# Sequence Analysis and Structural Prediction of the Severe Acute Respiratory Syndrome Coronavirus nsp5

Jia-Hai LU[#]*, Ding-Mei ZHANG[#], Guo-Ling WANG[#], Zhong-Min GUO[#], Juan LI[#], Bing-Yan TAN, Li-Ping OU-YANG, Wen-Hua LING, Xin-Bing YU, and Nan-Shan ZHONG[1]

*The Public Health School of Sun Yat-Sen University, Guangzhou 510080, China;*
[1] *Guangzhou Institute of Respiratory Diseases, Guangzhou Medical College, Guangzhou 510120, China*

**Abstract**      The non-structural proteins (nsp or replicase proteins) of coronaviruses are relatively conserved and can be effective targets for drugs. Few studies have been conducted into the function of the severe acute respiratory syndrome coronavirus (SARS-CoV) nsp5. In this study, bioinformatics methods were employed to predict the secondary structure and construct 3-D models of the SARS-CoV GD strain nsp5. Sequencing and sequential comparison was performed to analyze the mutation trend of the polymerase nsp5 gene during the epidemic process using a nucleotide-nucleotide basic local alignment search tool (BLASTN) and a protein-protein basic local alignment search tool (BLASTP). The results indicated that the nsp5 gene was steady during the epidemic process and the protein was homologous with other coronavirus nsp5 proteins. The protein encoded by the nsp5 gene was expressed in COS-7 cells and analyzed by sodium dodecylsulfate-polyacrylamide gel electrophoresis (SDS-PAGE). This study provided the foundation for further exploration of the protein's biological function, and contributed to the search for anti-SARS-CoV drugs.

**Key words**      severe acute respiratory syndrome (SARS); severe acute respiratory syndrome coronavirus (SARS-CoV); non-structural protein (nsp); replicase; secondary structure; 3-D structure

The severe acute respiratory syndrome (SARS) was first identified in Guangdong Province, southern China, and spread to 32 countries and regions in the world. A novel identified virus named the SARS coronavirus (SARS-CoV) has been established as the etiological agent [1]. There are three groups of coronaviruses: groups I and II contain mammalian viruses, while group III contains only avian viruses [2]. On the basis of unrooted phylogenetic analysis, SARS-CoV was proposed to represent a novel virus and was defined as a fourth group of coronaviruses [3,4]. But later, Snijder *et al.* suggested that the SARS-CoV lineage was an early split-off from the group II branch by rooted phylogenetic analysis of

coronavirus replicase genes [5]. The complete genome sequence of SARS-CoV is now available from GenBank. With the interpretation of the whole SARS-CoV gene sequence, the genes associated with SARS-CoV's function and etiology have become more and more familiar to us [3,4]. The SARS-CoV genome is a positive strand RNA of approximately 29,700 nucleotides and has a 5′ cap structure and 3′ poly(A) tail [4–6]. It is composed of at least 14 functional open reading frames (ORFs) that encode three classes of proteins: structural proteins (the S, M, E and N proteins), 16 non-structural proteins (the nsp or replicase proteins) involved in viral RNA synthesis, and proteins that are thought to be non-essential for replication in tissue culture but clearly provide a selective advantage *in vivo* (accessory proteins) [3,4,6,7].

The coronavirus replicases are encoded by two large 5′-proximal ORFs, which comprise approximately two-

