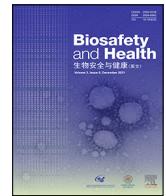




Since January 2020 Elsevier has created a COVID-19 resource centre with free information in English and Mandarin on the novel coronavirus COVID-19. The COVID-19 resource centre is hosted on Elsevier Connect, the company's public news and information website.

Elsevier hereby grants permission to make all its COVID-19-related research that is available on the COVID-19 resource centre - including this research content - immediately available in PubMed Central and other publicly funded repositories, such as the WHO COVID database with rights for unrestricted research re-use and analyses in any form or by any means with acknowledgement of the original source. These permissions are granted for free by Elsevier for as long as the COVID-19 resource centre remains active.



# Metagenomic evidence for the co-existence of SARS and H1N1 in patients from 2007-2012 flu seasons in France



Qi Liu<sup>a,b,c,1</sup>, Zhenglin Du<sup>b,c,d,1</sup>, Sihui Zhu<sup>b,c,d</sup>, Wenming Zhao<sup>b,c,d</sup>, Hua Chen<sup>a,b,c,\*</sup>, Yongbiao Xue<sup>b,c,d,\*</sup>

<sup>a</sup> CAS Key Laboratory for Genomic and Precision Medicine, Beijing Institute of Genomics, Chinese Academy of Sciences, Beijing 100101, China

<sup>b</sup> China National Center for Bioinformatics, Beijing 100101, China

<sup>c</sup> University of Chinese Academy of Sciences, Beijing 100049, China

<sup>d</sup> National Genomics Data Center, Beijing Institute of Genomics, Chinese Academy of Sciences, Beijing 100101, China

## ARTICLE INFO

### Article history:

Received 15 September 2021

Revised 29 October 2021

Accepted 3 November 2021

Available online 9 November 2021

### Keywords:

Metagenomics

SARS-CoV

Influenza A virus

Retrospective study

## ABSTRACT

By re-analyzing public metagenomic data from 101 patients infected with influenza A virus during the 2007–2012 H1N1 flu seasons in France, we identified 22 samples with SARS-CoV sequences. In three of them, the SARS genome sequences could be fully assembled out of each. These sequences are highly similar (99.99% and 99.70%) to the artificially constructed recombinant SARS-CoV (SARSr-CoV) strains generated by the J. Craig Venter Institute in the USA. Moreover, samples from different flu seasons have different SARS-CoV strains, and the divergence between these strains cannot be explained by natural evolution. Our study also shows that retrospective studies using public metagenomic data from past major epidemic outbreaks serve as a genomic strategy for the research of the origins or spread of infectious diseases.

© 2021 Chinese Medical Association Publishing House. Published by Elsevier BV. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

## 1. Introduction

Genome sequencing has been used to identify the pathogen, trace virus origin, and provide outbreak surveillance for infectious disease studies. For example, the availability of the complete genome of the severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2) in January 2020 sped up the identification of the pathogen and facilitated the development of effective vaccines. With the accumulation of extensive metagenomic sequence data in the past years, one potential application is to carry out genome-based retrospective studies on major historical outbreaks to understand the occurrence and development of viral epidemics.

## 2. Materials and Methods

### 2.1. Data collection

The metagenomic data of 101 patients infected by the influenza A virus were downloaded from the NCBI SRA database (project ID: PRJEB11406). These samples were collected by the National Influenza

\* Corresponding authors: National Genomics Data Center, Beijing Institute of Genomics, Chinese Academy of Sciences, Beijing 100101, China (Yongbiao Xue); CAS Key Laboratory for Genomic and Precision Medicine, Beijing Institute of Genomics, Chinese Academy of Sciences, Beijing 100101, China (Hua Chen).

E-mail addresses: [chenh@big.ac.cn](mailto:chenh@big.ac.cn) (H. Chen), [ybxue@big.ac.cn](mailto:ybxue@big.ac.cn) (Y. Xue).

<sup>1</sup> These authors contributed equally to this work.

