

# Efficient Prediction of Microplastic Counts from Mass Measurements

Shuyao Tan,<sup>†</sup> Joshua Taylor,<sup>‡</sup> and Elodie Passeport<sup>\*,†,¶</sup>

<sup>†</sup>*Department of Chemical Engineering and Applied Chemistry, University of Toronto, 200  
College Street, Toronto, ON M5S 3E5, Canada*

<sup>‡</sup>*The Edward S. Rogers Sr. Department of Electrical & Computer Engineering, University  
of Toronto, 10 King's College Road, Toronto, ON M5S 3G4, Canada*

<sup>¶</sup>*Department of Civil and Mineral Engineering, University of Toronto, 35 St. George  
Street, Toronto, ON M5S 1A4, Canada*

E-mail: [elodie.passeport@utoronto.ca](mailto:elodie.passeport@utoronto.ca)

## Abstract

Microplastics must be characterized and quantified to assess their impact. Current quantification procedures are time-consuming and prone to human error. This study evaluates the use of machine learning to estimate the number of microplastic particles based on aggregate particle weight measurements. Synthetic datasets are used to test the performance of linear regression, kernel ridge regression and decision trees. Kernel ridge regression achieves the strongest performance, and it is also tested with experimental datasets. The numerical results show that the algorithm is better at predicting the counts of larger and more homogeneous samples, and that contamination by organics does not significantly increase error. In mixed samples, prediction error is lower for heavier particles, with an error rate comparable to or better than that of manual counting. Overall, the proposed method is faster and easier than current approaches.

**Keywords:** Microplastics; Kernel Ridge Regression; Prediction; Machine Learning

## 14 1. Introduction

15 Plastic particles less than 5 mm in size are referred to as microplastics.<sup>1</sup> They are widespread  
16 in the environment, both in densely populated and remote areas.<sup>2-5</sup> Much remains unknown  
17 about their impact on organisms; while negative effects have been observed in several stud-  
18 ies, a clear, causal picture has not yet emerged.<sup>6</sup> This is in part due to the diversity of  
19 microplastics, i.e., shape, size, polymer identity, chemical mixture, color, and count, which  
20 are not accounted for in most studies on their ecological effects.<sup>6,7</sup>

21 Research on microplastics can have several different objectives, e.g., source identification,  
22 fate assessment, and (eco)toxicological impact evaluation, all of which require the charac-  
23 terization and quantification of microplastics. Depending on the subject, the quantity of  
24 microplastics can be referred to as count, concentration, abundance, or dose, expressed as  
25 particles per sampling area, per volume or per mass of sample.<sup>8-10</sup> Here, we will use the  
26 term count. The count of microplastics, though not sufficient alone to address all research  
27 questions, remains an important metric. Indeed, higher concentrations of microplastics are  
28 found near more populated areas, thus helping estimate proximity to sources. Microplastic  
29 counts are also used to evaluate the performance of water treatment systems.<sup>11,12</sup> Finally,  
30 the count of microplastics is also one of the drivers of observed effects on organisms.<sup>6</sup>

31 While currently there are no formal standards for microplastic characterization and quan-  
32 tification, the most common and non-destructive method involves some or all of the following  
33 steps: organic digestion, in which chemicals or enzymes are used to remove organics; den-  
34 sity separation, in which a solution of known density is used to separate lighter particles  
35 from heavier sediments; sieving, in which the particles are passed through sieves of different  
36 mesh sizes and grouped into different size ranges; visual identification, in which particles are  
37 manually sorted under a microscope for counting, sometimes particle size measurement, and  
38 classification for colors and shape categories.<sup>8-10,13-15</sup> The counting step is time-consuming  
39 and subject to human error. Usually a technique such as Raman or Fourier-transform in-  
40 frared spectroscopy (FTIR) is also used to verify that the particles are indeed plastic and

41 identify the polymer type, further increasing cost and time. As a result, characterization  
42 and quantification have become the rate-limiting steps in the analysis of microplastic con-  
43 tamination.

44 Recent efforts have attempted to make this process cheaper and faster by automation.  
45 Algorithms have been developed to automatically match  $\mu$ FTIR spectra to polymers, which  
46 take significantly less time than manual labour and involve less error<sup>3,16-18</sup>. Various machine  
47 learning algorithms have also been coupled with spectroscopic techniques to auto-identify  
48 microplastic particles.<sup>19</sup> For example, a convolutional neural network was trained to classify  
49 microbeads based on microscopic images, and support vector machines, adaptive boosting  
50 and random forests have been trained to identify particles from scanned images.<sup>20,21</sup> These  
51 methods are highly effective, but their applicability is limited for one or more of the following  
52 reasons: (i) they rely on expensive equipment, which require expertise to use, (ii) cannot  
53 reliably distinguish plastics from organic particles, and (iii) require the particles to be loosely  
54 distributed over a surface. Our objective is to develop an alternative that is cheap, fast, and  
55 simple enough that anyone can use it to reliably quantify microplastic contamination.

56 The objective of this study is to use machine learning to estimate the number of mi-  
57 croplastic particles, in different size ranges and morphology categories, from their aggregate  
58 weight measurements. Microplastic samples will be sieved into different size ranges, and  
59 the total particle weight of each size range will be measured. From the resulting set of  
60 weight measurements, we estimate the count of microplastic particles in each of the fol-  
61 lowing types: fragment, bead, film, fiber, and rubber. Koelmans *et al.* have proposed  
62 empirical relationships to calculate a microplastic mass-volume-number conversion factor  
63 based on data averaged over many microplastic studies.<sup>22</sup> In this study, we aim to develop  
64 a method for quantifying microplastic particles without assuming that all samples share  
65 similar compositions. We use three machine learning algorithms: linear regression, kernel  
66 ridge regression, and decision trees. We have chosen these techniques for their simplicity and  
67 low computational cost. While neural networks could potentially also work well, we prefer

68 these techniques for several reasons. First, the effectiveness of modern, deep neural networks  
69 comes from massive training datasets, which are not available in this application. Second,  
70 neural networks are more computationally intensive and often entail extensive parameter  
71 tuning, thus increasing the complexity of the procedure. In our numerical results, we will  
72 compare the performance of the algorithms to manual counting and evaluate the impact of  
73 mass measurement error on count prediction.

## 74 2. Methods

### 75 2.1. Experimental and Synthetic Datasets

76 Five sets of experimental data are taken from the literature. Isobe *et al.*<sup>23</sup> made samples by  
77 adding microplastic particles to sea water, and sent them to different laboratories around  
78 the world. These samples were sieved into multiple size ranges, weighed, and counted,  
79 resulting in 10 pairs of particle weight and count measurements. McIlwraith *et al.*<sup>24</sup> washed  
80 blankets in a washing machine, collected fibers from the effluent water, and reported 18  
81 pairs of total weight and count. Lavers *et al.*<sup>25</sup> reported 11 pairs of particle weight and  
82 count measurements obtained from sand samples collected in the Cocos islands. Puskic *et*  
83 *al.*<sup>26</sup> analysed microplastic particles in 52 sea birds (Short-Tailed Shearwaters), resulting in  
84 30 weight and count pairs. Corcoran and Arturo *et al.*<sup>27,28</sup> collected visible polymeric debris  
85 from 66 beaches across the Laurentian Great Lakes, measured the total weight of microplastic  
86 particles ranging from 1 mm to 5 mm, and reported particle count by morphology. These  
87 datasets will be referred to as *Interlab*, *Washing Machine*, *Cocos*, *Shearwaters*, and *Great*  
88 *Lakes*, respectively.

89 Due to the limited quantity and variation of published experimental data, a computa-  
90 tional procedure was developed to generate synthetic data. These synthetic datasets were  
91 used to evaluate the performance of the three algorithms under various circumstances. The  
92 procedure is summarized in Table S1. A particle's density depends on its polymer identity

93 and weathering state. The morphology (hereafter referred to as type) of a particle is defined  
94 by its largest, medium, and smallest dimensions. The number of each type of particles in  
95 each of the samples described below was randomly drawn from the ranges given in Table  
96 S1 in Supporting Information. On average, there can be around 10 rubbers, 40 fragments,  
97 50 films, 50 beads and 300 fibers in a sample. These values are approximates to analysed  
98 samples collected from freshwater systems that are subject to human activities.<sup>3,29-31</sup> The  
99 probability density plots of the number of each type of particle in the synthetic datasets can  
100 be found in Section S1 in Supporting Information. The dimensions of each particle were set  
101 as follows.

