

1 **TITLE GOES HERE**

2 **Running title:** Optimal resolution

3 Courtney R. Armour, (ande? begum? etc?), Patrick D. Schloss [†]

4 [†] To whom correspondence should be addressed:

5 pschloss@umich.edu

6 Department of Microbiology

7 University of Michigan

8 Ann Arbor, MI 48109

9 **observation format - max 1200 words, 2 figures, 25 ref**

¹⁰ **Abstract (max 250 words)**

¹¹ **Importance (max 150 words)**

12 Introduction

- 13 • CRC is one of the most common cancers and a leading cause of cancer related death
- 14 • There is evidence that the microbiome has a role in CRC development/progression and could be useful
15 for biomarker detection and diagnostics.
- 16 • Begum et al (mBio 2020) recently demonstrated effective application of machine learning (ML) to
17 microbiome based classification problems and developed a framework for applying ML practices in a
18 more reproducible way (mikropml).
- 19 • A common question when applying ML methods to microbiome data is which method and taxonomic
20 level should be use.
- 21 • This analysis utilizes the reproducible framework developed by Begum et al to quantify which ML
22 method and taxonomic level produce the best performing classifier for CRC data.

Results

- Across the five ML methods tested, model performance tends to increase with taxonomic level usually peaking around genus level and dropping off slightly with ASVs.
- Random forest was consistently the top performer at most taxonomic levels.
 - RF might be more appropriate for microbiome analysis since its more suitable for zero inflated data (need to look into literature)
- Within the RF model, the highest AUCs were observed for family, genus, and otu level data with no significant difference between the three (Figure 1). ASV performed significantly lower than OTU but not family and genus levels.

33 **Materials and Methods**

- 34 • data from prior study {baxter}
- 35 • mikropml package
- 36 • pvalues as previously described {begum}

