

1 **Phylogenomic disentangling of the *Bifidobacterium longum* subsp. *Infantis* taxon**

2 Running title: Phylogenomic dissection of *Bifidobacterium longum*

3
4 Key words: bifidobacteria, comparative genomics, longum, infantis, HGT

5
6 Chiara Tarracchini¹, Christian Milani^{1,2}, Gabriele Andrea Lugli¹, Leonardo Mancabelli¹, Federico
7 Fontana^{1,3}, Giulia Alessandri¹, Giulia Longhi^{1,3}, Rosaria Anzalone³, Alice Viappiani³, Francesca
8 Turrone^{1,2} Douwe van Sinderen⁴, and Marco Ventura^{1,2}

9
10 Laboratory of Probiogenomics, Department of Chemistry, Life Sciences, and Environmental
11 Sustainability, University of Parma, Parma, Italy¹; Microbiome Research Hub, University of Parma,
12 Parma, Italy²; GenProbio srl, Parma, Italy³; APC Microbiome Ireland and School of Microbiology,
13 Bioscience Institute, National University of Ireland, Cork, Ireland⁴

14
15
16
17
18
19
20
21
22
23
24
25
26
27
28 Correspondence. Mailing address for Marco Ventura Laboratory of Probiogenomics, Department of
29 Chemistry, Life Sciences, and Environmental Sustainability, University of Parma, Parco Area delle
30 Scienze 11a, 43124 Parma, Italy. Phone: ++39-521-905666. Fax: ++39-521-905604. E-mail:
31 marco.ventura@unipr.it

33 **Supplementary Text**

34 **Pan-genome and core genome of *B. longum* species.** Recently, pan-genome computation has allowed
35 to investigate genomic differences between a given (bifido)bacterial taxa as well as their evolutionary
36 development and phylogenomic relationships (1-3). In the framework of a species-level genomic
37 comparative analysis, the available genomes of *B. longum* were sent to pan-genome analysis, revealing
38 a total of 22,591 Clusters of Orthologous Groups (COGs), defined as CDSs showing identity > 50% with
39 alignment coverage > 80%. The pan-genome curve was built plotting the total number of COGs as a
40 function of the 272 *B. longum* genomes, showing an asymptotic behavior with a growth rate progressively
41 decreasing (Figure S1). More precisely, the number of new genes added by sequential inclusion of
42 genomes started from an average of 413.4 at the first three iterations, decreasing to an average of 49.7 at
43 the last three iterations. This trend is indicative of a pan-genome power trend line that has not yet reached
44 the plateau, suggesting that any additional genome included in the analysis will result in slight increases
45 of the pangenome size.

46 Moreover, the pan-genome analysis also revealed the core genome of this species as well as the unique
47 genes repertoire of each analyzed strain (Figure S1). A total of 510 COGs of *B. longum* pan-genome
48 were shared by all the strains, thus representing the core genome of this species. The Truly Unique Genes
49 (TUGs) for each *B. longum* strain, i.e., the genes present in just one strain, were also identified, revealing
50 a number of TUGs ranging from 296 genes for *B. longum* subsp. *longum* 1897B to 14 genes for two *B.*
51 *longum* subsp. *longum* strains, i.e., 9 and MC-42, with an average of 48.5 TUGs per genome. These
52 obtained average number of TUGs resulted comparable with that previously evidenced for
53 *Bifidobacterium pseudolongum*, *Bifidobacterium dentium*, showing an average of 41 and 60 TUGs per
54 genome, respectively (4, 5). These findings, coupled with the prediction of a pan-genome curve tending
55 to a plateau, suggested that the evolutionary pathway undertaken by this species has not led to obtaining
56 a high grade of strain-specific genomic variability, probably because of a high grade of specialization

57 toward a settlement in a limited range of ecological niches, e.g. mammalian gastrointestinal tract.
58 Nevertheless, the concomitant relatively small number of core genes, revealed a potential high grade of
59 intra-specific genomic variability, compared to that previously observed in members of *Bifidobacterium*
60 *bifidum* and *Bifidobacterium breve*, showing 1295 and 1307 of COGs, respectively (1, 6).

