

# FastAAI: Efficient Estimation of Genome Average Amino Acid Identity and Phylum-level relationships using Tetramers of Universal Proteins

Konstantinos Konstantinidis (✉ [kostas.konstantinidis@gatech.edu](mailto:kostas.konstantinidis@gatech.edu))

Georgia Institute of Technology

Carlos Ruiz-Perez

Georgia Institute of Technology

Kenji Gerhardt

Georgia Institute of Technology <https://orcid.org/0000-0003-0644-1862>

Luis Rodriguez-R

University of Innsbruck <https://orcid.org/0000-0001-7603-3093>

Chirag Jain

Indian Institute of Science <https://orcid.org/0000-0002-4300-0794>

James Tiedje

Michigan State University <https://orcid.org/0000-0002-8992-6218>

James Cole

Michigan State University

---

## Article

### Keywords:

**Posted Date:** March 23rd, 2022

**DOI:** <https://doi.org/10.21203/rs.3.rs-1459378/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

---

1 **FastAAI: Efficient Estimation of Genome Average Amino Acid**  
2 **Identity and Phylum-level relationships using Tetramers of Universal**  
3 **Proteins**

4 Carlos A Ruiz-Perez<sup>1\*</sup>, Kenji Gerhardt<sup>1\*</sup>, Luis M Rodriguez-R<sup>2,3\*</sup>, Chirag Jain<sup>4</sup>, James M. Tiedje<sup>5</sup>,  
5 James R Cole<sup>5</sup>, and Konstantinos T Konstantinidis<sup>1,2,\*</sup>

6 <sup>1</sup>School of Biological Sciences and <sup>2</sup>School of Civil and Environmental Engineering, Georgia Institute of Technology, Atlanta, GA  
7 30332, USA.

8 <sup>3</sup>Department of Microbiology and Digital Science Center (DiSC), University of Innsbruck, Innsbruck 6020, Austria.

9 <sup>4</sup>Department of Computational and Data Sciences, Indian Institute of Science, Bengaluru, KA 560012, India

10 <sup>5</sup>Center for Microbial Ecology, Michigan State University, East Lansing, MI 48824, USA.

11 \*Contributed equally

12 \*To whom correspondence should be addressed.

13

14 **Abstract**

15 Estimation of whole-genome relatedness and taxonomic identification are two important bioinformatics  
16 tasks in describing environmental or clinical microbiomes. The genome-aggregate Average Nucleotide  
17 Identity (ANI) is routinely used to derive the relatedness of closely related (species level) microbial and  
18 viral genomes, but it is not appropriate for more divergent genomes. Average Amino Acid Identity (AAI)  
19 can be used in the latter cases, but no current AAI implementation can efficiently compare thousands of  
20 genomes. Here we present FastAAI, a tool that estimates whole-genome pairwise relatedness using shared  
21 tetramers of universal proteins in a matter of microseconds, providing a speedup of up to 5 orders of  
22 magnitude when compared with current methods of calculating AAI or alternative whole-genome metrics.  
23 Further, FastAAI resolves distantly related genomes related at the phylum level with comparable accuracy  
24 to the phylogeny of ribosomal rRNA genes, substantially improving on a known limitation of current AAI

25 implementations. Therefore, Fast AAI uniquely expands the toolbox for microbiome analysis and allows it  
26 to scale to millions of genomes.

27

## 28 **Introduction**

29         The diversity of microbial and viral genomes on the planet is very large, e.g., estimated at over a  
30 billion species of (just) bacteria, and most of it remains undiscovered<sup>1</sup>. As genome sequencing can help  
31 characterizing this diversity and has recently become routine, most microbial scientists have been  
32 overwhelmed by the amount of available data. For many researchers, a complete and thorough analysis of  
33 the available genome data is not necessary. Instead, these users would be better served with straightforward  
34 methods that can identify their unknown DNA/RNA sequence(s) and be able to discriminate between well-  
35 understood taxa and those that are potentially novel. Thus, tools that can help researchers identify the most  
36 "interesting" or relevant genomes to their unknown/query genome(s) among thousands of candidates are  
37 important. Furthermore, metagenomics, which is the sequencing of environmental DNA, allows the genetic  
38 characterization of the majority of these microbes that have resisted cultivation in the laboratory. However,  
39 current tools to analyze metagenomic data are clearly lagging behind the development of sequencing  
40 technologies (and data), and do not scale with the number of metagenome-assembled genomes (MAGs),  
41 currently in the order of hundreds of thousands, that have become available<sup>2</sup>. This is a major limitation for  
42 better understanding, studying, and communicating about the biodiversity of uncultivated microorganisms  
43 that run the life-sustaining biogeochemical cycles on the planet, form critical associations with their plant  
44 and animal hosts, or produce products of biotechnological value.

45         One commonly used pair of methods to determine the level of novelty of a newly described taxon  
46 relative to the taxa already available in reference databases are genome-aggregate average nucleotide  
47 (ANI) and amino-acid (AAI) identity. ANI represents the average nucleotide identity of all genes shared  
48 between any two genomes and offers robust resolution between strains of the same or closely related  
49 species (i.e., showing 80-100% ANI). However, ANI is not appropriate for estimating genome relatedness

50 (and thus, classification) among more divergent genomes because nucleotide substitutions saturate and/or  
51 sequences cannot be reliably aligned at this level<sup>3</sup>. Instead, the average amino acid identity (AAI) should  
52 be used for moderately divergent genomes as it effectively delineates relationships at the genus and  
53 family levels, similarly to ANI for the species level. It is also important to note that the ANI/AAI-based  
54 approach can be applied to any unknown genome sequence, even if it is incomplete, that encodes at least  
55 a few genes shared with at least a few reference genomes<sup>4</sup>, whereas alternative approaches based on  
56 universal genes such as the small subunit ribosomal rRNA (SSU rRNA) gene are limited to sequences  
57 carrying the corresponding genes. Accordingly, the ANI/AAI approach has facilitated the analysis of an  
58 increasing number of MAG, single-cell amplified genome (SAG) or viral genome sequences and the need  
59 to perform whole-genome ANI or AAI calculations has grown exponentially in recent years.

