

# Optimization of a coronavirus genus recognition procedure based on the n-gene of prototypic strains

*Maria Chaley*<sup>1,\*</sup>, and *Vladimir Kutyrkin*<sup>2</sup>

<sup>1</sup>IMPB RAS – Branch of Keldysh Institute of Applied Mathematics RAS, Pushchino, Russia

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

**Abstract.** The article offers a solution to the problem of fast and efficient recognition of the coronavirus genus. For this purpose, the authors apply a virus genome targeting method based on the use of a sufficiently short and conserved N-gene of the nucleocapsid protein. Comparison of the codon frequency distributions in the N-gene of the analyzed genome and a set of 67 prototypical strains corresponding to the coronavirus subgenus allows us to recognize the genus of the coronavirus. This paper proposes optimization of the genus recognition of coronavirus by eliminating a significant number of codons from the 64 codons of the genetic code (26 in one case and 57 in the other). The authors achieved 100% genus recognition efficiency in a sample of 2,051 coronavirus genomes from the GenBank database with annotated subgenus in the optimized procedure. The authors also achieved 99% confidence when using the optimized coronavirus genus recognition procedure in a total sample of 3,242 genomes.

## 1 Introduction

Coronaviruses (Coronaviridae) are a large family of RNA-containing viruses capable of infecting humans and animals. The COVID-19 pandemic caused by SARS-CoV-2 has brought renewed attention to coronaviruses and their evolutionary potential [1]. High mutation and recombination rates contribute to the genotypic and phenotypic variability of coronaviruses [2]. Many scientists have studied the transmission of the virus from the main host population to other animal species and humans at the molecular genetic level [3]. Modern methods of analyzing viral genomes can create a system for monitoring gene pools of viral populations in unique ecosystems [4].

A reasonably fast and cost-effective targeted sequencing approach helps in tracking mutant strains of SARS-CoV-2 and other coronavirus species [5]. Here, only the sequence of the most variable fragment of the S-protein gene responsible for binding to cellular receptors during infection was determined. We can also propose a targeting approach for the rapid identification of coronavirus species, genus, and subgenus based on separately isolated fragments of the viral genome.

---

\* Corresponding author: [maramaria@yandex.ru](mailto:maramaria@yandex.ru)

Work [6] examined various approaches to the recognition of the coronavirus genus (CoVs) based on the genomes of prototypical coronavirus strains. The authors characterized the coronavirus genome by the distribution of codon frequencies of its individual structural and nonstructural genes: the M gene for membrane protein, the S gene for spike protein, the N gene for nucleocapsid protein, and the ORF1ab gene encoding several nonstructural proteins. Scientists divided each of the four genera ( $\alpha$ -CoV,  $\beta$ -CoV,  $\delta$ -CoV,  $\gamma$ -CoV) into subgenera that also have several prototypic strains. The coronavirus genus was characterized by an analytical, i.e. averaged, distribution of codon frequencies of all its prototypic strains. In the variant approach [6] to genus recognition, each variant (regarded as a target) relied on different combinations of structural and non-structural genes. The variant based on the N-gene of the nucleocapsid protein showed the best result in recognizing the genera  $\beta$ -CoV,  $\delta$ -CoV, and  $\gamma$ -CoV. Averaging the codon frequency distributions of prototypic strains of a significant number (12) subgenera in  $\alpha$ -CoV caused a low recognition efficiency of this genus.

This paper was aimed at fast and reliable identification of the genus from a fragment (target) of the coronavirus genome. The authors chose the N-gene among the previously identified targets as the most conservative and significantly shorter (~1200 nucleotides) compared to other genes. The N-gene encoded nucleocapsid protein has functions related to viral pathogenesis, transcription and replication. Scientists often use the N-gene for molecular diagnostics of CoVs [7-9].

We moved from recognizing the genus of coronavirus based on analytical (averaged) codon frequency distributions [6] to recognition based on individual frequency distributions in the N-gene of prototypic subgenus strains. Determining the subgenus of a coronavirus automatically determines its genus. We investigated the dependence of recognition efficiency on the use of different codon groups in the N-gene. The groups were formed according to selected levels of codon frequencies averaged over the N-gene of 67 prototypic strains of all four coronavirus genera. We identified the two most efficient groups of codons, including 38 codons in one case and 7 codons in the other. Reducing the number of codons used optimized the recognition procedure. This optimization improved the efficiency of coronavirus genus recognition compared to recognition based on the analytical mean distributions of the coronavirus genus of the paper [6].

