

Molecular evolution of the SARS-coronavirus during the course of the SARS epidemic in China

The Chinese SARS molecular epidemiology consortium

Supporting online materials

Materials and Methods

A. Epidemiological investigations

Official epidemiological records of the Guangdong Center for Disease Control and Prevention (GDCDCP), which represented an aggregate of the regular SARS epidemiology reports submitted by the local Centers for Disease Control and Prevention of individual cities, were reviewed. The contact and clinical histories of all of the early seemingly independent index cases and several key cases (*e.g.* HZS2-F) were reconfirmed either by review of hospital patient records or direct interview with the patients and/or the physicians-in-charge. In particular, eleven index cases from seven cities located in the Pearl River Delta region of Guangdong Province (Fig.1 and fig. S1), which occurred prior to the first superspread event of a Guangzhou hospital, HZS-2, were investigated in detail.

Majority of the specimens were collected by the virologists of GDCDCP, with the remaining samples collected by the staffs of local hospital or Guangzhou Center for Disease Control and Prevention.

B. Sequencing strategy and procedures

We isolated the viral RNA templates either from the culture supernatants of VeroE6 cells that showed cytopathic effects or directly from patients' specimens of SARS cases (including serum, stool, oropharyngeal swabs, nasal pharyngeal aspirates or autopsy lung tissues). RNA was extracted with the QIAamp viral RNA mini kit (Qiagen, Valencia, CA, USA) or TRIZOL Reagent (GIBCOBRL). The double-strand cDNA was synthesized with the SuperScript II cDNA system (Invitrogen, Carlsbad, CA, USA) or RNA PCR Kit (AMV) Ver.2.1 (Takara, Dalian China). To amplify the genomic sequences of the SARS-CoV, 53 sets of nested primers were designed based on the TOR2 sequence. The nested PCR fragments were directly sequenced in both forward and reverse directions on the ABI-3700 DNA sequencer (Applied Biosystems, Foster City, CA, USA) with 2- to 4-fold redundancy. For GZ02, PCR primers were designed to cover the whole genome in every 1kb interval with 200bp overlap with the adjacent fragment based on the TOR2 sequence. PCR products were sequenced using ABI BigDye Terminator Cycle Sequencing Kit on ABI-377. All of the nucleotide sequence variations of GZ02, which differ from that of the human SARS-CoV sequences available at GenBank as of June 2003, in particular, TOR2 and GZ01 (the sequence of an independent viral isolate from the same patient as GZ02 and currently renamed as GD01) sequences (including the 29-bp segment), were re-sequenced from RNA extractions from the same lung tissue specimen of that patient and the 5' end sequence was completed. The PHRED/PHRAP/CONSED software (University of Washington, Seattle, WA, USA; <http://www.phred.org>) was used for base calling, assembly, and editing. The assembled genome sequence was checked manually for accuracy and the regions with poor quality were re-sequenced. For data analysis, the nucleotide coordinate of GZ02 was used as a reference.

C. List of GenBank accession numbers for sequences mentioned in the text and SOM:

Sequences generated by this study	GenBank accession number	Sequences previously available	GenBank accession number
GD03T13 (S gene)	AY525636	SZ16 (palm civet)	AY304488
		SZ3 (palm civet)	AY304486
GZ02	AY390556	GD01 (GZ01)	AY278489
HGZ8L1-A	AY394981		
HSZ-A	AY394984		
HSZ-B (b, c)	AY394985, AY394994		
HSZ-C (b, c)	AY394986, AY394995		
ZS-A	AY394997	gz43 (S gene)	AY304490
ZS-B	AY394996	gz60 (S gene)	AY304491
ZS-C	AY395003		
GZ-A	AY394977		
JMD	AY394988		
HGZ8L1-B	AY394982		
HZS2-A	AY394983	CUHK-W1	AY278554
HZS2-Bb	AY395004	BJ04	AY279354
HZS2-C	AY394992	BJ01	AY278488
HZS2-D	AY394989	BJ02	AY278487
HZS2-E	AY394990	BJ03	AY278490
HGZ8L-2	AY394993		
HZS2-Fc	AY394991		
HZS2-Fb	AY394987		
CUHK-LC1	AY394998		
GZ-B	AY394978	TOR2	AY274119
GZ-C	AY394979	ZJ01	AY297028
GZ-D	AY394980		
CUHK-LC2	AY394999	CUHK-AG01	AY345986
CUHK-LC3	AY395000	CUHK-AG02	AY345987
CUHK-LC4	AY395001		
CUHK-LC5	AY395002		

Supporting Online Text S1

Statistical analysis for (A) the estimation of the neutral mutation rate and the date for the most recent common ancestor (MRCA); (B) calculation of the average Ka/Ks for three coding sequences (S, Orf1a, Orf1b) of the SARS-CoV genome within the three epidemic phases

(A) Estimation of the neutral mutation rate and the date for the most recent common ancestor (MRCA)

(1) Selection of samples

Culture artifacts may potentially introduce apparent sequence variations in the SARS-CoV genome. Therefore, we used sequences that are derived directly from the patients' clinical specimens for the present statistical analysis. Sequences generated from specimens collected more than 4 weeks after disease onset were also excluded. Among all of the available sequences, 10 (GZ02, CUHK-AG01, CUHK-AG02, GZ-C, GZ-D, HZS2-A, HZS2-Fb, HSZ-A, HSZ-Bb, HSZ-Cb) met all the criteria. We used GZ02 as the out-group, since it is the most divergent from all of the remaining 9 sequences (see Fig. 2 in the text).

(2) Estimation of the neutral mutation rate

The Pamilo-Bianchi-Li model (*S1–2*) was used to calculate the Ks for the 6 known concatenated coding sequences (orfla, orflb, S, E, M, and N) of the SARS-CoV genome. We plotted the Ks versus the sampling dates, defined as the number of days away from January 1, 2003 (fig S5).

There is a positive correlation between the Ks and the sampling dates (correlation

coefficient: 0.82). Thus, the synonymous substitution rate appears to be relatively constant throughout the sampling period. As there are less selective constraints for synonymous mutations in general, we proposed to use the synonymous substitution rate per site per day for the concatenated coding sequences to represent the neutral mutation rate of the SARS-CoV genome.

A simple linear regression model was used to estimate the neutral base substitution rate. The rate was estimated from the slope (β_1) of the fitted linear regression line and is found to be $8.26 \times 10^{-6} (\pm 2.16 \times 10^{-6}) \text{ nt}^{-1}\text{day}^{-1}$. This estimated mutation rate is quite similar to the values obtained for other known RNA viruses (S3–4). Specifically, the rate is about one third that for human immunodeficiency virus (S3).

The relationship can thus be described as, $Y = \beta_0 + \beta_1 X$ where Y is the Ks (using GZ02 as out-group) and X is the sampling date, which is measured by the number of days away from January 1, 2003.

(3) Estimation of the date of MRCA for the available samples

The intercept (β_0) of the fitted line is 1.055×10^{-3} , which equates to a sampling date of 0 and thus, corresponded to a calendar date in the end of year 2002. We used the GZ02 sequence as the out-group, which was sampled on February 11, 2003 (*i.e.*, 42 days after January 1, 2003). If T denotes the number of days before January 1, 2003 for the occurrence of the MRCA, then

$$T = \frac{\hat{\beta}_0 / \hat{\beta}_1 - 42}{2} = 46(\text{days}), \text{ which is equivalent to mid-November of 2002. The 95\% confidence}$$

interval for T is estimated to be 5.5 - 201.5 days (Supporting References and Notes: Appendix 1) (S5), meaning that the date for the MRCA is estimated to range from early June, 2002 to end of December, 2002. We know from the phylogenetic tree generated from all of the SARS-CoV sequences compared in this study (refer to Fig. 2 of the text and fig. S6) that the ancestor node of

the 10 samples we used to estimate the neutral mutation rate is the same as that for all of the other available human SARS-CoV sequences. Thus the estimated date of the MRCA should be applicable to all of the studied human SARS-CoV sequences.

