

SARS-CoV-2 and MERS-CoV Share the Furin Site CGG-CGG Genetic Footprint

Antonio R. Romeu¹

¹: Professor of Biochemistry and Molecular Biology. University Rovira i Virgili. Tarragona. Spain. Corresponding author. Email: antonioramon.romeu@iubilo.urv.cat

Abstract

The SARS-CoV-2 polybasic furin cleavage site is still a missing link. Remarkably, the two arginine residues of this protease recognition site are encoded by the CGG codon, which is rare in Betacoronavirus. However, the arginine pair is common at viral furin cleavage sites, but are not CGG-CGG encoded. The question is: Is this genetic footprint unique to the SARS-CoV-2? To address the issue, using Perl scripts, here I dissect in detail the NCBI Virus database in order to report the arginine dimers of the Betacoronavirus proteins. The main result reveals that a group of Middle East respiratory syndrome-related coronavirus (MERS-CoV) (isolates: camel/Nigeria/NVx/2016, host: *Camelus dromedarius*) also have the CGG-CGG arginine pair in the spike protein polybasic furin cleavage region. In addition, CGG-CGG encoded arginine pairs were found in the orf1ab polyprotein from HKU9 and HKU14 Betacoronavirus, as well as, in the nucleocapsid phosphoprotein from few SARS-CoV-2 isolates. To quantify the probability of finding the arginine CGG-CGG codon pair in Betacoronavirus, the likelihood ratio (LR) and a Markov model were defined. In conclusion, it is highly unlikely to find this genetic marker in betacoronaviruses wildlife, but they are there. Collectively, results shed light on recombination as origin of the virus CGG-CGG arginine pair in the S1/S2 cleavage site.

Key words

SARS-CoV-2, MERS-CoV, Arginine Pair, Polybasic Furin Cleavage Site, Arginine Codon, Markov Model, Bioinformatics.

Background

First of all, the structure and availability of the NCBI Virus database information (1), that makes this work possible, must be appreciated. Arginine is a polar and non-hydrophobic amino acid, with a positive charged group a physiological pH. Arginine participates in the binding of negatively charged substrates and/or protein actives sites (2). Consistently, arginine is involved in viral polybasic proteolytic cleavage sites, even as a dimers, as recognition motif of the ubiquitously expressed furin serine protease (3,4).

A notable characteristic of the SARS-CoV-2, that distinguishes from the rest of Sarbecovirus, is the acquisition of a polybasic furin cleavage site (PRRAR) at the S1-S2 boundary of the S glycoprotein (5). It greatly mediates the fusion of human cell and viral membranes, and the rapid human-to-human virus transmission (5-7). That acquisition was achieved through the insertion of four amino acids (PRRA). However, the furin protease recognition pattern is common in viral proteins, such as the hemagglutinin (H5) protein of the avian influenza viruses (3) or the spike glycoprotein of three of the seventh coronavirus known to infect humans (8): HCoV-HKU1 (RRKRR-760, coordinate based on GenBank: YP_173238.1), HCoV-OC43 (RRSRR-763, GenBank: AOL02453.1) and MERS-CoV (RSVRSV-753, GenBank: YP_009047204.1).

Another notable characteristic of the SARS-CoV-2 is the CGG-CGG coding sequence of the arginine dimer in that polybasic furin cleavage site. In the genetic code, arginine is encoded by six codons AGA, AGG, CGC, CGA,

CGG and CGT codons. CGG is a minority arginine codon in SARS-CoV-2 (9). Consistently, CGG-CGG encoded arginine dimers at viral polybasic furin cleavage sites have not been found (10). In this sense, SARS-CoV-2 has the most extreme CpG deficiency in all known Betacoronavirus genomes, probably to avoid the human antiviral defence, mediated by the zinc finger antiviral protein (ZAP) (11). On the other hand, the other thirteen SARS-CoV-2 proteome arginine dimers, which are strictly conserved in the closest Sarbecovirus strains, are not CGG-CGG encoded either (12).

