- 1 DPSN: standardizing the short names of amplicon-sequencing primers to avoid ambiguity
- 2
- 3 Yuxiang Tan<sup>#1</sup>, Yixia Tian<sup>#1</sup>, Junyu Chen<sup>2</sup>, Zhinan Yin<sup>1</sup>, Hengwen Yang<sup>\*1</sup>
- 4
- 5 <sup>1</sup> The First Affiliated Hospital, Biomedical Translational Research Institute, Guangdong Province
- 6 Key Laboratory of Molecular Immunology and Antibody Engineering, Jinan University,
- 7 Guangzhou 510632, China
- 8 <sup>2</sup> Shenzhen Institute of Synthetic Biology, Shenzhen Institutes of Advanced Technology, Chinese
- 9 Academy of Sciences, Shenzhen 518055, People's Republic of China
- 10 <sup>#</sup> These authors contributed equally
- 11
- 12 \* Corresponding authors: Dr. Hengwen Yang <u>benyang97@gmail.com</u>. Tel/ Fax: +86-020-85222787

### 14 Abstract

15	Background: Amplicon sequencing is the most widely used sequencing method to evaluate
16	microbial diversity in virtually all environments. Thus, appropriate and specific primers are needed
17	to amplify amplicon regions in amplicon sequencing. For this purpose, the community currently
18	uses probeBase, which curates rRNA-targeted probes and primers. Main Body: We found that 63.58%
19	of the primers in probeBase have problematic issues in the short name, full name, and/or position.
20	Furthermore, the current convention for short names causes ambiguity. We here introduce our new
21	Database of Primer Scientific Names (DPSN), which is a manually curated database for the 173
22	primers in probeBase complete with a new short name convention. Building on the work of
23	probeBase, we provide a more user-friendly and standardized system. The new short primer naming
24	convention has three basic components: 5' position on the sense strand, version, and direction. An
25	additional character for the name of the taxonomic group is also added in front of the name for
26	convenient use. Furthermore, DPSN contains primers for large subunit as well. In order to separate
27	them from the primers for small subunit, a header character is also recommended. Conclusion: All
28	173 primers in probeBase were corrected according to this new rule, and are stored in DPSN, which

29 is expected to facilitate accurate primer selection and better standardized communication in this

30 field.

- 31 Database URL: The DPSN database is available in a user-interactive website at
- 32 http://dpsn.gdimmunity.com
- 33 Keywords: database, scientific name, amplicon, primer

# 35 Background

36	Amplicon sequencing is a common sequencing method for microbial research from diverse
37	environmental or clinical samples [1, 2]. Amplicon sequencing is dependent on the choice of primers
38	for carrying out the amplification step. Thus, selection of the most appropriate primers is the
39	foundation of successful amplicon sequencing.
40	probeBase [3] is the only database currently available with updated lists of probes and primers,
41	along with links to other databases providing related information. At present, there is a total of 173
42	primers recorded in probeBase. In general, a primer is defined according to its short name (SN), full
43	name (FN), and sequence. However, in many cases, only the SN is used for the sake of convenience.
44	There are seven components of an FN [4]. Taking S-D-Bact-0338-a-A-18 as an example: "S" stands
45	for the target gene (Small Sub-Unit (SSU)), "D" represents the largest taxonomic group targeted
46	(Domain), "Bact" is the name of the taxonomic group (Bacteria), "0338" is the 5' position of the
47	sense strand, "a" presents the version, "A" denotes the identical strand ("S" for sense; "A" for
48	antisense), and "18" is the length of the primer. To avoid ambiguity, each primer should have a
49	unique SN; however, this is not the case. Different from FN, there is no guideline for how an SN

50	should be. Therefor	e, SNs were na	amed in a few	different ways,	such as P	rimer3, Bac9	27, 926r, and

- 51 934mcr. The most common ones were composited by the position and direction (for example, 926r),
- 52 or with an additional string for the name of the taxonomic group (for example, Arch 915r). The lack
- 53 of clear rules and sufficient information for accuracy leads to ambiguity of SNs.
- 54 In fact, there are 14 SNs that refer to multiple primer sequences, which could lead to confusion
- and cause several problems in application for users. For example, in the earth microbiome project
- 56 website [1], the author of the citation for a given primer is used along with the SN to better specify
- 57 the primer. Furthermore, the SN itself could be misleading. For example, primers 907r and 926r are
- 58 actually from the same region of the genome but with a difference of two bases in the sequence.
- 59 However, based on their SNs alone, a user would misinterpret these primers as being derived from
- 60 two different regions.
- 61 To resolve this problem, we here introduce Database of Primer Scientific Names (DPSN),
- 62 which is a database that has been manually curated to correct problematic and inconsistent features
- 63 (SN, FN, position, and length) of primers in probeBase according to an improved convention of
- 64 naming SNs. The new SNs still correspond to the old SNs and corrected FNs in a one-to-one relation.

