1   17.09.2016

4   _____

5

# Comparative Genomics of 19 *Porphyromonas gingivalis* Strains

7

8   *Tsute Chen[1]\*, Huma Siddiqui[2] and Ingar Olsen[2]*

9   *[1] Department of Microbiology, The Forsyth Institute, Cambridge, MA, USA*

10  *[2] Department of Oral Biology, University of Oslo, Oslo, Norway*

11

12  Running head: Genomics in *P. gingivalis* strains

13

14  *\*Correspondence:*

15  Ingar Olsen

16  ingar.olsen@odont.uio.no

17  Abstract contains 233 words

18  Article contains 1,088 words and 6 Figures

19    Currently, genome sequences of a total of 19 *Porphyromonas gingivalis* strains are available,
20    including eight completed genomes (strains W83, ATCC 33277[T], TDC60, HG66, A7436, AJW4,
21    381, and A7A1-28) and 11 high-coverage draft sequences (JCVI SC001, F0185, F0566, F0568,
22    F0569, F0570, SJD2, W4087, W50, Ando, and MP4-504) that are assembled into fewer than 300
23    contigs. The objective was to compare these genomes at both nucleotide and protein sequence
24    levels in order to understand their phylogenetic and functional relatedness. Four copies of *16S*
25    *rRNA* gene sequences were identified in each of the eight complete genomes and one in the other
26    11 unfinished genomes. These 43 *16S rRNA* sequences represent only 24 unique sequences and
27    the derived phylogenetic tree suggests a possible evolutionary history for these strains.
28    Phylogenomic comparison based on shared proteins and whole genome nucleotide sequences
29    consistently showed two close relations groups: one consisted of ATCC 33277, 381 and HG66,
30    another of W83, W50 and A7436. At least 1,037 core/shared proteins were identified in the 19 *P.*
31    *gingivalis* genomes based on the most stringent detecting parameters. Comparative functional
32    genomics based on genome-wide comparisons between NCBI and RAST annotations, as well as
33    additional approaches, revealed functions that are unique or missing in individual *P. gingivalis*
34    strains, or species-specific in all *P. gingivalis* strains, when compared to a neighboring species *P.*
35    *asaccharolytica*. All the comparative results of this study are available online for download at
36    ftp://www.homd.org/publication_data/20160425/

37

39    **INTRODUCTION**

40    The Gram-negative anaerobic rod-shaped bacterium *Porphyromonas gingivalis* is one of the
41    most important pathogens in chronic adult periodontitis (Socransky et al., 1998) and has been
42    called a keystone pathogen (Darveau et al., 2012; Hajishengallis et al., 2012). This description
43    implies that it can cause dysbiosis (imbalance in the relative abundance or influence of species
44    within a microbial community) even when present at a low colonization level. *P. gingivalis* has
45    also been found to be related to systemic diseases, e.g., cardiovascular diseases, rheumatoid
46    arthritis and Alzheimer's disease (Demmer and Desvarieux, 2006; Lundberg et al., 2010; Olsen
47    and Singhrao, 2015). It has become increasingly clear that strains of *P. gingivalis* differ in their
48    pathogenicity, e.g., in their ability to invade tissues and cells varying as much as three orders of
49    magnitude (Dorn et al., 2000; Lundberg et al., 2010; Dolgilevich et al., 2011; Olsen and
50    Progulske-Fox, 2015). Thus W83 is considered a virulent strain while ATCC 33277 is
51    considered less virulent.  The AJW4 strain had the lowest invasion ability of 27 strains tested
52    (Dolgilevich et al., 2011). A comparative genomics study focusing on differences that affect
53    virulence in a mouse model identified over 150 divergent genes (Chen et al., 2004). Dolgilevich
54    et al., (2011) suggested deficiency in multiple genes as a basis for the *P. gingivalis* non-invasive
55    phenotype. Actually, more than 100 genes were missing from the genome of a non-invading
56    strain. The interstrain genomic polymorphisms and the individual host response have been

57  suggested to be the key to disease initiation and progression (Dolgilevich et al., 2011). Genomic
58  arrangement may also play a key role in the difference in virulence. For example, Naito et al.,
59  (2008) found that although the genome size and GC content were almost the same in strain
60  ATCC 33277 and W83 there were extensive rearrangements between the two strains. *P.*
61  *gingivalis* has been suggested to harbor many genetic mobile elements such as insertion
62  sequence (IS), miniature inverted-repeat transposable element (MITE) and conjugative
63  transposons CTns (Duncan, 2003, Naito et al., 2008, Tribble et al., 2013, Klein et al., 2015).

64  Together they are responsible for the fluidic genomic structure of this species (Naito et al.
65  2008, Tribble et al., 2013). The structural changes of the *P. gingivalis* genomes caused by these
66  elements might have generated many strain-specific protein-coding sequences (CDs) and may
67  have resulted in difference in various phenotype including important virulence factors (Naito et
68  al. 2008).

69  To date, a total of 19 *P. gingivalis* genome sequences have been published including eight
70  completed (strains W83, ATCC 33277[T], TDC60, HG66, A7436, AJW4, 381, and A7A1-28);
71  and 11 high-coverage draft sequences (JCVI SC001, F0185, F0566, F0568, F0569, F0570, SJD2,
72  W4087, W50, Ando, and MP4-504) that are assembled into fewer than 300 contigs. These strains
73  were isolated from various sources including the well-studied laboratory cultures with different
74  degree of virulence, clinical samples from patients with different disease states, as well as an
75  environmental strain isolated from a hospital bathroom sink drain. Together these sequences
76  provide a great opportunity for a comparative genomics study and the results will provide
77  valuable information to better understand the disease mechanism of this important periodontal
78  pathogen. The aim of this study was to conduct in-silico genomics comparison for theses
79  genomes using various approaches in the areas of phylogenetics, phylogenomics, and functional
80  genomics. Results that we found most important and interesting are presented in this paper
81  whereas complete results derived from this study are also made available for download online for
82  further investigation.

83

84  **MATERIAL AND METHODS**

85  **Sequence Sources**

86  Genomic sequences used in this study were downloaded from the NCBI FTP site
87  (ftp.ncbi.nlm.nih.gov/genomes/all). The versions that were downloaded are also available online
88  at ftp://www.homd.org/publication_data/20160425. A summary of all the meta information for
89  each genome is available in the Excel file PG_Genome_Summary.xlsx in the above FTP folder.
90  This file lists all the detail information that are provided by NCBI, such as methods for
91  sequencing, assembling and annotation, as well as various IDs for the same genome including
92  GenBank Accession,  GenBank Assembly Accession, Refseq Accession, Refseq Assembly
93  Accession.  **Table 1** lists the basic information and sources of the sequence data of the 19 *P.*
94  *gingivalis* genomes analyzed in this report.

95

## Strain Information

97      **W83, ATCC 33277[T] and W50.** These most-studied laboratory cultures were among the first
98      *P. gingivalis* strains sequenced.  Strain W83 was isolated in the 1950s by H. Werner (Bonn,
99      Germany) from an undocumented human oral infection and was brought to The Pasteur Institute
100     by Madeleine Sebald during the 1960s. It was subsequently obtained by Christian Mouton
101     (Quebec, Canada) during the late 1970s. W83 was reported to be also known as strain HG66
102     (Nelson, Fleischmann et al. 2003), however it has been shown that the two are very different
103     strains based on data shown in this report. Strain W50 was originally isolated from a clinical
104     specimen by H. Werner and first studied for known virulence (Marsh et al., 1994). W50 is also
105     known as ATCC 53978 based on the description of the BioSample ID SAMN00792205
106     (http://www.ncbi.nlm.nih.gov/biosample/?term=SAMN00792205). The strain ATCC 33277 used
107     for genomic sequencing was directly obtained from the American Type Culture Collection
108     (ATCC) and was described as "has been kept for more than 20 years" by the authors (Naito et
109     al., 2008).

110     **TDC60.** This strain was isolated from a severe periodontal lesion at Tokyo Dental College in
111     Japan. Strain TDC60 exhibited higher pathogenicity in causing abscesses in mice than strains
112     W83 and ATCC 33277 and other strains tested in the college (Watanabe et al., 2011).

113     **JCVI SC001.** This strain was not isolated from the human oral cavity; instead the genomic
114     sequence was derived from single cells found in the biofilm of a hospital bathroom sink drain.
115     The sequence was the first report of a human pathogen sequenced from a single cell captured
116     from an environmental sample outside of the human host. An automated platform was used to
117     generate genomic DNA by the multiple displacement amplification (MDA) technique from
118     hundreds of single cells in parallel. Thus the bacterial culture or DNA source of the genomic
119     sequence obtained through MDA cannot be made available (Information source:
120     http://www.ncbi.nlm.nih.gov/biosample/SAMN02436407, also see reference (McLeanet al.,
121     2013).

122     **Strains sequenced by HMP.** A total of six strains (F0185, F0566, F0568, F0569, F0570,
123     and W4087) were sequenced by The Genome Institute of Washington University collaborated
124     with the Data Analysis and Coordination Center (DACC) of the Human Microbiome Project
125     (HMP) and the Human Oral Microbiome Database and were funded by a consortium of institutes
126     including the National Human Genome Research Institute (NHGRI)/National Institutes of Health
127     (NIH), and the National Institute of Dental and Craniofacial Research (NIDCR). Strain F0568
128     and F0569 were isolated in the 1980s in the USA from the subgingival plaque biofilm of black,
129     non-Hispanic male subjects (53 and 39 years old respectively) diagnosed with moderate
130     periodontitis. F0570 was isolated in the 1980s in the USA from a 39 years old non-Hispanic
131     white male diagnosed with moderate periodontitis. Strain F0185, F0566 and W4087 were
132     reported to be isolated from the oral cavity/mouth of human subjects. Information source:
133     GenBank records in **Table 1**.

134     **SJD2.** This strain was isolated from subgingival plaque of a patient in China with chronic
135     periodontitis. It was shown to have high virulent properties comparable with those of the strain
136     W83 in a mouse abscess model. It was reported to have a higher number of SJD2-specific genes
137     which suggests that strains isolated from a periodontal pocket of Chinese patients with chronic
138     periodontitis may have distinct genes (Liu et al., 2014).

139     **HG66.** HG66 (also known as DSM 28984) was isolated in Roland R. Arnold's laboratory at
140     the Emory School of Dentistry, Atlanta, GA in the 1960s and was maintained in Jan Potempa's
141     laboratory since 1989. This strain was of interest because it does not retain gingipains on the cell
142     surface, instead releases the majority of proteases in a soluble form. In fact HG66 secretes all
143     carboxy terminal domain-bearing proteins as soluble substances. Information source:
144     http://www.ncbi.nlm.nih.gov/biosample/SAMN02732406 and (Siddiqui et al., 2014).

145     **A7436.** This strain was isolated from the subgingival plaque of the tooth abscess of a
146     refractory periodontitis patient by V.R. Dowell, Jr., at the Centers for Disease Control and
147     Prevention in Atlanta, GA, in the mid-1980s.  Information source:
148     http://www.ncbi.nlm.nih.gov/biosample/SAMN03366764.

149     **AJW4.** This strain was isolated from the subgingival plaque of the tooth abscess of a
150     periodontitis patient by R.J. Genco and colleagues in 1988 at SUNY-Buffalo, and described by
151     A. Progulske-Fox and colleagues as a minimally invasive strain during *in vitro* cell culture
152     studies. Information source: http://www.ncbi.nlm.nih.gov/biosample/SAMN03372093.

153     **Ando.** This strain was isolated from the gingival sulcus of a human oral cavity in Japan in
154     1985. The genome of this strain was sequenced because it was reported to express a 53-kDa-type
155     Mfa1 fimbrium, a major fimbrilin variant of Mfa1 previously known in many *P. gingivalis*
156     strains.  Information source: http://www.ncbi.nlm.nih.gov/biosample/?term=SAMD00040429
157     and (Nagano et al., 2015).

158     **381.** Strain 381 was isolated from the subgingival plaque of the tooth abscess of a localized
159     chronic periodontitis patient by S. Socransky, A. Tanner, A. Crawford and colleagues at the
160     Forsyth Dental Center (currently The Forsyth Institute), in the early 1970s. Information source:
161     http://www.ncbi.nlm.nih.gov/biosample/SAMN03656156.

162     **A7A1-28.** A strain isolated from subgingival plaque of the tooth abscess of a periodontitis
163     patient, with non-insulin dependent diabetes mellitus, by M.E. Neiders and colleagues in the
164     mid-1987 at SUNY-Buffalo, and was described as a virulent strain with atypical fimbriae and
165     capsule phenotypes.  Information source:
166     http://www.ncbi.nlm.nih.gov/biosample/SAMN03653671.

167     **MP4-504.** This strain is a low-passage (fewer than five passages) clinical isolate sampled
168     from the periodontal pocket (8 mm probing depth) of a chronic periodontitis patient at the
169     University of Washington Graduate Periodontics Clinic in 1991. The important characteristics of
170     this strain  include stable adherence to oral streptococci, enhanced invasion of gingival epithelial

171  cells (GECs), strong inhibition of IL-8 production by GECs, and the ability to transfer DNA by
172  conjugation at high efficiencies (To et al., 2016).

173

174  **Data Analysis**

175    **16S rRNA phylogeny.** For the 16S rRNA gene phylogeny, *16S rRNA* gene sequences were
176  extracted from the genomes of the 19 *P. gingivalis* strains based on NCBI's annotation (the
177  *genomic.gff file in each of the downloaded genome folder). Sequences were pre-aligned with
178  MAFFT v6.935b (2012/08/21) (Katoh and Standley, 2013) and leading and trailing sequences
179  not present in all sequences were trimmed. The trimmed and aligned sequences, with an
180  alignment length of 1,425 bases and representing 20 unique sequences, were subjected to
181  QuickTree V 1.1 (Howe et al., 2002) using the "-kimura" option to calculate the substitution rate.
182  A copy of the 16S rRNA gene sequence from *Porphyromonas asaccharolytica* (PaDSM20707)
183  was used as the out-group during the phylogenetic tree construction.

184    **Core and unique proteins.** To study the phylogenetic relationship based on more
185  genes/proteins, protein sequences annotated by NCBI were used. Together with the outgroup
186  species PaDSM20707, a total of 41,625 proteins were annotated by NCBI, including 39,926
187  from the 19 *P. gingivalis* genomes and 1,699 from PaDSM20707. Of the 39,926 *P. gingivalis*
188  proteins, 37,667 are ≥ 50 amino acids in length and were searched for homologous clusters using
189  the "blastclust" software V.2.2.25
190  (http://www.ncbi.nlm.nih.gov/Web/Newsltr/Spring04/blastlab.html). Various sequence identity
191  cutoffs ranging from 10 to 90% and two minimal alignment length cutoffs 50% and 90% were
192  used for identifying the protein clusters. Proteins in each set of the identified clusters were
193  aligned with MAFFT and poorly aligned regions were filtered by Gblocks 0.91b (Talavera and
194  Castresana, 2007). Trees were constructed with FastTree 2.1.9 (Price et al., 2010) using the JTT
195  protein mutation model (Jones et al., 1992) and CAT+–gemma options to account for the
196  different rates of evolution at different sites. The reliability of tree splits were reported as "local
197  support values" based on the Shimodaira-Hasegawa test (Shimodaira and Hasegawa, 2001). For
198  comparison, all 41,625 proteins were also subject to the PhyloPhlAn software (Segata et al.,
199  2013) version 0.99 (8 May 2013).