DOI: 10.1111/j.1745-7270.2005.00066.x

third of the genome polyproteins. ORF1a and ORF1b are connected by a ribosomal frame shift site, and the ribosomal frame shift results in the translation of an ORF1a protein and a carboxyl-extended ORF1ab frame shift protein, which are also known as replicase polyproteins PP1a and PP1ab [4]. The ORF1a and ORF1ab translation products are polyprotein precursors, which are cleaved by viral proteinases, resulting in a minimum of 13 non-structural proteins, including a 3C-like proteinase (nsp5), an RNA-dependent RNA polymerase (nsp12, RdRp), a papain-like proteinase activity (PL2[pro], nsp3), and a superfamily 1-like helicase activity (HEL1, nsp13) [6] (following the nomenclature used by Snijder *et al.* [5]), and other function-unknown non-structural proteins. These replicase polyproteins play important roles in controlling the replication of viral genome RNA, the transcription of the sgmRNA, and the synthesis and translation of the mRNA of structural proteins (S, M, E, N) [4,8]. Structural proteins have been selected as candidates for anti-SARS drugs, but there is a chance of mutation at these proteins, leading to the inefficacy of anti-SARS drugs. The replicase proteins are relatively conserved and can be effective targets for anti-SARS drugs. In order to search for the targets of anti-SARS drugs and study the biological function of SARS-CoV nsp5, we carried out the following research. We compared the nsp5 of the SARS-CoV GD strain (isolated from Guangdong patient) with that of other SARS-CoV strains using a nucleotide-nucleotide basic local alignment search tool (BLASTN) to find out the variance tendency. In the same way we compared the amino acid sequence with that of other coronaviruses using a protein-protein basic local alignment search tool (BLASTP) to find the homologous protein. We know that the biological function of the gene is achieved by protein encoding, but the functions of the protein are closely related to its structure. In the process of biological evolution a protein's structure is more conserved than its sequence, so it is more useful for the study of the functions and mechanisms of a protein to master the protein structure information. To extract more information, we constructed the secondary structure and 3-D models of the nsp5 protein. To further study the function of nsp5, we cloned the nsp5 gene of SARS-CoV, isolated by molecular biological methods from a SARS patient in Guangdong. We constructed the eukaryotic expression plasmid and successfully expressed the nsp5 proteins. At present, few studies focused on the polymerase nsp5, hardly enough to elucidate its biological function. To deal with a possible recurrence of SARS, the determination of targets in the SARS-CoV is impor-

tant [7,8]. It provides the foundation for further exploration of the biological functions of this protein during transmission, and assists in the search for anti-SARS-CoV drugs.

## Materials and Methods

### SARS-CoV GD strain

The SARS-CoV GD strain was derived from a 47-year-old female who had returned from Hong Kong infected with SARS-CoV on 1 April, 2003. She was admitted to hospital on 9 April of the same year.

### Viral culture for SARS-CoV GD strain and extraction of SARS-CoV genomic of RNA

The virus was routinely cultivated in Vero E6 cells. The titer of virus was measured at approximately $10^{6.75}$ $TCID_{50}$/ml. The RNA was extracted according to the instructions of the Viral RNA extraction kit (Qiagen, Hilden, Germany) [11,12].

### Amplification of nsp5 gene by RT-PCR

The cDNA was reverse transcribed with random primers by RT-PCR (reverse transcription-polymerase chain reaction) kit (TaKaRa, Dalian, China) using the RNA extracted as the template. Reverse transcription was performed at 42 °C for 45 min, then at 95 °C for 5 min. The nsp5 gene is located within the nucleotides 12,022–12,615 of the SARS-CoV GD strain. The primers were designed according to the sequence of the SARS-CoV strain in Taipei (GenBank accession No. AY291451). The forward primer was 5′-GCCGTCGACATGGCTATT-GCTT-3′. The reverse primer was 5′-GCGCGGCCGCT-TACTGTAGTTTA-3′. GTCGAC is the *Sal*I site and GCGGCCGC is the *Not*I site. An aliquot (2 μl) of the cDNA product was amplified by PCR. PCR reactions were as follows: 95 °C for 1 min; 94 °C for 30 s, 50 °C for 50 s, and 72 °C for 1 min, 37 cycles; 72 °C for 10 min. Subsequently, PCR products were purified according to the instructions of the gel extraction kit (TaKaRa).

### Sequence analysis, prediction of secondary structure and 3-D model

The gene sequence was sent to Shanghai Genecore Bio Technologies Company Limited for DNA sequencing. The gene sequence and amino acid sequence were then compared with other homologous genes and amino acids by BLASTN and BLASTP to analyze the variance tendency

of nsp5. To extract more information about the nsp5 protein, we predicted the secondary structure using the Protein Structure Prediction Server (PSIPRED) [13]. A 3-D model of nsp5 was constructed using the Protein Fold Recognition (Threading) Server (3D-PSSM) [14]. RasMol and Chime plug-ins (http://www.umass.edu/microbio/rasmol/index.html) were used to view the 3-D models.

