

# Complete Microbial Genomes, Epigenomes, and Transcriptomes using Long-Read PacBio® Sequencing

Tyson A. Clark  
Pacific Biosciences, Menlo Park, CA



## Single Molecule, Real-Time (SMRT®) DNA Sequencing

### High Consensus Accuracy

Achieves >99.999% (QV50)  
Lack of systematic bias

### Lack of Sequence Context Bias

GC content  
Low complexity sequence

### Long Sequence Reads

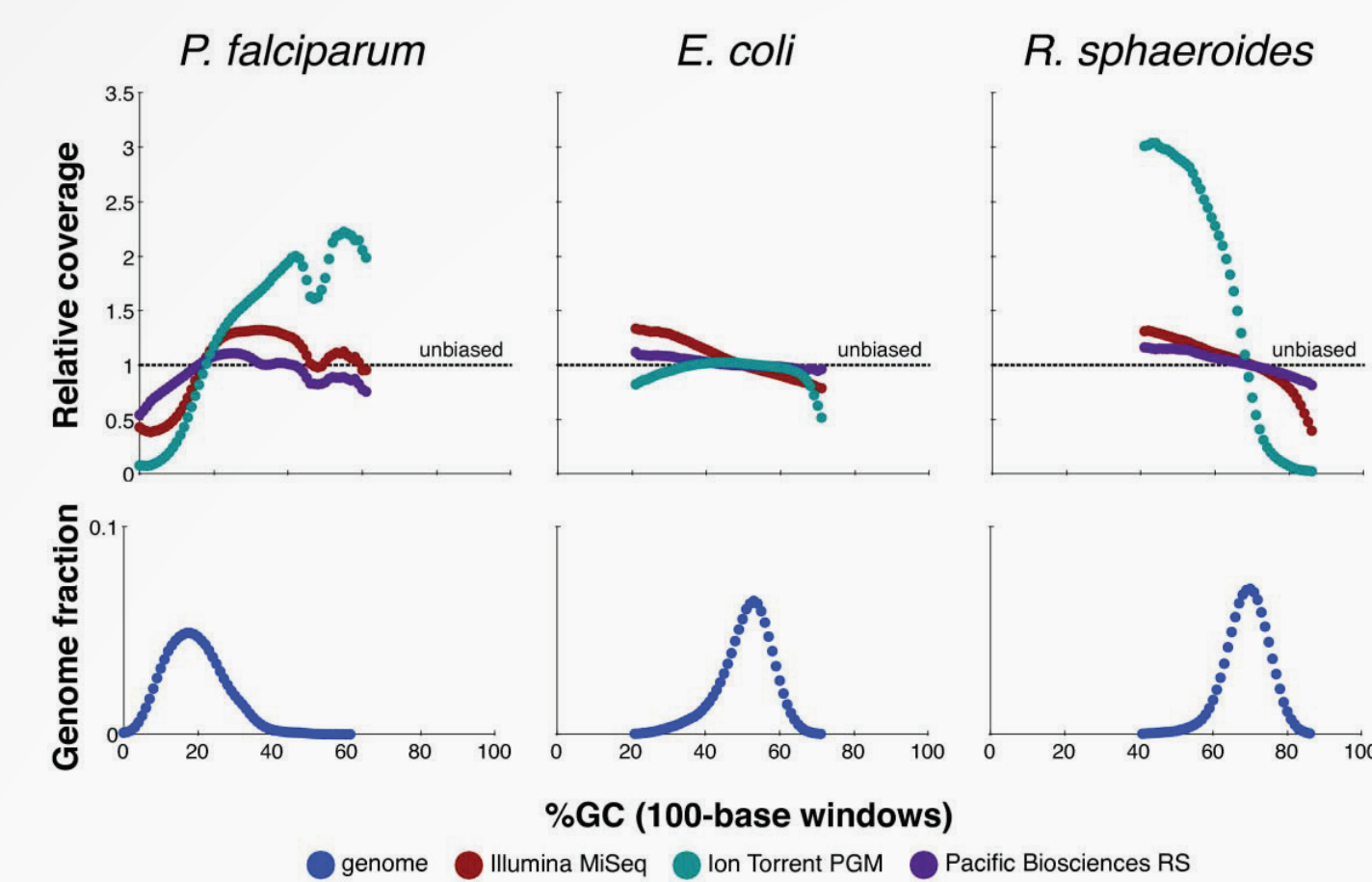
Finished microbial genomes  
Improve large, complex genomes  
Full-length cDNA sequencing  
Long-range haplotype phasing

