



OPEN

Analysis and forecasting of global real time RT-PCR primers and probes for SARS-CoV-2

Gowri Nayar[✉], Edward E. Seabolt, Mark Kunitomi, Akshay Agarwal, Kristen L. Beck, Vandana Mukherjee & James H. Kaufman

Rapid tests for active SARS-CoV-2 infections rely on reverse transcription polymerase chain reaction (RT-PCR). RT-PCR uses reverse transcription of RNA into complementary DNA (cDNA) and amplification of specific DNA (primer and probe) targets using polymerase chain reaction (PCR). The technology makes rapid and specific identification of the virus possible based on sequence homology of nucleic acid sequence and is much faster than tissue culture or animal cell models. However the technique can lose sensitivity over time as the virus evolves and the target sequences diverge from the selective primer sequences. Different primer sequences have been adopted in different geographic regions. As we rely on these existing RT-PCR primers to track and manage the spread of the Coronavirus, it is imperative to understand how SARS-CoV-2 mutations, over time and geographically, diverge from existing primers used today. In this study, we analyze the performance of the SARS-CoV-2 primers in use today by measuring the number of mismatches between primer sequence and genome targets over time and spatially. We find that there is a growing number of mismatches, an increase by 2% per month, as well as a high specificity of virus based on geographic location.

As the SARS-CoV-2 pandemic grows, an essential method for controlling its spread and determining readiness for the re-opening of public life is through rapid testing. Rapid tests for active SARS-CoV-2 infections are based on reverse transcription polymerase chain reaction (RT-PCR). These tests consist of a forward primer, reverse primer, and probe that together are used to amplify the signal from the targeted virus within a sample. The approach supports rapid and specific identification of the virus, and does not depend on tissue culture or animal cell models. However, RNA viruses evolve over time and a specific PCR test may lose sensitivity as the genotypic distribution of the virus changes or shifts. Phylodynamic studies suggest the mutation rate of SARS-CoV-2 is in the range 1.05×10^{-3} to 1.26×10^{-3} substitutions per site per year, approximately 1.5% variation increase per month¹, consistent with mutation rates reported for other *Coronaviridae*^{2–4}.

Sequence drift also leads to geospatial differences in the virus, resulting in varying test sensitivity by region. This study investigates the effectivity of current SARS-CoV-2 PCR tests over the development of the virus in space and time, and projects how the performance of each may change as the virus undergoes mutation. By taking a global perspective, using specific PCR protocols from several different countries together with genomic data from around the globe, our analysis shows how the existing tests respond differently over both time and location. By analyzing the number of mismatches of the PCR primers with respect to the sequenced SARS-CoV-2 genomes, we can measure how the targeted proteins are mutating. This provides an understanding of possible shortcomings of current tests, and suggests how often we may need to update those tests in the future. Through this work, we observe an average rate of amino acid sequence change of approximately 3% per month for the targeted proteins. Furthermore, we see that the virus genotype is spatially differentiated to the point that inter-country PCR testing already leads to a much higher rate of mismatches.

In support for global pandemic response, several countries have published their RT-PCR protocols. We have collected the primer sequences and protocols developed for six different regions—USA, Germany, China, Hong Kong, Japan, and Thailand—as provided by the WHO⁵. For all six protocols, we collect the forward, reverse, and probe sequences for each specific gene target. Table 1 details the different gene targets for each protocol. Most commonly, the PCR tests target the nucleoprotein (NP), followed by targets in the RNA-dependent RNA polymerase (RdRP) gene, and the envelope small membrane protein (E protein). NP is a structural protein that encapsidates the negative-stranded RNA genome. For other RNA viruses including influenza, the NP sequence is often used for species identification⁶. RNA-dependent RNA polymerase (RdRP) is an enzyme that catalyzes the replication of RNA from an RNA template. The membrane associated RdRP is an essential protein for

