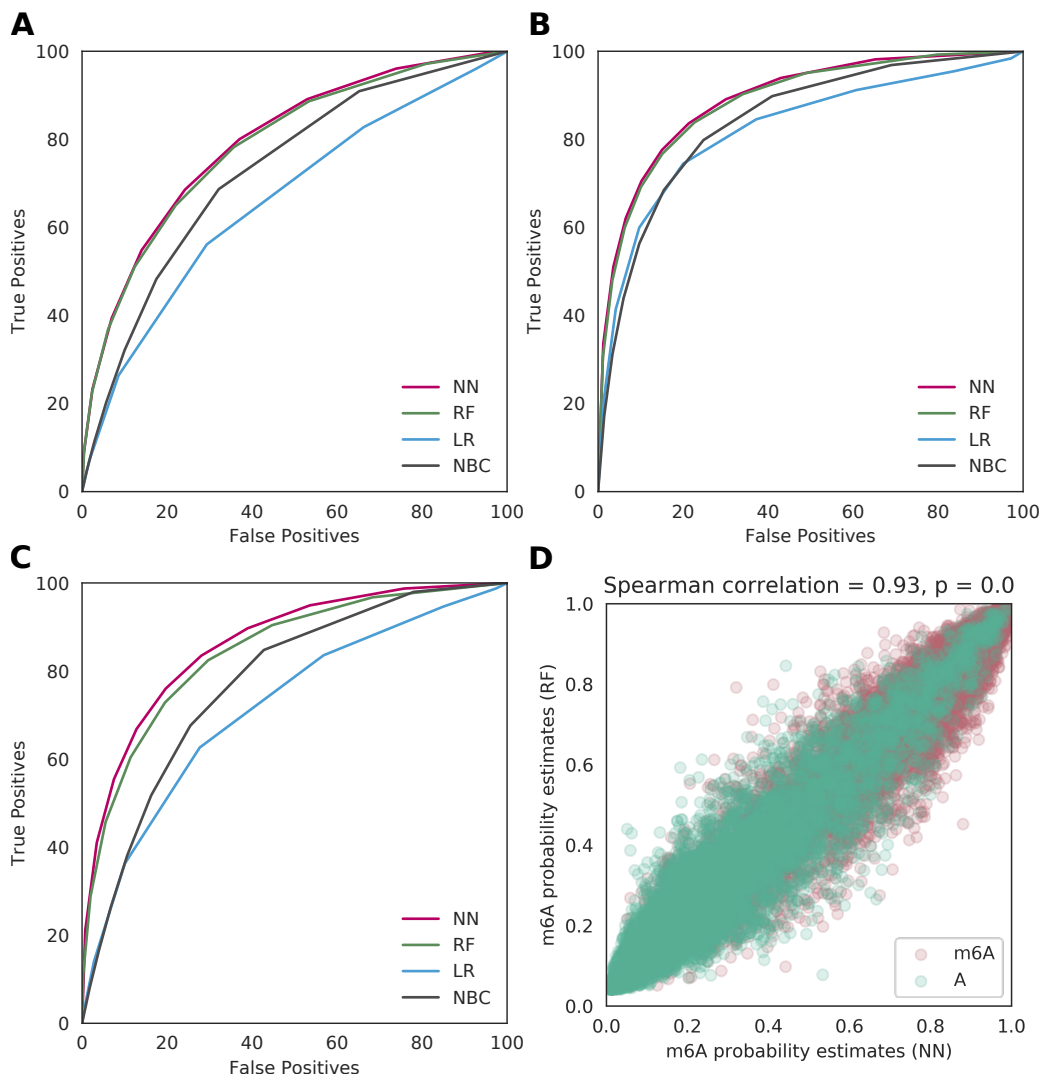


Supplementary Information

Single-molecule sequencing detection of *N*6-methyladenine in microbial reference materials

