

1 **Microbiota-based model improves the sensitivity for**
2 **detecting colonic lesions**

3

4 Nielson T. Baxter¹, Mack T. Ruffin IV², Mary A.M. Rogers³, and Patrick D. Schloss^{1*}

5

6 ¹Department of Microbiology and Immunology, University of Michigan, Ann Arbor, Michigan.

7 ²Department of Family Medicine, University of Michigan, Ann Arbor, Michigan.

8 ³Department of Internal Medicine, University of Michigan, Ann Arbor, Michigan.

9 *Correspondence to: pschloss@umich.edu

10

Colorectal cancer is the second leading cause of death among cancers in the United States¹. Although individuals diagnosed early have a greater than 90% chance of survival, more than one-third of individuals do not adhere to screening recommendations partly because the standard diagnostics, colonoscopy and sigmoidoscopy, are expensive and invasive¹⁻⁴. Thus, there is a great need to improve the sensitivity of non-invasive tests to detect early stage cancers and adenomas. Numerous studies have demonstrated a causal link between the formation of colonic lesions and the activity of the gut microbiota in tissue culture and animal models⁵⁻⁸. These findings have been complemented by studies in human populations identifying shifts in the composition of the gut microbiota associated with the progression of colorectal cancer⁹⁻¹³. These results suggest that the gut microbiota may represent a reservoir of biomarkers that would complement existing non-invasive methods such as the widely used fecal immunochemical test (FIT). Using stool samples from 490 patients we developed a cross-validated random forest classification model that detects colonic lesions using the relative abundance of gut microbiota and the concentration of hemoglobin in stool. The microbiota-based random forest model detected 95.0% of cancers and 57.1% of adenomas while FIT alone detected 75.0% and 15.7%, respectively. Of the colonic lesions missed by FIT, the model detected 80.0% of cancers and 49.1% of adenomas. These findings demonstrate the potential for microbiota analysis to complement existing screening methods to improve detection of colonic lesions. With a high sensitivity and low rate of false negatives, our model could be used to accurately identify asymptomatic individuals with preclinical disease and, thus, save lives.

Colorectal cancer mortality has steadily declined in recent decades, due in large part to increased screening¹. Yet current screening tests, the fecal immunochemical test (FIT) and the multitarget DNA test, have a sensitivity of 7.6% and 17.2%, respectively, for detecting non-advanced adenoma –

35 just the type of early lesion that screening is meant to identify¹⁴. Although structural exams
36 including colonoscopy and sigmoidoscopy are able to detect both adenomas and carcinomas, the
37 high cost and invasive nature are barriers for many people. Fear, discomfort, and embarrassment
38 are among the most cited reasons patients choose to forego CRC screening⁴. Likewise the large
39 disparity in screening rates between those with and without health insurance highlights the need
40 for inexpensive screening methods¹⁻³. Unfortunately cheaper, less invasive stool-based tests like
41 guaic fecal occult blood test and FIT are unable to reliably detect adenomas¹⁵. The newly introduced
42 stool DNA panel has improved accuracy compared to FIT, but is still limited in its ability to
43 accurately detect adenomas¹⁴. Thus there is need for novel screening methods that are inexpensive
44 and capable of detecting both cancer and adenomas.

45 The gut microbiota, the collection of microorganisms that inhabit the gastrointestinal tract, are one
46 potential source of biomarkers for detecting colonic lesions. Numerous studies have observed
47 alterations in the gut bacterial communities of patients with CRC⁹⁻¹³. Experiments in animal models
48 have demonstrated that such alterations have the potential to accelerate tumorigenesis⁵.
49 Furthermore, several members of the gut microbiota have been shown to potentiate both the
50 development and progression of CRC by a variety of mechanisms⁶⁻⁸. Although each of these
51 organisms may play a role in certain cases of CRC, none of them is present in every case. Therefore
52 we postulate that no one organism is an effective biomarker on its own and that focusing on a single
53 bacterial population excludes the potential that the microbial etiology of the disease is actually
54 polymicrobial.

55 We and others have shown that statistical models that take into account the abundances of multiple
56 bacterial species can be used to distinguish healthy individuals from those with CRC^{16,17}. In the
57 present study we expanded upon those findings by demonstrating the potential for microbiota

58 analysis to complement FIT for improved detection of colonic lesions, including adenomas. We
59 utilized the random forest algorithm, which is a decision tree-based machine learning algorithm for
60 classification that accounts for non-linear data and interactions among features and includes an
61 internal cross-validation to prevent overfitting¹⁸. By incorporating data on hemoglobin and
62 bacterial abundances into a single model (labeled the Multitarget Microbiota Test or MMT), we
63 were able to improve the sensitivity for adenomas and cancer compared to FIT alone.

64 We characterized the bacterial communities of stool samples from 490 patients using 16S rRNA
65 gene sequencing. Among these patients, 120 had CRC, 109 had advanced adenomas, 89 had non-
66 advanced adenomas, and 172 had no colonic lesions. We also tested each sample for the
67 concentration of hemoglobin using FIT. With these data we developed a random forest model that
68 incorporated the microbiota and FIT data and would differentiate normal individuals from those
69 with any type of colonic lesion (i.e. adenoma or carcinoma). We determined the optimal model
70 using the AUC-RF algorithm for maximizing the area under the curve (AUC) of the receiver
71 operating characteristic (ROC) curve for the MMT¹⁹. The optimal model combining hemoglobin
72 results and the microbiota used 23 bacterial populations, or operational taxonomic units (OTUs)
73 (Extended Data Fig. 1). Of those OTUs, 16 were members of the Firmicutes phylum, including 3 from
74 the Ruminococcaceae family and 10 from the Lachnospiraceae family, the predominant producers
75 of butyrate in the gut²⁰ (Extended Data Fig. 2). Three OTUs were associated with the genus
76 *Bacteroides*. The remaining OTUs were associated with *Porphyromonas*, *Parabacteroides*, *Collinsella*,
77 and Enterobacteriaceae. The OTU associated with *Porphyromonas* was most closely related to
78 *Porphyromonas asaccharolytica*, which has been previously shown to be predictive of CRC^{16,21}. Like
79 other studies^{13,16} we also observed an OTU associated with *Fusobacterium nucleatum* that was
80 enriched in cancer samples, however its relative abundance did not add sufficient information to be
81 included in the model. Interestingly the majority of OTUs used in the model, especially the

82 Lachnospiraceae, were enriched in normal patients, suggesting that a loss of beneficial organisms in
83 addition to the emergence of pathogens may be indicative of CRC development.

