



Highly multiplexed oligonucleotide probe-ligation testing enables efficient extraction-free SARS-CoV-2 detection and viral genotyping

Joel J. Credle¹ · Matthew L. Robinson² · Jonathan Gunn¹ · Daniel Monaco¹ · Brandon Sie¹ · Alexandra Tchir¹ · Justin Hardick^{2,3} · Xuwen Zheng¹ · Kathryn Shaw-Saliba³ · Richard E. Rothman^{2,3} · Susan H. Eshleman⁴ · Andrew Pekosz⁵ · Kasper Hansen⁶ · Heba Mostafa⁷ · Martin Steinegger⁸ · H. Benjamin Larman¹

Received: 19 August 2020 / Revised: 25 November 2020 / Accepted: 25 November 2020 / Published online: 3 February 2021
© The Author(s), under exclusive licence to United States & Canadian Academy of Pathology 2021

Abstract

There is an urgent and unprecedented need for sensitive and high-throughput molecular diagnostic tests to combat the SARS-CoV-2 pandemic. Here we present a generalized version of the RNA-mediated oligonucleotide Annealing Selection and Ligation with next generation DNA sequencing (RASL-seq) assay, called “capture RASL-seq” (cRASL-seq), which enables highly sensitive (down to ~1–100 pfu/ml or cfu/ml) and highly multiplexed (up to ~10,000 target sequences) detection of pathogens. Importantly, cRASL-seq analysis of COVID-19 patient nasopharyngeal (NP) swab specimens does not involve nucleic acid purification or reverse transcription, steps that have introduced supply bottlenecks into standard assay workflows. Our simplified protocol additionally enables the direct and efficient genotyping of selected, informative SARS-CoV-2 polymorphisms across the entire genome, which can be used for enhanced characterization of transmission chains at population scale and detection of viral clades with higher or lower virulence. Given its extremely low per-sample cost, simple and automatable protocol and analytics, probe panel modularity, and massive scalability, we propose that cRASL-seq testing is a powerful new technology with the potential to help mitigate the current pandemic and prevent similar public health crises.

These authors contributed equally: Matthew L. Robinson, Jonathan Gunn, Daniel Monaco

Supplementary information The online version of this article (<https://doi.org/10.1038/s41379-020-00730-5>) contains supplementary material, which is available to authorized users.

✉ Martin Steinegger
martin.steinegger@snu.ac.kr

✉ H. Benjamin Larman
hlarman1@jhmi.edu

¹ Institute for Cell Engineering, Immunology Division, Department of Pathology, Johns Hopkins University, Baltimore, MD, USA

² Division of Infectious Diseases, Department of Medicine, Johns Hopkins University School of Medicine, Baltimore, MD, USA

³ Department of Emergency Medicine, Johns Hopkins University School of Medicine, Baltimore, MD, USA

Introduction

Several RNA viruses have emerged in recent decades as threats to human health on a global scale (e.g., HIV, MERS, SARS, Ebola, Zika) [1, 2]. In each case, the impact of these viruses would likely have been substantially mitigated by more effective surveillance technologies and contact tracing programs. In the early stages of the COVID-19 pandemic, the crisis was exacerbated by lack of critical supplies in some regions, including the RNA extraction kits required

⁴ Division of Transfusion Medicine, Department of Pathology, Johns Hopkins University, Baltimore, MD, USA

⁵ W. Harry Feinstone Department of Molecular Microbiology and Immunology, Johns Hopkins Bloomberg School of Public Health, Baltimore, MD, USA

⁶ Department of Biostatistics, Johns Hopkins University, Baltimore, MD, USA

⁷ Division of Medical Microbiology, Department of Pathology, Johns Hopkins University, Baltimore, MD, USA

⁸ Biological Sciences & Institute of Molecular Biology and Genetics, Seoul National University, Seoul, South Korea

for reverse transcription polymerase chain reaction (RT-PCR) based molecular testing. As of January 1, 2021, the United States has reported over 20 million COVID-19 cases, over 350,000 deaths, and some of the worst economic losses on record. It is widely recognized that the development of a comprehensive testing infrastructure for large-scale diagnosis and surveillance, which is rapidly reconfigurable for emerging threats, is essential to ending the current pandemic and to preventing future crises due to emerging pathogens.

Current nucleic acid tests (NATs) for SARS-CoV-2 have key limitations. Traditional RT-PCR is relatively inexpensive, but requires a separate labor and time-intensive RNA extraction step prior to nucleic acid amplification. Cartridge-based nucleic-acid tests offer rapid results and minimal sample preparation, but production of cartridges and low-throughput instruments limit scalability; further, most current testing platforms feature single pathogen targets and do not provide strain or clade-level information. Multiplexed PCR platforms and metagenomic next-generation sequencing technologies have demonstrated some advantages for diagnosing infections compared with more traditional approaches, but cost, low sensitivity, and/or complicated informatics have limited adoption for routine use [3]. Platforms such as BioFire, Genmark ePlex and TaqMan array cards are able to identify up to 20–30 targets at a time, but their high per-sample costs (exceeding \$100/test), as well as their inherently low sample throughput, severely limit their utility in the setting of large-scale surveillance efforts [4–7]. In the midst of the COVID-19 pandemic, innovative techniques involving targeted use of NGS have been reported, such as “SwabSeq” [8] and “LAMP-Seq” [9]. While these methods are promising for detection of SARS-CoV-2, they may not be well suited to syndromic panel (multiplex) testing or generalized surveillance, and may provide only limited clade-level information.

Analysis of pathogen-associated RNA, versus DNA, can be valuable for several reasons. A large fraction of clinically important viruses, such as coronaviruses, have RNA genomes, and many have no DNA stage of their lifecycle. All of the viral NIAID Emerging Infectious Diseases Category A and B pathogens are RNA viruses [10]. Further, viral mRNA can also be detected from DNA viruses that cause disease, and may provide a diagnostic advantage over DNA testing by distinguishing between active and latent infections [11, 12]. For cellular pathogens, abundant RNA sequences, such as ribosomal RNA sequences, provide biological amplification compared to analysis of the organism’s genomic DNA, thereby enhancing detection sensitivity [13, 14]. In addition, RNA typically degrades rapidly outside of cells, permitting differentiation between living organisms and environmental/reagent contaminants.

Compared with DNA, RNA tends to be shorter and usually exists in single stranded form, making it more amenable to techniques involving probe hybridization. Finally, simultaneous analysis of viral and host mRNA expression has been shown to provide additional, clinically useful diagnostic and prognostic information about disease states [15–17]. Importantly, the RNA analysis method presented here avoids nucleic acid purification, which is an advantage since limited supplies of the reagents needed for this step of analysis have contributed to the disruption of large-scale SARS-CoV-2 testing efforts in the United States [18].

