

# 1 **Centrifuge+: improving metagenomic analysis upon Centrifuge**

2 Junfeng Liu<sup>1,2†\*</sup>, Yunran Ma<sup>1,2†</sup>, Yong Ren<sup>1,2</sup> and Hao Guo<sup>1,2\*</sup>

3 <sup>1</sup>State Key Laboratory of Translational Medicine and Innovative Drug Development,  
4 Jiangsu Simcere Diagnostics Co., Ltd., Nanjing, China

5 <sup>2</sup>Nanjing Simcere Medical Laboratory Science Co., Ltd., Nanjing, China

6 \*To whom correspondence should be addressed.

7 <sup>†</sup>The authors wish it to be known that, in their opinion, the first two authors should be  
8 regarded as Joint First Authors.

## 9 **Abstract**

10 **Summary:** Accurate abundance estimation of species is essential for metagenomic  
11 analysis. Although many methods have been developed for classification of  
12 metagenomic data and abundance estimation of species, the abundance estimation of  
13 species remains challenging due to the ambiguous reads that align equally well to  
14 more than one genome. Here, we present Centrifuge+, which introduces unique  
15 mapping rate to describe the influence of similarities among species in the reference  
16 database when analyzing ambiguous reads. In contrast to the popular Centrifuge,  
17 Centrifuge+ improved the accuracy of abundance estimation on simulated reads from  
18 4278 complete prokaryotic genomes.

19 **Availability and implementation:** The source code is available at  
20 <https://github.com/mNGSmethods/Centrifugep>.

21 **Contact:** [h.guo@foxmail.com](mailto:h.guo@foxmail.com) or [jlsjlf0101@126.com](mailto:jlsjlf0101@126.com)

22 **Supplementary information:** Supplementary data are available at *Bioinformatics*

23 online.

## 24 **1 Introduction**

25 Metagenomic sequencing has provided great improvements in microbiome analysis  
26 by metagenomic experiments that can be broadly categorized as either microbiome  
27 experiments or pathogen identification experiments (Knight *et al.*, 2018; Lu *et al.*,  
28 2022). In microbiome experiments, researchers focus on describing what is present in  
29 a given sample. For pathogen identification experiments, the focus of researcher is  
30 identifying one or few pathogenic microbes. In order to achieve the goal of  
31 metagenomic experiments, estimating the abundance of the species in a given sample  
32 becomes very important in metagenomic analysis. However, the ambiguous reads that  
33 align equally well to more than one genome make challenge for the abundance  
34 estimation of the species because it is very difficult to identify the taxon of ambiguous  
35 reads. There are two reasons for causing the ambiguous reads. The first reason is that  
36 closely related species are present in a given sample. The second reason is because of  
37 the nearly identical genomes in a reference database that is used for identifying the  
38 taxon of each read. In order to overcome the challenge of ambiguous reads, a separate  
39 abundance estimation algorithm is necessary for most metagenomic classification  
40 tools. To counter the ambiguous reads caused by closely related species in the same  
41 sample, Kim *et al.* (2016) defined a statistical model in the popular metagenomic  
42 classification tool Centrifuge (Kim *et al.*, 2016) and used it to estimate the abundance  
43 of species through an Expectation-Maximization (EM) algorithm. In the statistical  
44 model of Centrifuge, the probability is only depended on the abundance of species  $j$

45 and the length of the genomes of species  $j$  when the ambiguous read  $i$  is classified to  
46 species  $j$ . However, for the ambiguous reads caused by the nearly identical genomes  
47 in a reference database, the probability is also decided by the similarities between  
48 species  $j$  and the other species in the reference database if the ambiguous read  $i$  is  
49 classified to species  $j$ . Although the similarities among species in the reference  
50 database have been considered in the statistical model of Bracken (Lu *et al.*, 2017),  
51 which was developed to estimate species abundance in conjunction with Kraken  
52 (Wood and Salzberg, 2014), the statistical model of Bracken can be only used to  
53 analyze Kraken classification results and requires generating simulation data to  
54 estimate species abundance.

55 To address the above limitation, we introduce Centrifuge+, which modified the  
56 statistical model of Centrifuge and improved metagenomic analysis. In the modified  
57 statistical model, the influence of similarities among species in the reference database  
58 is described by unique mapping rate when analyzing the ambiguous reads. In addition,  
59 we use the modified statistical model to estimate species abundance through an  
60 Expectation-Maximization (EM) algorithm.

