## A: Pre-Trycycler assembly

Assembly A:
    contig_1: TCGGCGTGTGGTCTAAAGACTCCGGATGGGGCGTCATGGTTGATTCATCGATAATTTTC
    contig_2: AGCGTTGTACG

Assembly B:
    contig_1: GACGCCCCCATCCGGAGTCTTTAGGACCACTCGCCGAGAAAAATTATCGATGAATCACCA
    contig_2: TTGTAGCGAGCG
    contig_3: AAAAAA

Assembly C:
    contig_1: GCCGAGAAAAATTATCGATGAATCAACCATGACGCCCCATCCGGAGTCTTTAGACCACCGCC

Assembly D:
    contig_1: GATCCGGATGGGGCGTCATGGTTGATTCATCGATAATTTTTCTCGGCGGGTGGTCTAAA
    contig_2: AACGCCGCTACAAC

As input, Trycycler takes multiple different assemblies of the same genome. These can be generated using different assemblers and/or different read subsets.

## B: Clustering contigs

Cluster 1:
    A_contig_1: TCGGCGTGTGGTCTAAAGACTCCGGATGGGGCGTCATGGTTGATTCATCGATAATTTTC
    B_contig_1: GACGCCCCCATCCGGAGTCTTTAGGACCACTCGCCGAGAAAAATTATCGATGAATCACCA
    C_contig_1: GCCGAGAAAAATTATCGATGAATCAACCATGACGCCCCATCCGGAGTCTTTAGACCACCGCC
    D_contig_1: GATCCGGATGGGGCGTCATGGTTGATTCATCGATAATTTTTCTCGGCGGGTGGTCTAAA

Cluster 2:
    A_contig_2: AGCGTTGTACG
    B_contig_2: TTGTAGCGAGCG
    D_contig_2: AACGCCGCTACAAC

Cluster 3:
    B_contig_3: AAAAAA

cluster_2_A_contig_2
cluster_2_B_contig_3
cluster_2_D_contig_2
cluster_3_B_contig_2
cluster_1_A_contig_1
cluster_1_C_contig_1
cluster_1_D_contig_1
cluster_1_B_contig_1

Contigs from all assemblies are clustered based on their $k$-mer content. Trycycler makes a tree of the contig relationships to help users distinguish good clusters (which represent completely assembled replicons) vs bad clusters (which contain spurious, fragmented or incorrectly assembled sequences).

## C: Reconciling contigs

Normalise strands and fix circularisation:

Cluster 1:
    A_contig_1: GAAAATTATCGATGAATCAACCATGACGCCCCATCCGGAGTCTTTAGACCACACGCCGA
    B_contig_1: GACGCCCCCATCCGGAGTCTTTAGGACCACTCGCCGAGAAAAATTATCGATGAATCACCAT
    C_contig_1: GCCGAGAAAAATTATCGATGAATCAACCATGACGCCCCATCCGGAGTCTTTAGACCACC
    D_contig_1: TTTAGACCACCCGCCGAGAAAAATTATCGATGAATCAACCATGACGCCCCATCCGGATC

Cluster 2:
    A_contig_2: CGTACAACGCT
    B_contig_2: CGCTCGCTACAA
    D_contig_2: AACGCCGCTAC

Contig sequences are flipped to their reverse complement as necessary to ensure that all sequences within each cluster are on the same strand. For circular clusters, sequences are aligned to each other to repair circularisation issues: trimming overlapping bases or adding missing bases.