## 2 Materials

We used coronavirus genomes from four genera from the GenBank database: 3242 pieces, of which for 2051 genomes a subgenus was specified beside the genus. We used the codon distributions in the N genes of the prototypic strains from GenBank in the genus recognition procedure. Table 1 contains the access codes of the prototypical coronavirus strains in GenBank.

**Table 1.** GenBank access codes for prototypic coronavirus strains. Subgenus designations and indices of prototypic strains used in this work.

Genus	Subgenus	GenBank ID	Variant index k
$\alpha$ -CoV	<i>Colacovirus</i> ( $\alpha$ 01)	NC 022103	1
	<i>Decacovirus</i> ( $\alpha$ 02)	NC 028814	1
		NC 018871	2
	<i>Duvinacovirus</i> ( $\alpha$ 03)	NC 002645	1
	<i>Luchacovirus</i> ( $\alpha$ 04)	NC 032730	1
	<i>Minacovirus</i> ( $\alpha$ 05)	NC 023760	1
		KX512809	2
		KX512810	3

Continuation of Table 1.

$\beta$ -CoV	<i>Minunacovirus</i> ( $\alpha$ 06)	EU420138	1
		NC 010438	2
	<i>Miotacovirus</i> ( $\alpha$ 07)	NC 028811	1
	<i>Nyctacovirus</i> ( $\alpha$ 08)	NC 028833	1
	<i>Pedacovirus</i> ( $\alpha$ 09)	KT323979	1
		NC 009657	2
	<i>Rhinacovirus</i> ( $\alpha$ 10)	NC 009988	1
	<i>Setracovirus</i> ( $\alpha$ 11)	AY567487	1
		KY073745	2
	<i>Tegacovirus</i> ( $\alpha$ 12)	NC 038861	1
		KP981644	2
		FJ938051	3
		AY994055	4
		KR270796	5
	<i>Embecovirus</i> ( $\beta$ 01)	KF294357	1
		BCU00735	2
		KX432213	3
		EF446615	4
		AY391777	5
		NC 017083	6
		MF083115	7
		NC 026011	8
		AC 000192	9
		KF294371	10
		NC 012936	11
		NC 006577	12
	<i>Hibecovirus</i> ( $\beta$ 02)	NC 025217	1
<i>Merbecovirus</i> ( $\beta$ 03)	KF917527	1	
	JX869059	2	
	MG596803	3	
	MK679660	4	
	NC 009019	5	
	NC 009020	6	
<i>Nobecovirus</i> ( $\beta$ 04)	NC 030886	1	
	NC 009021	2	
<i>Sarbecovirus</i> ( $\beta$ 05)	MG772933	1	
	MG772934	2	
	AY278489	3	
	FJ588686	4	
	NC 045512	5	
	MT121216	6	
	MN996532	7	
$\delta$ -CoV	<i>Andecovirus</i> ( $\delta$ 01)	NC 016995	1
	<i>Buldecovirus</i> ( $\delta$ 02)	JQ065042	1
		KJ569769	2
		NC 016992	3
		NC 016991	4
		FJ376620	5
		NC 011550	6
		NC 016993	7
	<i>Herdecovirus</i> ( $\delta$ 03)	NC 016994	1
	<i>Moordecovirus</i> ( $\delta$ 04)	NC_016996	1

Continuation of Table 1.

$\gamma$ -CoV	<i>Cegacovirus</i> ( $\gamma$ 01)	EU111742	1
		KF793826	2
	<i>Igacovirus</i> ( $\gamma$ 02)	KF696629	1
		GQ504724	2
		NC 010800	3
		AY641576	4
		MK423877	5

Next to the name of the subgenus in parentheses is its designation in this paper (Table 1). We considered 67 prototypic strains, each of which is indexed in its subgenus. Work [10] provides details of species with prototypic coronavirus strains.

### 3 Methods

With the high mutation rate of viruses, it would be interesting to assess which N-gene codons determine whether a coronavirus belongs to a subgenus. This raises the question of a reduction in the number of codons in the N-gene of a coronavirus to recognize its subgenus. We selected the N-gene codons relative to their average frequency of occurrence in all subgenus of the prototype strains (see Table 1). Figure 1 shows a graph of the average frequency of codons in the N-gene of prototypic coronavirus strains. Thirty-eight codons occur with frequencies above the threshold shown by the dashed line. Seven codons occur with frequencies above the threshold shown by the dotted line.