Prior to analysis of RNA via qPCR or sequencing, purified RNA is typically first converted into complementary DNA (cDNA), via reverse transcription, for subsequent DNA polymerase-mediated amplification. Oligonucleotide probe-ligation assays use ligases to join two oligo probes together when they are hybridized adjacently on a nucleic acid template molecule. Ligated probes can be used similarly to cDNA in downstream molecular assays, such as qPCR or sequencing, with the added benefit that these cDNAs will contain common primer-binding sites for PCR amplification with a single pair of primers. Previous approaches to ligate DNA probes on RNA templates using the T4 DNA ligase were shown to be extremely inefficient [19]. The requirements for high enzyme concentrations, sensitivity to inhibitors and poor target fidelity, had limited the utility of oligo-probe based analysis of RNA. Several years ago, we and others described a modified RNA-mediated oligonucleotide Annealing, Selection, and Ligation with next-generation sequencing (RASL-seq) assay chemistry with enhanced sensitivity. This efficient reaction utilizes the T4 RNA Ligase 2 (Rnl2), which despite being an RNA ligase, efficiently catalyzes ligation of a DNA donor probe to a chimeric acceptor probe containing two bases of ribonucleotides at the ligation junction (Fig. 1A) [20–23]. In addition to the high sensitivity required for pathogen detection, RASL-seq also enables very high levels of probe set multiplexing, potentially providing the means for simultaneous analysis of pathogens, their ancestral lineages, and host immune response (Fig. 1B). By incorporating DNA barcodes into the primers used to amplify the ligation products, a high level of sample multiplexing is also achievable, which enables very high sample throughput and extremely low per-sample cost.

Immobilization of target molecules is commonly employed to permit removal of excess reagents and sample matrix via washing. In any probe ligation assay, it is important that excess probe be removed or destroyed in order to reduce the amount of non-specific background probe ligation [24–26]. In contrast to previously published methods, we have incorporated an oligonucleotide-mediated capture step for pathogen-associated RNA molecules, in an assay we refer to as “capture RASL-seq” or

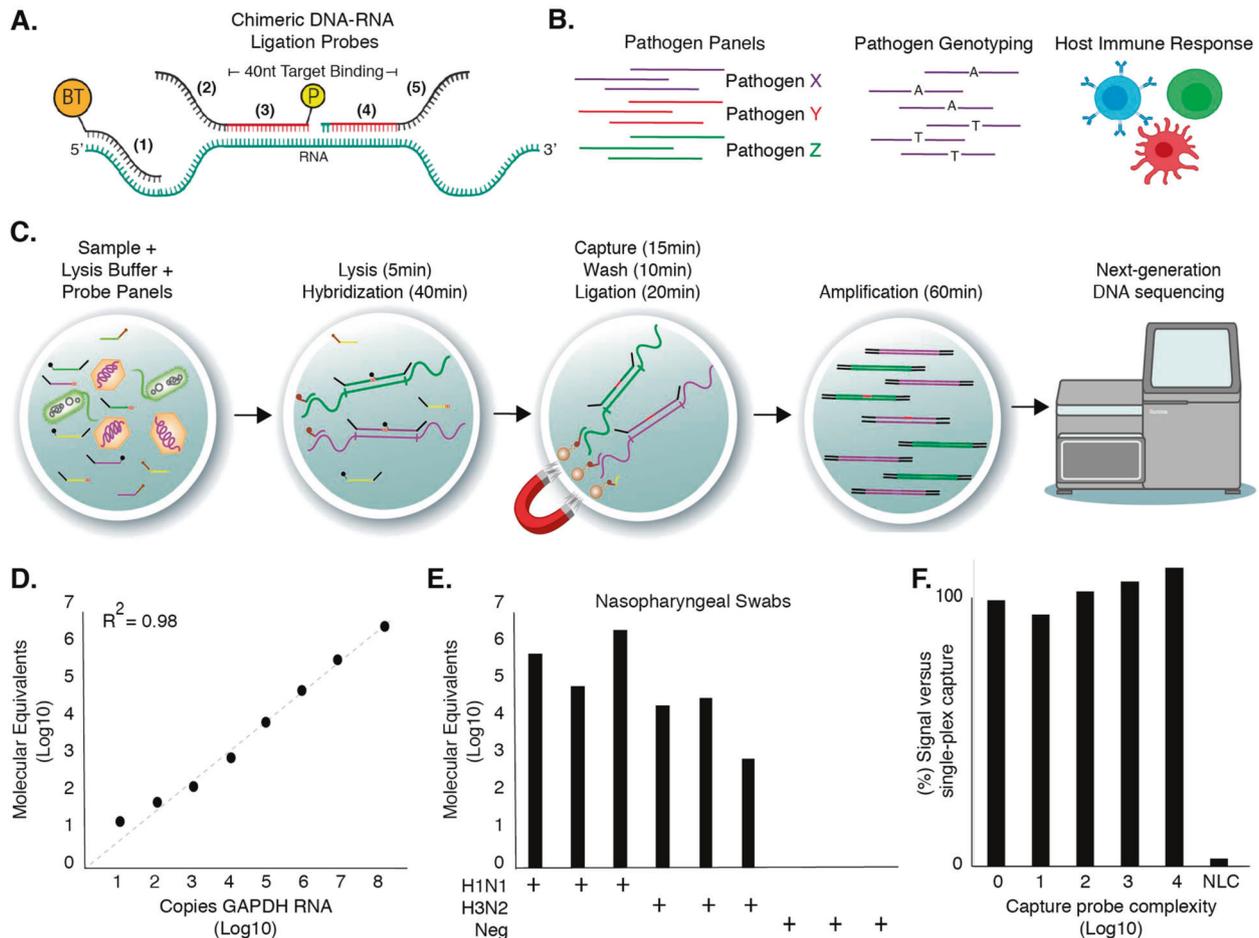


Fig. 1 The cRASL-seq assay. **A** Ligation probe set is composed of a chimeric DNA-RNA 3' acceptor probe (2 & 3) and a phosphorylated 5' donor probe (4 & 5). Each probe contains common PCR primer binding sequences (2 & 5). Two 20 nt target recognition sequences (3 & 4) bring these probes adjacent to one another on a target RNA, enabling their enzymatic ligation. A biotinylated capture probe (1) is used to separate the target sequences from irrelevant materials and excess ligation probes. **B** Complementary oligonucleotide probe ligation assays, which can be performed in a single reaction. **C** Sample (e.g., NP swab specimen) is added to lysis buffer containing cRASL

probes. After lysis and annealing, targets are captured for subsequent probe set ligation and sample-barcoding amplification, followed by amplicon pooling and NGS. **D** Amount of ligation product formed on transcribed GAPDH RNA as a function of input amount; analysis by qPCR of cRASL ligation product. **E** cRASL-seq test on a set of nine blinded NP swabs (unextracted) from six patients with influenza A and three negative controls. **F** Assay performed as in (**E**), with influenza capture probe doped into a background of irrelevant capture probe as indicated. For **D–F** molecular equivalents are calculated by normalization to a PCR spike-in sequence of defined copy number input.

“cRASL-seq” (Fig. 1C). By separating targeted from untargeted RNA molecules, and thus hybridized from unhybridized ligation probes, cRASL-seq permits extremely high assay specificity, which is especially important in the setting of diagnosing infectious diseases—particularly relevant in the early phases of an emerging pathogen threat when community prevalence is low. This method of target capture is distinct from RASL-seq analysis, which relies on immobilized oligo-dT for non-specific capture of polyadenylated mRNA.

We have previously demonstrated that libraries of ligation products are amplified with uniform efficiency [25], so that PCR spike-ins enable precise quantification of the copies of each ligation product formed prior to amplification. Quantification of target molecules has proven useful in

clinical settings, in determining the burden of organism(s) within a clinical specimen. To this end, we have developed several quality controls to assess sample integrity and ligation reaction efficiency. To monitor sample quality, we employed a probe set to measure host GAPDH mRNA. An internal synthetic RNA spike-in based on the M13 phage genome sequence was used to assess ligation efficiency. These quality control measures, combined with the PCR spike-in sequence ensures sample integrity, successful ligation, amplification and sufficient depth of sampling. Furthermore, we demonstrate how cRASL-seq probes can be used for simultaneous SARS-CoV-2 detection and SNP genotyping, which has utility for tracking chains of viral transmission and monitoring clade behavior. Recognizing the urgent need for large-scale testing at minimal cost, we

have optimized and characterized the performance of a streamlined, RNA purification-free protocol for direct analysis of nasopharyngeal (NP) swab and saliva specimens obtained from COVID-19 patients.