200    To identify proteins that are unique for each genome, all the 39,926 *P. gingivalis* proteins
201  were searched against each other using BLASTP 2.2.25 with default parameters (Altschul et al.,
202  1997). Those that did not match any other protein with expected *e* value ≤ 10 were considered
203  unique among the 19 genomes.

204    **Whole genome nucleotide comparisons.** Pairwise whole genome nucleotide to nucleotide
205  sequence alignment were plotted using NUCmer (NUCleotide MUMmer) version 3.1 (Delcher et
206  al., 2002). To compare the whole genome DNA similarity by the oligonucleotide frequency, all
207  possible 20-mer sequences present in the 20 genomes, including that of *P. asaccharolytica* strain
208  DSM 20707 used as an out-group, were categorized and the number of genomes in which a 20-

209 mer was present was recorded. Any given oligonucleotide can have a maximum of 20 (i.e.,
210 present in all 20 genomes) and a minimum of 1 (unique, found in only a single genome). To plot
211 the oligonucleotide frequencies, an overall frequency for every 500 bases across the entire
212 genome was calculated by recording the total number of genomes that all the possible 20-mer in
213 the 500 bases can be found in (maximal 20, minimal 1). Each of the 500 bases windows was
214 colored based on the genome frequency. Another plot was created similarly except that the non-
215 coding regions were masked with light blue color to highlight the oligonucleotide frequencies for
216 the areas that correspond to both forward (upper) and reverse-complement (lower) protein coding
217 sequences.

218     **Comparative functional genomics.** Three functional annotation systems were used and
219 compared in this study for all the 20 genomes – 1) the NCBI prokaryotic genome annotation
220 pipeline (Tatusova et al., 2016),  2) the SEED and RAST (Rapid Annotation using Subsystem
221 Technology) (Overbeek et al., 2014), and 3) the KOALA (KEGG Orthology And Links
222 Annotation) (Kanehisa et al., 2016). The NCBI annotation results were downloaded from the
223 NCBI FTP site described in the Sequence Sources above. The genomic DNA sequences were
224 sent to the SEED server (Aziz et al., 2012) using the Linux command-line and network-based
225 SEED API downloaded from the SEED server web site
226 (http://blog.theseed.org/servers/installation/distribution-of-the-seed-server-packages.html). The
227 NCBI annotated proteins were sent to the BLastKoala web site (http://www.kegg.jp/blastkoala)
228 to identify the KEGG Orthologs. The results of both NCBI and RAST annotations were
229 compared by several text based keyword searches. To identify more proteins in a particular
230 functional category that were somehow annotated in certain genomes but not in others, protein
231 sequences that were annotated in the same category from all 20 genomes were collected and used
232 as the query to search for more proteins of the same functional category. NCBI BLASTP was
233 used for this purpose and proteins with $\geq$ 95% sequence identity to and $\geq$ 95% coverage of the
234 query sequences were identified as highly similar proteins. The number of proteins related to the
235 IS5 transposase family was identified by the BlastKOALA program (Kanehisa et al., 2016) with
236 the matching to the KEGG Orthology (KO) number K07481. Additional functional comparison
237 results were also made available as several files in Excel format.

238 **Data and Results Availability**

239     To facilitate further comparison and future studies, all the data and results generated in this
240 study,  including the original downloaded  sequences, annotations, the comparative results
241 presented in this paper, as well as additional complete results that were not mentioned or
242 discussed,  are available for download from this FTP data repository site:
243 ftp://www.homd.org/publication_data/20160425 .

244

245 **RESULTS AND DISCUSSION**

246 **Summary of Genome Annotations**

247      The first *P. gingivalis* genome released was that of the strain W83 in 2003 and the latest one
248    was released in February 2016. Of the 19 genomes, eight were assembled into a single contig and
249    were considered complete and finished genomes; the remaining were released as various
250    numbers of sequence contigs assembled from whole genome shotgun (WGS) sequence reads.
251    The sequence of JCVI SC001 appears to have a 1-contig circular sequence under the Genbank
252    Accession number CM001843, however it is a pseudo-contig generated by ordering the 284
253    unassembled contigs (accession number APMB01000000) based on the homologous matches to
254    the genome of TDC60 (McLean et al., 2013) and joining the ordered contigs with 282 100-N
255    spacer sequences (total N length is 28,200 bps). Thus it is not considered a complete or finished
256    genome. Examining the sequences for the presence of Ns reveals the "completeness" of the
257    genomes. **Table 2** shows the reported length, non-N length, total number of Ns and the
258    distribution of the N fragments in the genomic sequences. Overall strain A7A1-28 is the smallest
259    of the completed *P. gingivalis* genomes with a size of 2,249,024 bps. HG66 has the largest size
260    of all the sequenced *P. gingivalis* genomes at 2,441,680 bps after removing the 100 Ns placed at
261    the end of the sequence. The placement of the 100 Ns at the end of the sequence was due to the
262    unsuccessful attempt to circularize the sequence with the minimus2 software used by the PacBio
263    sequencer at default settings (personal communication).  For this reason the HG66 genome
264    should not be considered complete. Almost all the unfinished draft genomes consist of various
265    numbers of Ns ranging from 698 Ns in SDJ2 to 7,200 Ns in F0569 (**Table 2**). It is likely that
266    some of these published contigs were assembled based on a reference genome and the Ns had
267    been filled in the gaps. Hence the true order of genes identified by the annotation process may
268    not be correct.

269      **Table 3** gives a numeric summary of the genome annotation results by the NCBI Prokaryotic
270    Genome Annotation Pipeline (released 2013,
271    http://www.ncbi.nlm.nih.gov/genome/annotation_prok/). The NCBI pipeline is capable of
272    identifying more than just the protein-coding genes, rRNAs and tRNAs, including several
273    interesting types of genes such as binding sites, repeat sequences, pseudo-genes, and several
274    types of non-coding RNAs (ncRNAs). However, since the NCBI pipeline is quite new, more
275    features are still being added and since some of the annotations of these *P. gingivalis* genomes
276    were done prior to 2013, the annotation results may not be comprehensive until the annotation is
277    updated again based on the latest NCBI pipeline.

278      In addition to the NCBI annotations, RAST (Rapid Annotations using Subsystems
279    Technology) is also a popular pipeline for annotating microbial genomes (Aziz et al., 2008). All
280    the 19 *P. gingivalis* genomes, as well as the chosen outgroup *P. asaccharolytica* DSM20707
281    were subjected to the RAST pipeline and the results were compared with those done by the
282    NCBI pipeline. As shown in **Table 4**, both the RAST and NCBI pipelines identified almost the
283    same number of rRNA and tRNA genes. However the numbers of protein-coding genes varied
284    quite significantly between the two pipelines. Although most of the genes were commonly
285    identified, up to hundreds of protein-coding sequences can be missed by either system.
286    Moreover, 86% (6,422 of 7,382 for all the 19 genomes) of these uniquely identified genes code
287    for hypothetical proteins and 80% are shorter than 100 amino acids in length (only 94 have
288    lengths ≤ 500 amino acids), thus the impact due to the annotation discrepancy may not be as

289     significant especially when drawing conclusions in genome-wide systematic analysis or
290     metabolic pathway capability.

291       A list of the 960 (7,382-6,422) non-hypothetical proteins is provided at the link
292     (ftp://www.homd.org/publication_data/20160425/2_Summary_of_Genome_Annotations/Non-
293     overlap_Non-hypothetical_protein_identified_by_NCBI_or_RAST.fasta).

294

**16S rRNA Phylogeny**

296       The 16S rRNA sequences have been used to infer the evolutionary relatedness of the
297     prokaryotes due to its slow rate of evolution (Woese et al., 1990). However multiple
298     rRNA genes including *16S rRNAs* are common in prokaryotic genomes  (Klappenbach et al.,
299     2000) and the genomic copy number of *16S rRNA* varies greatly among species from 1 to 15
300     (Vetrovsky and Baldrian, 2013). The number of *rRNA* genes was reported to correlate with the
301     rate at which phylogenetically diverse bacteria respond to resource availability (Klappenbach et
302     al., 2000). As shown in **Table 4**, all of the eight genomes which had been assembled to a single
303     contig contain four copies of *5S*, *16S* and *23S rRNA* genes respectively, thus it is reasonable to
304     believe that all *P. gingivalis* genomes have four copies of the rRNA operons. The lower number
305     of *rRNA* genes in the unfinished genomes is likely due to the incompleteness of the sequences
306     and is also likely due to the fact that genomes sequenced by short reads sequencing platforms
307     such as those of the Illumina sequencers cannot be easily assembled across the repeated regions
308     such as the highly conserved rRNA operons.

309       The *16S rRNA* sequences of all the 19 genomes annotated by NCBI were extracted and
310     aligned for the construction of a phylogenetic tree. Based on the annotation, there are a total of
311     24 unique *16S rRNA* gene sequences identified from the 19 genomes (**Table 5**, first column)
312     excluding the sequence of a close species *P. asaccharolytica* strain DSM 20707 (Accession
313     Number CP002689).  However, many of the sequence differences are due to different annotated
314     lengths. After aligning all the 24 unique sequences and trimming off the leading and trailing
315     sequences not present in all copies (trimmed aligned length = 1,425 bps), the aligned portion of
316     several sequences are identical and the number of unique sequences was reduced to 20 (second
317     column of **Table 5**). Strains 381, A7A1-28, ATCC 33277, and W83 all have four copies of
318     identical sequences and those of ATCC 33277 and 381, as well as three copies of HG66 shared
319     identical aligned/overlapping sequences. Strain A7436 shared three of its four copies of *16S*
320     *rRNA* sequences identically with those of W83. Together with the single copy from W50, they
321     formed an identical group of sequences. W50 has been known to be a close strain of W83, thus
322     the identical sequences between these two are not surprising. The explanation of identical copies
323     of the *16S rRNA* sequence in the genome is apparently due to the gene duplication event and the
324     fact that several strains shared identical duplicated sequences suggested that the duplication
325     event occurred after the speciation. Strains A7436, AJW4 and HG66 had three strain-specific
326     identical sequences with the 4[th] copy different from the other three.  Overall, all the *P. gingivalis*
327     *16S rRNA* gene sequences were extremely similar and often have only a single number of

328　nucleotide mismatches between any two strains (if not identical). Altogether only 16 loci on the
329　gene had nucleotide variation, with the exception of one copy in TDC60, which had a series of
330　A→C or G→C transversions between position 50 and 130. These aligned and trimmed
331　sequences, including the outgroup sequence from *P. asaccharolytica* strain DSM 20707, were
332　used to construct a phylogenetic tree based on Kimura's nucleotide substitution model and the
333　result is shown in **Figure 1**. The phylogenetic tree depicts a likely evolutionary path for these
334　different *P. gingivalis* strains. The strains 381, ATCC 33277 and HG66 appeared to be closer to
335　the potential common ancestor, based on the tree topology inferred with a close species as the
336　outgroup sequence. The other strains gradually diversified into deeper branching nodes with two
337　of the sequences from strains F0566 and TDC60 as the most deeply branched and mutated from
338　the common ancestor, which was inferred by using the sequence of a neighboring species
339　(PaDSM20707) as an outgroup.

340

**Core and Unique Proteins**

342　　The phylogenetic relationship inferred based on the *16S rRNA* gene sequences reported
343　above can only represent the evolution of this particular gene, hence a gene tree. A more
344　comprehensive way of studying the evolutionary relatedness of different genomes is to use as
345　much genomic information as possible in the analysis (i.e., phylogenomics). A popular approach
346　is to use the core proteins for the construction of a tree that may be closer to a true species tree, if
347　such a tree exists, or if there is no true species tree, may reflect more on the relatedness of these
348　strains at the genomic level. The concept of "core" proteins is ideally defined as proteins that are
349　present and required by all the genomes in study, however the identification of such group of
350　proteins, namely orthologues, is not straightforward and the results vary depending on the
351　criteria used. It is challenging, if not impossible, to identify all the orthologous proteins among a
352　group of genomes. In general, genomes of closer species or strains share more orthologues;
353　however any percent protein sequence identity chosen as the cutoff to test whether a group of
354　homologous proteins are truly orthologues (or paralogues) can always include some false
355　positive and negative orthologues. Nevertheless, one can still hypothesize that a more reliable
356　evolutionary relationship of a group of genomes can be obtained if the use of higher or lower
357　percent identity constrains does not affect the overall tree topology.

358　　To test this hypothesis, the "core" proteins were first identified among all *P. gingivalis*
359　genomes under different cutoffs. Based on the NCBI annotation, a total of 39,926 protein
360　sequences were identified. However for some unknown reasons, some of the annotated protein
361　lengths were as short as one or two amino acids. For example, proteins with Genbank IDs
362　GAP82138.1, GAP81676.1, and GAP81848.1 in strain Ando were identified with only 1, 2 and 2
363　amino acids in length respectively. These clearly were annotation errors caused by the
364　computational bugs in the annotation pipeline. In this analysis, only protein sequences with a
365　minimal length of 50 amino acids were used (a total of 37,667 proteins) for identifying core
366　proteins. They were subject to the "blastclust" program
367　(http://www.ncbi.nlm.nih.gov/Web/Newsltr/Spring04/blastlab.html) to identify clusters of

368 proteins that share a certain degree of sequence homology and with specified alignment length
369 coverage. In this analysis, if a protein is present (i.e., meets the % identity and alignment cutoffs
370 specified) in all 19 genomes (or 20 genomes if PaDSM20707 was used as the outgroup in some
371 results) it is considered as a core/shared protein, and if a protein is only present in a single
372 genome it is considered as a strain-specific unique protein.

373 **Figure 2** shows the potential numbers of both core and unique proteins in the 19 genomes
374 analyzed with "blastclust" by varying two parameters: sequence percent identity cutoffs (from
375 95% to 10%) and percent alignment length (90% and 50%). **Figure 2A** shows that regardless of
376 sequence identity cutoffs; the number of core proteins stays relatively constant around 1,000 with
377 90% as the alignment length cutoff. The number of core protein groups increased gradually from
378 1,037 at 95% identity, maximized at 1,045 at 60%, then decreased to 910 at 10%. The reason for
379 the increase from 95 to 60% was due to more core protein groups clustered together at lower %
380 identity. The decrease after 60% identity was due to the fact that different protein groups
381 identified with identity ≥ 60% began to merge into fewer groups. The 1,037 shared proteins were
382 detected under most stringent conditions thus it is reasonable to state that at least 1,037 core
383 proteins were detected based on 19 strains. This number is expectedly smaller than the 1,476
384 detected in the core genome based on 8 *P. gingivalis* strains (Brunner et al., 2010). It should also
385 be noted that the core genes/proteins are not the same as the "essential" genes, for which only ca.
386 400 were experimentally detected previously (Klein et al., 2012).

387 For the purpose of identifying a core/shared set of proteins for constructing a phylogenomic
388 tree, the 1,042 core proteins identified at 60% sequence identity and 90% alignment length
389 cutoffs were used for sequence alignment and tree building. This set of sequences is available for
390 download in the data repository FTP site mentioned in the Material and Methods. In addition, as
391 expected when the percent alignment length was decreased from 90% to 50%, more proteins
392 were identified as core proteins, e.g., from 1,289 at 95% identity cutoff to 1,301 at 60%, due to
393 the fact that more proteins share the same percent identity over shorter sequence length.

