

CoVDB: a comprehensive database for comparative analysis of coronavirus genes and genomes

Yi Huang¹, Susanna K. P. Lau^{1,2,3}, Patrick C. Y. Woo^{1,2,3,*} and Kwok-yung Yuen^{1,2,3}

¹Department of Microbiology, ²Research Centre of Infection and Immunology and ³State Key Laboratory of Emerging Infectious Diseases, The University of Hong Kong, Hong Kong

Received June 27, 2007; Revised August 24, 2007; Accepted September 11, 2007

ABSTRACT

The recent SARS epidemic has boosted interest in the discovery of novel human and animal coronaviruses. By July 2007, more than 3000 coronavirus sequence records, including 264 complete genomes, are available in GenBank. The number of coronavirus species with complete genomes available has increased from 9 in 2003 to 25 in 2007, of which six, including coronavirus HKU1, bat SARS coronavirus, group 1 bat coronavirus HKU2, groups 2c and 2d coronaviruses, were sequenced by our laboratory. To overcome the problems we encountered in the existing databases during comparative sequence analysis, we built a comprehensive database, CoVDB (<http://covdb.microbiology.hku.hk>), of annotated coronavirus genes and genomes. CoVDB provides a convenient platform for rapid and accurate batch sequence retrieval, the cornerstone and bottleneck for comparative gene or genome analysis. Sequences can be directly downloaded from the website in FASTA format. CoVDB also provides detailed annotation of all coronavirus sequences using a standardized nomenclature system, and overcomes the problems of duplicated and identical sequences in other databases. For complete genomes, a single representative sequence for each species is available for comparative analysis such as phylogenetic studies. With the annotated sequences in CoVDB, more specific blast search results can be generated for efficient downstream analysis.

INTRODUCTION

Coronaviruses are found in a wide variety of animals and are associated with respiratory, enteric, hepatic and neurological diseases of varying severity. Based on genotypic and serological characterization, coronaviruses

were divided into three distinct groups (1–3). As a result of the unique mechanism of viral replication, coronaviruses have a high frequency of recombination (2,4).

The recent severe acute respiratory syndrome (SARS) epidemic, the discovery of SARS coronavirus (SARS-CoV) and identification of SARS-CoV-like viruses from Himalayan palm civets and a raccoon dog from wild live markets in China have led to a boost in interest on discovery of novel coronaviruses in both humans and animals (5–9) (Figure 1). For human coronaviruses, a novel group 1 human coronavirus, human coronavirus NL63 (HCoV-NL63) was reported in 2004 (10,11), while we described the discovery, complete genome sequence and genetic diversity of a novel group 2 human coronavirus, coronavirus HKU1 (CoV-HKU1) in 2005 (4,12–14). As for animal coronaviruses, six group 1 (15–17), four group 2, including bat SARS-CoV and two new subgroups of group 2 coronaviruses (6,8,18,19), and 11 group 3 (20–23) coronaviruses have recently been described.

By July 2007, more than 3000 coronavirus sequence records, including a total of 264 complete genomes, are available in GenBank (24). Among the 25 coronavirus species with complete genome sequence available, six were sequenced by our group, including CoV-HKU1 and bat SARS-CoV (13,16,18,19). Furthermore, we defined two novel subgroups of group 2 coronavirus (18). During the process of batch sequence retrieval for comparative genome analysis of the coronavirus genomes that we sequenced, we encountered several major problems about the coronavirus sequences in GenBank as well as other coronavirus databases (Coronaviridae Bioinformatics Resource, <http://athena.bioc.uvic.ca/database.php?db=coronaviridae>; PATRIC <http://patric.vbi.vt.edu>) (25). First, in GenBank, the non-structural proteins in the polyprotein encoded by *orf1ab* were not annotated. Second, in all databases, for the non-structural proteins encoded by ORFs downstream to *orf1ab*, the annotations are often confusing because they are not annotated using a standardized system. Third, multiple accession numbers are often present for reference sequences (26). These problems often lead to confusion when sequence retrieval

*To whom correspondence should be addressed. Tel: 852 2855 4892; Fax: 852 2855 1241; Email: pcywoo@hkucc.hku.hk
Correspondence may also be addressed to Kwok-yung Yuen. Tel: 852 2855 4892; Fax: 852 2855 1241; Email: hkumicro@hkucc.hku.hk

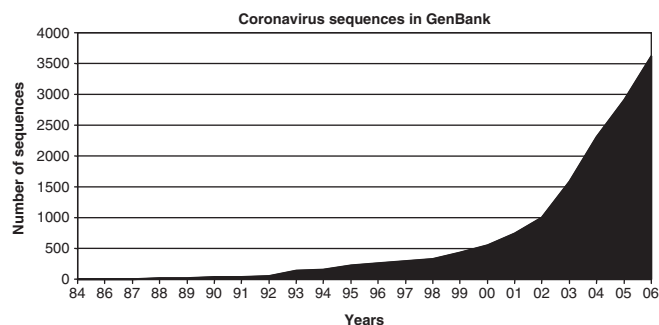


Figure 1. Number of coronavirus sequences in GenBank from 1984 to 2006.

is performed. Fourth, coronaviruses, especially SARS-CoV, amplified from different specimens may contain the same genome or gene sequences. These sequences usually lead to redundant work when they are analyzed.