### Base Modification Detection

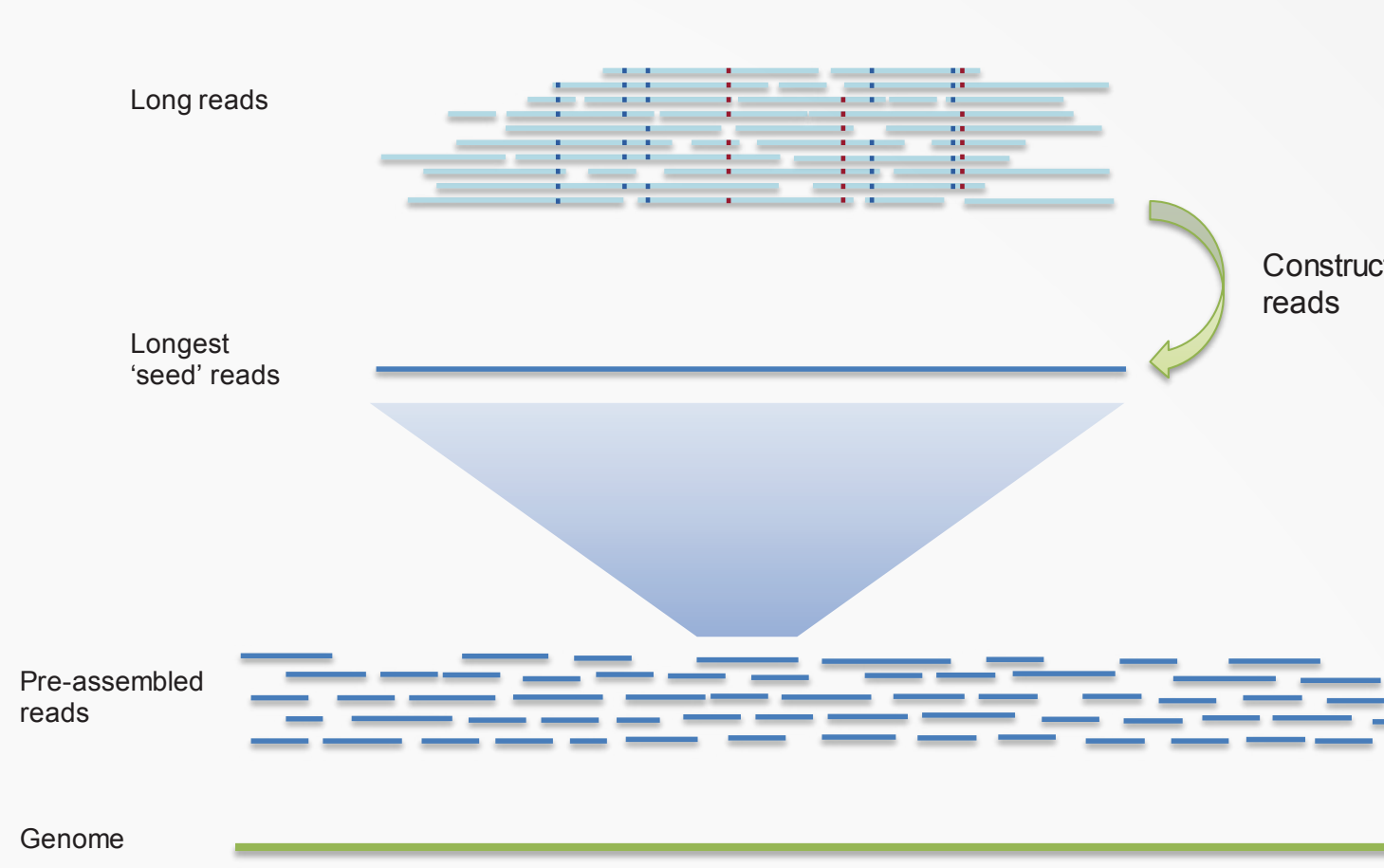
Epigenome characterization

Table 3 Genome assembly continuity and correctness comparison to secondary technologies									
Organism	Assembled with	Assembly bp	Contigs	N50	LAP	Discordant bases	QV		
<i>E. coli</i> K12	MSeq 100k 2x150 bp 300 bp (Illumina/CAI)	4,682,345	139	113,852	-69E+07	28	52.23		
	454 30x	4,589,757	93	117,600	-9.7E+07	17	54.26		
	<b>PacBio 200x</b>	<b>4,653,486</b>	<b>1</b>	<b>4,653,486</b>	<b>-9.6E+07</b>	<b>3</b>	<b>&gt;60</b>		
<i>E. coli</i> O157:H7	MSeq 100k 2x150 bp 500 bp (SPAdes/CAI)	5,433,737	413	133,641	-3.6E+07	62	49.43		
	454 22x + 8x 5 kbp + 10x 10 kbp	5,347,050	409	133,665	-3.7E+07	66	49.09		
	<b>PacBio 200x</b>	<b>5,611,389</b>	<b>9</b>	<b>4,324,437</b>	<b>-3.6E+07</b>	<b>0</b>	<b>&gt;60</b>		
<i>B. subtilis</i>	MSeq 100k 2x150 bp 500 bp (SPAdes/CAI)	2,377,594	83	222,446	-3.3E+07	10	53.76		
	454 30x	2,364,104	66	117,242	-3.3E+07	9	54.20		
	<b>PacBio 200x</b>	<b>2,411,068</b>	<b>1</b>	<b>2,411,068</b>	<b>-3.2E+07</b>	<b>0</b>	<b>&gt;60</b>		
<i>M. thermophilus</i>	MSeq 100k 2x150 bp 500 bp (MAsurCA/CAI)	2,721,965	89	84,094	-3.3E+07	47	47.63		
	<b>PacBio 200x</b>	<b>2,736,037</b>	<b>1</b>	<b>2,736,037</b>	<b>-3.3E+07</b>	<b>0</b>	<b>&gt;60</b>		
<i>F. tularensis</i>	MSeq 100k 2x250 bp 500 bp (SPAdes/CAI)	1,825,374	130	24,065	-1.3E+07	0	>60		
	454 30x	1,655,657	326	7,316	-1.3E+07	28	47.72		
	<b>PacBio 200x</b>	<b>1,877,407</b>	<b>3</b>	<b>573,021</b>	<b>-1.3E+07</b>	<b>0</b>	<b>&gt;60</b>		
<i>S. enterica</i> Newport	MSeq 50x 2x150 bp 500 bp (MAsurCA/CAI)	5,182,289	114	105,780	-2.2E+07	360	41.58		
	454 23x + 2x 10 kbp	5,005,089	172	127,513	-2.2E+07	39	51.08		
	<b>PacBio 200x</b>	<b>5,029,197</b>	<b>2</b>	<b>4,919,684</b>	<b>-2.2E+07</b>	<b>2</b>	<b>&gt;60</b>		