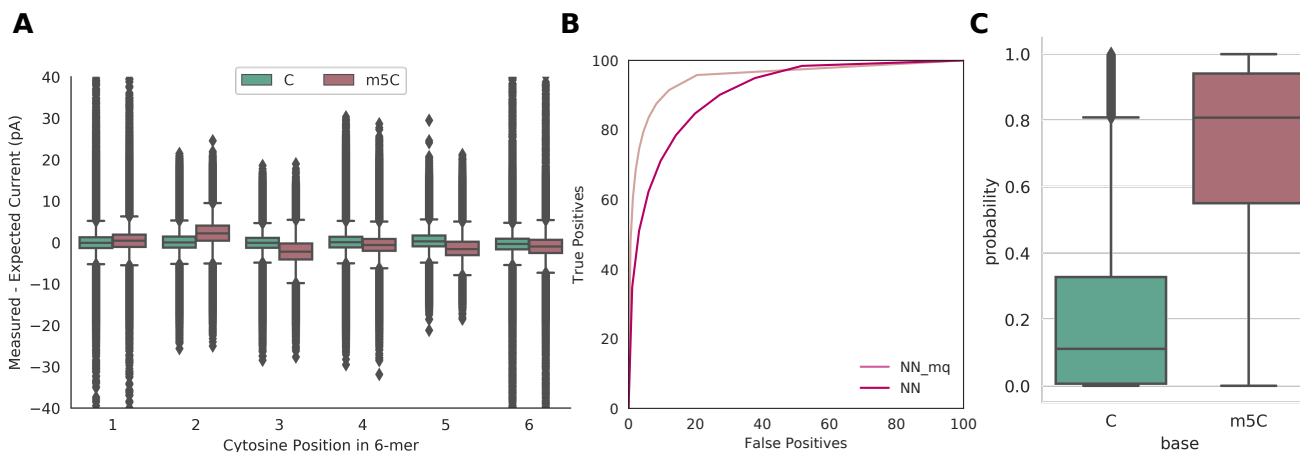
McIntyre et al.



Supplementary Figure 2. Model performance with different numbers of features. ROC plots for a support multi-layer perceptron neural network (NN), random forest classifier (RF), linear regression classifier (LR), and naïve Bayes classifier (NBC) using (A) 7-mer, (B) 11-mer, and (C) 15-mer contexts to classify. (D) A comparison of the probability scores for 10,000 randomly selected classifications per base using the top (x-axis, neural network) and second place (x-axis, random forest classifier) algorithms for 11-mer contexts. Models were generated using one R9 experiment and tested on data from a second.

# Features	Classifier	Accuracy	AUC
4	NN	72.2	0.798
4	RF	71.5	0.792
4	LR	63.4	0.671
4	NBC	65.4	0.739
6	NN	81.3	0.895
6	RF	80.8	0.890
6	LR	77.2	0.828
6	NBC	76.5	0.848
8	NN	78.2	0.865
8	RF	76.7	0.848
8	LR	67.5	0.726
8	NBC	67.7	0.782
6	NN_sk	77.8	0.862
6	NN_hq_sk	81.8	0.897
6	NN_hq	84.2	0.920
6	NN_pos	95.4	0.992