60         Despite these strengths, traditional ANI/AAI calculation is based on pairwise sequence  
61 alignments, e.g., BLAST<sup>5</sup>, and thus do not scale well with an increasing number of genomes. For  
62 instance, performing all vs. all ANI calculation among the ~13,000 reference genomes (~169 million  
63 pairwise comparisons) of all bacterial species taxonomically described to date<sup>6</sup> required ~1 month of  
64 processing time on a computer cluster of 500 CPUs. Performing a similar comparison among 100,000  
65 genomes (10 billion comparisons) is expected to take more than 1 year due to the exponential increase in  
66 the number of comparisons to perform. This is prohibitively expensive, even for larger computer clusters.  
67 While a faster ANI implementation based on the numbers of k-mers shared between genomes (FastANI)  
68 has been recently described by our team<sup>4</sup>, and has been widely used, current AAI implementations are not  
69 fast enough, especially considering that the majority of genomes in the reference database are moderately  
70 or distantly related to each other. Further, AAI does not provide robust resolution at the inter-domain and  
71 phylum levels (that is, between distantly related genomes) compared to the phylogenetic analysis of the  
72 SSU rRNA or other universal genes. Alternative approaches such as phylogenetic placement based on  
73 universal genes are accurate<sup>7,8</sup>, but are also computationally expensive and do not scale well with the  
74 exponentially increasing number of microbial genomes recovered from environmental samples. Thus,

75 there is an urgent need for a scalable tool to calculate genome relatedness, especially for deep-branching  
76 genomes with no close (e.g., same species) relatives available.

77 Here we present FastAAI, a bioinformatics tool aiming at solving the issues related to the speed  
78 and scalability of the AAI estimation among bacterial and archaeal genomes while providing higher  
79 resolution at the phylum and domain level. FastAAI is a standalone Python tool that leverages protein  
80 tetramer information in single-copy proteins (SCPs) “universally” present in microbial genomes to calculate  
81 genome similarities that are subsequently translated into “traditional” AAI values of relatedness. The  
82 current implementation of FastAAI can perform pairwise genome comparisons in a matter of microseconds  
83 and scales exceedingly well with a larger number of genomes compared.

84

## 85 **Results**

86 We focus our analysis below primarily on bacterial genomes because of their large numbers in the  
87 public datasets. FastAAI applies similarly to archaeal genomes but for genomes of other domains (e.g.,  
88 eukaryotic) adjustments in the collection of universal single-copy proteins will be required for best results.  
89 For our accuracy evaluations, we primarily focus on the comparison of FastAAI to 16S rRNA gene  
90 identities since the latter gene represents the backbone of bacterial taxonomy and has been most commonly  
91 used to assess genetic relatedness among taxa/genomes. We provide computational and performance  
92 comparisons to traditional alignment-based methods such as DIAMOND<sup>9</sup> and EzAAI<sup>10</sup> as well as  
93 alternative methods such as GTDB.

94

### 95 *FastAAI accurately estimates AAI*

96 FastAAI is intended to replace traditional AAI in resolving relationships above the species level;  
97 at and within the species level, ANI is sensitive and efficient enough (e.g., FastANI implementation) and  
98 thus, already serves as a complementary method to AAI. Therefore, we focused our performance evaluation  
99 of FastAAI on genomes with a relatedness between 30% and 90% AAI. Using the RefSeq dataset (see

100 Methods), we calculated a “reference” AAI value for each pair of genomes using DIAMOND and compared  
101 predictions made with FastAAI against these values. FastAAI’s estimator  $\widehat{AAI}$  (see Methods for details on  
102 the calculation of  $\widehat{AAI}$ ) achieves a highly accurate estimation of DIAMOND-based AAI between 30% and  
103 90%  $\widehat{AAI}$ , with an average prediction error of only 1.3% within this range (R-squared: 0.98, RMSE: 1.3%;  
104 Figure 1). The strong linear correlation between  $\widehat{AAI}$  and DIAMOND-based AAI values demonstrates that  
105 FastAAI can provide a similarly robust estimate of genome relatedness above the species level (Figure 2B  
106 and 2C). For the most distantly related genomes (e.g., related at around 30% AAI), the  $\widehat{AAI}$  and  
107 DIAMOND-based AAI values deviated the most from each other, which we investigated below with  
108 comparisons against the 16S rRNA gene and showed that this was due to higher accuracy of  $\widehat{AAI}$  at the  
109 domain and phylum levels.

110

### 111 *FastAAI is faster than alternative AAI implementations*

112 To demonstrate the gains in speed of FastAAI compared to alignment-based methods for  
113 calculating AAI, we performed an all-vs-all comparison of 100 randomly chosen RefSeq genomes (10,000  
114 pairwise comparisons) using FastAAI, DIAMOND-based AAI, and the recently published EzAAI. EzAAI  
115 uses MMseqs2<sup>11</sup> as a substitute for either BLAST or DIAMOND as its alignment engine, but calculates  
116 AAI in a fundamentally similar manner; that is, finding reciprocal best-match alignments between the  
117 proteins of each pair of genomes it processes and calculating AAI based on these best matches. As  
118 DIAMOND does not independently function as a pipeline for both protein prediction and alignment in the  
119 manner that FastAAI and EzAAI do, all three tools were supplied with the same set of predicted proteins  
120 generated from the original 100 RefSeq genomes using Prodigal. Timing covered all of the remaining  
121 comparable steps of each tool, including data formatting, database creation as needed, the calculation of  
122 AAI, and writing results.

123 The average time per pairwise genome comparison was 9.6 seconds (S.D. = 4.16) for DIAMOND,  
124 32.0 seconds (S.D. = 12.1) for EzAAI, and 0.031 seconds (S.D. non-applicable) for FastAAI. Although this

125 test already demonstrates more than two orders of magnitude speedup over both DIAMOND and EzAAI,  
126 it is in fact a nearly worst-case scenario for demonstrating the speed gains of FastAAI relative to these tools.  
127 As can be seen in Figure S1, the relative speed gain per comparison of FastAAI over DIAMOND and  
128 EzAAI increases as the datasets processed grow larger, with FastAAI becoming nearly 5 orders of  
129 magnitude faster per comparison than either DIAMOND or EzAAI in our largest test dataset (~2.29 billion  
130 pairwise comparisons) when a run started from predicted proteins.

131

### 132 ***FastAAI is faster and has a lower RAM requirement than GTDB-tk***

133 There are alternative methods to measure relatedness among genomes and/or taxonomically  
134 classify them aside from the traditional ANI/AAI; most notably, GTDB-tk<sup>12</sup> uses a combination of  
135 approximate, maximum-likelihood-based phylogenetic placement and FastANI to classify query genomes.  
136 Because GTDBtk can use only its own prebuilt database, we created an equivalent FastAAI database from  
137 the same set of genomes that are used in the latest GTDB database (rel. 202, last updated Jul. 9<sup>th</sup>, 2021),  
138 totaling 45,555 bacterial and 2,339 archaeal genomes. We randomly sampled three sets of 100 genomes  
139 from these 47,894 genomes and used these sets to test both GTDBtk and FastAAI by querying each set of  
140 100 selected genomes against the full complement of 47,894 target genomes. To ensure that the comparison  
141 between FastAAI and GTDBtk was fair, both tools were run with 10 threads, supplied with 256 GB of  
142 RAM, used their respective premade databases as search targets, and were supplied with query sets  
143 composed of genomes (not proteins or HMM results), as this is the starting point for a GTDBtk search. In  
144 addition to memory usage, we measured both the walltime time and CPU time for each tool, reflecting the  
145 user's waiting time from program start to program end and the sum of computer time across all 10 threads,  
146 respectively. On average, FastAAI completed in 9.8% of the runtime, 16.9% of the CPU time, and required  
147 only 0.4% of the RAM required by GTDBtk. All comparison results can be seen in Table 1. All times are  
148 reported as HH:MM:SS.