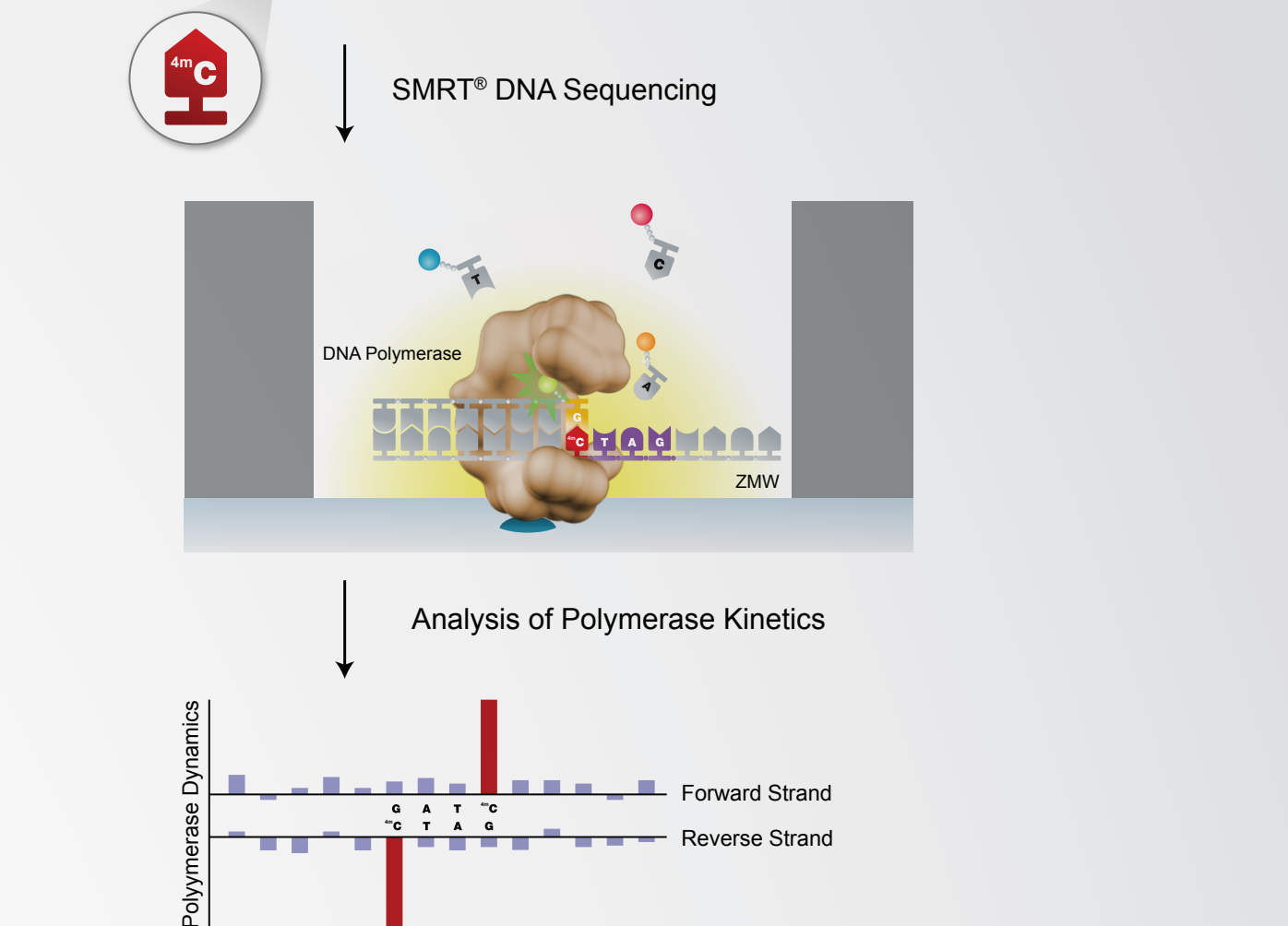
Reducing assembly complexity of microbial genomes with single-molecule sequencing, *Genome Biology*, Koren et al., [2013] **14**:R101 doi:10.1186/gb-2013-14-9-r101



Characterizing and measuring bias in sequence data, *Genome Biology*, Ross et al., [2013] **14**:R51 doi:10.1186/gb-2013-14-5-r51



Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data, *Nature Methods*, Chin et al., [2013] **10**: 563-569