Center near Paris, France, between 2007 and 2012, spanning five consecutive flu seasons. Sequencing data of these samples were submitted by Institute Pasteur, France, in 2018 (Table S1) [1].

### 2.2. Variant calling and consensus sequence generating

After removing sequencing adapters and trimming consecutive low-quality bases from the 5' and 3' read ends using cutadapt [2], clean reads were mapped to the SARS-CoV genome (NC\_004718.3) using BWA (V0.7.12) [3] with default parameters. Next, the Picard program (<http://picard.sourceforge.net>) was used to sort mapping results to BAM format and mark duplicates of PCR amplification. Then GATK (V4.1.6.0) [4] was used for SNP and indel calling. Finally, consensus sequences were generated by applying VCF variants to the reference sequence using bcftools (v1.9) [5].

### 2.3. Sequence alignment, phylogenetic and network analysis

Two hundred and fifty genomes of SARS-CoV with the genome size larger than 290,000 bases were downloaded from the NCBI Nucleotide database (<https://www.ncbi.nlm.nih.gov/nucleotide>). We then constructed a multiple sequence alignment of the 250 genomes and 3 assembled genomes, using the MAFFT v7.453 with parameter “--auto” [6], and the final alignment contains 30,327 nucleotides.

Neighbor-joining (NJ) phylogenetic trees of the 253 genome sequences were constructed using MEGA X 10.1.8 with a maximum

composite likelihood model and the default parameter settings [7]. In addition, phylogenetic relationships and mutations that occurred among unique genomes were further inspected from 253 genomes through median-joining networks [8], using the Network 10 (<http://www.fluxus-engineering.com/>) to examine changes of genetic variations across places and through times. For network analysis, an 81-bp block at the 5'-end including gaps and a 77-bp block at the 3'-end including gaps and the poly-A tails were trimmed out of the alignment, and the final alignment contains 30,169 nucleotides.

#### 2.4. Pairwise sequence alignment

We used the BLAST online tools with default parameters (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>) to align two sequences.

### 3. Results

We re-analyzed public metagenomic data from 101 patients infected with influenza A virus (H1N1), collected by the Pasteur Institute in France between 2007 and 2012, spanning five consecutive flu seasons [1]. In 22 (21.78%) out of the 101 patient samples, we identified genomic fragments of SARS-associated coronavirus (SARS-CoV) with different proportions (0.0003% – 0.6127%) (Table S1). For 3 samples with sequencing depth >30 and genome coverage >99% (ERR1091908, ERR1091910, and ERR1091914), we were able to obtain the assembled sequences using the variant calling method (Table S1, Fig. S1). These 22 samples have different numbers of mutations ranging from 10 to 116 compared with the SARS-CoV reference genome (NC\_004718.3, see Table S1). In addition, the samples collected during the same flu season share similar mutations, while samples from different flu seasons possess different sets of mutations (Fig. S2, Table 1), suggesting that distinct SARS-CoV strains existed in different flu seasons in France.

SARS-CoV once caused an outbreak of SARS in 2002, a life-threatening respiratory infectious disease [9], but disappeared in human populations after 2003. To investigate the origin of these SARS-CoV sequences detected in the patients infected with influenza A virus between 2007 and 2012, the three assembled SARS-CoV genome sequences (ERR1091908, ERR1091910, and ERR1091914) were pooled together with 250 complete SARS-CoV genomes downloaded from the NCBI database. We constructed a phylogeny of these sequences with the neighbor-joining approach and a haplotype network with the median-joining approach, respectively (Fig. 1). We found that ERR1091908 and ERR1091910 are clustered with the SARS-CoV ExoN1 strain (colored in blue in Fig. 1) while ERR1091914 is clustered with the SARS-CoV wtic-MB strain (colored in black in Fig. 1), consistent with the finding of differential SARS-CoV mutations in the patient samples from different flu seasons. Note that SARS-CoV ExoN1, SARS-CoV wtic-MB, SARS-CoV MA15, and SARS-CoV MA15 ExoN1 all belong to recombinant SARSr-CoV, a group of SARS-CoV sequences artificially constructed using the same infectious clone (ic) recombinant virus strain of SARS-CoV Urbani (AY278741) that was initially isolated from a patient with SARS [10–12]. Moreover, the three newly assembled SARS-CoV sequences along all the recombinant SARSr-CoV sequences are separated from other naturally occurring SARS-CoV sequences, including AY278741 (colored in yellow and cyan in Fig. 1) in the phylogeny and haplotype network. Therefore, the above-said results indicate that all the three SARS-CoV sequences (ERR1091908, ERR1091910, and ERR1091914) are more likely to be categorized as artificially constructed recombinant SARSr-CoV strains.

