

“CGG-CGG”: huella genética del origen de laboratorio del SARS-CoV-2

Antonio R. Romeu and Enric Ollé

AR: es Químico. Catedrático de Bioquímica y Biología Molecular de la Universidad Rovira i Virgili. Tarragona. España.

EO: es Veterinario y Bioquímico. Profesor Asociado del Departamento de Bioquímica y Biotecnología de la Universidad Rovira i Virgili. Tarragona. España.

El origen del virus es el gran debate, que se puede abordar desde muchos aspectos (1). En presente artículo lo abordamos exclusivamente desde la Ciencia, desde el análisis genético y de genomas.

El SARS-CoV-2 es un Coronavirus. En este grupo hay cuatro subgrupos (técnicamente, Géneros): Alfa-coronavirus, Beta-coronavirus, Gamma-coronavirus, Delta-coronavirus. El SARS-CoV-2 pertenece al grupo de Beta-coronavirus. Concretando más, dentro de los Beta-coronavirus, todavía hay cuatro subgrupos (subgéneros): Linaje A, Linaje B, Linaje C, y Linaje D. El SARS-CoV-2 pertenece al Linaje B (técnicamente, su subgénero se llama Sarbecovirus).

Es decir, el SARS-CoV-2 es un Coronavirus → Betacoronavirus → Sarbecovirus (2)

La característica fundamental del SARS-CoV-2 es que su proteína S, la que interacciona con las células humanas y produce la infección, ha adquirido un inserto con carga positiva en una región estratégica de la proteína, que facilita enormemente la infección y la propagación del virus de persona a persona (2,3,4). Es decir, un ancestro del virus, dentro del murciélago (o en otro animal intermedio) adquirió dicho inserto, y se convirtió en el actual SARS-CoV-2; después saltó del animal a las personas.

Así pues, de una forma reduccionista, se puede decir, que el origen del SARS-CoV-2 es el propio origen del inserto positivo en su proteína S (5).

Técnicamente, el inserto positivo se llama “sitio polibásico de escisión de furina”. La furina es una proteína de la membrana de nuestras células, que reconoce dicho inserto y colabora con los receptores correspondientes. En estas condiciones, la entrada del virus en nuestro organismo es por la puerta grande.

Estudiando el inserto del SARS-CoV-2, hay características que no encajan, y es lo que hizo pensar desde el principio que su origen pueda ser de laboratorio. El SARS-CoV-2, es el único miembro de su grupo (Sarbecovirus), que tiene este sitio de furina. Esto ya lo hace evolutivamente diferente, y desde el primer día, la pregunta del millón es ¿cómo le ha llegado?

En un análisis detallado de la composición del inserto, se observa que hay dos elementos positivos seguidos, llamados “arginina-arginina” (el símbolo de la arginina es una “R”; es decir, hay un dímero RR).

Un principio básico de la Biología o de la Genética es que cada elemento de las proteínas está codificado por 3 elementos (un triplete) del ADN o de los genes. Esto es el “Codigo Genético Universal”. En el caso concreto del elemento “arginina” (R), en los genes existen 6 grupos de 3 elementos, que codifican por igual a la R. Estos 6 elementos no se usan en la misma proporción en todas las especies. La Tabla 1 presenta los 6 tripletes que codifican R, y las frecuencias en que estos tripletes se utilizan en el SARS-CoV-2 (6) y en la especie humana (7). En la Tabla 1 se hace hincapié en que el uso de los 6 tripletes no es aleatorio como si se tratase de las caras de un dado.

Arginina (R)	AGA	AGG	CGA	CGC	CGG	CGT	
Frecuencia aleatoria:	1/6	1/6	1/6	1/6	1/6	1/6	
SARS-CoV-2 (%):	45	13	5	10	3	24	
Genoma humano (%):	20	20	11	19	21	9	

Tabla 1. Tripletes que codifican arginina, “R” (Código Genético Universal).

Volviendo al inserto positivo del SARS-CoV-2, su dímero RR está codificado por los tripletes CGG-CGG. Sorprendentemente, este triplete “CGG” sea el menos frecuente en el SARS-CoV-2 (3%) y el más utilizado en el genoma humano (21%). Aún sorprende más, que en el virus pandémico haya dos tripletes minoritarios seguidos. Por otra parte, y atendiendo a los datos de la Tabla 1, sea cual sea el origen de la secuencia CGGCGG, lo que es cierto es que está optimizado para las personas (8).

Es importante saber que el sitio de furina (inserto positivo), aunque no lo tienen los virus del grupo del SARS-CoV-2 (excepto él), es muy común en el mundo de los virus (9). Otros Coronavirus, lo tienen en su proteína S (por ejemplo, el MERS). Fuera del grupo de los Coronavirus, y dentro del grupo de los virus de la gripe; algunos virus de la gripe con la proteína H5, también tienen el inserto de furina. En general, muchos sitios de furina de los virus, comparten la presencia de un doblete RR (10).

En un estudio previo, estudiamos bioinformáticamente como estaban codificados los pares RR de los sitios de furina de toda clase de virus, en una muestra grande y representativa. Sorprendentemente, ningún par RR está codificado por la secuencia CGGCGG (6). De acuerdo con nuestros resultados, un estudio reciente describe que la secuencia CGG-CGG nunca se ha encontrado de forma natural en los Coronavirus (11).

En otro estudio (en proceso de publicación) hemos analizado los pares RR en las otras proteínas del SARS-CoV-2. El SARS-CoV-2 tiene 24 proteínas. En la proteína S, sólo hay un par RR (el del inserto). Pero, en las otras proteínas ¿tienen pares RR? Sólo en unas pocas, y no están codificados por la secuencia CGGCGG. La Tabla 2 presenta las proteínas del SARS-CoV-2 con dímeros RR, sus posiciones, y como están codificados.

Gen/proteína	Position en el genoma	RR	Posición	Código genético	Posición
<i>nsp3</i>	2720-8553	GVRR	1614	GGTGT T AGAAG G	4842
<i>nsp4</i>	8555-10054	RFRR	306	AGGTT T AGAAG A	918
		LKRR	401	CTAAAGAGAC G T	1203
<i>nsp6</i>	10973-11842	GARR	138	GGTGCTAGG A GA	414
<i>nsp13</i>	16237-18039	CIRR	22	TGCATAC G TAG A	66
		TCRR	443	ACTTGT C GG C GT	1329
		IPRR	595	ATTCCAC G TAG G	1785
<i>nsp14A2</i>	18040-19620	TYRR	53	ACCTATAG A AG A	159
		CDRR	213	TGTGATAG A CG T	639
S	21563-25384	SPRR	683	TCTCCTCGGCGG	2049
<i>ORF9 (N)</i>	28274-29533	KQRR	41	AAACAAC G T C GG	123
<i>Nucleocapsida</i>		YYRR	89	TACTAC C GAAG A	267
		ATRR	93	GCTACCAG A CG A	279
		FGRR	277	TTCGGCAG A CG T	831

Tabla 2. Proteínas del SARS-CoV-2 que tienen pares RR. La información relativa al inserto positivo de la proteína S se resalta en rojo.

La siguiente pregunta fue ¿los pares RR que se encuentran en el SARS-CoV-2 (a parte del RR del sitio de furina), se encuentran también en los virus de murciélago conocidos que son sus “parientes más cercanos”? La respuesta fue rotundamente sí. También fue sorprendente, que ninguno de ellos está codificado por la secuencia CGGCGG. Es más, muy pocos de los pares RR están codificados por el mismo triplete dos veces seguidas. Solo en unos pocos dobletes se repite el código “AGA” (AGAAGA), que es el mayoritario (45%) en el SARS-CoV-2 (Tabla 1). La Figura del Anexo 2 presenta, el detalle los dobletes RR de los Coronavirus estudiados y como están codificados.

Todo ellos nos lleva al origen del SARS-CoV-2. Frente a este gran enigma, hay que formularse la pregunta correcta: ¿cuál es el origen de la secuencia CGGCGG que codifica el inserto de la proteína S del Coronavirus?

Un origen natural debe contemplar dos posibles eventos (naturales): (i) mutaciones de inserción aleatorias en un antepasado próximo, en el murciélago (u otro animal), originando el SARS-CoV-2 que conocemos ; o (ii) relaciones de intercambio de material genético entre un antepasado próximo y otros virus, en el murciélago (u otro animal), dando lugar, también, al SARS-CoV-2 que conocemos.

Una inserción al azar, relativamente larga, como la incluye la secuencia CGGCGG que codifica el dímero RR en cuestión, es muy improbable. En la reproducción de estos Coronavirus (replicación) hay mecanismos de corrección que evitan grandes errores, como hubiese sido la inserción que aparece en el SARS-CoV-2. Es muy poco probable que por azar, dos tripletes muy poco frecuentes del virus aparezcan, nada más ni nada menos, que seguidos. Serían dos procesos de inserción independientes que tendrían que haber ocurrido simultáneamente en el tiempo y en el espacio. Otro proceso independiente, ocurrido también en el mismo tiempo y espacio, es que la secuencia insertada en el curso de la evolución de un antepasado del SARS-CoV-2, sea la secuencia más frecuente en las personas (Tabla 1). En cuestión de origen de las especies, no hay nada imposible, pero este mecanismo de mutación aleatoria, como origen del virus, parte de la comunidad científica no lo acepta; de lo contrario, no habría este debate sobre el origen.