(B) Calculation of the average Ka/Ks for three coding sequences (S, Orf1a, Orf1b) of the SARS-CoV genome within the three epidemic phases

The 61 human SARS-CoV genotypes were divided into three groups according to the three different phases of the epidemic (fig. S6: green for early phase, red and purple for middle phase and blue for late phase). The Ka/Ks for the S gene, Orf1a and Orf1b sequences were calculated for each group. The results are shown in table S3. All ratios were calculated in a similar manner and here we illustrate the method of calculation using the S gene. We included all the different S coding sequences in the calculation. For samples with identical sequences, one was randomly selected.

Within each group, Ka/Ks ratios were first calculated in a pairwise manner (2 sequences at a time) according to the Pamilo-Bianchi-Li model (S1-2). The average Ka/Ks and its standard error (table S3) were obtained using all these pairwise Ka/Ks values except those equal to infinite which is due to Ks=0. One-sided unpaired two-sample t-test was used to test whether the average Ka/Ks ratios for each of the studied coding region of the SARS-CoV sequences during the different epidemic phases were significantly different (table S3).

Supporting Online Text S2:

Acknowledgements

We appreciate the strong support from the Minister of Science and Technology of the Central Government of China, the Governments of Guangdong Province, Shanghai Municipality, and the Hong Kong Special Administrative Region. This work was supported by the State High Technology Development Program (863, Grant No. 2003AA208407), the State Key Program for Basic Research (973, Grant No. 2003CB514101), and a grant from the Guangdong Government for SARS research (2003FD02). RWKC, JST and YMDL are supported by the Hong Kong Research Grants Council Special Grant for SARS Research (CUHK4508/03M). HDS, GWZ, BWG, SJC, and ZC are partly supported by a SARS research grant from the BNP PARIBAS. PH, WZH and YXL are partly supported by the Shanghai Commission of Science and Technology. SYZ is supported by a special grant on SARS reservoir of the Institute of Zoology, Chinese Academy of Sciences, and a grant (30340035) from the National Science Foundation of China (NSFC). XNW and JLH are supported by the "973" grants (2001CB510008 and 2003CB514113) and the NSFC Fund for Distinguished Young Scholars. XJZ is supported by the "863" grants (2002AA214151 and 2003AA2Z347A). HT is supported by the NIH grant to C.I.W. We are grateful to the critical technical assistance supplied by Hui-Qiong ZHOU¹, Ping HUANG¹, Li-Mei DIAO¹, Qiu-Xia CHEN¹, Yan SHENG⁸, Yi CHEN⁸, Zheng RUAN⁸, Chun-Lei JIANG⁸, Yu LIU⁸, Wei SHEN¹⁴, and Li WANG¹⁴. We thank the Institute of Biochemistry and Cell Biology, SIBS, CAS for supplying VeroE6 cells. We thank Li Ruan, Wuchun Cao, and Ruifu Yang for sharing their unpublished information with us. We would also like to express our special appreciation to Shou-Yi Yu, Antione Danchin and Linfa Wang for helpful suggestions and comments on our research strategy and the manuscript preparation.

*All affiliations are the same as those listed in the main text.

Supporting References and Notes

Supporting references

- S1. P. Pamilo, N. O. Bianchi, *Mol. Biol. Evol.* **10**, 271 (1993).
- S2. W. H. Li, *J. Mol. Evol.* **36**, 96 (1993).
- S3. W. H. Li, M. Tanimura, P. M. Sharp, *Mol. Biol. Evol.* **5**, 313 (1988).
- S4. J. W. Drake, J. J. Holland, *Proc. Natl. Acad. Sci. U. S. A.* **96**, 13910 (1999).
- S5. J. A. Nelder, P. McCullagh, *Generalized Linear Models* (Chapman and Hall, London, ed. 2, 1989)
- S6. N. Saitou, M. Nei, *Mol. Biol. Evol.* **4**, 406 (1987).
- S7. M. Kimura, *J. Mol. Evol.* **16**, 111 (1980).

Supporting notes:

Appendix 1: Calculation of 95% confidence interval for T:

From the linear model, we obtained $\hat{\beta}_0 = 1.055 \times 10^{-3}$ and $\hat{\beta}_1 = 8.257 \times 10^{-6}$

Let's denote $\theta = \hat{\beta}_0 / \hat{\beta}_1$. So $\frac{\hat{\beta}_0 - \theta\hat{\beta}_1}{\hat{\sigma}\sqrt{Q(\theta)}} \sim t_{df=7}$, where $Q(\theta) = c_{11} - 2 \times \theta \times c_{12} + \theta^2 c_{22}$.

From the covariance matrix of the fitted linear model, we can get $\hat{\sigma}^2 \times c_{11} = 2.2658 \times 10^{-8}$, $\hat{\sigma}^2 \times c_{12} = -3.0558 \times 10^{-10}$ and $\hat{\sigma}^2 \times c_{22} = 4.7013 \times 10^{-12}$. So the 95% confidence of θ will satisfy,

$$\Pr\left\{ \left| \frac{1.055 \times 10^{-3} - \theta \times 8.257 \times 10^{-6}}{\sqrt{2.266 \times 10^{-8} + 2\theta \times 3.0558 \times 10^{-10} + \theta^2 \times 4.7013 \times 10^{-12}}} \right| < t_{7,0.025} \right\} = 0.05.$$

When we solved the inequality, we obtained the 95% confidence interval (CI) for θ ($\hat{\beta}_0 / \hat{\beta}_1$) which is (53, 445.1). Since T is equal to $\frac{\hat{\beta}_0 / \hat{\beta}_1 - 42}{2}$, the 95% CI for T will be estimated as (5.5, 201.55) (S5).