In view of these problems, we started to develop our own database for coronavirus gene and genome sequences in 2005. In this database, CoVDB, we sought to create a user-friendly platform for efficient batch sequence retrieval, which is crucial for comparative genome analysis. In this article, we describe this comprehensive database of annotated coronavirus genes and genomes, which provides a central source of information about coronaviruses. To further increase the usefulness of CoVDB, commonly used bioinformatics tools were also included for analysis of the sequence data.

MATERIALS AND METHODS

Database description

Sequence data. CoVDB is a web-based coronavirus database. Data of CoVDB is stored and managed by MySQL database management system. By July 2007, CoVDB contains 3982 coronavirus sequences and one torovirus genome sequence. Two hundred and sixty-four of them are complete genomes and the rest are partial genomes or genes. All data were retrieved from GenBank using modules of bioperl. We annotated sequences without gene information or non-structural protein boundary and labeled the 5' and 3' untranslated regions (UTRs) of the genomes. By July 2007, CoVDB contains 12 344 genes and UTRs.

Information on coronavirus genome characteristics. In addition to the two sequence retrieval pages, CoVDB collects information on coronavirus sequence characteristics, including genome organization, a brief description on each complete coronavirus genome, GC content, polyprotein cleavage sites, transcription regulatory sequences, acidic tandem repeat sequences and known RNA structures. These pieces of information can be accessed by clicking 'Genome' in the top menu bar of CoVDB. In the 'Tools' page, blast similarity search (27) against annotated coronavirus sequences in CoVDB can be performed and other commonly used tools are also provided.

Functionality of the database

Batch sequence retrieval. The main goal for setting up CoVDB is to provide a convenient and efficient platform for retrieving batches of coronavirus gene sequences. The interfaces of the database are simple and user friendly. All genes and genomes contain links to GenBank and/or pubmed. CoVDB contains two main pages for sequence retrieval. From the homepage, one can enter the first main page for retrieval of complete genomes and their genes by clicking 'CoVDB' (Figure 2a). From this page, users can obtain genes from specific coronavirus species by selecting the corresponding check boxes. We defined one representative genome from each species as the 'Type strain'. Most of the time, this 'Type strain' is the one assigned as the reference sequence in GenBank. By choosing the 'Type strain only' option, users can obtain one gene sequence per species and construct phylogenetic tree or perform other comparisons. An example of retrieving complete genome or a specific gene of complete genome of selected species is shown in Figure 2b and c.


From the page for retrieval of complete genomes and their genes, one can enter the second main page for retrieval of all complete and/or incomplete genes of a coronavirus (Figure 3a) by clicking 'From all groups of genes'. In this page, all the gene sequences are grouped vertically according to which coronavirus group and subgroup they belong to, and horizontally by the names of the genes. The option 'Exclude partial CDS' can be used if only complete genes are required. An example of retrieving all the sequence of a particular gene for a group of coronavirus is shown in Figure 3b. If the translated sequence of a selected gene has more than one stop codon which is probably due to sequencing error, the number in the 'Length' column of this gene will be marked in red.


Polyprotein annotation. In all coronavirus genomes, orf1ab occupies two-thirds of the genome and it is translated as a polyprotein. This polyprotein is post-translationally cleaved by 3C-like protease (3CL^{pro}) and papain-like protease (PL^{pro}) into 15–16 non-structural proteins. Some of the non-structural proteins, such as RNA-dependent RNA polymerase, helicase, 3CL^{pro} and PL^{pro} are essential for replication or virulence of the coronavirus, although the functions of others are still unclear. Due to the essentiality of the non-structural proteins, these sequences are often used for evolutionary analysis, primer design, etc. However, except for the reference sequences, detailed cleavage site information is not provided for the non-structural proteins in other sequences in GenBank. Since it has been shown that 3CL^{pro} and PL^{pro} of coronavirus cleave at conserved specific amino acids, the putative cleavage sites of the 15–16 non-structural proteins can be predicted by multiple sequence alignment. Using these pieces of information, we have annotated these non-structural proteins in all the coronavirus sequences for easy retrieval in CoVDB.

Protein/gene name unification. By convention, all non-structural proteins in the polyprotein encoded by orflab are named as 'nsp', with each protein numbered consecutively starting from the 5' end (nsp1–nsp16). The structural proteins after the polyprotein are hemagglutinin esterase (HE, in group 2a coronaviruses), spike glycoprotein (S), envelope protein (E), membrane protein (M) and nucleocapsid protein (N). However, there is

no unified naming system for the non-structural proteins encoded by ORFs downstream to orflab. This lack of a unified system greatly reduces the stability and accuracy of ortholog retrieval.