Results

We first determined whether the standard RASL-seq assay, which does not involve RNA purification, was compatible with analysis of COVID-19 patient mRNAs present in NP swab specimens, a matrix not previously analyzed using an oligonucleotide probe ligation technique. To this end, we utilized a large panel of RASL-seq probe sets designed to characterize human immune responses in a variety of settings (Table S1). A pool of 1,736 probe sets targeting 240 genes, 154 of which were assessed by analysis of exon-exon junction usage, were included in a standard RASL-seq assay with oligo-dT coated magnetic beads for capture of polyadenylated mRNA transcripts. In this experiment, an average of 727 (± 46) correctly paired probe sets, corresponding to 108 (± 4) genes, were sequenced at least ten times in a given sample. We observed very high reproducibility among technical replicates (average $R^2 = 0.95 \pm 0.03$, Fig. S1). High levels of housekeeping genes were detected as anticipated, and patterns of immunoglobulin gene expression could be reliably measured and were consistent among patients even with very different SARS-CoV-2 viral loads determined by RT-qPCR (Fig. S2). These findings indicated that cRASL-seq analysis of non-polyadenylated, pathogen-associated RNA molecules might also be possible using unextracted NP swab specimens.

We determined the dynamic range of the cRASL-seq assay for detection of 10^1 – 10^8 spiked-in target RNA molecules, observing exceptional linear performance over this range (Fig. 1D, $R^2 = 0.98$). To make the cRASL-seq protocol as fast and inexpensive as possible, we performed extensive optimization to minimize the time and reagents required for each step (Fig. 1C and S3), without compromising assay sensitivity.

We next tested the performance of the cRASL-seq assay in detecting influenza A virus in blinded, previously characterized NP swabs obtained by the Johns Hopkins Center of Excellence for Influenza Research and Surveillance (JHCEIRS). To detect influenza A, we designed a cRASL-seq probe set targeting a conserved sequence within the M-segment, according to our previously-established design principles [20]. A PCR spike-in standard was added at a known concentration to enable precise calculation of ligation product copy numbers. Upon specimen unblinding, we observed large numbers of reads mapping to the correctly paired M-segment probe set in all samples that contained either H1N1 or H3N2 influenza A virus (Fig. 1E). In

contrast, either zero or a small number of reads mapped to the negative control samples, providing a large signal-to-noise ratio that ranged from 10^3 to 10^6 . Saliva samples may be a more convenient sample type to collect in the setting of large-scale surveillance or research studies. We therefore assessed the performance of cRASL-seq on unextracted saliva samples spiked with 1,000 copies of influenza A virus. Undiluted or serially diluted saliva samples were added to the pre-mixed probe sets, virus and lysis buffer in a ratio of 1:1, and then subjected to cRASL-seq testing (Fig. S4). The undiluted saliva sample yielded the same results as the no-matrix control sample, indicating that the cRASL-seq protocol can indeed be used to analyze undiluted, unextracted saliva specimens without loss of detection sensitivity.

We wondered to what extent we could multiplex target capture probes, since we expected magnetic capture bead capacity to limit the level of multiplexing achievable. In order to model increasing probe pool complexity, we serially diluted the influenza M-segment biotinylated capture probe (maintained at a standard 5 pM concentration) into a background of an irrelevant biotinylated capture probe at a concentration increasing up to the binding capacity of the streptavidin coated magnetic beads used in the assay. We compared the signal of the M-segment probe in the singleplex assay (no additional capture probe) to that observed in the model multiplexed assays. We observed >90% of the singleplex signal even in a background of 10,000-fold excess irrelevant capture probe (Fig. 1F). While we have not explicitly tested higher levels of ligation probe multiplexing in this study, previous RASL-seq studies have employed panels of >5,000 probe sets [27]. With appropriate design of non-interfering capture and ligation probe sets, we thus anticipate that we could achieve multiplexing of up to 10,000 probe sets, without technical artifacts.

Having established that cRASL-seq could, in principle, be leveraged into a sophisticated infectious disease diagnostics platform, we next set out to determine whether a universal protocol could be employed for diverse classes of pathogens. Important human pathogens come from all kingdoms of life. We therefore tested the streamlined, extraction-free cRASL-seq protocol for detection of the following: fungal organisms (*Candida albicans* and *Cryptococcus neoformans*) using ITS and 26S/18S rRNA (Fig. 2A, B); acid fast bacteria (*Mycobacterium smegmatis*) (Fig. 2C), gram positive bacteria (*Staphylococcus aureus*) (Fig. 2D), and gram negative bacteria (*Pseudomonas aeruginosa* and *Haemophilus influenzae*) using 16S rRNA (Fig. 2E, F); DNA virus (Human cytomegalovirus) using pp65, US34, UL5, and UL22A mRNA (Fig. 2G); and RNA virus (Zika virus) using genomic RNA (Fig. 2H). Each organism was spiked into a separate reaction in serial dilution. The combined pool of 80 probes targeting all the RNAs were tested together in each reaction (Table S2).

