

APPENDIX

A MAIN DATA STRUCTURE

Successor tables (ST). The successor tables $\{ST_1, ST_2, \dots, ST_d\}$ of the sequence set $T = \{S_1, S_2, \dots, S_d\}$ are built to support the compression of the data and quick search for the immediate successors of the points. For a sequence $S_l = x_1, x_2, \dots, x_n$ from the sequence set T over a finite alphabet $\Sigma = (c_1, c_2, \dots, c_k)$, its successor table ST_l is a two-dimensional array, where $ST_l[i, j]$ (the element of the i th row and the j th column) is defined as

$$ST_l[i, j] = \min\{r | x_r = c_i, r \geq 1, r \geq j, 1 \leq i \leq |\Sigma|, 0 \leq j \leq n\} \quad (6)$$

From Eq. 6, we can see that $ST_l[i, j]$ denotes the minimal position r (the r th character position) of the sequence S_l with $x_r = c_i$ after position j . See the examples in Fig. 6.

$S_1 =$	T	G	A	C	G	A	T	C
0	1	2	3	4	5	6	7	8
A	3	3	3	6	6	6	—	—
C	4	4	4	4	8	8	8	—
G	2	2	5	5	5	—	—	—
T	1	7	7	7	7	7	7	—

(a) The Successor Table ST_1

$S_2 =$	A	T	G	C	T	C	A	G
0	1	2	3	4	5	6	7	8
A	1	7	7	7	7	7	—	—
C	4	4	4	4	6	6	—	—
G	3	3	3	8	8	8	8	—
T	2	2	5	5	5	—	—	—

(b) The Successor Table ST_2

$S_3 =$	C	T	A	G	T	A	C	G
0	1	2	3	4	5	6	7	8
A	3	3	3	6	6	6	—	—
C	1	7	7	7	7	7	—	—
G	4	4	4	4	8	8	8	—
T	2	2	5	5	5	—	—	—

(c) The Successor Table ST_3

Figure 6: The constructed successor tables ST_1, ST_2 and ST_3 corresponding to the sequences S_1, S_2 and S_3 (in the paper), where "—" indicates \emptyset .

The set S_{suc} of immediate successors of a d -dimensional point $p = (p_1, p_2, \dots, p_d)$ can be obtained efficiently in $O(d|\Sigma|)$ time. For a d -dimensional point p , the operation for producing its S_{suc} can be characterized by Eq. 7.

$$S_{suc} = \{(ST_1[i', p_1], ST_2[i', p_2], \dots, ST_d[i', p_d])\} \\ s.t. 1 \leq i' \leq |\Sigma|, \forall ST_l[i', p_l] \neq \emptyset, 1 \leq l, i \leq d \quad (7)$$

For example, for the dominant $(2, 3, 4)$ of the sequences S_1, S_2 and S_3 (in the paper), we can couple the corresponding rows 1-4 of the second, third and forth columns from the successor tables ST_1, ST_2 and ST_3 to obtain all its immediate successors $(3, 7, 6)$, $(4, 4, 7)$, $(5, 8, 8)$ and $(7, 5, 5)$ corresponding to the characters A, C, G, and T, respectively. There is no immediate successor for the dominant $(6, 7, 3)$ due to the coupling results $(_, _, 6)$, $(8, _, 7)$, $(_, 8, 4)$ and $(7, _, 5)$, which indicates none of the points is an immediate successor according to Eq. 7.

B EXPERIMENTAL RESULTS

We evaluate the performance of the approximate algorithms, whose performances vary in terms of not only efficiency but also precision. Here, the precision is measured by Eq. 3. The results for all the tested approximate algorithms, including the state-of-the-art *CRO*, *SA_MLCS* as well as our algorithm *HA_MLCS* are shown in Tables 3 and 4.

Table 3: Precisions (P) and running times (T) of *CRO* (A1), *SA_MLCS* (A2) and *HA_MLCS* (A3) for 5 sequences with various lengths ($|S_i|$).