61

62 **Phylogenetic analyses the *B. longum* taxon.** To further explore the overall genomic variability of *B.*
63 *longum* genome we observed that the pairwise percent Average Nucleotide Identity (ANI) ranging from
64 94.2 % to 98.9 (Table S2). Notably, previous *Bifidobacterium* phylogenomic studies showed that an ANI
65 threshold value of 94 % properly discriminates between bifidobacterial species (3, 7) as also reported by
66 literature for other genera belonging to *Bifidobacteriaceae* family (8). Accordingly, the ANI values above
67 94.1 % observed from this phylogenomic analysis revealed that all genome sequences correctly fall into
68 the same species, i.e., *B. longum*. Nonetheless, based on the ANI matrix encompassing all the 272
69 genomes (Table S2), it was possible to identify three subgroups corresponding to the three-known
70 subspecies of *B. longum*, within which the observed ANI values ranged from 96.3 % to 98.9 %, with an
71 average value of 98.26% (Table S2). Furthermore, in order to explore the phylogenetic relationships
72 between the strains of this species, we computed a phylogenetic tree based on the aminoacidic sequence
73 alignment of the 510 COGs constituting the core genome of this species (Figure S2). Notably, previous
74 literature showed that this approach allows precise high-resolution reconstruction of the phylogenetic
75 relationship between both distant and closely related taxa/strains (3, 8, 9). Due to the high number of
76 analyzed genomes belonging to the *B. longum* subsp. *longum* subspecies, we decided to generate an
77 additional tree encompassing just a pool of 21 representative genomes of this taxon in order to obtain a
78 better and more clear graphical visualization of the whole *B. longum* species phylogeny (Figure 1). In
79 particular, the selection of *B. longum* subsp. *longum* strains displayed in Figure 1 include the type strain,
80 i.e., DSM20219, as well as 12 genomes that had shown the lowest ANI values within the *B. longum*

81 subsp. *longum* subgroup respect to the type strain in order to maximize genomic variability among strains
82 selected for this subspecies (Table S2). Moreover, all the 11 strains isolated in the framework of this
83 study were also included in Figure 1. As expected, the resulting *B. longum*-based phylogenetic tree
84 revealed the presence of three main clades (Figure 1; Figure S2). These three clades encompassed
85 respectively 251, 11 and 10 *B. longum* genomes, constituting the *B. longum* subsp. *longum* taxonomic
86 group (*Bll*), the *B. longum* subsp. *infantis* taxonomic group (*Bli*) and the *B. longum* subsp. *suis* (*Bls*)
87 taxonomic group, respectively (Figure 1). An in-depth analysis of the tree revealed that four genomes
88 without subspecies classification as well as one that was presumed to belong to the *B. longum* subsp.
89 *longum* subspecies, i.e., JDM301, clustered in the *Bls* clade, thus suggesting a misclassification of this
90 latter strain (Figure 1; Figure S2), consistently with what previously observed through ANI analysis
91 (Table S2). Likewise, CCUG 52486 and 157F strains fall into *Bll* group, despite being previously
92 classified as *B. longum* subsp. *infantis*, thus indicating their mistaken taxonomic classification, also
93 confirmed by ANI analysis (Table S2). Notably, all the subsequent analyses performed in this study were
94 conducted exploiting taxonomic assignments derived from the phylogenomic and ANI analysis.
95 Interpretation of the phylogenomic tree suggests a clear phylogenetic separation between members of *B.*
96 *longum* subsp. *infantis* cluster and the other *B. longum* strains, indicative of earlier speciation respect to
97 *B. longum* subsp. *longum* and *B. longum* subsp. *suis*, which showed closer phylogenetic relationship.
98 Moreover, the phylogenomic-based approach, combined with ANI values assignment, was also exploited
99 to taxonomically classify the 11 newly isolated *B. longum* strains.

100

101 **The Pan- and Core- genome of the *B. longum* subspecies.** Pan-genome reconstruction can contribute
102 to deciphering these evolutionary events, by unveiling the genomic peculiarities as well as the shared
103 genetic traits characterizing a given bacterial taxon (10). In this framework, we separately conducted a
104 subspecies-specific pan-genome analyses, involving the 251 *B. longum* subsp. *longum*, 11 *B. longum*