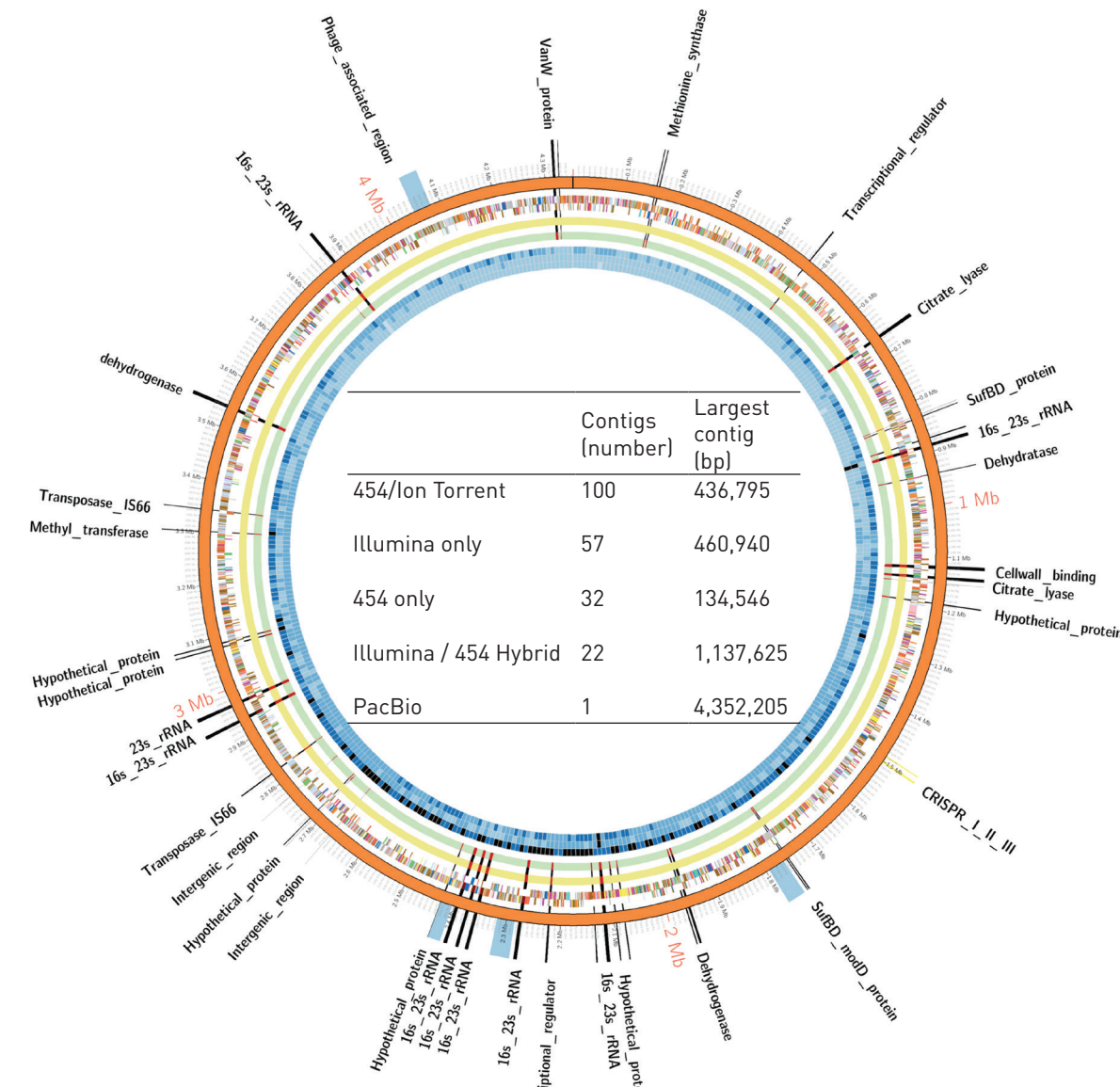


The methylomes of six bacteria, *Nucleic Acids Research*, Murray et al., [2012] **40**: 11450-11462

## Applications for Natural Product Discovery

### Industrial Microbiology

- Clostridium autoethanogenum* is capable of fermenting waste gases (CO, H<sub>2</sub>, CO<sub>2</sub>) into ethanol & commodity chemicals
- Previous draft genome in ~100 contigs
- Closed, high-quality genome generated from only PacBio data without manual finishing