394 **Figure 2B** shows the number of protein groups that are shared by 2 to 18 genomes (sub-core
395 proteins). The number decreases with the lower % identity because similar protein groups that
396 were identified as separate groups merged into a single (but larger) group due to the more
397 relaxed (lower) % identity (e.g., from 1,927 at 95% to 1,651 at 10%). However, contrary to the
398 core proteins above, which require proteins present in all 19 genomes, when the percent
399 alignment decreased, fewer sub-core proteins were identified. This is as expected because when
400 the percent alignment cutoff was lowered, a protein group which consists of only members from
401 for example 18 genomes, at higher cutoff, now may find a member in the 19[th] genome thus
402 disqualifying it as the 18-genome only sub-core group.

403 **Figure 2C** shows the number of strain-specific proteins that are present in only one genome
404 with a single copy. Similar to the sub-core groups, as the % identity decreases the number of
405 unique proteins becomes smaller because more proteins from different genomes were lumped
406 together as a homologous group under a lower % identity, resulting in the loss of the
407 "uniqueness". In addition, for example, at 60% sequence identity and 90% alignment cutoffs,

408    there were 2,289 proteins identified as present in a single genome, but the number was reduced
409    to 1,044 at 50% alignment cutoff – 1,245 proteins lost their uniqueness due to the presence of
410    more "similar" proteins found in other genomes.

411        For the unique proteins identified, it would be interesting to observe their distribution in the
412    19 genomes and the result may help understand which genomes possess more or fewer unique
413    proteins. **Figure 3** shows the distribution of the 1,044 unique proteins identified with the 50%
414    alignment cutoff (**Figure 3A**) and 2,289 with 90% cutoff (**Figure 3B**). Regardless of the
415    sequence identity and percent alignment cutoff, the results show that some strains possess
416    significantly more unique proteins than others. Four strains, F0566, F0568, F0569, and JCVI
417    SC001 have a significantly higher number of unique proteins under all identification conditions -
418    as high as 96–249 unique proteins (for percent identity 10% - 95% at 50% alignment length) in
419    the case of F0566 (**Figure 3A**). On the other hand, strains 381, ATCC 33277, A7436 and W83
420    are the four strains with the lowest number of unique proteins, only 10-15 unique proteins (10-
421    95% identity; 50% alignment) in the case of W83 (**Figure 3A**). Interestingly, strain W50, the
422    closest strain to W83, encodes more unique proteins (30 – 46) even though it is an unfinished
423    draft genome. Apparently the incompletes of the draft genomes are not the cause for the
424    difference in the number of unique proteins (JCVI SC001 and all the F strains are draft
425    genomes). This further suggests that the gaps in the draft genomes most likely contain only
426    repeated sequences that either do not encode for proteins, or encode repeated proteins that do not
427    contribute much to the genome's uniqueness.

428        Another noteworthy observation is that there is a consistent gap between the data points 95%
429    and 90% identity when searching for unique proteins in all strains under both alignment
430    conditions (**Figure 3**). This suggests that the more proteins identified as unique at 95% became
431    "similar" at 90%. Hence 90% sequence identity may be an ideal cutoff for differentiating
432    homologues and unique proteins, at least at the strain level.

433        **Table 6** lists the percentage of proteins that were annotated in NCBI as the hypothetical
434    proteins (functionally unknown) and the percentage of the unique proteins that were identified
435    with 80% as the sequence identity and 50% alignment length. The total percent hypothetical
436    proteins range from 26% (W50) to as high as 46% (F0566 and F0568), whereas the majority of
437    the unique proteins are hypothetical, from 68% (TDC60) to 100% (W83). Thus until more
438    functions of the hypothetical proteins are understood, it will still be challenging to understand
439    what each genome's overall "specialty" functions conferred by the unique genes. To give a
440    glimpse of what each genome's most unique functions are, based on the currently available
441    information, **Table 7** lists the functional annotations of the non-hypothetical unique proteins for
442    each of the 19 genomes (with default BLASTP parameter, i.e., expected e value ≤ 10) ( Altschul
443    et al., 1997). All of them are among the proteins identified above under the most stringent
444    parameters in terms of uniqueness – 50% sequence identity and 50% alignment length. Strain
445    JCVI SC001, an environment isolate from a hospital sink drain, has the most diverse functions
446    encoded by these unique proteins. Whether these annotations translate to unique functions of the
447    genome, require further investigation to ensure there are no other non-homologous proteins that

448 play similar functions.  All of the unique proteins identified are available by strain in the FTP
449 data repository.

450

451 **Phylogenomics by Homologous Proteins**

452     Once a group of putative core proteins is identified, they can be concatenated and aligned
453 together and used for compiling a phylogenomic tree to infer a possible evolutionary relationship
454 at a level closer to the species than just any single gene. In this analysis, the 1,045 proteins
455 shared by all 19 genomes at 60% sequence identity and 90% alignment cutoffs (**Figure 1A**),
456 were first aligned individually with the "mafft" software (Katoh and Standley, 2013). Each of the
457 1,045 protein sets contained exactly 19 aligned sequences, one from each of the 19 genomes. The
458 aligned proteins were concatenated in the same protein order. This generated a set of 19 mega
459 protein sequences with each consisting of 1,045 concatenated aligned sequences. The poorly
460 aligned sequence regions, including leading and trailing unaligned portions of the sequences, as
461 well as low-confidence parts of the alignment, such as positions that contain many gaps, were
462 removed with the "Gblock" tool (V 0.91) (Talavera and Castresana, 2007). After the Gblock
463 screening, a final set of 19 aligned protein sequences, each with a length of 395,174 amino acids
464 were used for constructing an unrooted tree. However among the 395,174 aligned amino acids,
465 only 17,389 positions had at least two different amino acids across proteins of all 19 *P. gingivalis*
466 genomes, the remaining 377,785 were all the same amino acids across all genomes. Thus only
467 those 17,389 informative or effective positions contributed to the pairwise distances calculated
468 among all genomes. **Figure 4A** is the result of the unrooted tree compiled based on the 1,045
469 shared proteins processed as described above. The overall topology is quite different from that of
470 the *16S rRNA* tree (**Figure 1**) with the exception of two very closely related groups of strains,
471 one consists of strains 381, ATCC 33277, and HG66 and another A7436, W50 and W83. This is
472 not surprising because both groups have members with identical *16S rRNA* sequences hence their
473 shared protein sequences are closer to each other in the group than other genomes.

474     To test whether including proteins from the outgroup species will result in a tree more similar
475 to that of *16S rRNA*, i.e., a tree that is rooted at a potential common ancestor for these strains,
476 orthologue candidates were first identified from the genome of *P. asaccharolytica* DSM 20707,
477 of which the *16S rRNA* sequence was also used for the *16S rRNA* tree. At 90% alignment length
478 cutoff, the number of homologous proteins in *P. asaccharolytica* decreases as the percent
479 sequence identity cutoff increases. The numbers of protein homologous to any of the 1,045 core
480 proteins used for the unrooted tree above are 0, 1, 7, 36, 146, 271 and 436 respectively for
481 percent identity cutoffs 95, 90, 85, 80, 70, 60 and 50%. **Figure 4B** and **C** are the two rooted
482 phylogenetic trees constructed based on the 36 (80% identity) and 436 (50% identity) proteins
483 shared between *P. asaccharolytica* DSM 20707 and all 19 *P. gingivalis* strains. After Gblocks
484 screening, the length of the aligned sequences were 12,646 (80% identity) and 177,272 (50%
485 identity) amino acids respectively and the number of effective amino acids positions are 154 and
486 4,771 respectively. In general, the branch lengths increased with more effective amino acids
487 positions which resulted in greater distances. Again, the only consistent close clusters were the

488  two grouped with identical *16S rRNA*, i.e., the group of 381, ATCC 32277 and HG66, and of
489  A7436, W50 and W83.

490  **Figure 4D** is the rooted tree constructed using the software PhyloPhlAn (Segata et al., 2013)
491  version 0.99 (8 May 2013). All 41,625 proteins annotated for the 20 genomes were subject to
492  PhyloPhlAn with the default parameters that excluded proteins shorter than 30 amino acids in
493  length. PhyloPhlAn finds among the input protein matches to a pre-set of the 400 most conserved
494  proteins for extracting the phylogenetic signals. A total of 264 query proteins were matched to
495  the 400 preset core but only 225 were present in all 20 genomes. These proteins were then
496  aligned individually and subsampled based on a sophisticated procedure provided by
497  PhyloPhlAn, which emphasizes regions both universally conserved and phylogenetically
498  discriminating. The final aligned, subsampled, and concatenated sequences had a length of 3,082
499  aligned amino acids with 840 effective positions. The PhyloPhlAn tree is shown in **Figure 4D**.
500  Similar to the two rooted trees (**Figure 4B** and **C**) and the 16S rRNA tree (**Figure 1**) the
501  PhyloPhlAn tree also placed the three strains ATCC 33277, 381 and HG66 closest (but much
502  closer) to the root and the remaining strains in a more linearly nested topology.

503  In summary, the only consensus based on interpretation of the three rooted protein trees and
504  the *16S rRNA* tree is that the group ATCC 33277, 381 and HG66 is less evolved and closest to
505  the common ancestor of this species (inferred based on the distance to the root). Strains W83,
506  W50 and A7436 consistently formed a close group regardless of how the trees were built, but
507  their exact phylogenetic position is inconclusive based on these analyses. The more
508  effective/informative aligned amino acid positions resulted in longer branches and pairwise
509  distances. To this end, the unrooted tree (**Figure 4A**) has the best resolution to reveal the
510  similarity/differences among these strains, in the most genome-wide manner. Until a group of
511  true orthologous proteins are identified (together with the outgroup) a true phylogenetic tree that
512  infers the evolutionary path for this species will not be accessible.

513

514  **Comparisons Based on Whole-Genome Nucleotide Sequences**

515  **1. MUMmer/NUCMER nucleotide plots**

516  Moving up the scale for comparison, one possible way is the whole genome nucleotide
517  alignment with a commonly used software MUMMER, which identified MUMs – minimal
518  unique matches between two genomic sequences (Delcher et al., 2002). **Figure 5** shows some of
519  the pairwise alignment results of the 19 *P. gingivalis* genomes. **Figure 5A** is the nucleotide
520  alignment between strains 381 and ATCC 33277 and the almost perfect diagonal high similarity
521  (red) match line indicates highly similar sequences, with only two visible exceptions – one
522  inversion and one insertion (to 381)/deletion (to ATCC 33277). Interestingly the inverted
523  sequence almost matches the inserted sequence; apparently the inverted sequence was duplicated
524  in the 381 genome and inserted somewhere else in the genome, where the ATCC 33277 genome
525  shows no counterpart. The high DNA sequence similarity between 381 and ATCC 33277 is also

526 supported by the identical *16S rRNA* gene sequence and copy numbers (**Figure 1**) as well as the
527 protein-based phylogenetic relationships (**Figure 4**), even though their genomes are not far from
528 identical. The phenomenon that a fairly large chunk of genomic sequence was duplicated and
529 inserted elsewhere in the genome is only observed in strain 381, as evidenced by the MUMMER
530 self-alignment of its genome (data available from the FTP site), but very similar to the alignment
531 between 381 and ATCC 33277). No duplication event was observed in the self-alignment of the
532 other 18 genomes.

533     Strain HG66 is the genome that is closest to 381 and ATCC 33277 based on *16S rRNA* genes
534 and protein sequences, on the other hand it shows the disconnected high similarity match lines,
535 which indicates more large-scale genomic arrangement between the two close strains – between
536 381 and HG66 (**Figure 5B**) and between ATCC 33277 and HG66 (**Figure 5C**).

537     The second closest groups of strains are A7436, W50, and W83 and their nucleotide
538 sequences are also highly similar based on the NUMMER plots (**Figure 5D** and **E**). However the
539 contigs of the unfinished draft genome of W50 were rearranged by MUMMER in the order
540 based on the similarity to the W83 sequence. Whether there is a large scale genomic
541 rearrangement between W83 and W50 cannot be known until the genome of W50 is completed.
542 Strain A7436, a finished genome, shows only one inversion of the genome when compared to
543 that of W83 (**Figure 5E**). The fact that A7436 is not as close to W50 and W83 as the distance
544 between HG66 and 381 (or ATCC 33277) based on *16S rRNA* and protein phylogeny (**Figure 1**
545 and **4**), suggests that the genomes of the group of HG66, 381 and ATCC 33277 have  higher
546 genomic sequence rearrangement activity than the A7436-W50-W83 group. The next genome
547 which is closest to the A7436-W50-W83 group is strain AJW4, with several visible (larger
548 fragments) of insertions/deletions and inversions when compared to A7436 (arrows heads in
549 **Figure 5F**). This relationship is also consistent with the *16S rRNA* gene tree (**Figure 1**).

550     Another interesting observation is the alignment between JCVI SC001 and TDC60. These
551 two strains are not among the closest groups based on the *16S rRNA* and protein sequences
552 (**Figure 1** and **4**). The MUMMER plot between these two genomes appears to be a straight
553 diagonal red line (**Figure 5G**), similar to that between 381 and ATCC 33277. However, since the
554 genomic sequence of JCVI SC001 was not really completed and closed to a circular
555 chromosomal format, the 284 *de novo* assembled contigs were mapped to the genome of TDC60
556 and the gaps were filled with Ns to form a single pseudo-contig (Genbank Accession
557 CM001843) (McLean et al., 2013). Thus the contig order in the published single contig genomic
558 sequence of JCVI SC001 may not be correct and the sequence similarity between JCVI SC001
559 and TDC60 may not be as "straight" as indicated in the MUMMER plot. In fact when the plot
560 was filtered to show only the region with percent identity $\geq 99\%$, the red line became fragmented
561 with large gaps (**Figure 5H**), indicating that  a large portion of the genomic sequences between
562 these two strains are under 99% similarity.

563     The complete pair-wise MUMMER plots of the 19 *P. gingivalis* genomes can be viewed in a
564 specifically designed interactive webpage at
565 http://bioinformatics.forsyth.org/publication/20160425. The web page provides interactive tools

566   to choose any two *P. gingivalis* genomes for the MUMMER results, as well as the possibility of
567   viewing the alignment at various percent sequence identity cutoffs.