In CoVDB, with the aim of facilitating gene retrieval, we tried to unify the naming of these non-structural proteins from different groups of coronaviruses. On the other hand, we have also tried to avoid radical changes

(a)  Home | [CoVDB](#) | [Genome](#) | [Tools](#) | [Publications](#) | [Contact](#) | [Links](#)

Sequence retrieve From all groups of genes 

Get genes from 264 completed genomes only

Group 1: HCoV-229E (1) HCoV-NL63 (9) TGEV (9) PEDV (1) PRCV (1) FIPV (3) BtCoV (1)
 Bat-CoV HKU2 (4)

Group 2a: HCoV-OC43 (5) CoV-HKU1A (13) CoV-HKU1B (3) CoV-HKU1C (6) BCoV (12) MHV (8) PHEV (1)
 Sable antelope (1) Giraffe (3)

SARS-CoV: Human (131) Palm civet (16) Bat HKU3 (8) Badger (1)

Group 2c: Bat-CoV HKU4 (5) Bat-CoV HKU5 (4)

Group 2d: Bat-CoV HKU9 (4)

Group 3: IBV (11) IBV-partridge (1) IBV-peafowl (1)

Torovirus: BoTV (1)

Select all Type strain only [Gene synonyms !\[\]\(a86311389280ae0b3c35372f357d29e9_img.jpg\)](#)

Note:

- Most of type strains here indicate reference sequences in genbank. But some of them which we have special interests are also included. If you hope to get all available sequences within one species, just unselect 'Type strain only'.
- The non-coding regions of 5' end and 3' end are also included.

Last updated: Aug 20, 2007. [Firefox !\[\]\(8f42d78483b86945861e530255975cad_img.jpg\)](#) is the recommended web browser for this site!
 Copyright (c) 2006 Department of Microbiology, The University of Hong Kong. All rights reserved.

(b) Thanks for searching COV db. Here are the results:

NCBIacc	Shortname	Length	Strain/isolate	Country:Region	Organism	Group	PMID	GC	GCskew
<input type="checkbox"/> NC_005147	HCoV-OC43	30738	OC43	Belgium	Human coronavirus OC43	G2a	15650185	0.368	0.176
<input type="checkbox"/> NC_006577	CoV-HKU1A	29926	N1	China:Hong Kong	Human coronavirus HKU1A	G2a	15613317	0.320	0.188
<input type="checkbox"/> NC_003045	BCoV	31028	BCoV-ENT	USA	Bovine coronavirus	G2a	11714968	0.371	0.174
<input type="checkbox"/> NC_001846	MHV	31357	MHV-A59	USA	Murine hepatitis virus	G2a	9426441	0.417	0.142
<input type="checkbox"/> NC_007732	PHEV	30480	VW572	Belgium	Porcine hemagglutinating encephalomyelitis virus	G2a	16809333	0.372	0.164
<input type="checkbox"/> NC_004718	SARS-human	29751	Tor2	Canada:Toronto	Human SARS coronavirus	G2b	15020242	0.407	0.020
<input type="checkbox"/> AY304488	SARS-civet	29731	SZ16	China:Shenzhen	Civet SARS coronavirus	G2b	12958366	0.408	0.020
<input type="checkbox"/> DQ022305	SARS-bat	29728	HKU3-1	China:Hong Kong	Bat SARS coronavirus	G2b	16169905	0.411	0.027

8 records!

Select all

Figure 2. Screenshots of CoVDB complete genome retrieval pages. **(a)** Specific gene can be retrieved using the pull-down list at the left lower corner. The number in brackets indicates the number of complete genomes for that coronavirus. **(b)** Example of showing genomes of selected species (some group 2a coronaviruses and SARS-CoV-related coronaviruses). Default is to show the "Type strain" for each species only. The columns NCBIacc and PMID link to GenBank and pubmed, respectively. **(c)** Example of showing S gene of selected species by choosing S in the pull-down list. For genes downstream to orflab, sequences upstream to the initiation codons can also be retrieved from this result page. This function is particularly useful for the detection of transcription regulatory sequences.

(c) Thanks for searching COV db. Here are the results:

ntacc	Shortname	Tag	Gene	From	To	Shift	Length	Protein_id	Strain/Isolate	Group	Country:Region
<input type="checkbox"/> NC_005147	HCoV-OC43	CDS	S	23644	27729	0	4086	NP_937950	OC43	G2a	Belgium
<input type="checkbox"/> NC_006577	CoV-HKU1A	CDS	S	22942	27012	0	4071	YP_173238	N1	G2a	China:Hong Kong
<input type="checkbox"/> NC_003045	BCoV	CDS	S	23641	27732	0	4092	NP_150077	BCoV-ENT	G2a	USA
<input type="checkbox"/> NC_001846	MHV	CDS	S	23929	27903	0	3975	NP_045300	MHV-A59	G2a	USA
<input type="checkbox"/> NC_007732	PHEV	CDS	S	23427	27476	0	4050	YP_459952	VW572	G2a	Belgium
<input type="checkbox"/> NC_004718	SARS-human	CDS	S	21492	25259	0	3768	NP_828851	Tor2	G2b	Canada:Toronto
<input type="checkbox"/> AY304488	SARS-civet	CDS	S	21477	25244	0	3768		SZ16	G2b	China:Shenzhen
<input type="checkbox"/> DQ022305	SARS-bat	CDS	S	21471	25199	0	3729	AAV88866	HKU3-1	G2b	China:Hong Kong

8 records!