### 65 **Construction and content**

#### 66 Data source

- 67 Information of all 173 primers in the probeBase dataset was manually extracted [3], including
- 68 the SN, FN, position, sequence, length, G+C content, and dissociation temperature.
- 69 The corresponding regions on the reference sequence of Escherichia coli K-12 substrain
- 70 MG1655 was extracted from the SILVA database [5] and served as the reference for confirming the
- 71 sequence position.

#### 72 Derivation of new SN naming convention

- 73 A unique SN should have at least three basic components to provide sufficient identifying
- 74 information: 5' position on the sense strand, version, and direction.
- 75 In the old SN, such as 907r, all reverse primers that start or end from position 907 will have
- the same name, which leads to ambiguity. Therefore, including additional information of the version,
- such as 907ar to indicate the version, could help to specify the primer sequence. Consistent with the
- old SN rule, "f" and "r" denote "forward primer" and "reverse primer," respectively. Moreover,

79	because the name of the taxonomic group provides helpful information for users to select
80	appropriate primers, DPSN also includes a shorthand for the name of the taxonomic group ("A" for
81	Archaea, "B" for Bacteria, "U" for universal, and "N" for nano) in front of the SN. Additionally, on
82	account of the need to separate SSU primers from large subunit (LSU) primers, the header represents
83	of target get from the FNs is retained. For instance, the SN 907ar represents the FN "S-D-Bact-
84	0907-a-A-20", which was corrected and recorded as S-B907ar in DPSN.

### 85 Amplification range validation of the primers

86	To validate the am	plification location of	primers according	g to the E.	coli K-12 refer	rence, BLAT
----	--------------------	-------------------------	-------------------	-------------	-----------------	-------------

87 of the National Center for Biotechnology Information [6] was employed as the aligner. However,

- 88 because of the presence of degenerate bases in primer sequences, the primers had to be converted
- 89 into expanded regular sequences, which was achieved using a customized Python script before
- 90 BLAT alignment. In particular, the additional parameters "-minMatch = 1-minScore = 8-minIdentity
- 91 = 70-stepSize =1-tileSize =8" were applied to BLAT, considering the length and mismatch tolerance

92 of primers.

93 To confirm the amplification location, the primers were also checked by the TestProbe function

94 in SILVA [5].

95

### 96 Utility and Discussion

#### 97 How to use DPSN

98	In order to make the search easy, DPSN supports fuzzy search on all the fields. This means
99	user can use any keyword to find the related primer(s), such as the intended 5' position and the sub-
100	string of the primer sequence. In return, DPSN will present the corrected information of the related
101	primers. As well as the original "Short Name" and "Full Name" from the probeBase, which will
102	help the user to connect the use of the primers in original papers. All the sequence in DPSN are the
103	same as the ones in probeBase.
104	Summary of Corrections and Discussion
105	Our careful review of probeBase identified five sequences with multiple primer names, 14
106	groups of SNs with multiple targets, and 91 SNs inconsistent with their FNs. Of the total 173 primers
107	in probeBase, the SNs for only 63 primers (36.42%) could direct the user to a unique sequence and

108 be considered as correct. Five sequences were multifold and had multiple primer names (Table 1).

109	Thirty SNs from 14 groups pointed to more than one sequence. The positions in 91 SNs were
110	different from the 5' position of their FNs. Overall, the positions of 35 primers in probeBase were
111	found to be incorrect.
112	In addition, a few FNs were found to be incorrect in probeBase, which have been manually
113	corrected in DPSN. Theoretically, the FNs of primers should be unique, since a single FN represents
114	a unique primer sequence. However, in probeBase, three FNs were duplicated and even represented
115	more than one sequence (Table 2). In the naming rule, the position in an FN is based on the $5'$
116	position; however, eight of the primers in probeBase violated this rule (Table 3). Even more
117	importantly, the length information of five FNs did not match the actual lengths of their sequences
118	(Table 4), and the directions of three primers were opposite to the actual direction of their sequences
119	(Table 5).