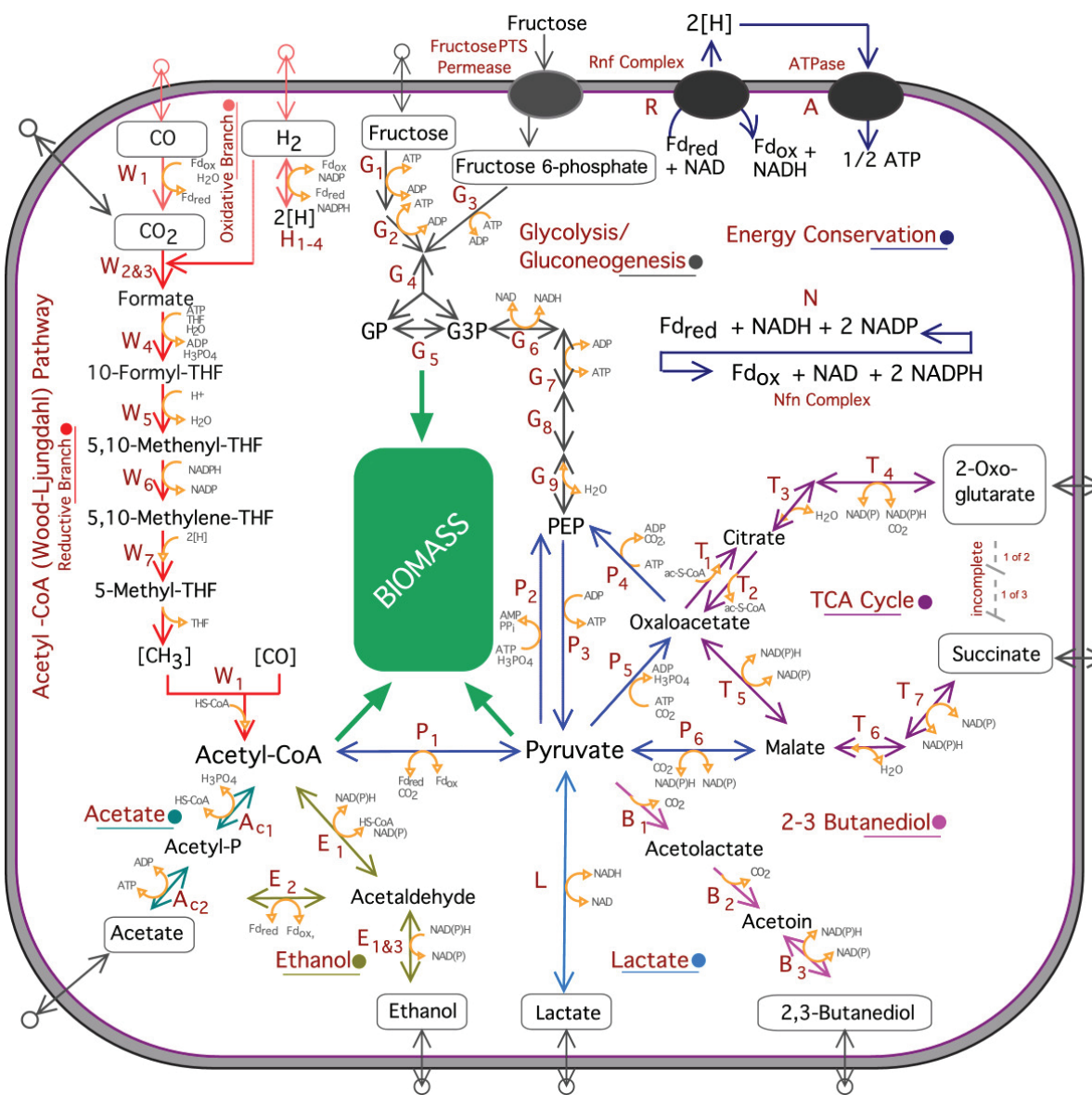


Genes missed by short-read assemblies, but present in PacBio assembly:

Table 3: Regions of low sequence-coverage

Locus tag	Start*	End*	Product description	PacBio coverage (x)	454 Coverage (x)	Illumina coverage (x)	454 Hybrid contig coverage	Draft assembly contig coverage
CAETHD_0602	663332	664234	Citrate lyase, beta subunit	111	29	45	None	None
CAETHD_0603	664234	664530	Citrate lyase, alpha subunit	107	29	43	None	None
CAETHD_0604	664530	664597	Citrate lyase, gamma subunit	109	23	43	None	Partial
CAETHD_0605	664597	664604	Malic protein NAD-binding protein	101	27	49	None	None
Intergenic	827340	827250	NA	106	30	53	None	None
rRNA	885055	887942	23S rRNA	87	77	147	None	None
rRNA	888206	889703	16S rRNA	102	54	145	None	None
CAETHD_1038	1114305	1114321	Cell wall-binding repeat 2-containing protein	127	49	Partial	None	None
CAETHD_1052	1136476	1138017	Citrate lyase, alpha subunit	107	22	53	Partial	None
CAETHD_1055	1139370	1140533	Malic protein NAD-binding protein	107	27	51	Partial	Partial
rRNA	2114155	2117942	23S rRNA	122	81	141	None	None
rRNA	2117334	2118031	16S rRNA	118	44	128	None	None
CAETHD_2076	2220149	2221596	Sigma54 specific transcriptional regulator, Fis family	122	32	85	Partial	Partial
CAETHD_2077	2221458	2221885	Transcriptional regulator, Fis family	126	21	92	Partial	None
CAETHD_2078	2222014	2222994	Putative sigma54 specific transcriptional regulator	135	30	77	Partial	Partial
rRNA	2271738	2272235	16S rRNA	145	10	26	None	None
rRNA	2273527	2274414	23S rRNA	158	10	26	None	None
rRNA	2355334	2356831	16S rRNA	145	11	24	None	None
rRNA	2357123	2360010	23S rRNA	136	13	23	None	None
rRNA	2372238	2372735	16S rRNA	128	13	21	None	None
rRNA	2374027	2374914	23S rRNA	126	14	19	None	None
rRNA	2397202	2397499	16S rRNA	134	12	20	None	None
rRNA	2397496	2397483	23S rRNA	142	11	21	None	None
CAETHD_2421	2823723	2824238	Transposase IS66	127	30	52	Partial	Partial
rRNA	2935186	2936483	16S rRNA	127	14	27	None	None
rRNA	2937473	2937945	rRNA, 16S	125	19	51	None	None
rRNA	2937053	2937126	rRNA, 16S	125	26	58	None	None
rRNA	2937443	2940330	23S rRNA	117	14	28	None	None
rRNA	2949792	2949889	16S rRNA	126	11	20	None	None
rRNA	2949779	2949851	rRNA, 16S	132	20	58	None	None
rRNA	2948859	2948932	rRNA, 16S	131	23	70	None	None
rRNA	2949722	2972109	23S rRNA	128	10	19	None	None
rRNA	3872016	3873511	16S rRNA	98	10	18	None	None
rRNA	3873937	3874824	23S rRNA	107	14	21	None	None