 Select all
Get -100 nt to +0 nt around start position

Figure 2. Continued.

in the names that may lead to confusion. In CoVDB, these non-structural proteins are named as NS2a, NS3x, NS4x, NS5x and NS7x ($x = a, b, c, \dots$). NS2a denotes the ORF between orf1ab and HE of group 2a coronaviruses. NS3x denotes the ORFs between S and E of groups 1, 2c, 2d and 3 coronaviruses. In most of these coronaviruses, there are two NS3x, named NS3a and NS3b. However, in group 1 coronaviruses, the genomes of some members (e.g. HCoV-NL63, PEDV) contain only one ORF between S and E. When we compared their putative amino acid sequences to the corresponding ones in other group 1 coronavirus genomes using BLAST, as well as searching for conserved domains using motifscan, results showed that the putative proteins encoded by these ORFs belonged to a protein family in Pfam originally assigned as 'Corona_NS3b' (accession number PF03053). Therefore, we named these ORFs as NS3b. NS4x denotes the ORFs between S and E of group 2a coronaviruses. NS5x denotes the ORFs between M and N of group 3 coronaviruses. One exception is NS5a of group 2a coronaviruses. Traditionally, this name denotes an ORF upstream of E in group 2a coronaviruses. Therefore, we have kept this name for that ORF in CoVDB. NS7x denotes the ORFs downstream of N gene. It is important to note that due to variations in genome organizations among different groups of coronaviruses (Table 1), NS genes with the same name in different coronavirus groups may not be orthologs of each other. The complete genome gene search page of CoVDB contains a link to a Gene synonyms page, which includes a list of synonymous names of the various genes in the coronavirus genomes.

Identical sequence labeling. Sequence redundancy is another problem of coronavirus sequences in public nucleotide databases. Different strains of the same species from samples collected in different locations or at different

times may possess completely or partially identical sequences. These sequences, though containing important epidemiological information, increase the workload during sequence analysis. In CoVDB, we compared all nucleotide sequences and labeled the identical ones to mitigate this problem. Users can choose to show or not to show strains with identical sequences by clicking on the check boxes to the left of the page (Figure 3b).

Blast similarity search. During the process of coronavirus gene sequences analysis, we encountered a major problem when coronavirus gene sequences, especially those of orf1ab, were used for blast search against GenBank or any other coronavirus databases. When part of the orf1ab gene (e.g. nsp5) is used as the query sequence, instead of getting the gene for the specific non-structural protein that the query sequence is homologous to, the results will only show that the hits are within orf1ab, or in some cases, shown to be within the entire coronavirus genome. Much time will be needed for further analyzing the results manually in order to locate the positions of the cleavage sites of the corresponding genes for the non-structural proteins, making it very inefficient for further downstream work.

This problem has been overcome by the annotated sequences in CoVDB. The blast search page of CoVDB is an interface for facilitating coronavirus similarity search. The background support program, blastall, is from the NCBI Blast package. The blast search page can be entered by clicking 'Tools' in the top menu bar in any page of CoVDB. Since all sequences in CoVDB are annotated, they can be grouped into different datasets for blast search. Users can choose one of the three nucleotide and two protein sequence datasets as the database for comparison (Figure 4). The three nucleotide sequence datasets are: CoV genes (nsp + genes after 1ab),