La segunda posibilidad de un origen natural se basa en un intercambio de material genético con otro virus (recombinación). Esta posibilidad es la más lógica. No obstante, si otro virus ha transferido el inserto; este virus donador tendría que tener el código CGGCGG, y de momento, no se conoce ningún dímero RR de sitio de furina de otros virus con dicho código. Es decir, la recombinación, como mecanismo molecular, se acepta, pero hasta que no se descubra un virus donador, queda abierta.

Finalmente, un origen de laboratorio de la secuencia CGGCGG del SARS-CoV-2 no es imposible. Los detractores alegan que un origen de laboratorio siempre deja huella, de ser así, la propia secuencia CGGCGG es la huella.

Referencias

1. Nicholas Wade. Origin of Covid — Following the Clues. Acceso el 03/07/2021. <https://nicholaswade.medium.com/origin-of-covid-following-the-clues-6f03564c038>.
2. Philip V'kovski, Annika Kratzel, Silvio Steiner, Hanspeter Stalder, Volker Thiel. Coronavirus biology and replication: implications for SARS-CoV-2. Nat. Rev. Microbiol. Oct 28;1-16, 2020. PMID: 33116300. doi: 10.1038/s41579-020-00468-6.
3. Britt Glaunsinger. Coronavirus biology. The second lecture in the COVID-19, SARS-CoV-2 and the Pandemic Series. University of California, Berkeley. 2020. Accessed June 23, 2021. <https://www.youtube.com/watch?v=r2mOU2qOCYs>.

4. Kristian G. Andersen, Andrew Rambaut, W Ian Lipkin, Edward C Holmes, Robert F Garry. The proximal origin of SARS-CoV-2. *Nat. Med.* 26:450-452, 2020. PMID: 32284615. doi: 10.1038/s41591-020-0820-9.
5. Manuel Ansede, Artur GAlocha y Mariano Zafra. ccu cgg cgg gca, thw 12 leters that changed the world. *El País*, 19-mayo-2020. Acceso el 03/07/2021.
https://elpais.com/elpais/2020/05/18/ciencia/1589818040_544543.html
6. Antonio R. Romeu, Enric Ollé. SARS-CoV-2 and the Secret of the Furin Site. *Preprints 2021*, 2021020264 (doi: 10.20944/preprints202102.0264.v1).
7. GenScript Codon Usage Frequency Table(chart) Tool. Accessed June 23, 2021.
<https://www.genscript.com/tools/codon-frequency-table>.
8. Antonio R. Romeu, Enric Ollé. SARS-CoV-2 Origin: An Affair of Codons?. *Preprints 2021*, 2021060121 (doi: 10.20944/preprints202106.0121.v1).
9. Elisabeth Braun, Daniel Sauter. Furin-mediated protein processing in infectious diseases and cancer. *Clin. Transl. Immunol.* E1073, 2019. PMID: 31406574. doi.org/10.1002/cti2.1073.
10. Imène Kara, Marjorie Poggi, Bernadette Bonardo, Roland Govers, Jean-François Landrier, Sun Tian, Ingo Leibiger, Robert Day, John W M Creemers, Franck Peiretti. The Paired Basic Amino Acid-cleaving Enzyme 4 (PACE4) Is Involved in the Maturation of Insulin Receptor Isoform B. *J. Biol. Chem.* 290:2812-2821, 2015. PMID: 25527501. doi: 10.1074/jbc.M114.592543.
11. Steven Quay, Richard Muller. The Science suggests a Wuhan lab leak. *The Wall Street Journal*, Monday, June 7, 2021. A17.

Anexos

Anexo 1. Metodología (lenguaje técnico)

Anexo 2 (en inglés). Figura de múltiples páginas

Se presentan fragmentos de los alineamientos múltiples de las secuencias de las proteínas y los genes de los Coronavirus del grupo del SARS-CoV-2, con presencia de dobles RR.

Para cada proteína, en primer lugar se presenta una figura apaisada a fin de localizarla en el genoma del virus.

En los fragmentos de los alineamientos múltiples, los dobles RR y sus códigos se resaltan en amarillo. Los números a la derecha indican las posiciones en la secuencia.

Anexo 1

Metodología

La fuente de información han sido las bases de datos NCBI GenBank (<https://www.ncbi.nlm.nih.gov/nucleotide/>) y GISAID (<https://www.gisaid.org/>).

El genoma de SARS-CoV-2 de referencia fue del Coronavirus WIV04, Wuhan-Hu-1 (NCBI, número de GenBank NC_045512.2 y GISAID, EPI_ISL_402124).

Coronavirus (Sarbecovirus) de la muestra de trabajo:

- Humanos: SARS-CoV-2 (MN996528.1, MT159709.2, MT066156.1), SARS-CoV (HKU-39849)
- De hurón: SARS-CoV (Tor2/FP1-10912)
- De murciélago: RatG13, RmYN02, RShsTT200, RShsTT182, RpYN06, RacCS203, PrC31, VZC45, ZXC21, BM48-31 y BtKY72.
- De pangolín: MP789, GX-P5L, GX-P4L, GX-P1E, GX-P5E y GX-P2V.

Plan de trabajo:

- Primero, identificación de las proteínas de SARS-CoV-2 con presencia de dímeros RR.
- Segundo, creación de una base de datos interna con dichas proteínas y sus genes, de los Coronavirus de la muestra de trabajo.
- Tercero, alineamientos múltiples, tanto de las secuencias de proteínas como de genes, mediante el programa EMBL Clustal Omega (v.1.2.4) (<https://www.ebi.ac.uk/Tools/msa/clustalo/>).

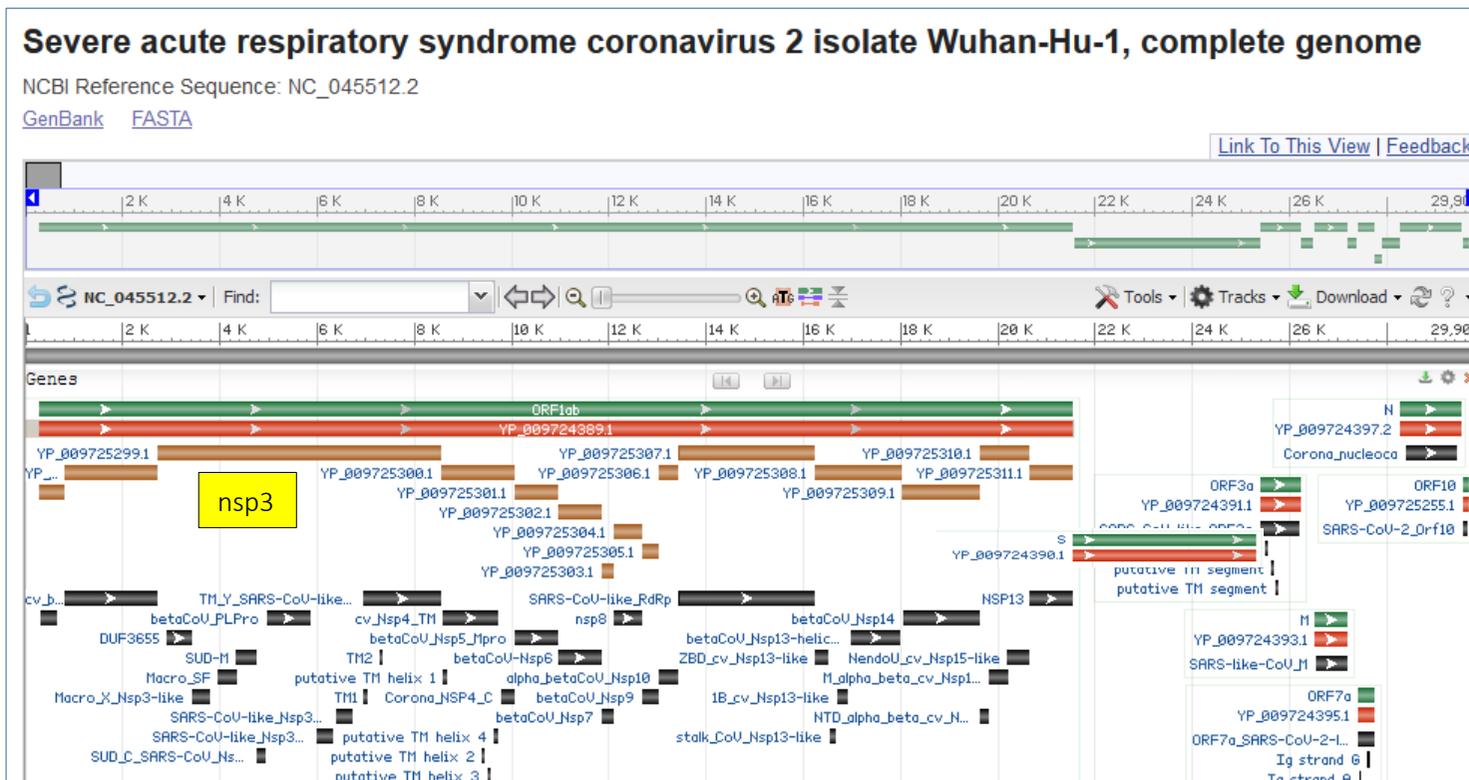
Dado que no todos los genomas se anotaron completamente, las secuencias de genes se obtuvieron mediante análisis BLASTn sucesivos, utilizando los genes de referencia del SARS-CoV-2 como query contra los genomas de Sarbecovirus de la muestra.