568   **2. Oligonucleotide frequency**

569       The MUMMER plots above are limited to viewing comparisons only between two given
570   genomes. To view and compare nucleotide difference/similarity for all genomes on the same
571   plot, the overall oligonucleotide composition and frequency can be measured along the entire
572   genome and the results can be plotted out and visually compared to each other. This analysis
573   started by collecting all the possible 20-mer sequences in all 20 genomes and then count for each
574   20-mer how many genomes have each particular sequence. The number of genomes (genomic
575   frequency) for each 20-mer thus ranges from 1 (unique) to 20 (universal). The frequencies can be
576   calculated and plotted along the entire genome by taking every 20-mer from the beginning to the
577   end of the genome. **Figure 6A** depicts the results of the 20-mer oligonucleotide frequencies
578   among all 20 genomes (including the out-group *P. asaccharolytica* DSM 20707).  If a region of a
579   genome is shared by all other 20 genomes, it is colored black; and if a region is unique to the
580   genome itself, it is colored bright yellow. In other words, a black region means that all the
581   possible 20-mer sequences appeared in all tested genomes, whereas the brightest yellow regions
582   have unique 20-mer sequences that are only found in one genome. For easy comparison, the
583   order or the genomes shown in **Figure 6A** and **B** were arranged according to that of the *16S*
584   *rRNA* tree (**Figure 1**) with a dendrogram reflecting similar tree topology. As expected, the two
585   closest strains ATCC 33277 and 381 share almost identical 20-mer frequency patterns, with the
586   exception of a small insertion at nucleotide position ca. 1,400,000, which is also detected by the
587   MUMMER plot in **Figure 5A**. The genome of strain 381 is ca. 24 Kbps longer than that of
588   ATCC 33277 due to this insert and the length difference is illustrated in **Figure 6A** because the
589   length of the bars were based on  actual genome sizes. Another interesting example observation
590   is that even though strain JCVI SC001 is closest to SDJ2 due to their identical *16S rRNA*
591   sequences (**Figure 1**), their oligonucleotide frequency patterns are quite different, with each
592   showing unique regions (brighter colors) at different places. This can be due to two possibilities:
593   1) the artificial order of the unfinished sequence contigs (in the plots, contig order was the same
594   as that in the downloaded sequences); and 2) the bona fide differences in sequence.  This is also
595   true for the other three genomes A7436, W83 and W50, which share identical *16S rRNA*
596   sequences but exhibit distinct frequency patterns.

597       The global view of the 19 *P. gingivalis* genomes reveals distinct areas of unique and semi-
598   unique frequencies (with frequencies from 2-18, i.e., red to orange colors in the plot). These
599   differences in nucleotide sequences are essentially reflected in the genes and then translated to
600   proteins, and ultimately accounting for differences in biological functions. This nucleotide-to-
601   protein projection is shown in **Figure 6B**, which is the same as **Figure 6A**, except the non-
602   protein coding regions were masked with a different color (light blur) and the frequency colors
603   were shown for open reading frames (ORFs) in their corresponding DNA coding strands. As
604   expected most of the ORFs are dark in color and are shared by the majority of the *P. gingivalis*
605   genomes. Stretches of ORFs with colors close to yellow on either strand should account for the
606   differences of the number of unique proteins previously identified (**Figure 3**).

607       Finally the out-group *P. asaccharolytica* DSM 20707 shows mostly yellow colors in the plot,
608    which is as expected and means that *P. asaccharolytica* does not share much of the 20-mer
609    oligonucleotide sequences with *P. gingivalis*. Interestingly, by lowering the oligomer size to 14
610    bases (14-mer), most of the DSM 20707 genome appears black, meaning that these two different
611    species share most of the 14-mer sequences (data not shown). When the plot was generated with
612    15-mer sequences, the *P. asaccharolytica* DSM 20707 genome started to show patches of yellow
613    area (data not shown), indicating that some unique 15-mer sequences are present between these
614    two species. With 15-mer all *P. gingivalis* genomes are black in the plot (data not shown),
615    meaning 15-mer is too short and have not enough resolution power to differentiate unique
616    regions among *P. gingivalis* genomes. Hence whether the oligomer frequency analysis can detect
617    unique/shared regions in a group of genomes, depends on the size of the oligomer. The choice of
618    20-mer was able to identify unique regions among the strains of *P. gingivalis*, as shown in
619    **Figure 6A** and **B**, yet is too sensitive for a different species.

620    **Comparative Functional Genomics**

621       The comparative genomics is less meaningful without association with biological functions.
622    Most functional genomic annotations rely on either DNA or protein sequence homology to other
623    sequences with known biological functions. The most popular genome annotation pipeline is
624    probably the NCBI Prokaryotic Genome Annotation Pipeline (Tatusova et al., 2016), which is
625    the current default annotation pipeline when a microbial genome sequence is deposited to and
626    published in NCBI. The NCBI pipeline is specifically designed to annotate bacterial and archaeal
627    genomes (chromosomes and plasmids). It is a multi-level process that includes prediction of
628    protein-coding genes, as well as other functional genome units such as structural RNAs, tRNAs,
629    small RNAs, pseudogenes, control regions, direct and inverted repeats, insertion sequences,
630    transposons and other mobile elements. However the NCBI pipeline has undergone major
631    changes since its first implementation in 2005 (Angiuoli et al., 2008). The annotation quality and
632    results, and the scope of genetic or functional elements identified may not be the same,
633    depending on when the genome was deposited to and annotated by NCBI.

634       In addition to the NCBI pipeline, there are also several other microbial genome annotation
635    pipelines published and commonly used, such as the RAST system - Rapid Annotation of
636    microbial genomes using Subsystems Technology (Overbeek et al., 2014); the BASys – Bacterial
637    Annotation System, a web server for automated bacterial genome annotation (Van Domselaar et
638    al., 2005). Above the gene level, there have also been tools and databases available for
639    constructing and comparing metabolic pathways of microbial genomes. Examples in this
640    category are IMG - the integrated microbial genomes comparative analysis system (Markowitz et
641    al., 2014) and BlastKOALA – a KEGG tool for functional characterization of genome sequences
642    (Kanehisa et al., 2016). Both systems provide annotation information beyond individual gene
643    and protein level, such as, in the case of IMG, conserved protein domain and groups COGs and
644    families (Pfam), as well as the enzymes and metabolic pathways inferred by BlastKOALA.

645       In this report we compared the *P. gingivalis* genomes at the functional level based on three
646    systems: the NCBI annotation, the RAST annotation, and the BlastKOALA inferred metabolic

647  pathways. The results of these analyses are too voluminous to be presented in text however the
648  complete results are provided in a central online site for download
649  (ftp://www.homd.org/publication_data/20150425). Here we summarized all the comparisons into
650  a single table (**Table 8**). Initially, functional comparisons were done based on simple text search
651  – by counting the number of genes with functional annotations containing several categories of
652  keywords listed  in the table  (in Italic font).

653      Interestingly marked differences were observed in the NCBI and RAST annotations either by
654  total protein count or by keyword searches. For example, as shown in **Table 4**, the difference of
655  the total number of protein encoding genes annotated by the two systems can be as large as 351
656  (for strain F0566). The difference in functional annotation is also quite noticeable (**Table 8**). For
657  example, only five of the 19 *P. gingivalis* genomes have proteins annotated as "gingipain" by
658  NCBI, whereas three other different genomes were annotated by RAST to have a single
659  gingipain gene.

660      To remedy the differences and apparent incompleteness of the two annotation systems, a
661  more effective way to detect most, if not all, proteins of the same function, is to perform
662  sequence similarity searches using protein sequences that had been identified. For example, the
663  19 proteins that were annotated as "gingipain" (16 by NCBI, three by RAST, **Table 8**) were
664  grouped together and used as the baits to search against all proteins of all 20 genomes identified
665  by both systems. A total of 84 proteins highly similar to the 19 gingipain proteins were detected
666  this way among all 19 *P. gingivalis* genomes ranging from two to seven gingipains per genome
667  (none was detected in *P. asaccharolytica*). These searches were conservative by setting a high
668  percent sequence identity and coverage, and so the numbers can be under-estimated. This
669  approach was done repeatedly for each of the seven functional categories that were deemed of
670  high interest by authors. All the proteins identified by either NCBI or RAST in each category
671  were collectively searched against all protein sequences in all 20 and the number of proteins with
672  $\geq$ 95% sequence identity and $\geq$ 95% alignment coverage to the query sequences were recorded.
673  The results of the BLAST searches were listed in the third row of each category in **Table 8**.
674  Unsurprisingly, the number of genes identified in all categories is higher than those provided by
675  either annotation system, and often higher than both systems combined. The fact that stringent
676  BLAST search identified more proteins of the same function indicates that the currently
677  microbial genome annotation pipelines are quite incomprehensive and are in need of
678  improvement.

679      For gingipains, using the 16 NCBI identified proteins and three RAST ones (**Table 8**), the
680  BLAST search of these sequences matched with many more proteins that are highly similar to
681  gingipains in all 19 *P. gingivalis* genomes. Examining the annotation for those proteins highly
682  matched with annotated gingipains, most of them were simply annotated as "hypothetical" or
683  "functionally unknown" proteins, while some were annotated as "peptidase".

684      Another notable observation is the high prevalence of the transposase proteins encoded in
685  this species, as high as 149 copes in strain A7436. The lower number of transposases detected in
686  those unfinished genomes is most likely due to the in-between-contig sequence gaps that may

687  contain highly repeated sequences such as the transposases and the IS elements. The completed
688  genome with lowest number of mobility related genes is strain A7A1-28 where only 68 were
689  detected in the genome.

690      Capsular polysaccharide  (CPS) has long been recognized as an important virulence factor
691  for *P. gingivalis* (Singh et al., 2011) and encapsulated strains are known to be more virulent than
692  the non-encapsulated ones (Laine and van Winkelhoff, 1998). When all the annotated capsule
693  related proteins were BLASTP searched against all genomes, the total number of capsule related
694  proteins ranged consistently between five and six copies (**Table 8**). For example, W83 is known
695  as an encapsulated strain and ATCC 33277 is non-encapsulated. However both strains encode six
696  copies of capsulated related genes. Of these, four were annotated as "CPS/capsule biosynthesis
697  proteins" by both NCBI or RAST. Interestingly, NCBI only identified three of these four CPS
698  biosynthesis protein. The $5^{th}$ one was annotated as "CPS transport protein" in W83 by NCBI
699  (Genbank ID AAQ65636.1) but was annotated as "conserved hypothetic protein" in ATCC 3327
700  by NCBI (BAG34043.1) or "tyrosine-protein kinase Wzc" in both W83 and ATCC 33277 by
701  RAST. The $6^{th}$ capsule related gene was annotated as "sugar isomerase" in ATCC 33277
702  (BAG34552.1 ) or "SIS domain protein" in W83 (AAQ65335.1) by NCBI. This same gene was
703  annotated as "arabinose 5-phosphate isomerase" in both W83 and ATCC 3327 by RAST and
704  "sugar phosphate isomerase involved in capsule formation" in several other strains by NCBI.
705  Taken together, this serves as an example of how inconsistent both annotations are, for genes
706  involved in a single biological function. By BLAST searching using proteins annotated as
707  capsule related genes annotated across all 19 *P. gingivalis* genomes, we were able to detect
708  consistently between five to six copies of genes involved in encapsulation for this species. The
709  fact the all *P. gingivalis* genomes contain a similar number of capsule related genes yet some are
710  encapsulated and others are not, indicates that these genes may subject to different gene
711  expression controls. It is thus likely that some non-encapsulated strains may become
712  encapsulated under certain specific *in vivo* conditions.

713      In a very different functional aspect, there is a high prevalence of the bacterial phage related
714  proteins, such as phage integrase/site-specific recombinase, phage tail component proteins, and
715  phage-related lysozyme. The number of phage related proteins detected in the 19 *P. gingivalis*
716  genomes ranged from 12 to 25. Functional bacteriophage have so far never been detected in this
717  species (Sandmeier et al., 1993) yet contrarily many proteins related to phage reproduction were
718  detected in all the 19 *P. gingivalis* strains. One most plausible explanation is the prevalence of
719  the CRISPR/Cas systems in this species (discussed below); another is also the presence of the
720  abortive phage infection proteins found in several strains (ATCC 33277, HG66, W83, AJW4,
721  SJD2, and MP4-504, data not shown).

722      As mentioned above, another very interesting category of enzymes reported in **Table 8** is the
723  prevalence of proteins associated with the CRISPR (clustered regularly interspaced short
724  palindromic repeats) elements. CRISPR, together with the Cas (CRISPR associated) proteins,
725  have been dubbed as the adaptive immune system for Bacteria and Archaea to ward off invading
726  foreign DNA (Horvath and Barrangou, 2010). However, although CRISPR arrays were detected
727  in all genomes (including outgroup *P. asaccharolytica*, 4th row in the CRISPR category of

728    **Table 8**), that is not the case for the Cas proteins. Cas was not detected in the genome of strain
729    JCVI SC001, and only one copy detected in strain AJW4 (3rd row in **Table 8** CRISPR category).
730    Strain F0569 has the highest number or CRISPR arrays detected using the online software
731    CRISPRfinger (http://crispr.i2bc.paris-saclay.fr/Server) but this strain does not have the highest
732    number of Cas proteins. Of all the CRISPR arrays detected, the length of the direct repeat (DR)
733    element ranged from 23 to 47 bps and the number of the DRs in the array ranged from five to
734    121 copies (data not shown but available from the online FTP site). Both ATCC 33277and strain
735    381 had three copies of nearly identical CRISPR arrays and both had one copy of the arrays with
736    121 DR sequences (and 120 spacer sequences). The high DR copy number may be an indication
737    for the CRISPR activity in the past. On the other hand, JSVI SC001 had three copies of CRISPR
738    arrays detected with DR of 31, 26, and 45 bps and repeat number 5, 7 and 6 respectively.
739    Whether or not this strain possesses a type of Cas protein that is very different from those in
740    other strains remains to be investigated. If this strain lacks any functional Cas protein, it is likely
741    to be susceptible to bacteriophage infection or the activation of the possible presence of
742    prophages as evidenced by the detection of 25 copies of phage related proteins (**Table 8**).

743    On the other side of the scale, at the metabolic pathway level, the KEGG pathways and
744    KEGG Orthology identified by BlastKOALA are BLAST-based, i.e., all the proteins sequences
745    regardless of their annotations, were BLAST-searched against the online protein database used
746    by BLastKOALA (http://www.kegg.jp/blastkoala/). Hence the comprehensiveness of the KEGG
747    pathways and the KO terms inferred by BlastKOALA depend on the completeness of the
748    proteins in the database. At any rate, several metabolic pathways identified to be unique to the
749    species *P. gingivalis* are: glycosphingolipid biosynthesis – globoseries; sphingolipid metabolism;
750    lysosome; glycosphingolipid biosynthesis - ganglio series; and glycosaminoglycan degradation.
751    These species specific pathways were suggested based on the fact that they were detected in all
752    of the 19 *P. gingivalis* genomes but not in *P. asaccharolytica* DSM 20707 (detailed data
753    available in the FTP site). When compared to *P. asaccharolytica* DSM 20707, BlastKOALA
754    determined that *P. gingivalis* lacks proteins involved in the following pathways: C5-branched
755    dibasic acid metabolism; AMPK signaling pathway; amoebiasis; thyroid hormone synthesis;
756    apoptosis; and arachidonic acid metabolism.

