

# <u>E. Ulzurrun<sup>1,2</sup>, D. del Hoyo<sup>2</sup>, M. Álvarez-Herrera<sup>3</sup>, P. Ruiz-Rodríguez<sup>3</sup>, B. Navarro-Domínguez<sup>3</sup>, S. Rodríguez Santana<sup>4</sup>,</u> D. Garcia Rasines<sup>4</sup>, R. Naveiro<sup>4</sup>, C. Gil<sup>1</sup>, J.M. Carazo<sup>2</sup>, M. Coscollá<sup>3</sup>, C.O.S. Sorzano<sup>2</sup>, N. Campillo<sup>1,4</sup>

1. Center for Biological Research Margarita Salas, CSIC. 2. National Biotechnology Center, CSIC. 3. Institute for Integrative Systems Biology, CSIC- University of Valencia. 4. Institute of Mathematical Sciences, CSIC

Background: Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2), initially reported in Wuhan (China) has spread worldwide. Like other viruses, SARS-CoV-2 accumulates mutations with each cycle of replication by continuously evolving a viral strain with one or more single nucleotide variants (SNVs). However, SNV coding for mutant residues that cause severe COVID-19 or lead to immune escape or vaccine failure are not well understood. We aim to identify haplotypes throughout the virus **proteome** associated with clinical phenotypes particularly with vaccine failure.

Methods: 386 and 4983 aligned SARS-Cov-2 consensus genome sequences and associated metadata from Gregorio Marañon Hospital (GMH386) and **GISAID** were used, respectively.

Pre-processing and analysis of haplotypes were carried out using scripts in R and Python programming languages, as well as python modules such as scipy.stats and others included in Biopython. Haplotype association was preformed by Chi-square test.

NCBI Reference Sequence NC\_045512.2 (Wuhan) was used to retrieve the coding sequences. Mutation is defined as mutational frequency for the minor nucleotide  $\geq 10\%$  for the major nucleotide at a particular gene locus. On the other hand, reference genome is made of the major nucleotide for each locus of the aligned sequences.

Results: A short description of whole population and haplotypes are shown in Tables 1 and 5. Using the nucleotide mutations, 85 and 45 non-synonymous amino acid mutations throughout the proteome of SARS-CoV-2 were identified for GMH386 and GISAID datasets, respectively. Associations with vaccine failure regardless of dose received were found in both set of sequences for two haplotypes. However, COVID-19 patients expected for haplotype 2 for GISAID was lower than expected significantly (p)  $= 5.88 \times 10^{-161}$  (Tables 2 and 6). Haplotypes of mutations are outlined in the Tables 3-4 and Tables 7-8.

Figure 1 shows colored in red the shared mutations for both datasets located at one homotrimer of the Spike protein. The receptor-binding domain (RBD) is colored in yellow. In contrast, the N-terminal domain (NTD) is highlighted in cyan. These loci are noteworthy with asterisk in Tables 4 and 8.

 Table 1. Brief description of the GM386 patients with a collection date from 2021-03-02

to 202204-24

Table 2. Vaccination distribution for the haplotypes with a frequency  $\geq$  10 %

	Female, n (%)	Male, n (%)	>= 65 y, n (%)	< 65 y, n (%)
GMH386 database, n = 386	199 (52)	187 (48)	250 (65)	136 (35)
h_GM386_1, n = 46	22 (48)	24 (52)	41 (89)	5 (11)
h_GM386_2, n = 39	24 (62)	15 (38)	35 (90)	4 (10)

Abbreviations: y = years; h\_GM386 = haplotype for Gregorio Marañon Hospital 386 dataset

	Vaccinated Yes	Expected frequencies Yes	Vaccinated No	Expected frequencies No	p value*
h_GM386_1, n = 46	37	26	9	20	1.41E-03
h_GM386_2, n = 39	31	22	8	17	5.83E-03