Anexo 2

Figura (varias páginas) de fragmentos de los alineamientos múltiples de las proteínas y genes de los Coronavirus de la muestra, siguiendo el orden de la Tabla 2.

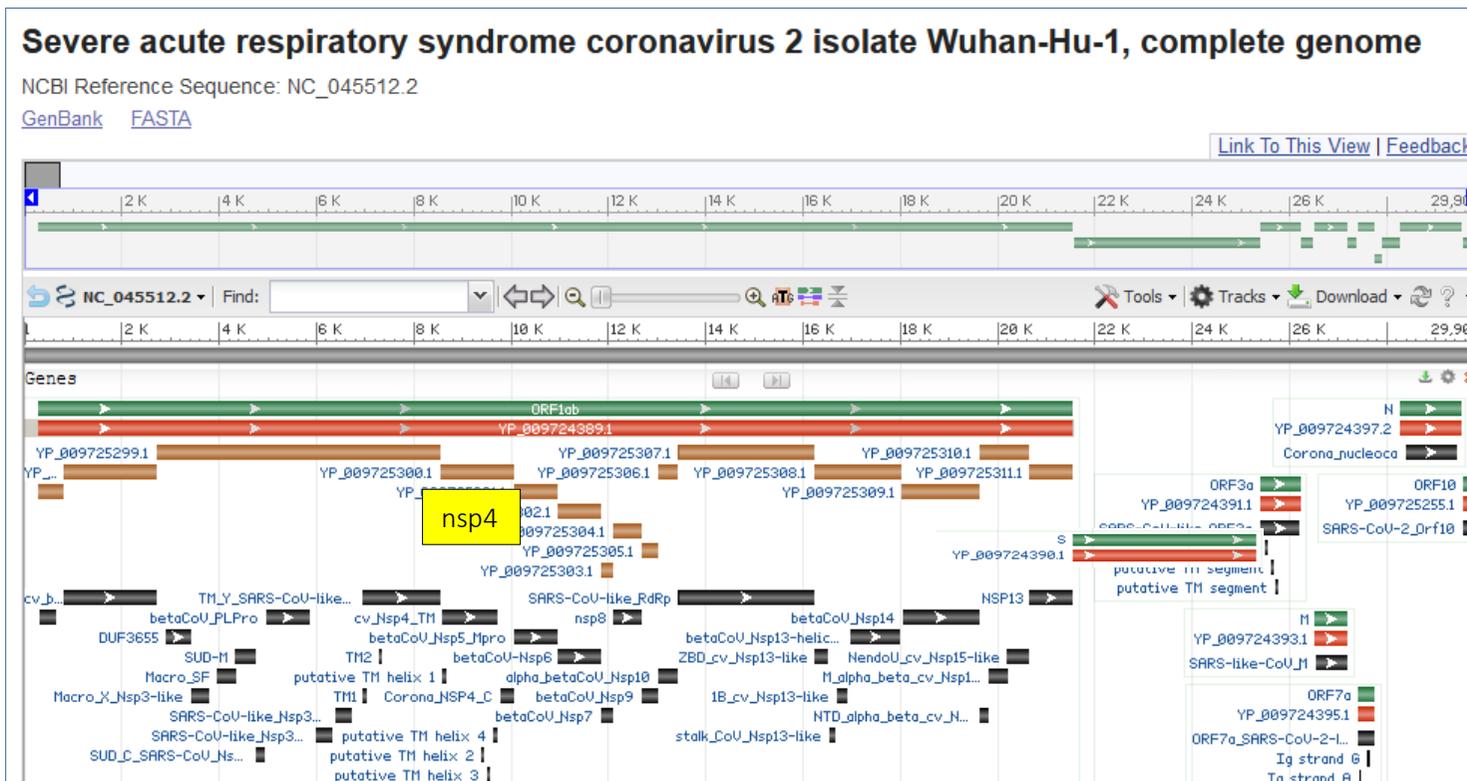
nsp3

Papain-like proteinase. Double-membrane vesicles formation. Macro domains



nsp4

Transmembrane domains. Double-membrane vesicles formation



Protein nsp4 (the RR pair 2)

BM48-31	AHLQWLAMFSPVFPFMITVYVICISTKHCHWFFSNYLRKRVVFN	358
PCoV_GX-P1E	AHVQWMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	397
PCoV_GX-P4L	AHVQWMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	397
PCoV_GX-P2V	AHVQWMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	397
PCoV_GX-P5L	AHVQWMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	397
PCoV_GX-P5E	AHVQWMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	397
ZC45	AHIQWVMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	420
ZXC21	AHIQWVMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	420
RacCS203	AHIQWVMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	420
PCoV-MP789	AHIQWVMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	420
RShSTT200	AHIQWVMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	420
RShSTT182	AHIQWVMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	420
INMI1	AHIQWVMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	420
CDC-CruiseA-12	AHIQWVMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	420
WIV04	AHIQWVMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	420
Wuhan-Hu-1	AHIQWVMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	420
RmYN02	AHIQWVMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	420
RaTG13	AHIQWVMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	420
RpYN06	AHIQWVMVMFTPLVPFWITIVYVICISTKHCHWFFSNYLRKRVVFN	420

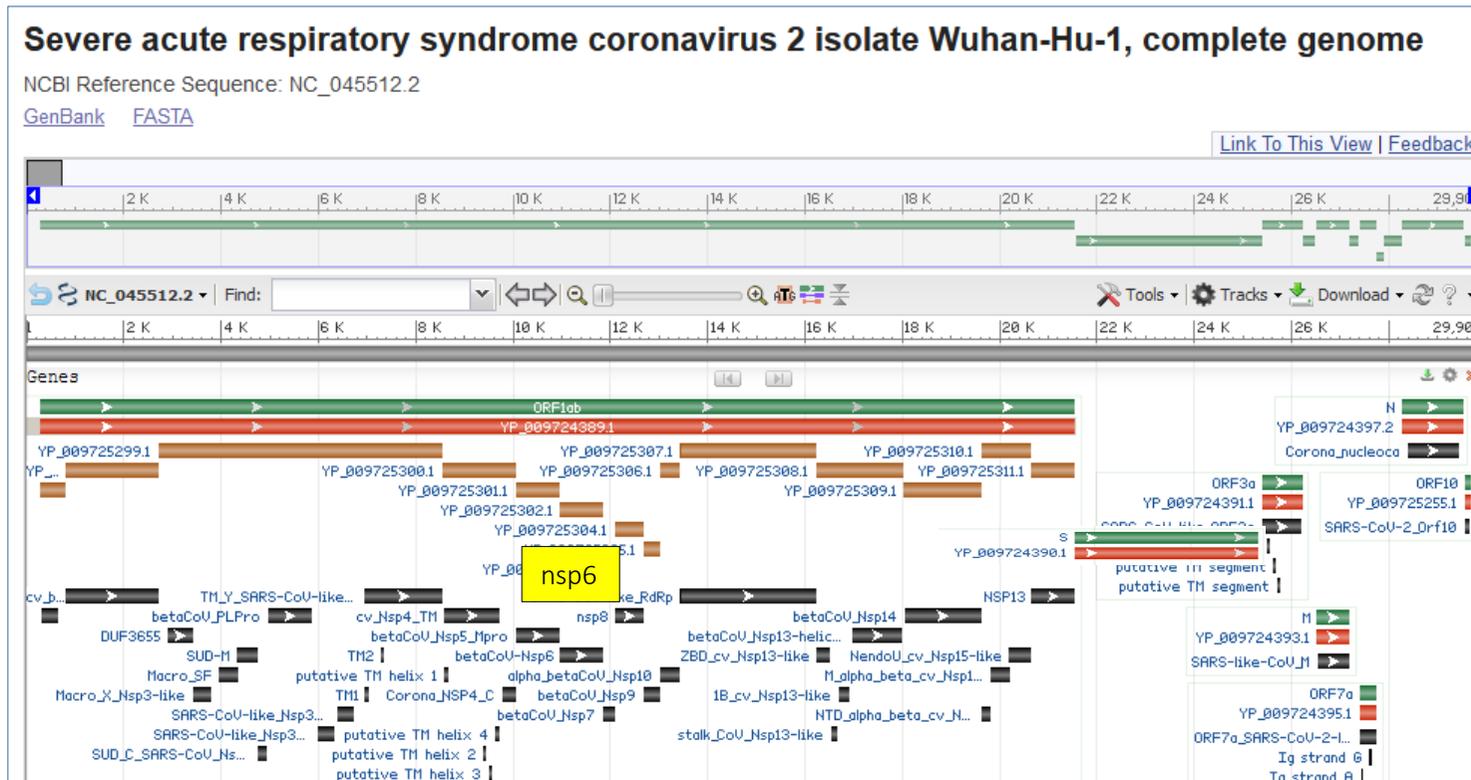
Gene nsp4 (the encoding of the RR pair 2)