84 To determine whether microbiota sequence data could be used to complement FIT, we compared
85 the performance of the MMT to FIT. The AUC for the MMT (AUC=0.755) was significantly higher
86 than FIT (AUC=0.639) for distinguishing adenoma from normal ($p<0.001$) or all lesions from
87 normal (FIT AUC=0.749, MMT AUC=0.829, $p<0.001$), but not cancer from normal (FIT AUC=0.929,
88 MMT AUC=0.952, $p=0.091$) (Fig. 1A). To generate a categorical prediction from the MMT, we
89 determined that the optimal threshold for the models's probability was 0.622 using Youden's J
90 statistic²². Samples scoring above this cutoff were classified as lesions, and those below the cutoff
91 were classified as normal. We then compared the sensitivity and specificity of the MMT to those of
92 FIT using the manufacturer recommended threshold of 100 ng/ml of hemoglobin. At these cutoffs
93 the MMT detected 95.0% of cancers and 57.1% of adenomas compared to 75.0% and 15.7% for FIT
94 (Table 1, Fig. 1B). When adenomas and cancers were pooled together, the MMT detected 71.4% of
95 lesions, while FIT only detected 38.1%. The MMT significantly improved sensitivity for both
96 advanced and non-advanced adenomas as well as multiple stages of cancer (Fig. 2). The increased
97 sensitivity of the MMT was accompanied by a decrease in specificity (83.7%) compared to FIT
98 (97.1%).

99 To better understand the relationship between the MMT and FIT, we compared the results of the
100 two tests for each sample (Fig. 3). All samples that tested positive by FIT also tested positive by the
101 MMT, indicating that the MMT did not miss any of the lesions that FIT was able to detect. However
102 the MMT was able to detect 80% of cancers and 49.1% of adenomas that FIT had failed to detect,
103 while maintaining a specificity of 86.2% (Extended Data Fig. 3). This result demonstrated that
104 incorporation of data from a subject's microbiota complemented FIT to improve its sensitivity.

105 The purpose of screening is to identify asymptomatic individuals with early stage disease (i.e., true
106 positives). Therefore, we estimated the number of true positives captured through FIT and MMT in
107 the recommended screening population in the United States (adults ages 50-74 years). The
108 prevalence of lesions in an average-risk population was obtained through a previously published
109 meta-analysis²³. Tests were utilized in series so that FIT, with a higher specificity (fewer false
110 positives), was applied first to minimize unnecessary diagnostic testing. MMT, with a higher
111 sensitivity (fewer false negatives), was then used to capture additional true positives in those with
112 negative FIT results (Extended Data Table 1). MMT was able to identify a large proportion of true
113 positives among individuals with a negative FIT result (55.1% for cancer, 72.0% for advanced
114 adenoma, 82.5% for non-advanced adenoma).

115 Previous studies have identified differences in diagnostic test performance for certain demographic
116 groups or for people taking certain medications²⁴⁻²⁶. Therefore we tested whether the MMT
117 performance differed between patient populations. We found no difference in model performance
118 according to age, BMI, NSAID usage, diabetes, smoking, or previous history of polyps (all $p>0.05$).
119 However the model was significantly better at differentiating normal from lesion for females than
120 for males ($p=0.016$; Extended Data Fig. 4). For females the model detected 73.5% of lesions with a
121 specificity of 89.2%. For males the model detected 69.9% of lesions with a specificity of 73.8%. This
122 difference was more pronounced for adenomas. The MMT detected 62.5% of adenomas in females
123 and 53.4% in males. Despite performing more poorly overall for males, the MMT did have a higher
124 sensitivity for cancer among males (98.5%) than females (90.4%). The difference in performance
125 between males and females seems to be due to differences in FIT results rather than differences in
126 the microbiome. After correcting for diagnosis, there was a significant effect of sex on FIT result
127 ($p=0.0057$, two-way ANOVA), but not on the overall structure of the microbiome($p=0.063$,
128 PERMANOVA).

129 It was recently shown that when FIT was combined with host-associated DNA biomarkers the
130 ability to detect adenomas and carcinomas was significantly improved over FIT alone¹⁴. The
131 sensitivity of the host-associated DNA screen was 92.3% for CRC and 42.4% for adenomas, which
132 are both slightly lower than what we observed with our MMT. Regardless of the relative
133 performance, such results support the assertion that because of the large interpersonal variation in
134 markers for adenomas and carcinomas, it is necessary to employ a panel of biomarkers and to use a
135 model that integrates the biomarkers. The accuracy of our model may be further improved by
136 incorporating additional biomarkers such as the host-associated biomarkers or those targeting
137 specific genes involved in the underlying mechanism of tumorigenesis such as toxins^{7,8,17}. More
138 generally, predictive and diagnostic models for other diseases with a microbial etiology may benefit
139 from a similar approach. For example, we recently demonstrated the ability to detect *Clostridium*
140 *difficile* infection based on the composition of the microbiota²⁷. Such models are likely to be useful
141 as microbiota sequencing gains traction as a tool for characterizing health.

142 Our findings demonstrate the potential for combining the analysis of a patient's microbiota with
143 conventional stool-based tests to improve CRC detection. Using the random forest algorithm it was
144 possible to interpret FIT results in the context of the microbiota. The MMT had significantly higher
145 sensitivity for lesions at almost all stages of tumorigenesis. Moreover the model detected the
146 majority of lesions that FIT was unable to detect. The shortcoming of the MMT is its lower
147 specificity but, by conducting the FIT and MMT in series, it is possible to maximize the number of
148 correctly identified individuals with preclinical lesions. The potential value of the MMT is its higher
149 sensitivity which, at its core, is the purpose of preventive screening – finding lesions earlier so that
150 cancer would be avoided.

151 **Methods Summary.** Fecal samples were collected from 490 subjects in 4 locations: Toronto
152 (Ontario, Canada), Boston (Massachusetts, USA), Houston (Texas, USA), and Ann Arbor (Michigan,
153 USA). Patient diagnoses were determined by colonoscopy and subsequent histopathological
154 examination of any biopsies taken. FIT was performed using OC FIT-CHEK sampling bottles and
155 processed using an OC-Auto Micro 80 automated system (Polymedco Inc.). The V4 region of the
156 bacterial 16S rRNA gene was amplified using custom barcoded primers, sequenced using an
157 Illumina MiSeq sequencer, and analyzed as described previously²⁸. A data analysis pipeline and all
158 necessary scripts to generate this paper are available at
159 github.com/SchlossLab/Baxter_gline007Modeling_2015. The sequence data are available in the
160 Sequence Read Archive under accession number SRP062005.