149

<b>Table 1. FastAAI vs GTDBtk runtime and resource usage; 100 genomes vs. GTDB Release 202</b>			
<b>Test</b>	<b>Walltime</b>	<b>CPU Time</b>	<b>Max RAM use (GB)</b>
FastAAI, Sample 1	00:07:46	00:46:28	0.81
GTDBtk, Sample 1	01:26:38	05:18:55	232.08
FastAAI, Sample 2	00:07:51	00:50:11	0.87
GTDBtk, Sample 2	01:17:41	04:12:42	231.67
FastAAI, Sample 3	00:08:31	00:54:35	0.98
GTDBtk, Sample 3	01:22:38	05:24:31	232.10

150

151 To further illustrate the minimal RAM required by FastAAI, we ran an additional all-vs-all tests  
 152 with FastAAI on varying sizes of random samples of the GTDB genomes, ranging in size from 50 vs. 50  
 153 to the entirety of GTDB release 202 (47,894 vs 47,894 genomes). The whole-GTDB all-vs-all test was the  
 154 longest-running and most memory intensive of these tests (~2.29 billion pairwise comparisons), which  
 155 FastAAI was able to complete in a walltime of 1:28:27, a CPU time of 10:33:07, and using only 4.21 GB  
 156 of RAM. Provided a preprocessed database (i.e., universal proteins identified and tetramers extracted), even  
 157 a comparison of this enormous scale could be accomplished with FastAAI within a day, on a typical  
 158 personal laptop. Performing the same task with GTDB-tk would take well in excess of 500 CPU hours and  
 159 require hundreds of gigabytes of RAM that are available only on supercomputers. The results for the other  
 160 tests can be seen in Table S1.

161

162 ***FastAAI vs. 16S rRNA gene identities and phylum-level resolution***

163 Our comparison of DIAMOND-based AAI and 16S rRNA gene identities using the bacterial  
 164 RefSeq genomes demonstrated a strong correlation (Figure S2A), with overlaps between adjacent  
 165 taxonomic ranks in terms of the AAI values (e.g., about 30% of the genomes grouped at the genus level as  
 166 their lowest shared rank show similar AAI values to genomes grouped at the family level), consistent with

167 the patterns observed in previous studies based on smaller genome datasets<sup>13</sup>. In these comparisons, we  
168 observed that the lowest values for 16S rRNA gene identities (70-80%) corresponded to (traditional) AAI  
169 values between 30-40%. However, the same range in AAI values (30-40%) also included genome pairs  
170 with higher 16S rRNA gene identities (80-85%), revealing substantial overlaps in AAI values (i.e., low  
171 resolution) at higher taxonomic levels, as noted also previously<sup>13</sup>. However, in the case of FastAAI, 16S  
172 rRNA gene identities between 70 and 80% corresponded to FastAAI's  $\widehat{AAI}$  values between 32.28 and 33.57  
173 with almost no overlap with same- vs. different-domain level comparisons (Figure S2B), revealing higher  
174 resolution of FastAAI at the phylum-domain level, comparable to that of 16S rRNA gene identities [see  
175 also Fig. S3 for discussion of a few outlier genome pairs observed.

176 More specifically, for 16S rRNA gene identities, we identified the inter-domain and inter-phylum  
177 valleys at 70.4 and 80.2%, respectively (Figure 2A). We observed that only 1.1% of same-domain  
178 comparisons were below 70.4%, while only 0.54% of different-domain comparisons were above this  
179 threshold, indicating limited overlap and good discrimination capacity at the domain level based on 16S  
180 rRNA gene identities. However, the same values for phylum-level comparisons revealed a lower  
181 discriminatory capacity; 34.7% of the comparisons were below the 80.2% threshold and 10.9% were above.  
182 Note that phyla are typically designated based on their branching patterns in the 16S rRNA gene phylogeny,  
183 not necessarily the 16S rRNA gene identity; hence, the results reported here do not necessarily reflect low  
184 performance of the 16S rRNA gene phylogeny approach. Our goal was to examine if convenient sequence  
185 identity thresholds for discriminating taxonomic ranks can be determined that could guide future taxonomic  
186 studies and description and to evaluate FastAAI's domain- and phylum-level resolution. Compared to 16S  
187 rRNA gene identity values, we found that the thresholds discriminating domain and phylum-level  
188 comparisons using FastAAI were 35.3% and 40.5%  $\widehat{AAI}$ . Interestingly, we found that only 0.07% of  
189 different-domain and 0.03% of same-domain comparisons were below and above the domain-threshold  
190 (35.3%), indicating similar -if not higher- discriminatory power at this level for  $\widehat{AAI}$  relative to the 16S  
191 rRNA gene. Phylum-level comparisons were nonetheless very similar to those observed for 16S rRNA with

192 20.4% same-phylum and 6.47% different-phylum comparisons below and above the identified threshold  
193 (40.5%), respectively. These results indicate that FastAAI offers a lower frequency of misclassifications  
194 (or higher resolving power) than (traditional) AAI and 16S rRNA gene identities at the domain and phylum  
195 level for deep-branching genomes. It should be also noted that the large number of genome pairs sharing  
196 the same phylum rank but showing identities below the discrimination threshold in all three metrics  
197 indicates that at least some of the corresponding genomes may deserve a distinct (novel) phylum  
198 designation based on genomic relatedness. We evaluated the remaining taxonomic ranks e.g., family and  
199 genus in a similar fashion, and showed that reliable FastAAI thresholds can be established for each rank,  
200 and that FastAAI's classification based on the best-matching genome largely (~95% of the cases) agreed  
201 with more sensitive tree placement classification approaches (see Supplementary Results for details).