Is the CGG-CGG encoded arginine dimer unique to SARS-CoV-2 polybasic furin cleavage site?

Based on the NCBI Virus database as a source of information, through a bioinformatics approach and using Perl scripts, all current Betacoronavirus arginine dimers and their coding regions are here reported. Full updated results are available in a Google Drive Folder (see below the Web address). Interestingly, arginine dimers were widely distributed in Betacoronavirus proteins, about 30% of them contained one or more of the amino acid pair. These proteins were mostly members of the non-structural orf1ab-polyprotein complex, and also in the structural S glycoprotein and nucleocapsid phosphoprotein. As regards the arginine codon usage focused on the Betacoronavirus arginine dimers, AGA was the majority (about 50%), followed by CGT (about 24%). CGG was minority (about 5%).

Table 1 summarizes the Betacoronavirus arginine dimers, that were encoded by CGG-CGG. The most remarkable discovery was the CGG-CGG arginine pair close to the furin recognition site of the spike glycoprotein from a group of MERS-CoVs (Table 1). Based on MERS-CoV spike glycoprotein structure (13), the S2 chain spans from arginine R-748 (coordinate based on UniProtKB – A0A023SFE5) to the C-terminal histidine H-1353, residues. In the case of human-infection, the MERS-CoV S glycoprotein is cleaved at R-748 generating the S1 and S2 subunits (8). However, it is worth noting that the CGG-CGG encoded arginine pair reported here (RR-700, coordinate based on GenBank AVN89376.1) is located 47 residues upstream the S1/S2 cleavage site (R-748), that creates, with a lysine residue, a true polybasic motif (KRR-700). Figure 1 shows sequence details. From the entire Betacoronavirus protein sample, there were 684 MERS-CoV spike glycoprotein sequences, of which 8 (1.17%) had the CGG-CGG encoded RR-700 dimer, in the rest was CGC-CGA encoded. In addition, the Betacoronavirus species MERS-CoV, *Rousettus* and *Eidolon helvum* bat coronavirus HKU9 and Rabbit coronavirus HKU14 also had CGG-CGG encoded arginine dimers in their orf1ab-polyprotein (Table 1). Within the SARS-CoV-2 species (apart from the S glycoprotein), only two SARS-CoV-2 isolates from North America showed a CGG-CGG arginine dimer in the orf1ab-polyprotein, and few SARS-CoV-2 isolates, also from North America, showed the first (out of four) nucleocapsid phosphoprotein arginine dimer encoded by CGG-CGG (Table 1).

CGG-CGG likelihood ratio (LR) and Markov model

Based on the structure of the NCBI Virus database, the results are grouped by Geographic Regions. The observed frequencies of the arginine codon pairs can be associated with probabilities. Also, based on the principles of forensic genetics (14), it is appropriate to ask for the LR value of the CGG-CGG genetic footprint, as a fundamental genetic marker of the pandemic virus. Given a Geographic Region, LR compares (ratio) the probability (P1) that if CGG-CGG encoding RR pair belongs to the SARS-CoV-2 (obviously, $P1 = 1$) with the probability (P2) that if CGG-CGG encoding RR pair belongs to a random Betacoronavirus isolate from the same SARS-CoV-2 Geographic Region (frequency). Only Africa and Asia Geographic Regions showed CGG-CGG frequencies other than zero, with the following LR values:

| Geographic Region | P1 | P2 | LR |
|-------------------|----|-----------|-----------|
| Africa | 1 | 1.82 E-04 | 5,492.90 |
| Asia | 1 | 9.27 E-05 | 10,790.15 |

In forensic genetics LR is used by juries or judges to draw inferences or conclusions and decide legal matters (14). So that LR should be large enough to allow that a genetic marker could be considered unique of a given forensic evidence. Here, the Africa and/or Asia Betacoronavirus LR were not excessively high, which agreed that arginine CGG-CGG is not unique SARS-CoV-2 genetic footprint.