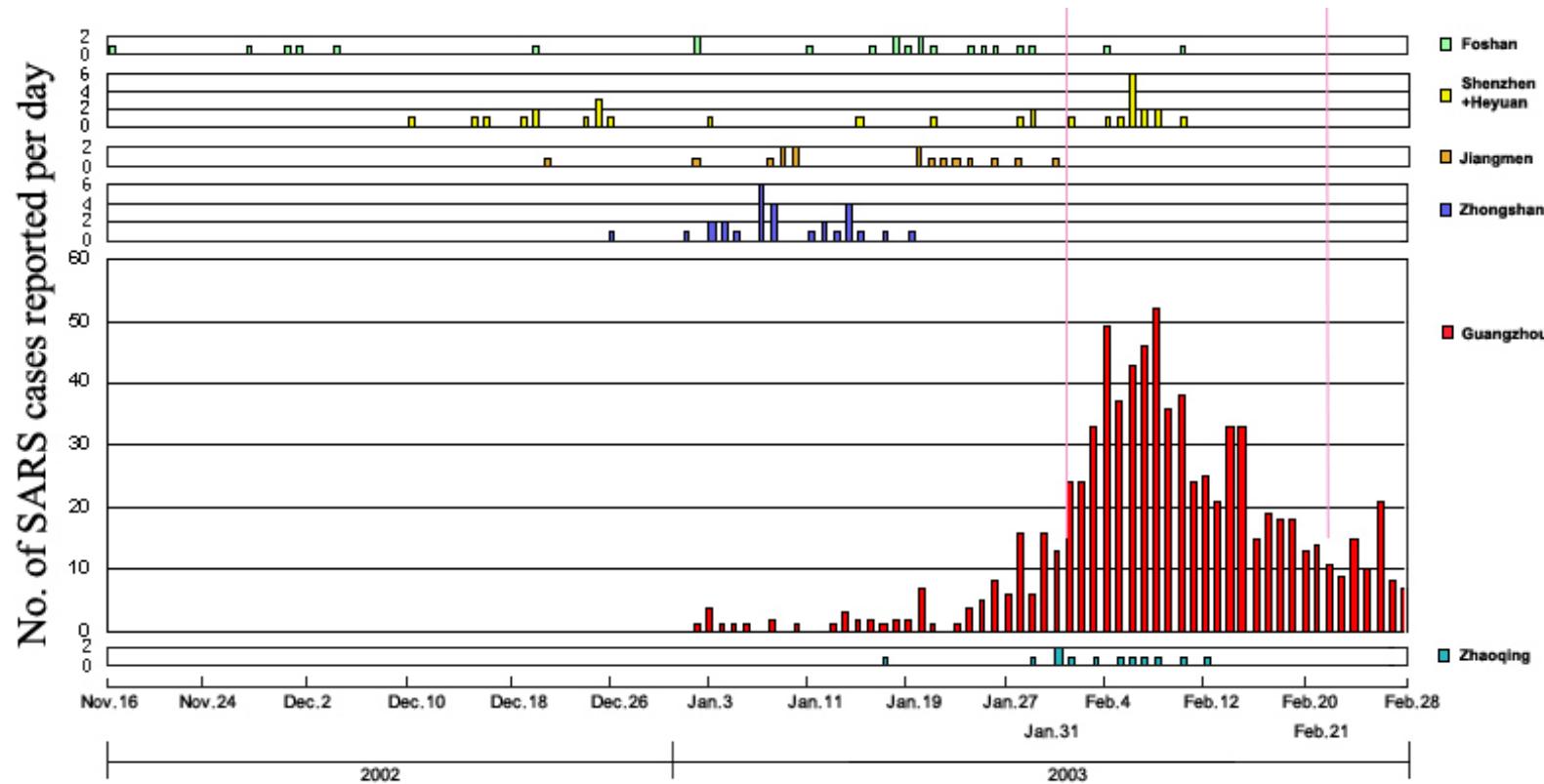


Fig. S1. Number of daily documented SARS cases reported from individual cities of the Guangdong Province, China, up to February 2003. Original epidemiological data were collected and analyzed by the Guangdong Center for Disease Control and Prevention. We combined the cases reported from the cities of Heyuan and Shenzhen because the Heyuan index case became infected in Shenzhen and after this nosocomial infection, no additional infections were reported in Heyuan. The order of the cities is arranged from top to bottom based on the disease onset date of their respective index cases, starting from the earliest to the latest dates of onset.

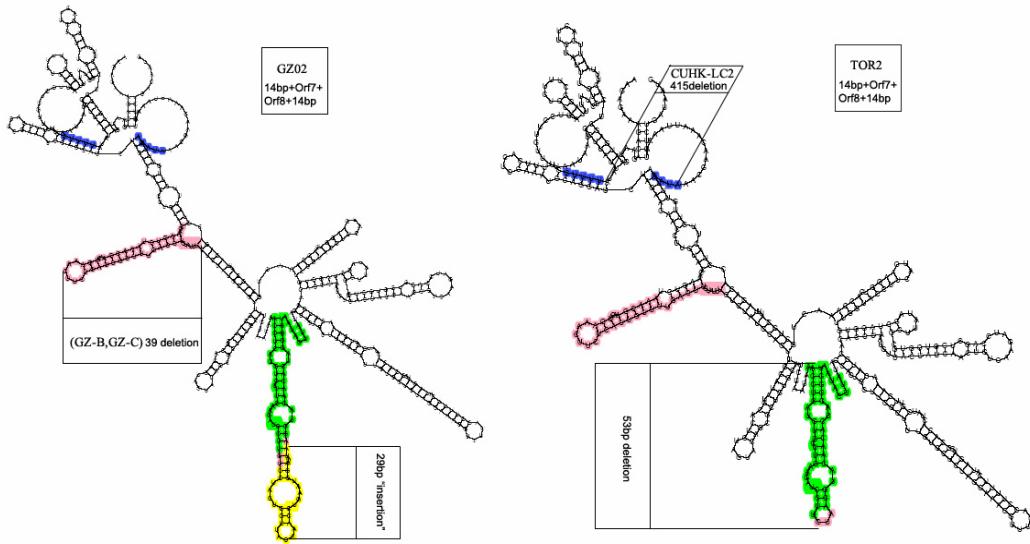
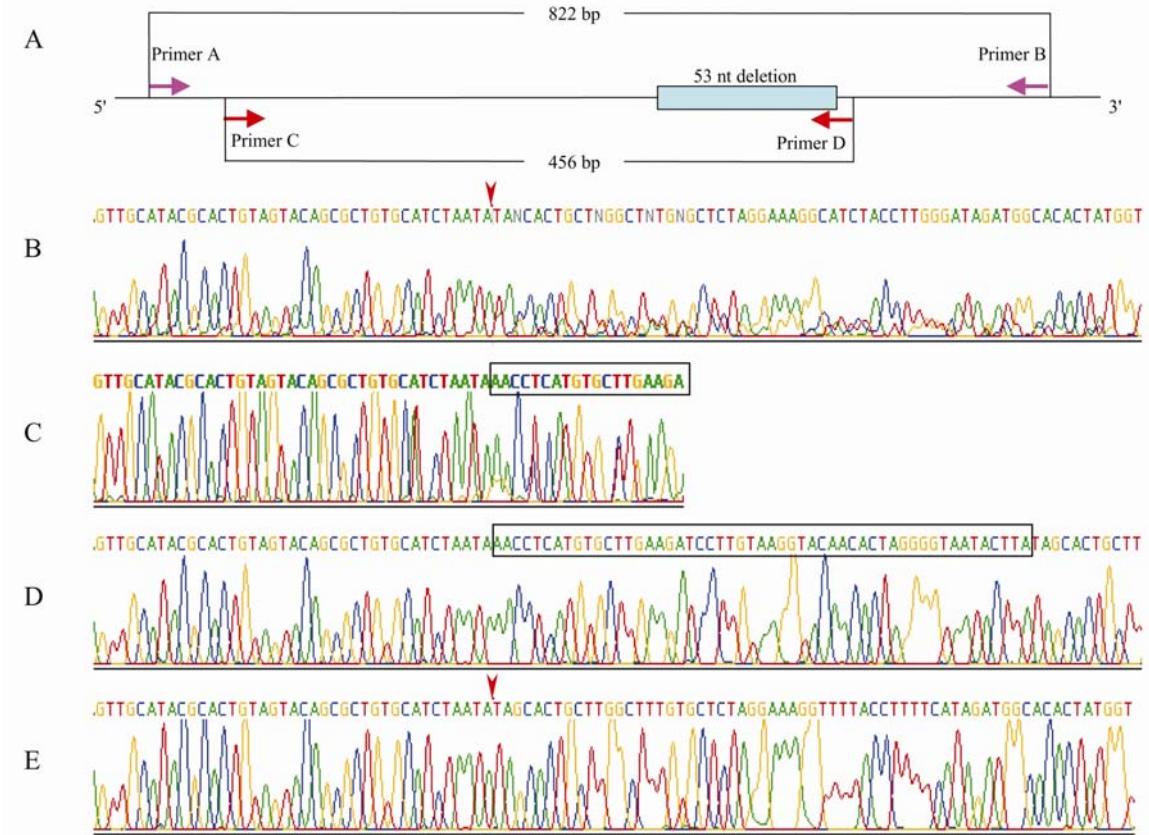
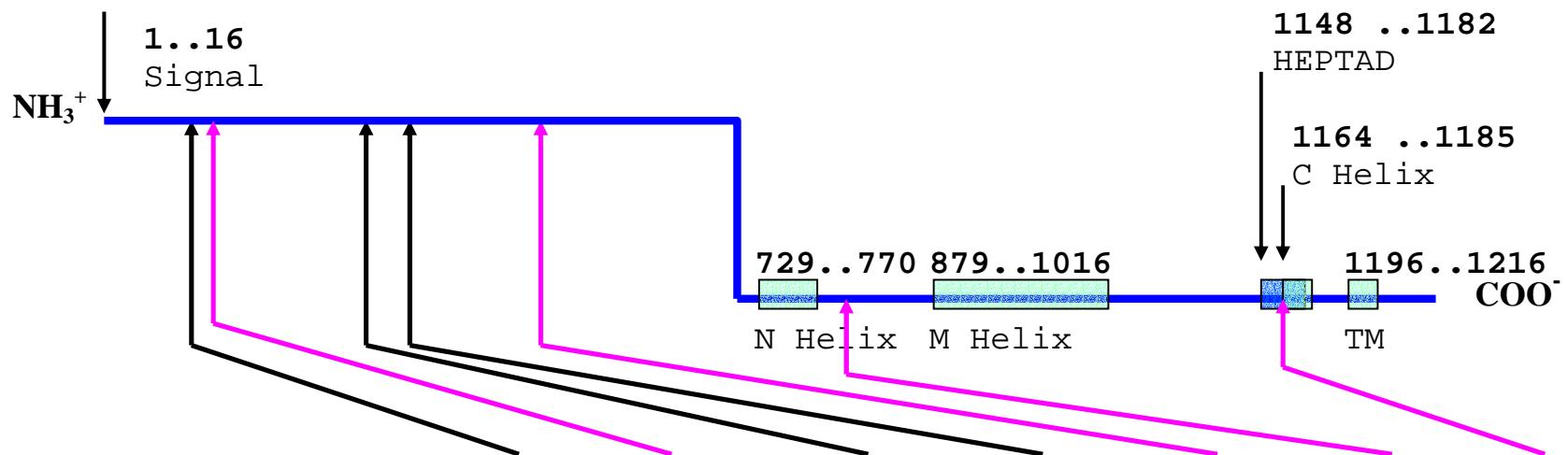