202

### 203 ***FastAAI is robust to varying levels of genome completeness and contamination***

204 Most studies using culture-independent genomic techniques filter recovered prokaryotic genomes  
205 according to their quality level based on estimated completeness and contamination<sup>14</sup>. Gene phylogenies or  
206 classification methods that rely on a single or few genes are often affected by genome incompleteness due  
207 to the absence of the required gene(s) for the analysis. In addition, contaminating gene sequences can be  
208 sources of error that can affect homology-based classification methods such as MEGAN or MyTaxa<sup>15, 16</sup>.  
209 In such cases, AAI is generally more robust because the number of proteins or genome fragments used in  
210 the calculations may be adequate even for relatively incomplete genomes (e.g., ~200 shared proteins are  
211 typically adequate<sup>4</sup>). Given that FastAAI relies on a reduced set of universal proteins, it is important to  
212 assess the consistency of the results at variable levels of genome quality. For this, we selected 500 genomes  
213 from the RefSeq dataset that were estimated to be 100% complete with 0% contaminated. We predicted the  
214 proteins for each selected genome and used FastAAI to calculate their genome relatedness against the  
215 smaller RefSeq dataset to create a baseline of reference  $\widehat{AAI}$  values. Each set of proteins per genome was  
216 then randomly subsampled and *in-silico* “contaminated” to create a set of genomes with varying levels of  
217 contamination and completeness (see methods, Figure S3A). Each of these incomplete and contaminated

218 genomes was then searched against the RefSeq dataset and the resulting  $\widehat{AAI}$  values of these searches were  
219 compared to the complete-genome results (baseline) to determine the degree of deviation from the baseline.  
220 Our results demonstrated that in the absence of contamination, FastAAI could recover almost identical  $\widehat{AAI}$   
221 values for genomes with varying levels of completeness, even down to 10% complete genomes, only  
222 deviating by, at most, 2.26% from the complete-genome-based  $\widehat{AAI}$  values in the genomes of the lowest  
223 completeness (Figure S3B). This result showed that FastAAI is robust even for incomplete genomes, as  
224 long as about 8 of the universal proteins are present (roughly corresponding to 10% completeness). As  
225 expected, the contamination resulted in an increased deviation of  $\widehat{AAI}$  values from the baseline. In the 100%  
226 completeness genomes, the deviation increased from 0% (100% completeness, 0% contamination) to 0.95%  
227 (100% completeness, ~37% contamination), a relatively small value considering the high level of  
228 contamination. This trend was also observed in lower completeness levels, where we observed the highest  
229 values of deviation, ranging from 2.26% to 3.98%  $\widehat{AAI}$  in genomes with 10% completeness and  
230 contamination of up to 37% (a very low-quality genome; Figure S3B). Therefore, considering that proposed  
231 classification standards for the quality of assembled genomes establish a low quality-draft as a genome with  
232 <50% completeness and <10% contamination<sup>14</sup>, we expect FastAAI to perform well, with less than 1%  $\widehat{AAI}$   
233 deviation from the true estimate, for such low quality genomes.

234

## 235 **Discussion**

236 This study presented FastAAI, a tool for the fast and accurate estimation of the whole-genome  
237 genetic relatedness and classification of microbial genomes that provides comparable results, if not better,  
238 with traditional methods such as AAI, and SSU rRNA gene identities, while dramatically improving in  
239 terms of speed and resolution. We also demonstrated substantial advantages over newer alternative tools  
240 for the classification of microorganisms in terms of speed and computational resources. The high accuracy  
241 and speed, especially in comparisons among deep-branching genomes (i.e., higher taxonomic levels), will  
242 be important for environmental and clinical samples that are expected to harbor substantial novel

243 diversity<sup>17</sup>. FastAAI can effectively replace the abovementioned methods for genome comparisons at the  
244 genus level and above, while reverting to the traditional AAI, or even better ANI, for closely related  
245 genomes with high FastAAI identities (species level). Thus, species-level comparisons should instead be  
246 resolved with FastANI or regular ANI comparisons, and particularly novel genomes with no apparent close  
247 relatives in the reference databases at the (same) genus level can be further processed with FastAAI and/or  
248 using phylogenetic methods when number of genomes is tractable.

249         FastAAI accelerates the calculation of AAI because it does not require the alignment of proteins  
250 between genomes, and instead estimates AAI between two genomes from the shared fraction of tetramers  
251 over a set of universal SCPs shared by the two genomes. This allows FastAAI to rapidly calculate AAI  
252 from massive datasets: even in all-vs-all comparisons of tens of thousands of genomes (i.e., hundreds of  
253 millions or billions of pairwise comparisons), FastAAI requires only a few hours of runtime and very  
254 modest computational resources to complete its estimation of AAI. Indeed, the overwhelming majority of  
255 runtime for FastAAI in all of our applicable tests is consumed in the preprocessing stage, i.e., protein  
256 prediction and HMM searches to identify SCPs, and the average time per comparison between a pair of  
257 genomes decreases as the number of comparisons performed increases (Figure S1). This speedup is due to  
258 the fact that even though the total number of pairwise genome comparisons increases with the square of the  
259 number of genomes, the slow pre-processing step grows linearly and only the number of fast tetramer  
260 pairwise comparisons increases with the square of the genome number. Similarly, comparison of a single  
261 novel genome to a large, preprocessed reference data set requires additional preprocessing of only the novel  
262 genome, which dominates comparison times against even our large data sets. It is important to note,  
263 however, that the preprocessing step is performed only once for each genome and can be stored in order to  
264 accelerate subsequent AAI calculations.

265         Notably, FastAAI retains its performance even in the presence of highly incomplete or  
266 contaminated genomes, which are commonly recovered by metagenomics surveys. Finally, from a database  
267 perspective, FastAAI retrieves the closest relative in most cases even using a highly incomplete database,  
268 which we expect that will only improve as databases become more comprehensive with genome

269 representatives from uncultured microorganisms. Therefore, we expect FastAAI to be a valuable tool for  
270 the scientific community aimed at the fast and accurate genome relatedness estimation and classification of  
271 microorganisms from clinical, industrial, or environmental sources.

272

## 273 **Online Methods**

### 274 *FastAAI implementation*

275 FastAAI is divided into four steps (Figure S4A). Briefly, FastAAI performs, (i) identification of  
276 marker proteins in each genome using a set of 122 universal single-copy proteins (SCP) represented by  
277 hidden Markov models (HMMs) (Table S2, and below), (ii) extraction of tetramers for each protein per  
278 genome, (iii) estimation of Jaccard indices based on the tetramers extracted from proteins shared by a  
279 genome pair, (iv) estimation of the average Jaccard index ( $\bar{J}$ ) followed by its translation to estimated AAI  
280 values (termed  $\widehat{AAI}$  to distinguish it from AAI values estimated based on alignment-based methods). The  
281 first step is performed by searching all proteins of a genome against a collection of HMMs with HMMer  
282 3.0<sup>18</sup> using profile trusted cutoffs (`--cut_tc`). These models were built from a set of 122 proteins typically  
283 found in most prokaryotic genomes in a single copy<sup>19</sup>. The second step consists of extracting all tetramers  
284 per marker protein. Tetramers are considered as present or absent regardless of the number of times they  
285 occur in a protein sequence. In the third step, for each pair of genomes, FastAAI identifies shared marker  
286 proteins; if a protein is only present in one genome, it is excluded from the comparison. Then, for each  
287 shared marker protein in the genome pair under comparison, a Jaccard similarity value is calculated based  
288 on the shared vs. total tetramers of the corresponding protein sequences (Figure S4B). Finally, FastAAI  
289 calculates the average of the previously estimated Jaccard values for all shared proteins ( $\bar{J}$ ), the standard  
290 deviation, and uses  $\bar{J}$  values to estimate AAI values ( $\widehat{AAI}$ ; read below for estimation formula). In cases  
291 where the inputs are DNA genome sequences, the program predicts protein sequences using Prodigal<sup>20</sup> and  
292 automatically selects the optimal translation table (11 or 4).