### Full metabolic reconstruction possible

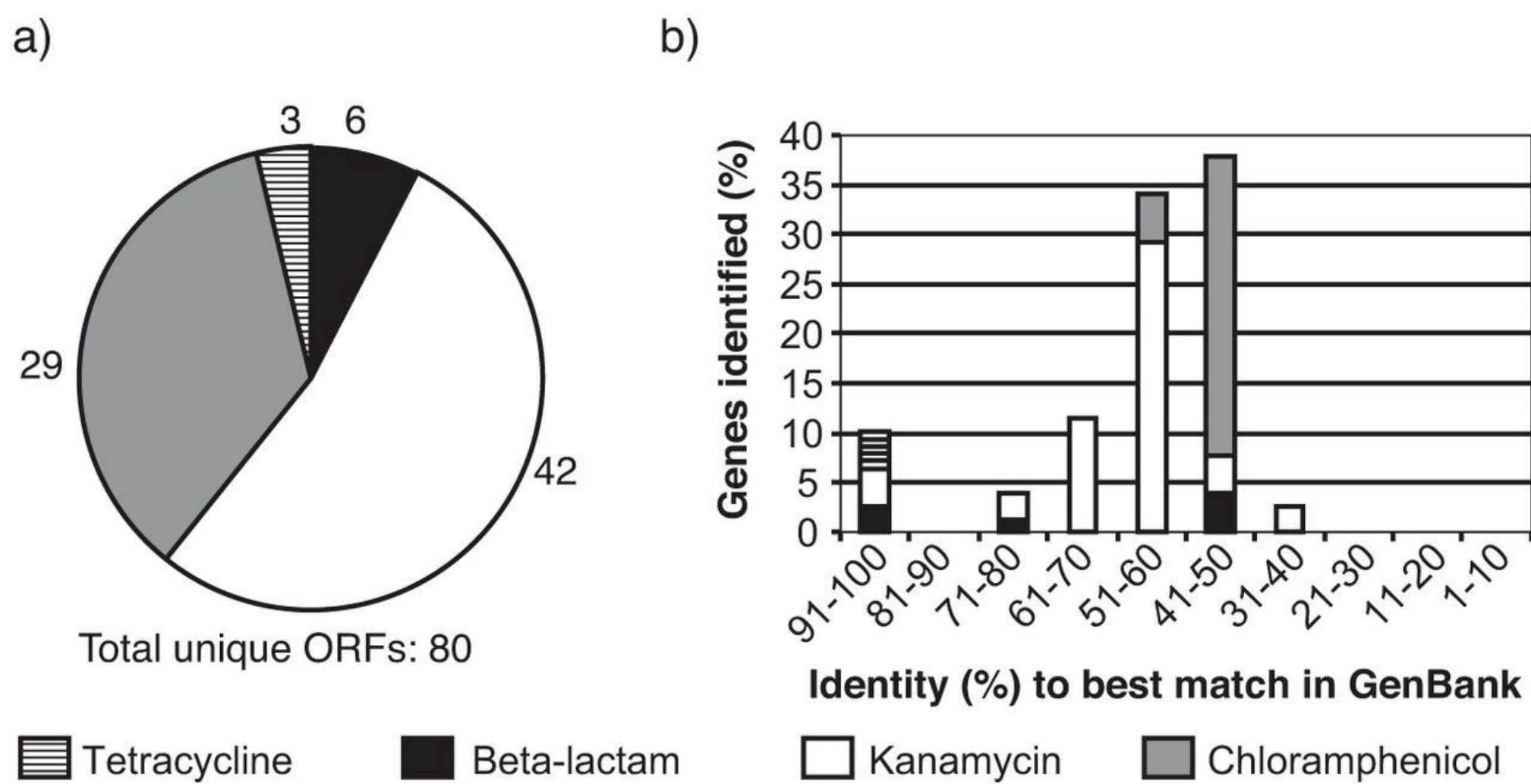


Brown et al. [2014] Comparison of single-molecule sequencing and hybrid approaches for finishing the genome of *Clostridium autoethanogenum* and analysis of CRISPR systems in industrial relevant Clostridia. *Biotechnology for Biofuels* **7**:40