161 **Methods (online only)**

162 ***Study Design/Patient sampling.*** Eligible patients for this study were at least 18 years old, willing
163 to sign informed consent, able to tolerate removal of 58 ml of blood, and willing to collect a stool
164 sample. Patient age at the time of enrollment ranged from 29 to 89 with a median of 60. All patients
165 were asymptomatic and were excluded if they had undergone surgery, radiation, or chemotherapy
166 for current CRC prior to baseline samples or had inflammatory bowel disease, known hereditary
167 non-polyposis CRC, or familial adenomatous polyposis. Colonoscopies were performed and fecal
168 samples were collected from subjects in 4 locations: Toronto (Ontario, Canada), Boston
169 (Massachusetts, USA), Houston (Texas, USA), and Ann Arbor (Michigan, USA). Patient diagnoses
170 were determined by colonoscopic examination and histopathological review of any biopsies taken.
171 Patients with an adenoma greater than 1cm, more than three adenomas of any size, or an adenoma
172 with villous histology were classified as advanced adenoma. Whole evacuated stool was collected
173 from each patient either prior to colonoscopy preparation or 1-2 weeks after colonoscopy. This has
174 been shown to be sufficient time for the microbiota to recover from colonoscopy preparation²⁹.

175 Stool samples were packed in ice, shipped to a processing center via next day delivery and stored at
176 -80°C. This study was approved by the University of Michigan Institutional Review Board and all
177 subjects provided informed consent.

178 ***Fecal Immunochemical Tests.*** Fecal material for FIT was collected from frozen stool aliquots using
179 OC FIT-CHEK sampling bottles (Polymedco Inc.) and processed using an OC-Auto Micro 80
180 automated system (Polymedco Inc.). Hemoglobin concentrations were used for generating ROC
181 curves for FIT and for building the MMT.

182 ***16S rRNA Sequencing.*** DNA was extracted from roughly 50 mg of fecal material from each subject
183 using the PowerSoil-htp 96 Well Soil DNA isolation kit (MO BIO Laboratories) and an epMotion
184 5075 automated pipetting system (Eppendorf). The V4 region of the bacterial 16S rRNA gene was
185 amplified using custom barcoded primers and sequenced as described previously using an Illumina
186 MiSeq sequencer²⁸. The 490 samples were divided into three sequencing runs to increase the per
187 sample sequencing depth. Although the same percentage of samples from the three groups were
188 represented on each sequencing run, samples were randomly assigned to the sequencing runs to
189 avoid confounding our analysis based on diagnosis or demographics.

190 ***Sequence Curation.*** The 16S rRNA gene sequences were curated using the mothur software
191 package, as described previously²⁸. Briefly, paired-end reads were merged into contigs, screened for
192 quality, aligned to SILVA 16S rRNA sequence database, and screened for chimeras. Curated
193 sequences were clustered in to operational taxonomic units (OTUs) using a 97% similarity cutoff
194 with the average neighbor clustering algorithm. The number of sequences in each sample was
195 rarefied to 10,000 per sample to minimize the effects of uneven sampling.

196 **Statistical Methods.** All statistical analyses were performed using R. Random Forest models were
197 generated using the AUCRF package¹⁹. The AUC of ROC curves was compared using the method
198 described by DeLong et al.³⁰. The optimal cutoff for the MMT was determined using Youden's *J*
199 statistic as implemented in the pROC package in R²². The sensitivities of FIT and the MMT were
200 compared using McNemar's chi-squared test. To control for diagnosis while testing the effects of sex
201 on the microbiome we used PERMANOVA as implemented in the adonis function in the vegan
202 package.

203 **Data Availability.** Raw fastq files and a MIMARKS file are available through the NCBI Sequence
204 Read Archive [SRP062005]. A data analysis pipeline and all necessary scripts are available at
205 github.com/SchlossLab/Baxter_gln007Modeling_2015.

206 **Literature cited**

- 207 1. Siegel, R., DeSantis, C. & Jemal, A. Colorectal cancer statistics, 2014. *CA: A Cancer Journal for*
208 *Clinicians* **64**, 104–117 (2014).
- 209 2. Disease Control, C. for, (CDC, P. & others. Vital signs: Colorectal cancer screening test use–United
210 states, 2012. *MMWR. Morbidity and Mortality Weekly Report* **62**, 881 (2013).
- 211 3. Hsia, J. *et al.* The importance of health insurance as a determinant of cancer screening: evidence
212 from the Women's Health Initiative. *Preventive Medicine* **31**, 261–270 (2000).
- 213 4. Jones, R. M., Devers, K. J., Kuzel, A. J. & Woolf, S. H. Patient-reported barriers to colorectal cancer
214 screening: a mixed-methods analysis. *American Journal of Preventive Medicine* **38**, 508–516 (2010).
- 215 5. Zackular, J. P. *et al.* The gut microbiome modulates colon tumorigenesis. *MBio* **4**, e00692–13
216 (2013).

- 217 6. Kostic, A. D. *et al.* Fusobacterium nucleatum potentiates intestinal tumorigenesis and modulates
218 the tumor-immune microenvironment. *Cell Host & Microbe* **14**, 207–215 (2013).
- 219 7. Wu, S. *et al.* A human colonic commensal promotes colon tumorigenesis via activation of T helper
220 type 17 T cell responses. *Nature Medicine* **15**, 1016–1022 (2009).
- 221 8. Arthur, J. C. *et al.* Intestinal inflammation targets cancer-inducing activity of the microbiota.
222 *Science* **338**, 120–123 (2012).
- 223 9. Wang, T. *et al.* Structural segregation of gut microbiota between colorectal cancer patients and
224 healthy volunteers. *The ISME Journal* **6**, 320–329 (2012).
- 225 10. Chen, H.-M. *et al.* Decreased dietary fiber intake and structural alteration of gut microbiota in
226 patients with advanced colorectal adenoma. *The American Journal of Clinical Nutrition* **97**, 1044–
227 1052 (2013).
- 228 11. Chen, W., Liu, F., Ling, Z., Tong, X. & Xiang, C. Human intestinal lumen and mucosa-associated
229 microbiota in patients with colorectal cancer. *PloS One* **7**, e39743 (2012).
- 230 12. Shen, X. J. *et al.* Molecular characterization of mucosal adherent bacteria and associations with
231 colorectal adenomas. *Gut Microbes* **1**, 138–147 (2010).
- 232 13. Kostic, A. D. *et al.* Genomic analysis identifies association of Fusobacterium with colorectal
233 carcinoma. *Genome Research* **22**, 292–298 (2012).
- 234 14. Imperiale, T. F. *et al.* Multitarget stool DNA testing for colorectal-cancer screening. *New England*
235 *Journal of Medicine* **370**, 1287–1297 (2014).
- 236 15. Hundt, S., Haug, U. & Brenner, H. Comparative evaluation of immunochemical fecal occult blood
237 tests for colorectal adenoma detection. *Annals of Internal Medicine* **150**, 162–169 (2009).

