

1 **Supplemental Information**

2 **Supplemental Methods**

3 **Oral Analysis Study Population:**

4 Three healthy adults contributed specimens to this study as participants in the healthy control
5 arm of the Oral Microbial Communities in States of Hyposalivation project (NIDCR Grant
6 Number: DE023113-01A1) (also see ref. 37). The study was approved by the Administrative
7 Panels on Human Subjects Research (Stanford IRB protocol #21586), and by the Human
8 Research Protection Program at UCSF (UCSF IRB protocol ##11-06283). Specimens were
9 collected between August and December 2014. All research subjects provided written informed
10 consent prior to specimen collection.

11 Exclusion criteria:

12 Study subjects were excluded if they had taken oral antibiotics or antifungals within six months
13 prior to screening date, if they used any medications other than oral contraceptives on a regular
14 basis, if they were under the age of 18, if they had evidence of active oral disease, or if they were
15 unable to adhere to the home- and clinic-collected sampling procedures.

16 **Sampling procedures:**

17 Individuals self-collected oral mucosal specimens every day for 30 days using Catch-All™
18 Sample Collection Swabs (Epicentre Biotechnologies, Madison, WI, USA). The seven sampling
19 sites included: anterior dorsal surface of the tongue, left and right buccal surfaces adjacent to the
20 first molars, upper and lower inner lips adjacent to the central incisors, the middle of the hard
21 palate, and the floor of the mouth adjacent to the attachment site of the lingual frenulum.
22 Subjects were instructed to sample each site by swabbing the area for 10 seconds while applying
23 moderate pressure and moving the swab in a circular motion, to collect specimens at the same
24 time of day, and to store the specimens in their home freezer immediately after collecting the
25 final specimen. Subjects were cautioned to prevent the swab from touching any surface other
26 than the intended sampling site, and to use a new, unopened swab if the swab touches a non-
27 target surface. Subjects also traveled to the Oral Medicine Faculty Clinic in the University of
28 California – San Francisco (UCSF) Dental Center for weekly appointments. There, a UCSF
29 dental clinician collected the same set of seven mucosal specimens following the same protocol.
30 Subjects were instructed to transport their home-collected specimens to the clinic appointments

31 in weekly batches using provided ice packs and freezer tote bags. Specimens were placed and
32 stored on dry ice immediately following collection. Home- and clinic-collected specimens were
33 then transported to the laboratory where they were stored at -80° C until time of DNA extraction.

34 **Detailed sample processing procedure**

35 Specimens were loaded into either a well of a 96-well plate or a 2mL Eppendorf tube, both
36 provided in the PowerSoil® kits described below. Swabs were either clipped into wells 1-2
37 centimeters above the foam head using a pair of sterile scissors, or bent until broken into
38 Eppendorf tubes. Whole genomic DNA was extracted from each specimen using the
39 PowerSoil®-HTP 96 well Soil DNA Isolation Kit (MO BIO Laboratories, Carlsbad, CA, USA)
40 or the PowerSoil® DNA Isolation Kit (2mL tube-based protocol) according to manufacturer's
41 instructions with two modifications: 1) inclusion of a 10-minute incubation step at 65° C
42 following addition of Solution C1 and 2) final elution in 125µL Solution C6. Genomic DNA was
43 PCR amplified using Golay error-correcting barcoded primers targeting the V4 hypervariable
44 region of the bacterial 16S rRNA gene [51]. The forward PCR primer (5' AAT GAT ACG GCG
45 ACC ACC GAG ATC TAC ACG CTN NNN NNN NNN NNT ATG GTA ATT GTG TGY
46 CAG CMG CCG CGG TAA 3') contains (from 5' to 3') an Illumina adapter, a unique 12-
47 nucleotide error-correcting Golay barcode (designated by 'N's), a forward primer pad, a 2-
48 nucleotide forward primer linker, and a broad-range bacterial primer, 515F. The reverse primer
49 (5' CAA GCA GAA GAC GGC ATA CGA GAT AGT CAG CCA GCC GGA CTA CNV GGG
50 TWT CTA AT 3') consists of an Illumina adapter, a reverse primer pad, a 2-nucleotide reverse
51 primer linker, and the broad-range bacterial primer 806R.

