

SUPPLEMENTARY TABLE AND FIGURES LEGENDS

Supplementary Figure S1. Comparative analysis of the CasCollect, complete assembly, and PRICE pipelines for isolate genomic datasets. (A) Schematic of the four steps for the CasCollect and complete assembly pipelines: files (blue) for trimming and stripping qualities; seed (orange) for the seed generation, expansion, and reassigning qualities; assembly (grey) for contig building; and annotate (yellow) for detecting CRISPR arrays and *cas* genes GFF3 output file. (B) Timings for each pipeline for each isolate genome and steps. (C) Percentage breakdown for time required for each step. (D) Number of contigs generated by each pipeline with (E) the length distribution for CasCollect, complete assembly, and PRICE. Colors for each step in (A) are used for the charts in (B) and (C).

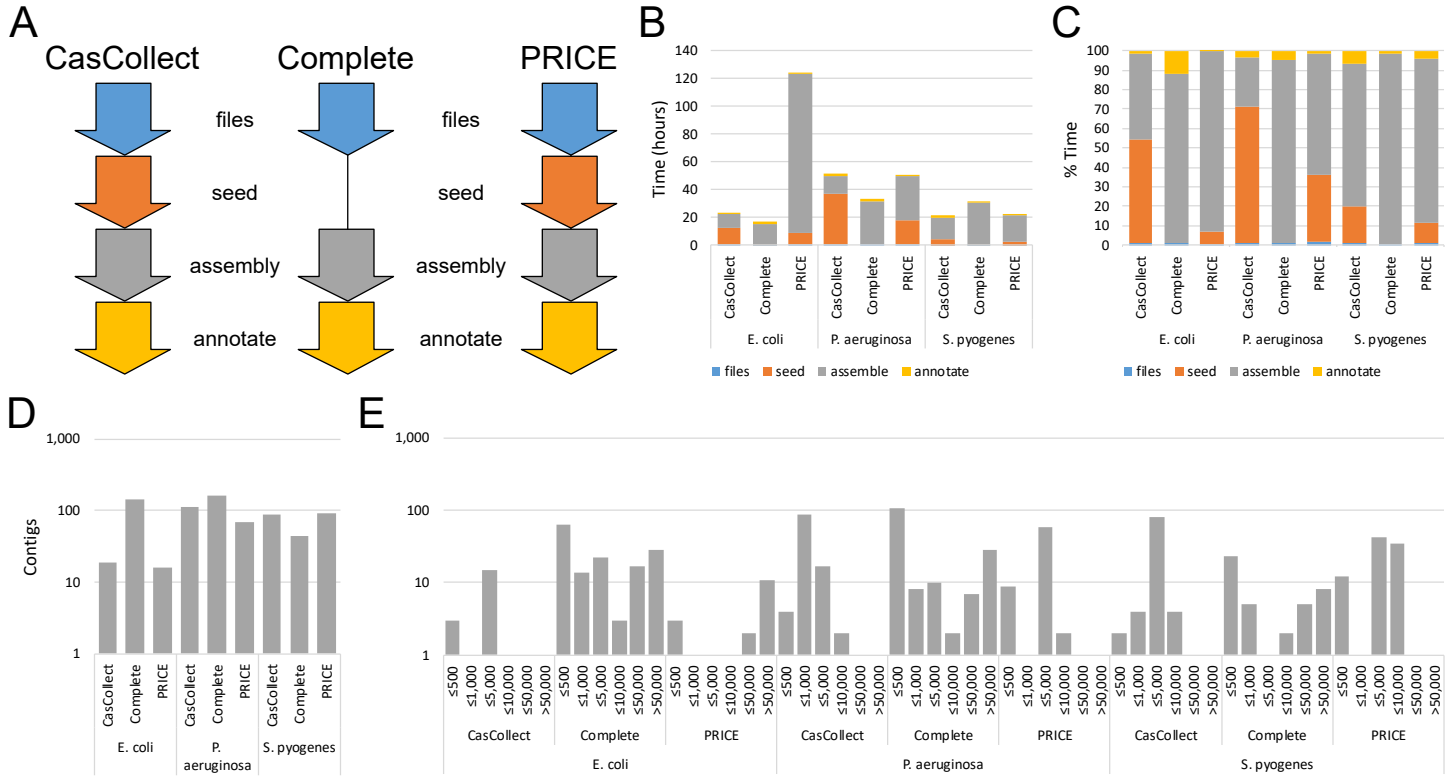
Supplementary Figure S2. Comparative analysis of the CasCollect and PRICE pipelines for metagenomic datasets. (A) Timings for each pipeline for each isolate genome and steps. (B) Percentage breakdown for time required for each step. The timings for PRICE were stopped after 200 cycles that corresponded to a 3-5 fold time increase compared to CasCollect. The timing for PRICE was insufficient to assemble the majority of *cas* operons.

