

344 Supplementary Materials and Methods.

345 Patient Recruitment:

346 In an effort to identify patients isolating in similar residential settings, the patient population
347 focused on University California San Diego (UCSD) students isolating in the established, on-
348 campus UCSD isolation dorm housing. Cases were identified through the UCSD Health system
349 as COVID-19 positive outpatients with a positive anterior nares clinical RT-qPCR assay from the
350 UCSD EXCITE (EXpedited COVID-19 IdenTification Environment) laboratory. Patients were
351 recruited to the study via phone call, enrolled into an IRB-approved study (UCSD protocol
352 200477), and confirmed to be active UCSD students isolating in the isolation dorms. Three
353 students were enrolled in the study: two of the students isolated in the on-campus isolation
354 dorms, and the third isolated in their on-campus residence (graduate housing with similar
355 architecture and design as the on-campus isolation dorms).

356 Surface Swabbing:

357 For each paired sample, two 1 mL sample collection tubes (ThermoFisher Scientific, 3740TS)
358 were prepared. One tube contained 800 μ L of 0.5% w/v sodium dodecyl sulfate (SDS) (Acros
359 Organics, 230420025) in water and was used for detecting SARS-CoV-2, and the second tube
360 contained 95% spectrophotometric-grade ethanol solution (Sigma-Aldrich #493511) which was
361 designated for 16S sequencing. To recover genetic material from the surfaces, a prewashed
362 cotton swab (Puritan, 806-WC) was pre-moistened with the ethanol solution and then used to
363 vigorously swab the surface. The cotton end of the swab was then placed back into the sample
364 collection tube and broken at the designated break point. The process was then immediately
365 repeated on an adjacent site of the same surface with a flocculated tip swab (Affordable IHC
366 Solutions) pre-moistened with the SDS solution, minimizing overlap between swabbed areas.
367
368

369 Viral Nucleic Acid Extraction and RT-qPCR:

370 Swabs stored in SDS were subjected to SARS-CoV-2 RT-qPCR detection following methods
371 previously described (1). Briefly, 150 μ L of the SDS solution were extracted with Omega
372 MagBind Viral DNA/RNA kit (Omega Bio-Tek, M6246) on Kingfisher Flex (ThermoFisher
373 Scientific) instruments. Viral gene detection was performed using a miniaturized TaqPath™
374 COVID-19 Combo Kit (ThermoFisher Scientific, A47814) assay on a QuantStudio 7 Pro with a
375 384-well sample block (ThermoFisher Scientific).
376

377 Microbial Nucleic Acid Extraction:

378 Sample plating and extractions of all surface swabs were carried out in a biosafety cabinet
379 Class II in a BSL2+ facility. Cotton tipped swabs suspended in 95% ethanol were plated into
380 bead plates from 96 MagMAX™ Microbiome Ultra Nucleic Acid Isolation Kits (A42357 Thermo
381 Fisher Scientific, USA). Following the KatharoSeq low biomass protocol (2), each sample
382 processing plate included eight positive controls consisting of 10-fold serial dilutions of a
383 microbial standard consisting of a gram negative *Paracoccus spp.* and gram positive *Bacillus*
384 *subtilis* ranging from 5 to 50 million cells per extraction, and 3 negative controls (Blanks,
385 sample-free lysis buffer). Nucleic acid extraction and purification was performed following
386 methods previously described (3). Briefly, samples were extracted in plates using the

387 MagMAX™ Microbiome Ultra Nucleic Acid Isolation Kit (Applied Biosystems™), following
388 manufacturer specifications, in KingFisher Flex™ robots (Thermo Fisher Scientific, USA),
389 including a bead beating step in a TissueLyser II (Qiagen, Germany) at 30 Hz for 2 min.