102 1. Rubber particles and fragments had similar dimensions but different densities. Their  
103 largest dimension was uniformly distributed between 153  $\mu\text{m}$  and 1500  $\mu\text{m}$ . The medium  
104 dimension was 70 to 100% of the largest dimension. The smallest dimension was 40 to 100%  
105 of the medium dimension. Rubber particles were lighter than fragments.

106 2. Films had similar densities to fragments, but their smallest dimensions were smaller.  
107 The largest and medium dimensions of films were drawn from the same range as those of  
108 fragments. The smallest dimension for films was 5 to 10% of the medium dimension.

109 3. Beads were assumed to be perfect spheres. Each bead's diameter was a uniformly  
110 drawn random number between 107  $\mu\text{m}$  and 1500  $\mu\text{m}$ .

111 4. Fibers were randomly divided into long and short fibers. Their diameters were between  
112 10 and 40  $\mu\text{m}$ . The length of each long fiber was 30 to 50 times of its diameter, and the  
113 length of each short fiber was 10 to 30 times of its diameter.

114 5. The largest dimension of organic particles was between 214 and 2000  $\mu\text{m}$ . The medium  
115 dimension was 50 to 100% of the largest dimension. The smallest dimension was between 40  
116 and 150  $\mu\text{m}$ .

117 The sieving process for non-fiber particles was simulated by comparing the size of each  
118 particle with the mesh size. If both the largest and medium dimensions were larger than  
119 the mesh size, the particle remained on the sieve; if not, the particle passed through the

120 sieve to the next one. In this study, the sieve mesh sizes were: 1000  $\mu\text{m}$ , 500  $\mu\text{m}$ , 300  $\mu\text{m}$   
121 and 106  $\mu\text{m}$ . Although fibers were thin enough to pass through every sieve, they could still  
122 become entangled with the mesh or with other particles. To approximate this randomness,  
123 the probability for long fibers to stay on each sieve was set to be 0.3, 0.4, 0.2, 0.1, and for  
124 short fibers, the distribution was 0.15, 0.35, 0.3, 0.2. These probabilities were estimated  
125 from storm water sample data collected from a parking lot using the same four sieve sizes  
126 specified above.<sup>29</sup>

127 After sieving, the total mass of all particles on each sieve was calculated by summing up all  
128 individual particle masses. Individual masses were obtained by multiplying the volume and  
129 density of each particle. For cuboid-shaped particles, volume was calculated by multiplying  
130 the three dimensions. Fiber volume was the product of its cross-sectional area and length.  
131 Spherical particles were treated as perfect spheres.

132 In this study, 10 datasets with different particle combinations were generated, each con-  
133 taining 100 samples (Section S2). These 10 combinations were:

- 134 1. Film and fiber.
- 135 2. Bead and fiber.
- 136 3. Rubber, fragment and fiber.
- 137 4. Rubber, fragment, fiber and half as many organic particles as fragments.
- 138 5. Rubber, fragment, fiber and as many organic particles as fragments.
- 139 6. Rubber, fragment, fiber and film.
- 140 7. Rubber, fragment, fiber, film and as many organic particles as fragments.
- 141 8. Rubber, fragment, fiber and bead.
- 142 9. Rubber, fragment, fiber, film and bead.
- 143 10. Rubber, fragment, fiber, film, bead and as many organic particles as fragments.

144 Fibers appeared in all samples because they are the most common type of particle found  
145 in water, sediment and air.<sup>29,32,33</sup> Fragments appeared in most samples because they are the  
146 second most common particle in water and sediment.<sup>32</sup> The above synthetic datasets form

147 the base case for the numerical experiments in later sections. To test the dependence of  
 148 algorithm performance on sample size, two more cases were generated with composition 9,  
 149 but with more particles, as described in Table S2.

## 150 2.2. Machine Learning Algorithms

151 Let  $x_i = [x_{i1}, x_{i2}, \dots, x_{ik}]$  be a row vector of the masses in each of  $k$  sieves for the  $i^{\text{th}}$  sample.  
 152 The matrix  $X = [x_1; x_2; x_3; \dots; x_n]$  contains the mass measurements for all  $n$  samples. Let  
 153  $y_i = [y_{i1}, y_{i2}, \dots, y_{im}]$  be a row vector of the number of each type of particle in each sieve for  
 154 the  $i^{\text{th}}$  sample, and  $Y = [y_1; y_2; y_3; \dots; y_n]$  be the matrix of all counts. Let  $\hat{Y} = f(X)$  be an  
 155 estimate of  $Y$  based on some function of  $X$ . In each of the following learning frameworks,  
 156 we assumed some simple structure for  $f$ , and then used training data consisting of known  
 157  $(X, Y)$  pairs to solve for the parameters of  $f$  that minimize prediction error.

### 158 Linear Regression (LR)

159 Linear regression assumes a relationship of the form

$$\underbrace{\begin{bmatrix} \hat{y}_{11} & \hat{y}_{12} & \dots & \hat{y}_{1m} \\ \hat{y}_{21} & \hat{y}_{22} & \dots & \hat{y}_{2m} \\ \vdots & \vdots & \vdots & \vdots \\ \hat{y}_{n1} & \hat{y}_{n2} & \dots & \hat{y}_{nm} \end{bmatrix}}_{\hat{Y}} = \underbrace{\begin{bmatrix} 1 & x_{11} & x_{12} & \dots & x_{1k} \\ 1 & x_{21} & x_{22} & \dots & x_{2k} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & x_{n1} & x_{n2} & \dots & x_{nk} \end{bmatrix}}_X \underbrace{\begin{bmatrix} b_1 & b_2 & \dots & b_m \\ a_{11} & a_{12} & \dots & a_{1m} \\ a_{21} & a_{22} & \dots & a_{2m} \\ \vdots & \vdots & \vdots & \vdots \\ a_{k1} & a_{k2} & \dots & a_{km} \end{bmatrix}}_A \quad (1)$$

160 Given training data  $(X, Y)$ , the matrix  $A$  that minimizes the mean squared error is given by

$$A = (X^T X)^{-1} X^T Y. \quad (2)$$

### 161 Kernelized Ridge Regression (KRR)

162 Ridge regression is a modification of linear regression.<sup>34</sup> Linear regression assumes that

163 there is no relationship between the independent variables. If there is some relationship  
164 among the independent variables, ridge regression often gives more reliable results than  
165 linear regression.<sup>35</sup> Ridge regression also has the form  $\hat{Y} = XA$ , but the matrix  $A$  is given  
166 by

$$A = (\lambda I + X^T X)^{-1} X^T Y, \quad (3)$$

167 where  $\lambda \geq 0$  is a complexity penalty and  $I$  is the identity matrix. The complexity penalty  
168 discourages large magnitude entries in  $A$ , hence reducing overfitting. This equation can also  
169 be written as

$$A = X^T (X X^T + \lambda I)^{-1} Y. \quad (4)$$

170 In KRR, the term  $X X^T$  is replaced by the Gaussian kernel function, given by

$$K(x_i, x_j) = \exp\left(-\frac{\|x_i - x_j\|^2}{2\sigma^2}\right), \quad (5)$$

171 where  $\sigma^2$  is referred to as the bandwidth of the kernel. The Gaussian kernel measures the  
172 similarity between any two points. Given a new input  $x$ , KRR identifies the portions of  
173 the training data with the highest similarity, which therefore have larger influences on the  
174 resulting estimate. Let  $K$  be an  $n$  by  $n$  matrix in which the  $ij^{\text{th}}$  entry is  $K(x_i, x_j)$ , and  
175 define

$$r = (K + \lambda I)^{-1} Y. \quad (6)$$

176 The parameter matrix  $A$  is given by

$$A = X^T r. \quad (7)$$

177 Given new input data  $x$ , the predicted output is given by

$$\hat{Y} = A^T x. \quad (8)$$



178 Note that this is equivalently written as

$$\hat{Y} = \sum_{i=1}^n r_i x_i^T x = \sum_{i=1}^n r_i K(x_i, x). \quad (9)$$

179 In this study, the Gaussian kernel parameters were found using grid search, as described  
180 in Section S2.

### 181 **Decision Tree (DT)**

182 Decision trees recursively partition a given dataset to minimize information entropy.<sup>34</sup>  
183 Information entropy measures the disorder of a dataset, and is given by

$$H = - \sum_{i=1}^n p(x_i) \log_2 p(x_i), \quad (10)$$

184 where  $p$  stands for probability. After each partition, the total entropy of the resulting subsets  
185 should be smaller than before. Training a decision tree only requires that one specifies its  
186 depth, that is, the level of partition. For example, a tree of depth two splits the original  
187 dataset into four subsets. Each final subset is called an end node. A prediction is made by  
188 passing a new data point through the decision tree, and outputting the average value of the  
189 training data at the resulting end node.