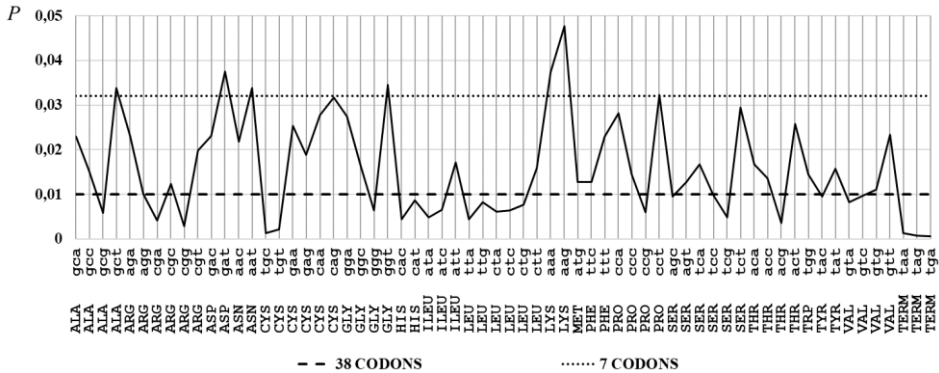


Fig. 1. Average frequency P of occurrence of codons in the N-gene of prototypic coronavirus strains.

Selecting a codon frequency threshold means that only those codons whose frequency exceeds the threshold are suitable for recognition. Figure 1 shows two such thresholds and the number of codons whose frequencies exceed each of them with different horizontal lines. Together with the codons, it shows the amino acids they encode along the horizontal axis.

Let the chosen threshold leave the set of  $C_r = \{c_1, c_2, \dots, c_r\}$  N-gene codons, where  $|C_r| = r$  – the number of codons in this population.  $\rho_i$  – the frequency of occurrence of a codon  $c_i$  in the N-gene of the coronavirus in question  $i = 1, r$ . The string  $\rho = (\rho_1, \rho_2, \dots, \rho_r)$  characterizes the coronavirus genome under analysis.

$S \in \{\alpha, \beta, \delta, \gamma\}$  is the symbol for the genus of the coronavirus,  $j$  is the two-digit subgenus number in the genus, and  $k$  is the index of the strain in this subgenus according to Table 1. For example,  $\alpha 12, 4$  is the notation of the prototype strain AY994055 with an index

$k = 4$  of the subgenus *Tegacovirus* from  $\alpha$ -CoV. Then  $\rho_i^{S_j,k}$  – the frequency of occurrence of codon  $c_i$  in the N-gene of the prototypic strain  $S_j,k$ , where  $S_j$  is the subgenus of the strain (and  $k$  is the index of the strain in that subgenus). Thus, the prototypic strain  $S_j,k$  is characterized by the string  $\mathbf{\rho}^{S_j,k} = (\rho_1^{S_j,k}, \rho_2^{S_j,k}, \dots, \rho_r^{S_j,k})$ .

The number  $\Delta^{S_j,k}(\mathbf{\rho})$  calculated using the formula (1) defines the deviation of the coronavirus genome with a frequency distribution of  $\mathbf{\rho}$  from the prototypical strain  $S_j,k$ .

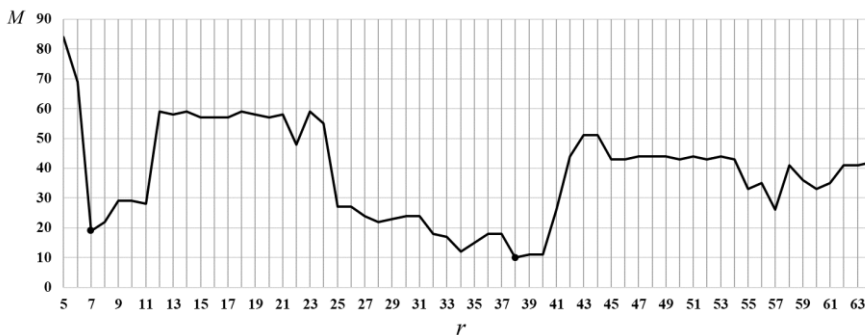
$$\Delta^{S_j,k}(\mathbf{\rho}) = \frac{1}{7} \sum_{i=1}^r \frac{|\rho_i - \rho_i^{S_j,k}|}{\rho_i^{S_j,k}}, \tag{1}$$

where the coefficient  $\frac{1}{7}$  is chosen for ease of presentation of results.