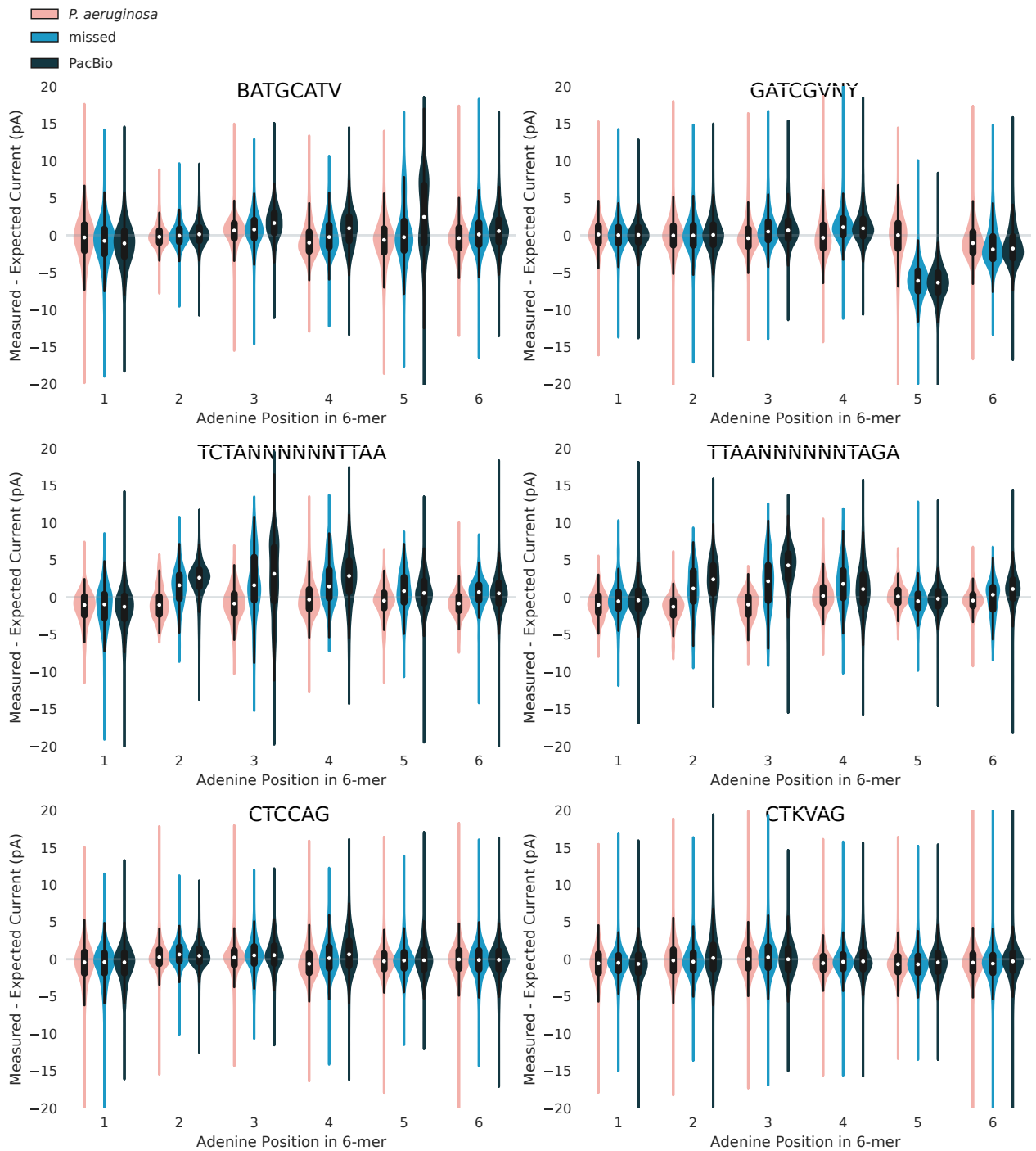
Supplementary Table 1. Model performance with different features and read qualities. Accuracy and area under the curve (AUC) for methods and parameters tested (number of features and filters) in Figure 2C and Supplementary Figure 2. Methods include neural network (NN), random forest (RF), linear regression (LR), and naïve Bayes classifier (NBC), and filters allowing up to two skips (sk), only high-quality reads (hq) of QV>9, and positions (pos) covered by 15 or more reads.