IBM Research, San Jose 95120, USA. ✉email: Gowri.Nayar@ibm.com

Country	Target		Sequence	
USA	Nucleoprotein 1	F	GACCCCAAATCAGCGAAAT	
		R	TCTGGTTACTGCCAGTTGAATCTG	
		P	ACCCCGCATTACGTTTGGTGGACC	
	Nucleoprotein 2	F	TTACAAACATTGGCCGCAAA	
		R	GCGCGACATTCCGAAGAA	
		P	ACAATTTGCCCCAGCGCTTCAG	
	Nucleoprotein 3	F	GGGAGCCTTGAATACACAAAA	
		R	TGTAGCACGATTGCAGCATTG	
		P	ACATTGGCACCCGCAATCCTG	
China	ORF1ab	F	CCCTGTGGGTTTTACACTTAA	
		R	ACGATTGTGCATCAGCTGA	
		P	CCGTCTGCGGTATGTGGAAAGTTATGG	
	Nucleoprotein	F	GGGGAACCTTCTCTGCTAGAAT	
		R	CAGACATTTGCTCTCAAGCTG	
		P	TTGCTGCTGCTTGACAGATT	
	Germany	RNA-dependent RNA polymerase	F	GTGARATGGTCATGTGTGGCGG
			R	CARATGTTAAASACACTATTAGCATA
			P1	CCAGGTGGAACRTCATCAGGTGATGC
P2			CAGGTGGAACCTCATCAGGAGATGC	
Envelope small membrane protein		F	ACAGGTACGTTAATAGTTAATAGCGT	
		R	ATATTGCAGCAGTACGCACACA	
		P	ACACTAGCCATCCTTACTGCGCTTCG	
Hong Kong		Orf1b	F	TGGGGCTTTACAGGTAACCT
			R	AACACGCTTAACAAAGCACTC
	P		TAGTTGTGATGCAATCATGACTAG	
	Nucleoprotein	F	TAATCAGACAAGGAAGTATTA	
		R	CGAAGGTGTGACTTCCATG	
		P	GCAAATTGTGCAATTTGCGG	
	Thailand	Nucleoprotein	F	CGTTTGGTGGACCTCAGAT
			R	CCCCACTGCGTTCTCCATT
			P	CAACTGGCAGTAACCA
France	RNA-dependent RNA polymerase IP2	F	ATGAGCTTAGTCCTGTTG	
		R	CTCCCTTTGTTGTGTTGT	
		P	AGATGTCTGTGCTGCCGGTA	
	RNA-dependent RNA polymerase IP4	F	GGTAACTGGTATGATTTCG	
		R	CTGGTCAAGGTTAATATAGG	
		P	TCATACAAACCACGCCAGG	
	Envelope small membrane protein	F	ACAGGTACGTTAATAGTTAATAGCGT	
		R	ATATTGCAGCAGTACGCACACA	
		P	ACACTAGCCATCCTTACTGCGCTTCG	
Japan	Nucleoprotein	F	AAATTTTGGGGACCAGGAAC	
		R	TGGCAGCTGTGTTAGGTCAAC	
		P	ATGTCGCGCATTGGCATGGA	

Table 1. Targeted genes by name by primers from the countries in the study.

Coronavirus replication⁷, and may be a primary target for the antiviral drug remdesivir⁸. The E protein is a small membrane protein involved in assembly, budding, envelope formation, and pathogenesis⁹. The SARS-CoV E protein also forms a Ca²⁺ permeable ion channel that alters homeostasis within cells which leads to the overproduction of IL-1beta^{10,11}.

Results

Primer comparison. We compare a corpus of 61,996 SARS-CoV-2 genomes to the set of published primer-probe sequences. Using the methods described below, we observed high sequence homology for at least 95% of all genomes for most of the PCRs, showing that each primer is able to detect most of the SARS-CoV-2 genomes sequenced at the time of this report. However, it is important to note that this result is biased by the nature of the study, as the genomes tested were the cause of a positive covid-19 infection. Thus to confirm the specificity of the assays, we tested each primer-probe set against a set of 5,000 *Alphacoronavirus*, *Gammacoronavirus*, and

PCR	Percent of hit genomes
America RP	0*
China ORF1ab	99.9
Japan NIID 2019-nCoV N	99.6
America 2019-nCoV N2	99.7
HongKong HKU-N	9.98
Thailand WH-NIC-N	99.7
China N	99.9
Germany E Sarbeco	99.9
France E Sarbeco	99.9
France nCoV IP2	99.8
America 2019-nCoV N1	99.9
France nCoV IP4	99.9
America 2019-nCoV N3	99.9
Germany RdRP-SARsR	99.9
HongKong HKU-ORFb-nsp14	99.7

Table 2. Percent of genomes that are hit by the described PCR test, identified by the country and target gene. *Indicates that the primer is designed to separate the any errant samples within the assay.

Deltacoronavirus. This test resulted in 0 matches, providing evidence for the specificity of the test sets to the *Betacoronavirus* genus, as no test matched the other coronaviruses. These three other genera of coronavirus are, by sequence, the closest in relation to betacoronavirus, and thus to the SARS-CoV-2 genome. Therefore, these results suggest a low probability in false positives for all the global primer-probe sets.

Table 2 shows the percent of genomes hit by each PCR test, labelled by the country and target gene region. The America RP is an additional housekeeping primer/probe set to detect the human RNase P gene to control for non-viral genes in the sample, and therefore, as expected, 0% of the SARS-CoV-2 genomes match with this set. Each primer-probe set will produce a positive assay for at least 95% of the genomes within the test set, suggesting that all the global primer-probe sets function similarly. However, when we study instead the number of mismatches for each primer-probe set and the hit genomes, we can see a disparity in the test sets. The Japan|NIID 2019-nCoV N primer-probe set has 1 mismatch with 99% of the genomes and the Germany|RdRP-SARsR set has 3 mismatches with 99% of genomes. Both these tests have at least 1 mismatch with every genome. A correction for the Japanese set was published, but this consistent mismatch persists with the corrected set as well. Several of the tests contain degenerate nucleotides in order to allow for a larger set of possible positive assays, but this could have the opposite effect. For example, the Germany|RdRP-SARsR reverse sequence contains an S, which denotes either an C or G, but every genome in the corpus contains a T at this locus, which causes the mismatch with every genome. Figure 1 shows the number of mismatches for all genomes created by each PCR, where we can see the range varying from 1796, created by the American N1 primer, to 42 mismatches, created by the French IP2 primer. Thus we observe that the measure of mismatches can be used as a proxy to identify the variation found within the specific gene fragments targeted the tests.