120	Table 2: Primer groups with the same FN in probeBase									
	Old SN	Old FN	Position	Sequence	GC%	TM*				
	Arch958f	S-D-Arch-0958-a-S-19	958–975	AATTGGANTCAACGCCGG	50	49				
	Arch958Bf	S-D-Arch-0958-a-S-19	958–976	AATTGGABTCAACGCCGGR	47	51.5				

b785	S-D-Bact-0785-a-A-19	785–803	CTACCAGGGTATCTAATCC	47	49
803r	S-D-Bact-0785-a-A-19	785–803	CTACCRGGGTATCTAATCC	47	50
518r	S-D-Bact-0518-a-A-17	518–534	ATTACCGCGGCTGCTGG	65	52
P518r	S-D-Bact-0518-a-A-17	518–534	ATTACCGCGGCTGCTGG	65	52

121 \*TM: dissociation temperature (°C)

122

#### Table 3: FNs of primers with inconsistent positions in probeBase

Old SN	Old FN	Position	Sequence	GC%	TM*
338	S-D-Bact-0338-a- A-19	337–355	TGCTGCCTCCCGTAGGA GT	63	58
1114mcr	S-P-Nano-1082- a-A-17	915–931	GGGTCTCGCCTGTTTCC	65	52
27F	S-D-Bact-0008-d- S-20	6–25	AGAGTTTGATCMTGGCT CAG	45	51
63F	S-D-Bact-0043-s- S-21	21–41	CAGGCCTAACACATGCA AGTC	52	52
Arch855R	S-D-Arch-0896- a-A-20	915–934	TCCCCCGCCAATTCCTTT AA	50	52

	bio-pJBS-V3	3.SER	S-D-Bact-094 A-20			GGTAAGGTT 5–965 TGC		CTTCG	CGT	55	53
	Primer3		S-D-Bact-0513 A-17	8-c-	340–3	357	GCCTACGGC G	GAGGCA	AGCA	72	57
	Primer2		S-D-Bact-034 A-18	0-a-	518-:	534	ATTACCGCG	GCTGC	TGG	65	52
123	*TM: dissocia	tion tem	perature (°C)								
124											
125	Table 4: FNs of p   Old SN Old FN		orime	rs wit	h the	wrong length i	n probel	Base			
			Position			Sequence	GC%	TM*	Correct	ted	
	Uni522r	S-*-U 15	Jniv-0517-a-A-	517–534		GWATTACCGCG GCKGCTG		61	54		18
	Primer2	S-D-H 18	3act-0340-a-A-	518-	518–534		ACCGCGGC CTGG	65	52		17
	U529r	S-*-U 18	Jniv-0522-a-A-	522–536 340–357		ACCGCGGCKGC TGGC		80	54.5		15
	Primer3	S-D-H 17	Bact-0518-c-A-				CTACGGGAG AGCAG	72	57		18
	Arch958f	S-D-4 19	Arch-0958-a-S-	958-	958–975		TGGANTCA GCCGG	50	49		18

126 \*TM: dissociation temperature (°C)

#### 127

128

129	Table 5: FNs of primers with the wrong strand in probeBase								
	Old SN	Old FN	Position	Sequence	GC%	TM *	Corrected FN S-D-Bact-0340-a- S-18		
	Primer3	S-D-Bact-0518- c-A-17	340–357	GCCTACGGGAGG CAGCAG	72	57			
	527f	S-D-Bact-0517- a-S-16	517–532	ACCGCGGGCCKGC TGGC	81	66	S-D-Bact-0517-a- A-16		
	536r	S-D-Bact-0519- a-A-18	519–536	CAGCMGCCGCG GTAATWC	61	54	S-D-Bact-0519-a- S-18		

130 \*TM: dissociation temperature (°C)

131 In DPSN, all of the SNs of the primers in probeBase have been updated according to the new

132 naming rule along with additional version information to provide a more unique identifier that is

133 still convenient to use. Overall, 110 problematic primers were corrected. Using the amended primer

134 name in DPSN, users can simply refer to the SN to specify a primer, because of the one-to-one

135 relation among the SN, FN, and sequence, and without the inconvenience of appending additional

136 information such as author name or sequence in the article.

# 137 Conclusion

138	In conclusion, because it is crucial to avoid vagueness in scientific research, the old SN system
139	of primers is problematic and should be replaced by the new naming rule as proposed herein. All of
140	these corrections have been curated in DPSN to improve searching convenience and accuracy.
141	Therefore, with DPSN, users can easily search an old name from probeBase or articles for its
142	amended name. For new articles, it is recommended that authors use the amended name to
143	accurately describe a primer.
144	DPSN currently focuses on only primers for amplicon sequencing on SSU and LSU, and thus
145	it can be assumed that the ambiguity problem still exists for primers in other amplicon regions, such
146	as ITS. Because the primers for these regions were not found in probeBase, we can collect and
147	import these primers into the naming system in DPSN in the future. To keep the database up to date,
148	DPSN accepts data submission of primers from researchers as well.