Fig. S2. Predicted RNA secondary structure of the Orf7b-Orf8 region of the SARS-CoV genome. SARS-CoV genotypic variations caused by major deletion events were observed on a number of occasions during the epidemic. All such deletions were confined to the Orf7b-Orf8 region. The genomic locations of the major deletions observed in this study are indicated on the predicted RNA secondary structures of the longest SARS-CoV genotype (left panel) and the genotype with the 29-nt deletion (right panel). The former genotype is represented by GZ02 while the latter is represented by TOR2. This latter genotype predominated the remainder of the epidemic from the middle phase onwards. For both panels, the illustrated region starts from 14 nucleotides upstream to the start of the predicted Orf7 to 14 nucleotides downstream to the end of Orf8. The illustrated region corresponds to nucleotide positions 27288 to 28161 on GZ02 and nucleotide positions 27259 to 28132 on TOR2. The prediction was made using the VIENNARNA:RNAfold software (<http://bioweb.pasteur.fr/>). GZ-B and GZ-C are two genotypes obtained from two Guangzhou patients with disease onset from mid-March but demonstrated a 39-nt deletion.





	GZ02 nucleotide coordinate	21715	21721	22207	22222	22422	23823	24978
Epidemic phase	S gene nucleotide coordinate	224	230	716	731	931	2,332	3,487
	S gene amino acid coordinate	75	77	239	244	311	778	1,163
Early	S gene codon (amino acid)	aGg (Arg)	gAc (Asp)	tTa (Leu)	aCt (Thr)	Aga (Arg)	Gat (Asp)	Gaa (Glu)
Late	S gene codon (amino acid)	aCg (Thr)	gGc (Gly)	tCa (Ser)	aTt (Ile)	Gga (Gly)	Tat (Tyr)	Aaa (Lys)

Fig. S4. The predicted amino acid residue alterations in the S protein caused by non-synonymous SNVs observed in the epidemic. The amino acid residue alterations predicted are listed and mapped to the approximate regions of the modeled S protein. Nucleotides 21721, 22222 and 23823 are the loci included in the 5-nt motif used for the classification of the major genotypes. These loci are shaded in the same colors as in Fig. 2 or Table S1.

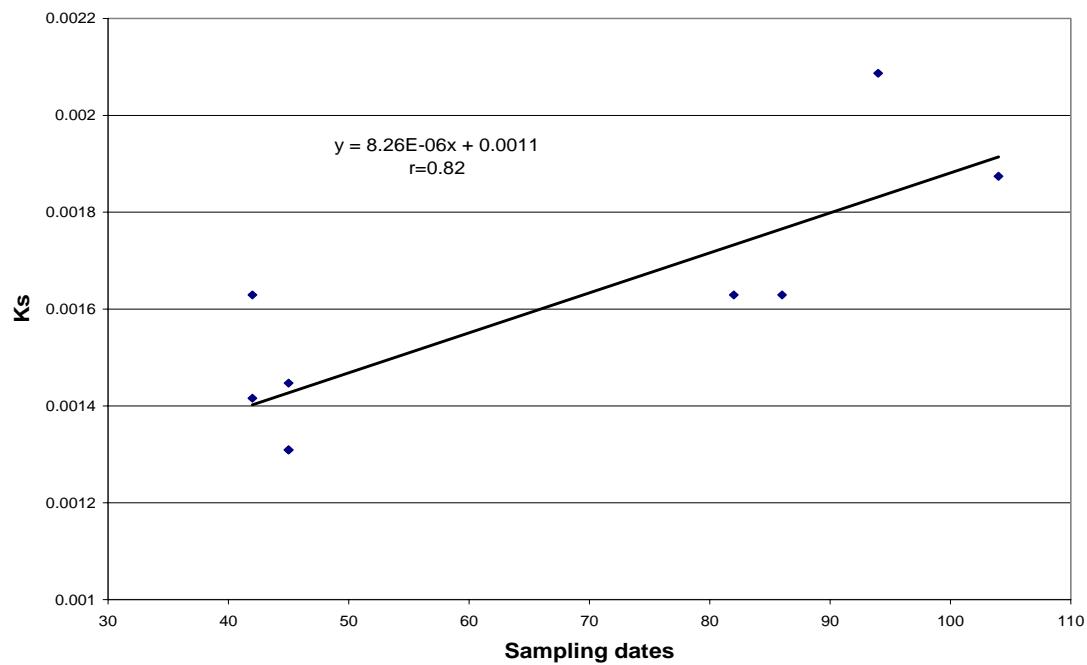


Fig. S5. Synonymous substitution rate, K_s , for the concatenated coding sequences versus sampling dates. The principles for sample selection and the statistical analysis methods are described above. The relatively most divergent sequence GZ02 was used as the out-group and K_s was determined for 9 representative human SARS-CoV sequences. The sampling dates are measured as the number of days away from January 1, 2003. A simple linear regression model was used to estimate the neutral base substitution rate.

Fig. S6: Neighbor-joining tree for the 60 SARS-CoV genomic sequences. ZJ01 was omitted for its low quality sequence and the tree was constructed using the Kimura two-parameter distance ($S = 7$).