757    **Concluding Remarks**

758    In this report 19 genomes of the species *P. gingivalis* as well as the outgroup species *P.*
759    *asaccharolytica* were compared at several different levels of information ranging from
760    nucleotide to genes to proteins and metabolic functions. Based on the single gene 16S rRNA
761    phylogeny and multi-gene pholygenomic approach using core/shared protein sequences, several
762    plausible evolutionary paths were suggested. Although there is no single evolutionary path
763    concluded by these analyses, two closely related groups were consistently observed throughout
764    the analyses. The first group consists of strains ATCC 33277, 381 and HG66 and the second of
765    W83, W50 and A7436. The group of ATCC 33277, 381 and HG66 is also closer to the possible
766    common ancestor inferred based on the use of an outgroup species *P. asaccharolytica*. We also
767    detected at least 1,037 core/shared proteins for this species based on 95% sequence similarity
768    and 90% alignment length. However the number of core proteins increases with the lowering of

the two detecting parameters. Functional and metabolic pathways were also compared and suggested several important functions of pathways that are unique to this species, to each strain, or missing in any particular strain. *P. gingivalis* has many genes encoding proteins related to or involved in gingipains, attachment (e.g., adhesins and fimbrins), capsules, and phages. These proteins were either missing or present in very few copies in the neighbor species *P. asaccharolytica.* Particularly intriguing observations were prevalence of many proteins related in phage productions and the equal prevalence of the CRISPR system in this species, with the exception of one strain lacking the Cas proteins.

Despite the large amount of comparative results generated in this study, there are still many different ways and software tools for analyzing and comparing a group of genomes. The complete results presented in this report, together with several other results that were only mentioned briefly here, are made available for download online at ftp://www.homd.org/publication_data/20150425. We hope these data are useful to the research community and more hypotheses can be formulated based on the current or future analyses in order to gain deeper understanding on this important periodontal pathogen.

**AUTHOR CONTRIBUTIONS**

TS: data acquisition, data analysis, data interpretation, writing of the manuscript, final approval of the version to be published; HS: data acquisition, data analysis, data interpretation, writing; IO: initiating the study, writing of the manuscript, revising the manuscript, final approval of the version to be published.

**SUPPLEMENTARY MATERIAL**

801     All the supplementary material for this article can be found online the FTP site at
802     ftp://www.homd.org/publication_data/20160425/)

803

804     **REFERENCES**

805     Altschul, S.F., Madden, T.L., Schäffer, A.A., Zhang, J., Zhang, Z., Miller, W., et al. (1997).
806         Gapped BLAST and PSI-BLAST: a new generation of protein database search programs.
807         *Nucleic Acids Res.* 25, 3389-3402.

808     Angiuoli, S.V., Gussman, A., Klimke, W., Cochrane, G., Field, D., Garrity, G., et al., (2008).
809         Toward an online repository of Standard Operating Procedures (SOPs) for (meta)genomic
810         annotation. *OMICS* 12, 137-141. doi: 10.1089/omi.2008.0017

811     Aziz, R. K., Devoid, S., Disz, T., Edwards, R.A., Henry, C.S., Olsen, G.J., et al., (2012). SEED
812         servers: high-performance access to the SEED genomes, annotations, and metabolic models.
813         *PLoS One* 7, e48053.

814     Aziz, R.K., Bartels, D., Best, A.A., DeJongh, M., Disz, T., Edwards, R.A., et al. (2008).The
815         RAST Server: rapid annotations using subsystems technology. *BMC Genomics* 9, 75. doi:
816         10.1186/1471-2164-9-75

817     Brunner, J., Wittink, F.R., Jonker, M.J., de Jong, M., Breit, T.M., Laine, M. L., et al., (2010).
818         The core genome of the anaerobic oral pathogenic bacterium Porphyromonas gingivalis.
819         *BMC Microbiol.* 10, 252.

820     Chen, T., Hosogi, Y., Nishikawa, K., Abbey, K., Fleischmann, R.D., Walling, J., et al. (2004).
821         Comparative whole-genome analysis of virulent and avirulent strains of *Porphyromonas*
822         *gingivalis. J. Bacteriol.* 186, 5473-5479.

823     Darveau, R.P., Hajishengallis, G., and Curtis, M.A. (2012). *Porphyromonas gingivalis* as a
824         potential community activist for disease. *J. Dent. Res.* 91, 816–820. doi:
825         http://dx.doi.org/10.1177/0022034512453589

826     Delcher, A.L., Phillippy, A., Carlton, J., and Salzberg, S.L. (2002). Fast algorithms for large-
827         scale genome alignment and comparison. *Nucleic Acids Res.* 30, 2478-2483.

828     Demmer, R.T., and Desvarieux, M. (2006). Periodontal infections and cardiovascular disease:
829         the heart of the matter. *J. Am. Dent. Assoc.* 137, 14S-20S.

830     Dolgilevich, S., Rafferty, B., Luchinskaya, D., and Kozarov, E. (2011). Genome comparison of
831         invasive and rare non-invasive strains reveals *Porphyromonas gingivalis* genetic
832         polymorphisms. *J. Oral Microbiol.* 3, 10.3402/jom.v3i0,5764

833     Dorn, B.R., Burks, J.N., Seifert, K.N., and Progulske-Fox, A. (2000). Invasion of endothelial and
834         epithelial cells by strains of *Porphyromonas gingivalis*. *FEMS Microbiol. Lett.* 187, 139-144.

835    Duncan,  M.J. (2003). Genomics of oral bacteria. *Crit. Rev. Oral Biol. Med.* 14,175-187.

836    Hajishengallis, G., Darveau, R.P., and Curtis, M.A. (2012).The keystone pathogen hypothesis.
837        *Nat. Rev. Microbiol.* 10, 717–725. doi: http://dx.doi.org/10.1038/nrmicro2873

838    Horvath, P., and Barrangou, R. (2010). CRISPR/Cas, the immune system of bacteria and
839        archaea. *Science* 327, 167-170.

840    Howe, K., Bateman, A., and Durbin, R. (2002). QuickTree: building huge Neighbour- Joining
841        trees of protein sequences. *Bioinformatics* 18, 1564-1567.

842    Jones, D.T., Taylor, W.R., and Thornton, J.M. (1992).The rapid generation of mutation data
843        matrices from protein sequences. *Comput. Appl. Biosci.* 8, 275-282.

844    Kanehisa, M., Sato, Y., and Morishima, K. (2016). BlastKOALA and GhostKOALA: KEGG
845        tools for functional characterization of genome and metagenome sequences. *J.        Mol.*
846        *Biol.* 428, 726-731.

847    Katoh, K., and Standley, D.M. (2013). MAFFT multiple sequence alignment software version
848        7: improvements in performance and usability. *Mol. Biol. Evol.* 30, 772-780. doi:
849        10.1093/molbev/mst010

850    Klappenbach, J.A., Dunbar, J.M., and Schmidt, T.M. (2000).  rRNA operon copy number
851        reflects ecological strategies of bacteria. *Appl. Environ. Microbiol.* 66, 1328-1333.

852    Klein, B.A., Chen, T., Scott. J.C., Koenigsberg, A.L., Duncan, M.J., and Hu, L.T. (2015).
853        Identification and characterization of a minisatellite contained within a novel miniature
854        inverted-repeat transposable element (MITE) of *Porphyromonas gingivalis*. *Mob. DNA* 6, 18.
855        doi: 10.1186/s13100-015-0049-1

856    Klein, B. A., Tenorio, E.L., Lazinski, D.W., Camilli, A., Duncan, M.J., and Hu, L.T. (2012).
857        Identification of essential genes of the periodontal pathogen Porphyromonas gingivalis. *BMC*
858        *Genomics* 13, 578.

859    Laine, M. L., and van Winkelhoff, A.J. (1998). Virulence of six capsular serotypes of
860        Porphyromonas gingivalis in a mouse model. *Oral Microbiol. Immunol.* 13, 322-325.

861    Liu, D., Zhou, Y., Naito, M., Yumoto, H., Li, Q., Miyake, Y., et al. (2014). Draft genome
862        sequence of *Porphyromonas gingivalis* strain SJD2, isolated from the periodontal pocket of a
863        patient with periodontitis in China. *Genome Announc.* 2. pii: e01091-13. doi:
864        10.1128/genomeA.01091-13

865    Lundberg, K., Wegner, N., Yucel-Lindberg, T., and Venables, P.J. (2010). Periodontitis in RA-
866        citrullinated enolase connection. *Nat. Rev. Rheumatol.* 6, 727-730.

867    Markowitz, V.M., Chen, I.M., Palaniappan, K., Chu, K., Szeto, E., Pillay, M., et al. (2014). IMG
868        4 version of the integrated microbial genomes comparative analysis system. *Nucleic Acids*
869        *Res.* 42 (Database issue), D560-567. doi: 10.1093/nar/gkt963

870  Marsh, P.D., McDermid, A.S., McKee, A.S., and Baskerville, A. (1994). The effect of growth
871      rate and haemin on the virulence and proteolytic activity of *Porphyromonas gingivalis* W50.
872      *Microbiology* 140, 861-865.

873  McLean, J.S., Lombardo, M.J., Ziegler, M.G., Novotny, M., Yee-Greenbaum, J., Badger, J.H., et
874      al. (2013). Genome of the pathogen *Porphyromonas gingivalis* recovered from a biofilm
875      in a hospital sink using a high-throughput single-cell genomics platform. *Genome Res.*
876      23, 867-877. doi: 10.1101/gr.150433.112

877  Nagano, K., Hasegawa, Y., Yoshida, Y., and Yoshimura, F. (2015). A major fimbrilin variant of
878      Mfa1fimbriae in *Porphyromonas gingivalis*. *J. Dent. Res.* 94, 1143-1148. doi:
879      10.1177/0022034515588275

880  Naito, M., Hirakawa, H., Yamashita, A., Ohara, N., Shoji, M., Yukitake, H., et al. (2008).
881      Determination of the genome sequence of *Porphyromonas gingivalis* strain ATCC 33277 and
882      genomic comparison with strain W83 revealed extensive genome  rearrangements in *P.*
883      *gingivalis*. *DNA Res.* 15, 215–225.  http://dx.doi.org/10.1093/dnares/dsn013

884  Nelson, K.E., Fleischmann, R.D., DeBoy, R.T., Paulsen, I.T., Fouts, D.E., Eisen, J.A., et al.
885      (2003). Complete genome sequence of the oral pathogenic Bacterium  *porphyromonas*
886      *gingivalis* strain W83. *J. Bacteriol.* 185, 5591-5601.

887  Olsen, I., and Progulske-Fox, A. (2015). Invasion of *Porphyromonas gingivalis* into vascular
888      cells and tissue. *J. Oral Microbiol.* 7, 28788. doi: 10.3402/jom.v7.28788

889  Olsen, I., Singhrao, S.K. (2015). Can oral infection be a risk factor for Alzheimer's disease? *J.*
890      *Oral Microbiol.*7, 29143. doi: 10.3402/jom.v7.29143

891  Overbeek, R., Olson, R., Pusch, G.D., Olsen, G.J., Davis, J.J., Disz, T., et al. (2014). The SEED
892      and the Rapid Annotation of microbial genomes using Subsystems Technology (RAST).
893      *Nucleic Acids Res.* 42 (Database issue), D206-214. doi:  10.1093/nar/gkt122

894  Price, M.N., Dehal, P.S., and Arkin, A.P. (2010). FastTree 2-approximately maximum-
895      likelyhood trees for large alignments. *PLoS One* 5, e9490.doi:
896      10.1371/journal.pone.0009490

897  Sandmeier, H., Bar, K., and Meyer, J. (1993). Search for bacteriophages of black-pigmented
898      gram-negative anaerobes from dental plaque. *FEMS Immunol. Med. Microbiol.* 6,193- 194.

899  Segata, N., Bornigen, D., Morgan, X.C., and Huttenhower, C. (2013). PhyloPhlAn is a new
900      method for improved phylogenetic and taxonomic placement of microbes. *Nat.  Commun.* 4,
901      2304.

902  Shimodaira, H., and Hasegawa, M. (2001). CONSEL: for assessing the confidence of
903      phylogenetic tree selection. *Bioinformatics* 17, 1246-1247.

904  Siddiqui, H., Yoder-Himes, D.R., Mizgalska, D., Nguyen, K.A., Potempa, J., and Olsen, I.
905      (2014). Genome sequence of *Porphyromonas gingivalis* strain HG66 (DSM 28984).
906      *Genome Announc.* 2. pii: e00947-14. doi: 10.1128/genomeA.00947-14

907  Singh, A., Wyant, T., Anaya-Bergman, C., Aduse-Opoku, J., Brunner, J., Laine, M.L., et al.
908      (2011). The capsule of Porphyromonas gingivalis leads to a reduction in the host
909      inflammatory response, evasion of phagocytosis, and increase in virulence. *Infect.  Immun.*
910      79, 4533-4542.

911  Socransky, S.S., Haffajee, A.D., Cugini, M.A., Smith, C., and Kent, R.L., Jr. (1998). Microbial
912      complexes in subgingival plaque. *J. Clin. Periodontol.* 2, 134-144.

913  Talavera, G., and Castresana, J. (2007). Improvement of phylogenies after removing divergent
914      and ambiguously aligned blocks from protein sequence alignments. *Syst. Biol.* 56, 564-577.

915  Tatusova, T., DiCuccio, M., Badretdin, A., Chetvernin, V., Nawrocki, E.P., Zaslavsky, L., et al.
916      (2016). NCBI prokaryotic genome annotation pipeline. *Nucleic Acids Res.* 19, 44,  6614-
917      6624. doi: 10.1093/nar/gkw569

918  Tatusova, T., DiCuccio, M., Badretdin, A., Chetvernin, V., Nawrocki, E.P., Zaslavsky, L., et al.
919      (2016). NCBI prokaryotic genome annotation pipeline. *Nucleic Acids Res.* 44, 6614- 6624.
920      doi: 10.1093/nar/gkw569

921  To, T.T., Liu, Q., Watling, M., Bumgarner, R.E., Darveau, R.P., and McLean, J.S. (2016). Draft
922      genome sequence of low-passage clinical isolate *Porphyromonas gingivalis* MP4-504.
923      *Genome Announc.* 4. pii: e00256-16. doi: 10.1128/genomeA.00256-16

924  Tribble, G.D., Kerr, J.E., and Wang, B.Y. (2013) Genetic diversity in the oral pathogen
925      *Porphyromonas gingivalis*: molecular mechanisms and biological consequences. *Future*
926      *Microbiol.*8, 607-620. doi: 10.2217/fmb.13.30

927  Van Domselaar, G.H., Stothard, P., Shrivastava, S., Cruz, J.A., Guo, A., Dong, X., et al. (2005).
928      BASys: a web server for automated bacterial genome annotation. *Nucleic Acids Res.*  33
929      (Web Server issue), W455-459.

930  Vetrovsky, T., and Baldrian, P. (2013). The variability of the 16S rRNA gene in bacterial
931      genomes and its consequences for bacterial community analyses. *PLoS One* 8, e57923. doi:
932      10.1371/journal.pone.0057923

933  Watanabe, T., Maruyama, F., Nozawa, T., Aoki, A., Okano, S., Shibata, Y., et al. (2011).
934      Complete genome sequence of the bacterium *Porphyromonas gingivalis* TDC60, which
935      causes periodontal disease. *J. Bacteriol.* 193, 4259-4260. doi: 10.1128/JB.05269-11

936  Woese, C.R., Kandler, O., and Wheelis, M.L. (1990). Towards a natural system of organisms:
937      proposal for the domains Archaea, Bacteria, and Eucarya. *Proc. Natl. Acad. Sci. USA* 87,
938      4576-4579.