- 238 16. Zackular, J. P., Rogers, M. A., Ruffin, M. T. & Schloss, P. D. The human gut microbiome as a
239 screening tool for colorectal cancer. *Cancer Prevention Research* **7**, 1112–1121 (2014).
- 240 17. Zeller, G. *et al.* Potential of fecal microbiota for early-stage detection of colorectal cancer.
241 *Molecular Systems Biology* **10**, 766 (2014).
- 242 18. Liaw, A. & Wiener, M. Classification and regression by randomForest. *R News* **2**, 18–22 (2002).
- 243 19. Calle, M. L., Urrea, V., Boulesteix, A.-L. & Malats, N. AUC-RF: A new strategy for genomic profiling
244 with random forest. *Human Heredity* **72**, 121–132 (2011).
- 245 20. Pryde, S. E., Duncan, S. H., Hold, G. L., Stewart, C. S. & Flint, H. J. The microbiology of butyrate
246 formation in the human colon. *FEMS Microbiology Letters* **217**, 133–139 (2002).
- 247 21. Rex, D. K. *et al.* American College of Gastroenterology guidelines for colorectal cancer screening
248 2008. *The American Journal of Gastroenterology* **104**, 739–750 (2009).
- 249 22. Youden, W. J. Index for rating diagnostic tests. *Cancer* **3**, 32–35 (1950).
- 250 23. Heitman, S. J. *et al.* Prevalence of adenomas and colorectal cancer in average risk individuals: a
251 systematic review and meta-analysis. *Clinical Gastroenterology and Hepatology* **7**, 1272–1278
252 (2009).
- 253 24. Symonds, E. L. *et al.* Factors affecting faecal immunochemical test positive rates: demographic,
254 pathological, behavioural and environmental variables. *Journal of Medical Screening*
255 0969141315584783 (2015).
- 256 25. Kapidzic, A. *et al.* Gender differences in fecal immunochemical test performance for early
257 detection of colorectal neoplasia. *Clinical Gastroenterology and Hepatology* (2015).

- 258 26. Levi, Z. *et al.* Sensitivity, but not specificity, of a quantitative immunochemical fecal occult blood
259 test for neoplasia is slightly increased by the use of low-dose aspirin, NSAIDs, and anticoagulants.
260 *The American Journal of Gastroenterology* **104**, 933–938 (2009).
- 261 27. Schubert, A. M., Sinani, H. & Schloss, P. D. Antibiotic-Induced Alterations of the Murine Gut
262 Microbiota and Subsequent Effects on Colonization Resistance against *Clostridium difficile*. *MBio* **6**,
263 (2015).
- 264 28. Kozich, J. J., Westcott, S. L., Baxter, N. T., Highlander, S. K. & Schloss, P. D. Development of a dual-
265 index sequencing strategy and curation pipeline for analyzing amplicon sequence data on the MiSeq
266 Illumina sequencing platform. *Applied and Environmental Microbiology* **79**, 5112–5120 (2013).
- 267 29. O’Brien, C. L., Allison, G. E., Grimpen, F. & Pavli, P. Impact of Colonoscopy Bowel Preparation on
268 Intestinal Microbiota. *PLoS ONE* **8**, e62815 (2013).
- 269 30. DeLong, E. R., DeLong, D. M. & Clarke-Pearson, D. L. Comparing the areas under two or more
270 correlated receiver operating characteristic curves: a nonparametric approach. *Biometrics* 837–845
271 (1988).

272

273 **Author Contributions**

274 All authors were involved in the conception and design of the study. NTB processed samples and
275 analyzed the data. All authors interpreted the data. NTB and PDS wrote the manuscript. All authors
276 reviewed and revised the manuscript.