On the other hand, to quantify the probability of the CGG-CGG presence, a First-Order Markov Chain was defined. The states were the arginine codons themselves. This Markov model allowed to determine the probability of the second arginine codon depending on the previous codon. Since arginine has six codons, in an arginine dimer there are 36 (6 x 6) chances of finding a codon pair (like a roll of two dice: 36 possible outcomes). By normalizing the codon pair frequencies, the stochastic matrix of the Markov chain could be created, whose elements are the transition probability between codons (states). As an example, Table 2 shows a stochastic Markov matrix, based on the arginine dimers found in a recent Asia Betacoronavirus protein sample. In this sample, if the first codon was AGG or CGA, the second was most likely AGA. If the first codon was CGG, the second was most likely CGT. The elements on the main diagonal mean the probability that the second codon was the same as the first. The significant presence of two arginine codons in a row occurred only in AGA-AGA.

Concluding remarks

In this work, about ten million of Betacoronavirus protein and coding sequences have been analysed, as well as, few million more of arginine pairs. It was a large sample which grow day by day. So, the present results are also updated (Google Drive Folder). Furthermore, analysis of arginine pairs from viruses of other taxonomic groups is going on. In conclusion, excluding the pair of the SARS-CoV-2 furin site (PRRAR), the arginine CGG-CGG encoding is highly unlikely in betacoronaviruses wildlife, but they are there. However, just because that presence, recombination may have operated into the origin of the virus S1/S2 protease recognition site. Recombination is the common method of viruses picking up new skills (15-19).

Full updated results:

https://drive.google.com/drive/folders/1Dp04BHDyMay1sBOGX000IFzfZTp_VrBu?usp=sharing

Acknowledgements

This work has not been awarded grants by any research-supporting institution.

Competing interest declaration

Author declare that he has no conflicts of interest.

References

1. National Center for Biotechnology Information (NCBI). NCBI Virus database. Accessed September 23, 2021. <https://www.ncbi.nlm.nih.gov/labs/virus/vssi/#/>. Eneida L Hatcher, Sergey A Zhdanov, Yiming Bao, Olga Blinkova, Eric P Nawrocki, Yuri Ostapchuck, Alejandro A Schäffer, J Rodney Brister. Virus Variation Resource-improved response to emergent viral outbreaks. *Nucleic Acids Res.* 45(D1):D482-D490, 2017. PMID: 27899678 doi: 10.1093/nar/gkw1065.
2. Michael J. Harms, Jamie L. Schlessman, Gloria R. Sue, and Bertrand García-Moreno. Arginine residues at internal positions in a protein are always charged. *Proc. Natl. Acad. Sci. U S A.* 108(47):18954-18959, 2011. PMID: 22080604. doi:10.1073/pnas.1104808108.