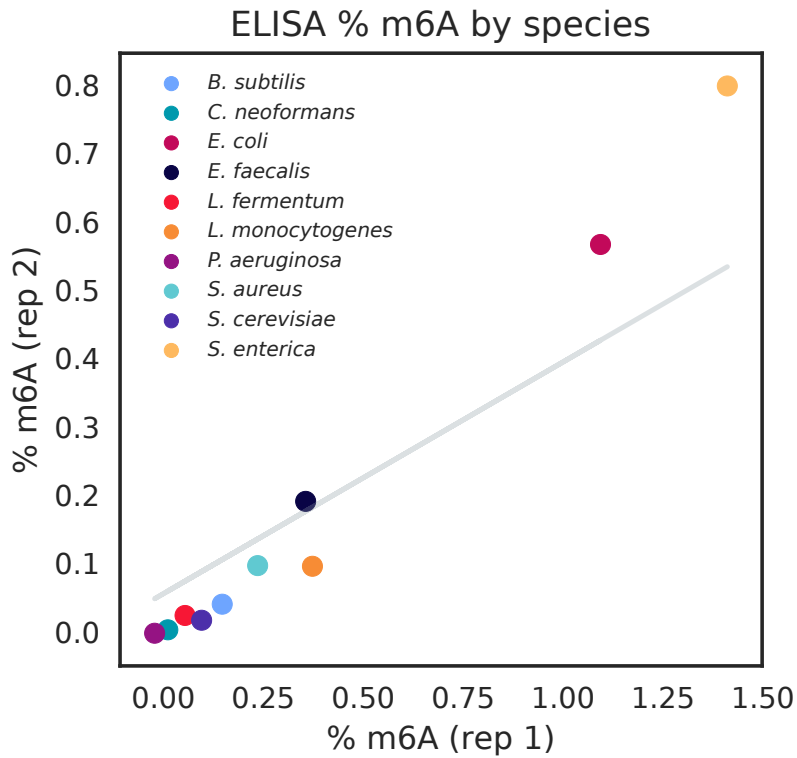
Supplementary Figure 3. Cytosine methylation detection. (A) Deviations between measured and expected current values across different positions in a 6-mer, depending on cytosine methylation status (outliers are truncated at +/- 40 pA). The variances for both methylated and unmethylated cytosines are higher than for m6A, likely in part because the reads were of lower average quality (5.64 for the M.Sssl dataset and 6.41 for the PCR dataset, vs. 7.65 for the Chiu dataset and 7.32 for the Mason dataset). **(B)** A ROC plot shows m5C classifications with per-read accuracy of 82.2%, increasing to 87.1% for “medium quality” data (“mq”, > mean base quality score of 7), similar to m6A despite including more sites with multiple methylated bases. **(C)** Probability scores for the detection of m5C using the same neural network classifier trained and tested using any quality reads from M. Sssl-modified data and whole genome amplified data. Boxplot centre lines show medians and whiskers 1.5x interquartile range.

Species (# putative m6A MTase genes)	Motif	Mean IPD Ratio	Total sites in genome	Detected in 2nd PacBio run?	PacBio + MeDIP-seq	PacBio + ONT	Sites detected by PB + ONT + MeDIP-seq	Sites missed by all methods
<i>B. subtilis</i> (2)	CNCANNNNNNRTGT	5.96	535	yes	516	291	288	3
	ACAYNNNNNNNTGNG	5.76	535	yes	516	496	482	3
<i>E. coli</i> (4)	none	2.01	6296		547	974	93	
	AAGANNNNNCTC	5.39	337	no data	153	191	97	0
	GAGNNNNNTCTT	6.03	337	no data	155	270	129	0
	GATC	5.67	39872	no data	24529	36805	12181	14
<i>E. faecalis</i> (4)	none	2.17	5054		1713	791	165	
	CAAYNNNNNNNTYG	5.40	861	yes	830	857	827	0
	CRAANNNNNNRTTG	5.17	861	yes	830	517	505	1
	CTKVAG	2.55	5909	CTBVAG	1788	1167	576	1246
	CTCCAG	2.80	478		205	89	54	48
<i>L. monocytogenes</i> (2)	none	1.84	7638		1540	1493	334	
	GCANNNNNNTGC	5.92	906	yes, also GCAAVDNNNTGC	886	570	562	0
	ANARAGTANYR	0.75	485	no	4	0	0	377
<i>P. aeruginosa</i> (1)	none	2.00	7438		626	1264	105	
	none	2.47	939		67	174	15	
<i>S. aureus</i> (4)	TCTANNNNNNTTAA	6.74	425	yes	408	267	259	3
	TTAANNNNNNTAGA	6.37	425	yes	408	3	3	5
	GAAGNNNNNTTRG	6.57	230	yes	197	104	90	1
	CYAANNNNNCTTC	5.63	230	yes	199	156	144	2
	GATCGVNY	2.23	399	GATCBVNYD	17	195	17	11
	none	2.03	4259		423	812	87	
	none	2.64	1714		233	183	28	
<i>S. cerevisiae</i>	none	2.64	1714		233	183	28	
	GATC	5.85	37910	yes	24007	33632	11885	172
<i>S. enterica</i> (3)	CAGAG	6.43	5819	yes	2048	3053	1111	110
	BATGCATV	3.81	730	BATGCAT	169	239	110	248
	none	2.47	2215		852	290	88	
	none	2.47	2215		852	290	88	