105 subsp. *infantis* and ten members of *B. longum* subsp. *suis* (Figure S3), These analyses also lead to the
106 definition of the *Bll*-, *Bls*- and *Bli*-Core Genome (CG) as the set of subspecies-specific core genes. In
107 detail, for building the *Bll*- *Bls*- and *Bli*-CG we take into account those COGs shared by at least 85 % of
108 the strains belonging to the same *B. longum* subspecies, and absent in the others. The *B. longum* subsp.
109 *longum* genomes and the members of *B. longum* subsp. *suis* showed an average genome sizes of 2.39 Mb
110 and 2.43 Mb, corresponding to an average of predicted CDS of 1916 and 1947, respectively (Table
111 S1). Moreover, such genome sizes were significantly reduced than those of members of *B. longum* subsp.
112 *infantis* (average of 2.65 Mb corresponding to 2170 CDS per genome) (Table S1) (ANOVA p-value <
113 0.001). Thus, these findings suggested that during evolution *B. longum* subsp. *infantis* taxon may have
114 obtained an increase in its genome size by progressive acquisition of new genetic materials (11).
115 Furthermore, analysis of the pangenome curves obtained for the three subspecies revealed that *B. longum*
116 subsp. *longum* tend toward reaching a plateau (Figure S3), as indicated by an average addition of 41.8
117 genes added to the pangenome in the last three iterations. In contrast, *B. longum* subsp. *suis* and *B. longum*
118 subsp. *infantis* pan-genome showed respectively an average of 191.4 and 99.1 genes added at the last
119 three iterations. These data demonstrated that both *B. longum* subsp. *infantis* and *B. longum* subsp. *suis*
120 were characterized by an open pan-genome (Figure S3), implying that further genomic sequencing efforts
121 will extend our knowledge regarding the genetic variability of these taxa.

122 Following subspecies-specific core genome investigation, a total of 59 genes was retrieved from *Bli*-CG,
123 while 23 and five core genes represented *Bll*-CG and *Bls*-CG, respectively. Notably, 20 of the 59 genes
124 constituting *Bli*-CG as well as two of the five genes forming *Bls*-CG were found within 100 % of the
125 strains belonging to the respective subspecies (Table 1). In contrast, no gene was found to be shared by
126 all the strains constituting the *B. longum* subsp. *longum* subspecies (Table 1), probably due to the higher
127 number of genomes retrieved for this subspecies, of which the majority were available as draft.

128

129 **References**

- 130 1. Lugli GA, Duranti S, Albert K, Mancabelli L, Napoli S, Viappiani A, et al. Unveiling Genomic Diversity
131 among Members of the Species *Bifidobacterium pseudolongum*, a Widely Distributed Gut Commensal of the
132 Animal Kingdom. *Appl Environ Microbiol.* 2019;85(8).
- 133 2. Duranti S, Milani C, Lugli GA, Mancabelli L, Turrone F, Ferrario C, et al. Evaluation of genetic diversity
134 among strains of the human gut commensal *Bifidobacterium adolescentis*. *Sci Rep.* 2016;6:23971.
- 135 3. Lugli GA, Milani C, Turrone F, Duranti S, Mancabelli L, Mangifesta M, et al. Comparative genomic and
136 phylogenomic analyses of the Bifidobacteriaceae family. *BMC Genomics.* 2017;18(1):568.
- 137 4. Bottacini F, O'Connell Motherway M, Kuczynski J, O'Connell KJ, Serafini F, Duranti S, et al. Comparative
138 genomics of the *Bifidobacterium breve* taxon. *BMC Genomics.* 2014;15:170.
- 139 5. Duranti S, Milani C, Lugli GA, Turrone F, Mancabelli L, Sanchez B, et al. Insights from genomes of
140 representatives of the human gut commensal *Bifidobacterium bifidum*. *Environ Microbiol.* 2015;17(7):2515-31.
- 141 6. Lugli GA, Tarracchini C, Alessandri G, Milani C, Mancabelli L, Turrone F, et al. Decoding the Genomic
142 Variability among Members of the *Bifidobacterium dentium* Species. *Microorganisms.* 2020;8(11).
- 143 7. Lugli GA, Milani C, Duranti S, Mancabelli L, Mangifesta M, Turrone F, et al. Tracking the Taxonomy of the
144 Genus *Bifidobacterium* Based on a Phylogenomic Approach. *Appl Environ Microbiol.* 2018;84(4).
- 145 8. Tarracchini C, Lugli GA, Mancabelli L, Milani C, Turrone F, Ventura M. Assessing the Genomic Variability
146 of *Gardnerella vaginalis* through Comparative Genomic Analyses: Evolutionary and Ecological Implications. *Appl*
147 *Environ Microbiol.* 2020;87(1).
- 148 9. Milani C, Lugli GA, Duranti S, Turrone F, Bottacini F, Mangifesta M, et al. Genomic encyclopedia of type
149 strains of the genus *Bifidobacterium*. *Appl Environ Microbiol.* 2014;80(20):6290-302.
- 150 10. Medini D, Donati C, Tettelin H, Massignani V, Rappuoli R. The microbial pan-genome. *Curr Opin Genet*
151 *Dev.* 2005;15(6):589-94.
- 152 11. Koskiniemi S, Sun S, Berg OG, Andersson DI. Selection-driven gene loss in bacteria. *Plos Genet.*
153 2012;8(6):e1002787.