BM48-31	TATGTTGTCTGTATTTCTATTAAGCATTGCCATTGGTTCTTTAGTAATTACCTCAAG	1015
PCoV_GX-P5E	TATGTCATTTGTATATCTACTAAGCATTGTTACTGGTTCTTTAGTAATTACCTTAGA	1131
PCoV_GX-P1E	TATGTCATTTGTATATCTACTAAGCATTGTTACTGGTTCTTTAGTAATTACCTTAGA	1131
PCoV_GX-P4L	TATGTCATTTGTATATCTACTAAGCATTGTTACTGGTTCTTTAGTAATTACCTTAGA	1131
PCoV_GX-P2V	TATGTCATTTGTATATCTACTAAGCATTGTTACTGGTTCTTTAGTAATTACCTTAGA	1131
PCoV_GX-P5L	TATGTCATTTGTATATCTACTAAGCATTGTTACTGGTTCTTTAGTAATTACCTTAGA	1131
ZC45	TATGTCATTTGCATTTCCACAAAGCATTGTTACTGGTTCTTTAGTAATTACCTAAAG	1200
ZXC21	TATGTCATTTGCATTTCCACAAAGCATTGTTACTGGTTCTTTAGTAATTACCTAAAG	1200
PCoV-MP789	TATGTCATTTGTATTTCCACAAAGCATTGTTACTGGTTCTTTAGCAACTACCTAAAG	1200
RacCS203	TATATCATTTGTATTTCCACAAAGCATTGTTACTGGTTCTTTAGCAATTACCTAAAG	1200
RaTG13	TATGTCATTTGTATTTCCACAAAGCATTGTTACTGGTTCTTTAGTAATTACCTAAA	1200
RShSTT200	TATATCATTTGTATTTCCACAAAGCATTGTTACTGGTTCTTTAGCAACTACCTAAA	1200
RShSTT182	TATATCATTTGTATTTCCACAAAGCATTGTTACTGGTTCTTTAGCAACTACCTAAA	1200
RmYN02	TATATCATTTGTATTTCCACAAAGCATTGTTACTGGTTCTTTAGCAACTACCTAAA	1200
INMI1	TATATCATTTGTATTTCCACAAAGCATTGTTACTGGTTCTTTAGTAATTACCTAAA	1200
CDC-CruiseA-12	TATATCATTTGTATTTCCACAAAGCATTGTTACTGGTTCTTTAGTAATTACCTAAA	1200
WIV04	TATATCATTTGTATTTCCACAAAGCATTGTTACTGGTTCTTTAGTAATTACCTAAA	1200
Wuhan-Hu-1	TATATCATTTGTATTTCCACAAAGCATTGTTACTGGTTCTTTAGTAATTACCTAAA	1200
RpYN06	TATATCATTTGTATTTCCACAAAGCATTGTTACTGGTTCTTTAGTAATTACCTAAA	1200

BM48-31	AGAGTTGCTTTAATGGTACTTCCTTTAGCACTTTTGAAGAAGCAGCTTTGTGTACATTC	1075
PCoV_GX-P5E	AGAGTTGCTTTAATGGTACTTCCTTTAGCACTTTTGAAGAAGCAGCTTTGTGTACATTC	1191
PCoV_GX-P1E	AGAGTTGCTTTAATGGTACTTCCTTTAGCACTTTTGAAGAAGCAGCTTTGTGTACATTC	1191
PCoV_GX-P4L	AGAGTTGCTTTAATGGTACTTCCTTTAGCACTTTTGAAGAAGCAGCTTTGTGTACATTC	1191
PCoV_GX-P2V	AGAGTTGCTTTAATGGTACTTCCTTTAGCACTTTTGAAGAAGCAGCTTTGTGTACATTC	1191
PCoV_GX-P5L	AGAGTTGCTTTAATGGTACTTCCTTTAGCACTTTTGAAGAAGCAGCTTTGTGTACATTC	1191
ZC45	CGTGTAGTCTTTAATGGTGTTCCTTTAGTACTTTTGAAGAAGCTGCCTTTATGCACCTTT	1260
ZXC21	CGTGTAGTCTTTAATGGTGTTCCTTTAGTACTTTTGAAGAAGCTGCCTTTATGCACCTTT	1260
PCoV-MP789	CGTGTAGTCTTTAATGGTGTTCCTTTAGTACTTTTGAAGAAGCTGCCTTTATGCACCTTT	1260
RacCS203	CGTGTAGTCTTTAATGGTGTTCCTTTAGTACTTTTGAAGAAGCTGCCTTTATGCACCTTT	1260
RaTG13	CGTGTAGTCTTTAATGGTGTTCCTTTAGTACTTTTGAAGAAGCTGCCTTTATGCACCTTT	1260
RShSTT200	CGTGTAGTCTTTAATGGTGTTCCTTTAGTACTTTTGAAGAAGCTGCCTTTATGCACCTTT	1260
RShSTT182	CGTGTAGTCTTTAATGGTGTTCCTTTAGTACTTTTGAAGAAGCTGCCTTTATGCACCTTT	1260
RmYN02	CGTGTAGTCTTTAATGGTGTTCCTTTAGTACTTTTGAAGAAGCTGCCTTTATGCACCTTT	1260
INMI1	CGTGTAGTCTTTAATGGTGTTCCTTTAGTACTTTTGAAGAAGCTGCCTTTATGCACCTTT	1260
CDC-CruiseA-12	CGTGTAGTCTTTAATGGTGTTCCTTTAGTACTTTTGAAGAAGCTGCCTTTATGCACCTTT	1260
WIV04	CGTGTAGTCTTTAATGGTGTTCCTTTAGTACTTTTGAAGAAGCTGCCTTTATGCACCTTT	1260
Wuhan-Hu-1	CGTGTAGTCTTTAATGGTGTTCCTTTAGTACTTTTGAAGAAGCTGCCTTTATGCACCTTT	1260
RpYN06	CGTGTAGTCTTTAATGGTGTTCCTTTAGTACTTTTGAAGAAGCTGCCTTTATGCACCTTT	1260

nsp6

Transmembrane domains. Double-membrane vesicles formation



Protein nsp6 (the RR pair)

BtKY72	ILLILMTARTVYDDARRVWTFMNVITLVYKVVYGNVLDQAIAMWALVISVTSNYSYSGVVT	173
ZC45	VLLILMTARTVYDDGARRVWTLMNVLTLVYKVVYGNALDQAISMWALIIISVTSNYSYSGVVT	180
ZXC21	VLLILMTARTVYDDSARRVWTLMNVLTLVYKVVYGNALDQAISMWALIIISVTSNYSYSGVVT	180
RmYN02	VLLILMTARTVYDDGARRVWTLMNVLTLVYKVVYGNALDQAISMWALIIISVTSNYSYSGVVT	180
RShSTT200	VLLILMTARTVYDDGARRVWTLMNVLTLVYKVVYGNALDQAISMWALIIISVTSNYSYSGVVT	180
RShSTT182	VLLILMTARTVYDDGARRVWTLMNVLTLVYKVVYGNALDQAISMWALIIISVTSNYSYSGVVT	180
PCoV-MP789	VLLILMTARTVYDDGARRVWTLMNVLTLVYKVVYGNALDQAISMWALIIISVTSNYSYSGVVT	180
RacCS203	VLLILMTARTVYDDGARRVWTLMNVLTLVYKVVYGNALDQAISMWALIIISVTSNYSYSGVVT	180
WIV04	VLLILMTARTVYDDGARRVWTLMNVLTLVYKVVYGNALDQAISMWALIIISVTSNYSYSGVVT	180
Wuhan-Hu-1	VLLILMTARTVYDDGARRVWTLMNVLTLVYKVVYGNALDQAISMWALIIISVTSNYSYSGVVT	180
INMI1	VLLILMTARTVYDDGARRVWTLMNVLTLVYKVVYGNALDQAISMWALIIISVTSNYSYSGVVT	180
CDC-CruiseA-12	VLLILMTARTVYDDGARRVWTLMNVLTLVYKVVYGNALDQAISMWALIIISVTSNYSYSGVVT	180
RpYN06	VLLILMTARTVYDDGARRVWTLMNVLTLVYKVVYGNALDQAISMWALIIISVTSNYSYSGVVT	180
RaTG13	VLLILMTARTVYDDGARRVWTLMNVLTLVYKVVYGNALDQAISMWALIIISVTSNYSYSGVVT	180
	:*****.:*****:*****:*****:*****:*****:*****:*****:*****:*****	

Gene nsp6 (the encoding RR pair)

BtKY72	ATACTGCTCATCCTCATGACAGCTAGAAGTGTCTATGACGATGCCACTAGACGAGTTTGG	399
ZC45	GTGTTATTAATCCTCATGACAGCAAGAACCGTATATGATGATGGTGCTAGAAGAGTTTGG	420
ZXC21	GTGTTATTAATCCTCATGACAGCAAGAAGTGTATATGATGATAGTGCTAGAAGAGTTTGG	420
PCoV-MP789	GTGTTATTAATCCTTATGACAGCAAGAAGTGTGTATGATGATGGTGCTAGAAGAGTTTGG	420
RaTG13	GTGTTACTAATCCTTATGACAGCAAGAAGTGTGTATGATGATGGTGCTAGGAGGGTGTGG	420
RpYN06	GTCTTACTAATCCTTATGACAGCAAGAAGTGTGTATGATGATGGTGCTAGGAGAGTGTGG	420
RmYN02	GTCTTACTAATCCTTATGACAGCAAGAAGTGTGTATGATGATGGTGCTAGGAGAGTGTGG	420
WIV04	GTGTTACTAATCCTTATGACAGCAAGAAGTGTGTATGATGATGGTGCTAGGAGAGTGTGG	420
Wuhan-Hu-1	GTGTTACTAATCCTTATGACAGCAAGAAGTGTGTATGATGATGGTGCTAGGAGAGTGTGG	420
INMI1	GTGTTACTAATCCTTATGACAGCAAGAAGTGTGTATGATGATGGTGCTAGGAGAGTGTGG	420
CDC-CruiseA-12	GTGTTACTAATCCTTATGACAGCAAGAAGTGTGTATGATGATGGTGCTAGGAGAGTGTGG	420
RacCS203	GTGTTACTAATCCTTATGACAGCAAGAAGTGTGTATGATGATGGTGCTAGGAGAGTGTGG	420
RShSTT200	GTGTTACTAATCCTTATGACAGCAAGAAGTGTGTATGATGATGGTGCTAGGAGAGTGTGG	420
RShSTT182	GTGTTACTAATCCTTATGACAGCAAGAAGTGTGTATGATGATGGTGCTAGGAGAGTGTGG	420
	* * * ***** ***** ***** ** ***** ** * ** * ** * ** *	