Time analysis

Following the methods described in “Time analysis methods” section, all genomes that fall within the 207 day range are segmented by date of collection and analyzed for mismatches to the various primer-probe sets. Using the metadata provided by GISAID, we were able to recover the date of collection for 50,687 genomes from the corpus. Figure 2 shows the average number of mismatches seen for all primers each day within this range, normalized by the number of genomes sampled in each day. From this analysis, we can see an average of 1.1 mismatches, with a 14% increase in mismatches over the 207 day time range. This corresponds to a ~2% increase per month. To estimate the mutation rate, from Fig. 2, we calculate the best-fit line using least squares, which results in an R^2 value of 0.6. This mutation rate is consistent with the expected rate of mutation of the SARS-CoV-2 virus¹⁻⁴. Thus on average, we can conclude that the regions targetted by the primer-probe sets undergo a similar mutation probability as found across the entire viral genome. Figure 3 shows the distribution of total mismatches and time averaged mismatches for each primer set over the defined time interval. Each violin shows the distribution of the mismatches, total and normalized, for each day and is colored by the gene targetted by that primer-probe set. These results suggest that the tests that target the nucleoprotein exhibit a larger distribution of total and normalized mismatches. However, it is important that it is the China and American tests only that are contributing to this larger distribution, while the Hong Kong and Thailand tests, which also target nucleoprotein, are consist with the other gene targets.

It is important to note that the total number of mismatches occurring is increasing and that many of these mismatches are being sustained in the evolving population. In order to identify a trend, genomes that occur close in time should have smaller change in mismatches than genomes that occur further apart in time. Figure 4 shows this comparison between delta time and delta mismatches for every pair of genomes for the France PCR targeting the RdRP gene (IP4). The graphs for the other PCRs may be found in the supplemental files. Each

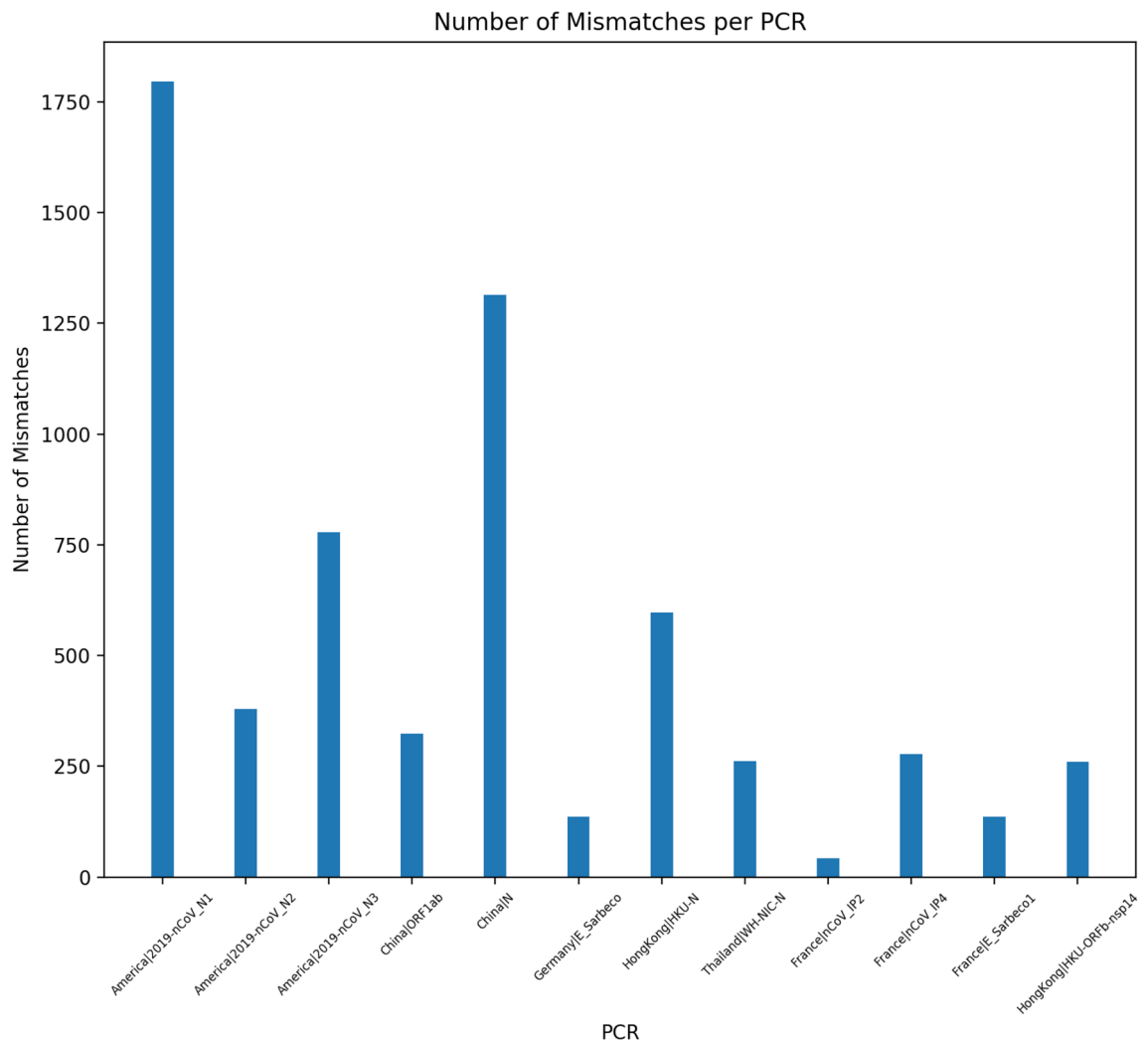


Figure 1. Total number of mismatches each PCR test creates when tested against the full corpus of SARS-CoV-2 genomes. Each PCR test is identified by the country of use and the targeted gene name.

point represents a pairwise comparison of the difference in mismatch plotted over the difference in time. We observe that the delta mismatches grows in variance as the genomes occur further apart in time. Furthermore, the Pearson coefficient is 0.99 between mismatches and the number of genomes sampled in a day for each PCR. This positive linear relationship between the number of genomes and the number of mismatches per day shows that the mismatches occur uniformly across the genomes sampled within a day (rather than a few genomes creating noise in the signal). The data indicates that the virus demonstrated sequence variability in the targeted gene regions and that this variability causes sequence mismatches to increase over time.