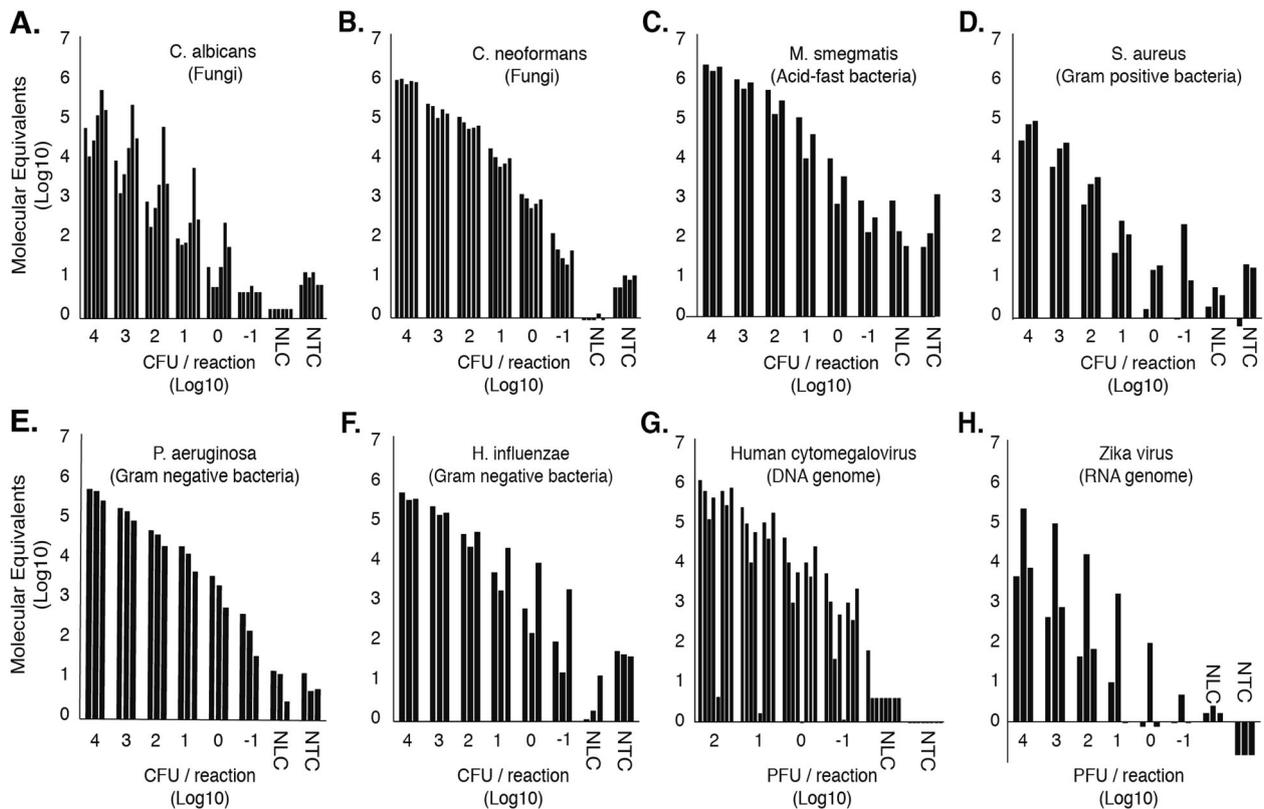


Fig. 2 Universal cRASL-seq assay for pathogen-associated RNA analysis. Each reference organism was serially diluted into PBS and added directly to the lysis buffer and probe pool. NLC, No Ligase Control; NTC, No Template Control. The extraction free protocol of Fig. 1C was performed with all 80 probe sets in a single pool.

Molecular Equivalents are calculated by normalizing read counts to a PCR spike-in sequence of known copy number. Detection of aggregate probes' read counts at a signal $>10\times$ the NTC was used to calculate the assay's limit of detection for each organism.

Negative control reactions included full reactions without any added organism ("no template control", NTC), as well as reactions containing the organisms but lacking the Rnl2 enzyme during the ligation step ("no ligase control", NLC). We considered an organism detected whenever the sum of the probes' normalized read counts was tenfold higher than the corresponding normalized read counts from the NTC sample. In each case, we observed a strong linear correlation between normalized read counts and organism input amount across several logs of abundance, down to limits of detection that ranged from ~ 1.5 to ~ 150 colony or plaque forming units per milliliter (Table 1). Organism reads detected in the NLC and NTC samples most likely arise due to index misassignment (also referred to as index "hopping" or "swapping"). The negative impact of index misassignment on specificity and sensitivity can be mitigated by using unique dual-indices and the use of non-patterned flow cells [28]. cRASL-seq can therefore be used to detect a broad range of pathogens with clinically relevant sensitivity, using a universal, nucleic acid purification-free protocol.

Table 1 Pan-pathogen detection using cRASL-seq.

Organism	Limit of detection
<i>Candida albicans</i>	150 cfu/mL
<i>Cryptococcus neoformans</i>	1.5 cfu/mL
<i>Mycobacterium smegmatis</i>	150 cfu/mL
<i>Staphylococcus aureus</i>	150 cfu/mL
<i>Pseudomonas aeruginosa</i>	1.5 cfu/mL
<i>Haemophilus influenzae</i>	1.5 cfu/mL
Human cytomegalovirus	1.5 pfu/mL
Zika virus	15 pfu/mL

Limits of detection using the multiplex cRASL-seq assay. Bacterial, fungal, or viral stocks were serially diluted in 1x-PBS, pH 7.4 and used in a cRASL-seq assay together with a multiplex probe pool targeting all organisms (*Candida albicans*, *Cryptococcus neoformans*, *Mycobacterium smegmatis*, *Staphylococcus aureus*, *Pseudomonas aeruginosa*, *Haemophilus influenzae*, *Cytomegalovirus*, *Zika virus*). Limits of detection were calculated as the serial dilution with a signal that was $10\times$ above no template controls.

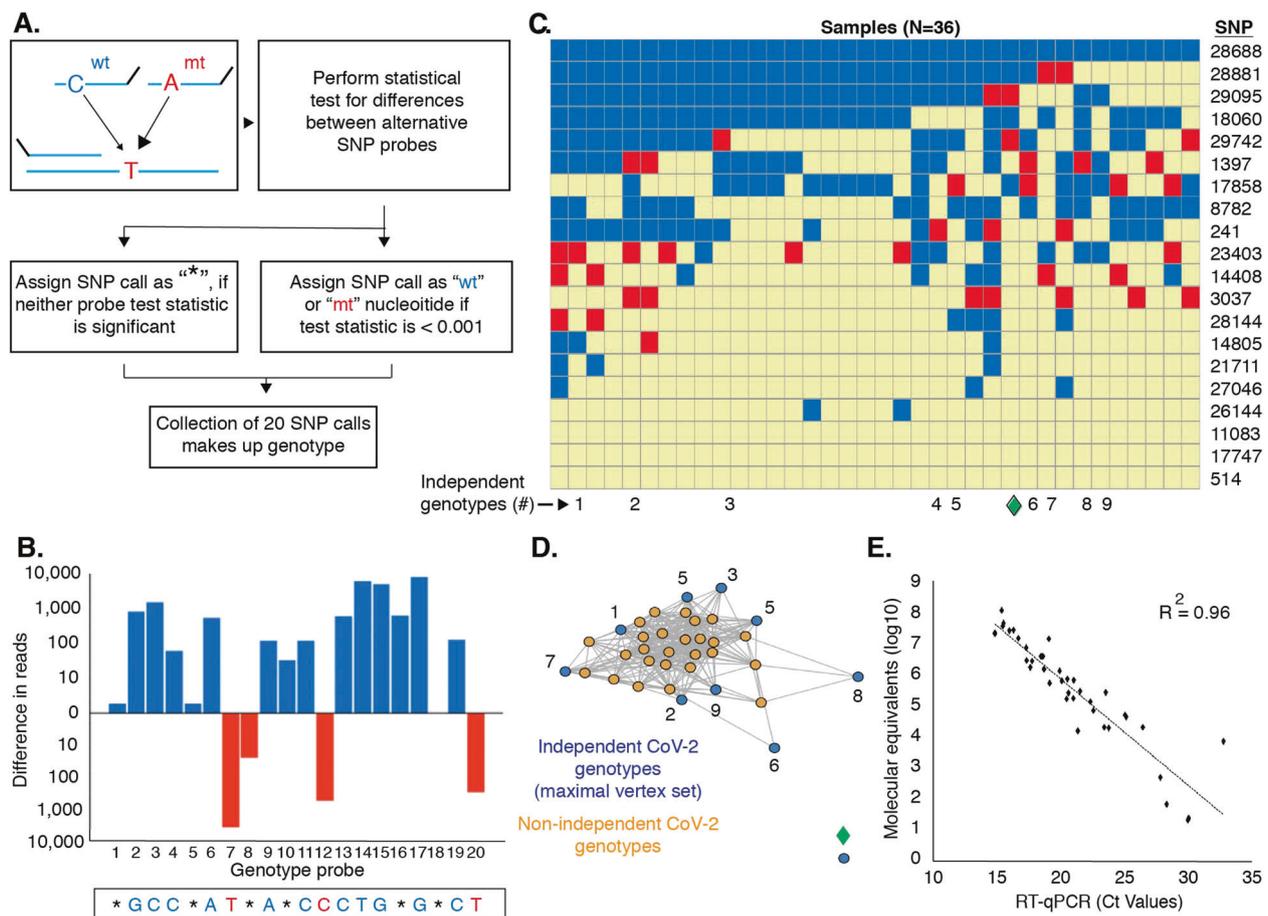


Fig. 3 Multiplexed SNP genotyping of SARS-CoV-2 gRNA directly from unextracted NP swabs. **A** A probe pair is designed with SNP position in the middle of the 5' phosphorylated donor probe. A base-calling algorithm is applied to the reads from each alternative probe. **B** 14 of 20 positions had a base call for the reference Washington isolate, which matched the known genotype without errors. **C** The 35 samples from the set of 40 PCR+ samples analyzed, which had 5 or more base calls, and the reference isolate (green diamond). Red indicates mutant, blue indicates wildtype versus the reference