Protein nsp13 (the RR pair 2)

BtKY72 YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 BM48 YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 PCoV_GX-P2V YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 PCoV_GX-P1E YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 PCoV_GX-P4L YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 PCoV_GX-P5L YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 PCoV_GX-P5E YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 ZXC21 YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 ZC45 YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 RacCS203 YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 RmYN02 YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 RpYN06 YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 HKU YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 Tor2/FP1-10912 YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 MP789 YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 INMI1 YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 CDC-CruiseA-12 YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 WIV04 YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 Wuhan-Hu-1 YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 RShSTT200 YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 RShSTT182 YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 RaTG13 YFNSVCRLMKTIGPDMFLGTCRRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI 480
 *****RRCPAEIVDTVSTLVYDNKLRHKDSSQCFKMFYKQVI *****

Gene nsp13 (the encoding of the RR pair 2)

PCoV_GX-P2V TGTAGAAGATGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAA 1380
 PCoV_GX-P5E TGTAGAAGATGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAA 1380
 PCoV_GX-P4L TGTAGAAGATGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAA 1380
 PCoV_GX-P5L TGTAGAAGATGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAA 1380
 PCoV_GX-P1E TGTAGAAGATGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAA 1380
 MP789 TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 RpYN06 TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 RmYN02 TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 RacCS203 TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 RShSTT200 TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 RShSTT182 TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 INMI1 TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 CDC-CruiseA-12 TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 WIV04 TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 Wuhan-Hu-1 TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 RaTG13 TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 HKU TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 Tor2/FP1-10912 TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 ZXC21 TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 ZC45 TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 BM48 TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 BtKY72 TGTGTCGGCTGTCCTGCTGAAATAGTTGACACTGTAAGTGCTCTAGTTTATGATAATAAG 1380
 ** * ***** * **

Protein nsp13 (the RR pair 3)

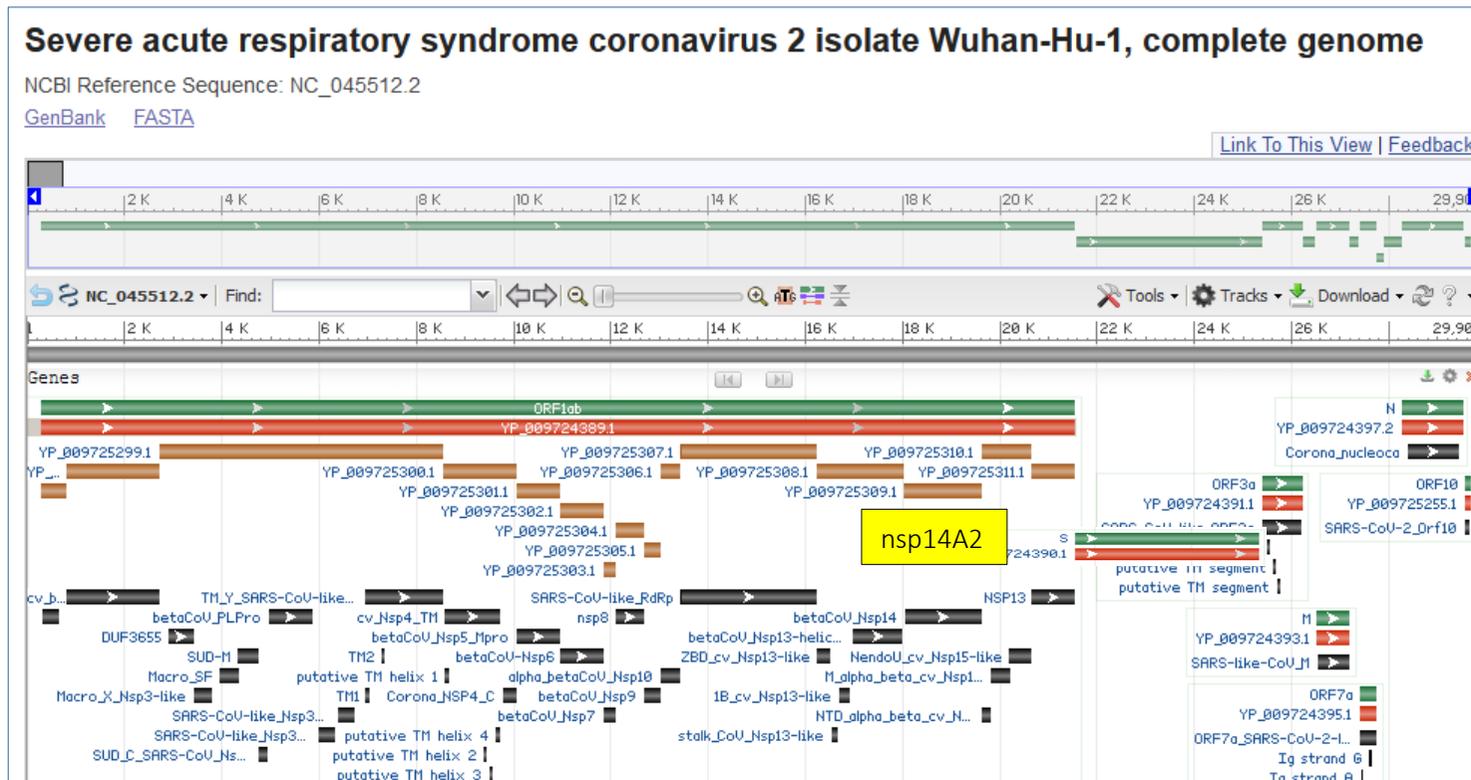
BtKY72	YDYVIFTQTETAHSCNVNRFNVAITRAKVGILCIMSDDRDLQFASLEVPRKSVAA--	598
BM48	YDYVIFAQTETAHSCNVNRFNVAITRAKVGILCIMSDDKLDYDKLQFTSLEVPRRSVA--	598
PCoV_GX-P2V	YDYVIFTQTETAHSCNVNRFNVAITRAKIGILCIMSDDRDLQFASLEVPRRVATL	600
PCoV_GX-P1E	YDYVIFTQTETAHSCNVNRFNVAITRAKIGILCIMSDDRDLQFASLEVPRRVATL	600
PCoV_GX-P4L	YDYVIFTQTETAHSCNVNRFNVAITRAKIGILCIMSDDRDLQFASLEVPRRVATL	600
PCoV_GX-P5L	YDYVIFTQTETAHSCNVNRFNVAITRAKIGILCIMSDDRDLQFASLEVPRRVATL	600
PCoV_GX-P5E	YDYVIFTQTETAHSCNVNRFNVAITRAKIGILCIMSDDRDLQFASLEVPRRVATL	600
ZXC21	YDYVIFTQTETAHSCNVNRFNVAITRAKIGILCIMSDDRDLQFASLEVPRRVATL	600
ZC45	YDYVIFTQTETAHSCNVNRFNVAITRAKIGILCIMSDDRDLQFMSLEVPRRVATL	600
RacCS203	YDYVIFTQTETSHSCNVNRFNVAITRAKVGILCIMSDDRDLQFASLEVPRRVATL	600
RmYN02	YDYVIFTQTETAHSCNVNRFNVAITRAKVGILCIMSDDRDLQFASLEVPRRVATL	600
RpYN06	YDYVIFTQTETAHSCNVNRFNVAITRAKVGILCIMSDDRDLQFASLEVPRRVATL	600
HKU	YDYVIFTQTETAHSCNVNRFNVAITRAKIGILCIMSDDRDLQFASLEVPRRVATL	600
Tor2/FP1-10912	YDYVIFTQTETAHSCNVNRFNVAITRAKIGILCIMSDDRDLQFASLEVPRRVATL	600
MP789	YDYVIFTQTETAHSCNVNRFNVAITRAKVGILCIMSDDRDLQFASLEVPRRVATL	600
INMI1	YDYVIFTQTETAHSCNVNRFNVAITRAKVGILCIMSDDRDLQFASLEVPRRVATL	600
CDC-CruiseA-12	YDYVIFTQTETAHSCNVNRFNVAITRAKVGILCIMSDDRDLQFASLEVPRRVATL	600
WIV04	YDYVIFTQTETAHSCNVNRFNVAITRAKVGILCIMSDDRDLQFASLEVPRRVATL	600
Wuhan-Hu-1	YDYVIFTQTETAHSCNVNRFNVAITRAKVGILCIMSDDRDLQFASLEVPRRVATL	600
RShSTT200	YDYVIFTQTETAHSCNVNRFNVAITRAKVGILCIMSDDRDLQFASLEVPRRVATL	600
RShSTT182	YDYVIFTQTETAHSCNVNRFNVAITRAKVGILCIMSDDRDLQFASLEVPRRVATL	600
RaTG13	YDYVIFTQTETAHSCNVNRFNVAITRAKVGILCIMSDDRDLQFASLEVPRRVATL	600
	*****.*****.*****.*****.*****.*****.***** ** :*. *	