Rotate to consistent start:

Cluster 1:
    A_contig_1: ATGAATCAACCATGACGCCCCATCCGGAGTCTTTAGACCACACGCCGAGAAAATTATCG
    B_contig_1: ATGAATCACCATGACGCCCCCATCCGGAGTCTTTAGGACCACTCGCCGAGAAAAATTATCG
    C_contig_1: ATGAATCAACCATGACGCCCCATCCGGAGTCTTTAGACCACCGCCGAGAAAAATTATCG
    D_contig_1: ATGAATCAACCATGACGCCCCATCCGGATCTTTAGACCACCCGCCGAGAAAAATTATCG

Cluster 2:
    A_contig_2: GCTCGTACAAC
    B_contig_2: GCTCGCTACAAC
    D_contig_2: GCCGCTACAAC

For each circular cluster, a starting sequence is identified (using a standard coding sequence, if possible) and the sequences are rotated to have a consistent start/end. Each cluster's sequences are now ready for global multiple sequence alignment.

## D: Multiple sequence alignment

Cluster 1:
    A_contig_1: ATGAATCAACCATGACGCCCC–ATCCGGAGTCTTTAG–ACCACACGCCGAGAAAA–TTATCG
    B_contig_1: ATGAATC–ACCATGACGCCCCCATCCGGAGTCTTTAGGACCACTCGCCGAGAAAAATTATCG
    C_contig_1: ATGAATCAACCATGACGCCCC–ATCCGGAGTCTTTAG–ACCAC–CGCCGAGAAAAATTATCG
    D_contig_1: ATGAATCAACCATGACGCCCC–ATCCGGA–TCTTTAG–ACCACCCGCCGAGAAAAATTATCG

Cluster 2:
    A_contig_2: GCTCG–TACAAC
    B_contig_2: GCTCGCTACAAC
    D_contig_2: GC–CGCTACAAC

Trycycler uses MUSCLE to produce a global multiple sequence alignment for each of the clusters.

## E: Partitioning reads

All reads:

    CTCGCC  AATTAT  AGAAAA  CTCGCT
    GAGAAA  TTAGAC  AACGCT  TCGCTA
    AGACCA  CGAGAA  CCGCCG  GACCAC
    TCTTTA  CACTCG  CGGAGT  CGCTCG
    ATCAAC  GCTCGC  GAAAAA  AACCAT
    GTCTTT  CCGCTA  GTACAA  CACCAT
    ACCACA  TACAAC  TGACGC  CCCATC
    ATGACG  CGCCGA  CTACAA  ACGCCG
    TCCGGA  AAAAAT  GCTACA  GGAGTC
    CATGAC  GCCCCA  ACAACG  GATGAA

Cluster 1 reads:
    CTCGCC  AATTAT  AGAAAA  GAGAAA
    TTAGAC  AGACCA  CGAGAA  CCGCCG
    GACCAC  TCTTTA  CACTCG  CGGAGT
    ATCAAC  GAAAAA  AACCAT  GTCTTT
    CACCAT  ACCACA  TGACGC  CCCATC
    ATGACG  CGCCGA  TCCGGA  AAAAAT
    GGAGTC  CATGAC  GCCCCA  GATGAA

Cluster 2 reads:
    CTCGCT  AACGCT  TCGCTA  CGCTCG
    GTACAA  GCTCGC  CCGCTA  TACAAC
    CTACAA  ACGCCG  GCTACA  ACAACG

Reads are aligned to each contig sequence and assigned to the cluster to which they best align.

## F: Generating a consensus

Divide alignment into chunks:

Cluster 1:
ATGAATC—A—ACCATGACGCCCC—C—ATCCGGA—G—TCTTTAG—G—ACCAC—A—T—C—CGCCGAGAAAA—A—TTATCG

Cluster 2:
GC—T—CG—C—TACAAC

The multiple sequence alignment is divided into chunks: "same" chunks where the sequences agree and "different" chunks where there are multiple possible options.

Choose best option for each chunk:

Cluster 1:
ATGAATC—A—ACCATGACGCCCC—C—ATCCGGA—G—TCTTTAG—G—ACCAC—T—C—CGCCGAGAAAA—A—TTATCG

Cluster 2:
GC—T—CG—C—TACAAC

For each "different" chunk, the most popular option is chosen (as defined by the minimum total Hamming distance to other options). When there is a tie, reads are aligned to each alternative to decide which option to keep (the one with the best total read alignment score).