277 **Author Information**

278 The authors declare no competing financial interests.

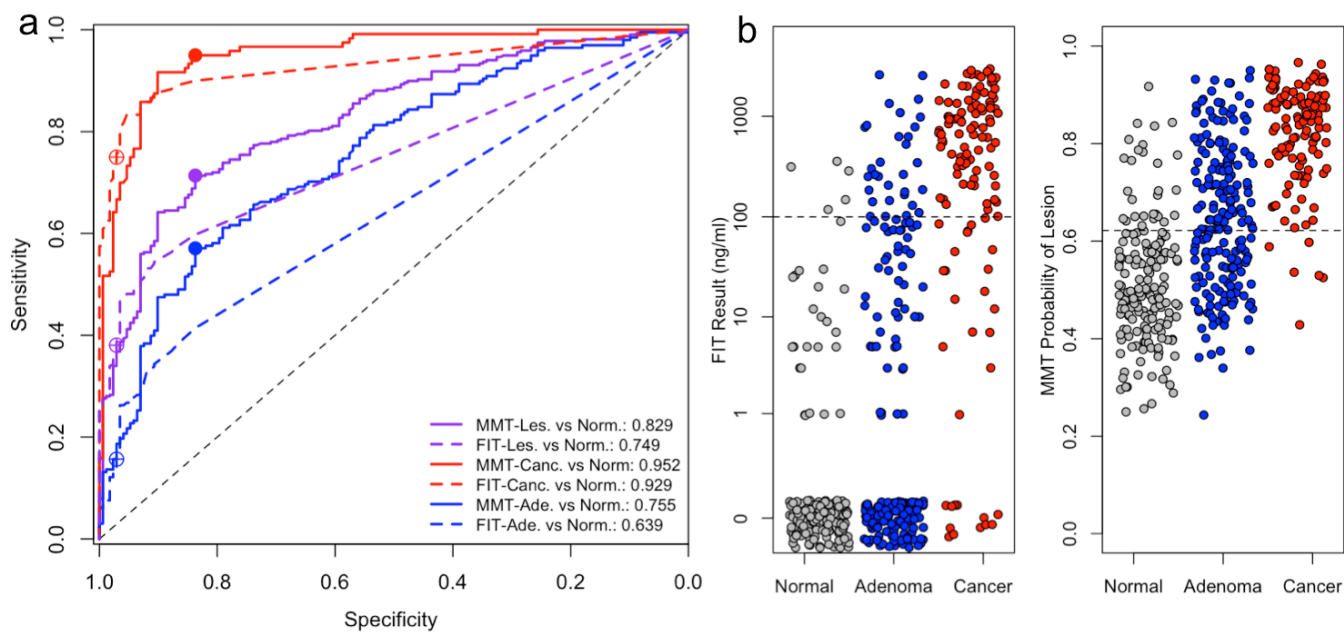
280 **Figures**

Diagnosis		FIT		Multitarget Microbiota Test	
		True Positives	Sensitivity (95% CI)	True Positives	Sensitivity (95% CI)
Cancer	n=120	90	75.0 (67.5-82.5)	114	95.0 (90.8-98.3)
Advanced Adenoma	n=109	21	19.3 (11.9-27.5)	64	58.7 (49.5-67.9)
Non Advanced Adenoma	n=89	10	11.2 (5.62-18)	49	55.1 (43.8-65.2)
Any Lesions	n=318	121	38.1 (33-43.4)	227	71.4 (66.4-76.4)
		True Negatives	Specificity (95% CI)	True Negatives	Specificity (95% CI)
Normal	n=172	167	97.1 (94.2-99.4)	144	83.7 (77.9-89)

281

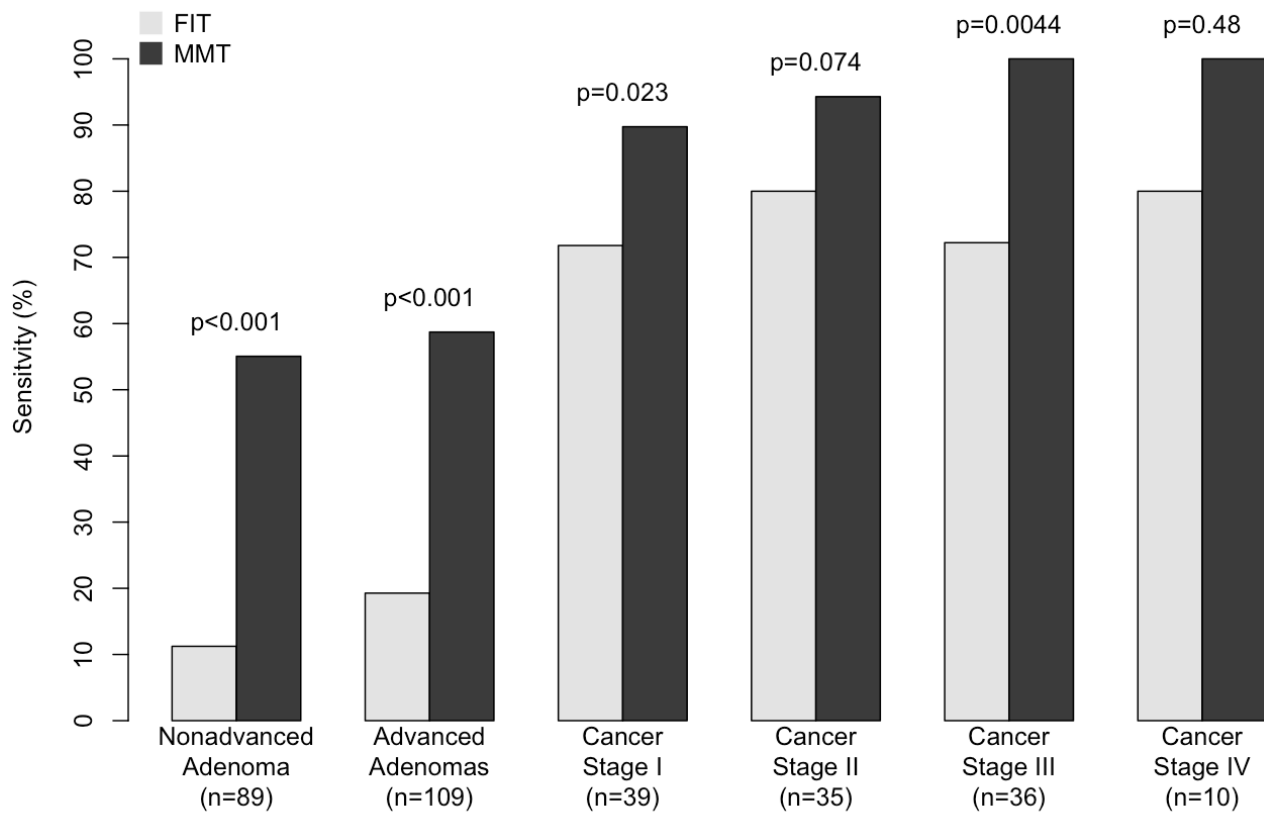
282 **Table 1.** Table of sensitivities and specificities for FIT and MMT. The 95% confidence intervals
 283 were computed with 2000 stratified bootstrap replicates.