52 Each specimen was PCR-amplified according to the following 75-µL PCR reaction set-up and
53 PCR cycling conditions. Replication ranged from 2-4 75-µL replicate reactions per specimen
54 (see sample mapping file for replicates per sample) depending on amplification efficiency, as
55 assayed by presence of a DNA band on agarose gel.

56

57

58

59

60

| Component | 1 rxn (μL) |
|--|-------------------|
| Molecular biology grade water (Sigma) | 26.7 |
| 2.5X HotMasterMix (5 PRIME, Gaithersburg, MD, USA) | 30 |
| 100 uM 806R primer (IDT, Coralville, IA, USA) | 0.3 |
| 10 uM 515F primer (barcoded) | 3 |
| Template DNA | 15 |

| Thermal cycling: | Total (μL) |
|-------------------------|-------------------|
| 94°C 3 minutes | |
| 30 cycles of: | |
| 94°C 45 seconds | 75 |
| 52°C 60 seconds | |
| 72°C 90 seconds | |
| 72°C 10 minutes | |
| 4°C HOLD | |

61 Following confirmation of PCR amplicon product by agarose gel electrophoresis, amplicons
62 were purified using the UltraClean®-htp 96 Well PCR Clean-Up Kit (MO BIO Laboratories,
63 Carlsbad, CA, USA) according to the manufacturer's instructions with one exception: amplicons
64 were eluted in 125μL Elution Buffer. Specimens' purified DNA content was quantified using the
65 Quant-iT™ High-Sensitivity dsDNA Assay Kit (Invitrogen, Carlsbad, CA, USA), and 82 ng
66 DNA per specimen was pooled into a single tube. Specimens whose amplicon DNA
67 concentrations were insufficient to be pooled in equimolar concentrations were pooled using
68 their entire available volume, approximately 100-115μL. Pooled amplicons were then ethanol
69 precipitated and resuspended in 400μL pH 8.0 TE buffer (Sigma). Resuspended amplicons were
70 then size-selected using agarose gel electrophoresis, and recovered using the QIAquick Gel
71 Extraction Kit (Qiagen, Valencia, CA, USA) according to manufacturer's instructions. Two
72 aliquots of 50μL amplicon were separately sequenced on an Illumina MiSeq v3 flowcell by the
73 W.M. Keck Center for Comparative Functional Genomics at the University of Illinois, Urbana-
74 Champaign, USA.

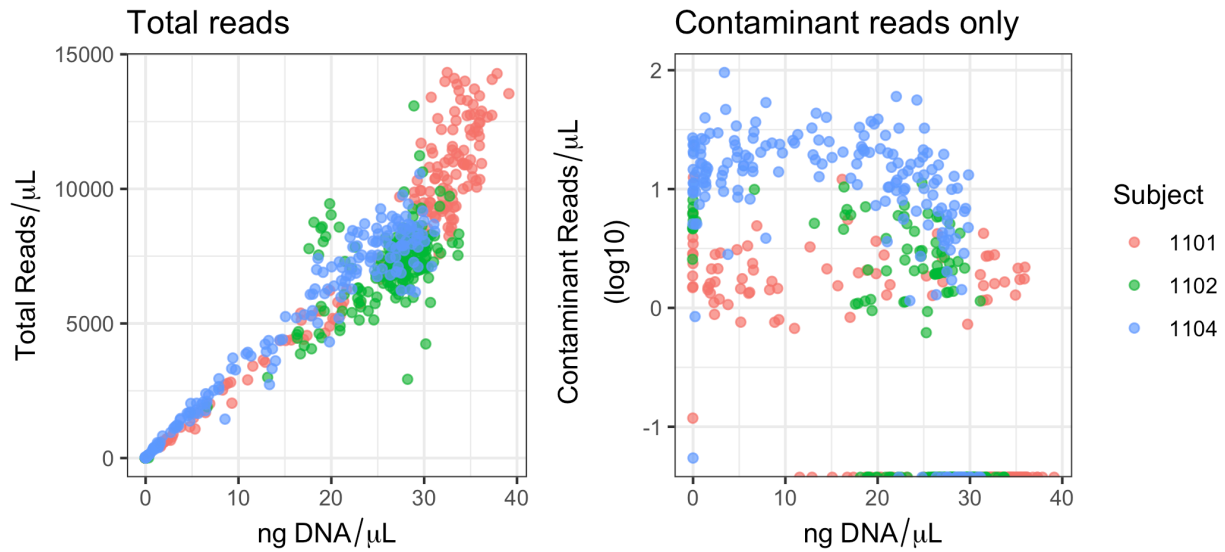
75 **Negative control description and processing**