Analysis of pathogen single nucleotide polymorphisms (SNPs) has utility for distinguishing closely related organisms (e.g., human and zoonotic brucellosis), in tracing chains of viral transmission, and for detecting strain-specific differences in transmission and/or virulence. We therefore tested the ability of cRASL-seq probes to directly genotype SNPs from SARS-CoV-2 RNA in COVID-19 patient NP swab specimens. Genotyping probe sets were designed to share a single 3' acceptor probe, which could pair with two alternative 5' phospho-donor probes corresponding to the alternative genotypes (Fig. 3A, Table S3). We placed the SNP recognition site in the center of the phospho-donor probe to maximally destabilize binding of the mismatched probe. A single base mismatch at this position caused an average of ~10× reduction in signal (range of 1.14–37×) when compared to

Wuhan seafood market isolate. Yellow indicates no call. **D** Network graph depicts each observed genotype (each individual node), two of which are linked if they do not have conflicting SNPs at any position. The blue nodes indicate a maximal vertex set of nine independent genotypes detected among the 35 patient samples that passed QC. **E** Comparison between reads from a SNP typing cRASL-seq probe set in the N gene, versus the Ct values from the corresponding RT-qPCR.

the perfect-match probe among the ten best sampled SNPs of the reference strain. In general, the impact of probe-to-RNA mismatches will depend on the position and type of mismatch, and will likely need to be determined empirically. Biotinylated capture probes were also designed to anneal within 200 nt of each genotyping probe set, to allow for some level of RNA degradation. In a proof-of-concept study, we designed genotyping probes for the 20 most entropic SARS-CoV-2 SNPs reported in the Nextstrain database (queried on 3/16/2020) [29]. These 20 SNPs span across the SARS-CoV-2 genome, ranging from nucleotide 241 to 29,095.

For each SNP, the number of reads from the reference ("wildtype") probe sequence is compared against the number of reads obtained from the non-reference ("mutant") probe sequence. If one of the probe sets is

preferentially incorporated into the ligation product (total reads mapped to probe pair > 100; fold-difference > 1.5; p value < 0.001, binomial test), the base is called and assigned to the position. If the probes do not have sufficient reads or they are not significantly different, no base is called and a “wildcard” is assigned to the position. The string of assigned bases and wildcards are then compared to the strings of corresponding bases from each SARS-CoV-2 genome deposited in the GISAID database. The genotyping assay was first tested using purified reference SARS-CoV-2 gRNA. At the highest input concentration, 2×10^5 copies per reaction, 14 of the 20 bases were called in both technical replicates, and these genotypes matched without error to the sequence of the isolate from which it originated (hCoV19/USA/WA1/2020/EPI_ISL_404895|20200119, Fig. 3B). As the input gRNA amount decreased, significantly fewer bases were called with high confidence. We assessed the reproducibility of the assay by testing eight NP swab specimens in duplicate and comparing the results (Fig. S5). All technical replicates agreed with each other at a cutoff of five SNPs called.

We used the 20-plex SARS-CoV-2 SNP genotyping panel to analyze 40 NP swab specimens obtained early in the pandemic from patients with RT-qPCR proven COVID-19. Of these, 5 or more SNPs could be called in 35 of the samples. These genotypes are displayed in Fig. 3C. To better understand the relationships among the observed genotypes, we used a network graph approach in which genotypes (nodes) are linked (share a connection) if they do not differ in any of the called SNPs (Fig. 3D, rows are variants and columns are patients). Using this network analysis, and by calculating the maximal independent vertex set, we were able to conclude that the infections among these 35 cases could be attributable to at least 9 distinct SARS-CoV-2 ancestral lineages. The reference Washington State isolate was not connected to our Baltimore network. The observed genotypes could additionally be associated with geographic locations, based on their matches to sequenced isolates in the GISAID database (Fig. S6), consistent with other reports of high viral diversity in early local circulation [30].

Finally, we wondered whether viral load could be simultaneously estimated using data from genotyping probe sets alone. After background subtraction using no ligase controls and seasonal coronavirus samples, a probe set targeting the synonymous 28688T>C SNP in the nucleocapsid gene produced positive signals in 36 of 38 PCR + COVID-19 samples that passed QC (out of 40 samples total) and none of the pre-pandemic controls (94.7% sensitivity, 100% specificity, 97.3% overall accuracy). Fig. 3E illustrates how well correlated these values are with RT-qPCR Ct values ($R^2 = 0.96$), indicating that cRASL-seq can be used to simultaneously determine viral load and viral

genotype at high sensitivity, high throughput and very low cost.

Discussion

We have developed a generalized version of the RASL-seq technology, called “capture RASL-seq” or “cRASL-seq”, and demonstrated its utility for highly multiplexed molecular analyses of pathogens and host responses directly from NP swab and saliva specimens, using a universal and streamlined protocol that does not require up-front nucleic acid purification or reverse transcription. The cRASL-seq protocol can easily be performed manually in a biosafety cabinet, does not require centrifugation or vortex steps that risk aerosolizing virus, does not rely on any specialized equipment, and incorporates easily scalable sample bar-coding to dramatically reduce per-sample sequencing cost and increase throughput [31–33]. Though not implemented here, the simple workflow can be easily automated using liquid handling instrumentation. For low complexity probe panels, such as a simple SARS-CoV-2 panel, tens of thousands of sequencing reads per-sample may ultimately provide sufficient sensitivity. A high output Illumina NovaSeq 6000 instrument run, which can generate up to 10^{10} single end reads, would therefore provide sufficient depth to analyze >100,000 samples at once. Another advantage of the cRASL-seq methodology is the extremely low per-probe assay concentration of 5 pM. A typical 100 nmol oligonucleotide synthesis scale will therefore yield a sufficient quantity of probe for tens of millions of 100 μ l cRASL-seq tests, at a per-test probe cost below one cent. The major cost-driving components of the reaction are the ligase, polymerase and magnetic beads. Enzyme production could be scaled up to reduce costs, while less expensive streptavidin capture matrices may be adapted to replace the magnetic beads. At the scale of millions of tests, it is therefore theoretically possible to reduce the per-test reagent and sequencing costs to below one dollar.