### Construction and identification of eukaryotic recombinant expression vector pCI-neo-nsp5

After the PCR products were purified, the PCR products and the expression vector pCI-neo (Promega, Madison, USA) were both digested with *Sal*I and *Not*I (TaKaRa). The ligation was performed between the target gene and the pCI-neo expression vector by T4 DNA ligase (TaKaRa), which was named pCI-neo-nsp5. Double restriction digestion and DNA sequencing was performed to identify the positive clone.

### Expression of the nsp5 protein and SDS-PAGE

The COS-7 cells (Promega) were cultured in Dulbecco's modified Eagle's medium (DMEM) supplemented with 10% fetal bovine serum in a humidified atmosphere of 5% $CO_2$ at 37 °C. The recombinant expression plasmid pCI-neo-nsp5 was transfected into COS-7 cells using lipofectin reagent according to the manufacturer's instructions and incubated for 72 h. The transfected cell clones (COS-7/pCI-neo-nsp5) were selected by the addition of G-418 (Geneticin, 800 μg/ml) and cultured under appropriate conditions. The cells were digested by 0.25% pepsin and washed with phosphate-buffered saline (PBS) three times, then sonicated. After sonication, the mixture was centrifuged at 10,000 *g* for 15 min at 4 °C. The supernatant and culture supernatant of COS-7 cells were both concentrated with PEG20000. The products were resolved by 7% sodium dodecyl-sulfate-polyacrylamide gel electrophoresis (SDS-PAGE) and stained with Coomassie brilliant blue (Sigma, USA).

## Results

### Sequence analysis, secondary structure and 3-D model

We compared the GD strain nsp5 sequence with those of 110 SARS-CoV strains obtained from the National Center for Biotechnology Information (NCBI) using BLASTN. The results demonstrated that the GD strain was completely homologous with 105 SARS-CoV strains (data not shown); only 5 of the 110 strains, when compared with the GD strain, had undergone mutation. As shown in **Table 1**, the GD01 strain had one mutation at the 129th base site, but the amino acid and secondary structure did not change. The HSZ-A strain had one mutation at the 267th base site. In ZJ01, mutation occurred at the 496th and the 497th base sites, which resulted in changes to the amino acid. ShanghaiQXC1 was completely in accord with ShanghaiQXC2, in which the t was substituted by c and V changed into A, but the secondary structure did not change (**Table 1**). The BLASTP result indicated that the nsp5 protein of the SARS-CoV GD strain is: 37% identical with replicase polyprotein 1ab of human coronavirus NL63 and ORF1a of human group I coronavirus associated with pneumonia; 35% identical with putative coronavirus nsp5 of human coronavirus 229E, coronavirus nsp5 of bovine coronavirus and RNA-directed RNA polymerase of murine hepatitis virus; 34%

**Table 1          Variance of SARS-CoV nsp5 protein**

| | Mutated base position | | | | | | | | | | | |
| | 129 | | | 267 | | | 496/497 | | | 584 | | |
| Strain | N | AA | S | N | AA | S | N | AA | S | N | AA | S |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| GD | aat | N | H | aca | T | H | att | I | H | gtt | V | H |
| GD01 | aa**c** | N | H | aca | T | H | att | I | H | gtt | V | H |
| HSZ-A | aat | N | H | ac**t** | T | H | att | I | H | gtt | V | H |
| ShanghaiQXC1 | aat | N | H | aca | T | H | att | I | H | g**c**t | *A* | H |
| ShanghaiQXC2 | aat | N | H | aca | T | H | att | I | H | g**c**t | *A* | H |
| ZJ01 | aat | N | H | aca | T | H | **ca**t | *H* | H | gtt | V | H |

AA, amino acid; N, nucleotide; H, helix; S, secondary structure. Bold letters, the nucleotide variation at the corresponding position of nsp5; italicized letters, the single nucleotide variations (SNVs) causing amino acid non-synonymous variations. GenBank accession No.: GD01, AY278489; HSZ-A, AY394984; ShanghaiQXC1, AY463059; ShanghaiQXC2, AY463060; ZJ01, AY297028.

identical with putative coronavirus nsp5 of porcine epidemic diarrhea virus; and 30% identical with replicase polyprotein 1a of avian infectious bronchitis virus (data not shown). The secondary structure of nsp5 of the GD strain is shown in **Fig. 1**. The 3-D model was constructed by 3D-PSSM using 1LKJ [PDB (Protein Data Bank) code] as the template (**Figs. 2** and **3**). The 3-D structure of 1LKJ is also shown in **Fig. 4**.