390

391 16S Sequencing:

392 16S rRNA gene amplification was performed according to the Earth Microbiome Project protocol
393 (4). Briefly, the V4 region of the 16S rRNA gene was targeted for amplification in a miniaturized
394 reaction (5) using the 515f-806r primers with Golay error-correcting barcodes. Amplicons were
395 pooled at equal volumes and the pool was purified with a QIAquick PCR purification kit
396 (QIAGEN). The pooled libraries were sequenced on a MiSeq (Illumina) instrument with a MiSeq
397 Reagent 300 cycle v2 Kit, with the appropriate sequencing primers.

398

399 Estimating genomic equivalents and microbial biomass:

400 To estimate viral genomic equivalents for each sample, we used published standard curves
401 relating average Cqs from RT-qPCR to known SARS-CoV-2 viral particle concentrations (in
402 GE's from digital droplet PCR) used to inoculate a variety of indoor surfaces (1). The equation
403 used depended on which qualitative category the surface materials belonged to: rough (carpet,
404 fabric) or smooth (e.g., acrylic, steel, glass, ceramic tile). The relationship between Cqs and
405 GEs for rough materials is $[GEs = -0.52 \times (Avg\ Cq) + 39.90]$ while for smooth materials the
406 equation used was $[GEs = -0.77 \times (Avg\ Cq) + 40.41]$.

407

408 We used an equivolumetric sequencing library pooling approach which allowed us to correlate
409 biomass to 16S amplicon counts (6).

410

411 Data processing:

412 16S sequences were demultiplexed, quality filtered, and denoised with Deblur (7) in Qiita (8)
413 using default parameters. Resulting feature tables were processed using QIIME2 (9).
414 Sequencing data available in Qiita study ID: 13957.

415

416 Katharoseq:

417 In addition to the 381 samples that underwent 16S sequencing, three negative controls (blanks)
418 and eight positive controls (a serially diluted bacterial stock, see Microbial Nucleic Acid
419 Extraction) were included in each 96-well extraction plate. The positive controls were used to
420 determine the threshold read count for which at least 80% of sequencing reads align to the
421 positive controls (10).

422

423 Alpha Diversity:

424 To explore the relationship between microbial biomass and SARS-CoV-2 status, we compared
425 the estimated SARS-CoV-2 viral load in GEs and the number of raw 16S reads for all samples.
426 The Pearson correlation coefficient was calculated to determine if the two measurements had a
427 linear relationship $[\log(16S\ Read\ Counts), \log(GE's)]$. The relationship between biomass (16S
428 read count) and SARS-CoV-2 detection status (Detected/Not Detected) for samples in the
429 same room type was tested with a Kruskal-Wallis H test. For the stringently filtered feature
430 tables, differences in Faith's Phylogenetic Diversity (Faith's PD) between SARS-CoV-2

431 detection status within each room were also tested using a Kruskal-Wallis H test. 2D Figures
432 were made using matplotlib (11).

433

434 Beta Diversity:

435 We used the unweighted Unifrac phylogenetic distance (12-13) to explore how the microbial
436 samples compare to each other. To quantify the effect size of different categorical variables on
437 our data, redundancy analysis (RDA) was applied to the unweighted Unifrac principal
438 coordinates. RDA estimates the contributions of individual and combined effects of multiple
439 covariates using the *varpart* function in R to perform linear constrained ordination (14). 2D
440 Figures were made using matplotlib (11) and EMPeror (15).

441

442 Differential Abundance:

443 To prepare the data for differential abundance we filtered the unrarefied feature table to exclude
444 features present in fewer than 10 samples and samples with depth less than 1000. This resulted
445 in a table of 258 samples and 1047 sOTUs. We performed multinomial regression using
446 Songbird (16) accounting for viral detection status, apartment, surface type, and indoor space
447 classifier as covariates. We used 5000 epochs and a learning rate of 0.0001 as
448 hyperparameters. Additionally, we specified a 3:1 split of training:testing samples for cross
449 validation. To ensure that our model was not overfitting we fit a null regression model with no
450 covariates using the same hyperparameters. Comparing the two models we found a positive
451 pseudo-Q² value of 0.059, indicating that our regression model outperformed the null model.