76 Negative controls consisted of two types: reagent-only and blank-swab controls. Reagent-only
77 controls consisted of empty PowerSoil® wells to which all DNA extraction, PCR, and
78 sequencing reagents were added. One reagent-only control was included for each of six plates in
79 either well B2 or B11. Blank-swab controls consisted of unopened foam swabs from the same lot
80 as those used for sample collection. Control swabs were not distributed to study subjects prior to
81 use. Four blank-swab controls were loaded into each of six PowerSoil® plates when true
82 samples were loaded, into either wells G1,G2,H1,H2, or wells G11,G12,H11,H12.

83

84

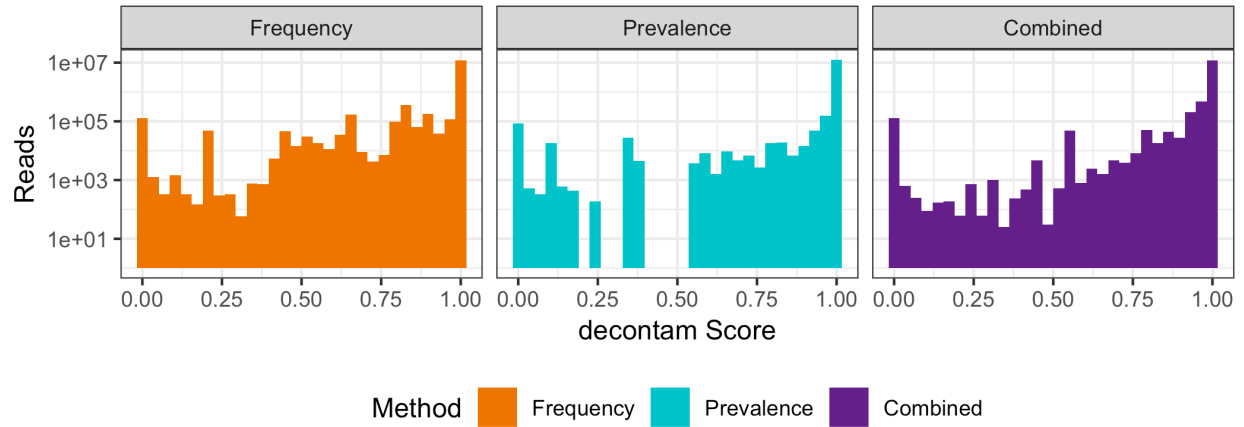
85 **Supplemental Figures**



86

87 **Figure S1. Concentration of contaminant features across oral samples.** ASVs from the oral dataset were
88 classified using decontam's prevalence method at the default threshold 0.1. Concentrations of total (left) or
89 contaminant (right) DNA (in reads per microliter of sample added to the DNA sequencing pool) are plotted against
90 the total post-PCR DNA concentration of the sample. Samples are colored by research subject. The total
91 concentration of contaminants is roughly even across samples, and independent of DNA concentration.