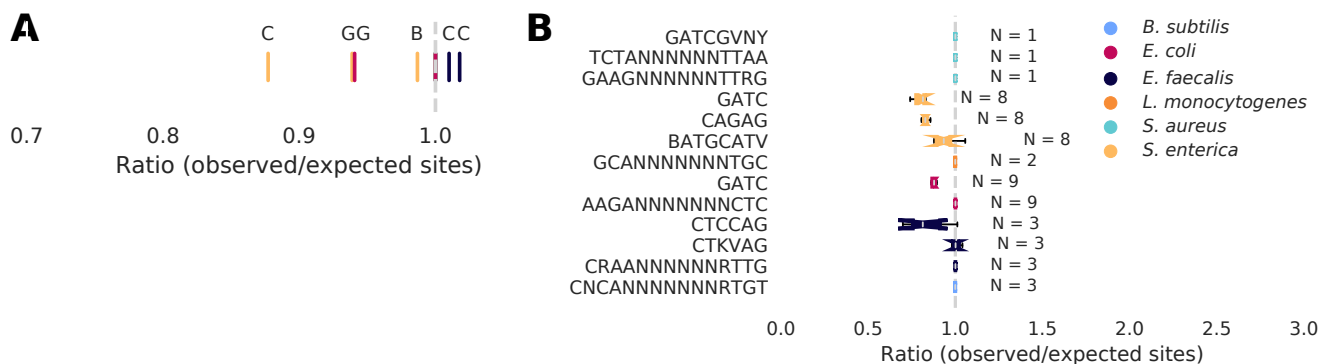
Supplementary Table 2. Predicted m6A motifs and non-motif sites predicted by PacBio across eight species. A total of fourteen motifs were predicted as methylated (nineteen including methylated reverse complements). IPD = inter-pulse duration, a measure of the rate of nucleotide addition during sequencing. mCaller was used for ONT predictions.



Supplementary Figure 4. Current differences among motifs. Plots of current differences from model values at motifs sites identified as methylated using PacBio sequencing, motif sites missed by PacBio in the same strain, and unmethylated sites for the same motif in *P. aeruginosa*. Inner boxplot centre lines show medians and whiskers 1.5x interquartile range.



Supplementary Figure 5. Comparison of ELISA replicates. The percent m⁶A in ten reference community strains (including two for which we did not have PacBio data – *C. neoformans* and *L. fermentum*), as measured using two replicates of a commercial ELISA kit. Pearson R = 0.60, p = 0.067; Spearman ρ = 0.16, p = 0.65.



Supplementary Figure 6. Motif depletion calculated with Kr. Ratio of observed over expected motif counts calculated by Kr in **(A)** assembled genomes, colored by species and labeled by first nucleotide where there were multiple motifs, and **(B)** in predicted prophages within assembled genomes. Boxplots centre lines show medians, notches confidence intervals, and whiskers 1.5x interquartile range.

Kr, expressed in counts instead of frequencies:

$$Kr(m) = L(g)^{-1^{L(m)}} \prod_{s \in S} N(s)^{-1^{(L(m)-L(s))}}$$

where

m = motif of interest

S = the set of all submotifs of m (including m)