452

453 Random Forest Classifier:

454 We performed machine learning analysis on the bacterial portion of the built environment
455 surface microbiome from 16S sequencing to predict the samples' SARS-CoV-2 status from
456 paired RT-qPCR detection results. Random forest classifiers were trained and tested following a
457 leave-one-site-out-cross-validation (LOSOCV) approach: the classifier was trained with samples
458 from N-1 sites and its performance was tested in the remaining site using a precision-recall
459 curve (Area Under the Precision Recall Curve (AUPRC), and Relative AUPRC). Classifiers were
460 trained on sOTU-level features with tuned hyperparameters as 20-time repeated, LOSOCV, with
461 sites resolved at the apartment_id (Fig. 2A) and room_type (Fig. 2B) levels using the R caret
462 package(17). The classifiers' performance was evaluated with AUPRC based on the samples'
463 SARS-CoV-2 status predictions of the holdout test site using the R PRROC package (18). The
464 importance of each sOTU for the prediction performance of the classifiers was estimated by the
465 built-in random forest scores in a 100-fold cross validation. We ranked the top 32 important
466 features by their average ranking of importance scores across the 100 classification models.
467 Relevant codebase for machine learning analysis is available at
468 <https://github.com/shihuang047/crossRanger> and is based on random forest implementation
469 from R ranger package (19).

470

471 Phylogenetic Tree visualization:

472 To identify phylogenetic clades important for the prediction of SARS-CoV-2 status from
473 environmental surface samples we visualized the top 32 important features identified by the

474 random forest classifier and the ranked differentially abundant features between SARS-CoV-2
475 status groups from multinomial regression using EMPress (20).

476

477 *3D Mapping:*

478 3D models were provided by UC San Diego's Housing, Dining, and Hospitality department. A
479 circular target was placed on all swabbed locations in each apartment. 3D coordinates were
480 picked following published methods (ref) (<https://github.com/MolecularCartography/ili>), and
481 merged with viral load (in GEs) data for visualization. 3D models and merged data (coordinates
482 and viral load) were visualized in ili (21).