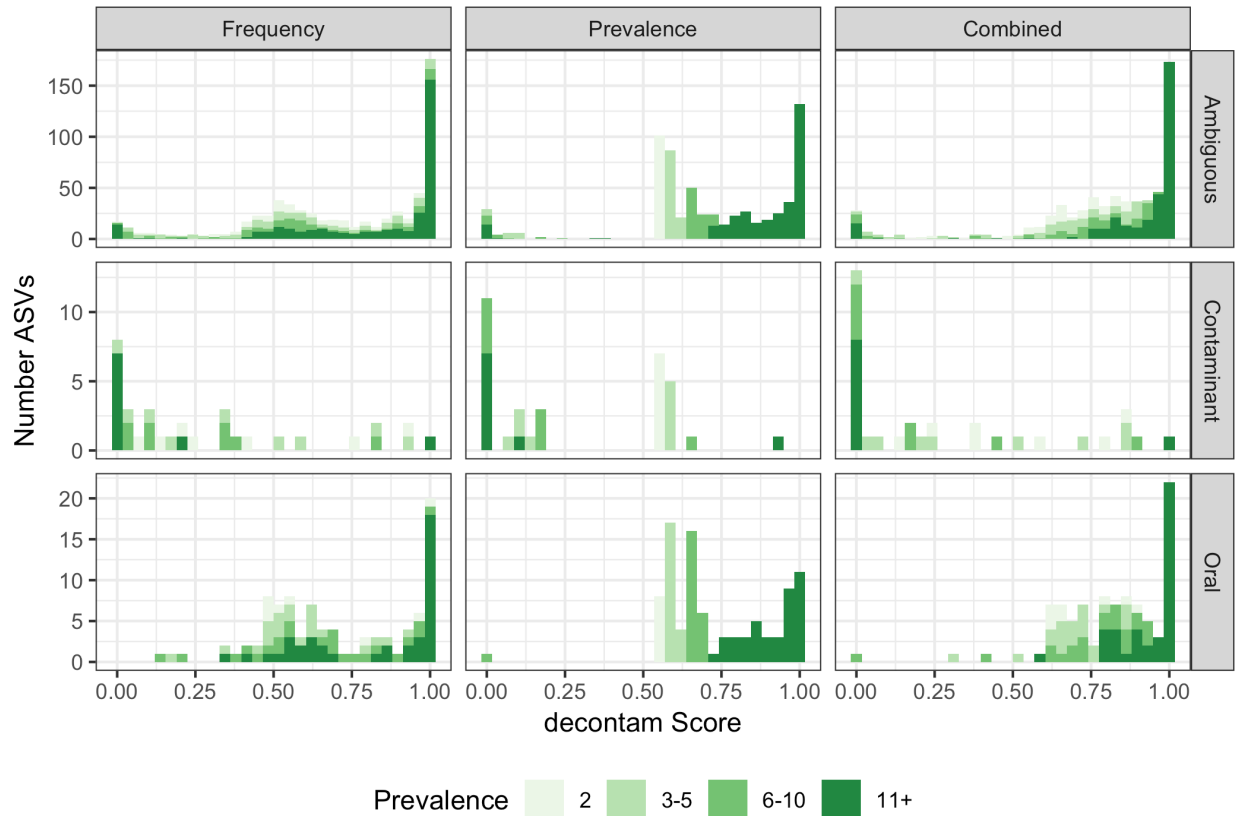
92



93

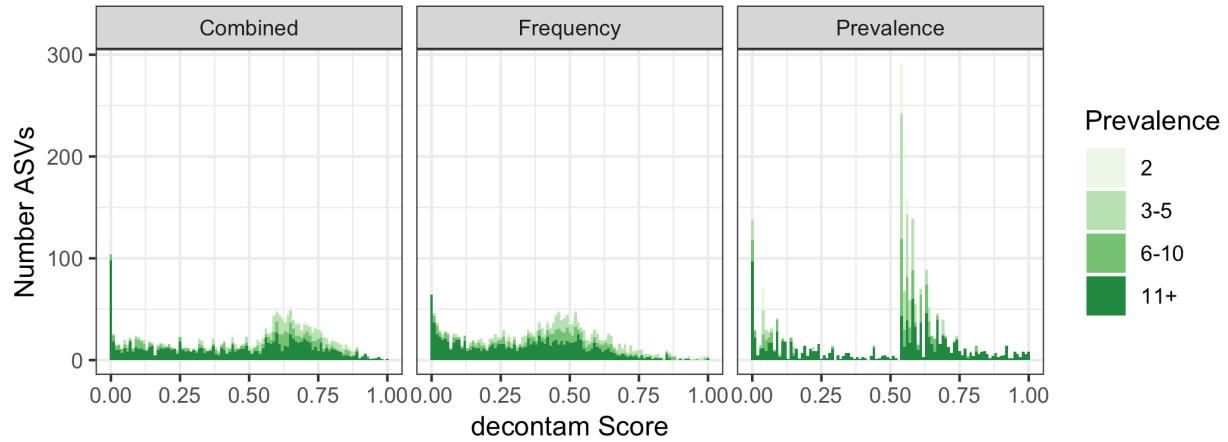
94 **Figure S2. Histogram of decontam scores in the oral mucosa dataset, weighted by ASV relative abundance.**