Geographical analysis

Geographical stratification is occurring as the SARS-CoV-2 mutates within each geographic location. Following the methods described in “[Geographical analysis methods](#)” section, geospatial analysis is conducted to identify patterns in mismatches found in genomes sequenced within versus outside the country of primer origin. Figure 5 shows that number of mismatches segmented by genomes that occur within the test’s country of origin and outside the test’s country of origin. We normalize the number of mismatches by the number of genomes seen within each category, in order to account for the bias in amount of sequencing performed various regions. For the majority of countries, the number of mismatches in the country is lower than the number of mismatches that occur with genomes sampled outside of the country. This shows that the virus displays localized tendencies within the targeted gene regions, in addition to the spike glycoprotein region that defines the common clade analysis. The two outliers, the Hong Kong and France primers, show a higher percent of mismatches within the country rather than from different countries. This however may be due to the bias in data, as we have the least number of sequences from within Hong Kong, resulting in a higher proportion of in-country mismatches. For the French primer, they perform among the best, with less than 50 mismatches in total, so we can conclude that these mismatches are negligible.

We can also analyze how the country’s primer-probe set change in performance over time, as they experience different variants emerging at different time points within the pandemic. Figure 6 shows the average number

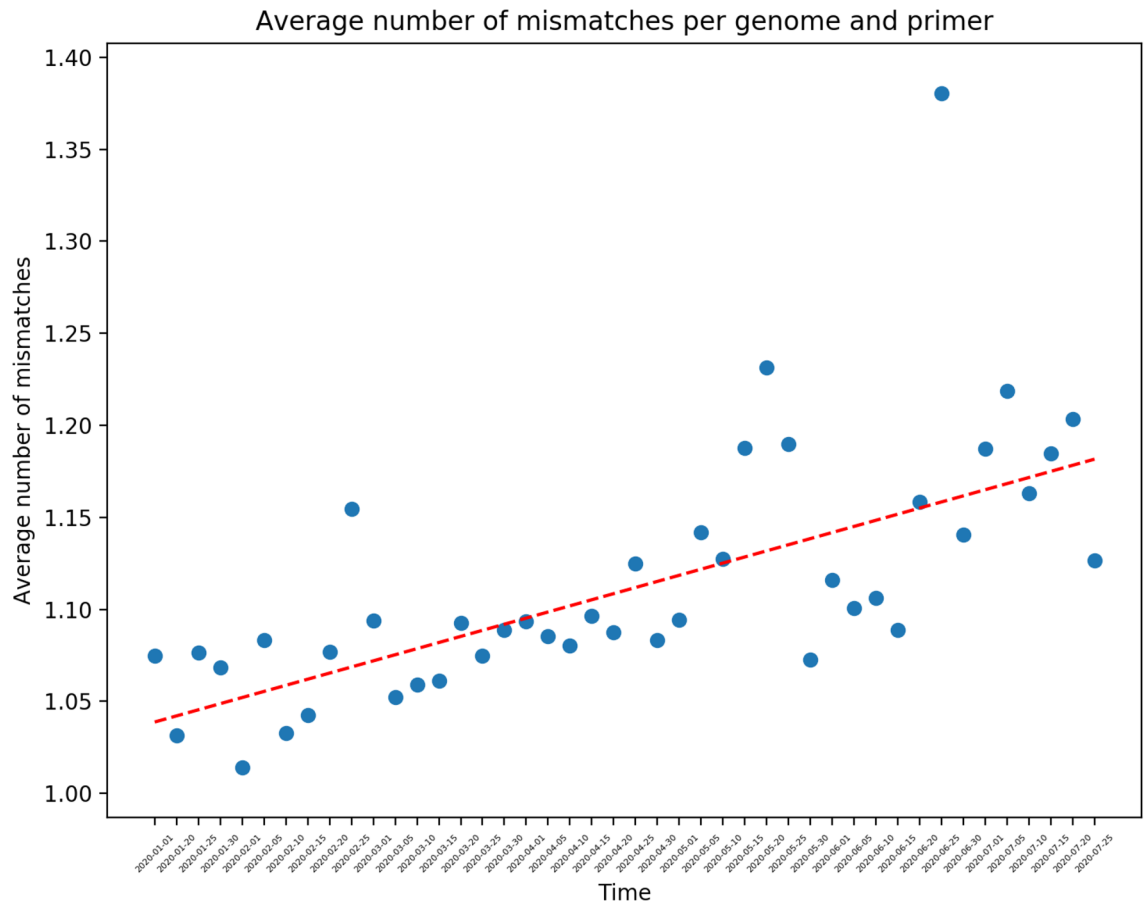


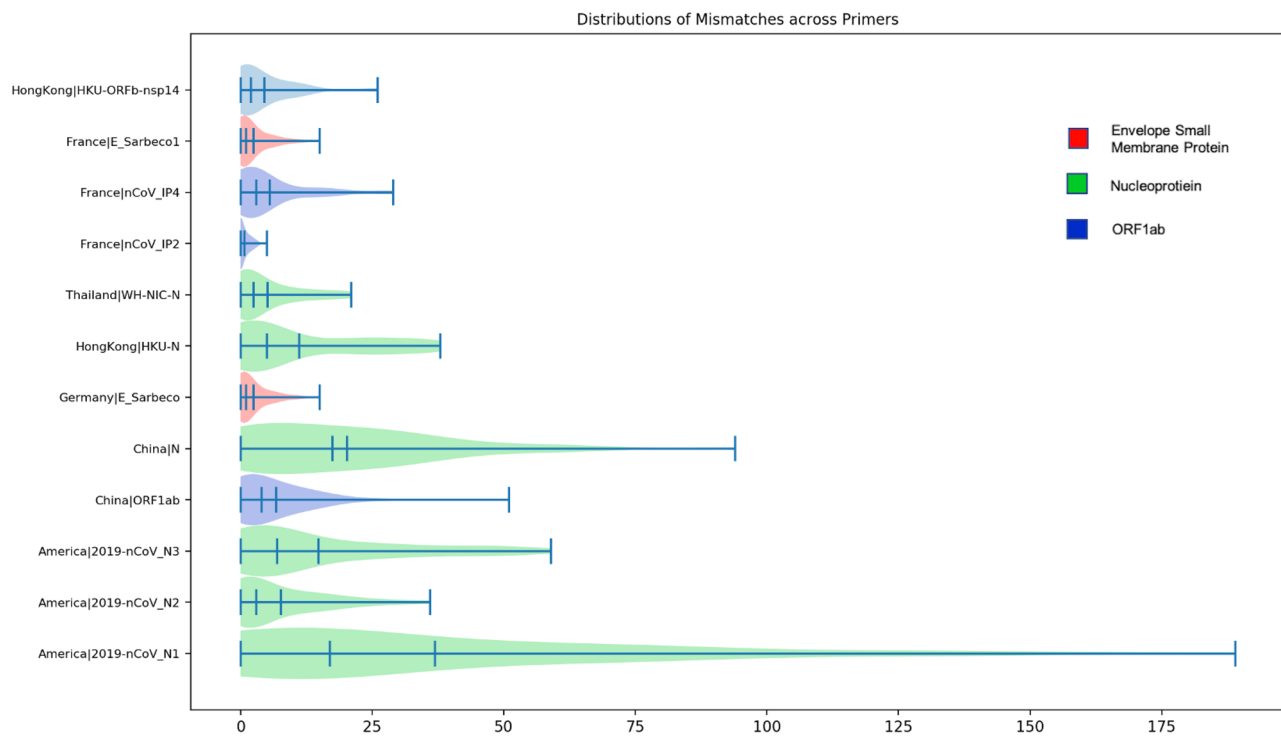
Figure 2. Average number of mismatches for all genomes and all PCR primers separated by the day on which the genome is collected. The dates shown are aggregated over every 5 day period.

of mismatches over time, grouped by the genomes sampled within and outside the country, for one American primer-probe. While the in-country average number of mismatches shows low variability, the out-country average number of mismatches show an increasing diversity in these targeted regions. This elucidates the increasing stratification of the virus depending on the location in which it is replicating. We can see the number of mismatches within the country growing linearly, while the number of mismatches outside the countries have much larger spikes. This is evidence that the high geospatial localization of this virus is reflected in the gene regions targeted by the primer-probe sets, and must be considered when administering the tests globally. The full set of graphs for each primer-probe set tested are available in the supplement.