### Livestock-Microbe Interactions

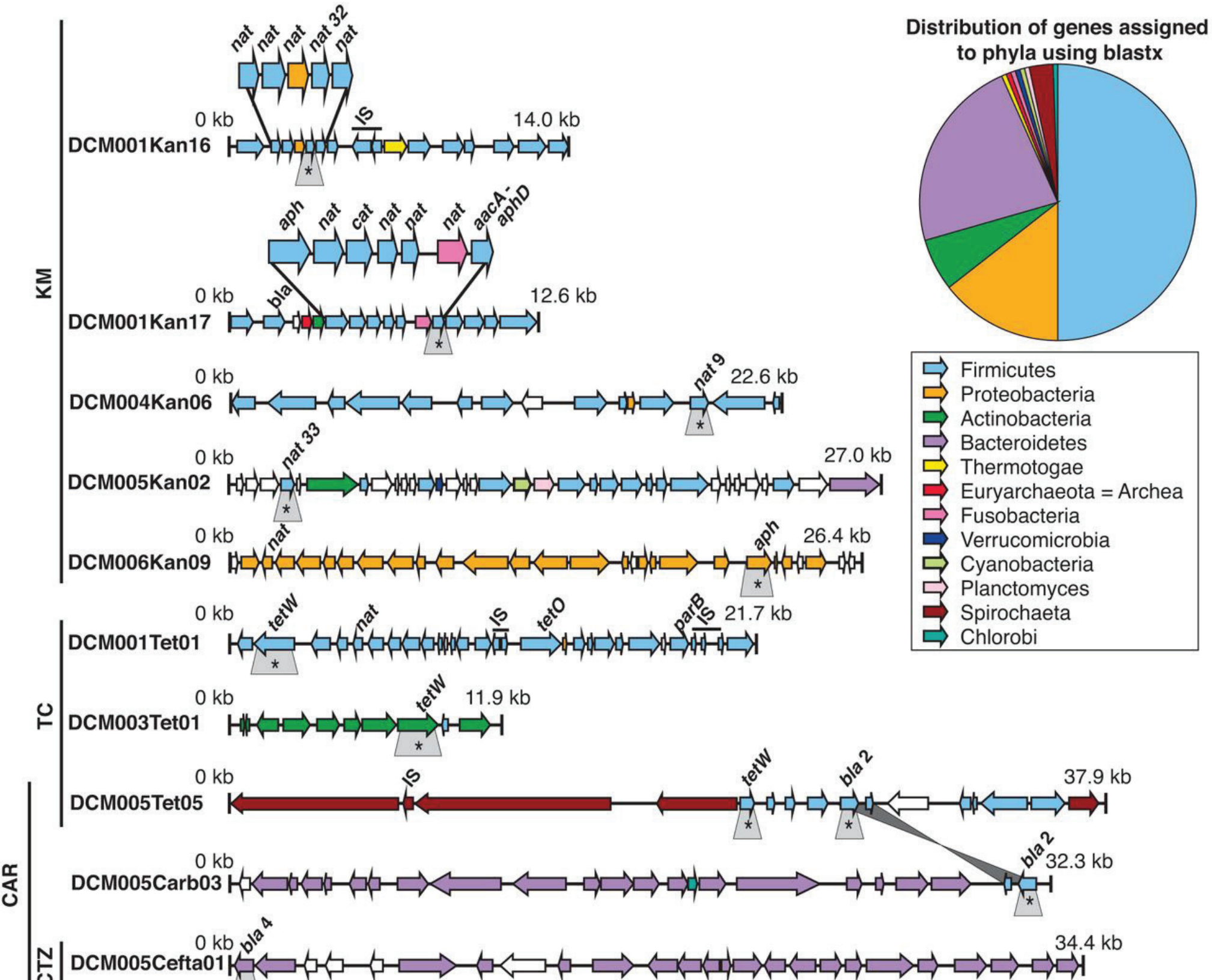
- Microbiomes of farm animals are reservoirs of antibiotic resistance genes
- Discovered a novel clade of chloramphenicol acetyltransferases
- Significantly extends the roster of functional antibiotic resistance genes found in animal gut bacteria

Antibiotic resistance genes from dairy cow manure.



Antibiotic resistance genes from dairy cow manure. [a] Distribution of 80 unique antibiotic resistance (AR) genes among four classes of antibiotics. [b] Distribution of similarity of 80 unique AR genes from manure compared to homologues in GenBank.

Organization of genes on ten metagenomic fosmid clones conferring resistance to kanamycin (KM), tetracycline (TC), carbencillin (CAR), or ceftazidime (CTZ) assessed using Pacific Biosciences RS sequencing technology.