ERR1091914 is almost identical (99.99%) to SARS-CoV wtic-MB (19 identical sequences in the data), with only one base pair different (bp) (R10626A, referred to the genome position of FJ882938.1) after

trimming the 20-bp at the 5'-end and 44-bp at the 3'-end of ERR1091914 (Table 1). Here, R represents A or G base. There are a total of 15 mutation differences between ERR1091914 and AY278741, the closest naturally occurring SARS-CoV, and the 15 mutations are shared by all the recombinant SARSr-CoV, indicating that ERR1091914 is indeed derived from the SARSr-CoV wtic-MB sequences instead of evolving independently from naturally occurring SARS-CoV sequences.

The other two newly identified SARS-CoV sequences (ERR1091908 and ERR1091910) are almost identical with only one base difference (99.99%). Both sequences are highly similar (99.70%) to the SARS-CoV ExoN1 sequence FJ882941.1, with 88 base differences after trimming 21-bp at the 5'-end and 65-bp at 3'-end (Fig. S3, Table 1). Annotation results showed that 84 of the 88 bases are located at coding regions of SARS-CoV and result in 45 amino acid (AA) substitutions. Most of these mutations are in the coding regions of open reading frame 1a (ORF1a) (13 substitutions), ORF1ab (7 substitutions), and the spike (S) glycoprotein (9 substitutions, Table 2), among which ORF1ab and S glycoprotein are functionally related to the viral replication, transmission and pathogenicity [13,14], suggesting potential gain-function effects of these mutations in ERR1091908 and ERR1091910. SARS-CoV ExoN1 strains are known to have a 21-fold increase in mutation rate during replication in previous research [12], which is consistent with the fact that the sequences within the SARS-CoV ExoN1 clades of the haplotype network are highly divergent compared with other recombinant SARSr-CoV clades (Fig. S3). Assuming a mutation rate of  $1.0 \times 10^{-3}$  per site per year for naturally evolving SARS-CoV and the SARS-CoV genome length of 29,751 bp, only 1.69 months (50.7 days) are needed to generate 88 site differences between the two SARS-CoV ExoN1 sequences, indicating continuous evolution of ERR1091908 and ERR1091910 after their recent divergence from FJ882941.1.

The sequencing data of 101 patients were submitted by the Pasteur Institute in France, which established a laboratory with a level-three biosafety standard and conducted research on SARS-CoV [15–17]. All the 116 recombinant SARSr-CoV sequences were submitted by the J. Craig Venter Institute (JCVI) from Tennessee, USA. The two institutes have collaborated and published their work on viral genome sequence [21].

### 4. Discussion

One possible explanation of the co-existence of SARS and H1N1 sequences in the patients is that the artificially constructed recombinant SARSr-CoV caused a co-infection outside the laboratory during 2007 - 2012, but did not result in the SARS-CoV epidemic; an alternative hypothesis is a contamination of the samples in the lab since the Pasteur Institute also conducted SARS-CoV studies. Moreover, samples from different flu seasons have different strains of SARS-CoV, and the divergence between these SARS-CoV strains cannot be explained by natural evolution. This intriguing finding warrants further efforts to sleuth out the culprit. In 2014, the Pasteur Institute France once lost vials containing patient samples collected during SARS (<https://www.sciencemag.org/news/2014/05/frances-institut-pasteur-under-fire-over-missing-sars-vials>). It raises a serious concern about laboratory biosafety in both institutions. Our study also shows that retrospective studies using public metagenomic data from past major epidemic outbreaks serve as a useful genomic strategy for the research of the origins or spread of infectious diseases.