284



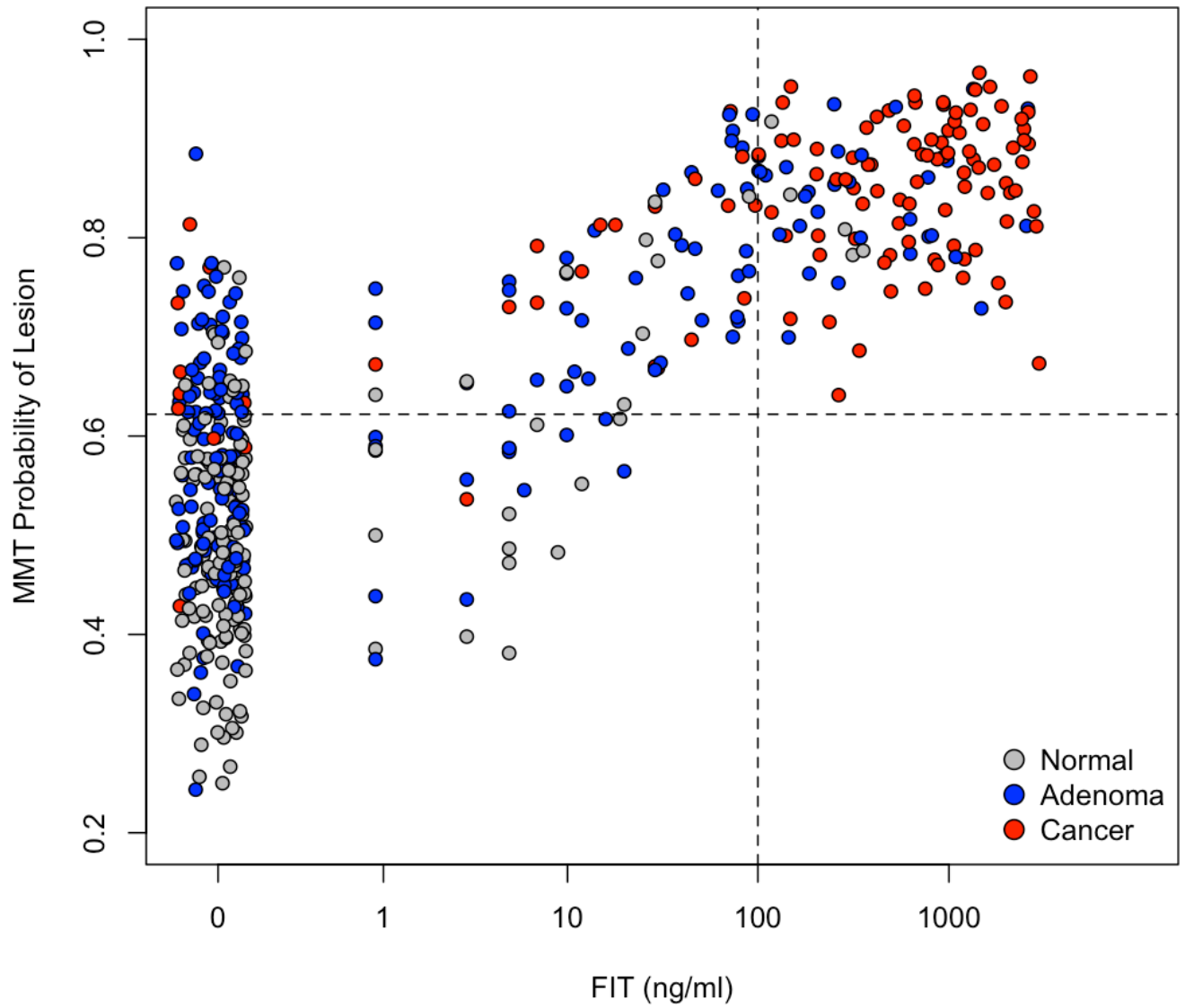
286

287 **Figure 1.** (a) ROC Curves for the MMT (solid lines) and FIT (dashed lines) for distinguishing normal
288 from any lesion (purple), normal from cancer (red) and normal from adenoma (blue). Filled dots
289 show the sensitivity and specificity of the MMT at the optimal cutoff (0.622). Open dots show the
290 sensitivity and specificity of FIT at the 100ng/ml cutoff. (b) Strip charts showing the results for FIT
291 and the MMT. Dashed lines show the cutoff for each test. Points with a FIT result of 0 are jittered to
292 improve visibility.



293

294 **Figure 2.** Barplot of sensitivities for FIT and MMT for each stage of tumor development. P-values
 295 based on McNemar's chi-squared test.



296
 297 **Figure 3.** Scatter plot of the results of FIT and MMT for each sample. Dashed lines show the cutoff
 298 for each test. Points with a FIT result of 0 are jittered to improve visibility.

299

300

301 **Extended Data Figures**

302

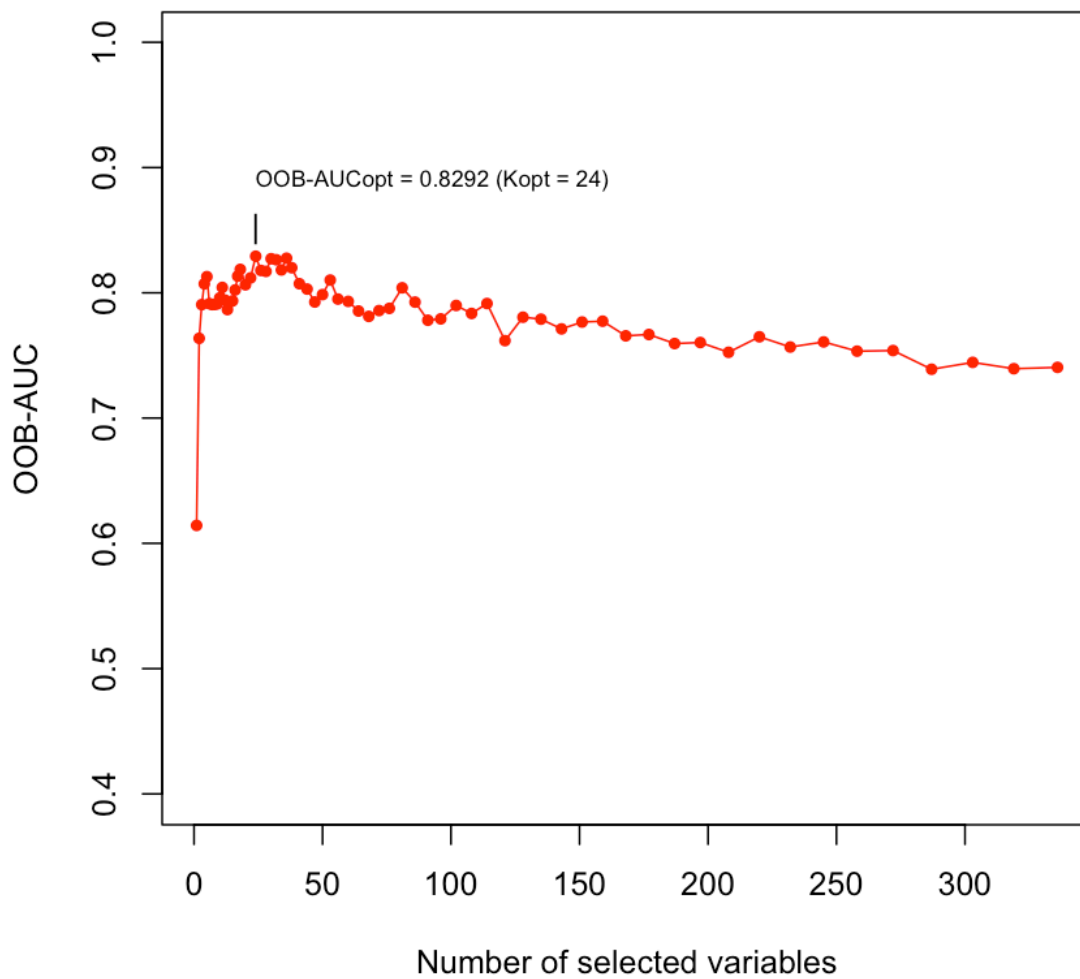
Number and Percentage of True Positives Identified through FIT and MMT

Condition	Prevalence	Number of Persons, ages 50-75 years, with Condition	True Positives identified through FIT	True Positives identified through additional MMT	Percentage of True Positives identified through additional MMT
Cancer	0.3%	241,483	181,112	222,260	55.1%
Advanced Adenoma	5.7%	4,588,174	885,518	2,276,159	72.0%
Non-advanced Adenoma	17.7%	14,247,488	1,595,719	7,507,376	82.5%

*Number of persons in the United States in 2010, 50-75 years of age, was 80,494,283.