483

484 References

- 485 1. Salido RA, Cantú VJ, Clark AE, Leibel SL, Foroughshafiei A, Saha A, Hakim A, Nouri
486 A, Lastrella AL, Castro-Martínez A, Plascencia A, Kapadia BK, Xia B, Ruiz CA,
487 Marotz CA, Maunder D, Lawrence ES, Smoot EW, Eisner E, Crescini ES, Kohn L,
488 Vargas LF, Chacón M, Betty M, Machnicki M, Wu MY, Baer NA, Belda-Ferre P,
489 Hoff P De, Seaver P, Ostrander RT, Tsai R, Sathe S, Aigner S, Morgan SC, Ngo
490 TT, Barber T, Cheung W, Carlin AF, Yeo GW, Laurent LC, Fielding-Miller R, Knight
491 R. 2021. Analysis of SARS-CoV-2 RNA Persistence across Indoor Surface
492 Materials Reveals Best Practices for Environmental Monitoring Programs.
493 *mSystems* <https://doi.org/10.1128/MSYSTEMS.01136-21>.
- 494 2. Minich JJ, Zhu Q, Janssen S, Hendrickson R, Amir A, Vetter R, Hyde J, Doty MM,
495 Stillwell K, Benardini J, Kim JH, Allen EE, Venkateswaran K, Knight R. 2018.
496 KatharoSeq Enables High-Throughput Microbiome Analysis from Low-Biomass
497 Samples. *mSystems* 3.
- 498 3. Shaffer JP, Marotz C, Belda-Ferre P, Martino C, Wandro S, Estaki M, Salido RA,
499 Carpenter CS, Zaramela LS, Minich JJ, Bryant M, Sanders K, Fraraccio S,
500 Ackermann G, Humphrey G, Swafford AD, Miller-Montgomery S, Knight R. 2021. A
501 comparison of DNA/RNA extraction protocols for high-throughput sequencing of
502 microbial communities. *Biotechniques* btn-2020-0153.
- 503 4. Thompson LR, Sanders JG, McDonald D, Amir A, Ladau J, Locey KJ, et al. A
504 communal catalogue reveals Earth's multiscale microbial diversity. *Nature*.
505 2017;551:457–63.
- 506 5. Minich JJ, Humphrey G, Benitez RAS, Sanders J, Swafford A, Allen EE, Knight R.
507 2018. High-Throughput Miniaturized 16S rRNA Amplicon Library Preparation
508 Reduces Costs while Preserving Microbiome Integrity. *mSystems* 3.
- 509 6. Cruz GNF, Christoff AP, de Oliveira LFV. 2021. Equivolumetric Protocol Generates
510 Library Sizes Proportional to Total Microbial Load in 16S Amplicon Sequencing.
511 *Front Microbiol* 12.
- 512 7. Amir A, McDonald D, Navas-Molina JA, Kopylova E, Morton JT, Zech Xu Z, Kightley
513 EP, Thompson LR, Hyde ER, Gonzalez A, Knight R. 2017. Deblur Rapidly
514 Resolves Single-Nucleotide Community Sequence Patterns. *mSystems* 2.
- 515 8. Gonzalez A, Navas-Molina JA, Kosciolk T, McDonald D, Vázquez-Baeza Y,
516 Ackermann G, DeReus J, Janssen S, Swafford AD, Orchanian SB, Sanders JG,
517 Shorenstein J, Holste H, Petrus S, Robbins-Pianka A, Brislawn CJ, Wang M,
518 Rideout JR, Bolyen E, Dillon M, Caporaso JG, Dorrestein PC, Knight R. 2018.
519 Qiita: rapid, web-enabled microbiome meta-analysis. *Nat Methods* 15:796–798.

520 9. Bolyen E, Rideout JR, Dillon MR, Bokulich NA, Abnet CC, Al-Ghalith GA, Alexander H,
521 Alm EJ, Arumugam M, Asnicar F, Bai Y, Bisanz JE, Bittinger K, Brejnrod A,
522 Brislawn CJ, Brown CT, Callahan BJ, Caraballo-Rodríguez AM, Chase J, Cope
523 EK, Da Silva R, Diener C, Dorrestein PC, Douglas GM, Durall DM, Duvallet C,
524 Edwardson CF, Ernst M, Estaki M, Fouquier J, Gauglitz JM, Gibbons SM, Gibson
525 DL, Gonzalez A, Gorlick K, Guo J, Hillmann B, Holmes S, Holste H, Huttenhower
526 C, Huttley GA, Janssen S, Jarmusch AK, Jiang L, Kaehler BD, Kang K Bin, Keefe
527 CR, Keim P, Kelley ST, Knights D, Koester I, Kosciulek T, Kreps J, Langille MGI,
528 Lee J, Ley R, Liu Y-X, Lofffield E, Lozupone C, Maher M, Marotz C, Martin BD,
529 McDonald D, McIver LJ, Melnik A V., Metcalf JL, Morgan SC, Morton JT, Naimey
530 AT, Navas-Molina JA, Nothias LF, Orchanian SB, Pearson T, Peoples SL, Petras
531 D, Preuss ML, Pruesse E, Rasmussen LB, Rivers A, Robeson MS, Rosenthal P,
532 Segata N, Shaffer M, Shiffer A, Sinha R, Song SJ, Spear JR, Swafford AD,
533 Thompson LR, Torres PJ, Trinh P, Tripathi A, Turnbaugh PJ, Ul-Hasan S, van der
534 Hooft JJJ, Vargas F, Vázquez-Baeza Y, Vogtmann E, von Hippel M, Walters W,
535 Wan Y, Wang M, Warren J, Weber KC, Williamson CHD, Willis AD, Xu ZZ,
536 Zaneveld JR, Zhang Y, Zhu Q, Knight R, Caporaso JG. 2019. Reproducible,
537 interactive, scalable and extensible microbiome data science using QIIME 2. *Nat*
538 *Biotechnol* 37:852–857.