293

294 **Testing Datasets**

295 To test performance and speed, we evaluated FastAAI on three primary datasets (Table 2). The first  
296 dataset (RefSeq) comprised all complete reference genomes from the NCBI RefSeq database with available  
297 taxonomic classification as of 11-13-2019. In total, this dataset was composed of 5,000 bacterial and 328  
298 archaeal genomes. The second dataset (TypeMat) included 10,573 genomes (10,200 *Bacteria* and 373  
299 *Archaea*) from type material with standing in nomenclature, which are available through the MiGA  
300 webserver with release number r2019-02<sup>6,21</sup>. Genomes in either dataset with completeness estimates below  
301 50% or contamination higher than 30% based on MiGA estimates were removed from further analyses as  
302 low-quality genomes. Finally, the third dataset (labeled PhylaLite) was composed of 39 complete genomes  
303 (36 *Bacteria* and 3 *Archaea*), one genome per phylum from the TypeMat collection. In addition to these  
304 primary datasets, we included simulated datasets derived from a subset of genomes in the RefSeq dataset  
305 (n=500) to evaluate FastAAI when genomes had different levels of completeness and contamination. For  
306 this, we created *in silico* genomes with varying degrees of completeness by randomly subsampling the  
307 single-copy marker proteins of these RefSeq genomes to a given level of completeness. Then, each of the  
308 resulting genomes was modified to increase contamination by randomly adding single-copy marker proteins  
309 from distantly related organisms to achieve a given level of contamination. All metrics used, including  
310 length, N50, completeness, and contamination, are shown for all datasets in Figure S5.

**Table 2. Datasets used to evaluate FastAAI**

<b>Name</b>	<b>No. of Genomes</b>	<b>Mean Completeness</b>	<b>Mean Contamination</b>
TypeMat	10,573	98.34	2.42
RefSeq	5,328	98.38	1.29
PhylaLite	39	95.98	2.15

311

312

313 **FastAAI's Method of Estimating AAI**

314 FastAAI utilizes HMMER to identify proteins in a genome's proteome as being one of the 122  
315 prokaryotic single-copy proteins (discussed above, Figure S4A). FastAAI considers a protein to be a  
316 particular SCP when HMMER identifies the protein's best match to be the SCP in question and the protein  
317 has no better match among the remaining SCPs. Consequently, a protein can only be classified as one SCP  
318 and each SCP can have only one protein as a representative within a single genome. No information about  
319 proteins that are not labeled as an SCP representative is retained by FastAAI; as a result, a genome can have  
320 at most 122 proteins considered in the process of estimating AAI. To estimate AAI for a pair of genomes,  
321 FastAAI determines which SCPs the two genomes share and collects the SCP representative proteins of  
322 each in matched-SCP pairs. FastAAI then calculates the Jaccard index between the tetramers of each pair  
323 and averages them to produce an initial estimator that we label as  $\bar{J}$  (Figure S4B).

324 Comparison of FastAAI to DIAMOND-based AAI using the RefSeq dataset revealed a linear  
325 correspondence between  $\bar{J}$  and AAI values ( $R^2$ : 0.96, RMSE: 2.22%; Figure S6A). However, despite the  
326 low root mean squared error (RMSE) and the high R-squared of the linear regression, AAI values predicted  
327 using the regression coefficients were consistently lower compared to the actual AAI values, especially in  
328 the inter-phylum and inter-class levels, as evidenced by the curvature over the regression line at these  
329 regions. To solve this and improve the correspondence between FastAAI's estimate and AAI values, we  
330 transformed  $\bar{J}$  by applying several transformations to linearize the relationship. The linearization  
331 coefficients were optimized using a gradient descent method aiming at decreasing the RMSE, excluding  
332 values with  $\bar{J} > 0.9$  (see below). Finally, to have a 1:1 correspondence, we applied a linear model on top of  
333 the initial transformation. The final transformation defined as the estimated AAI ( $\widehat{AAI}$ ) is described by the  
334 following formula:

$$335 \quad \widehat{AAI} = a + b * e^{(-(-c * \log(\bar{J}))^{1/d})}$$

$$336 \quad a = -0.3087057; b = 1.810741; c = 0.2607023; d = 3.435$$

337 The linear regression of AAI vs.  $\widehat{AAI}$  values in the 0-100% range shows a slightly lower R-squared  
338 of 0.95 and a high RMSE (22.02%; Figure S6B). This underperformance is due to data points with  $\bar{J}$  values  
339 higher than 0.9 (and thus, the corresponding values of  $\widehat{AAI}$ ) collapsing to the AAI range of ~95-100%,  
340 increasing the regression error and generating curvature in these portions of the plot (Figure S6B). The low  
341 correspondence between FastAAI's  $\bar{J}$  and AAI values in the upper range is presumably due to the high  
342 sequence conservation of the universal genes (thus, low resolving power) at this (the species) level.  
343 Therefore, FastAAI summarizes these high  $\bar{J}$  values as “>90% AAI” and we recommend using instead ANI  
344 (e.g., FastANI) or traditional AAI to compare genomes related at this level, as also suggested previously<sup>4</sup>.  
345 <sup>14</sup>. As AAI values tend to remain greater than 30% even across domains, FastAAI is also insensitive below  
346 30% AAI and so summarizes  $\widehat{AAI}$  estimates at or below 30% as “<30% AAI.”

347

#### 348 *Estimation of 16S rRNA gene identities, AAI, and GTDBtk-RED values*

349 We used INFERNAL (INFERENCE of RNA ALignment) v1.1.2<sup>22</sup> to identify and extract 16S rRNA  
350 gene sequence identities for all genomes in the RefSeq dataset. In cases where more than one 16S rRNA  
351 gene sequence was found in a single genome, the longest sequence was selected as the representative for  
352 the genome, or if multiple sequences had the same length, the representative was randomly chosen among  
353 them. We then performed pairwise alignments and determined the identities of each pair of sequences using  
354 a custom python script (available at [https://github.com/cruizperez/Bioinformatic\\_Tools](https://github.com/cruizperez/Bioinformatic_Tools)), which uses the  
355 Needleman–Wunsch algorithm with the same parameters as implemented in EMBOSS Needle<sup>23</sup>. We used  
356 the aai.rb tool from the Enveomics collection<sup>24</sup> with DIAMOND as the searching tool to calculate pairwise  
357 AAI values. Finally, we performed classifications and extracted Relative Evolutionary Diverge (RED)  
358 values based on the Genome Taxonomy Database (GTDB) using GTDBtk with release #202 of the  
359 database.