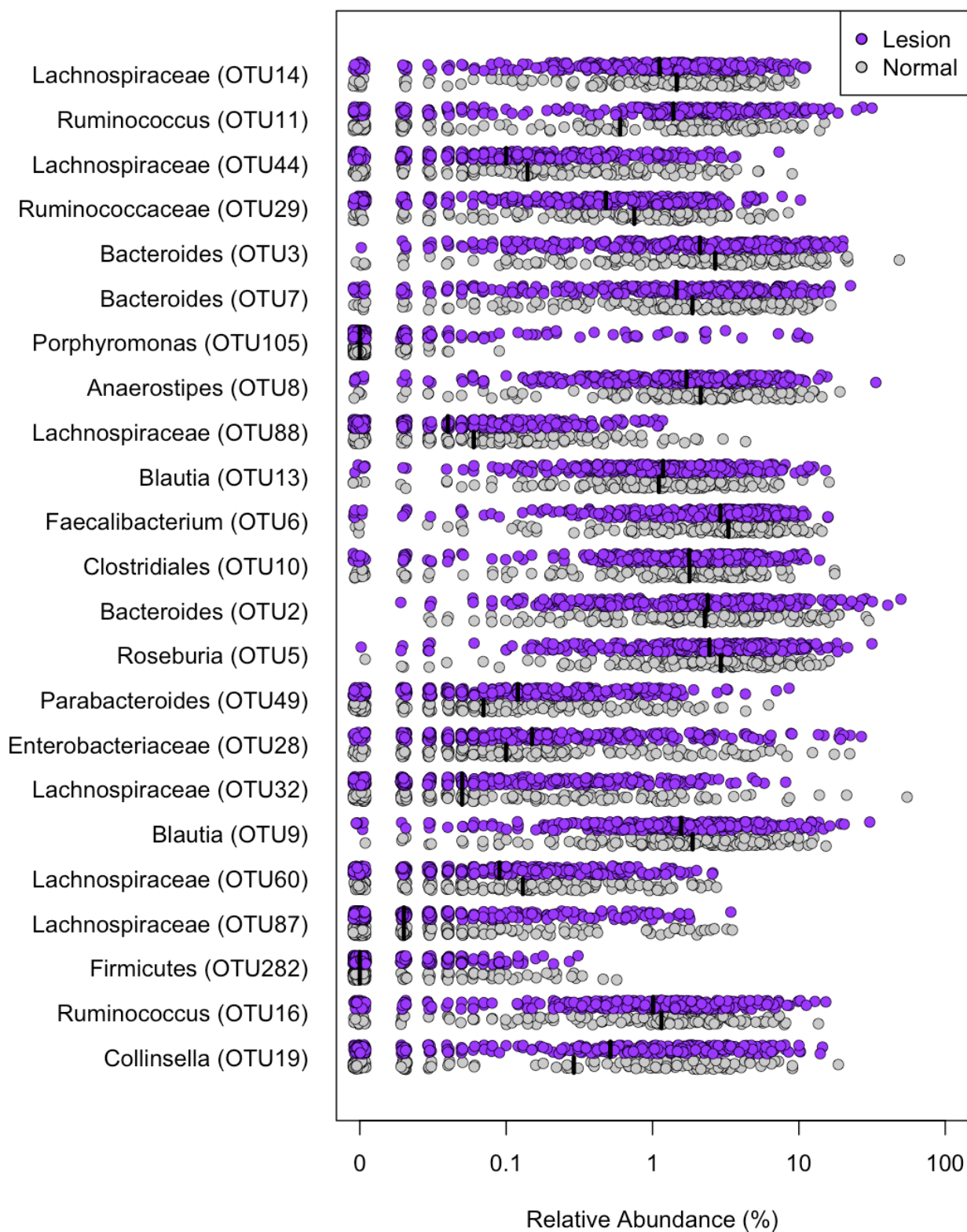
303

304 **Extended Data Table 1.** Number and proportion of true positives identified through FIT and MMT
305 in the United States in adults 50-75 years of age, based on published estimates of CRC prevalence.
306 Far right column shows percentage of true positives identified among individuals with a negative
307 FIT result.