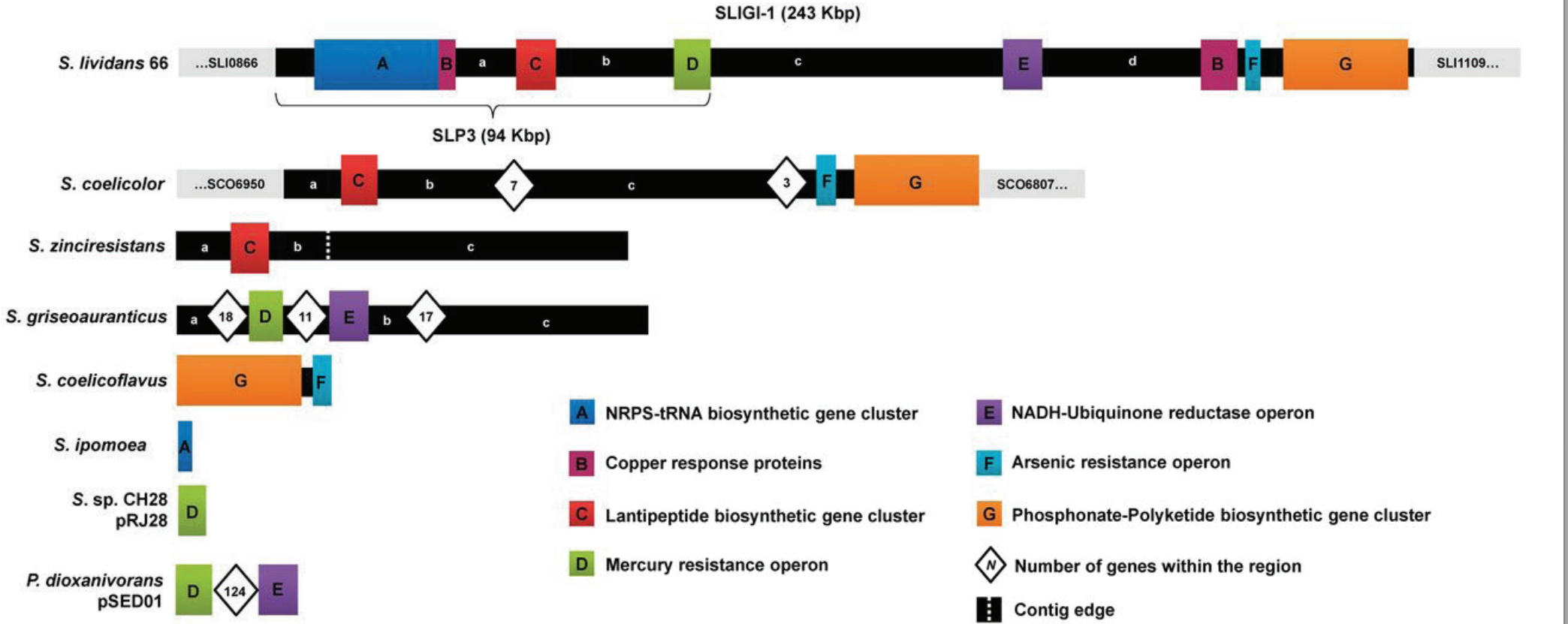


Organization of genes on ten metagenomic fosmid clones conferring resistance to kanamycin (KM), tetracycline (TC), carbencillin (CAR), or ceftazidime (CTZ) assessed using Pacific Biosciences RS sequencing technology. Genes are as follows: *nat*, N-acetyltransferase; *aph*, aminoglycoside phosphotransferase; *aacA-aphD*, bifunctional aminoglycoside-modifying enzyme; *tetW* and *tetO*, tetracycline resistance genes; *cat*, chloramphenicol acetyltransferase; *bla*, beta-lactamase. Genes with annotations related to mobile genetic elements are labeled as IS (insertion sequence elements/transposases) or *oripB* (plasmid-partitioning protein). Identical genes are connected with gray lines. Antibiotic resistance genes for which we demonstrated a function by subcloning are indicated with gray boxes and asterisks. All other labels were assigned based on sequence homology.

Wichmann et al. [2014] Diverse antibiotic resistance genes in dairy cow manure. *mBio* **5**:e01017-13

### Bioremediation

- High-quality genome sequence of a widely used actinomycete model organism, *Streptomyces lividans* 66
- High-quality genomes aid in identifying genetic diversity amongst closely related strains that otherwise would have been lost
- Resolves large mobile genomic island related to metal metabolism, which harbors the elusive plasmid SLP3
- SLP3 encodes enzymes that direct the synthesis of natural products that remain to be discovered, even in such well-studied strains as *S. lividans* 66



Cruz-Morales et al. [2013] The genome sequence of *Streptomyces lividans* 66 reveals a novel tRNA-dependent peptide biosynthetic system within a metal-related genomic island. *Genome Biology & Evolution* **5**:1165-1175

### Other Examples

#### Biocommodities:

Identification of Restriction-Modification Systems of *Bifidobacterium animalis* subsp. *lactis* CNCM I-2494 by SMRT Sequencing and Associated Methylation Analysis <http://dx.plos.org/10.1371/journal.pone.0094875.g006>  
Complete Genome Sequence of *Enterococcus mundtii* QJ 25, an Efficient L-Lysine-Lactic Acid-Producing Bacterium <http://dnaresearch.oxfordjournals.org/content/early/2014/03/10/dnaresearch.dsu003.full>

#### Biofuels:

The genome of the anaerobic fungus *Orpinomyces* sp. strain C1A reveals the unique evolutionary history of a remarkable plant biomass degrader <http://aem.asm.org/content/early/2013/05/20/AEM.00821-13.full.pdf+html> (fungus)  
Localized electron transfer rates and microelectrode-based enrichment of microbial communities within a phototrophic microbial mat <http://www.ncbi.nlm.nih.gov/pmc/articles/PMC3902354/> (16S)  
Long-term operation of microbial electrosynthesis systems improves acetate production by autotrophic microorganisms <http://pubs.acs.org/doi/abs/10.1021/es400341b> (16S)  
Capturing single cell genomes of active polysaccharide degraders: an unexpected contribution of Verrucomicrobia <http://www.plosone.org/article/info%3Adoi%2F10.1371%2Fjournal.pone.0035314>

#### Bioremediation:

Exploring the roles of DNA methylation in the metal-reducing bacterium *Shewanella oneidensis* MR-1 <http://jb.asm.org/content/195/21/4966.short>  
A crowd sourced funded project to sequence a fern. <https://experiment.com/projects/azolla-a-little-fern-with-massive-green-potential/community>  
Genome sequence of *Candidatus Microthrix parvicella* Bio1-7, a long-chain-fatty-acid-accumulating filamentous actinobacterium from a biological wastewater treatment plant <http://jb.asm.org/content/194/23/6670.long>

#### New bacteria with potential commercial applications:

Single cell genomic study of *Dehalococcoides* species from deep-sea sediments of the Peruvian Margin <http://www.wafergen.com/wp-content/uploads/2014/03/smej201424a.pdf>  
Genome Sequence of *Bacillus pumilus* MTCC B6033 <http://genome.asm.org/content/2/2/e00327-14.full>

#### Antibiotic discovery:

De Novo Assembly of the *Streptomyces* sp. Strain Mg1 Genome Using PacBio Single-Molecule Sequencing. <http://genome.asm.org/content/1/4/e00535-13.full.pdf+html>  
Value of a newly sequenced bacterial genome. <http://www.wjgnet.com/1949-8454/full/v5/i2/161.htm>

#### Livestock/plant microbiome interactions:

Genome Sequence and Methylation of Soil Bacterium *Gemmatirosa kalamazonensis* KBS708T, a Member of the Rarely Cultivated *Gemmatimonadetes* Phylum <http://genome.asm.org/content/2/2/e00226-14.full>  
Genome Sequence of *Pseudomonas* sp. Strain P482, a Tomato Rhizosphere Isolate with Broad-Spectrum Antimicrobial Activity. <http://www.ncbi.nlm.nih.gov/pubmed/24970823>

pacb.com

See over 200 customer publications at [www.pacb.com/publications](http://www.pacb.com/publications)

For Research Use Only. Not for use in diagnostic procedures. Pacific Biosciences, the Pacific Biosciences logo, PacBio, SMRT, SMRTbell and Iso-Seq are trademarks of Pacific Biosciences of California, Inc. All other trademarks are the property of their respective owners. © 2015 Pacific Biosciences of California, Inc. All rights reserved.