## 190 **2.3. Training and Evaluation Methods**

191 Training is the process of learning model parameters from given  $(X, Y)$  pairs. In this case, the  
192  $(X, Y)$  pairs that make up the training data could be obtained in several ways, for example,  
193 from traditional counting-based quantification, or from an existing dataset. Evaluation of  
194 a model is done by inputting data  $X$ , then comparing the estimated  $\hat{Y}$  with the known  $Y$ .  
195 In this study, the learning algorithms were trained and tested through leave-one-out cross-  
196 validation. In each dataset of 100 samples, 99 were used for training and the left-out sample  
197 was used to test the trained algorithm's performance. This process was repeated for every

198 sample in each dataset, hence yielding 100 predictions. The accuracy of each algorithm was  
199 quantified by mean absolute error (MAE), mean percentage error (MPE) and normalized  
200 root mean square error (NRMSE), which are defined as

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |\hat{y}_i - y_i| \quad (11)$$

$$\text{MPE} = \frac{1}{n} \sum_{i=1}^n \left| \frac{\hat{y}_i - y_i}{y_i} \right| \quad (12)$$

$$\text{NRMSE} = \frac{\sqrt{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - y_i)^2}}{y_{\max} - y_{\min}}, \quad (13)$$

201 where  $y$  and  $\hat{y}$  are respectively the true and predicted values. There were occasional samples  
202 in which the number of a certain type of particle was zero, and these cases were not included  
203 in MPE calculation. Another normalization factor,  $y_{\max} - y_{\min}$ , was used to calculate NRMSE  
204 to avoid zero denominators. We conducted our numerical experiments as described above,  
205 and also with measurement error, which reflects non-idealities such as balance roundoff. We  
206 simulated these errors by adding random numbers uniformly distributed between  $-30\%$  and  
207  $30\%$  of the original value to the two finer mesh sieve weights. Random numbers uniformly  
208 distributed between  $-10\%$  and  $10\%$  were added to the two coarser mesh sieve weights.

209 The 95% confidence intervals of the synthetic datasets were calculated using 1000 boot-  
210 strapped samples, in order to achieve a normal distribution. The statistical difference be-  
211 tween two sets of predictions were evaluated using the Wilcoxon Signed-Rank Test. A  
212 significance level of  $\alpha = 0.05$  was used.

213 To test the effect of training dataset size on prediction accuracy, the following number  
214 of training samples were also tested: 90, 80, 70, 60, 50, 40, 30, 20, 10, 8, 6, 4, and 2. The  
215 corresponding testing dataset sizes were 10, 20, 30 and so on. For each synthetic sample  
216 composition, the training-and-testing procedure was repeated 20 times, and each time the

217 training samples were randomly selected. The final result was reported as MPE for each  
218 type of particle against the number of training samples.

## 219 3. Results and Discussion

### 220 3.1. Training with Synthetic Data

Table 1: Mean Absolute Prediction Error

	KRR <sup>a</sup>	LR <sup>b</sup>	DT <sup>c</sup>	KRR-E <sup>d</sup>	LR-E <sup>d</sup>	DT-E <sup>d</sup>
Rubber	1.98 ± 0.07	2.14 ± 0.07	2.37 ± 0.10	2.02 ± 0.08	2.14 ± 0.07	2.34 ± 0.09
Fragment	2.06 ± 0.08	2.10 ± 0.07	2.99 ± 0.10	2.19 ± 0.08	2.23 ± 0.08	3.04 ± 0.10
Bead	3.25 ± 0.18	3.41 ± 0.18	4.17 ± 0.20	3.40 ± 0.18	3.60 ± 0.18	4.14 ± 0.20
Film	5.62 ± 0.25	5.79 ± 0.23	6.21 ± 0.25	5.72 ± 0.23	5.88 ± 0.25	6.50 ± 0.28
Fiber	20.2 ± 0.53	20.7 ± 0.53	22.3 ± 0.59	20.7 ± 0.49	21.1 ± 0.53	22.7 ± 0.62
Non-Fiber <sup>e</sup>	3.67 ± 0.12	3.74 ± 0.14	5.07 ± 0.18	3.89 ± 0.14	3.97 ± 0.14	5.10 ± 0.17

<sup>a</sup> Kernel Ridge Regression; <sup>b</sup> Linear Regression; <sup>c</sup> Decision Tree, implemented with Scikit-Learn<sup>36</sup>; <sup>d</sup> With weight measurement error; <sup>e</sup> Total number of microplastic particles excluding fibers

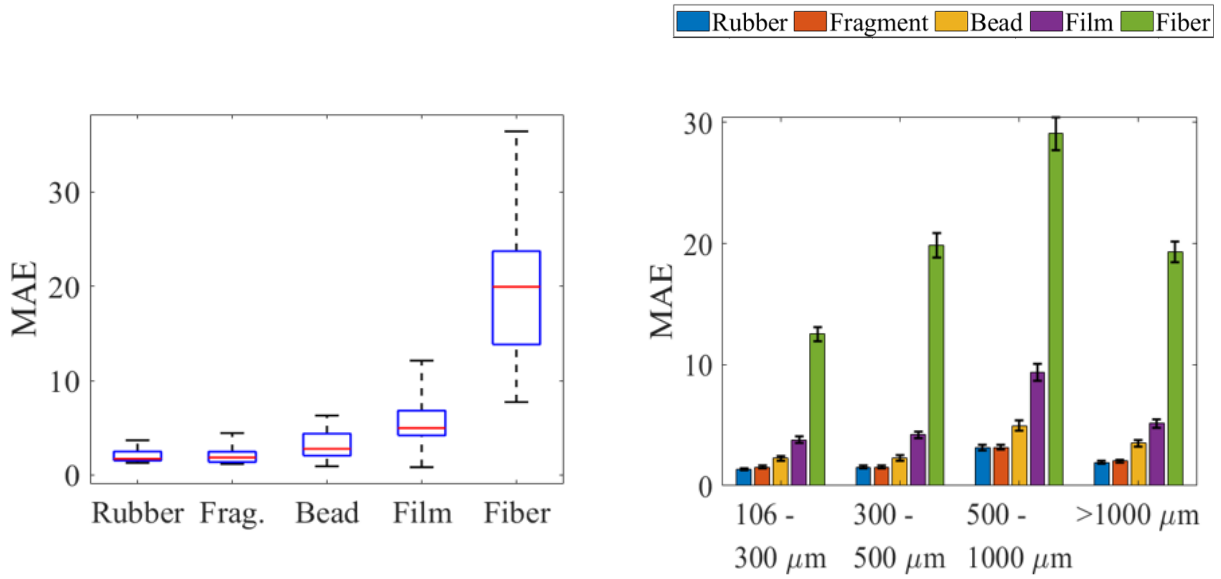
221 We test the algorithms on the synthetic data described in the Methods Section. The  
222 performances of the three algorithms were evaluated by comparing the MAE with 95% con-  
223 fidence interval for the synthetic samples. The results are reported in Table 1. The KRR  
224 prediction results were statistically different from LR and DT (p-values < 10<sup>-8</sup>). The results  
225 of LR and DT were similar (p-value = 0.57). Among the three algorithms, KRR always re-  
226 sulted in the lowest average MAE, both with and without incorporating weight measurement  
227 error. After adding weight measurement error, the performance of each algorithm did not  
228 significantly decrease (p-values were 0.71, 0.13 and 0.97 for KRR, LR and DT respectively).  
229 This suggests that all algorithms are somewhat robust to weight measurement error. We  
230 remark that DT is often poor at predicting continuous values because its output comes from  
231 a discrete set of values, the size of which is the number of end nodes.

232 KRR can be expected to outperform LR because, whereas LR assumes a linear relation-  
233 ship between the data, KRR can accommodate linear and some non-linear relationships.

234 While highly scalable, KRR is more complicated than LR in that one must tune its model  
235 parameters ( $\lambda$  and  $\sigma$ ). When simplicity is more important than accuracy, LR may be a  
236 superior choice.