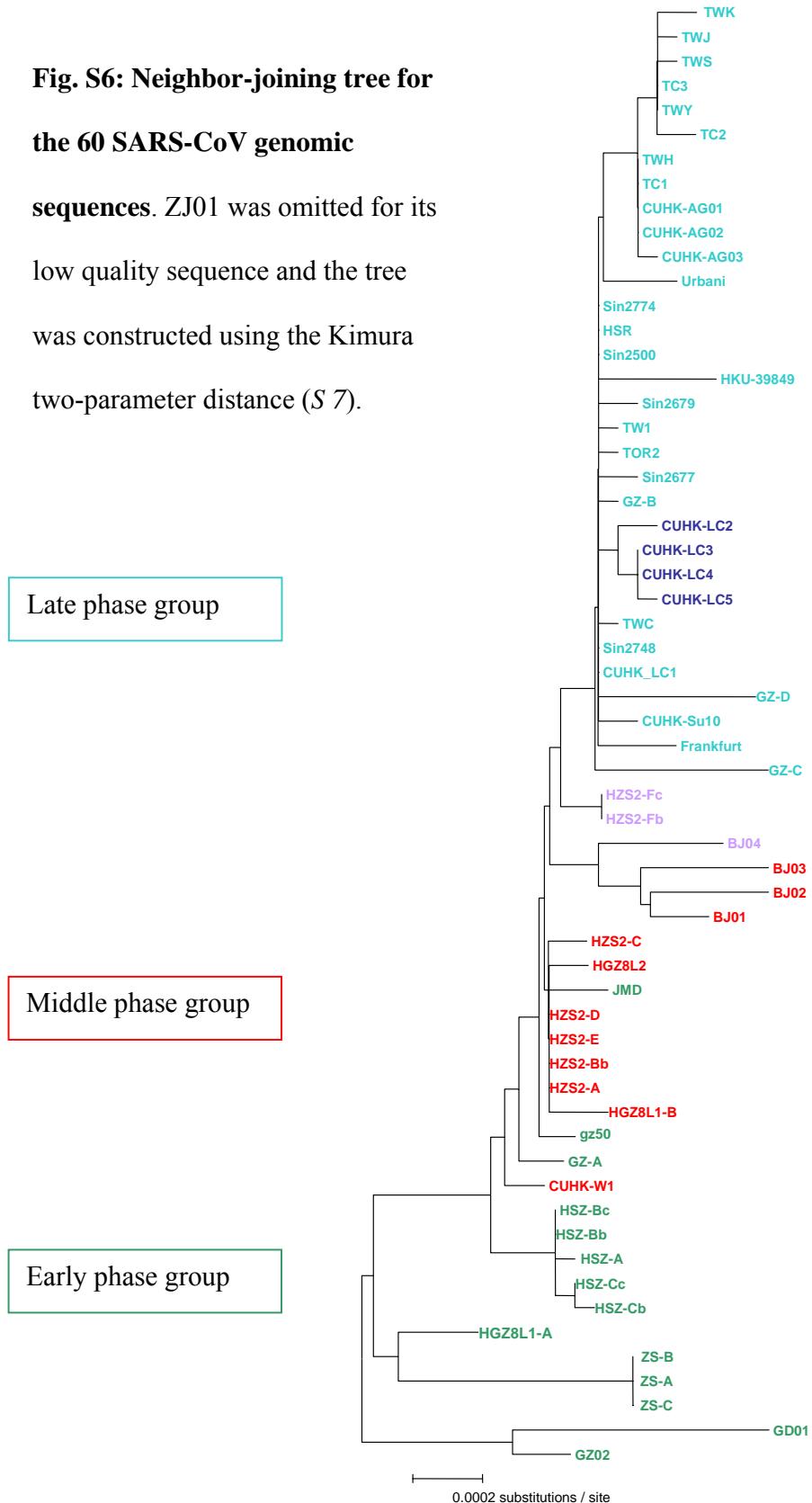


Fig. S7. Genotype clustering of the S genes from human SARS-CoV and palm civet SARS-like-CoV of the 2002/2003 epidemic and the 2003/2004 Guangdong index patient. The rooted phylogenetic tree for the nucleotide sequences of S genes from 2 palm civet SARS-like-CoV sequences (SZ16 and SZ3) and 64 human SARS-CoV sequences (61 as those used in the whole genome analysis of this project, in addition to two more S gene sequences, gz43 and gz60 of the 2002/2003 epidemic and the S gene sequence of the 2003/2004 Guangdong index patient, GD03T13 [Materials and Methods]). Only those variant sequences that were present in at least two independent samples were used for tree construction (total of 28 SNPs, synonymous and non-synonymous; Table S4). The map distance between individual sequences represents the extent of genotypic difference. A 6-nt motif that characterized the major phylogenetically-related genotypes are indicated in boxes. The nucleotide 24566, corresponding to the S gene nucleotide 3075 was not included in the characteristic motif because it only caused a synonymous variation. The sequences are named in concordance with their GenBank nomenclature.

S-nt	3739	C	C	C	C	C	C
SZ3	T	C	C	A	T	A	T
GD03T13	T	C	C	A	T	C	C
GZ02	T	C	C	A	C	T	C

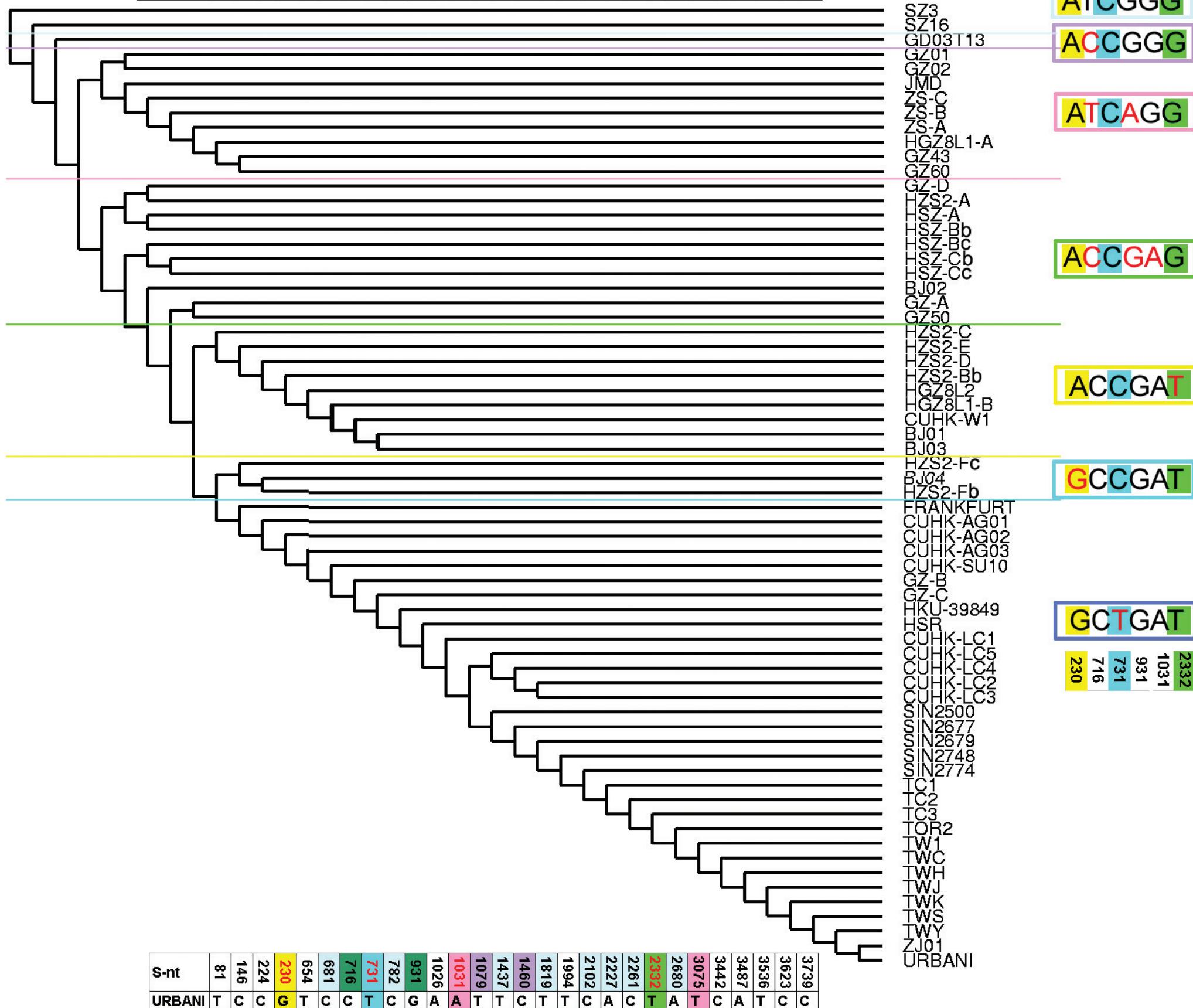


Table S1. The predicted coding sequence changes caused by the major deletion events in the Orf7b-Orf8 region of the SARS-CoV genome.