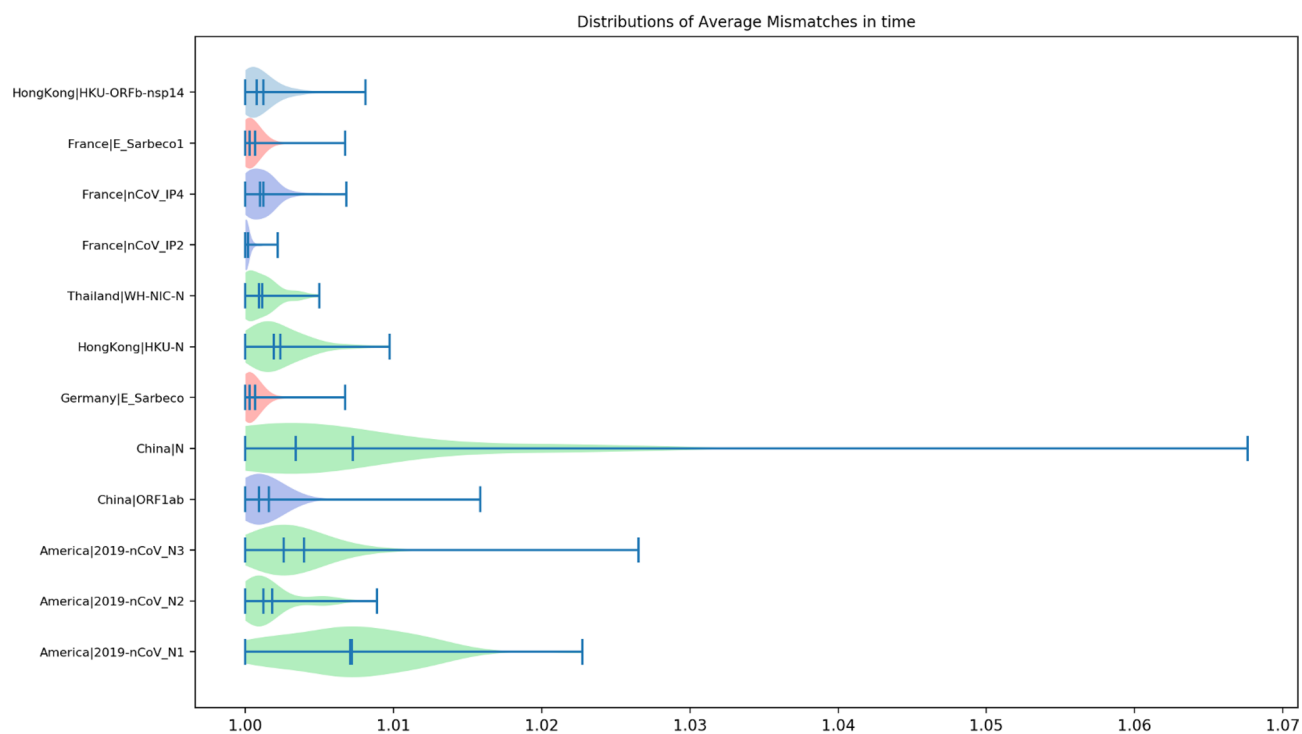
Clade analysis

Geospatial analysis leads to and overlaps with clade definitions, as the clades result in variants emerging in various localities. We can thus examine the number of mismatches segmented by the clades, as defined by Nextstrain, in order to determine if primer-probe sets are skewed in performance to a particular clade. 7125 genomes were placed in Clade 19A, 3706 genomes were placed in Clade 19B, 13437 genomes were placed in Clade 20A, 14,190 genomes were placed in Clade 20B, and 8852 genomes were placed in Clade 20C. Figure 7 shows the number of mismatches for each test per clade, normalized by the number of genomes in clade. This shows definite trends which confirm the geographic specificity of the virus; for example, the American nucleoprotein primers have the highest number of mismatches for clade 19A, which Nextstrain defines as originating from predominantly Asian genomes, while the Chinese primer has the lowest number of mismatches for this clade.

It is important to consider what mutations and gene regions define the clades, as this will impact the significance of a disproportionate number of mismatches. The clades are defined by specific mutations at nucleotide locations, detailed by Nextstrain documentation, which only overlaps with the primer binding regions for 3.37% of the genomes. 58% of the genomes in Clade 20B have a mismatch in the region defining Clade 20B and 1% of the genomes in Clade 19B have a mismatch in the region defining Clade 19B. The genomes in the other clades had no mismatches in the corresponding genomic regions. Therefore, the relationship between the primer mismatches and the genome clades are correlational rather than causal. Thus, the mismatches created by a primer-probe set do not define the clade, but rather are a result of mutations at different locations on the genome that is also being persisted. From this, we provide further evidence that the entire genome is undergoing mutation, and that the targeted gene regions are not resistant to variation.



A



B

Figure 3. Distribution of mismatches for each primer. **(A)** shows the total number of mismatches aggregated for each day within the time range. **(B)** shows the number of mismatches for each day averaged by the number of genomes that occur on a day within the time range.

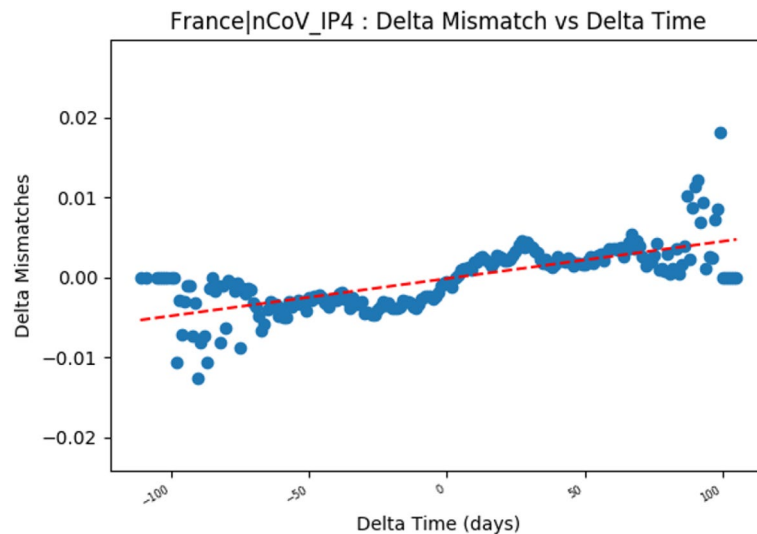


Figure 4. Change in number of mismatches between two occurrences over delta time between the two occurrences for the IP4 primer developed in France. The increasing slope shows that mutations are being sustained as we compare genomes that occur further apart in time. Graphs for all primers are included in the supplement.