Having got the distribution deviations (1) from all the distributions of the prototypical strains, we choose the strain  $S_j,k$  with the minimum deviation. We hypothesize that the coronavirus genome in question belongs to subgenus  $S_j$  and thus to genus  $S$  (see Table 1).

### 4 Results and discussion

We denote  $M(r)$  the number of errors in the coronavirus's recognition subgenus using a set of codons  $C_r$  and formula (1) for the 2051 coronavirus genome annotated by the subgenus. Figure 2 shows the dependence of  $M(r)$  the number of errors on the number of  $r$  codons used.



**Fig. 2.** Plot of the number of subgenus recognition errors in 2051 N-genes of coronaviruses as a function of the size  $r$  of the set of the genetic code codons used for recognition.

Analysis of the graph in Figure 2 reveals two particular minima on sets of 38 and seven codons. The set of 38 codons achieves the minimum number of errors ( $M = 10$ ). The analysis of Figure 2 reveals a sharp minimum ( $M = 19$ ) from a set of seven codons. Figure 3 shows the codons that are part of the aggregates  $C_{38}$  and  $C_7$ .

		<b>C<sub>38</sub></b>	<b>C<sub>7</sub></b>			<b>C<sub>38</sub></b>	<b>C<sub>7</sub></b>			<b>C<sub>38</sub></b>	<b>C<sub>7</sub></b>			<b>C<sub>38</sub></b>	<b>C<sub>7</sub></b>
<b>LYS</b>	<b>aaa</b>	1	1	<b>THR</b>	<b>aca</b>	1	0	<b>ALA</b>	<b>gct</b>	1	1	<b>LEU</b>	<b>ctt</b>	1	0
<b>LYS</b>	<b>aag</b>	1	1	<b>THR</b>	<b>acc</b>	1	0	<b>ALA</b>	<b>gca</b>	1	0	<b>LEU</b>	<b>cta</b>	0	0
<b>ASN</b>	<b>aat</b>	1	1	<b>THR</b>	<b>act</b>	1	0	<b>ALA</b>	<b>gcc</b>	1	0	<b>LEU</b>	<b>ctc</b>	0	0
<b>ASN</b>	<b>aac</b>	1	0	<b>THR</b>	<b>acg</b>	0	0	<b>ALA</b>	<b>gcg</b>	0	0	<b>LEU</b>	<b>ctg</b>	0	0
<b>ARG</b>	<b>cgt</b>	1	0	<b>MET</b>	<b>atg</b>	1	0	<b>VAL</b>	<b>gtt</b>	1	0	<b>LEU</b>	<b>tta</b>	0	0
<b>ARG</b>	<b>cgc</b>	1	0	<b>ILEU</b>	<b>att</b>	1	0	<b>VAL</b>	<b>gtg</b>	1	0	<b>LEU</b>	<b>tgt</b>	0	0
<b>ARG</b>	<b>cga</b>	0	0	<b>ILEU</b>	<b>ata</b>	0	0	<b>VAL</b>	<b>gta</b>	0	0	<b>PHE</b>	<b>ttt</b>	1	0
<b>ARG</b>	<b>cgg</b>	0	0	<b>ILEU</b>	<b>atc</b>	0	0	<b>VAL</b>	<b>gtc</b>	0	0	<b>PHE</b>	<b>ttc</b>	1	0
<b>ARG</b>	<b>aga</b>	1	0	<b>ASP</b>	<b>gat</b>	1	1	<b>CYS</b>	<b>caa</b>	1	0	<b>TYR</b>	<b>tat</b>	1	0
<b>ARG</b>	<b>agg</b>	1	0	<b>ASP</b>	<b>gac</b>	1	0	<b>CYS</b>	<b>cag</b>	1	0	<b>TYR</b>	<b>tac</b>	0	0
<b>SER</b>	<b>agt</b>	1	0	<b>CYS</b>	<b>gaa</b>	1	0	<b>HIS</b>	<b>cac</b>	0	0	<b>TERM</b>	<b>taa</b>	0	0
<b>SER</b>	<b>agc</b>	0	0	<b>CYS</b>	<b>gag</b>	1	0	<b>HIS</b>	<b>cat</b>	0	0	<b>TERM</b>	<b>tag</b>	0	0
<b>SER</b>	<b>tca</b>	1	0	<b>GLY</b>	<b>ggt</b>	1	1	<b>PRO</b>	<b>cct</b>	1	1	<b>TERM</b>	<b>tga</b>	0	0
<b>SER</b>	<b>tct</b>	1	0	<b>GLY</b>	<b>gga</b>	1	0	<b>PRO</b>	<b>cca</b>	1	0	<b>TRP</b>	<b>tgg</b>	1	0
<b>SER</b>	<b>tcc</b>	0	0	<b>GLY</b>	<b>ggc</b>	1	0	<b>PRO</b>	<b>ccc</b>	1	0	<b>CYS</b>	<b>tgc</b>	0	0
<b>SER</b>	<b>tcg</b>	0	0	<b>GLY</b>	<b>ggg</b>	0	0	<b>PRO</b>	<b>ccg</b>	0	0	<b>CYS</b>	<b>tgt</b>	0	0