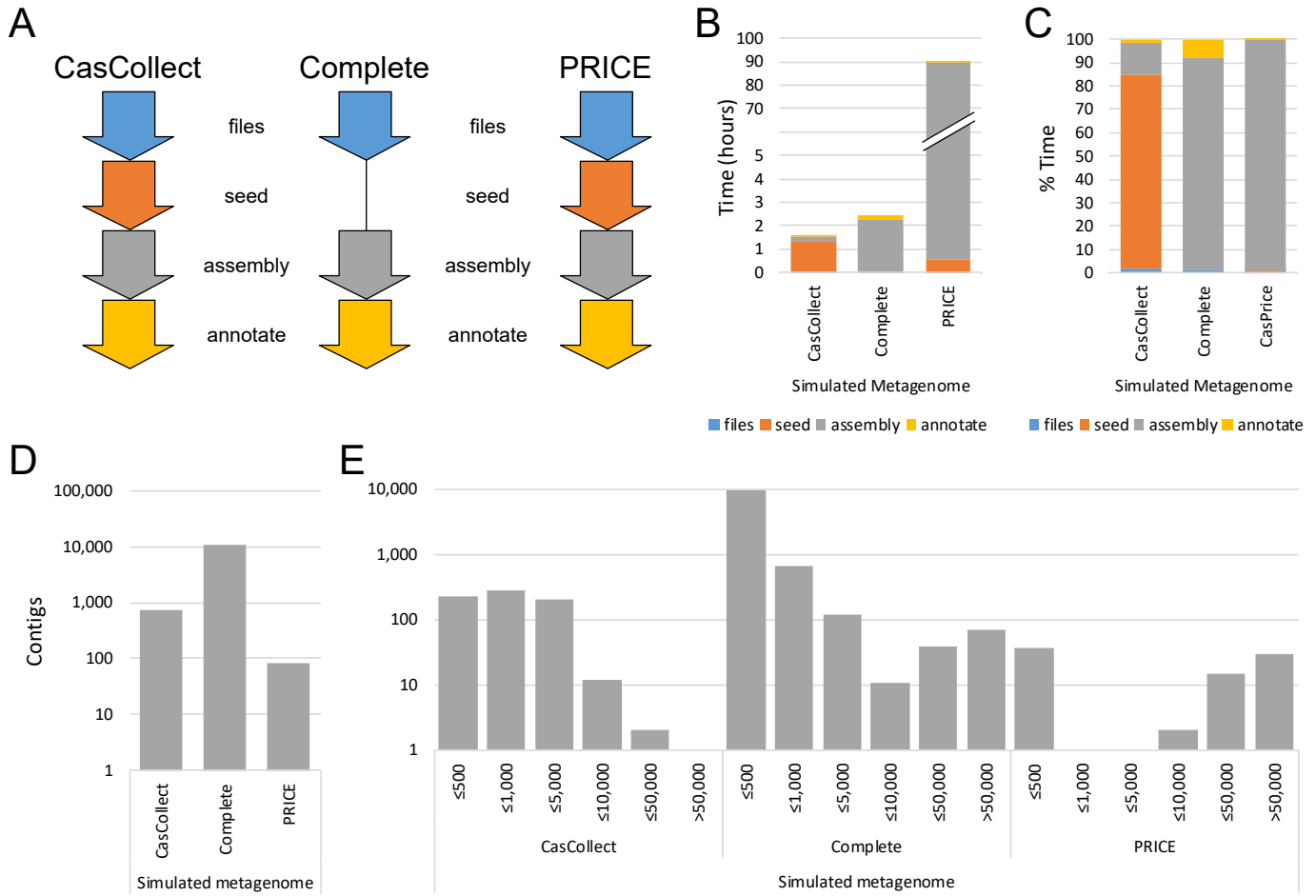
Supplementary Figure S3. Comparative analysis of the CasCollect and complete assembly pipelines for simulated metagenomic data. (A) Schematic of the four steps for the CasCollect and complete assembly pipelines: Files (blue) for trimming, Seed (orange) for the seed generation by protein mode, Assembly (grey) for contig building, and Annotate (yellow) for the general feature format file of CRISPR/Cas genes. (B) Timings for each pipeline with each step in the pipelines. (C) Percentage breakdown for time required for each step. (D) Number of contig generated by each pipeline with the length distribution for CasCollect (E) and complete assembly (F) characterized. Colors for each step in (A) used for the charts in (B) and (C).