	Orf7b	Orf8a	Orf8b	Nucleocapsid Protein
GZ02	27635:27769 MNE TLI D FYLCFLA FLLFLV LIMLII FWF SLE I QD LEEPCTKV *	27776: 28144 (Sars8) MKLL I VLTC I SLCSC I RTVVQRCASNKPHV LEDPCPT GYQPEW I RYNTRGNTY STA WL CALGKVL P FHR W TM VQT CP PNVT I NCQDPAGG AL I ARC W Y L HEGH QTA AFR DVL V LNKRTN*		28146: 29414 sars9a nucleocapsid protein MSDNGPQS N QR S APR I T F GG G TD S DN N QN G RNGARP K QRRPQ G LP N NTAS W FTALTQHGKEELRFPRGQ G VP I NTNSGPDDQ I GY Y R RAT TR RV R GG D GM K M K LS P R W Y F Y L GT G PEAS L PY G AN K E G I V W V A T E G AL N TP K D H I G TR R NP N NNAA T V L QL P Q G T T LP K G F YAE G SR GGSQASS R SS R SG N SR N STPG S RG N SP A RM S GG G ET A LA L LL D RL N Q L E S K V SG K QQ Q Q Q TV T KK S AA E ASK K PR Q K R T A TK QYN V TA Q FG R RG P EQ T Q G N F GD Q DL I R Q G T D Y K W WP Q I A Q F AP S AS A FF G MS R I G ME V TP S GT W LT Y H G A I K L DD K D P Q F K D N V I L LN K H I DAY K TF P TE P K D KK KK KT D EA Q PL P Q R Q K PT V TL P PA D MD D FS R QL Q NS M SG A ST T QA*
TOR2 29 dels	27635:27769 MNE TLI D FYLCFLA FLLFLV LIMLII FWF SLE I QD LEEPCTKV *	27776..27924 sars8a unknown (sars8b) MCL KI LV Y NR T R G NT Y STA WL CALGKVL P FHR W HTM V T CT P NV T I N CCQDPAGG AL I ARC W Y L HEGH QTA FRD V L V V LNKRTN*	27861: 28144 sars8b unknown (sars8b) MCL KI LV Y NR T R G NT Y STA WL CALGKVL P FHR W HTM V T CT P NV T I N CCQDPAGG AL I ARC W Y L HEGH QTA FRD V L V V LNKRTN*	28146: 29414 sars9a nucleocapsid protein MSDNGPQS N QR S APR I T F GG G TD S DN N QN G RNGARP K QRRPQ G LP N NTAS W FTALTQHGKEELRFPRGQ G VP I NTNSGPDDQ I GY Y R RAT TR RV R GG D GM K M K LS P R W Y F Y L GT G PEAS L PY G AN K E G I V W V A T E G AL N TP K D H I G TR R NP N NNAA T V L QL P Q G T T LP K G F YAE G SR GGSQASS R SS R SG N SR N STPG S RG N SP A RM S GG G ET A LA L LL D RL N Q L E S K V SG K QQ Q Q Q TV T KK S AA E ASK K PR Q K R T A TK QYN V TA Q FG R RG P EQ T Q G N F GD Q DL I R Q G T D Y K W WP Q I A Q F AP S AS A FF G MS R I G ME V TP S GT W LT Y H G A I K L DD K D P Q F K D N V I L LN K H I DAY K TF P TE P K D KK KK KT D EA Q PL P Q R Q K PT V TL P PA D MD D FS R QL Q NS M SG A ST T QA*
HGZ8L1-B 82 dels	27635:27769 MNE TLI D FYLCFLA FLLFLV LIMLII FWF SLE I QD LEEPCTKV *	27776..27963 MKLL I VLTC I SLCSC I IR T V Q RCAS N I ALL GF V L*	*	28146: 29414 sars9a nucleocapsid protein MSDNGPQS N QR S APR I T F GG G TD S DN N QN G RNGARP K QRRPQ G LP N NTAS W FTALTQHGKEELRFPRGQ G VP I NTNSGPDDQ I GY Y R RAT TR RV R GG D GM K M K LS P R W Y F Y L GT G PEAS L PY G AN K E G I V W V A T E G AL N TP K D H I G TR R NP N NNAA T V L QL P Q G T T LP K G F YAE G SR GGSQASS R SS R SG N SR N STPG S RG N SP A RM S GG G ET A LA L LL D RL N Q L E S K V SG K QQ Q Q Q TV T KK S AA E ASK K PR Q K R T A TK QYN V TA Q FG R RG P EQ T Q G N F GD Q DL I R Q G T D Y K W WP Q I A Q F AP S AS A FF G MS R I G ME V TP S GT W LT Y H G A I K L DD K D P Q F K D N V I L LN K H I DAY K TF P TE P K D KK KK KT D EA Q PL P Q R Q K PT V TL P PA D MD D FS R QL Q NS M SG A ST T QA*
CUHK-LC2 415 dels	27635:27769 MNE TLI D FYLCFLA FLLFLV LIMLII FWF SLE I QD LEEPCTKV *	*	*	27635: 29414 MNE TLI D FYLCFLA FL FLV LIMLII I NE Q I K M SD N G P Q S NR S AP R I T F GG G TD S DN N QN G RNGARP K QRRPQ G LP N NTAS W FTALTQHGKEELRFPRGQ G VP I NTNSGPDDQ I GY Y R T Q H G KE E LR F PR G Q G VP I NT N SG P DD Q I GY Y R A TR R RV R GG D GM K M K LS P R W Y F Y L GT G PEAS L PY G AN K E G I V W V A T E G AL N TP K D H I G TR R NP N NNAA T V L QL P Q G T T LP K G F YAE G SR V GG G SQASS R SS R SG N SR N STPG S RG N SP A RM S GG G ET A LA L LL D RL N Q L E S K V SG K QQ Q Q Q TV T KK S AA E ASK K PR Q K R T A TK V SG G K Q Q Q Q Q TV T KK S AA E ASK K PR Q K R T A TK Q Y N V T Q A F G R R PE Q T Q G N F G Q D DL I R Q G T D Y K W WP Q I A Q F AP S AS A FF G MS R I G ME V TP S GT W LT Y H G A I K L DD K D P Q F K D N V I L LN K H I DAY K TF P TE P K D KK KK KT D EA Q PL P Q R Q K PT V TL P PA D MD D FS R QL Q NS M SG A ST T QA*
GZ-C, GZ-B 39 dels +29 dels	27638:27844 MNE TLI D FYLCFLA FLLFLV LIMLII FWF SLE I QD LEEPCTKV LCSC I CT V VR CASN KPHV LED P	*	*	28146: 29414 sars9a nucleocapsid protein MSDNGPQS N QR S APR I T F GG G TD S DN N QN G RNGARP K QRRPQ G LP N NTAS W FTALTQHGKEELRFPRGQ G VP I NT N SG P DD Q I GY Y R RAT TR RV R GG D GM K M K LS P R W Y F Y L GT G PEAS L PY G AN K E G I V W V A T E G AL N TP K D H I G TR R NP N NNAA T V L QL P Q G T T LP K G F YAE G SR V GG G SQASS R SS R SG N SR N STPG S RG N SP A RM S GG G ET A LA L LL D RL N Q L E S K V SG K QQ Q Q Q TV T KK S AA E ASK K PR Q K R T A TK V QYN V TA Q FG R RG P EQ T Q G N F GD Q DL I R Q G T D Y K W WP Q I A Q F AP S AS A FF G MS R I G ME V TP S GT W LT Y H G A I K L DD K D P Q F K D N V I L LN K H I DAY K TF P TE P K D KK KK KT D EA Q PL P Q R Q K PT V TL P PA D MD D FS R QL Q NS M SG A ST T QA*

The amino acid sequences of the Orf7b, Orf8 (8a and 8b) and N protein as predicted for the major SARS-CoV deletion variants are

listed in the table. Amino acid sequence changes due to the deletion events are labeled in red. Corresponding nucleotide coordinates for each predicted open reading frame are based on the GZ02 sequence.

Table S2. List of all SNVs with multiple occurrences.

Sequence alignments were generated using CLUSTALW 1.83 with the Gonnet nuclear acid comparison matrix for the 63 SARS-CoV genomic sequences analyzed in this study (61 human SARS-CoV and 2 plam civet SARS-like-CoV sequences). The names of the genomic sequences (as from GenBank) are listed on the first column based on clusters determined by our scoring algorithm (see below). Characteristics of the SNVs observed are listed in the first row, including the variant nucleotides observed at the loci, the affected codon sequence, the resultant amino acids, the amino acid coordinate and the nucleotide coordinate (based on GZ02, shaded horizontally in blue). The relative genomic positions of the listed SNVs are indicated by the predicted open reading frames as shaded bars at the top of the Table.