Σ =4							Σ =20						
S _f	A1		A2		A3 (<i>m</i> = 100)		S _f	A1		A2		A3 (<i>m</i> = 100)	
	<i>P</i>	T(s)	<i>P</i>	T(s)	<i>P</i>	T(s)		<i>P</i>	T(s)	<i>P</i>	T(s)	<i>P</i>	T(s)
1.0E+3	0.524	0.08	0.822	0.97	0.988	0.71	1.0E+3	0.411	0.09	0.893	7.43	0.992	0.63
2.0E+3	0.496	0.11	0.818	1.94	0.981	1.42	2.0E+3	0.366	0.12	0.877	17.31	0.981	1.26
5.0E+3	0.391	1.01	0.792	2.41	0.985	1.35	5.0E+3	0.357	5.79	0.857	47.99	0.978	2.90
1.0E+4	0.347	3.48	0.743	9.33	0.975	5.50	1.0E+4	0.334	2.27	0.822	113.70	0.967	4.35
2.0E+4	0.322	5.86	0.717	32.59	0.972	9.22	2.0E+4	0.327	8.67	0.802	239.00	0.971	8.12
5.0E+4	0.314	33.01	0.701	87.06	0.968	27.56	5.0E+4	0.305	53.94	0.791	541.80	0.965	18.79
1.0E+5	0.296	133.20	0.678	156.79	0.959	52.86	1.0E+5	0.279	218.00	0.751	5272.00	0.952	39.87
1.0E+7	+	+	0.269	1609.00	0.951	557.20	1.0E+7	+	+	0.741	44190.00	0.962	403.60
1.0E+8	+	+	+	+	0.948	3430.00	1.0E+8	+	+	+	+	0.963	2985.00

Symbol '+' indicates the memory overflow leading to calculating failure.

Table 4: Precisions (P) and running times (T) of *CRO* (A1), *SA_MLCS* (A2) and *HA_MLCS* (A3) for d sequences with lengths ($|S_i|$) 1000 and 2000, respectively.

$ \Sigma =4, S_I =1000$							$ \Sigma =20, S_I =2000$						
d	A1		A2		A3 ($m = 100$)		d	A1		A2		A3 ($m = 100$)	
	P	$T(s)$	P	$T(s)$	P	$T(s)$		P	$T(s)$	P	$T(s)$	P	$T(s)$
10	0.689	0.03	0.792	0.62	0.988	0.41	10	0.698	0.03	0.801	0.12	0.987	0.08
50	0.667	0.06	0.756	0.82	0.977	0.51	50	0.671	0.06	0.811	0.11	0.971	0.20
100	0.645	0.08	0.758	0.97	0.975	0.63	100	0.655	0.07	0.801	0.10	0.978	0.13
400	0.632	0.17	0.742	1.02	0.973	0.74	400	0.589	0.17	0.812	0.23	0.975	0.19
800	0.617	0.19	0.738	1.33	0.973	0.88	1000	0.623	0.35	0.805	0.32	0.972	0.25
1000	0.602	0.26	0.717	2.26	0.969	1.24	5000	0.611	1.51	0.798	0.41	0.966	0.35
3000	0.587	0.45	0.703	3.11	0.967	1.64	10000	0.594	2.98	0.785	0.43	0.966	0.41
5000	0.577	0.70	0.695	3.57	0.965	1.76	14000	0.568	4.26	0.773	0.50	0.964	0.47
10000	0.534	1.11	0.686	3.83	0.959	2.18	20000	0.536	5.57	0.758	1.21	0.959	0.73

C THE EVOLUTION TREES OF COVID-19 VIRUSES FROM THE USA AND ENGLAND

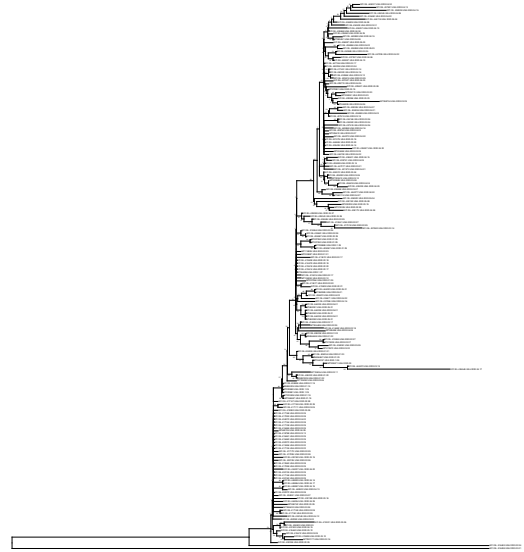


Figure 7: The evolution tree of COVID-19 viruses in human hosts from USA.