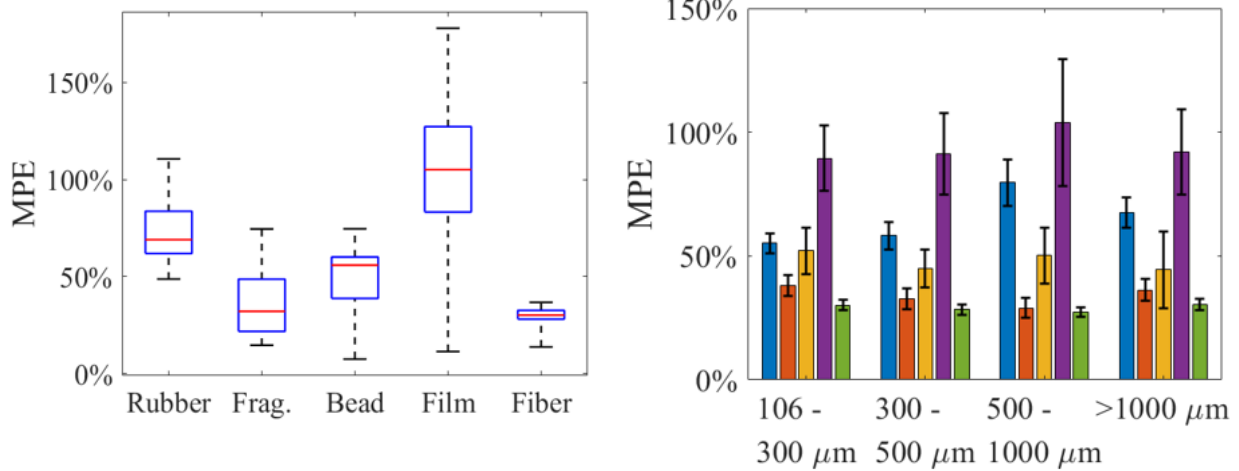
237 We hereon focus on KRR because it achieved the best performance. Its error for the  
238 synthetic datasets is summarized in Figure 1. On rubber, KRR achieves low MAE, high  
239 MPE and high NRMSE because rubber is usually present in smaller numbers than other  
240 particle types. As a result, even an absolute error of one or two particles can result in a  
241 percentage error that is more than 50%. On the contrary, fibers are the most numerous.  
242 This leads to a large MAE, which, when divided by the large number of fiber particles,  
243 leads to a small MPE. For other particles, despite the fact that they are present in similar  
244 numbers, films have larger median and average errors than beads, and beads have larger  
245 errors than fragments. MAE differs for each size fraction because each size fraction has  
246 different amounts of particles. The 500 – 1000  $\mu\text{m}$  size range has the most particles and thus  
247 the largest MAE. NRMSE for rubber, fragment, film and fiber is almost constant for all size  
248 ranges. Altogether, these results show that the larger the number of microplastics in a given  
249 category, the higher its MAE and the lower its MPE.

250 Figure 1g shows the stacked MAE for each of the ten synthetic datasets. The smallest  
251 prediction errors (MAE, MPE and NRMSE) for films and beads occur in their respectively  
252 purest samples, e.g., dataset 1 is composed entirely of films and fibers, and dataset 2 of beads  
253 and fibers. These values are even lower than that of fragments, because fragments always  
254 appear with rubbers in the synthetic datasets. MAE for beads and films increases dramati-  
255 cally with the addition of fragments and rubbers, while MAE for fragments only increases  
256 notably with the addition of beads. Also note that the smallest MAE for fibers occurs in  
257 dataset 1, made of only films and fibers. This is most likely because films are very light,  
258 so that fibers have a higher weight percentage in this dataset than in others. For mixture  
259 samples containing fragments, beads and films, MPEs for each type of particles increase  
260 with their decreasing weight percentage, as shown in Table S3: average MPE can increase



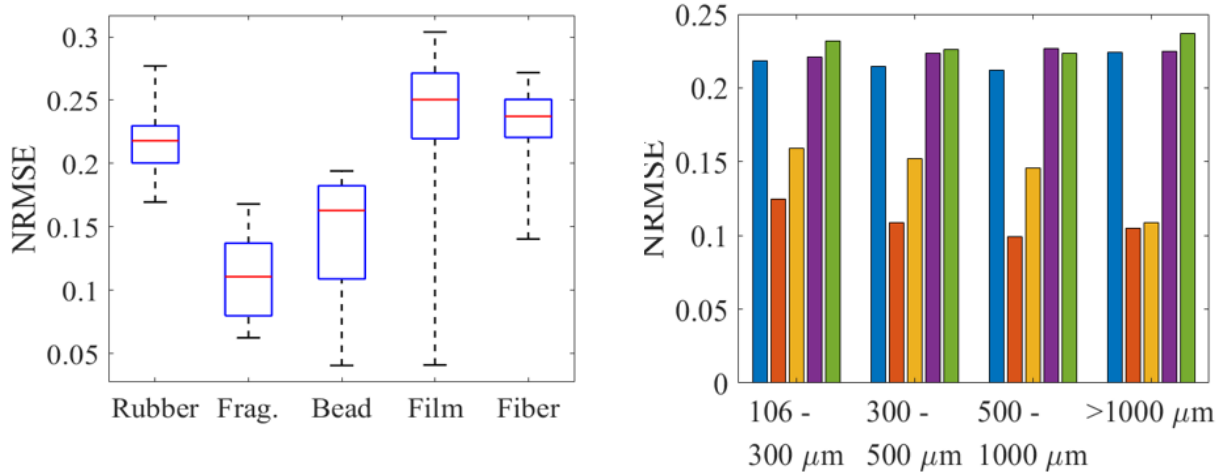
(a) MAE Boxplot

(b) MAE Distribution by Size



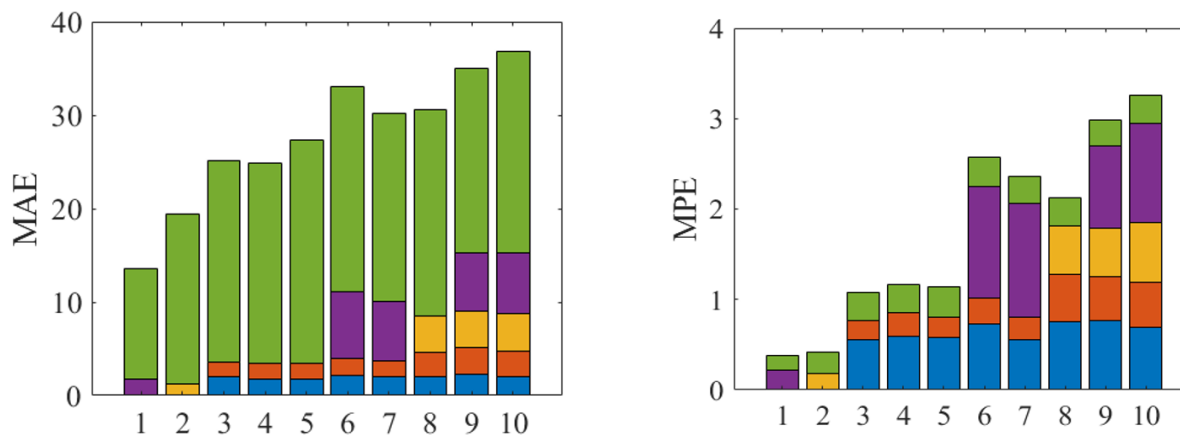
(c) MPE Boxplot

(d) MPE Distribution by Size



(e) NRMSE Boxplot

(f) NRMSE Distribution by Size



(g) Stacked MAE for each type of particle in each dataset.

(h) Stacked MPE for each type of particle in each dataset.

Figure 1: (Previous page.) Error of KRR(Kernel Ridge Regression) on synthetic data. In the boxplots, the upper and lower bars represent the upper and lower bounds of the data. The box represents the 25 and 75 percentile, and the middle line represents the median value. Error bars on the right hand side are 95% bootstrapped confidence interval. Frag. represents fragment. Sample complexity increases from composition 1 to 10. MAE: mean absolute error. MPE: mean percentage error. NRMSE: normalized root mean square error.

261 from 14% to 51% for fragments if their weight percentage drops from 80% to 20%. All these  
262 observations suggest that, when aggregate weight is used as the sole feature, heavier particle  
263 types have consistently lower prediction error because of their higher weight percentage in  
264 each sample. We can also see that MAE does not change significantly after adding various  
265 amounts of organics, most likely because the total weight of organics is smaller than that of  
266 the plastics, and the addition of organics does not change the order of weight percentages  
267 of other particles. For samples severely contaminated by organic particles, e.g., a lake with  
268 large amount of algae or plant debris sampled along with microplastics, an organic digestion  
269 step will be necessary prior to drying and measuring the sample weight, in order to limit the  
270 contribution of organics to the total sample weight.

271 To summarize, the procedure is accurate for particles that account for a significant frac-  
272 tion of a sample's weight. This is why it is relatively difficult to predict fiber counts in mixed  
273 samples, which are numerous but might account for only a few percent of the overall weight.