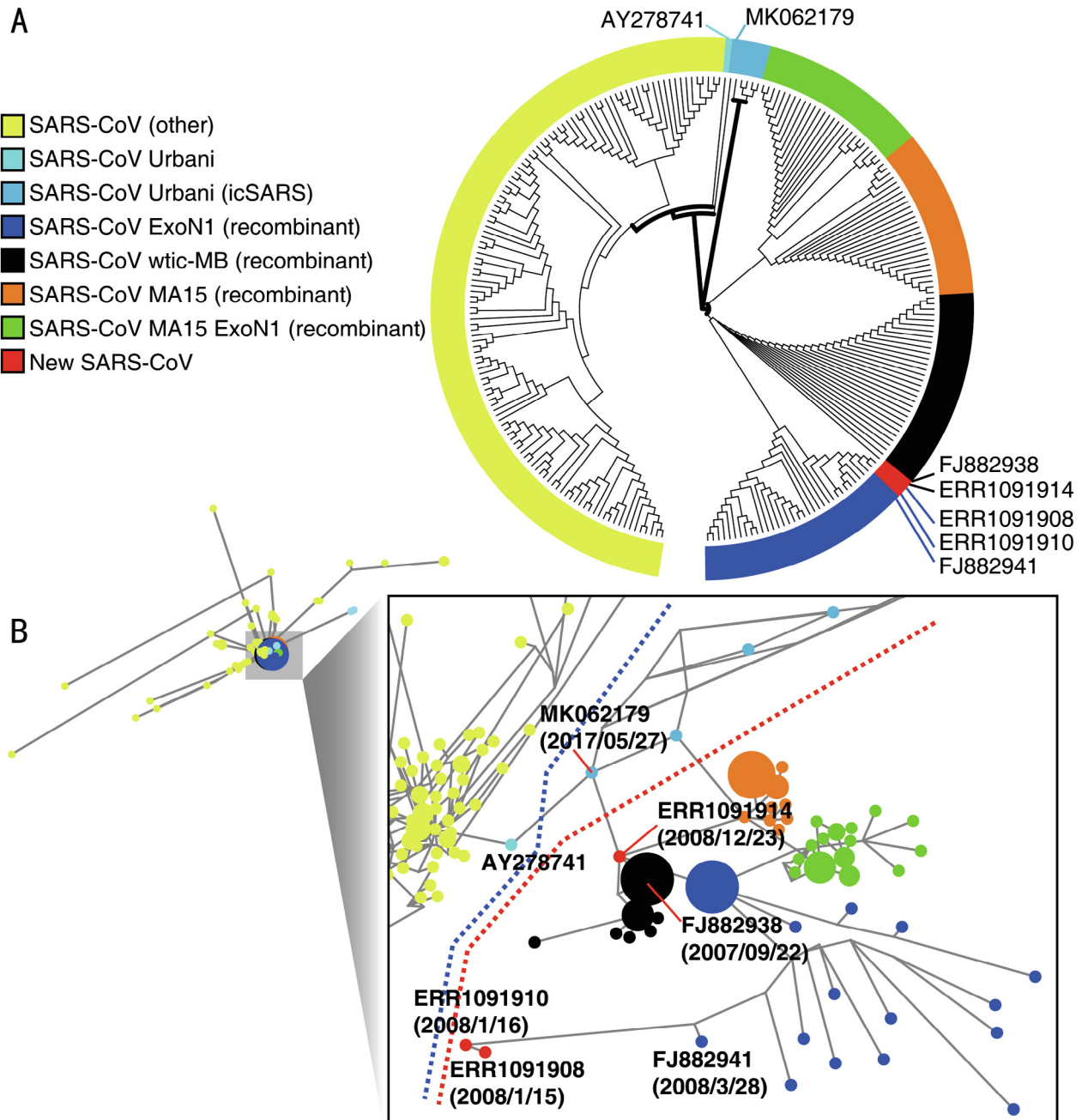
### Acknowledgements

This project was supported by the National Key R&D Program of China (2021YFC0863400), the Key Program of Chinese Academy of