308

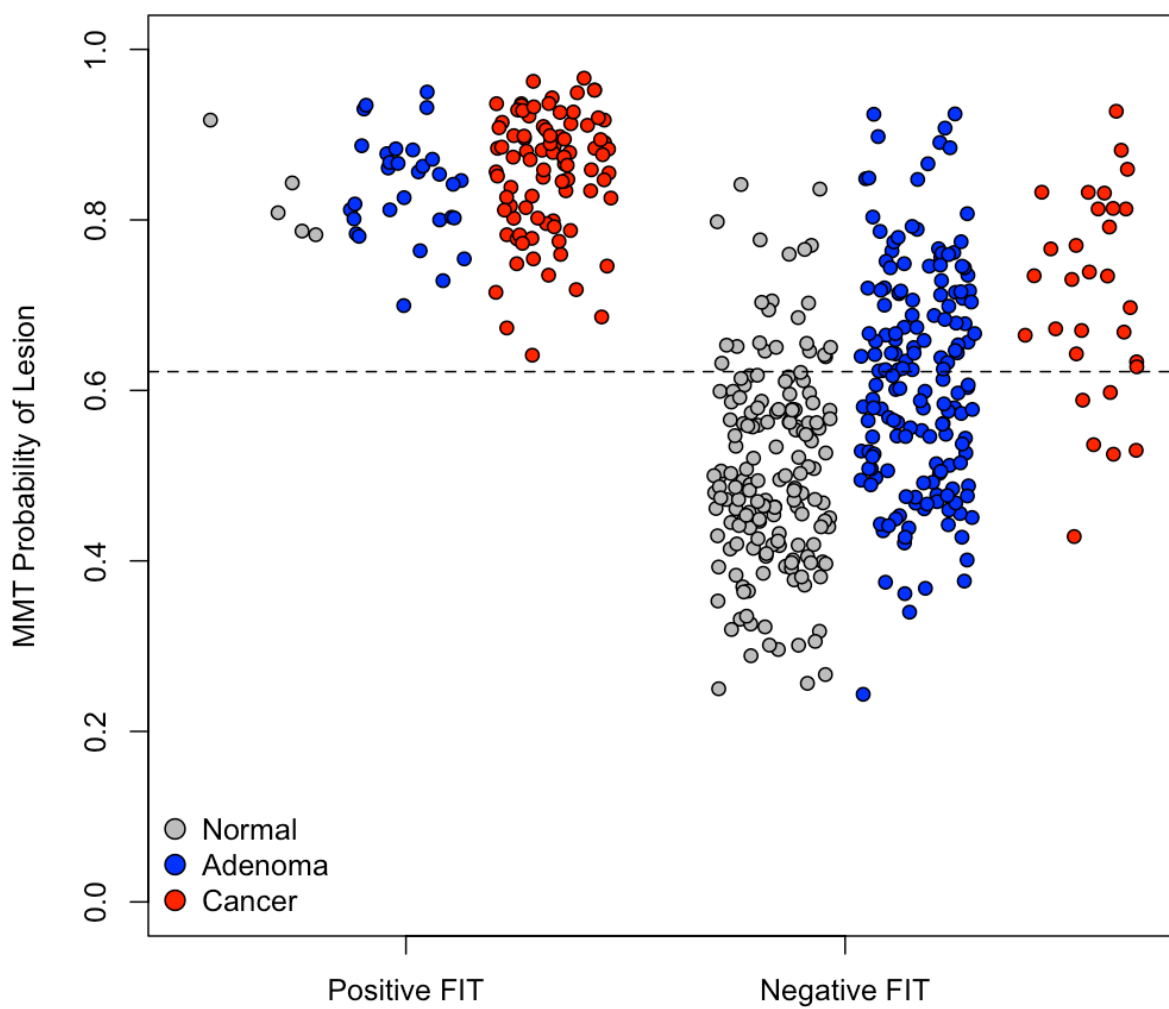
309 **Extended Data Figure 1.** Change in out-of-bag AUC with number of features in the MMT. The
 310 optimal model contains 24 features and has an AUC of 0.829.



311

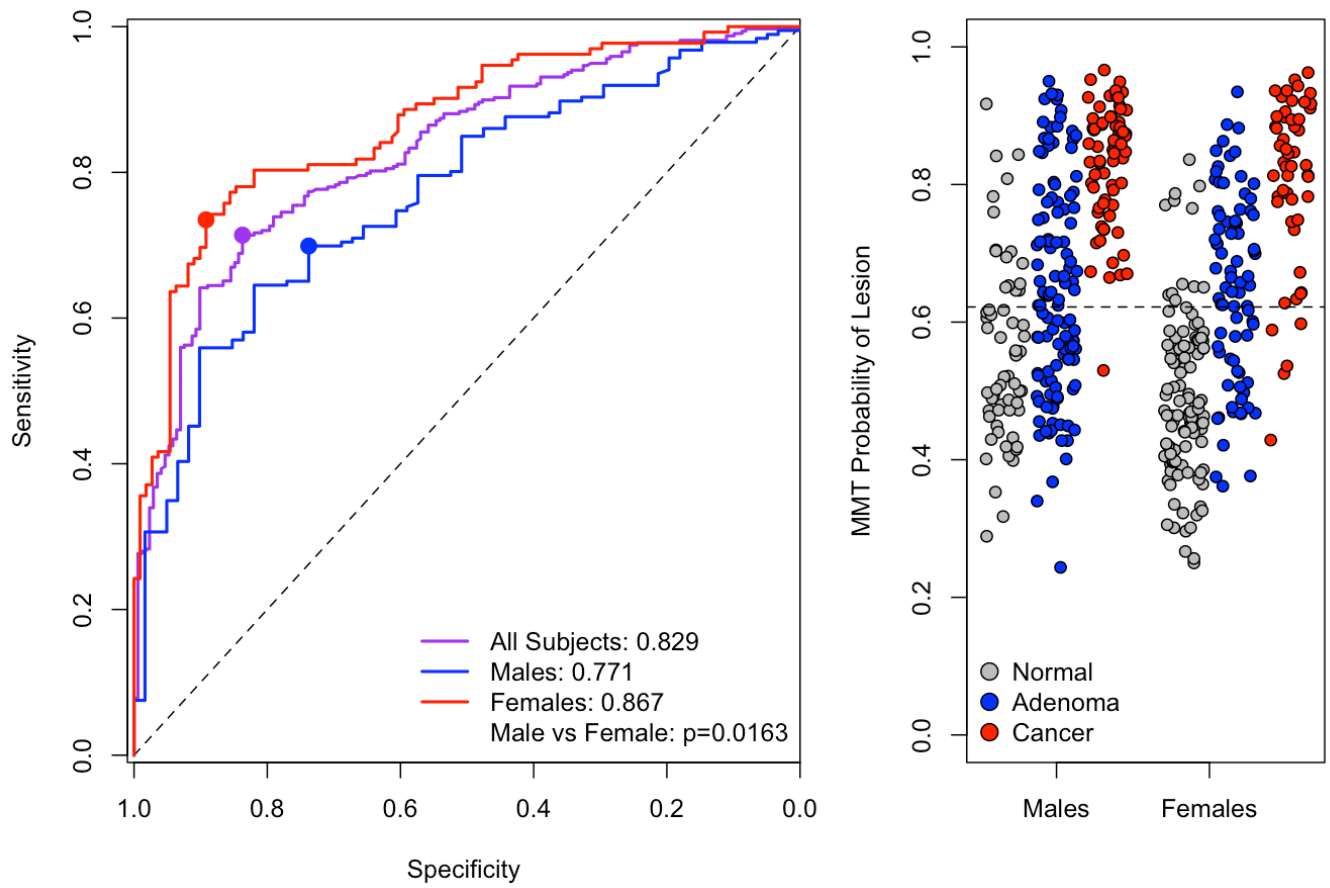
312 **Extended Data Figure 2.** Stripchart of the relative abundances of each OTU in the MMT with black

313 lines at the medians.



314

315 **Extended Data Figure 3.** Stripchart of MMT results for each sample based on FIT result.



316

317 **Extended Data Figure 4.** ROC curves (left) and stripchart (right) of MMT separated by sex.