Gene nsp13 (the encoding of the RR pair 3)

PCoV_GX-P2V	CTTTATGACAAATTACAATTTACAAGCCTTGAAGTTCCA	CGT	CGAAACGTGGCAACCTTA	1800
PCoV_GX-P5E	CTTTATGACAAATTACAATTTACAAGCCTTGAAGTTCCA	CGT	CGAAACGTGGCAACCTTA	1800
PCoV_GX-P4L	CTTTATGACAAATTACAATTTACAAGCCTTGAAGTTCCA	CGT	CGAAACGTGGCAACCTTA	1800
PCoV_GX-P5L	CTTTATGACAAATTACAATTTACAAGCCTTGAAGTTCCA	CGT	CGAAACGTGGCAACCTTA	1800
PCoV_GX-P1E	CTTTATGACAAATTACAATTTACAAGCCTTGAAGTTCCA	CGT	CGAAACGTGGCAACCTTA	1800
MP789	CTTTATGACAAGTTGCAATTTACAAGTCTTGAAATCCA	CGT	AGAAATGTGGCAACTTTA	1800
RpYN06	CTTTATGACAAGTTGCAATTTACAAGTCTTGAAATCCA	CGT	AGAAATGTGGCAACTTTA	1800
RmYN02	CTTTATGACAAGTTGCAATTTACAAGTCTTGAAATCCA	CGT	AGAAATGTGGCAACTTTA	1800
RacCS203	CTTTATGACAAGTTGCAATTTACAAGTCTTGAAATCCA	CGT	AGAAATGTGGCAACTTTA	1800
RShSTT200	CTTTATGACAAGTTGCAATTTACAAGTCTTGAAATCCA	CGT	AGAAATGTGGCAACTTTA	1800
RShSTT182	CTTTATGACAAGTTGCAATTTACAAGTCTTGAAATCCA	CGT	AGAAATGTGGCAACTTTA	1800
INMI1	CTTTATGACAAGTTGCAATTTACAAGTCTTGAAATCCA	CGT	AGAAATGTGGCAACTTTA	1800
CDC-CruiseA-12	CTTTATGACAAGTTGCAATTTACAAGTCTTGAAATCCA	CGT	AGAAATGTGGCAACTTTA	1800
WIV04	CTTTATGACAAGTTGCAATTTACAAGTCTTGAAATCCA	CGT	AGAAATGTGGCAACTTTA	1800
Wuhan-Hu-1	CTTTATGACAAGTTGCAATTTACAAGTCTTGAAATCCA	CGT	AGAAATGTGGCAACTTTA	1800
RaTG13	CTTTATGACAAGTTGCAATTTACAAGTCTTGAAATCCA	CGT	AGAAATGTGGCAACTTTA	1800
HKU	CTTTATGACAAACTGCAATTTACAAGTCTAGAAATACCA	CGT	CGCAATGTGGCTACATTA	1800
Tor2/FP1-10912	CTTTATGACAAACTGCAATTTACAAGTCTAGAAATACCA	CGT	CGCAATGTGGCTACATTA	1800
ZXC21	CTTTATGACAAGCTGCAATTTACGAGTCTAGAAGTACCG	CGT	CGTAATGTGGCTACTTTA	1800
ZC45	CTTTATGACAAGCTTCAATTTATGAGTCTAGAAGTACCG	CGT	CGAAATGTGGCTACTTTA	1800
BM48	CTCTATGATAAATTACAATTTACTAGTCTGGAAGTCCA	CGT	AGAGTGTGGC-----	1793
BtKY72	CTTTATGACAAACTCCAATTTGCTAGTCTAGAAGTCCA	CGT	AAAAGTGTGGC-----	1793
	** ***** ** * ***** ** ** * * **			

nsp14A2

Guanosine N7- methyltransferase. 3 to 5 exoribonuclease, proofreading



Protein:nsp14A2 (the RR pair 1)

PCoV_GX-P2V	AENVTLGFKDCSKVITGLHPTQAPTYLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
PCoV_GX-P5E	AENVTLGFKDCSKVITGLHPTQAPTYLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
PCoV_GX-P1E	AENVTLGFKDCSKVITGLHPTQAPTYLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
PCoV_GX-P4L	AENVTLGFKDCSKVITGLHPTQAPTYLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
PCoV_GX-P5L	AENVTLGFKDCSKVITGLHPTQAPTYLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
PCoV-MP789	AENVTLGFKDCSKVINGLHPTQALHTLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
RShSTT200	AENVTLGFKDCSKVITGLHPTQAPTHLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
RShSTT182	AENVTLGFKDCSKVITGLHPTQAPTHLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
RacCS203	AENVTLGFKDCSKVITGLHPTQAPTHLSVDTKFKTEGLCVDIPGIPKDLTYRRLISMMGF	60
RaTG13	AENVTLGFKDCSKVITGLHPTQAPTHLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
RmYN02	AENVTLGFKDCSKVITGLHPTQAPTHLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
INMI1	AENVTLGFKDCSKVITGLHPTQAPTHLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
CDC-CruiseA-12	AENVTLGFKDCSKVITGLHPTQAPTHLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
WIV04	AENVTLGFKDCSKVITGLHPTQAPTHLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
Wuhan-Hu-1	AENVTLGFKDCSKVITGLHPTQAPTHLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
RpYN06	AENVTLGFKDCSKVITGLHPTQAPTHLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
BtKY72	AENVTLGFKDCSKVINGLHPTQSPTYLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
BM48-31	-ENVTLGFKDCSKLITGLHPTQAPTYLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	59
Tor2_FP1-10912	AENVTLGFKDCSKIITGLHPTQAPTHLSVDIKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
HKU-39849	AENVTLGFKDCSKIITGLHPTQAPTHLSVDIKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
ZXC21	AENVTLGFKDCSKIITGLHPTQAPTHLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60
ZC45	AENVTLGFKDCSKIITGLHPTQAPTHLSVDTKFKTEGLCVDIPGIPKDMTYRRLISMMGF	60

*****.*.*****.*.:.*** *****.*****.*****.*****

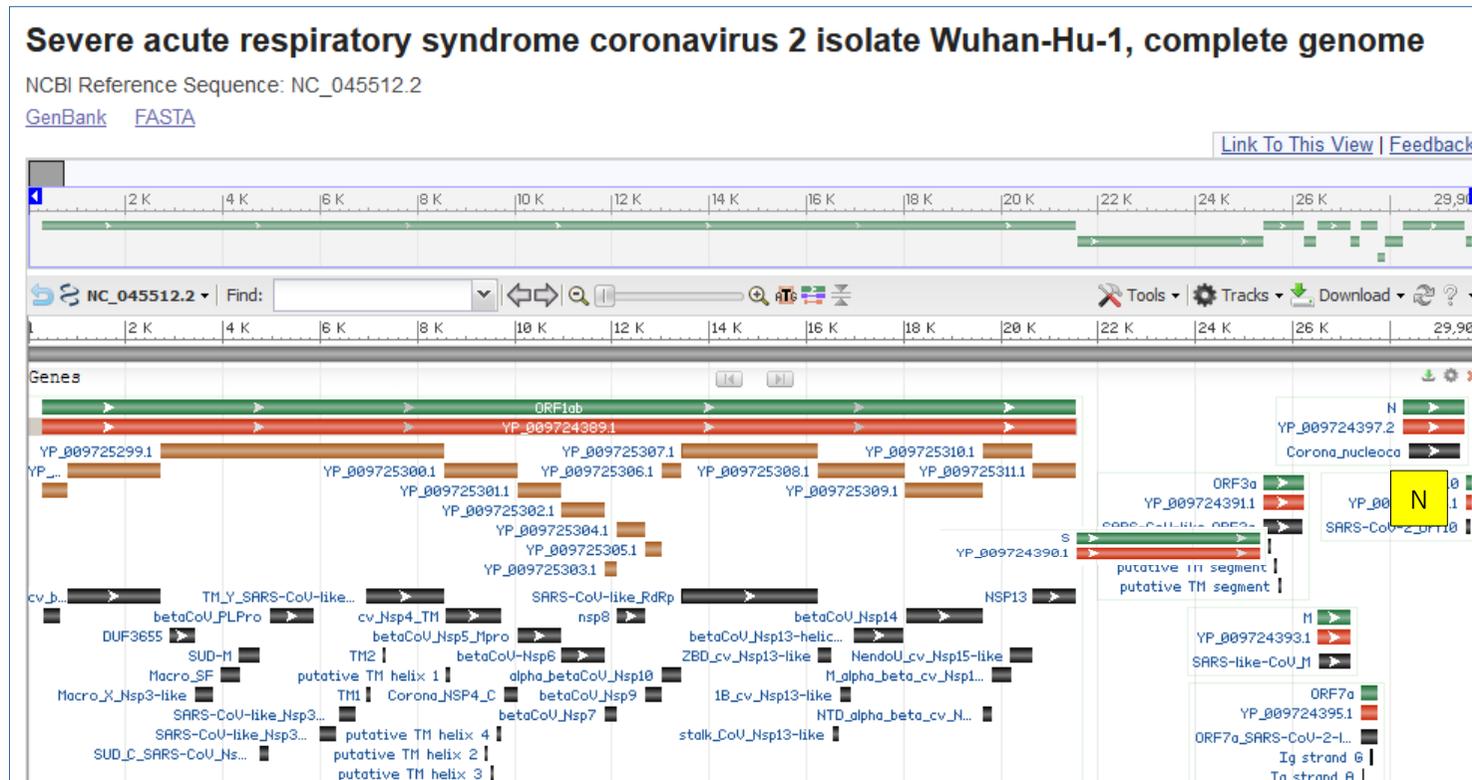
Gene nsp14A2 (the encoding of the RR pair 1)