274 On the other hand, it is relatively easy to predict fragment counts in mixed samples because  
275 they account for a moderate fraction of the overall weight.

276 The total error associated with microplastic count prediction by KRR remains compara-  
277 ble to or lower than that of manual counting, even for the more heterogenous samples. Isobe  
278 *et al.*<sup>23</sup> used artificial samples to estimate the relative uncertainty of microplastic counts  
279 performed by humans. In their study, samples with  $n_0 \approx 90$  fragments less than 2 mm were  
280 sent to 12 different laboratories. The relative error of a given laboratory's count,  $n$ , was  
281 calculated as  $(n - n_0)/n_0$ . The average error of all laboratories was around  $\pm 50\%$ , which we  
282 interpret as the error of a typical human count. In other visual sorting microplastic studies,  
283 researchers have also reported a false positive ratio ranging from 20% to 70%, meaning that,  
284 among the particles visually sorted as microplastics, 20% to 70% of them were determined  
285 via spectroscopic analysis to be non-plastic.<sup>8,29,37,38</sup>

286 In our synthetic datasets, the average number of fragments, beads and films is around  
287 50 for each, whenever present. Therefore, based on the similarity of *Interlab* samples to  
288 our synthetic data, it is reasonable to compare our results with theirs, with respect to each  
289 particle type. As discussed earlier, when these particles are present in their respective purest  
290 samples, their MPE can be less than 10%. In mixed samples, however, MPE for beads can go  
291 beyond 60% and films can go beyond 100%. Figure 1h shows MPE for each type of particle  
292 in the ten synthetic datasets. In mixed samples, MPE for films is always above 90%, due  
293 to their low contribution to the total weight, and that for rubbers is always above 50%, due  
294 to their small amount in each sample. MPE for fragments and fibers is below 35% except  
295 for datasets 8, 9, and 10, where the presence of beads increases the prediction error for  
296 fragments. When all particles are present, e.g., datasets 9 and 10, the prediction accuracy  
297 is better for fibers, comparable for beads and fragments, and worse for films and rubbers,  
298 compared to the accuracy of visual sorting. However, MPE decreases with increasing sample  
299 size, as indicated in Table S4. As the number of each type of particle increases, MPE for  
300 fragment, bead and film gradually decreases. MPE for fibers does not change significantly,

301 but MPE for rubber increases notably. This is because the number of rubber particles is a  
302 fraction of the number of fragments, as stated in Table S2. Then, there are several cases  
303 where the number of rubber particles is small, and the prediction is large, resulting in a large  
304 average MPE.

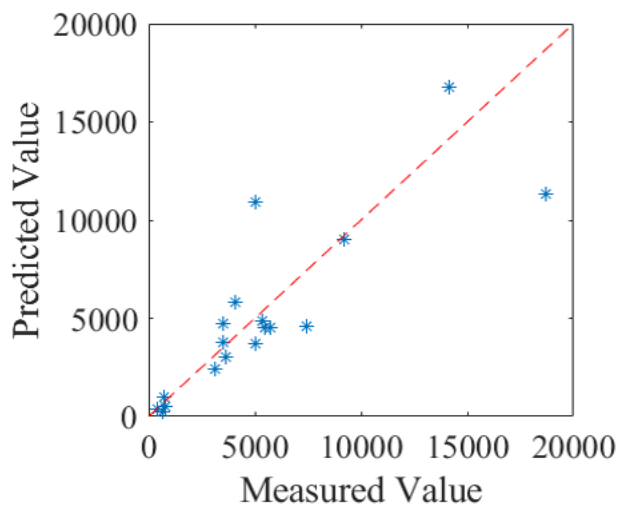
305 The above leave-one-out experiments correspond to training data with 99 samples, which  
306 may be impractically large in some cases. The effect of the number of training samples on  
307 KRR prediction accuracy can be found in Section S4. In general, MPE bottoms out when  
308 the number of training samples exceeds 20, regardless of sample composition and particle  
309 type. Therefore, in this case, there is little further accuracy to be gained from having more  
310 than 20 samples in the training data, and the leave-one-out results can be seen as the typical  
311 KRR performance.

312 These findings suggest that KRR, if given a moderate amount of accurate training data,  
313 can estimate counts more accurately than visual sorting on samples with homogeneous com-  
314 position, on heavier particles in mixed samples, and on larger samples, i.e., samples with  
315 larger particle number and total weight.

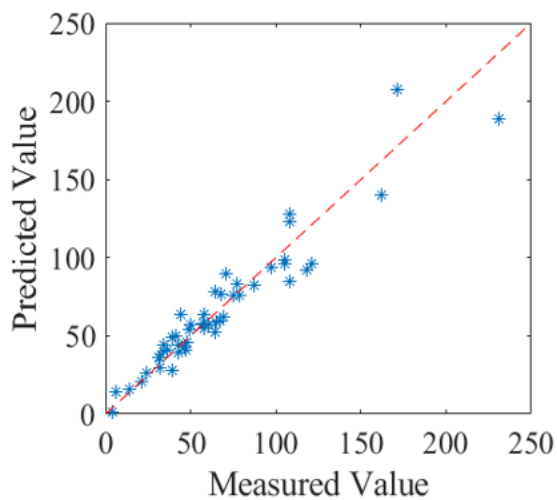
### 316 **3.2. Training with Experimental Data**

317 Figure 2 shows the performance of KRR on five datasets obtained from experiments: *Interlab*,  
318 *Washing Machine*, *Cocos*, *Shearwaters*, and *Great Lakes*. Training and testing was done  
319 through leave-one-out cross-validation. Data preprocessing was performed only for the *Great*  
320 *Lakes* dataset. For the latter, we first merged the 11 types of morphologies from the original  
321 report down to 6 (fiber, film, foam, fragment, pellet and rubber) in order to reduce data  
322 sparsity. Then, 3 outliers (outside three times the interquartile range of all data) were  
323 excluded from training and testing. Finally, training and testing was performed with the  
324 number of each type of particle, and with the total number of all particles, separately. Figures  
325 2a to 2f show the predicted values plotted against the measured values. In most cases, the  
326 outcomes scatter closely around the diagonal line, which represents perfect prediction. The