Obtaining SARS-CoV-2 genotype information as part of a large-scale surveillance effort would have key benefits. Capturing viral genetics could potentially identify chains of transmission, enabling more effective contact tracing and policymaking, while also detecting and tracking emerging strains with enhanced or diminished transmission potential and/or virulence. There is also intense investigation into the role that host genetics plays in COVID-19 disease severity. As these genotypes are defined, RASL-seq probes that distinguish host alleles could be additionally incorporated into the assay. In this study, we separated the cRASL-seq analysis of SARS-CoV-2 and the RASL-seq analysis of host immune responses into two separate reactions. However, by designing non-interfering probe panels and

balancing the proportion of streptavidin coated magnetic beads, versus oligo-dT coated magnetic beads, it should be straightforward to perform the two assays simultaneously in a single reaction.

cRASL-seq is not without limitations. Unlike unbiased metagenomic sequencing, target detection using cRASL-seq requires a priori selection of target organisms and knowledge of meaningful genomic variants. Another issue, which is associated with sample multiplexing on Illumina sequencing platforms has to do with index misassignment. Index misassignment can negatively impact sensitivity by increasing background, and reduce specificity by causing false positive detection. This problem can be greatly diminished via use of unique dual-indexing primers and non-patterned flow cells [28]. Alternative sequencing platforms (e.g., GenapSys) may be less prone to index misassignment. It is also possible that untemplated probe ligations can contribute to background signal. In addition to comparison with no template controls and no ligase controls, we can estimate the contribution of untemplated probe ligations within a sample by examining the rate of probe mispairing. For a typical COVID-19 sample, we find the probe mispairing rate to be below 2% of the correct probe pairing rate. This rate is particularly low when one considers that for N probe sets, there are N^2 possible mispairing combinations. The contribution of untemplated ligations to overall correctly paired probe ligation signal is therefore negligible.

While we have demonstrated a sensitivity comparable to single-plex RT-qPCR, the limits of detection are governed by overall sequencing depth, which can be reduced by consumption of reads from highly abundant ligation products (due to a high load of a co-infecting virus for example). However, since the ligation products are all amplified with a high degree of uniformity [25], simple RNA spike-in or PCR spike-in sequences can be used to determine the lower limit of detection sensitivity for each reaction. Analysis of host transcripts can also be used to assess sample acquisition sufficiency, a known source of false negative test results [34]. Another important consideration for COVID-19 molecular diagnostics is the turnaround time. When NGS is used to readout the cRASL-seq assay, turnaround time is unlikely to be less than ~24 h with currently available instrumentation. Per-sample sequencing cost considerations will favor analysis of large sample batches, which could further increase turnaround times. For large-scale regional and national level surveillance purposes however, 1–2 day turnaround time may be acceptable, given the costs and limitations associated with alternative methods. Faster, non-NGS-based readouts of cRASL probe ligation products may also be developed. For example, isothermal amplification, followed by array or test-strip hybridization may prove

more applicable in the point-of-care setting. Regardless of the readout, a robust, rapidly reconfigurable, multiplexed, inexpensive and high sample throughput platform for molecular surveillance, such as cRASL-seq, could be used to curb the COVID-19 pandemic and prevent future outbreaks from becoming pandemics.

Materials and methods

Probe design and synthesis

For each target sequence associated with the reference organisms (described in Results, and Supplemental Table S4), we identified 40-mer sites for ligation probes using CATCH [35]. To avoid overlapping probes, we set a stride of 40, allowing no mismatches and bypassing the cover extension in the design.py program (-pl 40 -l 40 -ps 40 -m 0 -e 0). We excluded probes that aligned against any target sequences from the other organisms, with an e value smaller than 10^{-3} , using MMseqs2 [36]. A similar design pipeline was employed for 20 mer capture probes with a final filter step to remove any overlapping ligation and capture probes. Finally, ligation and capture probes were filtered for binding properties using previously reported Primer3 conditions [20]. The design of the immunoglobulin gene expression probe panel was reported previously [20]. Ligation probes and capture probes (3'-diribonucleotide terminated acceptor probes, 5' phosphorylated probes, and biotinylated capture probes) were synthesized by Integrated DNA Technologies (Coraville, IA 52241, USA). Probes were diluted in water to 100 μ M, mixed in equimolar amounts to create multiplexed panels, and then aliquoted and stored at -80 °C (Table S1–S3).

Spike-ins and reference organisms

The synthetic PCR spike-in sequence used for determining molecular equivalents is a 74 nt oligo with a pseudo 40 nt ligated sequence flanked by the external 17 nt PCR1 primer binding sites: 5'-GGAGCTGTCGTTCACTCTGTCTCGGAGCTTACAGTrArU-TGACACTCAATCGGTCGCGTATGATCGGAAGAGCACAC-3'). The 40 nt irrelevant internal sequence is a scrambled version of a ligated GAPDH probe set. Reference organisms were purchased from American Type Culture Collection (ATCC, Manassas, VA) (Table S4). Organisms were reconstituted according to protocols provided by ATCC, aliquoted into single-use samples and stored at -80 °C until used. The Zika virus isolate was from a patient in Cali, Colombia, and was grown in Vero-E6 cells. The infectious titer of the virus (7.1×10^6 pfu/mL) was determined in the culture supernatant by a plaque assay. The full-length GAPDH RNA used in Fig. 1D

(RefSeq: NM_002046.6) was subcloned into a custom vector, linearized and transcribed in vitro using the HiScribe T7 High Yield RNA Synthesis Kit (NEB, Ipswich, MA). GAPDH IVT-RNA was purified by precipitation with lithium chloride followed by column purification using the RNeasy Mini Kit (Qiagen, Hilden, DE). Purified GAPDH IVT-RNA was quantified by nanodrop, aliquoted into single-use samples and stored at -80°C until used for the dynamic range experiments.

Nasopharyngeal swab and saliva specimens

NP swab specimens were collected from patients after informed written consent was obtained, under a protocol approved by the local governing human research protection committee. The Johns Hopkins Influenza Research and Surveillance (JH-CEIRS) program's human subjects protocol was approved by the Johns Hopkins School of Medicine Institutional Review Board (IRB): IRB90001667 and NIH Division of Microbiology and Infectious Diseases: Protocol 15-0103. Unextracted NP swab specimens ($n = 9$) were de-identified, blinded and provided for further analysis. Secondary use of all unextracted COVID-19 patient NP swab specimens ($n = 40$) was exempted by the Johns Hopkins University School of Medicine Institutional Review Board protocol (IRB00086059, and IRB00221396, Table S4). An RT-qPCR diagnostic test (RealStar[®] from Altona Diagnostics, Hamburg, Germany) [37] was performed on all COVID-19 patient NP swab specimens used in this study. NP swab specimens were stored either at -20°C or -80°C . Prior to saliva collection, donors were instructed to pre-rinse with H₂O for 30 s followed by abstaining from eating or drinking for an additional 10 min. Next, saliva (~2–4 mL) was collected in 15 mL conical tubes and centrifuged for 5 min at $3,750 \times g$. Clarified samples were then aliquoted and frozen at -20°C or -80°C until needed.

CRASL and RASL assays

Samples (39 μL) were added to 61 μL of a hybridization reaction mix containing 1X-SSC, 5 pM of each ligation and capture probe, 40 U of Protector RNase Inhibitor (Roche Diagnostics, Mannheim, Germany) and 50 μL of 2X DNA/RNA Shield (Zymo Research, Irvine CA) in a total reaction volume of 100 μL . Reactions were denatured for 5 min at 95°C followed by 40 min annealing at room temperature. 5 μL of a 50/50 slurry of streptavidin coated magnetic beads (Dynabeads MyOne C1, Thermo Fisher Scientific, Waltham MA) in 1x-PBS were added to each reaction and incubated with gentle shaking for 15 min at room temperature. 5 μL of Oligo-(dT)₂₅ beads (Dynabeads, Thermo Fisher Scientific, Waltham MA) were added to each reaction and incubated

with gentle shaking for 15 min at 25°C . Beads were collected on a magnet for 2 min and washed once with 1X-SSC buffer, followed by a final wash with 1X Rnl2 reaction buffer (50 mM Tris-HCl, 10 mM MgCl₂, 5 mM DTT, 1 mM ATP, pH 7.6). 10 μL of the ligation reaction containing 30 U of Rnl2 (Enzymatics, Beverly MA) in 1X Rnl2 buffer was incubated with the beads in suspension for 20 min at 37°C . Following ligation, beads were collected for 2 min on a magnet and then resuspended in 25 μL PCR master mix containing dual-indexing PCR primers with 8 nt i5/i7 barcodes and the P5/P7 Illumina adapters and Herculase-II DNA polymerase (Agilent, Santa Clara, CA). PCR cycling was as follows: an initial denaturing step at 95°C for 2 min, followed by 40 cycles of: 95°C for 20 s, 58°C for 30 s, 72°C for 30 s, with a final extension of 72°C for 3 min.