## 61 **2 Implementation**

62 Centrifuge+ is based on Centrifuge with the same methods of reference database  
63 sequence compression and classification of microbial sequences, but is different from  
64 Centrifuge on the statistical model, which considers the influence of similarities  
65 among species in the reference database on estimating species abundance. In order to  
66 implement the modified statistical model, we modified Centrifuge developed by Kim

67 et al. (2016) under the terms of the GNU General Public License and named it as  
68 Centrifuge+.

## 69 **2.1 The modified statistical model**

70 Similar to Centrifuge, the likelihood function is defined as follows:

$$71 \quad L(\partial|C) = \prod_{i=1}^R \sum_{j=1}^S \frac{\partial_j l_j}{\sum_{k=1}^S \partial_k l_k} C_{ij}$$

$$72 \quad C_{ij} = \begin{cases} p_j & \text{if read } i \text{ is uniquely mapped to species } j \\ 1 - p_j & \text{if read } i \text{ is multiply mapped to species } j \text{ and the other species} \\ 0 & \text{if read } i \text{ is not mapped to species } j \end{cases}$$

73 where  $R$  is the number of the reads,  $S$  is the number of species,  $\partial_j$  is the abundance

74 of species  $j$  and  $\sum_{j=1}^S \partial_j = 1$ ,  $l_j$  is the average length of genomes of species  $j$ , and  $p_j$  is

75 the unique mapping rate of species  $j$ . For species  $j$ , we count the number of reads that

76 are uniquely mapped to it,  $m$ . If the number of reads that can be classified to species  $j$

77 is  $n$ , the unique mapping rate of species  $j$  is  $m/n$ .

78 In the modified statistical model, we introduced the unique mapping rate ( $p_j$ ) to

79 describe the influence of similarities between species  $j$  and the other species in the

80 reference database when assigning a value to  $C_{ij}$ . However, in the statistical model of

81 Centrifuge, the value of  $C_{ij}$  is only 1 or 0 according to whether read  $i$  is mapped to

82 species  $j$ .

## 83 **2.2 Abundance analysis**

84 To estimate species abundance, the following EM procedure is implemented in

85 Centrifuge+.

86 Initialization step (I-step): the initial value of  $\partial_j$  is  $1/S$ .

87 Expectation step (E-step):

$$n_j = \frac{\sum_{i=1}^R \partial_j l_j C_{ij}}{\sum_{k=1}^S \partial_k l_k C_{ik}}$$

88

89 where  $n_j$  is the estimated number of reads mapped to species  $j$ .

90 Maximization step (M-step):

$$\partial'_j = \frac{n_j / l_j}{\sum_{k=1}^S n_k / l_k}$$

91

92 where  $\partial'_j$  is the updated estimation of species  $j$ 's abundance and used in the next  
93 iteration as  $\partial_j$ .

94 Centrifuge+ repeats E-step and M-step until  $\sum_{j=1}^S |\partial_j - \partial'_j| < 10^{-10}$ . The above EM

95 procedure is also implemented in Centrifuge except for E-step.