**Fig. 2  SARS-CoV GD strain nsp5 protein model**
Each amino acid is displayed in a specific color: pink, Val; blue, Gln; green, Ala; white, Gly; orange, Thr; red, Ser; brown, Met.



**Fig. 1  Predicted secondary structure of nsp5 protein**
Conf, confidence of prediction; Pred, predicted secondary structure; AA, target sequence.



**Fig. 3  N and C termini of SARS-CoV GD strain nsp5 protein**
The protein was colored according to the position of the amino acid in the molecule chain, from blue to green, to yellow, to orange, to red. Blue, C terminus of the protein; red, N terminus of the protein.



**Fig. 4  3-D structure of 1LKJ protein**
Red, alpha helices; yellow, beta strands; blue, turns; white, other amino acids.

## The identification of RT-PCR products and recombinant expressive plasmids pCI-neo-nsp5

As expected, a 594 bp gene fragment encoding the whole nsp5 protein of SARS-CoV was derived from total RNA by RT-PCR (**Fig. 5**, lane 1). The product of PCR was successfully cloned into pCI-neo vector, identified by double restriction digestion (*Sal*I and *Not*I) (**Fig. 5**, lane 4).
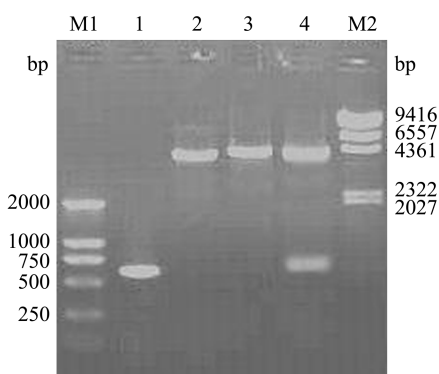


**Fig. 6      SDS-PAGE analysis of nsp5 protein from transfected COS-7 cells**
M, protein marker; 1, culture supernatant of COS-7 cells transfected by pCI-neo; 2, culture supernatant of COS-7 cells; 3, supernatant of the sonicated COS-7 cells transfected by pCI-neo; 4, culture supernatant of COS-7 cells transfected by pCI-neo-nsp5; 5, supernatant of the sonicated COS-7 cells transfected by pCI-neo-nsp5.



**Fig. 5      Identification of nsp5 gene PCR product and restriction enzyme analysis of pCI-neo-nsp5**
M1, DNA marker, DL2000; 1, PCR of the aimed gene (594 bp); 2, *Sal*I digestion of pCI-neo; 3, *Sal*I digestion of pCI-neo-nsp5; 4, *Sal*I and *Not*I digestion of pCI-neo-nsp5; M2, DNA marker, λ/*Hin*dIII digest.

## Expression of the nsp5 protein

The recombinant expression vector was successfully transfected into COS-7 cells using lipofectin reagent. The cells carrying the nsp5 gene were obtained and the nsp5 protein was successfully expressed. The expressed protein was analyzed by SDS-PAGE (**Fig. 6**). The protein was approximately 21 kDa, and the study of its function is underway.

## Discussion

The SARS-CoV genome is approximately 29.7 kb long and composed of at least 14 functional ORFs which encode the S, N, M, E and other proteins [4]. The S, N, M and E proteins have been successfully expressed and purified [2–4]. The SARS outbreak has inspired a myriad of studies into virt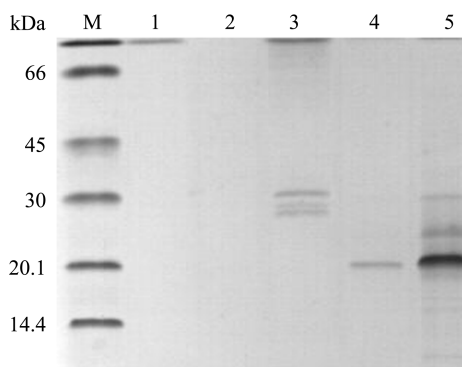ually every aspect of SARS-CoV biology, including viral pathogenesis, tissue tropism, genome structure expression and replication, as well as SARS-CoV structural and nonstructural proteins.