Sequence retrieve From complete genomes Exclude partial CDS

Gene	G1	G2a	G2b	G2c	G2d	G3	Gene	G1	G2a	G2b:sars	G2c	G2d	G3
S'UTR	<input type="checkbox"/> 36	<input type="checkbox"/> 77	<input type="checkbox"/> 148	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 21	NS3b	<input type="checkbox"/> 67	-	<input type="checkbox"/> 3b: 151	<input type="checkbox"/> 8	-	<input type="checkbox"/> 43
1ab	<input type="checkbox"/> 26	<input type="checkbox"/> 52	<input type="checkbox"/> 148	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 14	NS3c	<input type="checkbox"/> 11	-	-	<input type="checkbox"/> 9	-	E
nsp1	<input type="checkbox"/> 29	<input type="checkbox"/> 58	<input type="checkbox"/> 156	<input type="checkbox"/> 9	<input type="checkbox"/> 4	-	NS3d	<input type="checkbox"/> 4	-	-	<input type="checkbox"/> 9	-	-
nsp2	<input type="checkbox"/> 29	<input type="checkbox"/> 52	<input type="checkbox"/> 168	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 35	NS4	-	<input type="checkbox"/> 9	-	-	-	-
nsp3	<input type="checkbox"/> 110	<input type="checkbox"/> 63	<input type="checkbox"/> 196	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 34	NS4a	-	<input type="checkbox"/> 41	-	-	-	-
nsp4	<input type="checkbox"/> 28	<input type="checkbox"/> 51	<input type="checkbox"/> 180	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 14	NS4b	-	<input type="checkbox"/> 35	-	-	-	-
nsp5	<input type="checkbox"/> 29	<input type="checkbox"/> 52	<input type="checkbox"/> 160	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 15	NS5a	-	<input type="checkbox"/> 82	-	-	-	-
nsp6	<input type="checkbox"/> 28	<input type="checkbox"/> 51	<input type="checkbox"/> 158	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 14	E	<input type="checkbox"/> 59	<input type="checkbox"/> 110	<input type="checkbox"/> 197	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 88
nsp7	<input type="checkbox"/> 28	<input type="checkbox"/> 51	<input type="checkbox"/> 156	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 14	M	<input type="checkbox"/> 97	<input type="checkbox"/> 104	<input type="checkbox"/> 197	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 162
nsp8	<input type="checkbox"/> 28	<input type="checkbox"/> 51	<input type="checkbox"/> 159	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 14	NS5a	-	-	-	-	-	<input type="checkbox"/> 59
nsp9	<input type="checkbox"/> 97	<input type="checkbox"/> 51	<input type="checkbox"/> 156	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 14	NS5b	-	-	-	-	-	<input type="checkbox"/> 63
nsp10	<input type="checkbox"/> 29	<input type="checkbox"/> 52	<input type="checkbox"/> 155	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 14	-	-	-	<input type="checkbox"/> 6: 177	-	-	-
nsp11	<input type="checkbox"/> 29	<input type="checkbox"/> 52	<input type="checkbox"/> 154	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 14	-	-	-	<input type="checkbox"/> 7a: 193	-	-	-
nsp12	<input type="checkbox"/> 111	<input type="checkbox"/> 97	<input type="checkbox"/> 182	<input type="checkbox"/> 39	<input type="checkbox"/> 4	<input type="checkbox"/> 96	-	-	-	<input type="checkbox"/> 7b: 198	-	-	-
nsp13	<input type="checkbox"/> 65	<input type="checkbox"/> 52	<input type="checkbox"/> 179	<input type="checkbox"/> 19	<input type="checkbox"/> 4	<input checked="" type="checkbox"/> 15	-	-	-	<input type="checkbox"/> 8: 40	-	-	-
nsp14	<input type="checkbox"/> 34	<input type="checkbox"/> 55	<input type="checkbox"/> 186	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 15	-	-	-	<input type="checkbox"/> 8a: 121	-	-	-
nsp15	<input type="checkbox"/> 30	<input type="checkbox"/> 53	<input type="checkbox"/> 179	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 16	-	-	-	<input type="checkbox"/> 8b: 117	-	-	-
nsp16	<input type="checkbox"/> 34	<input type="checkbox"/> 54	<input type="checkbox"/> 185	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 17	N	<input type="checkbox"/> 142	<input type="checkbox"/> 113	<input type="checkbox"/> 203	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 221
NS2a	-	<input type="checkbox"/> 42	-	-	-	-	N2	-	<input type="checkbox"/> 55	<input type="checkbox"/> 9b: 161	<input type="checkbox"/> 9	-	-
HE	-	<input type="checkbox"/> 157	-	-	-	-	-	-	-	<input type="checkbox"/> 9c: 158	-	-	-
S	<input type="checkbox"/> 360	<input type="checkbox"/> 314	<input type="checkbox"/> 354	<input type="checkbox"/> 10	<input type="checkbox"/> 4	<input type="checkbox"/> 943	NS7a	<input type="checkbox"/> 68	-	-	-	<input type="checkbox"/> 4	<input type="checkbox"/> 5
NS3	<input type="checkbox"/> 9	-	-	-	<input type="checkbox"/> 4	-	NS7b	<input type="checkbox"/> 68	-	-	-	<input type="checkbox"/> 4	<input type="checkbox"/> 8
NS3a	<input type="checkbox"/> 47	-	<input type="checkbox"/> 3a: 169	<input type="checkbox"/> 9	-	<input type="checkbox"/> 57	3'UTR	<input type="checkbox"/> 31	<input type="checkbox"/> 58	<input type="checkbox"/> 155	<input type="checkbox"/> 9	<input type="checkbox"/> 4	<input type="checkbox"/> 73

Last updated: May 30, 2007. Firefox is the recommended web browser for this site!

Copyright (c) 2006 Department of Microbiology, The University of Hong Kong. All rights reserved.

Figure 3. Screenshots of all gene retrieval pages. (a) Gene sequences are grouped vertically according to which coronavirus group and subgroup they belong to, and horizontally by the name of the genes. The numbers next to each checkbox indicates the number of that gene in CoVDB. The option 'Exclude partial CDS' can be used if only complete genes are required. (b) Example of showing the 15 sequences of nsp13 in group 3 coronaviruses. The first column is CoVDB gene id. In the Uniq column, 'Uniq' will be shown if there is no other identical sequence in CoVDB. Otherwise, gene id of the sequences identical to it will be shown.

CoV genes (1ab + genes after 1ab) and CoV GenBank strains, which are the original sequences retrieved from GenBank. The two protein sequence datasets are the translated sequences of the first two nucleotide datasets: CoV proteins (nsp + aa after 1ab) and CoV proteins (1ab + aa after 1ab).

MyBlast. 'MyBlast' employs the same blast program as the Blast page mentioned above. However, instead of selecting a predefined nucleotide or amino acid sequence database, multiple sequences can be pasted into the second sequence input box to generate a temporary sequence database. One or more query sequences can be pasted into the first sequence input box for blastn or blastp search against the temporary sequence database.