Variant loci that are characteristic of and allow the segregation of SARS-CoV genotypes into major groups were determined by a method developed for this work. First, for a given mutation site i, sequences are sorted into two groups according to the nucleotide on this site. Second, The number of mutations on the other mutation site k of group j is counted to give $N_k^j(i)$. Then the resulting counts are summed up to give the score S (i) for site i.

$$S(i) = \sum_{k=1}^C \left\{ \sum_{j=1}^2 N_k^j(i) \right\}$$

Where C is the number of mutation sites. Then the scores for all mutation sites are obtained by reiterating the above steps. Finally, the mutation site m that carries the smallest score is chosen to be the primary clustering marker and the genotypes are initially clustered based on this primary marker.

$$m = \arg \min_{1 \leq i \leq C} \{S(i)\}$$

Once the first level of clustering is determined, the above process can be applied recursively until all of the sub-clustering is completed. Further fine adjustment in the resultant clusters / groups of genotypes was performed by integrating information of the epidemiological relationship of genotypes and sequence quality. The same data set was used to generate an unrooted phylogenetic tree with the PHYLIP software package (Fig. 2).

The table lists all of the 107 SNVs that are observed in more than one of the 63 genomic sequences. Only 85 of these SNVs are seen in more than one of the human-derived sequences. Fifty-two of these SNVs were predicted to cause amino acid changes (non-synonymous variations).

SNVs that contribute to the grouping of genotypes based on the predominant clustering criteria described above were further highlighted in the first row with different color shading, except for:

1. The 22 nucleotide sites exhibiting a sequence only observed in the SARS-like-CoV sequences from palm civets or with a single variation observed in human SARS-CoV sequences. These SNVs are shaded with orange color.
2. Synonymous variations.
3. Non-synonymous variations causing similar amino acid alterations, with the exception of the 2 SNVs that were shown to be significantly associated with the transition between the epidemic phases, namely, nt 9404 and nt 17564. The nt 22522 was also highlighted for its association with some minor phase variation events.
4. Some SNVs with the following features:

i. Variations that were only shown in samples of the same patient (*e.g.*, nts 508, 17131, and 28089 for GZ02/GZ01[GD01]; nt 25320 for HSZ-Cc/HSZ-Cb).

ii. The variations were only shared by two sequences, one or both of the sequences belonged to a minor genotype group (*e.g.*, nts 17421, 21637, and 25521 for the genotypes gz50 and GZ-A).

iii. Variations that were not consistent within one transmission path (*e.g.*, nt 9095 for GZ02/GD01 and gz50; nt 25844 for GZ02/GD01 and the ZS group).

The positions of the 5-nt motif used to classify the major genotypes are shaded, including the 2 loci external to the S gene (17564 and 27827, shaded in grey) and the three loci in the S gene, namely positions 21721 (yellow), 22222 (blue) and 23823 (green). These S gene residues are shaded in the same color scheme as used in Figs. 2 and S4. All other SNVs that contribute to the clustering of genotypes are highlighted in pink.

The listed genomic sequences are clustered into groups based on our scoring method. Major genotype clusters were demarcated by solid horizontal lines, where a green line separates the early and middle phases, while a blue line separates the middle and late phases. The further sub-classification of genotypes within the major groups is demarcated by dashed horizontal lines.

	orf1ab polyprotein (pp1ab)								S				sars3a			E	M	X	N	
	orf1a polyprotein (pp1a)				nonstructural polyprotein				S				sars3a		E	M	X	N		
SNV															E	M	X	N		
SNV	Codon	AA switch	AA residue #	nt coordinate																
SZ3	GCC	T	GCA	CTT	CAA	A	C	GTA	C	CCC	T	G	C	G	T	GG	CAG	AT	G	
SZ16	GCCT	T	GCA	CTT	CAA	A	C	GTA	C	CCC	T	G	C	G	T	GG	CAG	AT	G	
GZ02	TCCG	G	ACCT	CGA	CAT	A	T	CCCT	G	GTA	T	CC	C	A	AT	T	CC	AC	CG	
GD01	TTC	GGT	ACCT	CGC	CA	T	A	CCCT	T	CGT	A	CC	T	C	AG	TT	CGT	AT	CCTACG	
HGZ8L1-A	GCCT	T	GCAC	CT	CAA	A	C	GTA	C	CCC	T	G	C	G	T	GG	CAG	AT	G	
HSZ-Cc	GCCG	GT	AT	CTT	ACT	A	GGT	CCCT	T	TGCT	C	T	G	C	TA	CC	AC	CC	ACG	
HSZ-A	GCCG	GT	AT	CTT	ACCA	N	NN	CCCT	T	TGCT	C	T	G	C	TA	CC	NN	AA	ACG	
HSZ-Bb	GCCG	GT	AT	CTT	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	NN	AA	ATACG	TT	
HSZ-Cb	GCCG	GT	AT	CTT	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	NN	AA	ATACG	TT	
HSZ-Bc	GCCG	GT	AT	CTT	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT	
GZ50	GTC	CG	TTT	CTC	ACCA	AGGT	CCCT	T	TGCT	T	CTT	T	G	C	TA	CC	GA	AA	ATACG	TT
GZ-A	GTC	CG	TTT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT	
JMD	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
HGZ8L1-B	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
ZS-A	GCT	GT	TT	ACCT	CAA	A	CAGA	GATT	CCA	CGCT	T	CC	CT	CG	TA	CA	CC	GT	CCACCCATA	
ZS-B	GCT	GT	TT	ACCT	CAA	A	CAGA	GATT	CCA	CGCT	T	CC	CT	CG	TA	CA	CC	GT	CCACCCATA	
ZS-C	GCT	GT	TT	ACCT	CAA	A	CAGA	GATT	CCA	CGCT	T	CC	CT	CG	TA	CA	CC	GT	CCACCCATA	
BJ04	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
BJ03	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
BJ02	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
BJ01	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
CUHK-W1	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
HZS2-D	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
HZS2-E	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
HZS2-C	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
HGZ8L2	GCCG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT	
HZS2-Bb	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
HSZ2-A	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
HZS2-Fc	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
HZS2-Fb	GTC	CG	GT	AT	CTC	ACCA	AGGN	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
TWC	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
Sin2679	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
ZJ01	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
HSR	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
TW1	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
HKU-39849	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT
GZ-D	GTC	CG	GT	AT	CTC	ACCA	AGGT	CCCT	T	TGCT	C	T	G	C	TA	CC	GA	AA	ATACG	TT