According to the knowledge of coronaviruses, the replicase polyproteins play important roles in controlling the replication of viral genome RNA, the transcription of the sgmRNA, and the synthesis and translation of the mRNA of structural proteins (S, M, E, N) [15]. In this study the nsp5 gene of the SARS-CoV GD strain was compared with other SARS-CoV strains obtained from NCBI to find the variance tendency. The results demonstrated that the GD strain was completely homologous with 105 strains and only five strains had mutations (**Table 1**). The table shows that the GD01 strain had one mutation at the 129th base site, but the amino acid and secondary structure did not change. We know that GD01 was the early case in the transmission. The result suggested that the 129th base site mutated in the middle and late period of transmission, but the corresponding amino acid and secondary structure did not change. The HSZ-A strain, which was also submitted by Guangzhou, had one mutation at the 267th base site, also without amino acid or secondary structure change. In ZJ01, submitted 12 May 2003 in Zhejiang Province, mutation occurred at the 496th and 497th base sites, which resulted in change to the amino acid. ShanghaiQXC1, submitted 11 November 2003, was completely in accord with ShanghaiQXC2, in which the t was substituted by c, and V changed into A, but the secondary structure did not change. From the 110

SARS-CoV strains, there are only five mutations at five base sites in five strains, two amino acid substitutions and no structural change. The results indicated that nsp5 is very stable during the transmission and evolutionary processes, suggesting it is an appropriate target for anti-SARS drugs. The mutations may also be important for adaptation to host conditions, and the mechanism of these changes still needs to be explored. The secondary structure of nsp5 of the GD strain is shown in **Fig. 1**. Clearly it is mostly composed of helices. The BLASTP result indicated that the nsp5 protein of the SARS-CoV GD strain is homologous with the nsp5 protein of other coronaviruses. We can base our exploration of the function of the SARS-CoV nsp5 protein, and our design of anti-SARS drugs, on other coronavirus nsp5 proteins. Interestingly, the nsp5 protein of the SARS-CoV GD strain possessed the highest identification of replicase polyprotein 1ab of human coronavirus NL63, which was also identified in January 2003 [16]. Further exploration is needed to establish whether a relationship exists between them.

To gain further insight, we constructed a rudimentary 3-D model of the nsp5 protein by 3D-PSSM, using 1LKJ as the template. This method determines the possible fold type of a particular protein using the folds of a different protein, whose structure is already known to us, as a template. This model is generated directly from the alignment, with no subsequent refinement. The results revealed the three proteins with the highest PSSM *E*-value: 1LKJ (PDB code), which is a chain of calmodulin from the yeast *Saccharomyces cerevisiae*, with an *E*-value of 4.8; 1ABV, which is the delta subunit of the $F_1F_o$-ATP synthase, with an *E*-value of 7.44; 1J1EC, which is troponin, with an *E*-value of 7.51. The 3-D structure of 1LKJ is also shown in **Fig. 4**. The results demonstrated that 1LKJ is mostly composed of 8 alpha helices. The secondary structure of GD nsp5 also demonstrated that it is mostly made up of alpha helices, which indicated that it was reasonable to construct the 3-D model of the SARS-CoV nsp5 protein using 1LKJ as the template. Calmodulin is a ubiquitous $Ca^{2+}$-binding regulatory protein that mediates the action of $Ca^{2+}$ in a diverse array of biological events, including cell division, microtubular depolymerization and cyclic nucleotide metabolism, $Ca^{2+}$ transport, and protein phosphorylation [17]. The identification between nsp5 and calmodulin led us to consider whether nsp5 had a similar function to that of calmodulin, which needed further exploration. The results also showed the protein in the fold library with the highest sequence identity with nsp5 was 1UD0 (PDB code, heat-shock-like protein a chain). As reported above, 3D-PSSM is a method of constructing a 3-D model of a submitted protein using the protein in the fold library with the most homologous folds. The models generated by 3D-PSSM are simple mappings from the coordinates of the template structure and the query sequence residues aligned to them. The *E*-value is the measure of confidence in the prediction, and the closer this score is to zero, the better the textual match. In this result the lowest *E*-value was 4.8, which indicated that certainty was less than 50%. Consequently the generated 3-D model of nsp5 is only rudimentary. An exact model can only be obtained through further study.