ORF finder for coronavirus. This ORF finder is specifically designed for coronavirus genome analysis. The result page shows the positions and lengths of each putative

ORF and the position of the putative ribosomal frame-shift site for translation of orf1ab. The nucleotide or amino acid sequences of the ORFs can be shown by selecting the corresponding check boxes. To facilitate genome comparison and annotation, the most closely related coronavirus, which had been annotated in CoVDB, can be chosen from a pull-down list for comparison using blast search. This function is particularly useful for determining the range of nsp in orf1ab.

DISCUSSION

Rapid and accurate batch sequence retrieval is both the cornerstone and bottleneck for comparative gene or genome analysis. During the process of complete genome sequencing and comparative analysis of the various novel human and animal coronavirus genomes in the past 2 years, we have developed a comprehensive

(b) Thanks for searching COV db. Here are the results:

- NA: not available
- Shift: -1 frameshift position.
- Type: c-complete CDS, p-partial CDS.
- Length: Red number means this sequence may contain multiple stop codons due to sequencing or mutation.
- Country:Region: If no virus location is available, this indicates the submitters' country and region.
- Please let me know if there is any incorrect message. Thanks!

15!

	ntacc	Gene	From	To	Shift	Length	Type	Protein_id	Short	Strain/Isolate	Country:Region	Uniq	PMID
<input checked="" type="checkbox"/>	5307	AY319651	nsp13	15163	16965	0	1803	c	IBV	BJ	China:Beijing	Uniq	0
<input checked="" type="checkbox"/>	5332	AY514485	nsp13	15136	16935	0	1800	c	IBV	Cal99	USA	Uniq	16927130
<input checked="" type="checkbox"/>	5407	AY851295	nsp13	15138	16937	0	1800	c	IBV	Mass41;M41	USA	Uniq	0
<input checked="" type="checkbox"/>	5432	DQ001338	nsp13	15132	16931	0	1800	c	IBV	IBV-EP3	Singapore	Uniq	16137658
<input checked="" type="checkbox"/>	5457	DQ001339	nsp13	15129	16928	0	1800	c	IBV	IBV-p65	Singapore	Uniq	16137658
<input checked="" type="checkbox"/>	5532	DQ288927	nsp13	15125	16924	0	1800	c	IBV	SAIBK	China:Sichuan	Uniq	0
<input checked="" type="checkbox"/>	5557	AY338732	nsp13	15079	16878	0	1800	c	IBV	LX4	China:Heilongjiang	Uniq	15223561
<input checked="" type="checkbox"/>	5582	AY692454	nsp13	15132	16931	0	1800	c	IBV	Beaudette(VC)	USA	Uniq	0
<input checked="" type="checkbox"/>	9959	DQ834384	nsp13	15139	16938	0	1800	c	IBV	M41	USA	Uniq	0
<input checked="" type="checkbox"/>	5357	AY641576	nsp13	15075	16874	0	1800	c	IBV-peafowl	Peafowl/GD/KQ6/2003	China:Guangdong	Uniq	0
<input checked="" type="checkbox"/>	5382	AY646283	nsp13	14964	16763	0	1800	c	IBV-partridge	Partridge/GD/S14/2003	China:Guangdong	Uniq	0
<input checked="" type="checkbox"/>	5482	NC_001451	nsp13	15132	16931	0	1800	c	NP_740630	IBV	Beaudette	United Kingdom	5482 3027249
<input type="checkbox"/>	5507	AJ311317	nsp13	15132	16931	0	1800	c	IBV	BeaudetteCK	United Kingdom	5482	11711626
<input type="checkbox"/>	10006	M94356	nsp13	15132	16931	0	1800	c	IBV	Beaudette(M42)	United Kingdom	5482	3027249
<input type="checkbox"/>	10985	Z30541	nsp13	3843	5642	0	1800	c	IBV	Beaudette	United Kingdom	5482	3027249

Select all

Figure 3. Continued.

Table 1. Genome organization of different groups of coronavirus

Group	Organizations
1	5'UTR-nsp1-16-S-NS3x-E-M-N-(NS7x)-3'UTR
2a	5'UTR-nsp1-16-(NS2a)-HE-S-(NS4x)-NS5a-E-M-N-3'UTR
2b	5'UTR-nsp1-16-S-sars3x-E-M-sars6-sars7x-sars8x-N-3'UTR
2c	5'UTR-nsp1-16-S-NS3x-E-M-N-3'UTR
2d	5'UTR-nsp1-16-S-NS3x-E-M-N-(NS7x)-3'UTR
3	5'UTR-nsp1-16-S-NS3x-E-M-NS5x-N-(NS7x)-3'UTR

database, CoVDB, of annotated coronavirus genes and genomes, which offers efficient batch sequence retrieval and analysis. As shown by our experience in using CoVDB for comparative genome analysis of novel coronaviruses we have discovered (4,13,16,18,19), we find that CoVDB is more rapid and efficient than other existing coronavirus databases for batch sequence retrieval for the following reasons. First, we have performed annotation on all non-structural proteins in the polyprotein encoded by orflab of every single sequence. Second, annotation was performed for the non-structural proteins encoded by ORFs downstream to orflab using a standardized system, with some exceptions given to some names that have been used for a long time so as to minimize confusion. Third, all sequences with identical nucleotide sequences were labeled where one can choose to show or not to show strains with identical sequences. Fourth, CoVDB contains not

only complete coronavirus genome sequences, but also incomplete genomes and their genes. Some genes of coronaviruses, such as *pol*, spike and nucleocapsid are sequenced much more frequently than others because they are either most conserved or least conserved. These gene sequences are particularly important for evolutionary analysis, single nucleotide polymorphism studies and design of primers for RT-PCR or quantitative RT-PCR amplification.