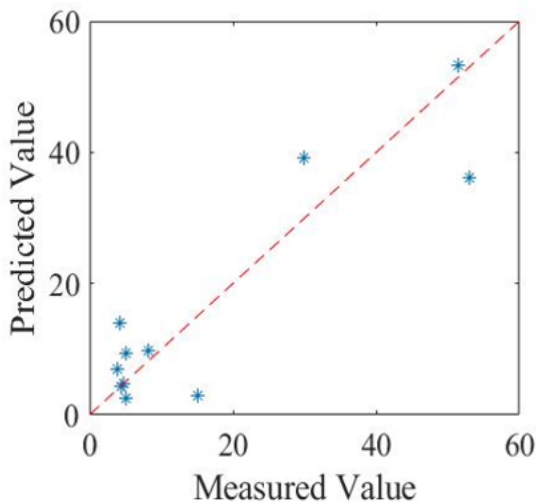




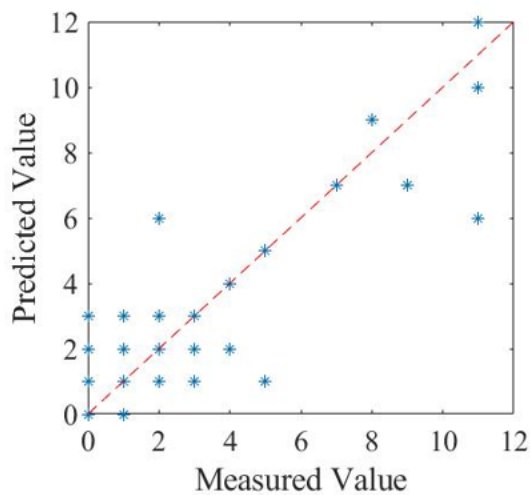
(a) *Washing Machine*



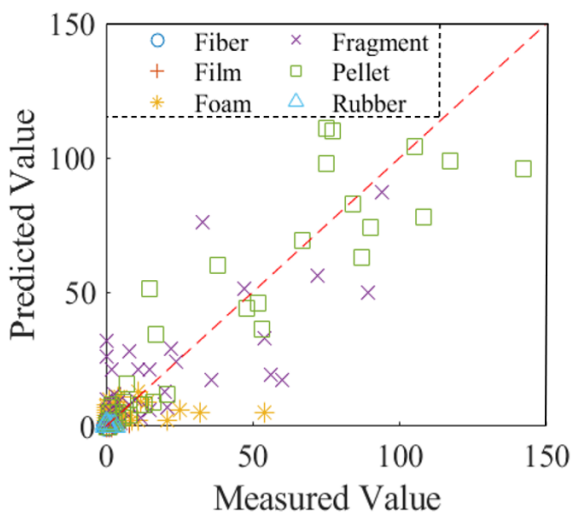
(b) *Interlab*



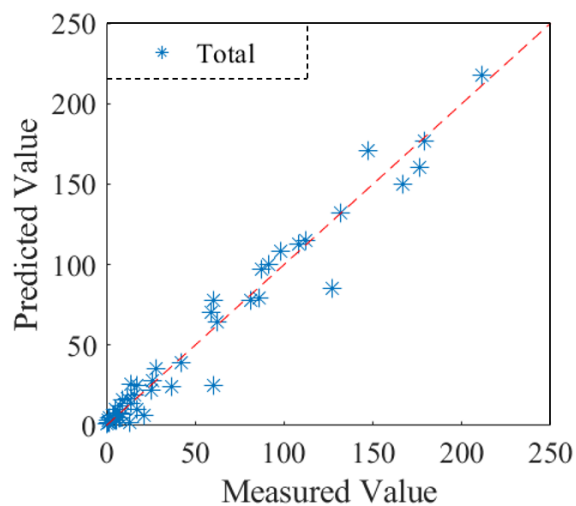
(c) *Cocos*



(d) *Shearwaters*



(e) *Great Lakes*, predicting the number of each type of particles



17

Figure 2: KRR (Kernel Ridge Regression) Performance Summary on Real Data

(f) *Great Lakes*, predicting the total number of all microplastic particles

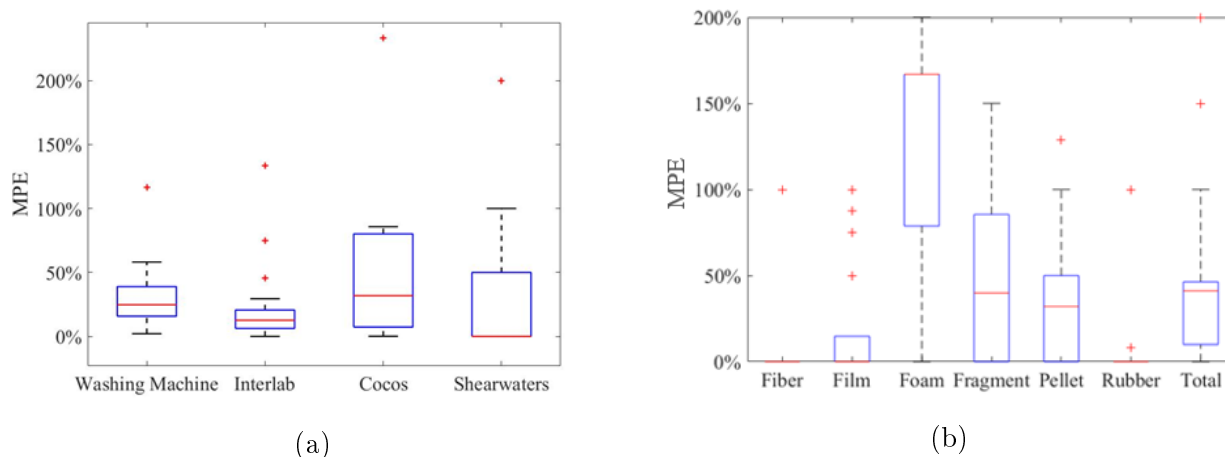


Figure 3: MPE (Mean Percentage Error) of KRR (Kernel Ridge Regression) prediction on: (a) the *Interlab*, *Washing Machine*, *Cocos*, and *Shearwaters* datasets, and (b) the *Great Lakes* dataset, with the display of y-axis limited at 200%. Original version is included as Figure S16(a). In the boxplots, the upper and lower bars represent the upper and lower bounds of the data. The box represents the 25 and 75 percentile, and the middle line represents the median value.

327 MPE for all datasets are mainly below 50%, with median values at 25% for *Washing Maching*,  
328 13% for *Interlab*, 32% for *Cocos* and 0% for *Shearwaters*, as shown by Figure 3a. The MPE  
329 for the *Great Lakes* dataset varies for different types of particles (Figure 3b), but in most  
330 cases the median value is less than 50%. MPEs for fiber, film and rubber are low in this  
331 case because the numbers of these particles show small variances, as demonstrated by Figure  
332 S16(b). MPEs for pellet, fragment and foam show a similar trend as found with artificial  
333 datasets: prediction error is smaller for heavier particles. These results demonstrate that  
334 predictions from a well-trained KRR can be more accurate than the results of traditional  
335 counting. However, there are a few large prediction errors, which could be the result of  
336 patterns undetected by KRR, or errors in the measured counts.

337 If the original count error is Gaussian with zero mean and constant variance, it will not  
338 affect the expected value of the model parameters.<sup>39</sup> In the case of microplastic quantifi-  
339 cation, it is possible for human counts to be biased, violating the zero-mean requirement,  
340 e.g., by always over-counting cotton fibers from air contamination or ignoring a type of  
341 particle thought to be non-plastic. To reduce the effect of biased measurement errors, train-

342 ing datasets should be as accurate as possible, for example, verified by FTIR or Raman  
343 spectroscopy and corrected to report plastic-only particles. Most studies that used such  
344 methods to confirm particle identities (plastic vs. non-plastic), showed that between 20% to  
345 70%<sup>8,29,37,38</sup> of the total number of microparticles are not plastics, and for sediment samples  
346 this rate can be as high as 98%.<sup>40</sup> As a consequence, accurate training data is necessary for  
347 making accurate predictions.

348 To optimize prediction quality, a minimum training dataset size is required. The recom-  
349 mended training dataset size increases with the number of predicted variables. When the  
350 number of predicted quantities is small, e.g., only fiber counts were predicted for the *Wash-*  
351 *ing Machine* data, a training dataset with 8 or more samples can be adequate, as suggested  
352 by Jenkins *et al.*<sup>41</sup> and shown in Figure S15. In cases like synthetic dataset 9, which has 4  
353 size fractions and 5 types of particle, a training dataset of 20 samples is recommended, as  
354 demonstrated by our numerical experiments. The training samples should cover as much the  
355 data range as possible, meaning that, the samples should be representative of the full range  
356 of inputs and outcomes.

### 357 **3.3. Implementation**

358 Assessing the extent and impact of microplastic contamination requires the collection of sam-  
359 ples at high temporal and spatial resolutions. This section briefly compares our approach  
360 to standard quantification procedures and demonstrates its power to quickly collect large  
361 amounts of data. Figure 4 shows the two quantification procedures after organic digestion,  
362 density separation, and sieving. This comparison mainly focuses on time, cost and level  
363 of expertise required. The equipment needed for traditional manual counting method costs  
364 approximately 1,000 USD for visual sorting. The cost of prediction based on sieved mass is  
365 approximately 5,000 USD to 7,000 USD in equipment for sample mass determination. Com-  
366 mercial laboratories charge between 1,000 to 5,000 USD for microplastic sample analysis,  
367 and around 50 to 250 times less for total suspended solid concentration measurement ( $\sim 20$