3. Elisabeth Braun, Daniel Sauter. Furin-mediated protein processing in infectious diseases and cancer. *Clin. Transl. Immunol.* E1073, 2019. PMID: 31406574. doi.org/10.1002/cti2.1073.
4. Imène Kara, Marjorie Poggi, Bernadette Bonardo, Roland Govers, Jean-François Landrier, Sun Tian, Ingo Leibiger, Robert Day, John W M Creemers, Franck Peiretti. The Paired Basic Amino Acid-cleaving Enzyme 4 (PACE4) Is Involved in the Maturation of Insulin Receptor Isoform B. *J. Biol. Chem.* 290:2812-2821, 2015. PMID: 25527501. doi: 10.1074/jbc.M114.592543.
5. Kristian G. Andersen, Andrew Rambaut, W Ian Lipkin, Edward C Holmes, Robert F Garry. The proximal origin of SARS-CoV-2. *Nat. Med.* 26(4):450-452, 2020. PMID: 32284615. doi: 10.1038/s41591-020-0820-9.
6. Philip V'kovski, Annika Kratzel, Silvio Steiner, Hanspeter Stalder, Volker Thiel. Coronavirus biology and replication: implications for SARS-CoV-2. *Nat. Rev. Microbiol.* 19(3):155-170, 2021. PMID: 33116300. doi: 10.1038/s41579-020-00468-6.
7. Markus Hoffmann, Hannah Kleine-Weber, Stefan Pöhlmann. Multibasic Cleavage Site in the Spike Protein of SARS-CoV-2 Is Essential for Infection of Human Lung Cells. *Mol. Cell* 78(4):779-784, 2020. PMID: 32362314. doi: 10.1016/j.molcel.2020.04.022.
8. Bruno Coutard, Camille Valle, Xavier de Lamballerie, Bruno Canard, Nabil G Seidah, Etienne Decroly. The spike glycoprotein of the new coronavirus 2019-nCoV contains a furin-like cleavage site absent in CoV of the same clade. *Antiviral Res.* 176:104742, 2020. PMID: 32057769. doi: 10.1016/j.antiviral.2020.104742.
9. Xuhua Xia. Domains and Functions of Spike Protein in Sars-Cov-2 in the Context of Vaccine Design. *Viruses* 13(1):109, 2021. PMID: 33466921. doi: 10.3390/v13010109.
10. Antonio R. Romeu, Enric Ollé. SARS-CoV-2 and the Secret of the Furin Site. Preprints 2021, 2021020264. doi: 10.20944/preprints202102.0264.v1.
11. Shuai Xia, Qiaoshuai Lan, Shan Su, Xinling Wang, Wei Xu, Zezhong Liu, Yun Zhu, Qian Wang, Lu Lu, Shibo Jiang. The role of furin cleavage site in SARS-CoV-2 spike protein-mediated membrane fusion in the presence or absence of trypsin. *Signal Transduct. Target Ther.* 5(1):92, 2020. PMID: 32532959. doi: 10.1038/s41392-020-0184-0.
12. Antonio R. Romeu, Enric Ollé. The SARS-CoV-2 arginine dimers. Research Square Preprints 2021. doi:10.21203/rs.3.rs-770380/v1.
13. Spike glycoprotein. Middle East respiratory syndrome-related coronavirus (MERS-CoV). UniProtKB - AOA023SFE5 (AOA023SFE5_MERS). Accessed September 23, 2021. <https://www.uniprot.org/uniprot/AOA023SFE5>.
14. John M. Butler. *Fundamentals of Forensic DNA Typing*. 2009. Academic Press. ISBN 978-0-12-374999-4. Maryland, USA. doi.org/10.1016/C2009-0-01945-X
15. Simon-Loriere, E. & Holmes, E.C. Why do RNA viruses recombine? *Nat. Rev. Microbiol.* **9**, 617-26 (2011). PMID: 21725337. doi: 10.1038/nrmicro2614.
16. Zhou, H. et al. Identification of novel bat coronaviruses sheds light on the evolutionary origins of SARS-CoV-2 and related viruses. *Cell* **184**, 4380-4391.e14 (2021). PMID: 34147139. doi: 10.1016/j.cell.2021.06.008.
17. Li, L.-L. et al. A novel SARS-CoV-2 related coronavirus with complex recombination isolated from bats in Yunnan province, China. *Emerg. Microbes Infect.* **10**, 1683-1690 (2021). PMID: 34348599. doi: 10.1080/22221751.2021.1964925.

18. Bernd Kaina. On the Origin of SARS-CoV-2: Did Cell Culture Experiments Lead to Increased Virulence of the Progenitor Virus for Humans? *In Vivo* 35(3):1313-1326, 2021. PMID: 33910809. doi: 10.21873/invivo.12384.
19. Smriti Mallapaty. Closest known relatives of virus behind COVID-19 found in Laos. *Nature* 597, 603 (2021). PMID: 34561634. doi: <https://doi.org/10.1038/d41586-021-02596-2>.
20. Fábio Madeira, Young Mi Park, Joon Lee, Nicola Buso, Tamer Gur, Nandana Madhusoodanan, Prasad Basutkar, Adrian R N Tivey, Simon C Potter, Robert D Finn, Rodrigo Lopez. The EMBL-EBI search and sequence analysis tools APIs in 2019. *Nucl. Acids Res.* 47(W1):W636-W641, 2019. PMID: 30976793. doi: 10.1093/nar/gkz268