95 Scores for each amplicon sequence variant (ASV) present in two or more samples were computed by the frequency,
 96 prevalence and combined methods as implemented in the *isContaminant* function in the decontam R package. The
 97 histogram of scores for each method is shown, after weighting by the number of reads of each ASV. The y-axis is
 98 log-scaled. Most reads are assigned high scores, indicating non-contaminant origin.



99
 100 **Figure S3. Histogram of decontam scores in the oral mucosa dataset, stratified by reference-based**
 101 **classification.** Scores for each amplicon sequence variant (ASV) present in two or more samples were computed by
 102 the frequency, prevalence and combined methods as implemented in the *isContaminant* function in the decontam R
 103 package. The histogram of scores is shown, with color intensity depending on the number of samples (or prevalence)
 104 in which each ASV was present. Separate histograms are shown for ASVs classified as Ambiguous, Contaminant
 105 and Oral sequences based on comparison to curated reference databases (Methods). Reference-based contaminants
 106 are predominantly assigned low scores, reference-based oral ASVs are predominantly assigned high scores, while
 107 ambiguous ASVs consist of a mixture of high and low scores.

108



109

110

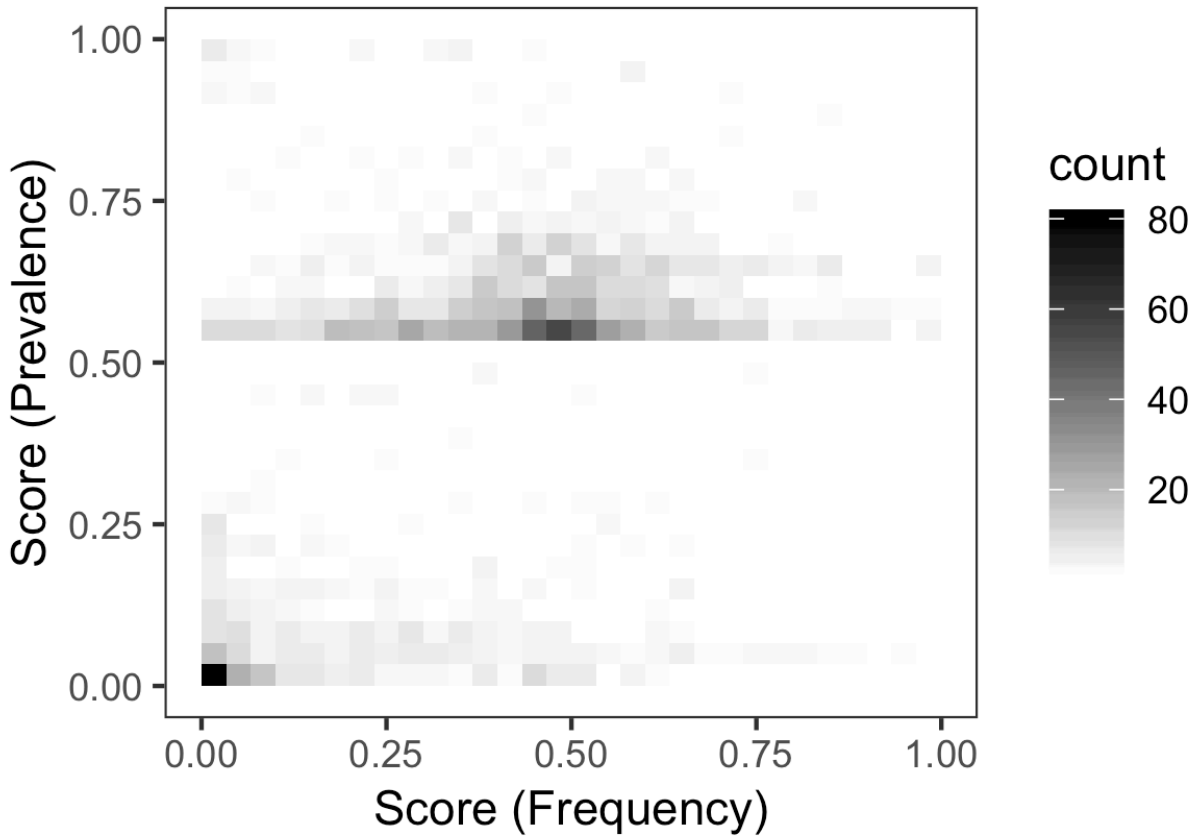
Figure S4. Histogram of decontam scores in the preterm birth dataset. Scores for each amplicon sequence variant (ASV) present in two or more samples were computed by the frequency, prevalence and combined methods as implemented in the *isContaminant* function in the decontam R package. The histogram of scores is shown, with color intensity depending on the number of samples (or prevalence) in which each ASV was present. A sharp peak in the score histogram is evidence at around $P = 0.1$.