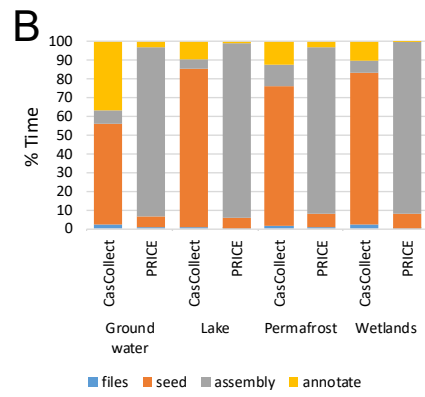
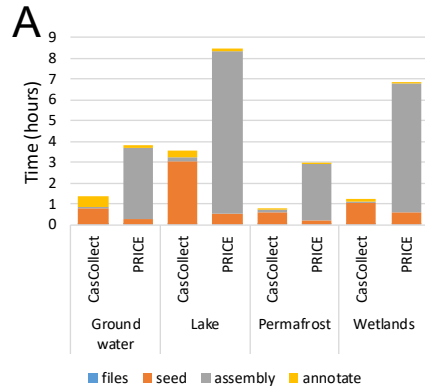
Supplementary Figure S4. Discrepancy in spacer number within CRISPR arrays. (A) Schematic of the CRISPR array from AR110 (SRR3112345) with CasCollect targeted and complete assembly. (B) The sequence of each direct repeat (uppercase) and spacer (lowercase). (C) Read coverage from mapping the raw sequencing reads onto the CRISPR arrays from targeted and complete assembly. Coverage across the targeted assembled was consistent, while the 180 bp region comprising three duplicates of the direct repeat and spacer in the complete assembly had low coverage.

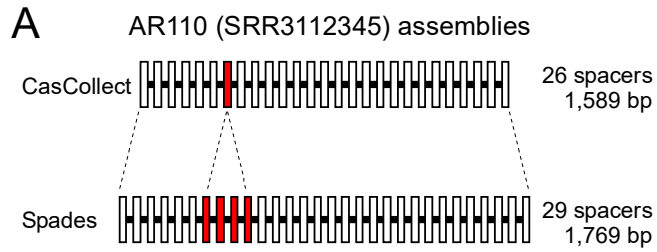
Supplementary Table S1. List of species with CRISPR repeats available for the DNA mode of CasCollect.

Supplementary Table S2. Identification of *Pseudomonas aeruginosa* CRISPR arrays lacking *cas* operons with CasCollect.