Table 1. CGG-CGG encoded arginine dimer from Betacoronavirus protein sequences. The list is limited to records that exclude those from the SARS-CoV-2 polybasic furin cleavage site (PRRAR). The data is grouped by Geographic Regions

| Species* | Isolate | Protein** | Acession | Length | Dimer | Protein Position | Coding | Gene/Genome Position | Host |
|----------|-----------------------------|-------------|--------------|--------|-------|------------------|--------|----------------------|------------------------------|
| Africa | | | | | | | | | |
| MERS-CoV | camel/Nigeria/NV1787/2016 | S protein | AVN89398 | 1353 | RR | 700 | CGGCGG | 2100 | <i>Camelus dromedarius</i> |
| MERS-CoV | camel/Nigeria/NV1712/2016 | S protein | AVN89453 | 1353 | RR | 700 | CGGCGG | 2100 | <i>Camelus dromedarius</i> |
| HKU9 | PREDICT-GVF-CM-ECO06464 | RdRp | ATU79936 | 112 | RR | 85 | CGGCGG | 255 | <i>Rousettus aegyptiacus</i> |
| MERS-CoV | camel/Nigeria/NV1673/2016 | S protein | AVN89387 | 1353 | RR | 700 | CGGCGG | 2100 | <i>Camelus dromedarius</i> |
| MERS-CoV | camel/Nigeria/NV1657/2016 | S protein | AVN89442 | 1353 | RR | 700 | CGGCGG | 2100 | <i>Camelus dromedarius</i> |
| HKU9 | PREDICT-GVF-CM-ECO06646 | RdRp | ATU79938 | 112 | RR | 85 | CGGCGG | 255 | <i>Eidolon helvum</i> |
| MERS-CoV | camel/Nigeria/NV2020/2016 | S protein | AVN89409 | 1353 | RR | 700 | CGGCGG | 2100 | <i>Camelus dromedarius</i> |
| MERS-CoV | camel/Nigeria/NV2040/2016 | S protein | AVN89420 | 1353 | RR | 700 | CGGCGG | 2100 | <i>Camelus dromedarius</i> |
| MERS-CoV | camel/Nigeria/NV1989/2016 | S protein | AVN89431 | 1353 | RR | 700 | CGGCGG | 2100 | <i>Camelus dromedarius</i> |
| MERS-CoV | camel/Nigeria/NV1405/2016 | S protein | AVN89376 | 1353 | RR | 700 | CGGCGG | 2100 | <i>Camelus dromedarius</i> |
| Asia | | | | | | | | | |
| HKU14 | | polyprotein | AFE48811 | 7151 | RR | 6763 | CGGCGG | 20289 | <i>Oryctolagus cuniculus</i> |
| HKU9 | Rousettus spp/Jinghong/2009 | ORF1ab | AVP25405 | 6920 | RR | 4951 | CGGCGG | 14853 | <i>Rousettus sp.</i> |
| HKU9 | | ORF1ab | ADM33573 | 6923 | RR | 2569 | CGGCGG | 7707 | <i>Chiroptera</i> |
| HKU14 | | polyprotein | AFE48810 | 7151 | RR | 6763 | CGGCGG | 20289 | <i>Oryctolagus cuniculus</i> |
| HKU14 | | polyprotein | AFE48822 | 7112 | RR | 6724 | CGGCGG | 20172 | <i>Oryctolagus cuniculus</i> |
| HKU14 | | nsp15 | YP_009924422 | 375 | RR | 286 | CGGCGG | 858 | <i>Oryctolagus cuniculus</i> |
| HKU14 | | ORF1ab | YP_005454239 | 7151 | RR | 6763 | CGGCGG | 20289 | <i>Oryctolagus cuniculus</i> |
| HKU9 | | ORF1ab | ABN10926 | 6923 | RR | 2569 | CGGCGG | 7707 | <i>Chiroptera</i> |
| MERS-CoV | | ORF1ab | AHY61336 | 7179 | RR | 6174 | CGGCGG | 18522 | <i>Vespertilio sinensis</i> |
| HKU9 | | ORF1ab | ADM33557 | 6923 | RR | 2569 | CGGCGG | 7707 | <i>Chiroptera</i> |