Library preparation and sequencing

Barcoded PCR products were analyzed individually or as a pool on a 3% agarose gel to confirm amplicon size and purity. Barcoded PCR products were pooled and purified using NucleoSpin PCR Clean-up columns (Mackery Nagel, Duren DE). Pooled libraries were sequenced on an Illumina NextSeq 500 instrument (Illumina, La Jolla CA), using a single-end 50-cycle protocol with a custom read 1 sequencing primer and custom i5/i7 sequencing primers as previously described [20, 25].

Quantification of ligation products

Relative quantification of ligated products (molecular equivalents) were calculated using a synthetic PCR spike-in (described above), which was added to PCR1 reactions at 3,000 or 5,000 molecules/reaction. For samples measured by qPCR, molecular equivalents were calculated using the following equation: $N_{\text{spike}} * 2^{(C_{t\text{spike}} - C_{t\text{ligation}})}$ where N_{spike} is the molecules of spike-in added to PCR1 reactions. For samples analyzed by sequencing, a pseudocount was added to each probe set's read count and the molecular equivalents were calculated by taking the ratio of ligation product read count to spike-in read count, multiplied by the molecules of spike-in added to the PCR1 reaction. For comparison of read counts with Ct values, we performed baseline subtraction on spike-in normalized probe values by subtracting the maximum normalized value of each probe, which occurred in the four negative samples (two no ligase controls and two seasonal coronavirus samples). The corrected values of the best performing genotyping probe set, which targets the N gene at genome position 28,688, was plotted against clinically-determined SARS-CoV-2 RT-qPCR Ct values. We fit an exponential regression to the data, and graphed the resulting data points and regression on a semilog plot. For analysis of host gene expression,

psuedocounted and GAPDH-normalized values were obtained for each probe set. For the analysis of Ig gene expression, the normalized values of all probe sets corresponding to each Ig (3 probes each for IGHA, IGHD, IGHE, and IGHM, and 8 probes for IGHG1–4) were summed to obtain the final values plotted in Fig. S2.

Analysis of sequencing data

Sequencing reads were trimmed to the first 40 bases, demultiplexed and aligned against a reference database of the intended ligation products using exact matching. Genotyping data was analyzed as follows. Genotyping probe pairs were evaluated using an exact binomial test with a null model of equal abundance between the wildtype probe and the mutant probe. A SNP was called if the binomial p value was <0.001 , the total reads mapped to the probe pair was >100 , and the fold change between the probes was >1.5 . A locus with read counts that failed any of these criteria was not called, and thus considered a “wildcard”. Each set of 20 SNP calls and wildcards is referred to as the SARS-CoV-2 genotype associated with a given sample. We utilized a network graph-based approach to visualize the relationships of genotypes detected in the Baltimore COVID-19 cohort. Genotypes were represented as nodes and were connected if there was no disagreement between the pair of genotypes (wildcards could match anything). We utilized the R package *igraph* to determine which set of samples would result in the largest number of unique genotypes (the “maximal independent vertex set”).

Acknowledgements This work was made possible by support from the Johns Hopkins Medicine Discovery Fund in the form of an Innovation Award, a grant from the National Cancer Institute’s Innovative Molecular Analysis Technology (IMAT) Program and an administrative supplement from the Cancer Moonshot initiative (CA202875, CA18099607), a Prostate Cancer Foundation Young Investigator Award, the JHCEIRS (contract number HHSN272201400007C) and a grant from the National Institute of Allergy and Infectious Disease to the HIV Prevention Trials Network (HPTN) Laboratory Center (AI068613). Additional support was provided through the National Research Foundation of Korea (NRF) grant by the Korean government (MEST) (No. 2019R1A6A1A10073437) and the Creative-Pioneering Researchers Program through Seoul National University. We are grateful to Beatriz Parra for providing the Zika virus isolate used in the limit of detection studies and to Priya Duggal and Rebecca Munday for assistance designing host mRNA targeting probes. Many thanks to Haiping Hao and Linda Orzolek at the Johns Hopkins Transcriptomics and Deep Sequencing Core Facility for critical assistance with sequencing.

Author contributions JJC and JG performed the COVID-19 NP swab specimen experiments and assisted with data analysis. MLR assisted in experimental design, probe design, and helped write the paper. SHE helped write the paper. DM analyzed the genotyping data sets. JJC, JG, and BS performed the optimization experiments. AT and JJC performed the reference organisms limit of detection studies. JH, KSS, RER, AP and HM provided clinical specimens and served as technical

consultants. XZ assisted with data analysis. KH designed the human splicing RASL-seq probe panel. MS developed software to design cRASL-seq probe sets, including the genotyping probes. HBL conceived the project and oversaw all aspects of the experimental design and data analyses.

Compliance with ethical standards

Conflict of interest HBL and JJC are inventors on a patent application filed by Johns Hopkins University that covers the use of the cRASL-seq technology to identify infectious organisms and have founded Portal Bioscience to commercialize probe ligation technologies. HBL is a founder of Alchemab and ImmuneID, and is an advisor to CDI Laboratories and TScan Therapeutics.

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

References

1. Morens DM, Fauci AS. Emerging Infectious Diseases: threats to Human Health and Global Stability. *PLOS Pathog.* 2013;9:e1003467.
2. Paules CI, Eisinger RW, Marston HD, Fauci AS. What Recent History Has Taught Us About Responding to Emerging Infectious Disease Threats. *Ann Intern Med.* 2017;167:805–11.
3. Kandathil AJ, Breitwieser FP, Sachithanandham J, Robinson M, Mehta SH, Timp W, et al. Presence of Human Hepatitis E Virus in a Cohort of People Who Inject Drugs. *Ann Intern Med.* 2017;167:1–7.
4. Liu J, Ochieng C, Wiersma S, Stroher U, Towner JS, Whitmer S, et al. Development of a TaqMan Array Card for acute-febrile-illness outbreak investigation and surveillance of emerging pathogens, including Ebola virus. *J Clin Microbiol.* 2016;54:49–58.
5. Hercik C, Cosmas L, Mogeni OD, Wamola N, Kohi W, Houghton E, et al. A Combined Syndromic Approach to Examine Viral, Bacterial, and Parasitic Agents among Febrile Patients: a Pilot Study in Kilombero, Tanzania. *Am J Trop Med Hyg.* 2018;98:625–32.
6. Leber, AL, Everhart, K, Balada-Llasat, JM, Cullison, J, Daly, J, Holt, S, et al. Multicenter Evaluation of BioFire FilmArray Meningitis/Encephalitis Panel for Detection of Bacteria, Viruses, and Yeast in Cerebrospinal Fluid Specimens. *J Clin Microbiol* (2016).
7. Buss SN, Leber A, Chapin K, Fey PD, Bankowski MJ, Jones MK, et al. Multicenter evaluation of the BioFire FilmArray gastrointestinal panel for etiologic diagnosis of infectious gastroenteritis. *J Clin Microbiol.* 2015;53:915–25.
8. Srivatsan S, Han PD, van Raay K, Wolf CR, McCulloch DJ, Kim AE, et al. Preliminary support for a “dry swab, extraction free” protocol for SARS-CoV-2 testing via RT-qPCR. *bioRxiv.* 2020.
9. Schmid-Burgk JL, Li D, Feldman D, Słabicki M, Borrajo J, Strecker J, et al. LAMP-Seq: Population-scale COVID-19 diagnostics using a compressed barcode space. *bioRxiv.* 2020.
10. NIAID Emerging Infectious Diseases/ Pathogens | NIH: National Institute of Allergy and Infectious Diseases. <https://www.niaid.nih.gov/research/emerging-infectious-diseases-pathogens> (2018).
11. Xu M, Arku B, Jartti T, Koskinen J, Peltola V, Hedman K, et al. Comparative Diagnosis of Human Bocavirus 1 Respiratory Infection With Messenger RNA Reverse-Transcription Polymerase Chain Reaction (PCR), DNA Quantitative PCR, and Serology. *J Infect Dis.* 2017;215:1551–7.
12. Graf EH, Simmon KE, Tardif KD, Hymas W, Flygare S, Eilbeck K, et al. Unbiased detection of respiratory viruses by use of RNA sequencing-based metagenomics: a systematic comparison to a commercial PCR panel. *J Clin Microbiol.* 2016;54:1000–7.