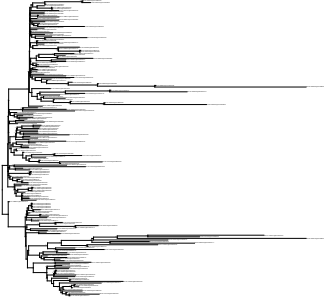


Figure 8: The evolution tree of COVID-19 viruses in human hosts from England.

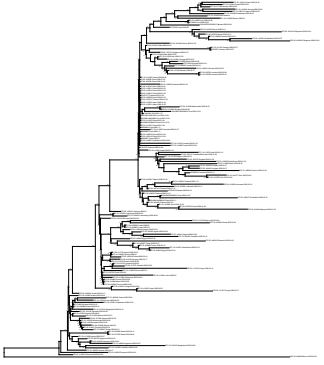


Figure 9: The evolution tree of COVID-19 viruses in human hosts from different 79 countries.

D ALGORITHM PSEUDOCODE

Algorithm 1 HA-MLCS(T, Σ)

```

1:  $ST \leftarrow$  Construct Successor Tables of sequence set  $T$  with  $\Sigma$ ;
2:  $k \leftarrow 0$ ;  $MLCS\text{-}ODAG \leftarrow \emptyset$ ;  $D^k \leftarrow \{(0, 0, \dots, 0)\}$ ;
3:  $D_{init}^{k+1} \leftarrow \text{Successor}(D^k, ST)$ ; //calculate successor point set  $D_{init}^{k+1}$  of  $D^k$ 
4: while  $D_{init}^{k+1} \neq \emptyset$  do //Constructing  $MLCS\text{-}ODAG$  with layer by layer
5:    $(D_{init}^{k+1})_{1st} \leftarrow \text{BestNondominatedSorting}(D_{init}^{k+1})$ ;
6:   Calculate  $\text{score}(p)$  by Eq. 1,  $p \in (D_{init}^{k+1})_{1st}$ ;
7:    $D^{k+1} \leftarrow$  Keep top  $m$  points with the minimum score in  $(D_{init}^{k+1})_{1st}$ ;
8:    $MLCS\text{-}ODAG \leftarrow MLCS\text{-}ODAG \cup D^{k+1}$ ;  $k \leftarrow k + 1$ ;
9:    $D_{init}^{k+1} \leftarrow \text{Successor}(D^k, ST)$ ;
10: end while
11:  $\text{maxlevel} \leftarrow k - 1$ ;  $k \leftarrow 0$ ;  $D^0 \leftarrow \{(\infty, \infty, \dots, \infty)\}$ ;
12: while  $D^k \neq \emptyset$  do //Algorithm BackwardTopSorting
13:    $D^{k+1} \leftarrow \emptyset$ ;
14:   for  $q \in D^k$  do
15:     for  $p \in \text{precursor}[q]$  do
16:       if  $\text{tlevel}[p] + k \neq \text{maxlevel}$  then
17:         Delete  $p$  from  $MLCS\text{-}ODAG$ ;
18:       else
19:          $D^{k+1} \leftarrow D^{k+1} \cup \{p\}$ ;
20:       end if
21:     end for
22:   end for
23:    $k \leftarrow k + 1$ ;
24: end while
25: return all of the  $MLCS$ s of sequence set  $T$ ;

```

The proposed algorithm *HA-MLCS* is implemented in Java JDK1.8. Where m is a user-customized parameter ($1 \leq m \in \mathbb{Z}$), which represents how many number of key points to be retained in each layer

when constructing *MLCS-ODAG*. The source code of algorithm *HA-MLCS* is available at: <https://github.com/HA-MLCS/HA-MLCS>

E THE SIMILARITY BETWEEN COVID-19 AND RELATED VIRUSES

Table 5: The similarity matrices.