939

**TABLE 1. Summary of all the *P. gingivalis* genome sequences compared in this report [1].**

| Strain | Sequence Release Date [2] | Genome Size (bps) | Contigs | GenBank Accession | Bioproject | Biosample [3] | Submitter |
|---|---|---|---|---|---|---|---|
| W83 | 2003-09-02 | 2,343,476 | 1 | AE015924 | PRJNA48 | SAMN02603720 | *Porphyromonas gingivalis* Genome Project |
| ATCC_33277 | 2008-05-20 | 2,354,886 | 1 | AP009380 | PRJDA19051 | | Kitasato Univ. |
| TDC60 | 2011-05-23 | 2,339,898 | 1 | AP012203 | PRJDA66755 | | Tokyo Medical and Dental Univ. |
| W50 | 2012-06-25 | 2,242,062 | 104 | AJZS01000000 | PRJNA78905 | SAMN00792205 | J. Craig Venter Institute |
| JCVI_SC001 | 2013-04-24 | 2,426,396 | 1 / 284 | CM001843 [4] / APMB01000000 | PRJNA167667 | SAMN02436407 | J. Craig Venter Institute |
| F0568 | 2013-09-16 | 2,334,744 | 154 | AWUU01000000 | PRJNA173937 | SAMN02436723 | Washington Univ. |
| F0569 | 2013-09-16 | 2,249,227 | 111 | AWUV01000000 | PRJNA173938 | SAMN02436724 | Washington Univ. |
| F0570 | 2013-09-16 | 2,282,791 | 117 | AWUW01000000 | PRJNA173939 | SAMN02436747 | Washington Univ. |
| F0185 | 2013-09-16 | 2,246,368 | 113 | AWVC01000000 | PRJNA198891 | SAMN02436815 | Washington Univ. |
| F0566 | 2013-09-16 | 2,306,092 | 192 | AWVD01000000 | PRJNA198892 | SAMN02436881 | Washington Univ. |
| W4087 | 2013-09-16 | 2,216,597 | 114 | AWVE01000000 | PRJNA198893 | SAMN02436749 | Washington Univ. |
| SJD2 | 2013-12-04 | 2,329,548 | 117 | ASYL01000000 | PRJNA205615 | SAMN02470968 | Shanghai Jiao Tong Univ. School of Medicine |
| HG66 | 2014-08-14 | 2,441,780 | 1 | CP007756 | PRJNA245225 | SAMN02732406 | Univ. of Louisville |
| A7436 | 2015-08-11 | 2,367,029 | 1 | CP011995 | PRJNA276132 | SAMN03366764 | Univ. of Florida |
| AJW4 | 2015-08-26 | 2,372,492 | 1 | CP011996 | PRJNA276132 | SAMN03372093 | Univ. of Florida |
| Ando | 2015-09-17 | 2,229,994 | 112 | BCBV01000000 | PRJDB4201 | SAMD00040429 | Lab. of Plant Genomics and Genetics, Dept. of Plant Genome Research, Kazusa DNA Research Institute |
| 381 | 2015-10-14 | 2,378,872 | 1 | CP012889 | PRJNA276132 | SAMN03656156 | Univ. of Florida |
| A7A1-28 | 2015-11-17 | 2,249,024 | 1 | CP013131 | PRJNA276132 | SAMN03653671 | Univ. of Florida |
| MP4-504 | 2016-02-09 | 2,373,453 | 92 | LOEL01000000 | PRJNA305025 | SAMN04309157 | Univ. of Washington |

[1] For a more detailed list of this table please follow this web link: ftp://www.homd.org/publication_data/20160425/

[2] Genomes of this table are sorted by the original sequence release date.

[3] Unassembled raw sequence reads from which the assembly that was done can be traced back by the Biosample ID, if available.

[4] This Genbank number shows the sequence as "circular", however it is a single pseudo-contig with many Ns filling the gaps. Thus it should not be considered as a complete genome.

**TABLE 2. Effective (non-Ns) sizes of the genomes.**

| Strain | Contigs | Size(bps) | Non-N Size(bps) [1] | Ns (bps) | N Fragment Size Range (Fragment Count) |
|---|---|---|---|---|---|
| HG66 | 1 | 2,441,780 | 2,441,680 | 100 | 100 (1) |
| JCVI_SC001 | 1 | 2,426,396 | 2,398,196 | 28,200 | 100 (282) |
| 381 | 1 | 2,378,872 | 2,378,872 | 0 | None |
| MP4-504 | 92 | 2,373,453 | 2,373,453 | 0 | None |
| AJW4 | 1 | 2,372,492 | 2,372,492 | 0 | None |
| A7436 | 1 | 2,367,029 | 2,367,029 | 0 | None |
| ATCC_33277 | 1 | 2,354,886 | 2,354,886 | 0 | None |
| W83 | 1 | 2,343,476 | 2,343,476 | 0 | None |
| TDC60 | 1 | 2,339,898 | 2,339,897 | 1 | 1 (1) |
| SJD2 | 117 | 2,329,548 | 2,328,850 | 698 | 4-256 (23) |
| F0568 | 154 | 2,334,744 | 2,328,244 | 6,500 | 100 (65) |
| F0566 | 192 | 2,306,092 | 2,300,992 | 5,100 | 100 (51) |
| F0570 | 117 | 2,282,791 | 2,278,391 | 4,400 | 100 (44) |
| A7A1-28 | 1 | 2,249,024 | 2,249,024 | 0 | None |
| W50 | 104 | 2,242,062 | 2,242,060 | 2 | 1 (2) |
| F0569 | 111 | 2,249,227 | 2,242,027 | 7,200 | 100 (72) |
| F0185 | 113 | 2,246,368 | 2,240,268 | 6,100 | 100 (61) |
| Ando | 112 | 2,229,994 | 2,227,972 | 2,022 | 10-100 (61) |
| W4087 | 114 | 2,216,597 | 2,212,597 | 4,000 | 100 (40) |

[1] Genomes are ordered based on the non-N size.

**TABLE 3. Summary of the NCBI Annotation[1].**

| Strain | Protein coding | tRNA | rRNA | tmRNA[3] | Repeat region | Binding site | Pseudo-gene | ncRNA[2] | | | | Other |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | | | | | | | | Antisense-RNA | RNase-P-RNA | Auto-catalytically spliced intron | Other ncRNA | |
| W83 | 1,909 | 53 | 12 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 41 |
| ATCC_33277 | 2,090 | 53 | 12 | 0 | 210 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| TDC60 | 2,220 | 53 | 12 | 1 | 380 | 7 | 0 | 0 | 0 | 0 | 1 | 34 |
| W50 | 2,016 | 48 | 3 | 1 | 0 | 8 | 0 | 1 | 1 | 0 | 0 | 0 |
| JCVL_SC001 | 2,354 | 45 | 3 | 1 | 0 | 8 | 0 | 1 | 1 | 0 | 0 | 0 |
| F0568 | 2,410 | 46 | 3 | 1 | 0 | 7 | 0 | 0 | 1 | 0 | 0 | 0 |
| F0569 | 2,297 | 46 | 3 | 1 | 0 | 7 | 0 | 0 | 1 | 1 | 0 | 0 |
| F0570 | 2,315 | 44 | 3 | 1 | 0 | 7 | 0 | 0 | 1 | 1 | 0 | 0 |
| F0185 | 2,233 | 45 | 3 | 1 | 0 | 7 | 0 | 0 | 1 | 0 | 0 | 0 |
| F0566 | 2,392 | 45 | 3 | 1 | 0 | 7 | 0 | 0 | 1 | 1 | 0 | 0 |
| W4087 | 2,202 | 45 | 3 | 1 | 0 | 7 | 0 | 0 | 1 | 1 | 0 | 0 |
| SJD2 | 2,012 | 48 | 3 | 0 | 3 | 0 | 62 | 0 | 0 | 0 | 0 | 0 |
| HG66 | 1,958 | 53 | 12 | 0 | 3 | 5 | 38 | 0 | 1 | 0 | 0 | 0 |
| A7436 | 2,004 | 53 | 12 | 1 | 4 | 0 | 3 | 0 | 1 | 0 | 0 | 0 |
| AJW4 | 2,002 | 53 | 12 | 1 | 2 | 0 | 2 | 0 | 1 | 0 | 0 | 0 |
| Ando | 1,770 | 47 | 4 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 381 | 1,968 | 53 | 12 | 1 | 3 | 0 | 9 | 0 | 1 | 1 | 0 | 0 |
| A7A1-28 | 1,841 | 53 | 12 | 1 | 5 | 0 | 37 | 0 | 1 | 0 | 0 | 0 |
| MP4-504 | 1,889 | 47 | 3 | 0 | 3 | 2 | 99 | 0 | 1 | 0 | 0 | 0 |

[1] Data analyzed based on the gff files of each genome generated by the NCBI annotation pipeline.

[2] non-coding RNA

3 Trans-messenger RNA: a bacterial RNA molecule with dual tRNA-like and mRNA-like properties

**TABLE 4. Comparison of NCBI and RAST genome annotations[1].**

| Strain | Total NCBI | Total RAST | Common / Unique[2] | 5S rRNA | 16S rRNA | 23S rRNA | tRNA |
|---|---|---|---|---|---|---|---|
| W83 | 1,909 | 2,163 | 1,784 / 80 / 334 | 4 | 4 | 4 | 53 |
| ATCC_33277 | 2,090 | 2,092 | 1,911 / 154 / 144 | 4 | 4 | 4 | 53 |
| TDC60 | 2,220 | 2,090 | 1,880 / 286 / 167 | 4 | 4 | 4 | 53 |
| W50 | 2,016 | 2,036 | 1,887 / 102 / 123 | 1 | 1 | 1 | 48 |
| JCVI_ SC001 | 2,354 | 2,136 | 2,030 / 276 / 78 | 1 | 1 | 1 | 45/42 |
| F0568 | 2,417 | 2,096 | 1,939 / 403 / 111 | 1 | 1 | 1 | 46 |
| F0569 | 2,297 | 1,982 | 1,845 / 377 / 92 | 1 | 1 | 1 | 46 |
| F0570 | 2,316 | 2,063 | 1,912 / 338 / 107 | 1 | 1 | 1 | 44 |
| F0185 | 2,236 | 2,005 | 1,862 / 319 / 107 | 1 | 1 | 1 | 45 |
| F0566 | 2,395 | 2,044 | 1,885 / 428 / 112 | 1 | 1 | 1 | 45 |
| W4087 | 2,204 | 1,973 | 1,850 / 303 / 92 | 1 | 1 | 1 | 45 |
| SJD2 | 2,020 | 2,166 | 1,845 / 136 / 271 | 1 | 1 | 1 | 48/47 |
| HG66 | 1,958 | 2,215 | 1,881 / 58 / 298 | 4 | 4 | 4 | 53 |
| A7436 | 2,004 | 2,173 | 1,898 / 84 / 239 | 4 | 4 | 4 | 53 |
| AJW4 | 2,002 | 2,139 | 1,884 / 104 / 226 | 4 | 4 | 4 | 53 |
| Ando | 1,788 | 1,989 | 1,674 / 76 / 275 | 2 | 1 | 1 | 47 |
| 381 | 1,968 | 2,108 | 1,853 / 91 / 221 | 4 | 4 | 4 | 53 |
| A7A1-28 | 1,841 | 2,039 | 1,736 / 89 / 269 | 4 | 4 | 4 | 53 |
| MP4-504 | 1,891 | 2,181 | 1,806 / 68 / 347 | 1 | 1 | 1 | 47 |

[1] Only protein-coding, rRNA and tRNA genes were compared since these are the only types of genes annotated by RAST.

[2] The three numbers shown (X / Y / Z) are X: common genes, genes with $\geq$ 80 % overlapped based on the annotated start and end postion; Y: RAST unique genes, gene annotated by RAST without overlap of any NCBI gene; Z: NCBI unique genes, genes annotated by NCBI without overlap to any RAST gene. There are genes that are partially overlapping to each other with < 80% of the length not included.

**TABLE 5. Unique *16S rRNA* gene sequences in *P. gingivalis* genomes.**

| Original Sequence | Trimmed Sequence[1] | Copy Number | Original Length (bps) | Strains (Copy Number) [2] |
|---|---|---|---|---|
| Unique Seq 1 | Unique Trimmed Seq 1 | 4 | 1422 | 381 (4) |
| Unique Seq 2 | | 4 | 1475 | ATCC33277 (4) |
| Unique Seq 3 | | 3 | 1538 | HG66 (3) |
| Unique Seq 4 | Unique Trimmed Seq 2 | 1 | 1538 | HG66 |
| Unique Seq 5 | Unique Trimmed Seq 3 | 3 | 1422 | A7436 (3) |
| Unique Seq 6 | | 5 | 1475 | W50; W83 (4) |
| Unique Seq 7 | Unique Trimmed Seq 4 | 1 | 1422 | A7436 |
| Unique Seq 8 | Unique Trimmed Seq 5 | 4 | 1422 | A7A1-28 (4) |
| Unique Seq 9 | Unique Trimmed Seq 6 | 3 | 1422 | AJW4 (3) |
| Unique Seq 10 | Unique Trimmed Seq 7 | 1 | 1422 | AJW4 |
| Unique Seq 11 | Unique Trimmed Seq 8 | 1 | 1521 | TDC60 |
| Unique Seq 12 | | 1 | 1520 | TDC60 |
| Unique Seq 13 | Unique Trimmed Seq 9 | 1 | 1522 | TDC60 |
| Unique Seq 14 | Unique Trimmed Seq 10 | 1 | 1520 | TDC60 |
| Unique Seq 15 | Unique Trimmed Seq 11 | 1 | 1475 | JCVI SC001 |
| Unique Seq 16 | Unique Trimmed Seq 12 | 1 | 1538 | SJD2 |
| Unique Seq 17 | Unique Trimmed Seq 13 | 1 | 1475 | Ando |
| Unique Seq 18 | Unique Trimmed Seq 14 | 1 | 1520 | W4087 |
| Unique Seq 19 | Unique Trimmed Seq 15 | 1 | 1520 | F0569 |
| Unique Seq 20 | Unique Trimmed Seq 16 | 1 | 1520 | F0568 |
| Unique Seq 21 | Unique Trimmed Seq 17 | 1 | 1520 | F0185 |
| Unique Seq 22 | Unique Trimmed Seq 18 | 1 | 1520 | F0566 |
| Unique Seq 23 | Unique Trimmed Seq 19 | 1 | 1542 | MP4-504 |
| Unique Seq 24 | Unique Trimmed Seq 20 | 1 | 1520 | F0570 |
| Unique Seq 25[3] | Unique Trimmed Seq 21 | 2 | 1517 | PaDSM20707 (2) |

[1] Sequences were pre-aligned with the software MAFFT v6.935b (2012/08/21) (Katoh and Standley, 2013) with default setting; after trimming the leading and trailing sequences not present for all genomes, the trimmed aligned sequence length is 1,425 bps in length.

[2] If multiple copies of identical sequences are present, the copy number is indicated in the parenthesis.