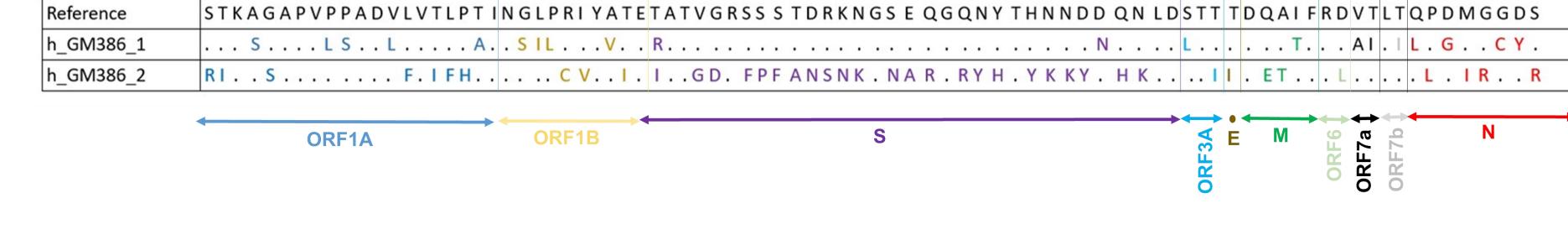
Abbreviations: h\_GM386 = haplotype for Gregorio Marañon Hospital 386 dataset \*Chi-square test

# Table 3. Haplotypes with 85 non-synonymous mutations for GM386

	and the second		
faranca		DT INCLODIVATETATVCDCC	ON LDCTT TDOAL CDDVTLTODDMCC

### Table 4. Protein mutations of SARS-Cov-2 for GM386

ORF1a:135	ORF1a:2710	ORF1b:591	<b>S:19</b>	S:405	S:501 <sup>*</sup>	S:981 <sup>*</sup>	M:112	N:203
ORF1a:842	ORF1a:2907	ORF1b:662	S:67 <sup>*</sup>	<b>S:408</b>	S:505 <sup>*</sup>	S:1146	<b>ORF6:20</b>	N:204
			*					



ORF1a:856	ORF1a:2930	ORF1b:829	S:95	S:417	S:547	<b>ORF3a:26</b>	ORF6:61	N:215
ORF1a:1306	ORF1a:3027	ORF1b:1000	S:213	<b>S:440</b>	S:655 <sup>*</sup>	ORF3a:64	ORF7a:82	N:377
ORF1a:1307	ORF1a:3053	ORF1b:1315	S:339 <sup>*</sup>	<b>S:446</b>	S:679	ORF3a:223	ORF7a:120	N:413
ORF1a:1352	ORF1a:3090	ORF1b:1566	S:346	S:477	S:764	<b>E:9</b>	ORF7b:18	
ORF1a:1803	ORF1a:3201	ORF1b:1759	S:371 <sup>*</sup>	S:484 <sup>*</sup>	S:796 <sup>*</sup>	M:3	ORF7b:40	
ORF1a:1887	ORF1a:3395	ORF1b:1918	S:373 <sup>*</sup>	S:493 <sup>*</sup>	S:950	M:19	N:9	
ORF1a:2046	ORF1a:3646	ORF1b:2163	S:375 <sup>*</sup>	S:496 <sup>*</sup>	S:954 <sup>*</sup>	M:63	N:13	
ORF1a:2287	ORF1a:3758	ORF1b:2196	S:376	S:498 <sup>*</sup>	S:969 <sup>*</sup>	M:82	N:63	

Table 5. Brief description of the GISAID patients with a collection date from 2020-11-03 and 2022-02-03

Table 6. Vaccination distribution for the haplotypes with a frequency ≥ 10 %

	Female, n (%)	Male, n (%)	>= 65 y, n (%)	< 65 y, n (%)
GISAID database, n = 4983	2772 (56)	2200 (44)	1368 (28)	3615 (72)
h_GISAID_1, n = 902	496 (55)	405 (45)	268 (30)	634 (70)
h_ GISAID_2, n = 927	577 (62)	350 (38)	327 (35)	600 (65)
h_ GISAID_3, n = 607	347 (57)	260 (43)	132 (22)	475 (78)