360

#### 361 *FastAAI taxonomic classification accuracy*

362 We evaluated the ability of FastAAI  $\widehat{AAI}$  values to discriminate genomes at different taxonomic  
363 levels by using the TypeMat dataset as the ground truth and reference taxonomy. For this, we determined  
364 threshold values delimiting taxonomic ranks by finding valleys in the distribution of  $\widehat{AAI}$  values between  
365 inter- vs. intra-rank comparisons for the TypeMat genome dataset. For example, genomes assigned to the  
366 same domain typically show values above 35.3%  $\widehat{AAI}$  and those assigned to different domain typically have  
367  $\widehat{AAI}$  values below 35.3%. Genomes in the RefSeq dataset (and their associated taxonomy, assumed to be  
368 correct) were treated as queries, and their taxonomic classification based on the estimated  $\widehat{AAI}$  values  
369 against TypeMat genomes was compared to their available taxonomy (from RefSeq). For each rank, we  
370 built a confusion matrix where we labeled the classifications obtained based on  $\widehat{AAI}$  value thresholds and  
371 compared them to the available taxonomic classification of the RefSeq query genome from NCBI as  
372 follows: True positive (TP) for pairs with  $\widehat{AAI} >$  taxonomic rank threshold (or just “threshold” for  
373 simplicity) and where query and reference shared the same taxon, true negative (TN) for pairs with  $\widehat{AAI} <$   
374 threshold and where the pair belonged to a different taxon in the rank, false positive (FP) for pairs with  $\widehat{AAI}$   
375  $>$  threshold and where the pair belonged to a different taxon in the rank, and false negatives (FN) for pairs  
376 with  $\widehat{AAI} <$  threshold and where the pair belonged to the same taxon. From this confusion matrix, we  
377 calculated the accuracy, recall, precision, and F1 score for each taxonomic rank as follows:

$$378 \quad Accuracy = \frac{TP + TN}{TP + TN + FP + FN}$$

$$379 \quad Recall = \frac{TP}{TP + FN}$$

$$380 \quad Precision = \frac{TP}{TP + FP}$$

$$381 \quad F1 = 2 \times \frac{Precision \times Recall}{Precision + Recall}$$

### 382 ***FastAAI availability***

383 FastAAI is an open-source tool freely available at <https://github.com/KGerhardt/FastAAI>.

384

## 385 Acknowledgments

386 This work has been supported by US NSF under awards # DBI1356288, -1356380, -1759831, and -  
387 1759892.

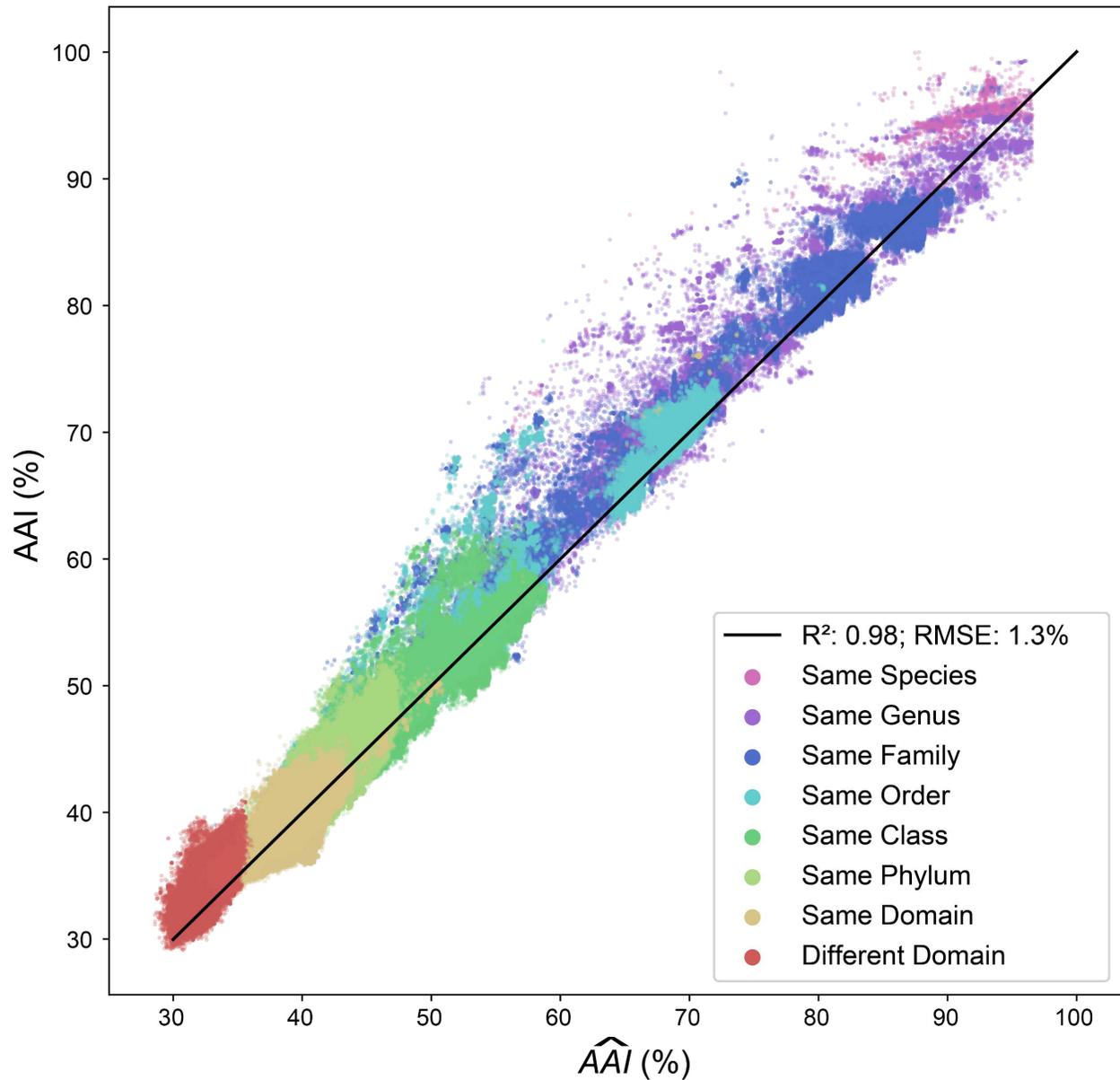
388

## 389 References

- 390 1. Amann, R. and R. Rossello-Mora (2016). "After All, Only Millions?" mBio **7**(4).
- 391 2. Nayfach, S., S. Roux, R. Seshadri, D. Udway, N. Varghese, F. Schulz, D. Wu, D. Paez-  
392 Espino, I. M. Chen, M. Huntemann, K. Palaniappan, J. Ladau, S. Mukherjee, T. B. K.  
393 Reddy, T. Nielsen, E. Kirton, J. P. Faria, J. N. Edirisinghe, C. S. Henry, S. P. Jungbluth,  
394 D. Chivian, P. Dehal, E. M. Wood-Charlson, A. P. Arkin, S. G. Tringe, A. Visel, I. M. D.  
395 Consortium, T. Woyke, N. J. Mouncey, N. N. Ivanova, N. C. Kyrpides and E. A. Elie-  
396 Fadrosch (2021). "A genomic catalog of Earth's microbiomes." Nat Biotechnol **39**(4): 499-  
397 509.
- 398 3. Konstantinidis, K. T. and J. M. Tiedje (2007). "Prokaryotic taxonomy and phylogeny in  
399 the genomic era: advancements and challenges ahead." Curr Opin Microbiol **10**(5): 504-  
400 509.
- 401 4. Jain, C., R. L. Rodriguez, A. M. Phillippy, K. T. Konstantinidis and S. Aluru (2018). "High  
402 throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries."  
403 Nat Commun **9**(1): 5114.
- 404 5. Altschul, S. F., T. L. Madden, A. A. Schaffer, J. Zhang, Z. Zhang, W. Miller and D. J.  
405 Lipman (1997). "Gapped BLAST and PSI-BLAST: a new generation of protein database  
406 search programs." Nucl. Acids. Res. **25**(17): 3389-3402.
- 407 6. Rodriguez, R. L., S. Gunturu, W. T. Harvey, R. Rossello-Mora, J. M. Tiedje, J. R. Cole  
408 and K. T. Konstantinidis (2018). "The Microbial Genomes Atlas (MiGA) webserver:  
409 taxonomic and gene diversity analysis of Archaea and Bacteria at the whole genome level."  
410 Nucleic Acids Res **46**(W1): W282-W288.
- 411 7. Sunagawa, S., D. R. Mende, G. Zeller, F. Izquierdo-Carrasco, S. A. Berger, J. R. Kultima,  
412 L. P. Coelho, M. Arumugam, J. Tap, H. B. Nielsen, S. Rasmussen, S. Brunak, O. Pedersen,  
413 F. Guarner, W. M. de Vos, J. Wang, J. Li, J. Dore, S. D. Ehrlich, A. Stamatakis and P. Bork  
414 (2013). "Metagenomic species profiling using universal phylogenetic marker genes." Nat  
415 Methods **10**(12): 1196-1199.
- 416 8. Parks, D. H., M. Chuvochina, D. W. Waite, C. Rinke, A. Skarszewski, P. A. Chaumeil and  
417 P. Hugenholtz (2018). "A standardized bacterial taxonomy based on genome phylogeny  
418 substantially revises the tree of life." Nat Biotechnol **36**(10): 996-1004.