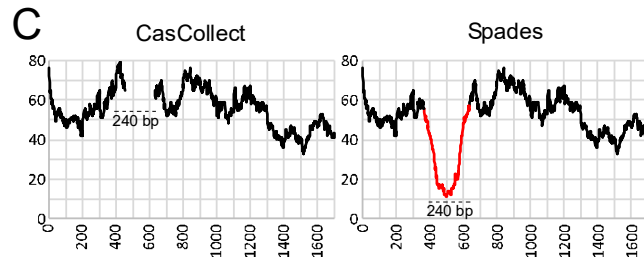
B

1 28 60

```

GTTCACTGCCGTGTAGGCAGCTAAGAAA tggctgatcaggctccagaa cggatcgtagac
GTTCACTGCCGTGTAGGCAGCTAAGAAA ttgatatgccggtagaacgtcggcgagacat
GTTCACTGCCGTGTAGGCAGCTAAGAAA tcgaacgctgctgagcgcgaacgcatagatg
GTTCACTGCCGTGTAGGCAGCTAAGAAA agaccgagggcgtcgaaaactcgatgatc
GTTCACTGCCGTGTAGGCAGCTAAGAAA ttcgacggccacgctcagcccggcccaggcc
GTTCACTGCCGTGTAGGCAGCTAAGAAA ttattgaaatccatcagcggctcgcactgtctc
GTTCACTGCCGTGTAGGCAGCTAAGAAA tgatcaccacaagcgtgcgta tcgcccattcc
GTTCACTGCCGTGTAGGCAGCTAAGAAA tgatcaccacaagcgtgcgta tcgcccattcc
GTTCACTGCCGTGTAGGCAGCTAAGAAA tgatcaccacaagcgtgcgta tcgcccattcc
GTTCACTGCCGTGTAGGCAGCTAAGAAA tgatcaccacaagcgtgcgta tcgcccattcc
GTTCACTGCCGTGTAGGCAGCTAAGAAA tgggtgtccaacatcgacgggtcgaactgtctc
GTTCACTGCCGTGTAGGCAGCTAAGAAA tggctgata tggccgga tcatagcgcgacctac
GTTCACTGCCGTGTAGGCAGCTAAGAAA ttoggtacttctgaaccatacgtcgcgcgata
GTTCACTGCCGTGTAGGCAGCTAAGAAA agtcatcgatgaaacgacgagccggtcagtgcc
GTTCACTGCCGTGTAGGCAGCTAAGAAA agaagctggagcagcggctggcggcaattcgt
GTTCACTGCCGTGTAGGCAGCTAAGAAA ccggacgttcacgctggtggtgagacca tccg
GTTCACTGCCGTGTAGGCAGCTAAGAAA tggctgtcgtcgtcgtcgtcgtcgtcgtcgtat
GTTCACTGCCGTGTAGGCAGCTAAGAAA aggtggtccagagcgggtcgaacggcacggtc
GTTCACTGCCGTGTAGGCAGCTAAGAAA gaaccgcccgttcatctgtgaaaggccatcgtc
GTTCACTGCCGTGTAGGCAGCTAAGAAA acatcagcgcgcggtagccgatgocgatatac
GTTCACTGCCGTGTAGGCAGCTAAGAAA acca tcccggccacgggttgcggacacctg
GTTCACTGTCGTGTAGGCAGCTAAGAAA gtccatccgggtagg tcaactcaccgtcgtat
GTTCACTGCCGTGTAGGCAGCTAAGAAA tggagagtgaaccgctcaagaccgagcccgag
GTTCACTGCCGTGTAGGCAGCTAAGAAA tga tgcgggaca tgggacgtttcgcgggaacc
GTTCACTGCCGTGTAGGCAGCTAAGAAA ccggacgcccetaatc tggagggc tccctggca
GTTCACTGCCGTGTAGGCAGCTAAGAAA ttccccgcggagcatagcagggatattctgt
GTTCACTGCCGTGTAGGCAGCTAAGAAA tacaagcagatcggcagctgagcggcaaggga
GTTCACTGCCGTGTAGGCAGCTAAGAAA tcgtactggtcagat tccgattcggaaagcc
GTTCACTGCCGTGTAGGCAGCTAAGAAA ttogacgcgcgtagggttgcgcgcatcgcctc
GTTCACTGCCGTGTAGGCAGCTAAGAAA