13. Murphy SC, Prentice JL, Williamson K, Wallis CK, Fang FC, Fried M, et al. Real-Time Quantitative Reverse Transcription PCR for Monitoring of Blood-Stage *Plasmodium falciparum* Infections in Malaria Human Challenge Trials. *Am J Trop Med Hyg.* 2012;86:383–94.
14. Backstedt BT, Buyuktanir O, Lindow J, Wunder EA Jr, Reis MG, Usmani-Brown S, et al. Efficient detection of pathogenic leptospires using 16S ribosomal RNA. *PLoS ONE.* 2015;10:e0128913.
15. Tsalik EL, Henaio R, Nichols M, Burke T, Ko ER, McClain MT, et al. Host gene expression classifiers diagnose acute respiratory illness etiology. *Sci Transl Med.* 2016;8:322ra311.
16. Woods CW, McClain MT, Chen M, Zaas AK, Nicholson BP, Varkey J, et al. A host transcriptional signature for presymptomatic detection of infection in humans exposed to influenza H1N1 or H3N2. *PLoS ONE.* 2013;8:e52198.
17. Ahn SH, Tsalik EL, Cyr DD, Zhang Y, van Velkinburgh JC, Langley RJ, et al. Gene expression-based classifiers identify *Staphylococcus aureus* infection in mice and humans. *PLoS ONE.* 2013;8:e48979.
18. Esbin MN, Whitney ON, Chong S, Maurer A, Darzacq X, Tjian R. Overcoming the bottleneck to widespread testing: a rapid review of nucleic acid testing approaches for COVID-19 detection. *RNA.* 2020;ma-076232.
19. Nilsson M, Antson DO, Barbany G, Landegren U. RNA-templated DNA ligation for transcript analysis. *Nucleic Acids Res.* 2001;29:578–81.
20. Larman HB, Scott ER, Wogan M, Oliveira G, Torkamani A, Schultz PG. Sensitive, multiplex and direct quantification of RNA sequences using a modified RASL assay. *Nucleic Acids Res.* 2014;42:9146–57.
21. Zhang P, Liu Y, Zhang Y, Liu C, Wang Z, Li Z. Multiplex ligation-dependent probe amplification (MLPA) for ultrasensitive multiplexed microRNA detection using ribonucleotide-modified DNA probes. *Chem Commun (Camb).* 2013;49:10013–5.
22. Nandakumar J, Shuman S. How an RNA ligase discriminates RNA versus DNA damage. *Mol Cell.* 2004;16:211–21.
23. Nandakumar J, Shuman S, Lima CD. RNA ligase structures reveal the basis for RNA specificity and conformational changes that drive ligation forward. *Cell.* 2006;127:71–84.
24. Li H, Qiu J, Fu XD. RASL-seq for massively parallel and quantitative analysis of gene expression. *Curr Protoc Mol Biol Chapter 4, Unit 4.* 2012;13:11–19.
25. Credle JJ, Itoh CY, Yuan T, Sharma R, Scott ER, Workman RE, et al. Multiplexed analysis of fixed tissue RNA using Ligation in situ Hybridization. *Nucleic Acids Res.* 2017;45:e128.
26. Yeakley JM, Shepard PJ, Goyena DE, VanSteenhouse HC, McComb JD, Seligmann BE. A trichostatin A expression signature identified by TempO-Seq targeted whole transcriptome profiling. *PLoS ONE.* 2017;12:e0178302.
27. Yi J, Shen HF, Qiu JS, Huang MF, Zhang WJ, Ding JC, et al. JMJD6 and U2AF65 co-regulate alternative splicing in both JMJD6 enzymatic activity dependent and independent manner. *Nucleic Acids Res.* 2017;45:3503–18.
28. Costello M, Fleharty M, Abreu J, Farjoun Y, Ferriera S, Holmes L, et al. Characterization and remediation of sample index swaps by non-redundant dual indexing on massively parallel sequencing platforms. *BMC Genom.* 2018;19:332.
29. Hadfield J, Megill C, Bell SM, Huddleston J, Potter B, Callender C, et al. Nextstrain: real-time tracking of pathogen evolution. *Bioinformatics.* 2018;34:4121–3.
30. Thielen PM, Wohl S, Mehoke T, Ramakrishnan S, Kirsche M, Falade-Nwulia O, et al. Genomic Diversity of SARS-CoV-2 During Early Introduction into the United States National Capital Region. *medRxiv.* 2020;2020.08.13.20174136.
31. Abdalhamid B, Bilder CR, McCutchen EL, Hinrichs SH, Koepsell SA, Iwen PC. Assessment of Specimen Pooling to Conserve SARS CoV-2 Testing Resources. *Am J Clin Pathol.* 2020;153:715–8.
32. Aragón-Caqueo D, Fernández-Salinas J, Laroze D. Optimization of group size in pool testing strategy for SARS-CoV-2: A simple mathematical model. *J Med Virol.* 2020;92:1988–94.
33. Hogan CA, Sahoo MK, Pinsky BA. Sample Pooling as a Strategy to Detect Community Transmission of SARS-CoV-2. *JAMA.* 2020;323:1967–9.
34. Tang YW, Schmitz JE, Persing DH, Stratton CW. Laboratory Diagnosis of COVID-19: Current Issues and Challenges. *J Clin Microbiol.* 2020;58:e00512–20.
35. Metsky HC, Siddle KJ, Gladden-Young A, Qu J, Yang DK, Brehio P, et al. Capturing sequence diversity in metagenomes with comprehensive and scalable probe design. *Nat Biotechnol.* 2019;37:160–8.
36. Steinegger M, Soding J. MMseqs2 enables sensitive protein sequence searching for the analysis of massive data sets. *Nat Biotechnol.* 2017;35:1026–8.
37. Uhteg K, Jarrett J, Richards M, Howard C, Morehead E, Geahr M, et al. Comparing the analytical performance of three SARS-CoV-2 molecular diagnostic assays. *J Clin Virol.* 2020;127:104384.