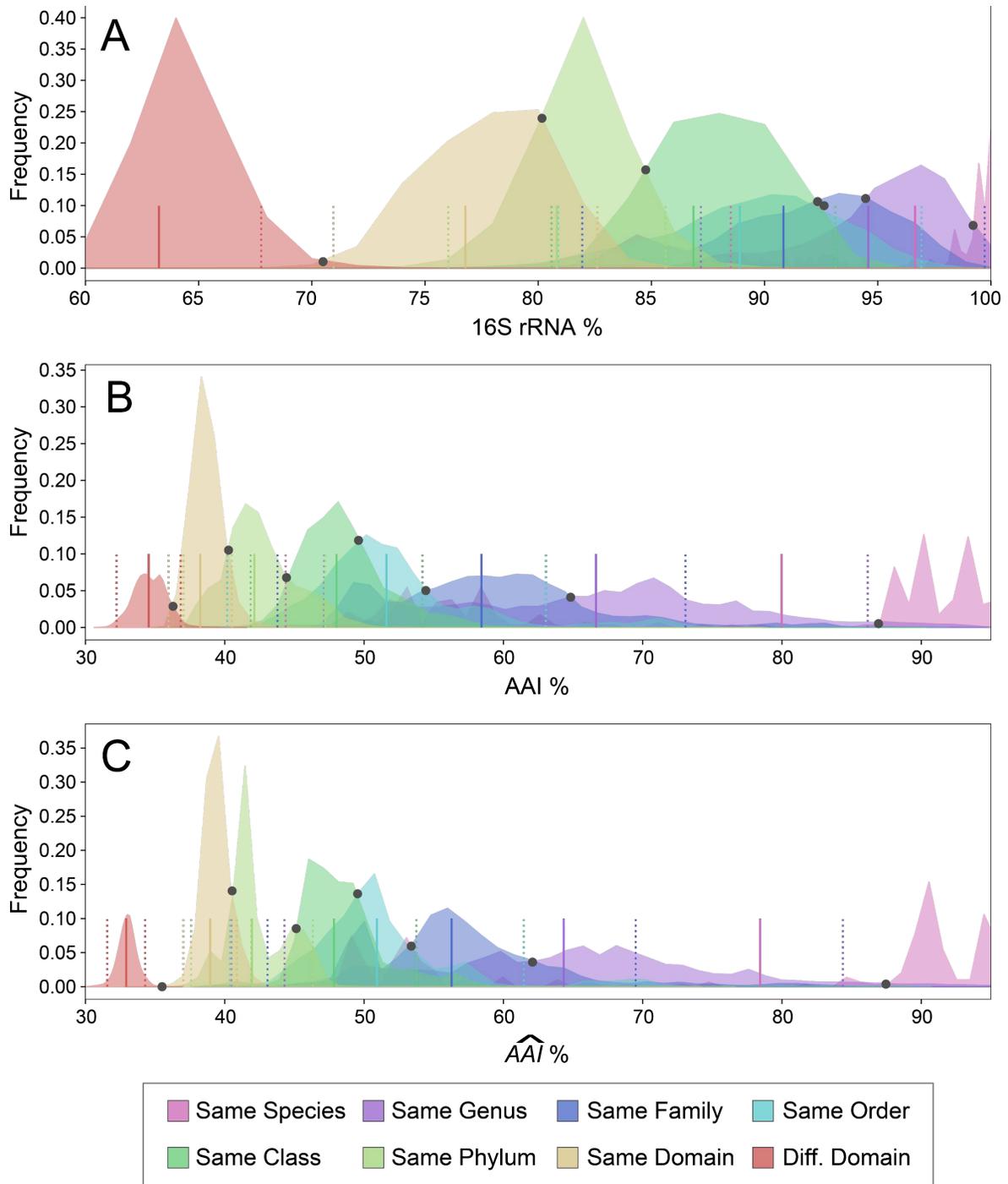
- 419 9. Buchfink B, Xie C, Huson DH. "Fast and sensitive protein alignment using DIAMOND" (2015).  
420 Nat Methods. **12**(1):59-60.
- 421 10. Kim, D., S. Park and J. Chun (2021). "Introducing EzAAI: a pipeline for high throughput  
422 calculations of prokaryotic average amino acid identity." J Microbiol **59**(5): 476-480.
- 423 11. Steinegger, M. and J. Soding (2017). "MMseqs2 enables sensitive protein sequence  
424 searching for the analysis of massive data sets." Nat Biotechnol **35**(11): 1026-1028.
- 425 12. Chaumeil, P. A., A. J. Mussig, P. Hugenholtz and D. H. Parks (2019). "GTDB-Tk: a toolkit  
426 to classify genomes with the Genome Taxonomy Database." Bioinformatics **36**(6):1925-  
427 1927.
- 428 13. Konstantinidis, K. T. and J. M. Tiedje (2005). "Towards a genome-based taxonomy for  
429 prokaryotes." J Bacteriol **187**(18): 6258-6264.
- 430 14. Bowers, R. M., N. C. Kyrpides, R. Stepanauskas, M. Harmon-Smith, D. Doud, T. B. K.  
431 Reddy, F. Schulz, J. Jarett, A. R. Rivers, E. A. Elie-Fadrosh, S. G. Tringe, N. N. Ivanova,  
432 A. Copeland, A. Clum, E. D. Becraft, R. R. Malmstrom, B. Birren, M. Podar, P. Bork, G.  
433 M. Weinstock, G. M. Garrity, J. A. Dodsworth, S. Yooseph, G. Sutton, F. O. Glockner, J.  
434 A. Gilbert, W. C. Nelson, S. J. Hallam, S. P. Jungbluth, T. J. G. Ettema, S. Tighe, K. T.  
435 Konstantinidis, W. T. Liu, B. J. Baker, T. Rattei, J. A. Eisen, B. Hedlund, K. D. McMahon,  
436 N. Fierer, R. Knight, R. Finn, G. Cochrane, I. Karsch-Mizrachi, G. W. Tyson, C. Rinke, C.  
437 Genome Standards, A. Lapidus, F. Meyer, P. Yilmaz, D. H. Parks, A. M. Eren, L. Schriml,  
438 J. F. Banfield, P. Hugenholtz and T. Woyke (2017). "Minimum information about a single  
439 amplified genome (MISAG) and a metagenome-assembled genome (MIMAG) of bacteria  
440 and archaea." Nat Biotechnol **35**(8): 725-731.
- 441 15. Huson, D. H., S. Mitra, H. J. Ruscheweyh, N. Weber and S. C. Schuster (2011).  
442 "Integrative analysis of environmental sequences using MEGAN4." Genome Res **21**(9):  
443 1552-1560.
- 444 16. Luo, C., R. L. Rodriguez and K. T. Konstantinidis (2014). "MyTaxa: an advanced  
445 taxonomic classifier for genomic and metagenomic sequences." Nucleic Acids Res **42**(8):  
446 e73.
- 447 17. Rinke, C., P. Schwientek, A. Sczyrba, N. N. Ivanova, I. J. Anderson, J. F. Cheng, A.  
448 Darling, S. Malfatti, B. K. Swan, E. A. Gies, J. A. Dodsworth, B. P. Hedlund, G. Tsiamis,  
449 S. M. Sievert, W. T. Liu, J. A. Eisen, S. J. Hallam, N. C. Kyrpides, R. Stepanauskas, E. M.  
450 Rubin, P. Hugenholtz and T. Woyke (2013). "Insights into the phylogeny and coding  
451 potential of microbial dark matter." Nature **499**(7459): 431-437.
- 452 18. Eddy, S. R. (2011). "Accelerated Profile HMM Searches." PLoS Comput Biol **7**(10):  
453 e1002195.
- 454 19. Lee, M. D. (2019). "GToTree: a user-friendly workflow for phylogenomics."  
455 Bioinformatics **35**(20): 4162-4164.

- 456 20. Hyatt, D., G. L. Chen, P. F. Locascio, M. L. Land, F. W. Larimer and L. J. Hauser (2010).  
457 "Prodigal: prokaryotic gene recognition and translation initiation site identification." BMC  
458 Bioinformatics **11**: 119.
- 459 21. Rodriguez-R, L. M., W. T. Harvey, R. Rosselló-Mora, J. M. Tiedje, J. R. Cole and K. T.  
460 Konstantinidis (2020). Classifying prokaryotic genomes using the Microbial Genomes  
461 Atlas (MiGA) webserver. Bergey's Manual of Systematics of Archaea and Bacteria: 1-11.
- 462 22. Nawrocki, E. P. and S. R. Eddy (2013). "Infernal 1.1: 100-fold faster RNA homology  
463 searches." Bioinformatics **29**(22): 2933-2935.
- 464 23. Madeira, F., Y. M. Park, J. Lee, N. Buso, T. Gur, N. Madhusoodanan, P. Basutkar, A. R.  
465 N. Tivey, S. C. Potter, R. D. Finn and R. Lopez (2019). "The EMBL-EBI search and  
466 sequence analysis tools APIs in 2019." Nucleic Acids Res **47**(W1): W636-W641.
- 467 24. Rodriguez-R LM, Konstantindis K. (2016). "The enveomics collection: a toolbox for  
468 specialized analyses of microbial genomes and metagenomes." PeerJ Preprints **4**:e1900v1.



470

471 **Figure 1.** Identity correspondence between AAI and  $\widehat{AAI}$  derived from the FastAAI value  $\bar{J}$  transformation  
 472 using the RefSeq dataset. Each dot represents a pairwise comparison colored by the lowest taxonomic rank  
 473 (figure key) that the two genomes in the pair share. Note the linear correspondence between both values  
 474 and the low error associated with the linear regression.



475

476 **Figure 2.** Distribution of 16S rRNA gene identities, AAI, and  $\widehat{AAI}$  values, according to the lowest  
 477 taxonomic level shared by the pairs of genomes analyzed. Underlying data are based on all vs. all  
 478 comparisons of the TypeMat genomes. Note the large overlaps between taxonomic ranks using all  
 479 methods, but also the virtually no overlap between the inter- (red distribution) and intra-domain (yellow  
 480 distribution) groups for  $\widehat{AAI}$ , which is consistent to the picture based on 16S rRNA gene identities. Solid  
 481 lines represent the means of the distribution while dotted lines represent one standard deviation around  
 482 the mean.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [25FastAAISupplementalText.docx](#)