	orf1ab polyprotein (pp1ab)								S								sars3a				E	M	N			
	orf1a polyprotein (pp1a)				nonstructural polyprotein				S				sars3a				E		M		N					
SNV	aa		codon		aa		codon		aa		codon		aa		codon		aa		codon		aa		codon			
	ct	gt	at	tc	ct	gt	at	tc	ct	gt	at	tc	ct	gt	at	tc	ct	gt	at	tc	ct	gt	at	tc	ct	
SNV	ag	aaa3	K-K	376	29247	G				ag	aaa3	K-K	376	29247	G				ag	aaa3	K-K	376	29247	G		
Codon	ct	act2	T-T	25	28193				ct	act2	T-T	25	28193				ct	act2	P-L	111	27243					
AA switch	ac	aat1	N-H	76	28193				ac	aat1	N-H	76	28193				ac	aat1	F-L	76	28089					
AA residue #	ct	ttt1	R-C	17	27827				ct	ttt1	R-C	17	27827				ct	ttt1	I-V	86	26653					
nt coordinate	at	tgc3	S-F	11	27810				at	tgc3	S-F	11	27810				ct	tgc2	A-V	68	26600					
Urbani	ta	tgg1	W-R	193	25844				ta	tgg1	W-R	193	25844				ta	tgg1	W-R	193	25844					
Sin2748	ga	ggt1	G-S	381	26410				ga	ggt1	G-S	381	26410				ga	ggt1	G-S	381	26410					
Sin2677	ct	gcc2	A-V	1208	25114				ct	gcc2	A-V	1208	25114				ct	gcc2	A-V	1208	25114					
Sin2500	ac	ccs3	P-P	261	26050				ac	ccs3	P-P	261	26050				ac	ccs3	P-P	261	26050					
Frankfurt	at	aat3	A-A	255	26032				at	aat3	A-A	255	26032				at	aat3	A-A	255	26032					
Sin2774	tg	tgc3	C-C	63	26586				tg	tgc3	C-C	63	26586				tg	tgc3	C-C	63	26586					
CUHK-Su10	ta	tta1	L-F	1265	25289				ta	tta1	L-F	1265	25289				ta	tta1	L-F	1265	25289					
CUHK-LC1	ct	cic1	L-F	1247	25230				ct	cic1	L-F	1247	25230				ct	cic1	L-F	1247	25230					
CUHK-AG01	ct	tgc3	C-C	1025	24566				ct	tgc3	C-C	1025	24566				ct	tgc3	C-C	1025	24566					
CUHK-AG02	ct	ttt2	F-S	360	22570				ct	ttt2	F-S	360	22570				ct	ttt2	F-S	360	22570					
CUHK-AG03	tc	tta2	L-S	665	23485				tc	tta2	L-S	665	23485				tc	tta2	L-S	665	23485					
TWH	ca	aac3	N-K	227	2227				ca	aac3	N-K	227	2227				ca	aac3	N-K	227	2227					
TC1	ct	act2	T-I	244	22222				ct	act2	T-I	244	22222				ct	act2	T-I	244	22222					
TWY	tc	tta2	L-S	239	22207				tc	tta2	L-S	239	22207				tc	tta2	L-S	239	22207					
TWS	ca	aac3	N-K	227	2227				ca	aac3	N-K	227	2227				ca	aac3	N-K	227	2227					
TWK	ct	ttt2	F-S	218	22145				ct	ttt2	F-S	218	22145				ct	ttt2	F-S	218	22145					
TWJ	tg	tgc2	G-D	77	21721				tg	tgc2	G-D	77	21721				tg	tgc2	G-D	77	21721					
TC3	cg	agg2	T-R	75	21715				cg	agg2	T-R	75	21715				cg	agg2	T-R	75	21715					
TC2	ct	tca2	S-L	49	21637				ct	tca2	S-L	49	21637				ct	tca2	S-L	49	21637					
GZ-B	tg	gtc3	A-A	27	21572				tg	gtc3	A-A	27	21572				tg	gtc3	A-A	27	21572					
GZ-C	ta	aat3	N-N	2694	21479				ta	aat3	N-N	2694	21479				ta	aat3	N-N	2694	21479					
TOR2	tc	ttt2	F-V	3062	9448				tc	ttt2	F-V	3062	9448				tc	ttt2	F-V	3062	9448					
CUHK-LC2	tg	tgc3	C-C	2746	8505				tg	tgc3	C-C	2746	8505				tg	tgc3	C-C	2746	8505					
CUHK-LC3	ta	tta2	L-S	2269	7070				ta	tta2	L-S	2269	7070				ta	tta2	L-S	2269	7070					
CUHK-LC4	ct	ttt2	F-V	2222	6929				ct	ttt2	F-V	2222	6929				ct	ttt2	F-V	2222	6929					
CUHK-LC5	tg	tgc3	C-C	2116	6612				tg	tgc3	C-C	2116	6612				tg	tgc3	C-C	2116	6612					
gt	ggg2	A-S	549	1909				gt	ggg2	A-S	549	1909				gt	ggg2	A-S	549	1909						
gt	aat3	A-A	1102	3570				gt	aat3	A-A	1102	3570				gt	aat3	A-A	1102	3570						
ct	ttt2	A-V	1233	3962				ct	ttt2	A-V	1233	3962				ct	ttt2	A-V	1233	3962						
tc	tgc3	S-S	1196	3852				tc	tgc3	S-S	1196	3852				tc	tgc3	S-S	1196	3852						
ct	ttt2	I-T	1136	3671				ct	ttt2	I-T	1136	3671				ct	ttt2	I-T	1136	3671						
ct	ttt2	T-I	1121	3626				ct	ttt2	T-I	1121	3626				ct	ttt2	T-I	1121	3626						
ct	ttt2	R-K	1319	4220				ct	ttt2	R-K	1319	4220				ct	ttt2	R-K	1319	4220						
ct	ttt2	A-V	1021	3326				ct	ttt2	A-V	1021	3326				ct	ttt2	A-V	1021	3326						
tc	tgc3	C-C	2526	7842				tc	tgc3	C-C	2526	7842				tc	tgc3	C-C	2526	7842						
ct	ttt2	F-V																								

Table S3. Statistical analysis for the change of Ka/Ks ratios for different coding regions of the SARS-CoV sequences during the different epidemic phases.

Proteins	Epidemic phases	$\bar{K}s(10^{-3})$	$\bar{K}a(10^{-3})$	$\bar{Ka / Ks}$	s.e.($\bar{Ka / Ks}$)	H_1^*	P-value†
Spike	early	1.698	1.997	1.248	0.081	Ka/Ks (early) > Ka/Ks (middle)	2.3×10^{-7}
	middle	2.377	0.898	0.410	0.087		
	late	1.355	0.257	0.219	0.043	Ka/Ks (middle) > Ka/Ks (late)	0.034
Orf1b	early	1.494	0.570	0.562	0.145	Ka/Ks (early) > Ka/Ks (late)	0.091
	middle	1.048	0.240	0.315	0.108		
	late	0.577	0.159	0.344	0.048		
Orf1a	early	1.264	1.030	0.923	0.124	Ka/Ks (early) > Ka/Ks (late)	7.4×10^{-5}
	middle	0.526	0.565	1.293	0.202		
	late	0.557	0.139	0.369	0.060		

* H_1 means the alternative hypothesis.

†One-sided unpaired two-sample t-test was used

Table S4. List of all SNVs in S genes with multiple occurrences.

Sequence alignments were generated by the method described in the notes for Table S2. Sequences used were described in the legend of Fig. S7. SNVs that contribute to the grouping of genotypes based on the predominant clustering criteria previously described were further highlighted in the rows with different color shading:

1. Green is for the SNV characteristically distinguishable between the early versus the middle phase viral isolates. Yellow is for the SNV that was characteristic for the motif transition between the middle phase and the late phase. Blue is for the SNV characteristically distinguishable between the middle versus the late phase viral isolates. These 3 coloring systems are consistent with that used for the whole genome analysis (Table S2 and Fig. 2)
2. Pink is for the SNVs characteristically distinguishable between groups of the early phase isolates, including synonymous and non-synonymous variations. Dark green is for the SNVs causing non-synonymous variations among significant portions of the early isolates. These variations were less systematic than the pink shaded SNVs with respect to their correlation with the epidemiological data.
3. Light blue indicates the palm civet specific SNVs known so far, while purple indicates the SNVs shared by palm civet SARS-like-CoV, SZ16 and SZ3, and the most recent human SARS-CoV, GD03T13.

	25230	ct	ctc1	L-F	1247	3739	C	C	C
nt coordinate	25114	ct	gcc2	A-V	1208	3623	C		
SNV	25027	tc	ctc2	L-P	1179	3536	T		
Codon	24978	ga	aaa1	E-K	1163	3487	G		
AA switch	24932	ct	ctt1	D-D	1147	3441	C		
AA residue #	24566	ct	tgt3	C-C	1025	3075	C	C	C
S-nt	24171	ag	acc1	T-A	894	2680	G	G	G
	23823	gt	tat1	D-Y	778	2332	G	G	G
	23752	ct	gct2	A-V	754	2261	T	T	T
	23718	ag	aca1	T-A	743	2227	G	G	G
	23593	ct	tca2	S-L	701	2102	T	T	T
	23485	tc	tta2	L-S	665	1994	C	C	C
	23310	tc	tct1	S-P	607	1819	T	T	T
	22951	cg	act2	T-S	487	1460	G	G	G
	22928	ta	aat3	N-K	479	1437	A	A	A
	22570	tc	ttt2	F-S	360	1079	C	C	C
	22522	ga	aaa2	R-K	344	1031	G	G	G
	22422	ag	ggaa1	R-G	311	931	G	G	G
	22273	ca	aac3	N-K	227	681	A	A	A
	22145	ct	cc3	P-P	218	654	C	C	C
	21721	ga	ggc2	G-D	77	230	T	T	T
	21715	cg	acg2	T-R	75	224	C	C	C
	21637	ct	tca2	S-L	49	146	C	C	C
	21572	tg	gct3	A-A	27	81	T	T	T