# 149 Abbreviations

150 **DPSN**: Database of Primers' Scientific Names **SN**: short name FN: full name

# 151 **Declarations**

- 152 Ethics approval and consent to participate
- 153 Not applicable
- 154 Consent for publication
- 155 Not applicable
- 156 Availability of data and material
- 157 <u>Database Name:</u> Database of Primers' Scientific Names (DPSN)
- 158 <u>Database URL:</u> http://dpsn.gdimmunity.com

#### **159** Competing interests

160 The authors declare that they have no competing interests.

#### 161 Funding

162 This work was supported by the Science and Technology Department of Guangdong Province of

163 China (grant no. 2017A030310179 to Dr. Yuxiang Tan); The Major International Joint Research

- 164 Program of China (grant no. 31420103901), and the "111 project" (grant no. B16021) to Dr. Zhinan
- 165 Yin and Hengwen Yang. The three grants supported the whole study through design and collection,
- analysis and interpretation of data and in writing the manuscript.
- 167 Authors' contributions
- 168 YuT conceived of the idea, conducted the analysis, and wrote the manuscript. YiT performed the
- 169 data collection. JC provided data of LSU primers. HY and ZY supervised the project and
- 170 participated in the revision of the manuscript. All authors read and approved the final manuscript.

#### 171 Acknowledgements

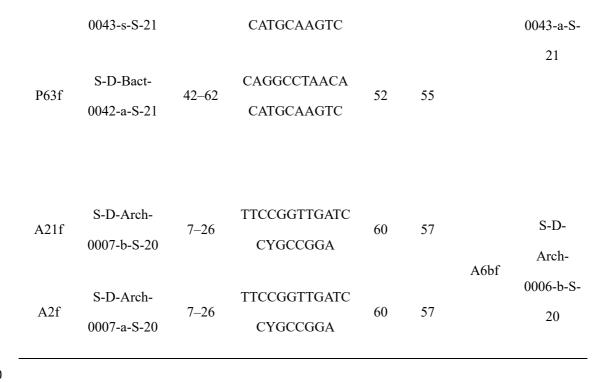
- 172 We would like to thank Becky Kusko for editing suggestions and Editage (www.editage.com) for
- 173 English language editing.

#### 174 References

- 175 1. Gilbert JA, Jansson JK, Knight R. The Earth Microbiome project: successes and aspirations.
- 176 BMC Biol. 2014;12:69.
- 177 2. The Human Microbiome Project Consortium. A framework for human microbiome research.
- 178 Nature 2012;486:215-21.

179	3.	Greuter D, Loy A, Horn M, Rattei T. probeBase-an online resource for rRNA-targeted
180		oligonucleotide probes and primers: new features. Nucleic Acids Res. 2016;44:D586-9.
181	4.	Klindworth A, Pruesse E, Schwwer T, Peplies J, Quast C, Horn M, et al. Evaluation of general
182		16S ribosomal RNA gene PCR primers for classical and next-generation sequencing-based
183		diversity studies. Nucleic Acids Res. 2013;41:e1.
184	5.	Quast C, Pruesse E, Yilmaz P, Gerken J, Schwwer T, Yarza P, Peplies J, Glöckner FO. The
185		SILVA ribosomal RNA gene database project: improved data processing and web-based tools.
186		Nucleic Acids Res. 2013;41:D590-6.
187	6.	Kent WJ. BLATthe BLAST-like alignment tool. Genome Res. 2002;12:656-4.

Table 1: Primer groups with the same sequence in probeBase											
Old SN	Old FN	Position	Sequence	GC%	TM*	Corr. SN	Corr. FN				
Arch95 8Bf	S-D-Arch- 0938-b-S-19	938–956	AATTGGABTCA ACGCCGGR	47	51.5	A958bf	S-D- Arch- 0958-b-S- 19				
Arch95 8Bf	S-D-Arch- 0958-a-S-19	958–976	AATTGGABTCA ACGCCGGR	47	51.5						
U519f	S-*-Univ- 0519-a-S-18	519–536	CAGCMGCCGCG GTAATWC	61	54	U519bf	S-*-Univ- 0519-a-S- 18				
536r	S-D-Bact- 0519-a-A-18	519–536	CAGCMGCCGCG GTAATWC	61	54						
518r	S-D-Bact- 0518-a-A-17	518–534	ATTACCGCGGCT GCTGG	65	52						
P518r	S-D-Bact- 0518-a-A-17	518–534	ATTACCGCGGGCT GCTGG	65	52	B518ar	S-D-Bact 0518-a-A 17				
Primer2	S-D-Bact- 0340-a-A-18	518–534	ATTACCGCGGGCT GCTGG	65	52						
63f	S-D-Bact-	21-41	CAGGCCTAACA	52	52	B43af	S-D-Bact				



190

191 \*TM: dissociation temperature (°C)