## G: Post-Trycycler polishing

Trycycler assembly:
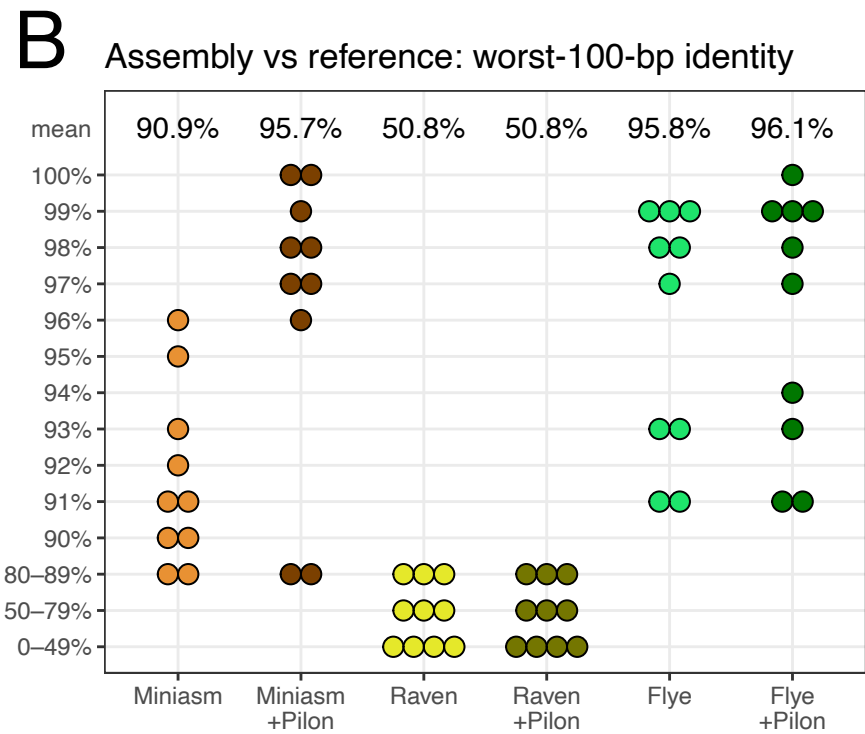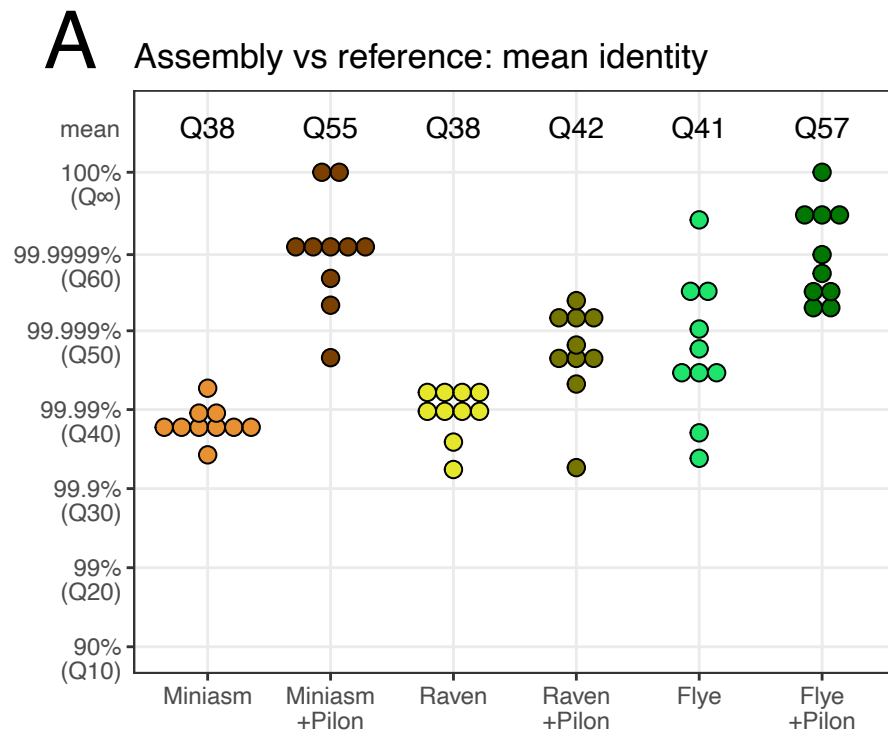ATGAATCAACCATGACGCCCCATCCGGAGTCTTTAGGACCACTCGCCGAGAAAAATTATCG          GCTCGCTACAAC