| | | | | | | | | | |
|-------|--|-------------|----------|------|----|------|--------|-------|------------------------------|
| HKU14 | | polyprotein | AFE48800 | 7151 | RR | 6763 | CGGCGG | 20289 | <i>Oryctolagus cuniculus</i> |
| HKU9 | | ORF1ab | ADM33565 | 6903 | RR | 2041 | CGGCGG | 6123 | <i>Chiroptera</i> |

| | |
|-------|---------|
| North | America |
|-------|---------|

| | | | | | | | | | |
|------------|----------------------|--------|----------|------|----|------|--------|-------|---------------------|
| SARS-CoV-2 | NC-CDC-LC0027271 | ORF9 | QTG55296 | 419 | RR | 41 | CGGCGG | 123 | <i>Homo sapiens</i> |
| SARS-CoV-2 | JW0066 | ORF9 | QXX31736 | 419 | RR | 41 | CGGCGG | 123 | <i>Homo sapiens</i> |
| SARS-CoV-2 | FL-BPHL-5767 | ORF9 | QZW21341 | 419 | RR | 41 | CGGCGG | 123 | <i>Homo sapiens</i> |
| SARS-CoV-2 | NC-CDC-STM-000025458 | ORF9 | QTC13492 | 419 | RR | 41 | CGGCGG | 123 | <i>Homo sapiens</i> |
| SARS-CoV-2 | NC-SLPH-0070 | ORF9 | QTX13191 | 419 | RR | 41 | CGGCGG | 123 | <i>Homo sapiens</i> |
| SARS-CoV-2 | NC-CDC-STM-000028277 | ORF9 | QTX74547 | 419 | RR | 41 | CGGCGG | 123 | <i>Homo sapiens</i> |
| SARS-CoV-2 | MI-CDC-STM-000045686 | ORF9 | QTW56377 | 419 | RR | 41 | CGGCGG | 123 | <i>Homo sapiens</i> |
| SARS-CoV-2 | NC-CDC-STM-000032315 | ORF9 | QTX83255 | 419 | RR | 41 | CGGCGG | 123 | <i>Homo sapiens</i> |
| SARS-CoV-2 | JW0864 | ORF9 | QXX38166 | 419 | RR | 41 | CGGCGG | 123 | <i>Homo sapiens</i> |
| SARS-CoV-2 | NC-CDC-STM-000026863 | ORF9 | QTC16696 | 419 | RR | 41 | CGGCGG | 123 | <i>Homo sapiens</i> |
| SARS-CoV-2 | MI-CDC-STM-000046692 | ORF9 | QTW59960 | 419 | RR | 41 | CGGCGG | 123 | <i>Homo sapiens</i> |
| SARS-CoV-2 | CA-CDC-ASC210119629 | ORF1ab | UBF72283 | 7096 | RR | 5767 | CGGCGG | 17301 | <i>Homo sapiens</i> |
| SARS-CoV-2 | MD-MDH-4405 | ORF1ab | UAB59607 | 7096 | RR | 5767 | CGGCGG | 17301 | <i>Homo sapiens</i> |

* Full name and taxonomic identifier of the Betacoronavirus species: SARS-CoV-2, Severe acute respiratory syndrome coronavirus 2 (taxid:2697049); MERS-CoV, Middle East respiratory syndrome-related coronavirus (taxid:1335626); HKU9, Roussettus bat coronavirus HKU9 (taxid:694006); HKU14, Rabbit coronavirus HKU14 (taxid:1160968).

** Protein name: RdRp, RNA-dependent RNA polymerase; ORF1ab, orf1ab polyprotein; nsp15, non structural protein 15; ORF9, nucleocapsid phosphoprotein.