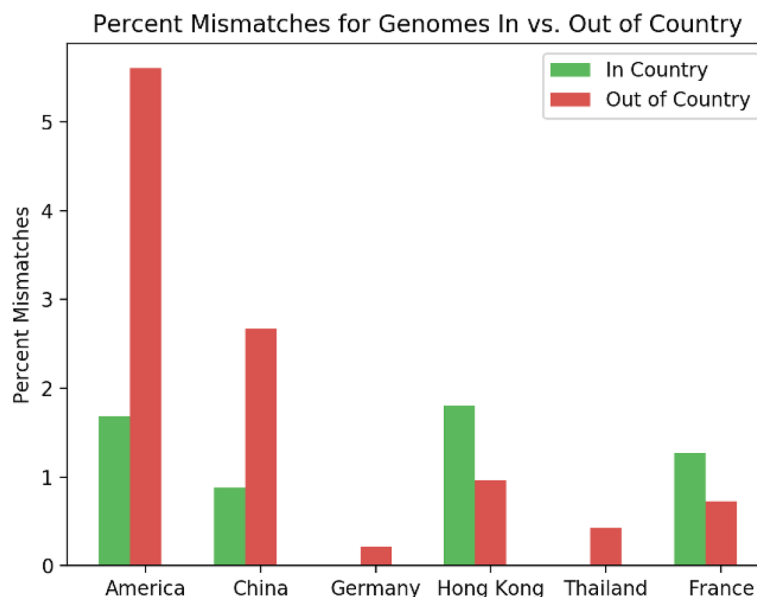


Figure 5. Number of mismatches for each PCR test tested on all SARS-CoV-2 genomes, split between genomes collected within the same country as the test and outside the country. For Japan, 100% of genomes, both in and out of the country, have 1 mismatch, and therefore not shown in the figure. For 9 out of the 11 PCR tests, there are a higher number of mismatches for total genomes that occur outside the country than genomes that occur inside the country.

Discussion

By taking a global perspective on both the SARS-CoV-2 genomes and the common RT-PCR protocols, we are able to highlight important trends within the data. We observe an increasing number of mismatches between the primer and target genome sequence as time progresses. We can also see that the number of mismatches is higher when we compare genomes sampled outside of the country that designed the test compared to within the country. While these metrics do not quantify the performance of the test, they demonstrate a growing divergence between the targeted gene sequences and the test primers.

As shown by D. Bru et al.¹², a single mutation can result in an underestimation of the gene copy number by up to 1000-fold. Our results reveal, today, an average of 1.1 mismatches between the primer and target sequences, with a growth of 2% each month. Understanding copy number is critical to correct interpretation of a PCR assay.



Figure 7. Average number of mutations for each PCR test that occur within each clade, as defined by NextStrain.

target. This indicates that variations in the primer target sequences have not yet reached large enough statistical significance to define a new clade in the Nextstrain phylogeny, although the variants that are present in the primer region may cause a decrease in amplification signal within the assay.

With the emergence of specific mutations that are spreading at faster rates, this analysis becomes more important in evaluating the possible need for primer re-design. As shown with emerging tools, such as PrimerScan and CoV-GLUE, the community is supporting the findings that as the virus mutates, the primer-probe sets must be kept up-to-date. This work further elucidates this as we perform statistical analysis on this data, showing the increasing rate of mismatches across time and location and that these mismatches are being sustained as the virus spreads. The emergence of the B.1.1.7 strain contains mutation in the regions encoding for the envelope small membrane protein and the nucleoprotein, both targeted by the current primers. With the number of cases of SARS-CoV-2 globally, it is highly probable that the genome will mutate in the primer target regions.

Methods

Data description. GISAID has emerged as a leading source of SARS-CoV-2 genomes, containing the largest number of genomes sequences around the world with metadata about the location and time of collection¹⁵. SARS-CoV-2 genomes from the GISAID repository were curated, collecting high quality genomes within the date range Aug 24, 2017–July 31, 2020¹⁶. While this date range precedes the start of the current outbreak, the genome sequences from the earlier points and time serve as a control for comparison. We define high quality genomes as those with less than 1% N within the sequence and less 0.05% unique non-synonymous mutation. By taking these measures, we reduce the noise generated from random mutations or sequencing errors found within the genome. This resulted in a set of 61,996 SARS-CoV-2 genomes, for which we evaluated primer homology.

The WHO has published primers from six countries - China, France, USA, Japan, Germany, Hong Kong, and Thailand⁵. Each protocol published is a RT-PCR assay method, and for each primer set, a forward, reverse and probe sequence is provided⁵. For this study, we use the sequences as provided with no modifications made.

PCR primer comparison. Using the primer sequences and SARS-CoV-2 genomes described above, we perform a sequence comparison. Specifically, we used BLASTN with parameters similar to Primer-BLAST¹⁷. This procedure was verified to account for full alignments of the forward, reverse, and probe sequences of primers¹⁸. The BLAST results are then parsed, ensuring that the forward, reverse, and probe sequences match

a given genome and that the probe sequence is matched spatially in the forward and reverse directions on the genome, and the number of mismatches is aggregated for each PCR sequence and genome. This metric does not necessarily predict whether the PCR test would generate a positive or negative outcome for the particular genome, but rather measures variability within the targeted gene region. Since all genomes included in this corpus are associated with SARS-CoV-2, it can be assumed that they were collected by a positive assay. Mutations in the targeted gene region, over time, can affect the sensitivity of the primers.

Time analysis methods. For each regional test, the primers each target a particular section of the genome derived from various reference genomes. However, as replication and mutation of the virus occurs, these targeted regions of circulating virus genomes accumulate sequence differences from the reference. Thus, the efficacy of the primer may decrease over time. As more mutations accumulate, it is important to measure the rate of mismatch growth between primer sequence and targeted section as a function of time. From this rate it is possible to anticipate when target sequences used in a regional test should be updated. To estimate the mutation rate of the targeted genes over time, we group the genomes by their date of sampling and aggregate the number of mismatches for each day. In order to reduce noise from days with few genomes collected, for any time-based analysis, we consider only those days that have over 100 unique genomes sequenced. With this restriction data is available for a time range between Jan 1, 2020–July 25, 2020, for a total of 207 days. This process removes outlier data that was sequenced prior to the start of the pandemic, including sequences that were collected from non-human hosts.