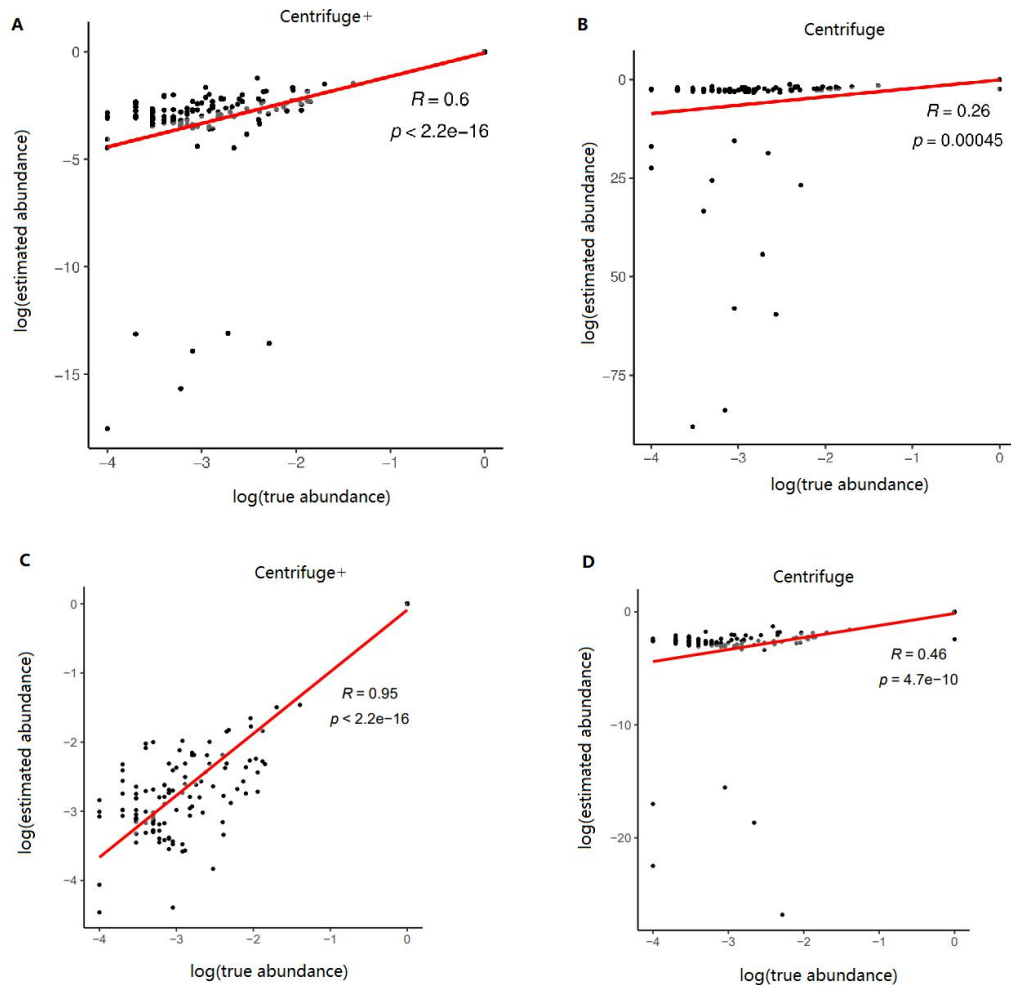
### 96 **3 Results and discussion**

97 We compared Centrifuge+ and Centrifuge by assessing the match between the  
98 estimated abundance and the true abundance distribution of genomes in the simulated  
99 reads at the species level. The simulated read data set was created from the 4278  
100 complete prokaryotic genomes in RefSeq (Pruitt *et al.*, 2014) by Kim *et al.* (2016).  
101 They used the Mason simulator (Luke *et al.*, 2005) to generate 10 million 100-bp  
102 reads and the resulting file was named `bacteria_sim10M.fa` (Kim *et al.*, 2016). Then,  
103 they randomly down-sampled the datasets to 10 thousand reads (`bacteria_sim10K.fa`)  
104 without replacement. We used this dataset for the performance comparison of  
105 Centrifuge+ and Centrifuge (Supplementary Materials). Pearson's correlation

106 coefficient between the true abundance and the estimated abundance of Centrifuge+  
107 was 0.6 at the species level based on 10 thousand simulated reads (Fig. 1A). However,  
108 Pearson's correlation coefficient was only 0.26 at the species level when comparing  
109 the true abundance and the estimated abundance of Centrifuge (Fig. 1B). If the top  
110 five percent worst abundance estimates were omitted, the correlation coefficient of  
111 Centrifuge+ can improve to 0.95 (Fig. 1C). But, the correlation coefficient of  
112 Centrifuge only can improve to 0.46 when omitting the top five percent worst  
113 abundance estimates (Fig. 1D). The above results show that the abundance estimates  
114 of Centrifuge+ are more closely matched to the true abundance than Centrifuge. Even  
115 for the top five percent worst abundance estimates, Centrifuge+ is still significantly  
116 better than Centrifuge (Supplementary Fig. S1). Moreover, the more accurate  
117 abundance estimates make Centrifuge+ to have a higher recall (75.86% VS 65.51%),  
118 which is the proportion of true positive species divided by the number of distinct  
119 species actually in the sample (Supplementary Table S1). Though Centrifuge+ was a  
120 little lower than Centrifuge on the the precision (98.33% VS 100%), which is the  
121 proportion of true positive species identified in the sample divided by the number of  
122 total identified species, Centrifuge+ has a higer F1 score (0.86 VS 0.79) that is the  
123 harmonic mean of recall and precision (Supplementary Table S1).