$N(s)$ = count of submotif s in the genome or sequence of interest

g = genome or sequence of interest

$L(x)$ = length of x

Species	PacBio		Illumina (WGS)		Illumina (MeDIP-seq)		Illumina (WGBS)	ONT	
	FDA Site (RSII)	UF (Sequel)	FDA Site (HiSeq4000)	WCM (iSeq)	WCM (iSeq)	WCM (NextSeq)	WCM (HiSeq)	UCSF	WCM
<i>B. subtilis</i>	X	X	X		X		X		X
<i>E. faecalis</i>	X	X	X			X	X		X
<i>E. coli</i>	X		X			X	X		X (x2)
<i>L. fermentum</i>				X		X	X		X
<i>L. monocytogenes</i>	X	X	X			X	X		X
<i>P. aeruginosa</i>	X		X			X	X		X
<i>S. enterica</i>	X	X	X			X	X		X
<i>S. aureus</i>	X	X	X			X	X		X
<i>C. neoformans</i>				X (x2)		X	X		X
<i>S. cerevisiae</i>	X	X		X (x2)		X	X		X
NASA <i>E. coli</i> (K12)						X	X	X	X
Pool				X					

Supplementary Table 3. Summary of data generated for the samples discussed in the manuscript. The NASA *E. coli* K12 ONT sequencing included lambda phage and mouse DNA.

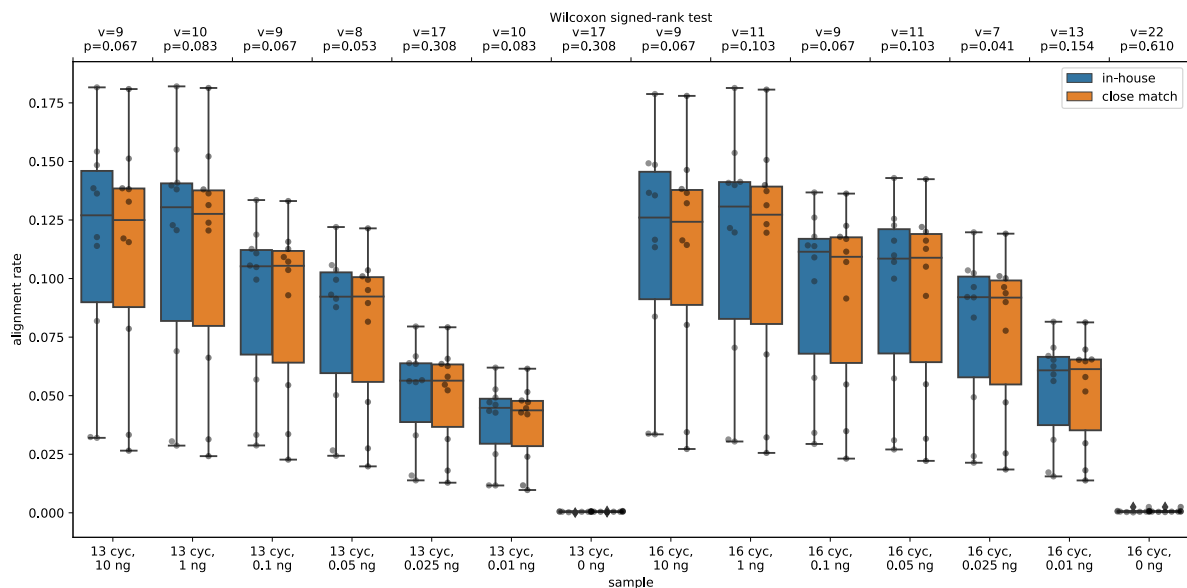
Assembly, contig	Identified match	Pearson's r
<i>B. subtilis</i> , contig 1	<i>B. subtilis</i> subsp. <i>spizizenii</i> W23	0.99866
<i>B. subtilis</i> , contig 2	<i>B. subtilis</i> subsp. <i>spizizenii</i> TU-B-10	0.89454
<i>C. neoformans</i>	<i>C. neoformans</i> var. <i>neoformans</i> B-3501A	0.97109
<i>E. coli</i> , contig 1	<i>E. coli</i> str. K-12 substr. MG1655	0.99977
<i>E. coli</i> , contig 2	<i>E. coli</i> O104:H4 str. 2009EL-2050 plasmid p09EL50	0.97142
<i>E. faecalis</i>	<i>E. faecalis</i> OG1RF	0.99957
<i>L. fermentum</i>	<i>L. fermentum</i> F-6	0.99843
<i>L. monocytogenes</i>	<i>L. monocytogenes</i> SLCC2378	0.99990
<i>P. aeruginosa</i>	<i>P. aeruginosa</i> SCV20265	0.99990
<i>S. aureus</i> , contig 1	<i>S. aureus</i> subsp. <i>aureus</i> ECT-R 2	0.99948
<i>S. aureus</i> , contig 2	<i>S. aureus</i> subsp. <i>aureus</i> TW20 plasmid pTW20_1	0.60116
<i>S. cerevisiae</i>	<i>S. cerevisiae</i> BY4742	0.99955
<i>S. enterica</i> , contig 1	<i>S. enterica</i> subsp. <i>enterica</i> serovar Choleraesuis str. SC-B67	0.99974
<i>S. enterica</i> , contig 2	<i>S. enterica</i> subsp. <i>enterica</i> serovar Choleraesuis str. SC-B67 plasmid pSCV50	0.99122