Abbreviations: y = years; h\_GISAID = haplotype for GISAID dataset dataset

# Table 7. Haplotypes with 45 non-synonymous mutations for GISAID

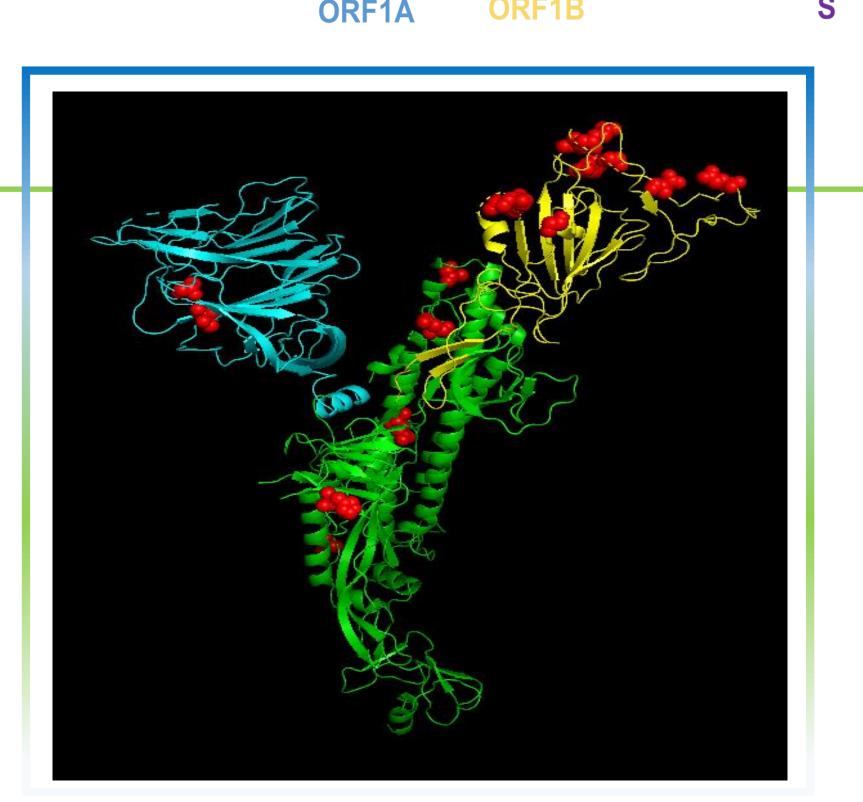
Reference	KAAPP ADVI GPA ATGSSSEQGQNYHDQNLS TT D Q R VTL T PDMGGD
h_GISAID_1	. <b>S</b> . L <b>S</b> L. <mark>S L V</mark>
h_GISAID_2	RTVVIDLPFAR SR YH YYH KF I GEL. I R
h_GISAID_3	. <b>S</b> . L <b>S</b> L. <mark>S L V</mark> . I
	$\leftarrow \qquad \qquad$

	Vaccinated Yes	Expected frequencies Yes	Vaccinated No	Expected frequencies No	p value*
h_GISAID_1, n = 902	477	391	425	511	1.82E-10
h_GISAID_2, n = 927	33	402	894	526	5.88E-161
h_GISAID_3, n = 607	328	263	279	344	1.61E-08

Abbreviations: h\_GISAID = haplotype for GISAID dataset dataset \*Chi-square test

# Table 8. Protein mutations of SARS-Cov-2 for GISAID

ORF1a:856	ORF1b:1000	S:493 <sup>*</sup>	ORF3a:26	N:13
ORF1a:1306	ORF1b:1918	S:496 <sup>*</sup>	<b>ORF3a:64</b>	N:63
ORF1a:1707	S:67 <sup>*</sup>	S:498 <sup>*</sup>	<b>E:9</b>	N:203
ORF1a:2046	S:95 <sup>*</sup>	<b>S:501</b> *	M:3	N:204
ORF1a:2287	S:339 <sup>*</sup>	<b>S:505</b> *	<b>M:19</b>	N:215
ORF1a:2710	S:371 <sup>*</sup>	<b>S:655</b> *	<b>ORF6:20</b>	N:377
ORF1a:2907	S:373 <sup>*</sup>	S:796 <sup>*</sup>	ORF7a:82	
ORF1a:2930	S:375 <sup>*</sup>	<b>S:954</b> *	ORF7a:120	
ORF1a:3758	S:484 <sup>*</sup>	S:969 <sup>*</sup>	<b>ORF7b:18</b>	
ORF1b:662		S:981 <sup>*</sup>	ORF7b:40	



Conclusions: Haplotype analysis suggest the relevance of the genetic variability throughout the SARS-CoV-2 whole-proteome in the response of **COVID-19** patients to the vaccine.

This research work was funded by the European Commission-NextGenerationEU (Regulation EU 2020/2094), through CSIC's Global Health Platform (PTI Salud Global)



*Figure 1.* Common mutations located at the Spike protein