124 When describing the probability of observed read, the statistical model in  
125 Centrifuge does not distinguish between species in the processing of unique and  
126 multiple mapping reads, that is, no matter to which species the read is classified to,  
127 the unique mapping rates of different species are the same. Centrifuge's above

128 processing method implies the following assumption: the probability of the  
129 occurrence of unique and multiple mapping reads of different species is only  
130 determined by species abundance. However, due to the influence of reference genome  
131 similarity, the probability of unique mapping reads and multiple mapping reads of  
132 different species will also be different. For example, an observation sample contains  
133 20 reads with two species, A and B, whose abundance ratio is 1:1 and genome length  
134 ratio is 1:1. If 6 reads are unique mapped to species A, 4 reads are unique mapped to  
135 species B, and the remaining 10 reads are mapped to both species A and species B,  
136 then according to the statistical model of Centrifuge, in which the value of  $C_{ij}$  is only  
137 1 or 0 according to whether read  $i$  is mapped to species  $j$ , the estimated abundances of  
138 species A and species B are 0.6 and 0.4 respectively. When the influence of reference  
139 genome similarity is considered in the statistical model of Centrifuge+ by introducing  
140 the unique mapping rate, the estimated abundances of species A and species B are  
141 0.57 and 0.43 respectively and more closer to true abundances. Therefore,  
142 Centrifuge+ can improve the accuracy of abundance estimates than Centrifuge  
143 according to the above discussion.



144

145 **Fig. 1.** Comparison of log-scaled true abundance and estimated abundance at species level based on 10

146 thousand simulated reads.  $R$  and  $p$  are Pearson's correlation coefficient and  $p$ -value, respectively. (A)

147 Comparison of log-scaled true abundance and Centrifuge+ abundance estimates; (B) Comparison of

148 log-scaled true abundance and Centrifuge abundance estimates; (C) Comparison of log-scaled true

149 abundance and Centrifuge+ abundance estimates when the five percent worst abundance estimates of

150 Centrifuge+ were omitted; (D) Comparison of log-scaled true abundance and Centrifuge abundance

151 estimates when the five percent worst abundance estimates of Centrifuge were omitted.

## 152 **4 Conclusion**

153 Because Centrifuge can analyze not only short reads, but also long reads, Centrifuge

154 has a wide range of application scenarios, such as Pavian (Breitwieser and Salzberg,



155 2020) and minoTour (Munro *et al.*, 2022). Centrifuge is especially applied for ONT  
156 shotgun sequencing analysis and is now included as a step in WIMP, which is a  
157 quantitative analysis tool for real-time species identification based on the MinIon  
158 released by Oxford Nanopore Technologies. In contrast to Centrifuge, Centrifuge+  
159 improved the accuracy of abundance estimates by modifying the statistical model in  
160 Centrifuge. The more accurate abundance estimates will be benefit to improve the  
161 precision-recall analysis for species identification. Hence, Centrifuge+ will be more  
162 widely applied for metagenomic analysis, particularly for real-time species  
163 identification.

#### 164 **Funding**

165 This research was supported by China's National Key R&D Program (Grant No.  
166 2018YFE0102100 and 2022YFC2505100) and the Collaborative Innovation Major  
167 Project of Zhengzhou (Grant No. 20XTZX08017).

168 *Conflict of Interest:* none declared.

#### 169 **References**

- 170 Breitwieser F P, Salzberg S L. Pavian: interactive analysis of metagenomics data for microbiome  
171 studies and pathogen identification[J]. *Bioinformatics*, 2020, 36(4): 1303-1304.
- 172 Kim D, Song L, Breitwieser F P, et al. Centrifuge: rapid and sensitive classification of metagenomic  
173 sequences[J]. *Genome research*, 2016, 26(12): 1721-1729.
- 174 Knight R, Vrbanac A, Taylor B C, et al. Best practices for analysing microbiomes[J]. *Nature Reviews*  
175 *Microbiology*, 2018, 16(7): 410-422.
- 176 Luke S, Cioffi-Revilla C, Panait L, et al. Mason: A multiagent simulation environment[J]. *Simulation*,  
177 2005, 81(7): 517-527.
- 178 Lu J, Breitwieser F P, Thielen P, et al. Bracken: estimating species abundance in metagenomics

- 179 data[J]. PeerJ Computer Science, 2017, 3: e104.
- 180 Lu J, Rincon N, Wood D E, et al. Metagenome analysis using the Kraken software suite[J]. Nature  
181 protocols, 2022: 1-25.
- 182 Munro R, Santos R, Payne A, et al. minoTour, real-time monitoring and analysis for nanopore  
183 sequencers[J]. Bioinformatics, 2022, 38(4): 1133-1135.
- 184 Pruitt K D, Brown G R, Hiatt S M, et al. RefSeq: an update on mammalian reference sequences[J].  
185 Nucleic acids research, 2014, 42(D1): D756-D763.
- 186 Wood D E, Salzberg S L. Kraken: ultrafast metagenomic sequence classification using exact  
187 alignments[J]. Genome biology, 2014, 15(3): 1-12.
- 188