Supplementary Table 4. Similarity to published genomes. Closely related genomes and plasmids inferred with ANI value Pearson's correlations. Note the lower Pearson's r for the *S. aureus* plasmid, which implies it is significantly different from previously known plasmids in this species.

Supplementary note: the *L. fermentum* genome was assembled using Illumina iSeq100 and nanopore data with a custom pipeline. Error corrected Nanopore and HiSeq reads were assembled using SPAdes¹. Scaffolding of the assembled contigs was done using SSpace². and gap filling was executed using GapFiller³.

Species	Assembly statistics			Related strain statistics		
	BUSCOs, %	GC, %	Length, bp	BUSCOs, %	GC, %	Length, bp
<i>B. subtilis</i>	100.0	44.0	4 090 859	100.0	43.9	4 027 676
<i>C. neoformans</i>	62.1 (58.3)	48.3	27 181 434	96.9	45.6	19 699 782
<i>E. coli</i>	98.0 (93.9)	50.5	4 895 096	98.7 (98.0)	50.7	4 750 926
<i>E. faecalis</i>	97.3 (95.9)	37.6	2 865 554	97.3	37.8	2 739 625
<i>L. fermentum</i>	99.3	52.3	1 905 487	98.0 (97.3)	51.7	2 064 620
<i>L. monocytogenes</i>	100.0 (99.3)	38.0	3 012 754	100.0	38.0	2 941 360
<i>P. aeruginosa</i>	100.0	66.2	6 811 589	100.0	66.3	6 725 183
<i>S. aureus</i>	98.6	32.8	2 785 696	98.6	32.9	2 759 125
<i>S. cerevisiae</i>	96.2 (87.6)	38.1	12 843 354	99.0 (92.1)	38.1	12 147 068
<i>S. enterica</i>	98.6	52.2	4 852 168	98.7 (98.0)	52.2	4 805 258

Supplementary Table 5. Benchmarking Universal Single-Copy Orthologs. Percentages of BUSCOs identified in the genomes, GC content and total assembly lengths . The primary value is the percentage of all found BUSCOs, the value in the parentheses is the percentage of contiguous single-copy BUSCOs (where absent, is equal to the primary value).



Supplementary Figure 7. Illumina sequencing for the pooled microbial reference community aligned to assembled genomes. Bwa mem alignment rates for microbial community data at different dilutions and PCR amplification cycles, generated using the Illumina iSeq100. Reads were aligned genomes for each of the ten strains in the reference community, using either in house assemblies or closely related strains. Alignment to in house assemblies tended to be slightly better, but not significantly so. Boxes span from the first to the third quartile, with medians denoted by a centre line, and whiskers represent the 1.5x interquartile range.

Supplementary References

1. Bankevich, A. *et al.* SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *Journal of computational biology* **19**, 455–477 (2012).
2. Boetzer, M., Henkel, C. V., Jansen, H. J., Butler, D. & Pirovano, W. Scaffolding pre-assembled contigs using SSPACE. *Bioinformatics* **27**, 578–579 (2010).
3. Nadalin, F., Vezzi, F. & Policriti, A. GapFiller: a de novo assembly approach to fill the gap within paired reads. *BMC bioinformatics* **13**, S8 (2012).