**Table 1**  
Summary of the three SARS-CoV sequences and their closest sequences.

Cluster	Sample ID	Sample collection location	Collection date	Organism	Submitter	Identity with SARS-CoV sequences (%)
1	ERR1091914	Haute Normandie, France	2008/12/23	Influenza A virus	Institute Pasteur	99.99
	FJ882938*	Tennessee, USA	2007/9/22	SARS-CoV wtic-MB	J. Craig Venter Institute	
2	ERR1091908	Lorraine, France	2008/1/15	Influenza A virus	Institute Pasteur	99.70
	ERR1091910	Picardie, France	2008/1/16	Influenza A virus		
	FJ882941	Nashville, Tennessee, USA	2008/3/28	SARS-CoV ExoN1	J. Craig Venter Institute	

\*One sequence closed to ERR1091914 is shown and the remaining 18 sequences closed to ERR1091914 are presented in Table S2.



**Fig. 1.** Phylogeny tree and haplotype network of 253 SARS-CoV genome sequences. (A) A phylogeny of the 253 SARS-CoV genome sequences is constructed with the neighbor-joining approach. (B) Haplotype network of the 253 SARS-CoV genome sequences is constructed with the median-joining approach. Size of each circle represents the number of identical sequences.



**Table 2**

The positions of 84 base differences between ERR1091908 and FJ882941.1

ORF	CDS (position in NC_004718.3)	Mutations found in ERR1091908 compared to FJ882941.1		Function [ref]
		Positions of nucleotide change (number)	Positions of amino acid change in SARS-CoV protein (number)	
ORF 1a	265-13413	654, 707, 1771, 1905, 2976, 3229, 3491, 3603, 3845, 4731, 4808, 5015, 5061, 5236, 5412, 6087, 6265, 6459, 6476, 7484, 8004, 8922, 10119, 10658, 12411, 13149 (26)	148, 503, 989, 1076, 1194, 1489, 1515, 1584, 1658, 2001, 2071, 2407, 3465 (13)	Involved in viral replication and transcription, and virus pathogenesis [13,18]
ORF 1ab	265-21485	13874, 13925, 14178, 14630, 14876, 15497, 15605, 15740, 15821, 15905, 16356, 16386, 17269, 17602, 18238, 18239, 18244, 18245, 18749, 18860, 19082, 19814, 19917, 20528, 20555, 20789, 21038 (27)	987, 997, 1291, 1402, 1614, 1616, 2174 (7)	
S	21492-25259	21860, 22206, 22352, 22423, 23243, 23374, 23468, 23518, 23823, 24249, 24873, 24910, 24957 (13)	239, 311, 628, 676, 778, 920, 1128, 1140, 1156 (9)	Associated with cell entry of SARS-CoV and viral transmission [14,19]
ORF 3a	25268-26092	25550, 25626, 25783, 25800, 26049 (5)	120, 178, 261 (3)	Playing roles in virus uptake and release, viral-related apoptosis, and formation of viral envelope [20]
ORF 3b	25689-26153	25783, 25800, 26049, 26121 (4)	32, 38, 121, 145 (4)	Involving in immunomodulation, and acting as interferon antagonist [20]
E	26117-26347	26121, 26226, 26241, 26335 (4)	2, 37, 42 (3)	A small integral membrane proteins with roles in virus morphogenesis, assembly, budding, and replication [18]
M	26398-27063	NA	NA	NA
ORF 6	27074-27265	27167, 27248 (2)	32, 59 (2)	Acting as a $\beta$ -interferon antagonist and contribute to virulence [20]
ORF 7a	27273-27641	27290, 27639 (2)	123 (1)	Involving in virus-host interaction and contribute to SARS-CoV pathogenesis [20]
ORF 7b	27638-27772	27639, 27648 (2)	1, 4 (2)	A potential attenuating factor [20]
ORF 8a	27779-27898	NA	NA	Potential roles in the host
ORF 8b	27864-28118	27917 (1)	NA	ubiquitin-proteasome system [20]
N	28120-29388	28557, 29271, 29324 (3)	402 (1)	Playing role in virus replication and transcription, and acting as an interferon antagonist [18,19]
ORF 9b	28130-28426	NA	NA	Inducing caspase-dependent apoptosis [20]
Total number	NA	84	45	NA