539

540 10. Minich JJ, Zhu Q, Janssen S, Hendrickson R, Amir A, Vetter R, Hyde J, Doty MM,
541 Stillwell K, Benardini J, Kim JH, Allen EE, Venkateswaran K, Knight R. 2018.
542 KatharoSeq Enables High-Throughput Microbiome Analysis from Low-Biomass
543 Samples. *mSystems* 3.

544 11. Hunter JD. 2007. Matplotlib: A 2D graphics environment. *Comput Sci Eng* 9:90–95.

545 12. Lozupone C, Lladser ME, Knights D, Stombaugh J, Knight R. 2011. UniFrac: an
546 effective distance metric for microbial community comparison. *ISME J* 5:169.

547 13. McDonald D, Vázquez-Baeza Y, Koslicki D, McClelland J, Reeve N, Xu Z, Gonzalez
548 A, Knight R. 2018. Striped UniFrac: enabling microbiome analysis at
549 unprecedented scale. *Nat Methods* 2018 1511 15:847–848.

550 14. Falony G, Joossens M, Vieira-Silva S, Wang J, Darzi Y, Faust K, Kurilshikov A,
551 Bonder MJ, Valles-Colomer M, Vandeputte D, Tito RY, Chaffron S, Rymenans L,
552 Verspecht C, Sutter L De, Lima-Mendez G, D'hoë K, Jonckheere K, Homola D,
553 Garcia R, Tigchelaar EF, Eeckhaut L, Fu J, Henckaerts L, Zhernakova A,
554 Wijmenga C, Raes J. 2016. Population-level analysis of gut microbiome variation.
555 *Science* (80-) 352:560–564.

556

- 557 15. Vázquez-Baeza Y, Pirrung M, Gonzalez A, Knight R. 2013. EMPeror: a tool for
558 visualizing high-throughput microbial community data. *Gigascience* 2:16.
- 559 16. Morton JT, Marotz C, Washburne A, Silverman J, Zaramela LS, Edlund A, Zengler K,
560 Knight R. 2019. Establishing microbial composition measurement standards with
561 reference frames. *Nat Commun* 10:2719.
- 562 17. Kuhn M. 2008. Building Predictive Models in R Using the caret Package. *J Stat Softw*
563 28:1–26.
- 564 18. Keilwagen J, Grosse I, Grau J. 2014. Area under Precision-Recall Curves for
565 Weighted and Unweighted Data. *PLoS One* 9:e92209.
- 566 19. Wright MN, Ziegler A. 2017. ranger: A Fast Implementation of Random Forests for
567 High Dimensional Data in C++ and R. *J Stat Softw* 77:1–17.
- 568 20. Cantrell K, Fedarko MW, Rahman G, McDonald D, Yang Y, Zaw T, Gonzalez A,
569 Janssen S, Estaki M, Haiminen N, Beck KL, Zhu Q, Sayyari E, Morton JT,
570 Armstrong G, Tripathi A, Gauglitz JM, Marotz C, Matteson NL, Martino C, Sanders
571 JG, Carrieri AP, Song SJ, Swafford AD, Dorrestein PC, Andersen KG, Parida L,
572 Kim H-C, Vázquez-Baeza Y, Knight R. 2021. EMPress Enables Tree-Guided,
573 Interactive, and Exploratory Analyses of Multi-omic Data Sets. *mSystems* 6.
- 574 21. Protsyuk I, Melnik A V., Nothias LF, Rappez L, Phapale P, Aksenov AA, Bouslimani
575 A, Ryazanov S, Dorrestein PC, Alexandrov T. 2017. 3D molecular cartography
576 using LC–MS facilitated by Optimus and 'ili software. *Nat Protoc* 2017 131 13:134–
577 154.