PCoV_GX-P2V	GACATACCAGGAATACCAAAAGACATGACCTATAGGAGACTCATCTCTATGATGGGTTTC	180
PCoV_GX-P5E	GACATACCAGGAATACCAAAAGACATGACCTATAGGAGACTCATCTCTATGATGGGTTTC	180
PCoV_GX-P1E	GACATACCAGGAATACCAAAAGACATGACCTATAGGAGACTCATCTCTATGATGGGTTTC	180
PCoV_GX-P4L	GACATACCAGGAATACCAAAAGACATGACCTATAGGAGACTCATCTCTATGATGGGTTTC	180
PCoV_GX-P5L	GACATACCAGGAATACCAAAAGACATGACCTATAGGAGACTCATCTCTATGATGGGTTTC	180
PCoV-MP789	GACATACCAGGTATACCCAAGGACATGACCTATAGGAGACTCATTCCATGATGGGTTTC	180
RShSTT200	GACATACCTGGCATACTAAGGACATGACTTATAGAAGACTCATCTCTATGATGGGTTTC	180
RShSTT182	GACATACCTGGCATACTAAGGACATGACTTATAGAAGACTCATCTCTATGATGGGTTTC	180
RacCS203	GACATACCTGGCATACTAAGGACTTGACCTATAGAAGACTCATCTCTATGATGGGTTTC	180
RpYN06	GACATACCTGGCATACTAAGGACATGACCTATAGAAGACTCATCTCTATGATGGGTTTT	180
RaTG13	GACATACCTGGTATACCTAAGGACATGACCTATAGAAGACTCATCTCTATGATGGGTTTC	180
INMI1	GACATACCTGGCATACTAAGGACATGACCTATAGAAGACTCATCTCTATGATGGGTTTT	180
CDC-CruiseA-12	GACATACCTGGCATACTAAGGACATGACCTATAGAAGACTCATCTCTATGATGGGTTTT	180
WIV04	GACATACCTGGCATACTAAGGACATGACCTATAGAAGACTCATCTCTATGATGGGTTTT	180
Wuhan-Hu-1	GACATACCTGGCATACTAAGGACATGACCTATAGAAGACTCATCTCTATGATGGGTTTT	180
RmYN02	GACATACCTGGCATACTAAGGACATGACCTATAGAAGACTCATCTCCATGATGGGTTTC	180
BtKY72	GACATACCTAGCATACCTAAGGACATGACTTATCGTAGACTCATCTCTATGATGGGCTTC	180
BM48-31	GACATACCAGGAATACCAAAAGGACATGACCTATCGTAGGCTCATCTCTATGATGGGTTTT	177
Tor2/FP1-10912	GACATACCAGGCATACCAAAAGGACATGACCTACCGTAGACTCATCTCTATGATGGGTTTC	180
HKU-39849	GACATACCAGGCATACCAAAAGGACATGACCTACCGTAGACTCATCTCTATGATGGGTTTC	180
ZXC21	GACATACCAGGAATACCAAAAGACATGACCTATCGTAGACTCATCTCTATGATGGGTTTT	180
ZC45	GACATACCAGGAATACCAAAAGGACATGACCTATCGTAGACTCATCTCTATGATGGGCTTC	180

***** * ***** ** ** * ***** * * * ***** * * ***** **

Orf9 N. Nucleocapsid

Packages the positive strand viral genome RNA. Virion assembly



Protein N orf9 Nucleoprotein (the RR pair 1)

HKU-39849	MSDNGPQSNQRSAPRITFFGGPTDSTDNNQNGGRNGARPKQRRPQGLPNNTASWFTALTQH	60
Tor2/FP1-10912	MSDNGPQSNQRSAPRITFFGGPTDSTDNNQNGGRNGARPKQRRPQGLPNNTASWFTALTQH	60
PcoV_GX-P5E	MSDNGPQ-N--RAPRITFFGGPSDSTDNNQNGDRSGARPKQRRPQGLPNNTASWFTALTQH	57
PcoV_GX-P2V	MSDNGPQ-X--RAPRITFFGGPSDSTDNNQNGDRSGARPKQRRPQGLPNNTASWFTALTQH	57
PcoV_GX-P4L	MSDNGPQ-N--RAPRITFFGGPSDSTDNNQNGDRSGARPKQRRPQGLPNNTASWFTALTQH	57
PCoV_GX-P5L	MSDNGPQ-N--RAPRITFFGGPSDSTDNNQNGDRSGARPKQRRPQGLPNNTASWFTALTQH	57
PcoV_GX-P1E	MSDNGPQ-N--RAPRITFFGGPSDSTDNNQNGDRSGARPKQRRPQGLPNNTASWFTALTQH	57
ZC45	MSDNGPQ-NQRSAPRITFFGGPSDSSDNNQNGERNGARPKQRRPQGLPNNTASWFTALTQH	59
ZXC21	MSDNGPQ-NQSSAPRITFFGGPSDSSDNNQNGERNGARPKQRRPQGLPNNTASWFTALTQH	59
RShSTT182	MSDNGPQ-NQRNAPRITFFGGPSDSSDNNQNGERSGARPKQRRPQGLPNNTASWFTALTQH	59
RShSTT200	MSDNGPQ-NQRNAPRITFFGGPSDSSDNNQNGERSGARPKQRRPQGLPNNTASWFTALTQH	59
RacCS203	MSDNGPQ-NQRNAPRITFFGGPSDSSDNNQNGERSGARPKQRRPQGLPNNTASWFTALTQH	59
PanCoV-MP789	MSDNGPQ-NQRNAPRITFFGGPSDSSDNNQNGERSGARPKQRRPQGLPNNTASWFTALTQH	59
RmYN02	MSDNGHQ-SQRNAPRITFFGGPSDSTGNSQNGERSGARPKQRRPQGLPNNTASWFTALTQH	59
RaTG13	MSDNGPQ-NQRNAPRITFFGGPSDSTGNSQNGERSGARPKQRRPQGLPNNTASWFTALTQH	59
WIV04	MSDNGPQ-NQRNAPRITFFGGPSDSTGNSQNGERSGARSKQRRPQGLPNNTASWFTALTQH	59
CDC-CruiseA-12	MSDNGPQ-NQRNAPRITFFGGPSDSTGNSQNGERSGARSKQRRPQGLPNNTASWFTALTQH	59
WIV04	MSDNGPQ-NQRNAPRITFFGGPSDSTGNSQNGERSGARSKQRRPQGLPNNTASWFTALTQH	59
INMI1	MSDNGPQ-NQRNAPRITFFGGPSDSTGNSQNGERSGARSKQRRPQGLPNNTASWFTALTQH	59
RpYN06	MSDNGPQ-SQRNAPRITFFGGPSDSTGNSQNGERSGARPKQRRPQGLPNNTASWFTALTQH	59
BM48-131	MTDNGQ-SNSRNAPRITFFGV-SDTSDNNQNAERAGARPKQRRPQGLPNNTASWFTALTQH	58
BtKY72	MTDNGQ-QGPRNAPRITFFGV-SDNFDNNQNGDRTGARPKHRRPQGLPNNTASWFTALTQH	58

*.*** ***** :*. . . .:*. * ** * :***** * .***** **

Gene N orf9 Nucleoprotein (the encoding of the RR pair 1)

