thousand years of evolution. *J. Mol. Biol.* 2006. V. 364. P. 964–973. doi: [10.1016/j.jmb.2006.09.074](https://doi.org/10.1016/j.jmb.2006.09.074)
  15. Su S., Wong G., Shi W., Liu J., Lai A.C.K., Zhou J., Liu W., Bi Y., Gao G.F. Epidemiology, Genetic Recombination, and Pathogenesis of Coronaviruses. *Trends in Microbiology*. 2016. V. 24. No. 6. P. 490–502. doi: [10.1016/j.tim.2016.03.003](https://doi.org/10.1016/j.tim.2016.03.003)

16. Edara V-V., Pinsky B.A., Suthar M.S., Lai L., Davis-Gardner M.E., Floyd K., Flowers M.W., Wrammert J., Hussaini L., Rose Ciric C. et al. Infection and vaccine-induced neutralizing-antibody responses to the SARS-CoV-2 B.1.617 Variants. *N. Engl. J. Med.* 2021. V. 385. No. 7. P. 664–666. doi: [10.1056/NEJMc2107799](https://doi.org/10.1056/NEJMc2107799)
17. Long S.W., Olsen R.J., Christensen P.A., Subedi S., Olson R., Davis J.J., Ojeda Saavedra M., Yerramilli P., Pruitt L., Reppond K. Sequence Analysis of 20,453 Severe Acute Respiratory Syndrome coronavirus 2 genomes from the Houston metropolitan area identifies the emergence and widespread distribution of multiple isolates of all major variants of concern. *Am. J. Pathol.* 2021. V. 191. No. 6. P. 983–992. doi: [10.1016/j.ajpath.2021.03.004](https://doi.org/10.1016/j.ajpath.2021.03.004)
18. Борисова Н.И., Котов И.А., Колесников А.А., Каптелова В.В., Сперанская А.С., Кондрашева Л.Ю., Тиванова Е.В., Хафизов К.Ф., Акимкин В.Г. Мониторинг распространения вариантов SARSCoV-2 (*Coronaviridae: Coronavirinae: Betacoronavirus; Sarbecovirus*) на территории Московского региона с помощью таргетного высокопроизводительного секвенирования. *Вопросы вирусологии.* 2021. Т. 66. № 4. С. 269–278. doi: [10.36233/0507-4088-72](https://doi.org/10.36233/0507-4088-72)
19. Чалей М.Б., Тюлько Ж.С., Кутыркин В.А. Распознавание видов флавивирусов на основе кодирующих последовательностей полипротеинов. *Математическая биология и биоинформатика.* 2019. Т. 14. № 2. С. 533–542. doi: [10.17537/2019.14.533](https://doi.org/10.17537/2019.14.533)
20. Чалей М.Б., Тюлько Ж.С., Кутыркин В.А. Исследование структуры кодирования ORF1ab, S, M и N генов коронавирусов. *Математическая биология и биоинформатика.* 2020. Т. 15. № 2. С. 441–454. doi: [10.17537/2020.15.441](https://doi.org/10.17537/2020.15.441)
21. Lai M.M.C. Recombination in large RNA viruses: Coronaviruses. *Seminars in Virology.* 1996. V. 7. No. 6. P. 381–388. doi: [10.1006/smy.1996.0046](https://doi.org/10.1006/smy.1996.0046)
22. Tao Y., Shi M., Chommanard C., Queen K., Zhang J., Markotter W., Kuzmin I.V., Holmes E.C., Tong S. Surveillance of bat coronaviruses in Kenya identifies relatives of human coronaviruses NL63 and 229E and their recombination history. *Journal of Virology.* 2017. V. 91. No. 5. P. e01953–16. doi: [10.1128/JVI.01953-16](https://doi.org/10.1128/JVI.01953-16)
23. Luk H.K.H., Li X., Fung J., Lau S.K.P., Woo P.C.Y. Molecular epidemiology, evolution and phylogeny of SARS coronavirus. *Infection, Genetics and Evolution.* 2019. V. 71. P. 21–30. doi: [10.1016/j.meegid.2019.03.001](https://doi.org/10.1016/j.meegid.2019.03.001)
24. *GenBank.* URL: <https://ftp.ncbi.nlm.nih.gov/genbank/> (accessed 20.12.2021).

Рукопись поступила в редакцию 28.01.2022, переработанный вариант поступил 09.03.2022.  
Дата опубликования 15.03.2022.

# Coronavirus Genus Recognition Based on Prototype Virus Variants

Chaley M.B.<sup>1</sup>, Kutyrkin V.A.<sup>2</sup>

<sup>1</sup>*Institute of Mathematical Problems of Biology RAS, Pushchino, Russia*

<sup>2</sup>*Moscow State Technical University n.a. N.E. Bauman, Moscow, Russia*

**Abstract.** Method named as variant approach to recognizing genus of coronavirus that is based on frequency of codon distribution in viral ORF1ab and genes of structural proteins (S, M and N) was proposed in the work. This method uses modified statistics whose efficiency was demonstrated earlier for flavivirus species recognition. To recognize genus of coronavirus the variant approach considers both various combinations of several structural coronavirus genes and individual structural genes. Finally, coronavirus genus is determined in the result of analysis of all variants considered. The method proposed was developed with the help of learning sample from prototype viral variants of *Alphacoronavirus*, *Betacoronavirus*, *Deltacoronavirus* and *Gammacoronavirus* genus. Application of the variant approach to recognizing genus of coronavirus has demonstrated the approach high assurance at level of 95 %. Among all variants of joint analysis, the most reliability (98 %) in recognizing genus has been achieved if codon frequency of the ORF1ab was used. Variant approach has revealed a phenomenon of mosaic structure in coronavirus genomes, i.e., when the results of genus recognition for a few genes differ from final conclusion about coronavirus genus. It seems that such phenomenon reflects homologous recombinations of the genes between various species of the coronaviruses and plasticity of their genomes in evolutionary processes.

**Key words:** *coronavirus genome, ORF1ab, S-gene, M-gene, N-gene, statistical analysis, variant approach to recognizing coronavirus genus.*