Sciences (KJZD-SW-L14), and the National Natural Science Foundation of China (Grant No. 31571370 and 91731302).

### Conflict of interest statement

The authors declare that there are no conflicts of interest.

### Author contributions

**Qi Liu:** Formal Analysis, Visualization, Writing - Original Draft. **Zhenglin Du:** Investigation, Data Curation, Visualization, Writing - Original Draft. **Sihui Zhu:** Formal Analysis. **Wenming Zhao:** Investigation. **Hua Chen:** Conceptualization, Project Administration, Methodology, Writing - Review & Editing. **Yongbiao Xue:** Conceptualization, Supervision.

### Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.bsheal.2021.11.002>.

### References

- [1] I. Pelletier, D. Rousset, V. Enouf, et al, Highly heterogeneous temperature sensitivity of 2009 pandemic influenza A(H1N1) viral isolates, northern France, Euro. Surveill. 16 (2011), 19999. <https://doi.org/10.2807/ese.16.43.19999-en>.
- [2] M. Martin, Cutadapt removes adapter sequences from high-throughput sequencing reads, EMBnet J. 17 (1) (2011) 10, <https://doi.org/10.14806/ej.17.110.14806/ej.17.1.200>.
- [3] H. Li, R. Durbin, Fast and accurate short read alignment with Burrows-Wheeler transform, Bioinformatics 25 (14) (2009) 1754–1760, <https://doi.org/10.1093/bioinformatics/btp324>.
- [4] A. McKenna, M. Hanna, E. Banks, A. Sivachenko, K. Cibulskis, A. Kernytzky, K. Garimella, D. Altshuler, S. Gabriel, M. Daly, M.A. DePristo, The Genome Analysis Toolkit: A MapReduce Framework for Analyzing Next-Generation DNA Sequencing Data, Genome Res. 20 (9) (2010) 1297–1303, <https://doi.org/10.1101/gr.107524.110>.
- [5] P. Danecek, J.K. Bonfield, J. Liddle, J. Marshall, V. Ohan, M.O. Pollard, A. Whitwham, T. Keane, S.A. McCarthy, R.M. Davies, H. Li, Twelve years of SAMtools and BCFtools, GigaScience 10 (2021), giab008. <https://doi.org/10.1093/gigascience/giab008>.
- [6] K. Katoh, K. Misawa, K.-i. Kuma, T. Miyata, MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform, Nucl. Acids Res. 30 (2002) 3059–3066, <https://doi.org/10.1093/nar/gkf436>.
- [7] S. Kumar, G. Stecher, M. Li, C. Knyaz, K. Tamura, MEGA X: Molecular Evolutionary Genetics Analysis across Computing Platforms, Mol. Biol. Evol. 35 (2018) 1547–1549, <https://doi.org/10.1093/molbev/msy096>.
- [8] H.J. Bandelt, P. Forster, A. Rohl, Median-joining networks for inferring intraspecific phylogenies, Mol. Biol. Evol. 16 (1) (1999) 37–48, <https://doi.org/10.1093/oxfordjournals.molbev.a026036>.
- [9] T. Ksiazek, G., et al, A Novel Coronavirus Associated with Severe Acute Respiratory Syndrome, N. Engl. J. Med. 348 (2003) 1953–1966, <https://doi.org/10.1056/NEJMoa030781>.
- [10] B. Yount, K.M. Curtis, E.A. Fritz, L.E. Hensley, P.B. Jahrling, E. Prentice, M.R. Denison, T.W. Geisbert, R.S. Baric, Reverse genetics with a full-length infectious cDNA of severe acute respiratory syndrome coronavirus, Proc. Nat. Acad. Sci. U.S.A. 100 (22) (2003) 12995–13000, <https://doi.org/10.1073/pnas.1735582100>.
- [11] A. Roberts, D. Deming, C.D. Paddock, A. Cheng, K. Subbarao, A Mouse-adapted SARS-coronavirus causes disease and mortality in BALB/c mice, PLoS Pathog. 3 (1) (2007), e5. <https://doi.org/10.1371/journal.ppat.0030005>.
- [12] L.D. Eckerle, M.M. Becker, R.A. Halpin, et al., Infidelity of SARS-CoV Nsp14-Exonuclease Mutant Virus Replication Is Revealed by Complete Genome Sequencing, PLoS Pathog. 6 (5) (2010), e1000896. <https://doi.org/10.1371/journal.ppat.1000896>.
- [13] R.L. Graham, J.S. Sparks, L.D. Eckerle, A.C. Sims, M.R. Denison, SARS coronavirus replicase proteins in pathogenesis, Virus Res. 133 (1) (2008) 88–100, <https://doi.org/10.1016/j.virusres.2007.02.017>.