In this study, we constructed the eukaryotic recombinant expression plasmidpCI-neo-nsp5. The recombinant expression vector pCI-neo-nsp5 was transfected into COS-7 cells, and the protein was then expressed and analyzed by SDS-PAGE. The successful construction of the 3-D structure model, the preliminary function prediction and stable expression of nsp5 protein established the foundation for the further exploration of its actual biological function and contributed to the search for anti-SARS drugs.

## Acknowledgements

## References

1  Enserink M. Infectious diseases. Second suspect in the global mystery outbreak. Science 2003, 299: 1963

2  Gao F, Ou HY, Chen LL, Zheng WX, Zhang CT. Prediction of proteinase cleavage sites in polyproteins of coronaviruses and its application in analyzing SARS-CoV genomes. FEBS Lett 2003, 553: 451–456

3  Rota PA, Oberste MS, Monroe SS, Nix WA, Campagnoli R, Icenogle JP, Peñaranda S *et al*. Characterization of a novel coronavirus associated with severe acute respiratory syndrome. Science 2003, 300: 1394–1399

4  Marra MA, Jones SJ, Astell CR, Holt RA, Brooks-Wilson A, Butterfield YSN, Khattra J *et al*. The genome sequence of the SARS- associated coronavirus. Science 2003, 300: 1399–1404

5  Snijder EJ, Bredenbeek PJ, Dobbe JC, Thiel V, Ziebuhr J, Poon LLM, Guan Y *et al*. Unique and conserved features of genome and proteome of SARS-coronavirus, an early split-off from the coronavirus group 2 lineage. J Mol Biol 2003, 331: 991–1004

6  Sutton G, Fry E, Carter L, Sainsbury S, Walter T, Nettleship J, Berrow N *et al*. The nsp9 replicase protein of SARS-coronavirus, structure and functional insights. Structure 2004, 12: 341–353

7  Thiel V, Ivanov KA, Putics A, Hertzig T, Schelle B, Bayer S, Weiβbrich B *et al*. Mechanisms and enzymes involved in SARS coronavirus genome expression. J Gen Virol 2003, 84: 2305–2315

8  Tijms MA, Dinten LC, Gorbalenya AE, Snijder EJ. A zinc finger-containing papain-like protease couples subgenomic mRNA synthesis to genome translation in a positive-stranded RNA virus. Proc Natl Acad Sci USA 2001, 98: 1889–1894

9  Bermejo Martin JF, Jimenez JL, Munoz-Fernandez A. Pentoxifylline and severe acute respiratory syndrome (SARS): A drug to be considered. Med Sci Monit 2003, 9: 29–34

10 Anand K, Ziebuhr J, Wadhwani P, Mesters J R, Hilgenfeld R. Coronavirus main proteinase (3CL-Pro) structure: Basis for design of anti-SARS Drugs. Science 2003, 300: 1763–1767

11 Lu JH, Yan XG, Guo ZM, Zheng HY, Zhang X, Wan ZY, Zhang RL *et al*. Establishment of SARS virus vaccine line. Guangdong Med J (Chin) 2003, 24: 194–195

12 Zhang CH, Guo ZM, Zheng HY, Lu JH, Wang YF, Yan XG, Zhao Y *et al*. Humoral immune responses in rabbits induced by an experimental inactivated severe acute respiratory syndrome coronavirus vaccine prepared from F69 strain. Chin Med J 2004, 117: 1625–1629

13 McGuffin LJ, Bryson K, Jones DT. The PSIPRED protein structure prediction server. Bioinformatics 2000, 16: 404–405

14 Kelley LA, MacCallum RM, Sternberg MJE. Enhanced genome annotation using structural profiles in the program 3D-PSSM. J Mol Biol 2000, 299: 499–520

15 Lu Y, Denison MR. Determinants of mouse hepatitis virus 3C-like proteinase activity. Virology 1997, 230: 335–342

16 Ahmad K. New human coronavirus isolated. Lancet Infect Dis 2004, 4: 255

17 Davidkova G, Zhang SP, Nichols RA, Weiss B. Reduced level of calmodulin in PC12 cells induced by stable expression of calmodulin antisense RNA inhibits cell proliferation and induces neurite outgrowth. Neuroscience 1996, 75: 1003–1019

Edited by
**Bing SUN**