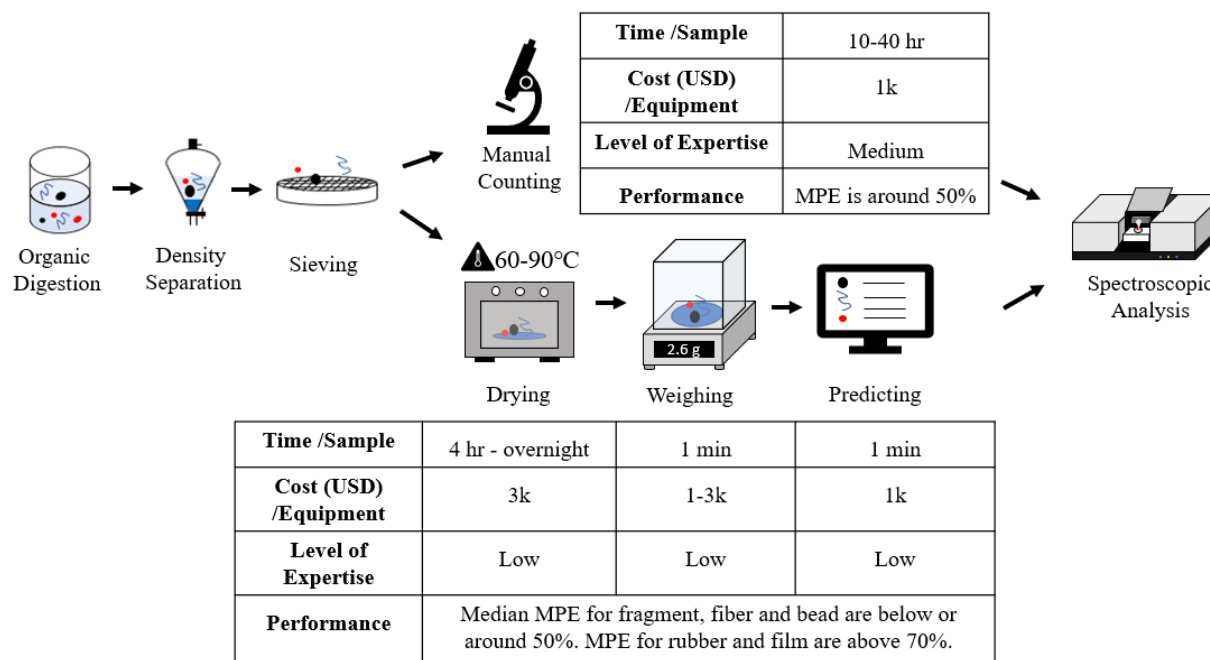


Figure 4: Comparison between the current and the proposed method. k stands for thousand. The time needed to manually count microplastic particles using a microscope was estimated from laboratory experiences, at around 10 to 40 hours per sample, with each sample containing around 50 to 100 particles. Predicting counts from sample weights consists of three parts: 1) drying, which ranges from 4 hours to overnight depending on the temperature used, 2) weighing, estimated as 1 minute, and 3) prediction, estimated as 1 minute for an already trained algorithm.<sup>24,42</sup> To this, one can add the time needed to process the required number of samples through the traditional manual counting method to generate a training dataset, i.e., an extra 10 to 40 hours per training sample. The equipment cost for manual counting is dominated by the microscope, estimated as 1,000 USD. The cost for the prediction method is estimated as the total cost of the oven (~ 3,000 USD), balance (~ 1,000 to 3,000 USD) and computer (~ 1,000 USD). The cost for the balance can be lower for heavier samples where two decimals are enough; whereas, for lighter samples a four-decimal balance will be preferred.

368 USD), which is comparable to the proposed method on determining the weight of microplas-  
369 tics. The level of expertise required was estimated based on the fact that visual sorting of  
370 microplastic particles requires a certain amount of training; whereas, drying, weighing, and  
371 running a trained algorithm are all relatively straightforward tasks. Finally, our method  
372 requires between 10% to 40% the time of traditional methods. It can achieve better MPE  
373 than manual counting results reported by Isobe *et al.* for heavier particles in a complicated  
374 sample, and our MPE for lighter particles can be reduced to below 30% when the sample  
375 size increases.<sup>23</sup> In conclusion, by using machine learning to predict microplastic counts, we  
376 can reduce the experimental labour by 50% to 80%, without significantly increasing the cost  
377 or sacrificing reliability.

378 We now discuss three hypothetical use cases. Graphical representations can be found in  
379 Section S6. Depending on the context, additional features that might influence the count  
380 can be added. When assessing environmental contamination by microplastics, these features  
381 can include time of the year, season, longitude, latitude, distance to the nearest urban or  
382 industrial center, and so on. We note, however, that a feature will only improve performance  
383 if it is actually related to the count. For example, location may be particularly useful if some  
384 of sampling sites are near a highway.

385 In the first hypothetical use case, microplastic contamination is monitored at a single  
386 site, e.g., a lake, over the course of a year (Figure S17). If one collects several samples  
387 each month, one sample from each collection period can be counted using the traditional  
388 method and verified by spectroscopic analysis to form a training dataset. This training data  
389 would be used to train KRR. The rest of the samples would be sieved and weighed, and the  
390 microplastic counts would be estimated using the trained KRR. Here, the time of the year  
391 could be added as an extra feature if it influences microplastic counts.

392 In the second hypothetical use case, one wishes to monitor a large geographical area  
393 (Figure S18). One could perform traditional counting on samples from a small number of  
394 locations, use this data to train KRR, and then use the algorithm to quickly estimate counts

395 at a larger number of remaining locations. In this case, location might be included as an  
396 additional feature in the training data.

397 The third hypothetical use case involves samples with homogeneous composition, e.g.,  
398 atmospheric or washing machine effluent samples composed of mainly fibers (Table S5).  
399 As demonstrated before, KRR would perform better if the sample contains only one type  
400 of particle. Ideally, one would weigh all collected samples and choose from them several  
401 samples that depict the range of the measurements, count them using traditional manual  
402 method and use them to train KRR. The rest of the samples can be predicted using the  
403 trained model.

## 404 4. Conclusion

405 This study demonstrated the possibility of using regression methods for microplastic quan-  
406 tification based on sieved weight measurements. The method works better for pure samples  
407 made of one type of particles, or for heavier particles in mixed samples. Note that these  
408 results are based on a set of mass measurements, and that further improvement may be  
409 attainable by incorporating features such as geographic location and time of measurement.  
410 This method is faster and easier than current visual quantification methods. As shown in  
411 Figure 4, the new method takes less than half of the time needed for manual counting, and  
412 it only involves fundamental laboratory equipment. We therefore believe that it could be  
413 widely useful to both experts and non-experts. The accuracy of predictions from KRR are  
414 better than that attained by human counting, suggesting that this approach has the po-  
415 tential to accelerate experimental research on this topic without sacrificing the quality of  
416 microplastic quantification.

## 417 Acknowledgements and Data Availability

418 This research was supported by the University of Toronto and the Natural Sciences and  
419 Engineering Research Council of Canada, Canada Research Chairs, grant no. 950-230892.  
420 Raw data were generated at the University of Toronto. Derived data supporting the findings  
421 of this study are available from the corresponding author E.P. on request.

## References

- (1) Moore, C. J. Synthetic polymers in the marine environment: A rapidly increasing, long-term threat. *Environmental Research* **2008**, *108*, 131 – 139, The Plastic World.
- (2) Wright, S. L.; Ulke, J.; Font, A.; Chan, K. L.; Kelly, F. J. Atmospheric microplastic deposition in an urban environment and an evaluation of transport. *Environment International* **2020**, *136*, 105411.
- (3) Liu, F.; Olesen, K. B.; Borregaard, A. R.; Vollertsen, J. Microplastics in urban and highway stormwater retention ponds. *Science of the Total Environment* **2019**, *671*, 992–1000.
- (4) Allen, S.; Allen, D.; Phoenix, V. R.; Le Roux, G.; Durántez Jiménez, P.; Simonneau, A.; Binet, S.; Galop, D. Atmospheric transport and deposition of microplastics in a remote mountain catchment. *Nature Geoscience* **2019**, *12*, 339–344.
- (5) González-Pleiter, M.; Velázquez, D.; Edo, C.; Carretero, O.; Gago, J.; Barón-Sola, Á.; Hernández, L. E.; Yousef, I.; Quesada, A.; Leganés, F.; Rosal, R.; Fernández-Piñas, F. Fibers spreading worldwide: Microplastics and other anthropogenic litter in an Arctic freshwater lake. *Science of the Total Environment* **2020**, *722*, 137904.
- (6) Bucci, K.; Tulio, M.; Rochman, C. M. What is known and unknown about the effects

- of plastic pollution: A meta-analysis and systematic review. *Ecological Applications* **2020**, *30*.
- (7) Rochman, C. M. et al. Rethinking microplastics as a diverse contaminant suite. *Environmental Toxicology and Chemistry* **2019**, *38*, 703–711.
- (8) Hidalgo-Ruz, V.; Gutow, L.; Thompson, R. C.; Thiel, M. Microplastics in the Marine Environment: A Review of the Methods Used for Identification and Quantification. *Environmental Science & Technology* **2012**, *46*, 3060–3075, PMID: 22321064.
- (9) Li, J.; Liu, H.; Paul Chen, J. Microplastics in freshwater systems: A review on occurrence, environmental effects, and methods for microplastics detection. *Water Research* **2018**, *137*, 362–374.
- (10) Methods for sampling and detection of microplastics in water and sediment: A critical review. *TrAC - Trends in Analytical Chemistry* **2019**, *110*, 150–159.
- (11) Eerkes-Medrano, D.; Thompson, R. In *Microplastic Contamination in Aquatic Environments*; Zeng, E. Y., Ed.; Elsevier, 2018; pp 95 – 132.
- (12) Sun, J.; Dai, X.; Wang, Q.; van Loosdrecht, M. C.; Ni, B.-J. Microplastics in wastewater treatment plants: Detection, occurrence and removal. *Water Research* **2019**, *152*, 21 – 37.
- (13) Silva, A. B.; Bastos, A. S.; Justino, C. I.; da Costa, J. P.; Duarte, A. C.; Rocha-Santos, T. A. Microplastics in the environment: Challenges in analytical chemistry - A review. *Analytica Chimica Acta* **2018**, *1017*, 1–19.
- (14) Andrady, A. L. Microplastics in the marine environment. *Marine Pollution Bulletin* **2011**, *62*, 1596 – 1605.
- (15) Mai, L.; Bao, L.-J.; Shi, L.; Wong, C. S.; Zeng, E. Y. A review of methods for measuring