154

155

156

157

158

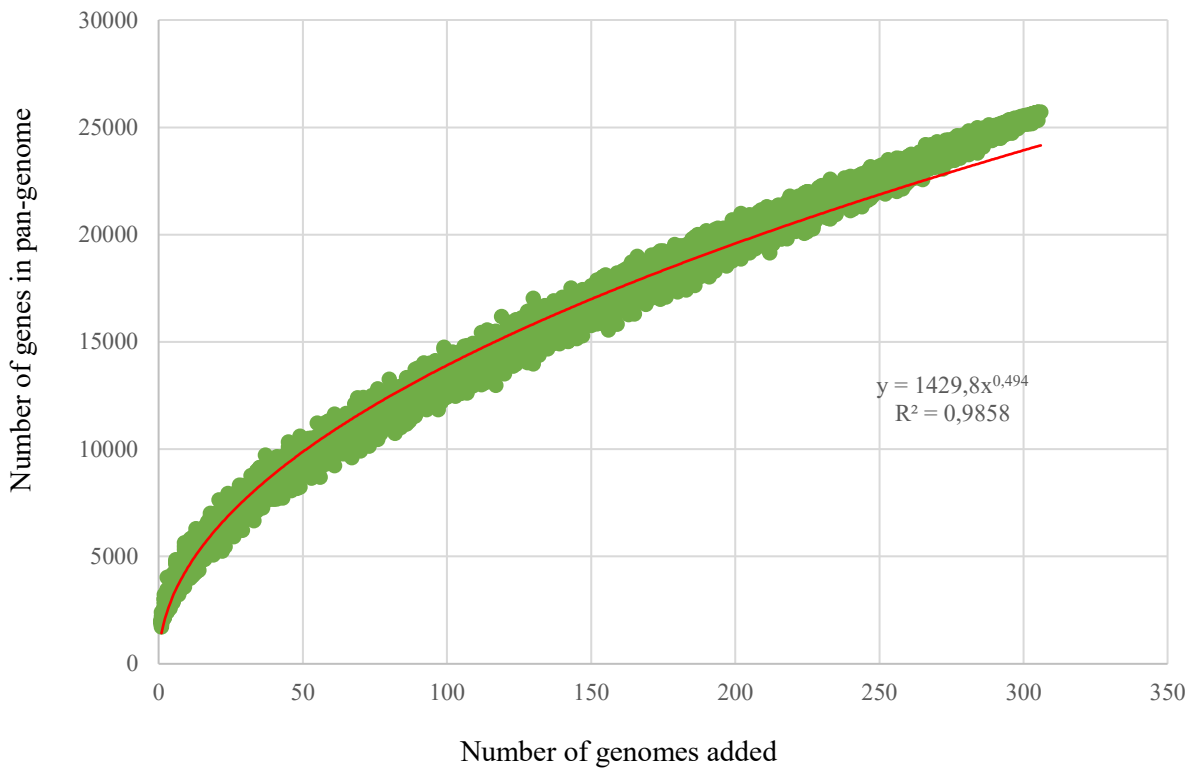
159 **Supplementary figure legends**

160 **Figure S1. Pan-genome of *B. longum* species.** Panel a shows the pangenome size based on the
161 sequential addition of the 289 *B. longum* genomes. Panel b displays the number of core genes (light red),
162 dispensable genes (light blue), and unique genes (purple) identified in the pangenome analysis.

163 **Figure S2. Complete *B. longum*-based phylogenomic tree.** The phylogenomic tree, showing all the
164 289 genomes of *B. longum* included in this study, was based on the concatenation of the 501 *B. longum*
165 core genes and was built through the neighbor-joining method. Bootstrap percentages above 50 are
166 shown at node points, based on 1,000 replicates. Phylogenetic clusters are highlighted with similarly
167 colored branches.