HKU-39849	CGCCGACCCCAAGGTTTACCCAATAAATACTGCGTCTTGGTTCACAGCTCTCACTCAGCAT	180
Tor2/FP1-10912	CGCCGACCCCAAGGTTTACCCAATAAATACTGCGTCTTGGTTCACAGCTCTCACTCAGCAT	180
PcoV_GX-P5E	CGAAGGCCCCAGGGATTACCCAATAAATACTGCGTCTTGGTTCACCGCTCTCACTCAACAT	171
PcoV_GX-P2V	CGAAGGCCCCAGGGATTACCCAATAAATACTGCGTCTTGGTTCACCGCTCTCACTCAACAT	171
PcoV_GX-P4L	CGAAGGCCCCAGGGATTACCCAATAAATACTGCGTCTTGGTTCACCGCTCTCACTCAACAT	171
PcoV_GX-P5L	CGAAGGCCCCAGGGATTACCCAATAAATACTGCGTCTTGGTTCACCGCTCTCACTCAACAT	171
PcoV_GX-P1E	CGAAGGCCCCAGGGATTACCCAATAAATACTGCGTCTTGGTTCACCGCTCTCACTCAACAT	171
ZC45	CGTCGACCCCAAGGCTTACCCAATAAATACTGCATCTTGGTTCACCGCTCTCACTCAACAT	177
ZXC21	CGTCGACCCCAAGGCTTACCCAATAAATACTGCATCTTGGTTCACCGCTCTCACTCAACAT	177
PanCoV-MP789	CGTCGTCCCCAAGGTTTACCCAATAAATACTGCGTCTTGGTTCACCGCTCTCACTCAACAT	177
Wuhan-Hu-1	CGTCGGCCCCAAGGTTTACCCAATAAATACTGCGTCTTGGTTCACCGCTCTCACTCAACAT	177
CDC-CruiseA-12	CGTCGGCCCCAAGGTTTACCCAATAAATACTGCGTCTTGGTTCACCGCTCTCACTCAACAT	177
WIV04	CGTCGGCCCCAAGGTTTACCCAATAAATACTGCGTCTTGGTTCACCGCTCTCACTCAACAT	177
INMI1	CGTCGGCCCCAAGGTTTACCCAATAAATACTGCGTCTTGGTTCACCGCTCTCACTCAACAT	177
Bat-CoV-RatG13	CGTCGGCCCTCAAGGTTTACCCAATAAATACTGCGTCTTGGTTCACCGCTCTCACTCAACAT	177
RmYN02	CGTCGTCCCCAAGGTTTACCCAATAAATACTGCGTCTTGGTTCACCGCTCTCACTCAACAT	177
RpYN06	CGTCGTCCCCAAGGTTTACCCAATAAATACTGCGTCTTGGTTCACCGCTCTCACTCAACAT	177
RShSTT200	CGTCGTCCACAAGGCCTACCCAATAAATACTGCGTCTTGGTTCACCGCTCTCACTCAGCAT	177
RShSTT182	CGTCGTCCACAAGGCCTACCCAATAAATACTGCGTCTTGGTTCACCGCTCTCACTCAGCAT	177
RacCS203	CGTCGTCCCCAAGGTTTACCCAATAAATACTGCGTCTTGGTTCACCGCTCTCACTCAGCAT	177
BM48-131	AGAAGACCGCAAGGCCCTCTAACAAACACAGCATCCTGGTTCACAGCTCTCACTCAGCAT	174
BtKY72	AGAAGACCGCAAGGCCCTCTAACAAACACGGCATCCTGGTTCACCGCTCTCACTCAACAT	174

* * * * * * * * * * *

Protein N orf9 Nucleoprotein (the RR pair 2 and 3)

HKU-39849	GKEELRFPRGQGVPIINTNSGDDQIGYYRRATRRVRGGDGKMKELSPRWYFYLLGTGPEA	120
Tor2/FP1-10912	GKEELRFPRGQGVPIINTNSGDDQIGYYRRATRRVRGGDGKMKELSPRWYFYLLGTGPEA	120
PcoV_GX-P5E	GKEDLRFPRGQGVPIINTNSTKDDQIGYYRRATRRVRGGDGKMKDLSPRWYFYLLGTGPEA	117
PcoV_GX-P2V	GKEDLRFPRGQGVPIINTNSTKDDQIGYYRRATRRVRGGDGKMKDLSPRWYFYLLGTGPEA	117
PcoV_GX-P4L	GKEDLRFPRGQGVPIINTNSTKDDQIGYYRRATRRVRGGDGKMKDLSPRWYFYLLGTGPEA	117
PCoV_GX-P5L	GKEDLRFPRGQGVPIINTNSTKDDQIGYYRRATRRVRGGDGKMKDLSPRWYFYLLGTGPEA	117
PcoV_GX-P1E	GKEDLRFPRGQGVPIINTNSTKDDQIGYYRRATRRVRGGDGKMKDLSPRWYFYLLGTGPEA	117
ZC45	GKENLTFPRGQGVPIINTNSKDDQIGYYRRATRRIRGGDGKMKELSPRWYFYLLGTGPEA	119
ZXC21	GKENLTFPRGQGVPIINTNSKDDQIGYYRRATRRIRGGDGKMKELSPRWYFYLLGTGPEA	119
RShSTT182	GKENLTFPRGQGVPIINTNSTKDDQIGYYRRATRRIRGGDGKMKDLSPRWYFYLLGTGPEA	119
RShSTT200	GKENLTFPRGQGVPIINTNSTKDDQIGYYRRATRRIRGGDGKMKDLSPRWYFYLLGTGPEA	119
RacCS203	GKENLTFPRGQGVPIINTNSTKDDQIGYYRRATRRIRGGDGKMKDLSPRWYFYLLGTGPEA	119
PanCoV-MP789	GKEDLRFPRGQGVPIINTNSPDDQIGYYRRATRRIRGGDGKMKDLSPRWYFYLLGTGPEA	119
RmYN02	GKEDLKFPRGQGVPIINTNSRDDQIGYYRRATRRIRGGDGKMKDLSPRWYFYLLGTGPEA	119
RaTG13	GKEDLKFPRGQGVPIINTNSPDDQIGYYRRATRRIRGGDGKMKDLSPRWYFYLLGTGPEA	119
WIV04	GKEDLKFPRGQGVPIINTNSPDDQIGYYRRATRRIRGGDGKMKDLSPRWYFYLLGTGPEA	119
CDC-CruiseA-12	GKEDLKFPRGQGVPIINTNSPDDQIGYYRRATRRIRGGDGKMKDLSPRWYFYLLGTGPEA	119
WIV04	GKEDLKFPRGQGVPIINTNSPDDQIGYYRRATRRIRGGDGKMKDLSPRWYFYLLGTGPEA	119
INMI1	GKEDLKFPRGQGVPIINTNSPDDQIGYYRRATRRIRGGDGKMKDLSPRWYFYLLGTGPEA	119
RpYN06	GKEDLKFPRGQGVPIINTNSPDDQIGYYRRATRRIRGGDGKMKDLSPRWYFYLLGTGPEA	119
BM48-31	GKEGLSFPRGQGVPIINTNSRDDQIGYYRRATRRVRGGDGKMKELSPRWYFYLLGTGPEA	118
BtKY72	GKETLTFPRGQGVPIINTNSGKDDQIGYYRRASRRVRGGDGKMKELSPRWYFYLLGTGPEA	118
	*** * ***** . ***** ***** . ** . ***** . ***** *****	

Gene N orf9 Nucleoprotein (the encoding of the RR pair 2 and 3)

HKU-39849	CCAGATGACCAAATTTGGCTACTACCGAAGAAGCTACCCGACGAGTTTCGTGGTGGTGACGGC	300
Tor2/FP1-10912	CCAGATGACCAAATTTGGCTACTACCGAAGAAGCTACCCGACGAGTTTCGTGGTGGTGACGGC	300
PcoV_GX-P5E	AAAGATGACCAAATTTGGCTACTACCGAAGAAGCTACCCAGACGAGTTTCGTGGTGGTGACGGT	291
PcoV_GX-P2V	AAAGATGACCAAATTTGGCTACTACCGAAGAAGCTACCCAGACGAGTTTCGTGGTGGTGACGGT	291
PcoV_GX-P4L	AAAGATGACCAAATTTGGCTACTACCGAAGAAGCTACCCAGACGAGTTTCGTGGTGGTGACGGT	291
PcoV_GX-P5L	AAAGATGACCAAATTTGGCTACTACCGAAGAAGCTACCCAGACGAGTTTCGTGGTGGTGACGGT	291
PcoV_GX-P1E	AAAGATGACCAAATTTGGCTACTACCGAAGAAGCTACCCAGACGAGTTTCGTGGTGGTGACGGT	291
ZC45	AAAGATGACCAAATTTGGCTACTACCGTAGAGCTACCCAGACGAATTCGTGGTGGTGACGGT	297
ZXC21	AAAGATGACCAAATTTGGCTACTACCGTAGAGCTACCCAGACGAATTCGTGGCGGTGACGGT	297
PanCoV-MP789	CCAGATGACCAAATTTGGCTACTACCGAAGAAGCTACCCAGACGAATTCGTGGTGGTGACGGT	297
Wuhan-Hu-1	CCAGATGACCAAATTTGGCTACTACCGAAGAAGCTACCCAGACGAATTCGTGGTGGTGACGGT	297
CDC-CruiseA-12	CCAGATGACCAAATTTGGCTACTACCGAAGAAGCTACCCAGACGAATTCGTGGTGGTGACGGT	297
WIV04	CCAGATGACCAAATTTGGCTACTACCGAAGAAGCTACCCAGACGAATTCGTGGTGGTGACGGT	297
INMI1	CCAGATGACCAAATTTGGCTACTACCGAAGAAGCTACCCAGACGAATTCGTGGTGGTGACGGT	297
Bat-CoV-RatG13	CCAGATGACCAAATTTGGCTACTACCGAAGAAGCTACCCAGACGAATTCGTGGTGGTGACGGT	297
RmYN02	CGAGATGACCAAATTTGGCTACTACCGAAGAAGCTACCCAGACGAATTCGTGGTGGTGACGGT	297
RpYN06	CCAGATGATCAAATTTGGCTACTACCGAAGAAGCTACCCAGACGAATTCGTGGTGGTGACGGT	297
RShSTT200	AAAGATGACCAAGATTGGCTACTACCGAAGAAGCTACCCAGACGAATTCGTGGTGGTGACGGT	297
RShSTT182	AAAGATGACCAAGATTGGCTACTACCGAAGAAGCTACCCAGACGAATTCGTGGTGGTGACGGT	297
RacCS203	AAAGATGACCAAGATTGGCTACTACCGAAGAAGCTACCCAGACGAATTCGTGGTGGTGACGGT	297
BM48-31	AGGGACGACCAAATTTGGCTACTATCGCAGAGCTACCCGACGAGTTTCGTGGTGGTGATGGT	294
BtKY72	AAAGATGACCAAATTTGGCTACTATAGAAGAAGCTTCCCGACGAGTTTCGTGGTGGTGACGGA	294
	*** ** * ***** ***** * ***** * ***** ***** **	