**Fig. 3.** The C<sub>38</sub> and C<sub>7</sub> aggregates of the genetic code codons with the lowest number of subgenus recognition errors by the N-gene of coronavirus. The number 1 shows codon entry into the aggregate, 0 shows no codon in the aggregate. To the left of the codon is the three-letter designation of the amino acid encoded. Stop codons are designated as TERM.

Note that there are no genus recognition errors out of 10 for subgenus recognition based on 38 codons. When subgenus recognition is based on 7 codons out of 19 errors, there are 14 errors in genus recognition of coronavirus.

If we consider a combined sample of genus and subgenus annotated (2051 genomes) and genus-only annotated (1191 genomes) coronavirus genomes, recognition of coronavirus genus based on 38 codons will lead to 10 errors and recognition based on 7 codons to 34 errors.

The proposed methods for subgenus and genus recognition of coronaviruses yield at least 99% confidence.

## 5 Conclusions

When using the above two approaches for genus recognition based on N genes, out of 2051 genomes annotated for the subgenus in the first case (with averaging), we noted 43 errors (of which 19 were incorrectly determined by the genus coronavirus) and in the second (without averaging); we noted 36 errors, of which only three were incorrectly determined by the genus. This result shows that the use of analytical (averaged) codon frequency distributions of prototypic strains is less effective for recognizing both the subgenus and genus of coronaviruses.

## References

1. D. K. Lvov, S. V. Alkhovsky, L. V. Kolobukhina, E. I. Burtseva, *Problems of Virology* **65**, 6 (2020)
2. E. Salem, V. Dhanasekaran, H. Cassard, B. Hause, S. Maman, G. Meyer, et al., *Viruses* **12**, 534 (2020)
3. E. Corrales-Aguilar, M. Schwemmler Eds., *Bats and Viruses: Current Research and Future Trends* (Academic Press, Caister, 2020)
4. D. K. Lvov, S. V. Borisevich, S. V. Alkhovsky, E. I. Burtseva, *Infectious Diseases: News, Opinions, Training* **8**, 96 (2019)
5. N. I. Borisova, I. A. Kotov, A. A. Kolesnikov, V. V. Kaptelova, A. S. Speranskaya, L. Yu. Kondrasheva, et al., *Problems of Virology* **66**, 269 (2021)
6. M. B. Chaley, V. A. Kutyrkin, *Math. Biol. Bioinf.* **17**, 10 (2022)

7. L. J. Saif, *Vet. Clin. North Am. Food Anim. Pract.* **26**, 349 (2010)
8. A. N. Vlasova, L. J. Saif, *Front. Vet. Sci.* **8**, 643220 (2021)
9. A. G. Glotov, A. V. Nefedchenko, A. G. Yuzhakov, S. V. Koteneva, T. I. Glotova, A. K. Komina, N. Yu. Krasnikov, *Problems of Virology* **67**, 465 (2022)
10. M. Yu. Shchelkanov, A. Yu. Popova, V. G. Dedkov, V. G. Akimkin, V. V. Maleev, *Russian Journal of Infection and Immunity* **10**, 221 (2020)