```



CRISPR array repeats	
<i>Aeropyrum camini</i>	<i>Microclunatus phosphovorus</i> NM-1
<i>Alicyclobacillus acidocaldarius</i> subsp <i>acidocaldarius</i> Tc-4-1	<i>Mycobacterium canettii</i> CIPT 140010059
<i>Azotobacter vinelandii</i> CA	<i>Mycobacterium tuberculosis</i> RGTB327
<i>Bdellovibrio exovorus</i> JSS	<i>Mycoplasma gallisepticum</i> S6
<i>Bibersteinia trehalosi</i> USDA-ARS-USMARC-192	<i>Myxococcus stipitatus</i> DSM 14675
<i>Bifidobacterium longum</i> subsp <i>longum</i> KACC 91563	<i>Neisseria meningitidis</i> Z2491
<i>Brachyspira hyodysenteriae</i> WA1	<i>Nocardia cyriacigeorgica</i> GUH-2
<i>Burkholderia</i> sp Y123	<i>Nocardiopsis alba</i> ATCC BAA-2165
<i>Calothrix</i> sp PCC 6303	<i>Nostoc</i> sp PCC 7120
<i>Campylobacter</i> sp 03-427	<i>Pasteurella multocida</i> 36950
<i>Candidatus Arthromitus</i> sp SFB-mouse-Japan	<i>Pectobacterium carotovorum</i> subsp <i>carotovorum</i> PCC21
<i>Chamaesiphon minutus</i> PCC 6605	<i>Persephonella marina</i> EX-H1
<i>Clostridium</i> sp SY8519	<i>Phycisphaera mikurensis</i> NBRC 102666
<i>Corynebacterium terpenotabidum</i> Y-11	<i>Pleurocapsa</i> sp PCC 7327
<i>Crinalium epipsammum</i> PCC 9333	<i>Porphyromonas gingivalis</i> TDC60
<i>Cronobacter sakazakii</i> ES15	<i>Propionibacterium acidipropionici</i> ATCC 4875
<i>Cycloclasticus zancles</i> 7-ME	<i>Providencia stuartii</i> MRSN 2154
<i>Dehalococcoides mccartyi</i> DCMB5	<i>Pseudomonas aeruginosa</i> DK2
<i>Desulfocapsa sulfexigens</i> DSM 10523	<i>Pyrobaculum</i> sp 1860
<i>Desulfovibrio hydrothermalis</i> AM13	<i>Pyrococcus</i> sp ST04
<i>Enterococcus faecalis</i> D32	<i>Ralstonia solanacearum</i> Po82
<i>Escherichia coli</i> str K-12 substr MDS42 DNA	<i>Rhodospirillum photometricum</i> DSM 122
<i>Feridicoccus fontis</i> Kam940	<i>Rhodospirillum rubrum</i> F11
<i>Flavobacterium branchiophilum</i>	<i>Riemerella anatipestifer</i> RA-CH-1
<i>Frankia</i> sp Ccl3	<i>Saccharothrix espanaensis</i> DSM 44229
<i>Gallibacterium anatis</i> UMN179	<i>Salinispora arenicola</i> CNS-205
<i>Geitlerinema</i> sp PCC 7407	<i>Salmonella enterica</i> subsp <i>enterica</i> serovar Typhimurium str 798
<i>Geobacillus thermoleovorans</i> CCB US3 UF5	<i>Sorangium cellulosum</i> So0157-2
<i>Gloeobacter</i> sp JS	<i>Spiroplasma apis</i> B31
<i>Gloeocapsa</i> sp PCC 7428	<i>Stanieria cyanosphaera</i> PCC 7437
<i>Gordonia polyisoprenivorans</i> VH2	<i>Staphylococcus aureus</i> 08BA02176
<i>Granulibacter bethesdensis</i> CGDNIH1	<i>Streptococcus equi</i> subsp <i>zooepidemicus</i> ATCC 35246
<i>Halorhabdus tiamatea</i> SARL4B	<i>Streptococcus lutetiensis</i> 033
<i>Helicobacter cinaedi</i> PAGU611	<i>Streptococcus macedonicus</i> ACA-DC 198
<i>Hydrogenobaculum</i> sp 3684	<i>Streptococcus salivarius</i> CCHSS3
<i>Hyphomicrobium nitratorans</i> NL23	<i>Streptococcus suis</i> YB51
<i>Kibdelosporangium phytohabitans</i> strain KLBMP1111	<i>Streptococcus thermophilus</i> JIM 8232
<i>Klebsiella oxytoca</i> KCTC 1686	<i>Streptomyces hygrosopicus</i> subsp <i>jinggansensis</i> 5008
<i>Lactobacillus brevis</i> KB290 DNA	<i>Sulfobacillus acidophilus</i> TPY
<i>Lactobacillus ruminis</i> ATCC 27782	<i>Sulfolobus acidocaldarius</i> N8
<i>Leptothrix cholodnii</i> SP-6	<i>Synechococcus</i> sp PCC 7502
<i>Leuconostoc gelidum</i> JB7	<i>Synechocystis</i> sp PCC 6803
<i>Listeria ivanovii</i> subsp <i>londoniensis</i> strain WSLC 30167	<i>Taylorella asinigenitalis</i> MCE3
<i>Listeria monocytogenes</i> M7	<i>Tepidanaerobacter acetatoxydans</i> Re1
<i>Mannheimia haemolytica</i> USDA-ARS-SAM-185	<i>Thermacetogenium phaeum</i> DSM 12270
<i>Megasphaera elsdenii</i> DSM 20460	<i>Thermoanaerobacterium saccharolyticum</i> JW-SL-YS485
<i>Melioribacter roseus</i> P3M	<i>Thermococcus</i> sp 4557
<i>Melissococcus plutonius</i> DAT561	<i>Thermogladus</i> sp 1633
<i>Metallosphaera cuprina</i> Ar-4	<i>Thermotoga elfii</i> NBRC 107921
<i>Methanocella</i> sp HZ254	<i>Thermus</i> sp CCB US3 UF1
<i>Methanoculleus bourgensis</i> MS2T	<i>Thioalkalivibrio nitratireducens</i> DSM 14787
<i>Methanosaeta harundinacea</i> 6Ac	<i>Tistrella mobilis</i> KA081020-06
<i>Microcoleus</i> sp PCC 7113	