Geographical analysis methods. As the virus has spread throughout the world, we see particular mutations that are specific to outbreaks by geospatial location. As studies using Bayesian coalescent analysis have shown, high evolutionary rates and fast population growth of the SARS-CoV-2 virus results in increasing diversification of the virus by geographic location¹⁹. To understand how the PCR tests respond differently for genomes collected by country, we first extract the country of sampling for each genome from the fasta header provided by GISAID and then group the number of mismatches found in the genome by in country versus out of country.

Clade analysis methods. SARS-CoV-2 genomes have been categorized into clades to define groups of mutations. For this analysis, we use the clades as indicated by NextStrain, which are defined by frequency and geographic spread. Their script to categorize genomes within the specific clade definitions was used to classify each genome within the dataset²⁰. Furthermore, NextStrain publishes the genome locus that defines each clade, and these loci were compared to the genome location the primer targets bind to. By grouping the number of mismatches for each PCR by the genomes' clade we see how different genetic variations affect the PCR test performance.

Received: 6 January 2021; Accepted: 31 March 2021

Published online: 26 April 2021

References

- Hill, V. & Rambaut, A. Phylodynamic analysis of sars-cov-2 | update 2020-03-06. *virological.org* (2020).
- Gytis, D. *et al.* Mers-cov spillover at the camel-human interface. *eLife* **7**, e31257 (2018).
- Cotten, M. *et al.* Spread, circulation, and evolution of the middle east respiratory syndrome coronavirus. *MBio* **5**, (2014).
- Bacic, R. S. *et al.* Episodic evolution mediates interspecies transfer of a murine coronavirus. *J. Virol.* **71**, 1946–1955 (1997).
- Who in-house assays (2020).
- Burger, H. *et al.* Sequence of the nucleoprotein gene of influenza a/parrot/ulster/73. *Virus Res.* **3**, 35–40. [https://doi.org/10.1016/0168-1702\(85\)90039-5](https://doi.org/10.1016/0168-1702(85)90039-5) (1985).
- Gao, Y. *et al.* Structure of the rna-dependent rna polymerase from covid-19 virus. *Science* **779–82**, (2020).
- Elfiky, A. A. Ribavirin, Remdesivir, Sofosbuvir, Galidesivir, and Tenofovir against SARS-CoV-2 RNA dependent RNA polymerase (RdRp): A molecular docking study. *Life sciences* **253**, 117592 (2020).
- Schoeman, D. *et al.* Coronavirus envelope protein: current knowledge. *Viol. J.* **16**, 1–22 (2019).
- Surya, W. *et al.* Mers coronavirus envelope protein has a single transmembrane domain that forms pentameric ion channels. *Virus Res.* **201**, 61–66 (2015).
- Nieto-Torres, J. *et al.* Severe acute respiratory syndrome coronavirus e protein transports calcium ions and activates the nlrp3 inflammasome. *Virology* (2015).
- Bru, D., Martin-Laurent, F. & Philippot, L. Quantification of the detrimental effect of a single primer-template mismatch by real-time pcr using the 16s rrna gene as an example. *Appl. Environ. Microbiol.* <https://doi.org/10.1128/AEM.02403-07> (2008).
- Christopherson, C., Sninsky, J. & Kwok, S. Phylodynamic analysis of sars-cov-2 genomes. *Nucleic Acids Res.* (2020).
- Carter, L. *et al.* Assay techniques and test development for covid-19 diagnosis. *ACS Cent. Sci.* <https://doi.org/10.1021/acscentsci.0c00501> (2020).
- Shu, Y. & McCauley, J. GISAID: global initiative on sharing all influenza data—from vision to reality. *Eurosurveillance* **22**, 30494 (2017).
- Seabolt, E. *et al.* Ibm functional genomics platform, a cloud-based platform for studying microbial life at scale. *IEEE/ACM Trans. Comput. Biol. Bioinform.* <https://doi.org/10.1109/TCBB.2020.3021231> (2020).
- Camacho, C. *et al.* Blast+: architecture and applications. *BMC Bioinform.* <https://doi.org/10.1186/1471-2105-10-421> (2009).
- Ye, J. *et al.* Primer-blast: a tool to design target-specific primers for polymerase chain reaction. *BMC Bioinform.*, <https://doi.org/10.1186/1471-2105-13-134> (2012).
- Castells, M. *et al.* Evidence of increasing diversification of emerging sars-cov-2 strains. *J. Med. Virol.* <https://doi.org/10.1002/jmv.26018> (2020).
- Hadfield, J. *et al.* Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics* <https://doi.org/10.1093/bioinformatics/bty407> (2018).

Acknowledgements

The authors would like to acknowledge the GISAID Initiative and NCBI for the provision of data.

Author contributions

G.N. conceived the experiment and analysis, M.K. verified the results, E.S. was the architect of the platform used, A.A and K.L.B. performed genome quality analysis, J.H.K and V.M. provided scientific guidance and domain specific knowledge.

Competing interests

The authors declare no competing interests.

Additional information

Correspondence and requests for materials should be addressed to G.N.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021