Availability

CoVDB is constructed by the Department of Microbiology, the University of Hong Kong. It is available at no charge at <http://covdb.microbiology.hku.hk>.

ACKNOWLEDGEMENTS

We are grateful to the generous support of Mr Hui Hoy and Mr Hui Ming in the genomic sequencing platform. This work is partly supported by the Research Grant Council Grant; University Development Fund and Outstanding Young Researcher Award, The University of Hong Kong; The Tung Wah Group of Hospitals Fund for Research in Infectious Diseases; the HKSAR Research Fund for the Control of Infectious Diseases of the Health, Welfare and Food Bureau; and the Providence Foundation Limited in memory of the late Dr Lui Hac Minh. Funding to pay the Open Access publication



Blast

Myblast

Formatseq

ORF finder

Translation

Fragment

Blast (Documents of Blast)

Fasta format sequences (you can paste more than one sequence):

Database: Program: E-value:

Filter: Number of Alignments:

Format:

Last updated: Feb 09, 2007. is the recommended web browser for this site!

Copyright (c) 2006 Department of Microbiology, The University of Hong Kong. All rights reserved.

Figure 4. Screenshot of blast similarity search page. Five datasets can be chosen as the database for comparison.

charges for this article was provided by Research Grant Council Grant.

Conflict of interest statement. None declared.

REFERENCES

- Brian,D.A. and Baric,R.S. (2005) Coronavirus genome structure and replication. *Curr. Top. Microbiol. Immunol.*, **287**, 1–30.
- Lai,M.M. and Cavanagh,D. (1997) The molecular biology of coronaviruses. *Adv. Virus Res.*, **48**, 1–100.
- Ziebuhr,J. (2004) Molecular biology of severe acute respiratory syndrome coronavirus. *Curr. Opin. Microbiol.*, **7**, 412–419.
- Woo,P.C., Lau,S.K., Yip,C.C., Huang,Y., Tsoi,H.W., Chan,K.H. and Yuen,K.Y. (2006) Comparative analysis of 22 coronavirus HKU1 genomes reveals a novel genotype and evidence of natural recombination in coronavirus HKU1. *J. Virol.*, **80**, 7136–7145.
- Guan,Y., Zheng,B.J., He,Y.Q., Liu,X.L., Zhuang,Z.X., Cheung,C.L., Luo,S.W., Li,P.H., Zhang,L.J. *et al.* (2003) Isolation and characterization of viruses related to the SARS coronavirus from animals in southern China. *Science*, **302**, 276–278.
- Marra,M.A., Jones,S.J., Astell,C.R., Holt,R.A., Brooks-Wilson,A., Butterfield,Y.S., Khattri,J., Asano,J.K., Barber,S.A. *et al.* (2003) The Genome sequence of the SARS-associated coronavirus. *Science*, **300**, 1399–1404.
- Peiris,J.S., Lai,S.T., Poon,L.L., Guan,Y., Yam,L.Y., Lim,W., Nicholls,J., Yee,W.K., Yan,W.W. *et al.* (2003) Coronavirus as a possible cause of severe acute respiratory syndrome. *Lancet*, **361**, 1319–1325.
- Rota,P.A., Oberste,M.S., Monroe,S.S., Nix,W.A., Campagnoli,R., Icenogle,J.P., Penaranda,S., Bankamp,B., Maher,K. *et al.* (2003) Characterization of a novel coronavirus associated with severe acute respiratory syndrome. *Science*, **300**, 1394–1399.
- Woo,P.C., Lau,S.K., Tsoi,H.W., Chan,K.H., Wong,B.H., Che,X.Y., Tam,V.K., Tam,S.C., Cheng,V.C. *et al.* (2004) Relative rates of non-pneumonic SARS coronavirus infection and SARS coronavirus pneumonia. *Lancet*, **363**, 841–845.
- Fouchier,R.A., Hartwig,N.G., Bestebroer,T.M., Niemeyer,B., de Jong,J.C., Simon,J.H. and Osterhaus,A.D. (2004) A previously undescribed coronavirus associated with respiratory disease in humans. *Proc. Natl Acad. Sci. USA*, **101**, 6212–6216.
- van der Hoek,L., Pyrc,K., Jebbink,M.F., Vermeulen-Oost,W., Berkhout,R.J., Wolthers,K.C., Wertheim-van Dillen,P.M., Kaandorp,J., Spaargaren,J. *et al.* (2004) Identification of a new human coronavirus. *Nat. Med.*, **10**, 368–373.
- Woo,P.C., Huang,Y., Lau,S.K., Tsoi,H.W. and Yuen,K.Y. (2005) In silico analysis of ORF1ab in coronavirus HKU1 genome reveals a unique putative cleavage site of coronavirus HKU1 3C-like protease. *Microbiol. Immunol.*, **49**, 899–908.
- Woo,P.C., Lau,S.K., Chu,C.M., Chan,K.H., Tsoi,H.W., Huang,Y., Wong,B.H., Poon,R.W., Cai,J.J. *et al.* (2005) Characterization and complete genome sequence of a novel coronavirus, coronavirus HKU1, from patients with pneumonia. *J. Virol.*, **79**, 884–895.
- Woo,P.C., Lau,S.K., Tsoi,H.W., Huang,Y., Poon,R.W., Chu,C.M., Lee,R.A., Luk,W.K., Wong,G.K. *et al.* (2005) Clinical and molecular epidemiological features of coronavirus HKU1-associated community-acquired pneumonia. *J. Infect. Dis.*, **192**, 1898–1907.
- Woo,P.C., Lau,S.K., Li,K.S., Poon,R.W., Wong,B.H., Tsoi,H.W., Yip,B.C., Huang,Y., Chan,K.H. *et al.* (2006) Molecular diversity of coronaviruses in bats. *Virology*, **351**, 180–187.
- Lau,S.K., Woo,P.C., Li,K.S., Huang,Y., Wang,M., Lam,C.S., Xu,H., Guo,R., Chan,K.H. *et al.* (2007) Complete genome sequence of bat coronavirus HKU2 from Chinese horseshoe bats revealed a much smaller spike gene with a different evolutionary lineage from the rest of the genome. *Virology*, doi: 10.1016/j.virol.2007.06.009.
- Tang,X.C., Zhang,J.X., Zhang,S.Y., Wang,P., Fan,X.H., Li,L.F., Li,G., Dong,B.Q., Liu,W. *et al.* (2006) Prevalence and genetic