- [14] M. Bolles, E. Donaldson, R. Baric, SARS-CoV and emergent coronaviruses: viral determinants of interspecies transmission, *Curr. Opin. Virol.* 1 (6) (2011) 624–634, <https://doi.org/10.1016/j.coviro.2011.10.012>.
- [15] O.D. Mantke, H. Schmitz, Z. Herve, P. Heyman, A. Papa, M. Niedrig, Quality assurance for the diagnostics of viral diseases to enhance the emergency preparedness in Europe, *Euro. Surveill.* 10 (2005) 102–106, <https://doi.org/10.2807/esm.10.06.00545-en>.
- [16] Y.L. Yap, X.W. Zhang, Andonov., Structural analysis of inhibition mechanisms of Aurintricarboxylic Acid on SARS-CoV polymerase and other proteins, *Comput. Biol. Chem.* 29 (3) (2005) 212–219, <https://doi.org/10.1016/j.compbiolchem.2005.04.006>.
- [17] A. Fontanet, Cross-species transmission: last obstacle before pandemic, *Transfus. Clin. Bio.* 14 (2007) 16, <https://doi.org/10.1016/j.traci.2007.04.012>.
- [18] S. Perlman, J. Netland, Coronaviruses post-SARS: update on replication and pathogenesis, *Nat. Rev. Microbiol.* 7 (6) (2009) 439–450, <https://doi.org/10.1038/nrmicro2147>.
- [19] R. Hilgenfeld, M. Peiris, From SARS to MERS: 10 years of research on highly pathogenic human coronaviruses, *Antiviral Res.* 100 (1) (2013) 286–295, <https://doi.org/10.1016/j.antiviral.2013.08.015>.
- [20] D.X. Liu, T.S. Fung, K.K. Chong, A. Shukla, R. Hilgenfeld, Accessory proteins of SARS-CoV and other coronaviruses, *Antiviral Res.* 109 (2014) 97–109, <https://doi.org/10.1016/j.antiviral>.
- [21] M. Eppinger, M.J. Rosovitz, W.F. Fricke, et al, The complete genome sequence of *Yersinia pseudotuberculosis* IP31758, the causative agent of Far East Scarlet-Like Fever, *PLoS Genet.* 3 (8) (2007) e142, <https://doi.org/10.1371/journal.pgen.0030142>.