168 **Figure S3. Prediction of the *B. longum* subspecies pangenome.** For each *B. longum* subspecies, Panel
169 a shows the pan-genomes curve representing the variation of pan-genome size resulting from the
170 sequential genomes addition. Panel b represents a pie chart of the number of core genes (red), dispensable
171 genes (light blue) and unique genes (purple).

a)



b)

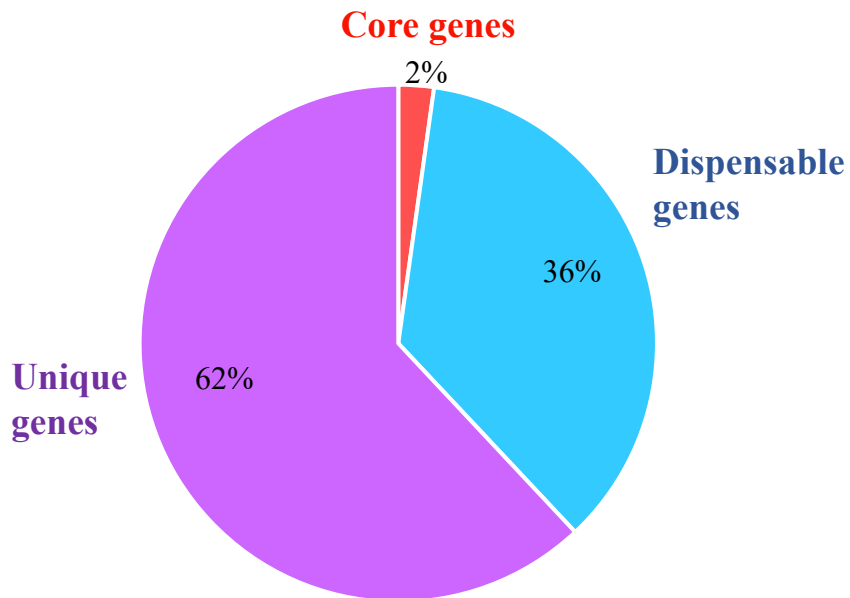
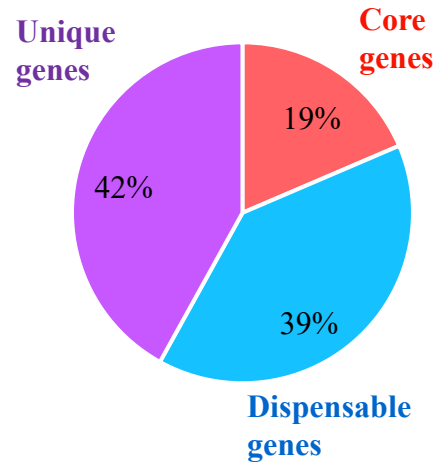
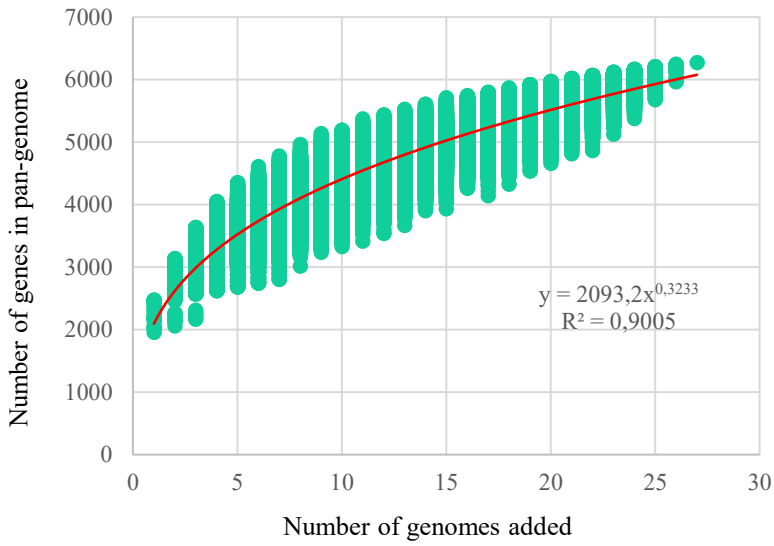
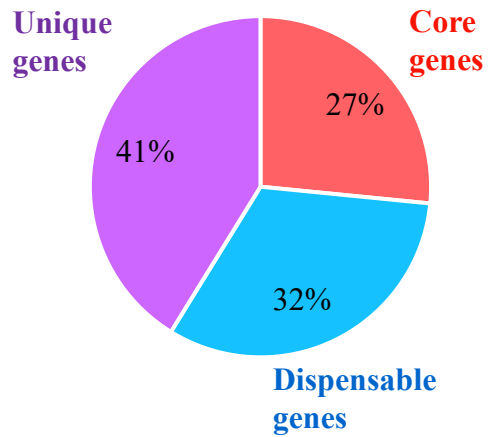
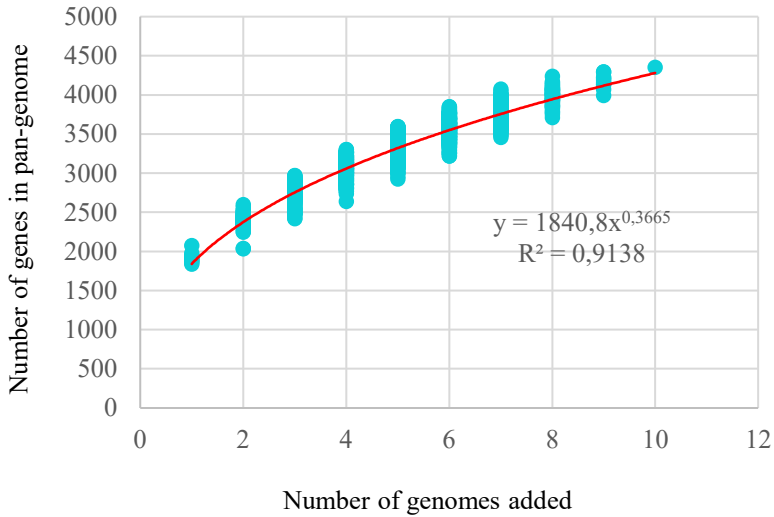