[2] Sequence of *P. asaccharolytica* strain DSM 20707 (from Genbank ID: CP002689) was included as outgroup

**TABLE 6. Percent hypothetical proteins 19 *P. gingivalis* genomes.**

| Strain [1] | Total | % Total hypothetical | Total unique (80% identity) | % Unique hypothetical (80% identity) |
|---|---|---|---|---|
| HG66 | 1,958 | 28% | 53 | 81% |
| 381 | 1,968 | 27% | 13 | 85% |
| ATCC_33277 | 2,090 | 42% | 14 | 79% |
| A7A1-28 | 1,841 | 28% | 46 | 78% |
| MP4-504 | 1,891 | 27% | 34 | 85% |
| Ando | 1,788 | 29% | 61 | 70% |
| F0568 | 2,417 | 46% | 114 | 88% |
| F0569 | 2,297 | 45% | 125 | 86% |
| W4087 | 2,204 | 43% | 94 | 78% |
| F0185 | 2,236 | 43% | 72 | 88% |
| F0570 | 2,316 | 44% | 96 | 90% |
| JCVI_ SC001 | 2,354 | 30% | 172 | 72% |
| SJD2 | 2,020 | 35% | 79 | 82% |
| AJW4 | 2,002 | 29% | 45 | 76% |
| A7436 | 2,004 | 28% | 25 | 80% |
| W50 | 2,016 | 26% | 34 | 91% |
| W83 | 1,909 | 35% | 13 | 100% |
| F0566 | 2,395 | 46% | 161 | 86% |
| TDC60 | 2,220 | 41% | 78 | 68% |

[1] The strains were ordered somewhat according to the 16S rRNA phylogenetic tree shown in Fig. 1.

[2] The unique proteins were identified by "blastclust" program with parameters 80% as the sequence identity and 50% alignment length.

**TABLE 7. Non-hypothetical unique[1] proteins in 19 *P. gingivalis* genomes.**

| Strain | Annotation |
| --- | --- |
| HG66 | glyoxalase |
| A7A1-28 | beta-galactosidase; putative hydrolase or acyltransferase of alpha/beta superfamily |
| Ando | DNA polymerase III subunits gamma and tau, partial<br>external scaffolding protein D<br>replication-associated protein A<br>major spike protein G |
| F0568 | DGQHR domain protein |
| F0569 | toxin-antitoxin system,<br>toxin component, Fic domain protein |
| W4087 | CAAX amino terminal protease family protein<br>phage portal protein, SPP1 family<br>phage uncharacterized protein |
| F0185 | peptidase S24-like protein |
| JCVI_SC001 | thioesterase family protein, partial<br>starch-binding protein, SusD-like domain protein, partial<br>spermine/spermidine synthase, partial<br>phage portal protein, lambda family, partial<br>head to-tail joining protein W<br>serine carboxypeptidase domain protein, partial<br>NYN domain protein<br>imidazoleglycerol-phosphate dehydratase domain protein, partial<br>carbohydrate kinase, PfkB domain protein<br>PF13785 domain protein, partial<br>DNA-binding helix-turn-helix protein |
| SJD2 | transposase ISPsy14 |
| AJW4 | geranylgeranyl pyrophosphate synthase<br>T5orf172 domain-containing protein |
| A7436 | transposase |
| W50 | transposase, mutator-like family protein |
| TDC60 | terminase |

[1] These proteins were searched against all the proteins in the 19 genomes and matched none but itself at the default BLASTP 2.2.25 parameter (i.e., with expected $e$ value ≤10) ( Altschul et al., 1997)

**TABLE 8. Comparative functional genomics of *P. gingivalis* genomes[1,2,3].**

| Annotation Source | ATCC33277 | HG66 | 381 | W83 | W50 | A7436 | AJW4 | F0570 | JCVISC001 | SJD2 | F0568 | F0569 | Ando | F0185 | W4087 | MP4-504 | A7A1-28 | F0566 | TDC60 | PaDSM20707 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **Gingipain**: *rgp, kgp, gingipain* | | | | | | | | | | | | | | | | | | | | |
| NCBI | 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 0 | 0 | 0 | 2 | 0 | 1 | 0 | 0 | 0 | 5 | 0 |
| RAST | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| BLAST[4] | **7** | **6** | **7** | **5** | **3** | **5** | **6** | **3** | **4** | **2** | **3** | **3** | **5** | **3** | **4** | **4** | **4** | **4** | **6** | **0** |
| **Attachment**: *adhesin, fim, pili, pilus, fimbriae, fimbrilin* | | | | | | | | | | | | | | | | | | | | |
| NCBI | 7 | 6 | 5 | 1 | 8 | 4 | 4 | 1 | 12 | 3 | 1 | 1 | 6 | 1 | 1 | 10 | 3 | 1 | 5 | 0 |
| RAST | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| BLAST | **16** | **17** | **18** | **14** | **11** | **14** | **17** | **12** | **16** | **11** | **13** | **13** | **13** | **13** | **12** | **15** | **15** | **11** | **14** | **0** |
| **Heme**: *heme, haga, hagb, hagc, hemaglu, hemoglo* | | | | | | | | | | | | | | | | | | | | |
| NCBI | 4 | 1 | 2 | 4 | 4 | 1 | 2 | 3 | 3 | 1 | 4 | 4 | 4 | 4 | 5 | 1 | 2 | 4 | 4 | 1 |
| RAST | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 1 | 3 | 1 | 2 | 1 | 1 | 1 | 3 | 1 | 1 | 1 | 4 |
| BLAST | **10** | **11** | **10** | **8** | **8** | **10** | **10** | **6** | **6** | **7** | **5** | **6** | **7** | **5** | **8** | **8** | **8** | **7** | **10** | **5** |
| **Gene Mobility**: *transposon, ISPg, transposase, conjugation, insertion element* | | | | | | | | | | | | | | | | | | | | |
| NCBI | 118 | 68 | 94 | 73 | 20 | 98 | 73 | 25 | 30 | 13 | 35 | 23 | 14 | 26 | 14 | 26 | 25 | 45 | 65 | 32 |
| RAST | 46 | 50 | 56 | 48 | 35 | 64 | 69 | 42 | 57 | 48 | 57 | 38 | 28 | 40 | 27 | 87 | 46 | 61 | 51 | 24 |
| BLAST | **131** | **133** | **139** | **138** | **46** | **149** | **132** | **46** | **71** | **54** | **66** | **45** | **40** | **43** | **34** | **110** | **68** | **72** | **89** | **37** |
| **Transposase, IS5 family**; K07481[5] | | | | | | | | | | | | | | | | | | | | |
| KEGG Orthology | 47 | 45 | 45 | 13 | 3 | 27 | 16 | 0 | 1 | 1 | 1 | 0 | 3 | 2 | 2 | 2 | 14 | 1 | 22 | 0 |
| **Capsule**: *capsul* | | | | | | | | | | | | | | | | | | | | |
| NCBI | 2 | 3 | 3 | 3 | 1 | 4 | 4 | 1 | 1 | 3 | 1 | 1 | 3 | 1 | 1 | 4 | 3 | 1 | 2 | 1 |
| RAST | 3 | 3 | 3 | 3 | 3 | 3 | 2 | 2 | 3 | 3 | 2 | 2 | 2 | 3 | 3 | 2 | 3 | 3 | 3 | 0 |
| BLAST | **6** | **6** | **6** | **6** | **6** | **6** | **5** | **5** | **6** | **6** | **5** | **5** | **5** | **6** | **6** | **6** | **6** | **6** | **6** | **1** |
| **CRISPR**: *crispr* | | | | | | | | | | | | | | | | | | | | |
| NCBI | 4 | 12 | 11 | 11 | 15 | 15 | 1 | 5 | 0 | 0 | 5 | 12 | 2 | 5 | 5 | 6 | 14 | 10 | 11 | 7 |
| RAST | 12 | 12 | 12 | 12 | 12 | 12 | 0 | 6 | 0 | 3 | 5 | 12 | 3 | 5 | 5 | 5 | 11 | 8 | 12 | 7 |
| BLAST | **14** | **14** | **14** | **15** | **15** | **15** | **1** | **6** | **0** | **3** | **6** | **13** | **3** | **6** | **5** | **6** | **14** | **11** | **15** | **7** |
| CRISPR arrays[6] | 3 | 3 | 3 | 4 | 5 | 5 | 2 | 3 | 3 | 3 | 7 | 22 | 3 | 15 | 4 | 3 | 4 | 7 | 5 | 2 |
| **Phage**: *phage* | | | | | | | | | | | | | | | | | | | | |
| NCBI | 1 | 1 | 3 | 1 | 6 | 1 | 2 | 10 | 13 | 1 | 9 | 5 | 8 | 8 | 12 | 3 | 1 | 6 | 3 | 1 |
| RAST | 3 | 2 | 3 | 4 | 4 | 4 | 4 | 2 | 6 | 5 | 2 | 2 | 6 | 4 | 7 | 3 | 2 | 1 | 1 | 2 |
| BLAST | **13** | **13** | **15** | **20** | **18** | **19** | **22** | **19** | **25** | **19** | **18** | **13** | **17** | **18** | **25** | **17** | **13** | **13** | **12** | **3** |

[1] Results were compiled based on the NCBI or RAST genome annotations. Total number of proteins containing any of the keywords shown in each category were recorded for each genome and for NCBI and RAST annotations separately. The detail results are provided in the Supplemental Files X-X and can be downloaded from the FTP site: ftp://bioinformatics.forsyth.org/publication_data/20160425/

[2] The keyword search was perform in a case-insensitive manner and allow matching of the partial word.

[3] The order of genomes was based on that similar to the 16S rRNA phylogenetic tree.

[4] BLAST: all the proteins identified by NCBI and RAST were collected and the sequences searched against all the proteins of all 20 genomes using BPLSTP. The numbers indicated for each genome are the number of proteins with ≥ 95% sequence identity and ≥ 95% coverage of the query sequences. The numbers were calculated separated for NCBI and RAST annotated proteins, and the larger number of the two are shown in this table.

[5] The number of proteins related to the IS5 transposase family was identified by the BlastKOALA program (Kanehisa et al., 2016) with the matching to the KEGG Orthology (KO) number K0748

[6] The number of CRISPR arrays detected by the online software  CRISPRfinger (http://crispr.i2bc.paris-saclay.fr/Server/)  ; only  the number of  "confirmed" candidates  were reported   thus excluding those "questionable" ones, which only have two DR and one spacer sequences.

**Figure Legends**

**FIGURE 1| Phylogenetic tree of *P. gingivalis 16S rRNA* gene sequences.**

A total of 24 unique *16S rRNA* gene sequences were extracted from the genomes of 19 *P. gingivalis* strains annotated by NCBI.  Sequences were pre-aligned with MAFFT v6.935b (2012/08/21) (Katoh and Standley, 2013) and leading and trailing sequences not present in all sequences were trimmed. The trimmed aligned sequences represent 20 unique sequences and were subject to QuickTree V 1.1 (Howe et al., 2002) using the "-kimura" option to calculate the substitution rate. Sequence of *P. asaccharolytica* strain DSM 20707 (PaDSM20707) was used as out-group. The branch length of the out-group was truncated to fit the tree in the figure and the substitution rate is indicated with the blue number. The red numbers next to the branching point are the bootstrap values based on 100 iterations. Sequences of different strains were separated by semicolons and the number of sequences were indicated in the parentheses in the format of (x – y / z), where x and y are the start and end IDs and z the total number in the strain.

**FIGURE 2| Core and unique genes in *P. gingivalis* surveyed by sequence identity and alignment length.**

All 37,667 protein sequences that were annotated by NCBI and with length ≥ 50 amino acids were searched for homologous clusters using the "blastclust" software V.2.2.25 (http://www.ncbi.nlm.nih.gov/Web/Newsltr/Spring04/blastlab.html) with various % identity and sequence alignment length as parameters.

**FIGURE 3| Unique proteins in 19 *P. gingivalis* strains.**

Unique proteins of each of the 19 *P. gingivalis* genomes were identified as proteins found in only one genome without any similar counterpart in any other. The cutoffs for defining the similar counterpart were from 10 to 90% sequence identities and two alignment length cutoffs at 50% (**A**) and 90% (**B**) respectively. Software "blastclust" V.2.2.25 (http://www.ncbi.nlm.nih.gov/Web/Newsltr/Spring04/blastlab.html) was used with varying % identity and alignment cutoffs as described. The strains were ordered somewhat according to the 16S rRNA phylogenetic tree showed in **Figure 1**.

**FIGURE 4|*P. gingivalis* phylogenomic trees based on core proteins identified at various percent sequence identities.**

**A**) unrooted tree based on the 1,045 shared proteins identified by "blastclust" with 60% as the sequence identity and 90% as the alignment length cutoffs; the alignment generated a total of 17,389

effective (non-identical) protein sequence positions across all 19 genomes and the tree was constructed based on these positions; **B**) rooted tree based on 436 proteins (out of 1,045) that are also found in *P. asaccharolytica* strain DSM 20707 (PaDSM20707) with ≥ 50% sequence identity and ≥ 90% alignment length; the alignment generated 4,771 effective protein sequence positions; **C**) rooted tree based on 36 proteins shared among 20 genomes with ≥ 80% sequence identity and ≥ 90% alignment length. Proteins were aligned with MAFFT v6.935b (2012/08/21) (Katoh and Standley, 2013) and poorly aligned regions were filtered by Gblocks 0.91b (Talavera et al., 2007). Trees were constructed with FastTree 2.1.9 ( Price et al., 2010) using the JTT protein mutation model (Jones et al., 1992)  and CAT+–gemma options to account for the different rates of evolution at different sites. The reliability of tree splits were reported as "local support values" based on Shimodaira-Hasegawa test (Shimodaira and Hasegawa, 2001) and are printed in blue on the split. The branch length (substitution rate) of the outgroup PaDSM20707 was truncated and the length were printed in black (**B** and **C**); **D)** Rooted tree constructed using PhyloPhlAn (Segata et al. 2013) by directly subjecting all NCBI annotated proteins of the 20 genomes to the software, resulting in 840 effective protein positions from 225 aligned proteins.

**FIGURE 5|** DNA-DNA sequence alignment between *P. gingivalis* genomes. Genomic sequence alignment between several pairs of *P. gingivalis* strains were plotted using NUCmer (NUCleotide MUMmer) version 3.1 (Delcher et al., 2002). The sequence percent identities of detected homologous fragments were plotted in gradient colors based on the percentage. The axes are the nucleotide coordination in the genomes. The orders of the contigs in the unfinished genomes were rearranged based on the reference genome (genome on X- axis).