- diversity of coronaviruses in bats from China. *J. Virol.*, **80**, 7481–7490.
18. Woo, P.C., Wang, M., Lau, S.K., Xu, H., Poon, R.W., Guo, R., Wong, B.H., Gao, K., Tsoi, H.W. *et al.* (2007) Comparative analysis of twelve genomes of three novel group 2c and group 2d coronaviruses reveals unique group and subgroup features. *J. Virol.*, **81**, 1574–1585.
 19. Lau, S.K., Woo, P.C., Li, K.S., Huang, Y., Tsoi, H.W., Wong, B.H., Wong, S.S., Leung, S.Y., Chan, K.H. *et al.* (2005) Severe acute respiratory syndrome coronavirus-like virus in Chinese horseshoe bats. *Proc. Natl Acad. Sci. USA*, **102**, 14040–14045.
 20. Cavanagh, D., Mawditt, K., Welchman Dde, B., Britton, P. and Gough, R.E. (2002) Coronaviruses from pheasants (*Phasianus colchicus*) are genetically closely related to coronaviruses of domestic fowl (infectious bronchitis virus) and turkeys. *Avian Pathol.*, **31**, 81–93.
 21. East, M.L., Moestl, K., Benetka, V., Pitra, C., Honer, O.P., Wachter, B. and Hofer, H. (2004) Coronavirus infection of spotted hyenas in the Serengeti ecosystem. *Vet. Microbiol.*, **102**, 1–9.
 22. Jonassen, C.M., Kofstad, T., Larsen, I.L., Lovland, A., Handeland, K., Follestad, A. and Lillehaug, A. (2005) Molecular identification and characterization of novel coronaviruses infecting graylag geese (*Anser anser*), feral pigeons (*Columba livia*) and mallards (*Anas platyrhynchos*). *J. Gen. Virol.*, **86**, 1597–1607.
 23. Liu, S., Chen, J., Chen, J., Kong, X., Shao, Y., Han, Z., Feng, L., Cai, X., Gu, S. *et al.* (2005) Isolation of avian infectious bronchitis coronavirus from domestic peafowl (*Pavo cristatus*) and teal (*Anas*). *J. Gen. Virol.*, **86**, 719–725.
 24. Benson, D.A., Karsch-Mizrachi, I., Lipman, D.J., Ostell, J. and Wheeler, D.L. (2007) GenBank. *Nucleic Acids Res.*, **35**, D21–D25.
 25. Snyder, E.E., Kampanya, N., Lu, J., Nordberg, E.K., Karur, H.R., Shukla, M., Soneja, J., Tian, Y., Xue, T. *et al.* (2007) PATRIC: the VBI PathoSystems Resource Integration Center. *Nucleic Acids Res.*, **35**, D401–D406.
 26. Pruitt, K.D., Tatusova, T. and Maglott, D.R. (2007) NCBI reference sequences (RefSeq): a curated non-redundant sequence database of genomes, transcripts and proteins. *Nucleic Acids Res.*, **35**, D61–D65.
 27. Altschul, S.F., Gish, W., Miller, W., Myers, E.W. and Lipman, D.J. (1990) Basic local alignment search tool. *J. Mol. Biol.*, **215**, 403–410.