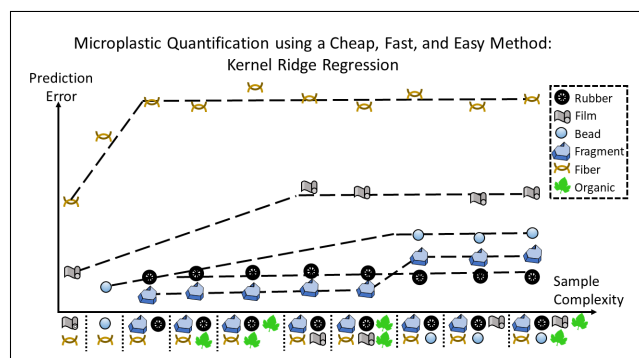
- microplastics in aquatic environments. *Environmental Science and Pollution Research* **2018**, *25*, 11319–11332.
- (16) Primpke, S.; Lorenz, C.; Rascher-Friesenhausen, R.; Gerdtts, G. An automated approach for microplastics analysis using focal plane array (FPA) FTIR microscopy and image analysis. *Anal. Methods* **2017**, *9*, 1499–1511.
- (17) Primpke, S.; Dias, P. A.; Gerdtts, G. Automated identification and quantification of microfibrils and microplastics. *Analytical Methods* **2019**, *11*, 2138–2147.
- (18) Primpke, S.; Cross, R. K.; Mintenig, S. M.; Simon, M.; Vianello, A.; Gerdtts, G.; Vollertsen, J. Toward the Systematic Identification of Microplastics in the Environment: Evaluation of a New Independent Software Tool (siMPle) for Spectroscopic Analysis. *Applied Spectroscopy* **2020**, *74*, 1127–1138.
- (19) Michel, A. P. M.; Morrison, A. E.; Preston, V. L.; Marx, C. T.; Colson, B. C.; White, H. K. Rapid Identification of Marine Plastic Debris via Spectroscopic Techniques and Machine Learning Classifiers. *Environmental Science & Technology* **2020**, *54*, 10630–10637, PMID: 32697577.
- (20) Yurtsever, M.; Yurtsever, U. Use of a convolutional neural network for the classification of microbeads in urban wastewater. *Chemosphere* **2019**, *216*, 271 – 280.
- (21) Lorenzo-Navarro, J.; Castrillón Santana, M.; Gómez, M.; Herrera, A.; Marín-Reyes, P. Automatic Counting and Classification of Microplastic Particles. 2018.
- (22) Koelmans, A. A.; Redondo-Hasselerharm, P. E.; Mohamed Nor, N. H.; Kooi, M. Solving the Nonalignment of Methods and Approaches Used in Microplastic Research to Consistently Characterize Risk. *Environmental Science & Technology* **2020**,
- (23) Isobe, A. et al. An interlaboratory comparison exercise for the determination of mi-

- croplastics in standard sample bottles. *Marine Pollution Bulletin* **2019**, *146*, 831 – 837.
- (24) McIlwraith, H. K.; Lin, J.; Erdle, L. M.; Mallos, N.; Diamond, M. L.; Rochman, C. M. Capturing microfibers – marketed technologies reduce microfiber emissions from washing machines. *Marine Pollution Bulletin* **2019**, *139*, 40 – 45.
- (25) Lavers, J. L.; Dicks, L.; Dicks, M. R.; Finger, A. Significant plastic accumulation on the Cocos (Keeling) Islands, Australia. *Scientific Reports* **2019**, *9*, 1–9.
- (26) Puskic, P. S.; Lavers, J. L.; Adams, L. R.; Bond, A. L. Ingested plastic and trace element concentrations in Short-tailed Shearwaters (*Ardenna tenuirostris*). *Marine Pollution Bulletin* **2020**, *155*, 111143.
- (27) Corcoran, P. L.; de Haan Ward, J.; Arturo, I. A.; Belontz, S. L.; Moore, T.; Hill-Svehla, C. M.; Robertson, K.; Wood, K.; Jazvac, K. A comprehensive investigation of industrial plastic pellets on beaches across the Laurentian Great Lakes and the factors governing their distribution. *Science of The Total Environment* **2020**, *747*, 141227.
- (28) Arturo, I. Plastic debris in the Laurentian Great Lakes System, North America: Analysis of types, abundances, and sources. *Electronic Thesis and Dissertation Repository* **2021**, *7758*.
- (29) Smyth, K.; Drake, J.; Li, Y.; Rochman, C.; Van Seters, T.; Passeport, E. Bioretention cells remove microplastics from urban stormwater. *Water Research* **2021**, *191*.
- (30) Olesen, K. B.; Stephansen, D. A.; van Alst, N.; Vollertsen, J. Microplastics in a stormwater pond. *Water (Switzerland)* **2019**, *11*.
- (31) Lima, A.; Costa, M.; Barletta, M. Distribution patterns of microplastics within the plankton of a tropical estuary. *Environmental Research* **2014**, *132*, 146 – 155.

- (32) Burns, E. E.; Boxall, A. B. Microplastics in the aquatic environment: Evidence for or against adverse impacts and major knowledge gaps. *Environmental Toxicology and Chemistry* **2018**, *37*, 2776–2796.
- (33) Zhang, Y.; Kang, S.; Allen, S.; Allen, D.; Gao, T.; Sillanpää, M. Atmospheric microplastics: A review on current status and perspectives. *Earth-Science Reviews* **2020**, *203*, 103118.
- (34) Murphy, K. P. *Machine Learning: A Probabilistic Perspective*; The MIT Press, 2012.
- (35) El-Dereny, M.; Rashwan, N. Solving multicollinearity problem using ridge regression models. *International Journal of Contemporary Mathematical Sciences* **2011**, *6*.
- (36) Pedregosa, F.; Varoquaux, G.; Gramfort, A.; Michel, V.; Thirion, B.; Grisel, O.; Blondel, M.; Prettenhofer, P.; Weiss, R.; Dubourg, V., et al. Scikit-learn: Machine learning in Python. *Journal of machine learning research* **2011**, *12*, 2825–2830.
- (37) Eriksen, M.; Mason, S.; Wilson, S.; Box, C.; Zellers, A.; Edwards, W.; Farley, H.; Amato, S. Microplastic pollution in the surface waters of the Laurentian Great Lakes. *Marine Pollution Bulletin* **2013**, *77*, 177 – 182.
- (38) Lenz, R.; Enders, K.; Stedmon, C. A.; Mackenzie, D. M.; Nielsen, T. G. A critical assessment of visual identification of marine microplastic using Raman spectroscopy for analysis improvement. *Marine Pollution Bulletin* **2015**, *100*, 82 – 91.
- (39) Cawley, G.; Talbot, N.; Chapelle, O. Estimating Predictive Variances with Kernel Ridge Regression. 2005; pp 56–77.
- (40) Löder, M. G. J.; Gerdtz, G. In *Marine Anthropogenic Litter*; Bergmann, M., Gutow, L., Klages, M., Eds.; Springer International Publishing: Cham, 2015; pp 201–227.
- (41) Jenkins, D.; Quintana-Ascencio, P. A solution to minimum sample size for regressions. *PLOS ONE* **2020**, *15*, e0229345.

- (42) Rodrigues, S.; R. Almeida, C. M.; Ramos, S. Adaptation of a laboratory protocol to quantity microplastics contamination in estuarine waters. *MethodsX* **2019**, *6*, 740 – 749.

## Graphical TOC Entry



## Highlights

- Predict microplastic counts based on aggregate mass measurements
- Better prediction accuracy for homogeneous samples and heavier particles
- Mass measurement error did not affect prediction accuracy
- Faster, cheaper and easier microplastic quantification method