**FIGURE 6| Genomic DNA similarity of 19 *P. gingivalis* genomes compared by oligonucleotide frequency.**

All possible 20-mer sequences present in all genomes, including that of *P. asaccharolytica* strain DSM 20707 (PaDSM2070) used as an out-group, were categorized and the number of genomes in which a 20-mer is present, was recorded. Panel **A** was generated by first calculating the average number of genomes for all the 20 mers present in every 500-nucleotide windows across the entire genome and then color each window based on the genome frequency (minimum 1 in yellow and maximum 20 in black). Panel **B** was similar to **A** but the non-coding regions were masked with light blue color to highlight the oligonucleotide frequencies for the areas that correspond to both forward (upper) and reverse-complement (lower) protein coding sequences. The order of the unfinished genomic contigs was arranged in the same order as appeared in the sequences downloaded from NCBI. The genomes in the plot were ordered based on the 16S rRNA phylogenetic tree (**Figure 1**) with a dendrogram derived from the same tree to show the relatedness.
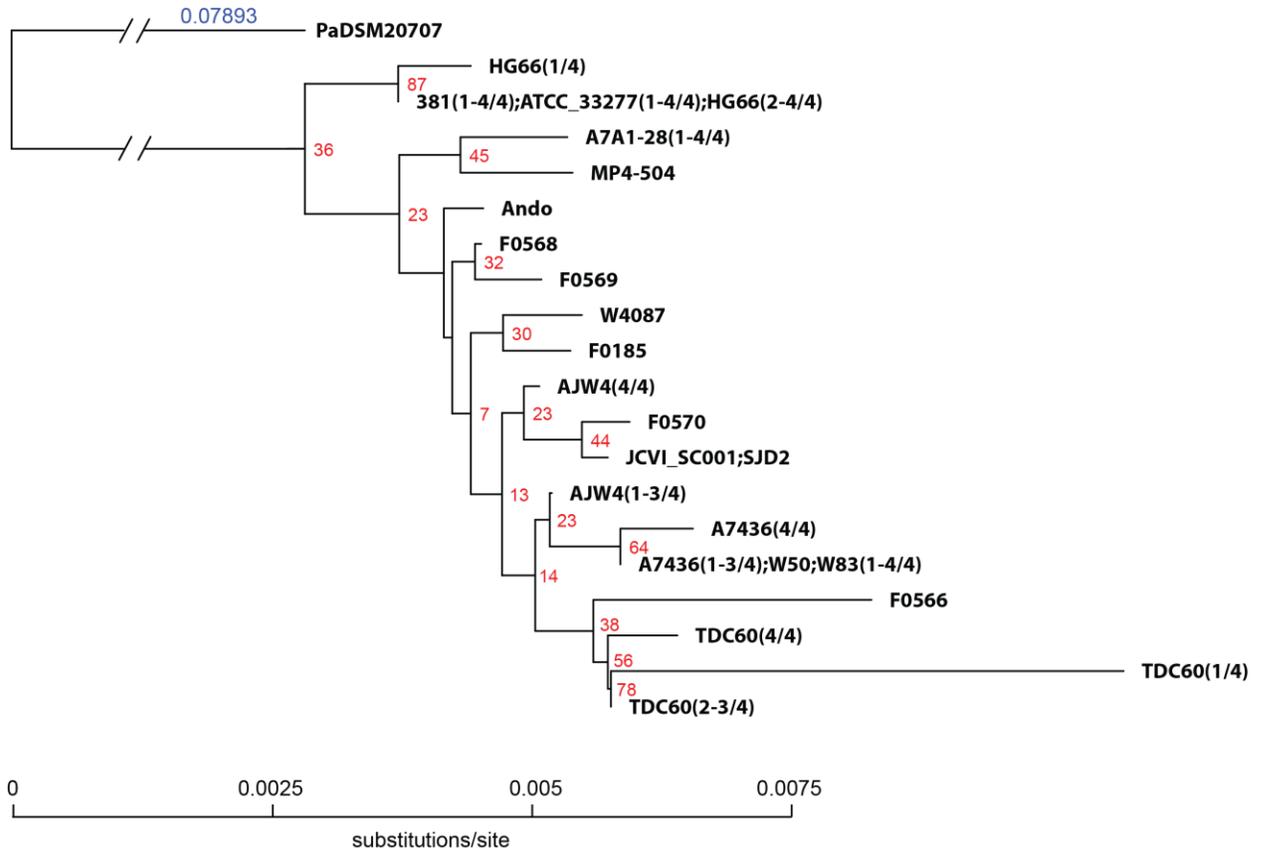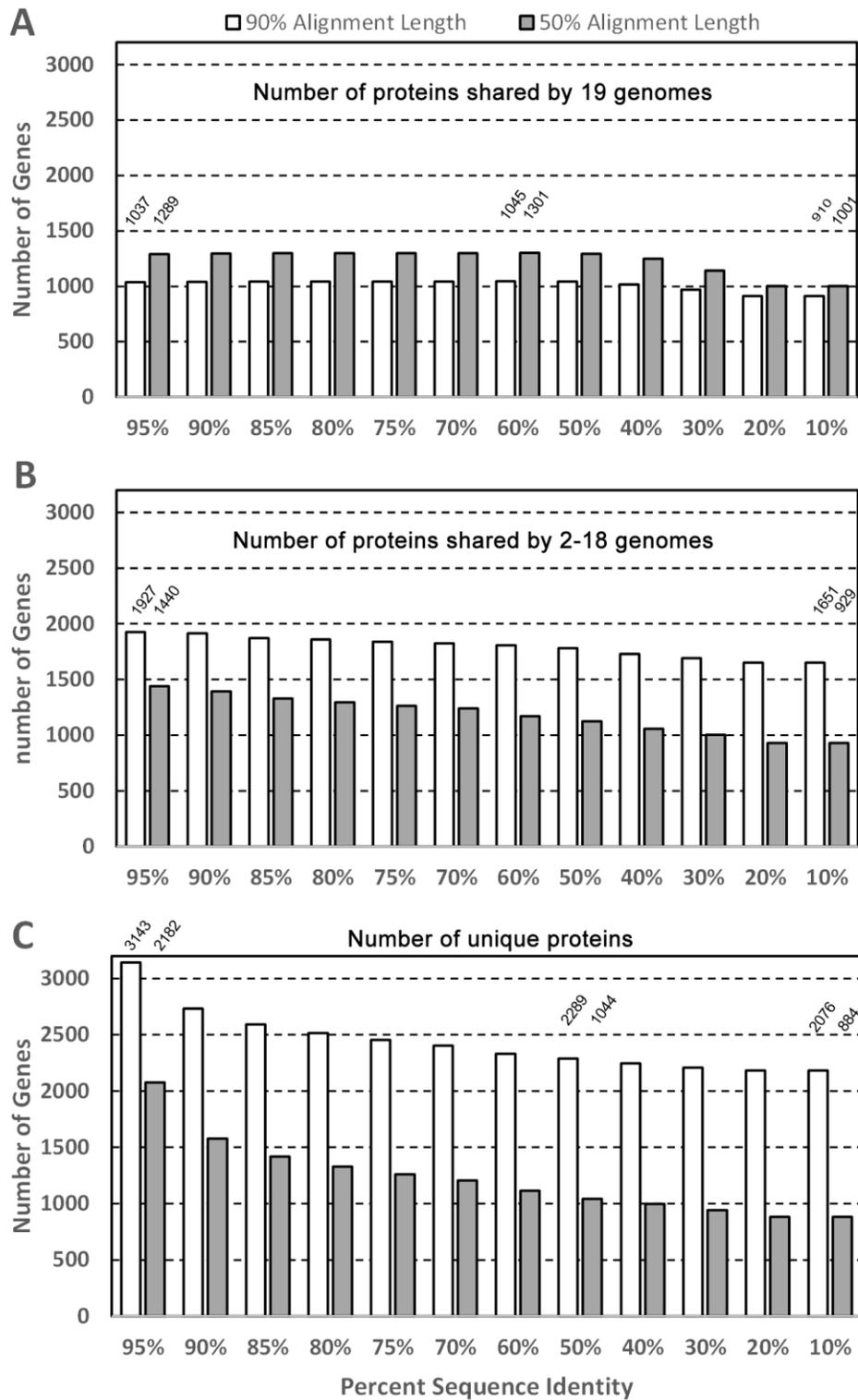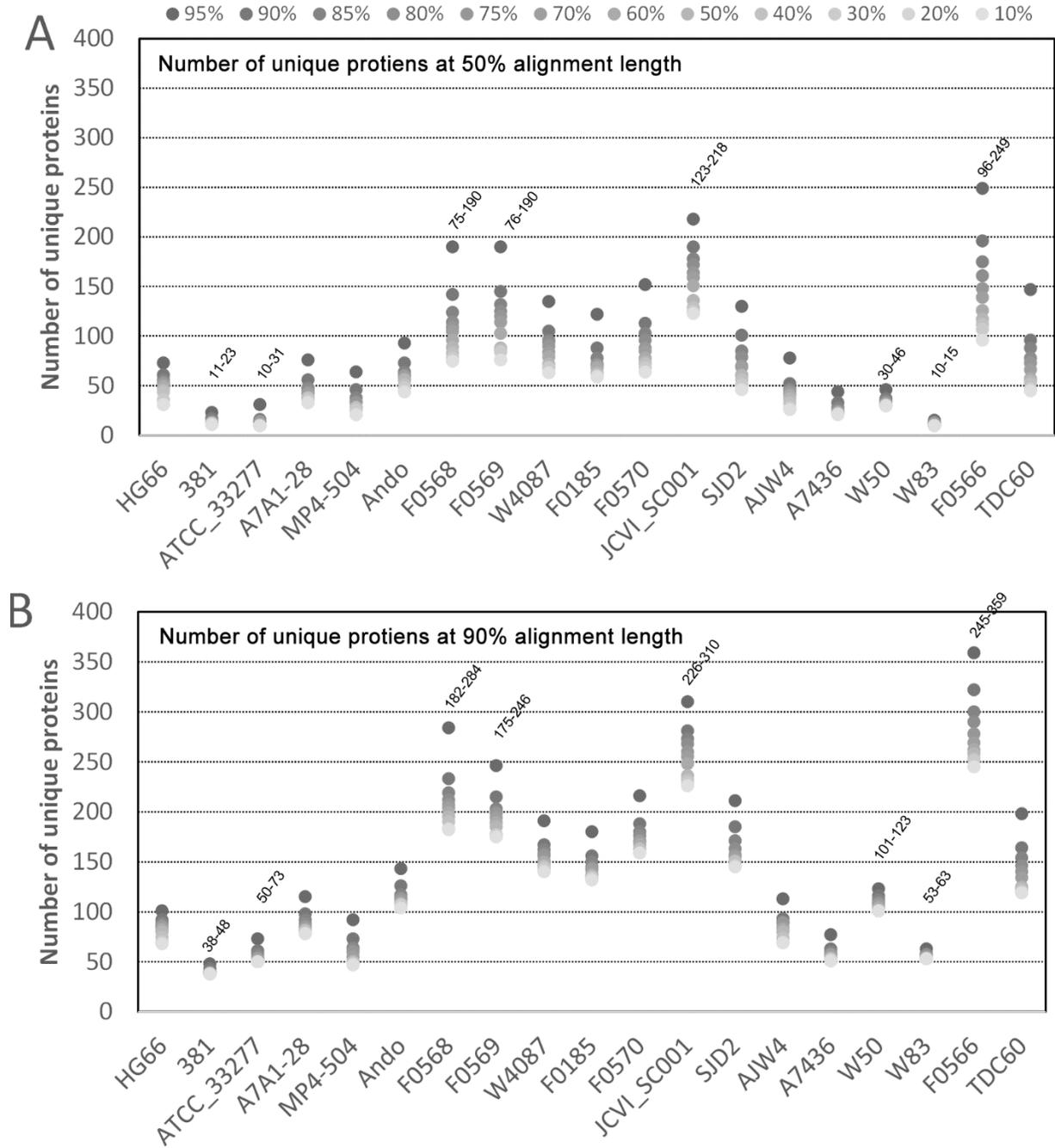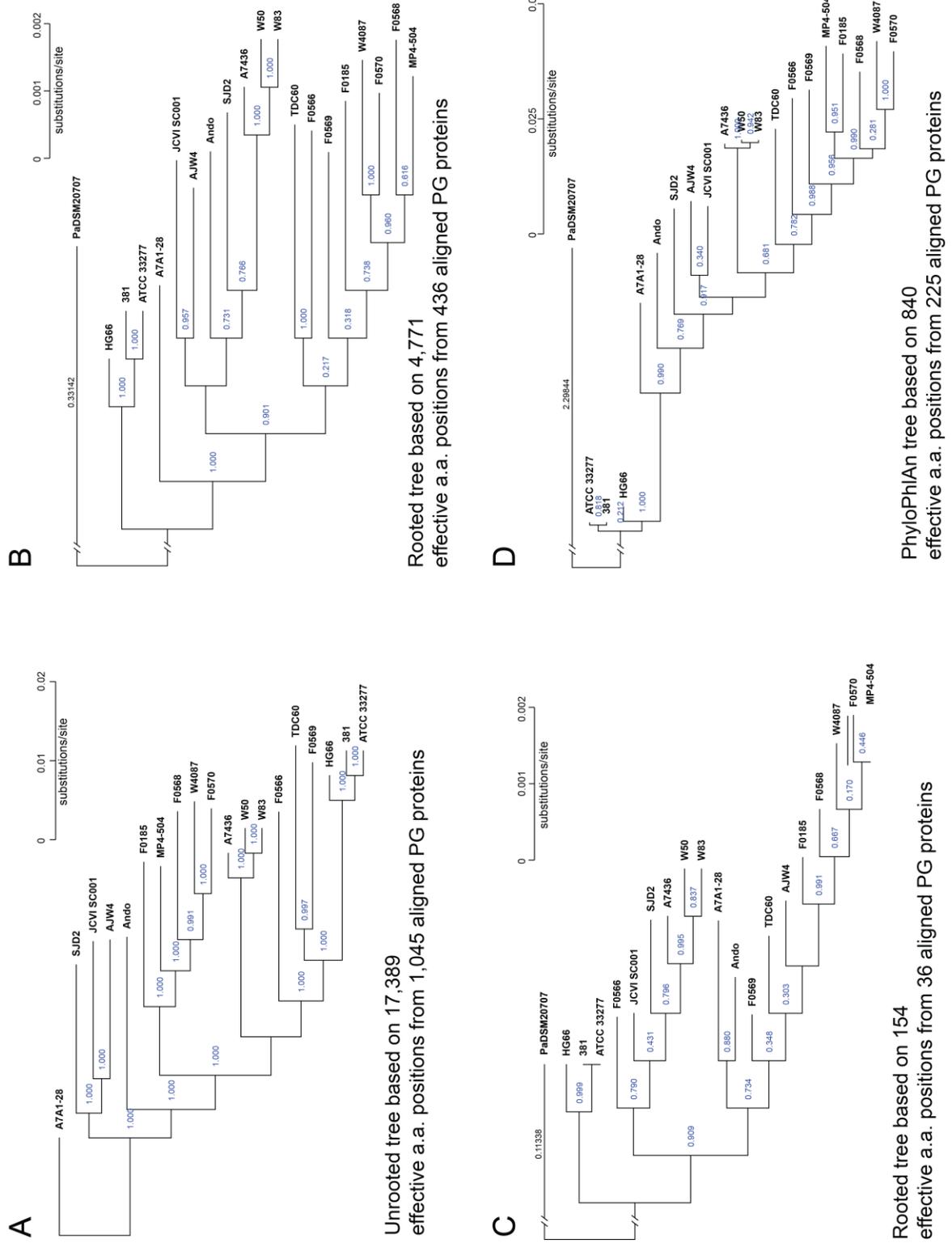
**FIGURE 1**

**FIGURE 2**

**FIGURE 3**

**FIGURE 4**



A

Unrooted tree based on 17,389
effective a.a. positions from 1,045 aligned PG proteins

B

Rooted tree based on 4,771
effective a.a. positions from 436 aligned PG proteins

C

Rooted tree based on 154
effective a.a. positions from 36 aligned PG proteins

D

PhyloPhlAn tree based on 840
effective a.a. positions from 225 aligned PG proteins

**FIGURE 5**

**FIGURE 6A**



**FIGURE 6B**