LCS-Based											
	HCov229E	HCovHKU1	HCovNL63	SARS	Lassa	MERS	Victoria	Yamagata	Ebola	Dengue	Environment
2019.12	0.652594	0.691539	0.658972	0.678047	0.626388	0.340651	0.694165	0.425849	0.425918	0.499862	0.340251
2020.01	0.652732	0.691539	0.659161	0.680932	0.626833	0.340785	0.694567	0.426070	0.426050	0.499955	0.340377
2020.02	0.652867	0.691513	0.659455	0.681191	0.626881	0.341032	0.694733	0.426254	0.426232	0.500041	0.340410
2020.03	0.654415	0.689552	0.660277	0.681593	0.627379	0.342154	0.694857	0.427542	0.427520	0.501517	0.341793
2020.04	0.656404	0.679613	0.651387	0.687941	0.627384	0.342254	0.694740	0.428367	0.428367	0.503008	0.342201
2020.05	0.658049	0.680861	0.658063	0.688351	0.628885	0.344117	0.695210	0.430231	0.429576	0.503550	0.343054
Average	0.653533	0.688998	0.658004	0.683167	0.626620	0.341829	0.694634	0.427387	0.427247	0.501289	0.341379
LD-Based											
	HCov229E	HCovHKU1	HCovNL63	SARS	Lassa	MERS	Victoria	Yamagata	Ebola	Dengue	Environment
2019.12	0.558677	0.561025	0.560186	0.561579	0.560521	0.338617	0.574539	0.418381	0.418381	0.471776	0.338617
2020.01	0.558679	0.561025	0.560202	0.561598	0.560544	0.338750	0.575088	0.418465	0.418508	0.471802	0.338617
2020.02	0.558685	0.561048	0.560880	0.561874	0.560624	0.339674	0.576760	0.418670	0.418719	0.472045	0.338617
2020.03	0.558188	0.561446	0.560093	0.562361	0.560700	0.341071	0.571322	0.419872	0.419881	0.472709	0.340589
2020.04	0.558583	0.561800	0.560262	0.562394	0.560590	0.341193	0.571710	0.422083	0.422847	0.472724	0.341857
2020.05	0.558189	0.560186	0.558221	0.562513	0.560599	0.342833	0.571981	0.423131	0.423201	0.472724	0.342833
Average	0.558297	0.561133	0.560212	0.562070	0.560525	0.340740	0.571133	0.420991	0.420257	0.471872	0.340486

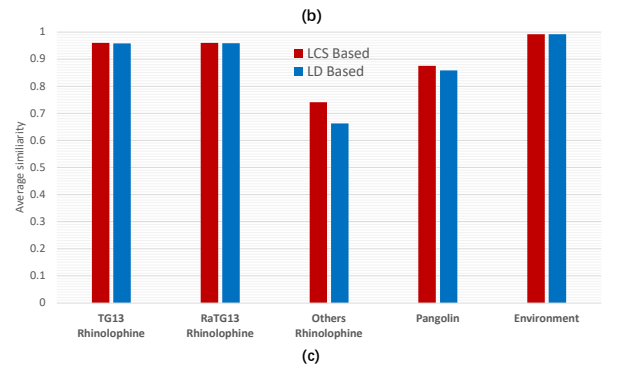
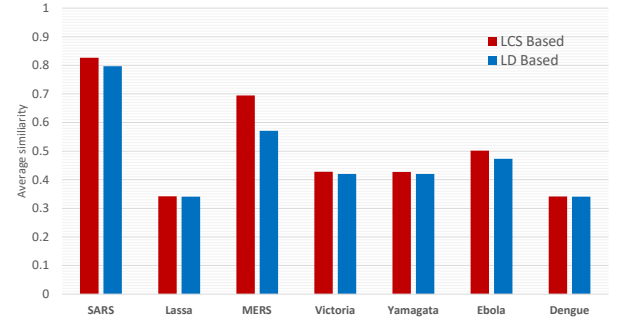
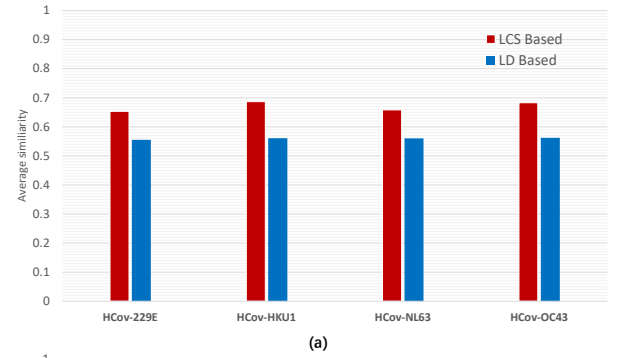


Figure 10: The schematic diagram between COVID-19 strains and other viruses.