Strain	SRA	Reads	Type	CRISPR		Cas Genes	
				Arrays	Spacers	#	Operon
AR090	SRR3112335	2,126,644	n/a	0	0	0	n/a
AR092	SRR3112336	1,225,103	n/a	0	0	0	n/a
AR094	SRR3112337	1,657,533	n/a	0	0	0	n/a
AR100	SRR3112340	2,087,105	n/a	0	0	0	n/a
AR105	SRR3112342	1,989,806	n/a	0	0	0	n/a
AR231	SRR4417542	478,617	n/a	1	8	0	n/a
AR232	SRR4417553	375,904	n/a	0	0	0	n/a
AR233	SRR4417558	615,859	n/a	0	0	0	n/a
AR234	SRR4417559	477,110	n/a	0	0	0	n/a
AR236	SRR4417561	488,238	n/a	0	0	0	n/a
AR237	SRR4417562	1,175,855	n/a	0	0	0	n/a
AR239	SRR4417532	1,008,367	n/a	1	4	0	n/a
AR240	SRR5122331	1,647,468	n/a	0	0	0	n/a
AR244	SRR4417533	856,508	n/a	0	0	0	n/a
AR247	SRR4417534	856,694	n/a	0	0	0	n/a
AR251	SRR4417535	882,484	n/a	0	0	0	n/a
AR252	SRR5122325	1,141,358	n/a	0	0	0	n/a
AR253	SRR4417536	884,811	n/a	0	0	0	n/a
AR257	SRR4417540	951,684	n/a	0	0	0	n/a
AR258	SRR4417541	975,106	n/a	0	0	0	n/a
AR259	SRR4417543	881,366	n/a	0	0	0	n/a
AR260	SRR4417544	953,363	n/a	1	12	0	n/a
AR263	SRR4417547	918,300	n/a	0	0	0	n/a
AR264	SRR4417548	877,917	n/a	1	10	0	n/a
AR265	SRR4417549	766,911	n/a	1	13	0	n/a
AR266	SRR4417550	857,564	n/a	0	0	0	n/a
AR269	SRR4417554	1,011,282	n/a	0	0	0	n/a
AR270	SRR4417555	1,012,411	n/a	0	0	0	n/a
AR271	SRR4417556	896,578	n/a	0	0	0	n/a
AR272	SRR4417557	777,887	n/a	0	0	0	n/a
AR352	SRR6985662	7,982,657	n/a	1	14	0	n/a
AR353	SRR6799371	3,957,778	n/a	1	10	0	n/a
AR354	SRR6799376	4,153,625	n/a	0	0	0	n/a
AR357	SRR6799382	4,035,637	n/a	0	0	0	n/a
AR359	SRR6799387	3,194,525	n/a	0	0	0	n/a