Figure S1

***B. longum* subsp. *infantis* pangenome**



***B. longum* subsp. *suis* pangenome**



***B. longum* subsp. *longum* pangenome**

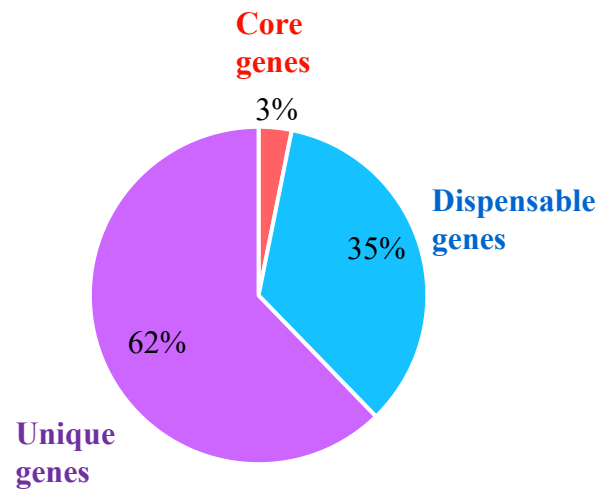
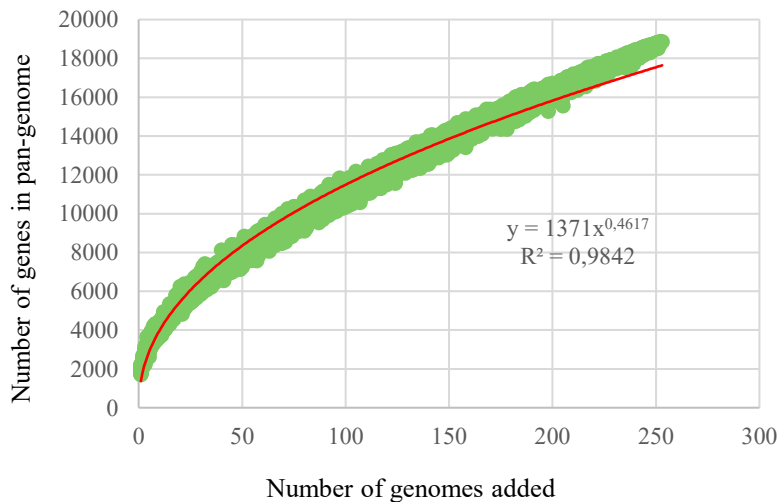


Figure S3