Figure 1. Spike glycoprotein furin S1/S2 crecognition region

| | | | |
|------------|----------------|--|-----|
| SARS-CoV-2 | YP_009724390.1 | AIHADQL--TPTWRVYSTGSNVFQTRAGCLIGAEHVN-NSYECDIPIGAGICASYQTQTN | 679 |
| MERS-CoV | AVN89376.1 | TMSQYSRSTRSML KRR DSTYGPLQTPVGCVLGLVNSSLFVEDCKLPLGQSLCALPDTPST | 744 |
| | | * * * * * | |
| SARS-CoV-2 | YP_009724390.1 | -S PRRAR SVASQSI---IAYTMSLGAENSVAYSNNNSIAIPTNF ⁷³⁵ TISVTTEILPVSM ⁷³⁵ TKTS | 735 |
| MERS-CoV | AVN89376.1 | LT PRSVR SVPGEMRLASIAFNHPIQV-DQLNSSFYFKLSIPTNF ⁸⁰³ SFGVTQEYIQ ⁸⁰³ TTIQKVT | 803 |
| | | * * * * * | |

Fragment of spike glycoprotein pairwise alignment from SARS-CoV-2 and MERS-CoV, corresponding to the S1/S2 cleavage site region. The sequences were from the following Betacoronavirus: SARS-CoV-2, isolate Wuhan-Hu-1, NCBI Reference Sequence NC_045512.2; and MERS-CoV, isolate MERS-CoV camel/Nigeria/NV1405/2016, GenBank MG923474.1. Sequence alignment was created by Clustal Omega (v.1.2.4) using default parameters (19). Strictly conserved amino acids are denoted by *, gaps are denoted by -. Positions of sequence amino acid residues are indicated by the numbers on the right. The MERS-CoV **CGG-CGG** encoded arginine doublet (**RR-700**), located 47 residues upstream of the S1/S2 cleavage site is highlighted in bold and red, within the polybasic motif (**KRR**), highlighted in yellow. The specific SARS-CoV-2 and MERS-CoV furin protease recognition pattern and the S1/S2 cleavage positions **R-685** and **R-748**, respectively, are also highlighted in bold and yellow.

Table 2. Stochastic matrix of the First-Order Markov Chain defined by the arginine codon sequence encoding the arginine dimers of the Asia Betacoronavirus protein sample.

| | AGA | AGG | CGA | CGC | CGG | CGT | Sum |
|-----|------------|------------|------------|------------|------------|------------|-----|
| AGA | 0.29218655 | 0.18973403 | 0.11623226 | 0.01059752 | 0.00056244 | 0.39068721 | 1 |
| AGG | 0.80859016 | 0.00471785 | 0.07812021 | 0.00251210 | 0.08050977 | 0.02554991 | 1 |
| CGA | 0.81517094 | 0.00356125 | 0.01709402 | 0.05864198 | 0.00356125 | 0.10197056 | 1 |
| CGC | 0.21081081 | 0.35091892 | 0.04345946 | 0.07221622 | 0.01167568 | 0.31091892 | 1 |
| CGG | 0.08980123 | 0.00070989 | 0.10151443 | 0.01123994 | 0.00153810 | 0.79519640 | 1 |
| CGT | 0.39124861 | 0.27172256 | 0.00256375 | 0.02674612 | 0.30058204 | 0.00713692 | 1 |

The stochastic matrix is a square matrix of transition probability between arginine codons (states). The rows are probabilistic vectors. An element of the matrix means the probability that the second arginine codon would be that of the column if the first is that of the row. Consequently, the sum of the elements of a row is 1. Data used to create this Betacoronavirus-Asia stochastic Markov matrix: Total number of analysed Betacoronavirus protein sequences, 93,977; total number of protein sequences having arginine dimer(s), 34,346 (36.55%); total number of arginine dimers in the sample, 134,249; total number of SARS-CoV-2 (CGG-CGG) polybasic furin cleavage site arginine dimers, 6,859 (5.11%). To avoid distortions in calculations of transition probabilities between arginine codons, the arginine dimers of the SARS-CoV-2 furin site were excluded.