113

114

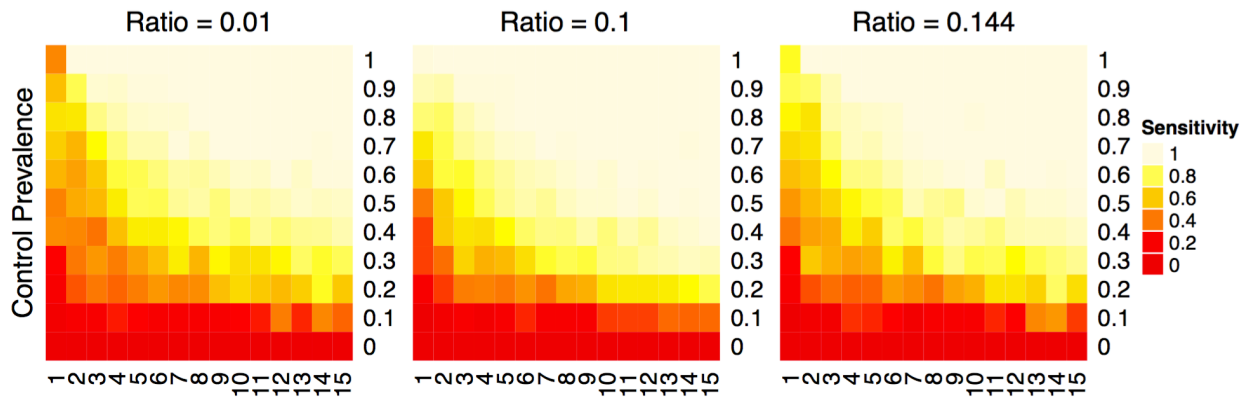
115

116



117
 118
 119
 120
 121
 122
 123
 124
 125

Figure S5. Density plot of scores assigned by the frequency and prevalence methods in the preterm birth dataset. Scores for each amplicon sequence variant (ASV) present in two or more samples were computed by the frequency and prevalence methods as implemented in the *isContaminant* function in the *decontam* R package. The density of joint score assignments is shown, with dark colors indicating higher numbers of ASVs assigned scores corresponding to that x-y position. Scores assigned by the frequency and prevalent methods are fairly consistent, especially at low values.



126 **Figure S6. Simulation analysis of prevalence method sensitivity.** Simulated prevalence distributions were
 127 generated based on a varying prevalence of contaminants in negative controls, a varying number of negative control
 128 samples, and three sample:control prevalence ratios of 0.01, 0.1, and 0.144. The mean prevalence ratio of
 129 contaminants in the oral data is 0.144. The number of controls is shown on the x-axis, and the prevalence of
 130 contaminants in negative controls on the y-axis. Sensitivity over 100 replicate simulations is shown by the color-
 131 scale. The number of true samples was fixed to 60 in these simulations, but its precise value has little impact on
 132 sensitivity as long as it is much greater than the number of controls. 5-6 negative control samples is sufficient to
 133 detect most contaminants under these scenarios.