After long-read polishing:
ATGAATCAACCATGACGCCCCCATCCGGAGTCTTTAGGACCACTCGCCGAGAAAAATTATCG          GCTCGCTAGAAC

After short-read polishing:
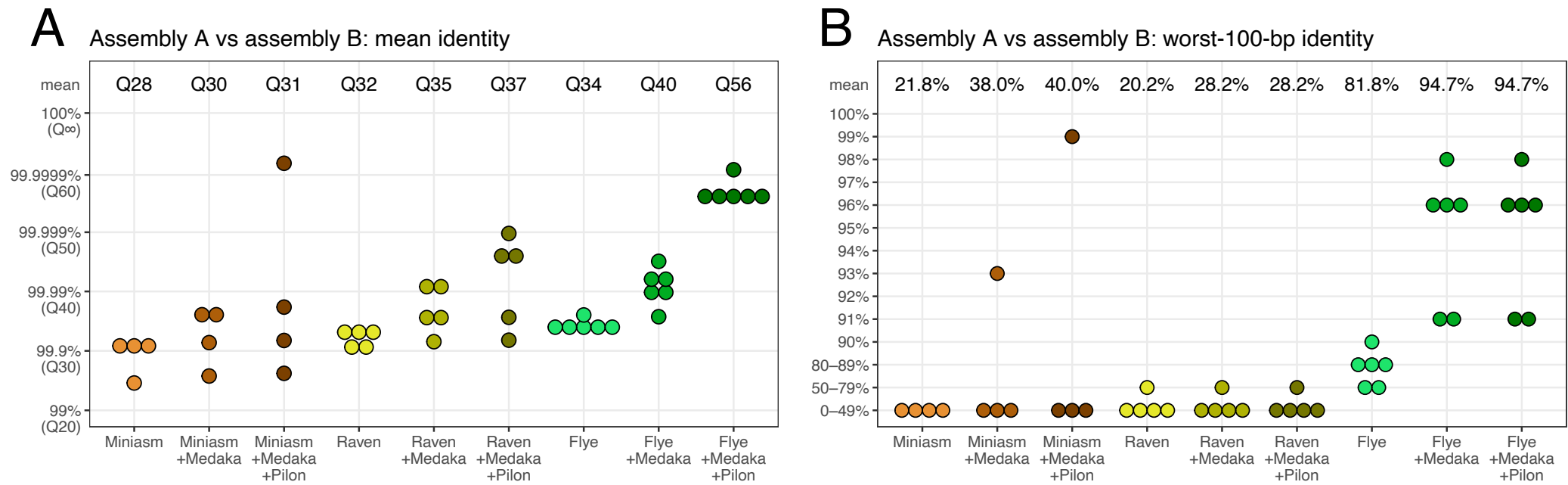ATGAATCAACCATGACGCCCCCATCCGGAGTCTTTAGGACCACTCGCCGAGAAAAATTATCG          GCTCGCTAGAAC

After Trycycler is finished, platform-specific long-read polishing (e.g. Medaka for ONT sequencing) can reduce the number of small-scale errors in the assembly. If available, short-read polishing (e.g. with Pilon) can further reduce small-scale errors.

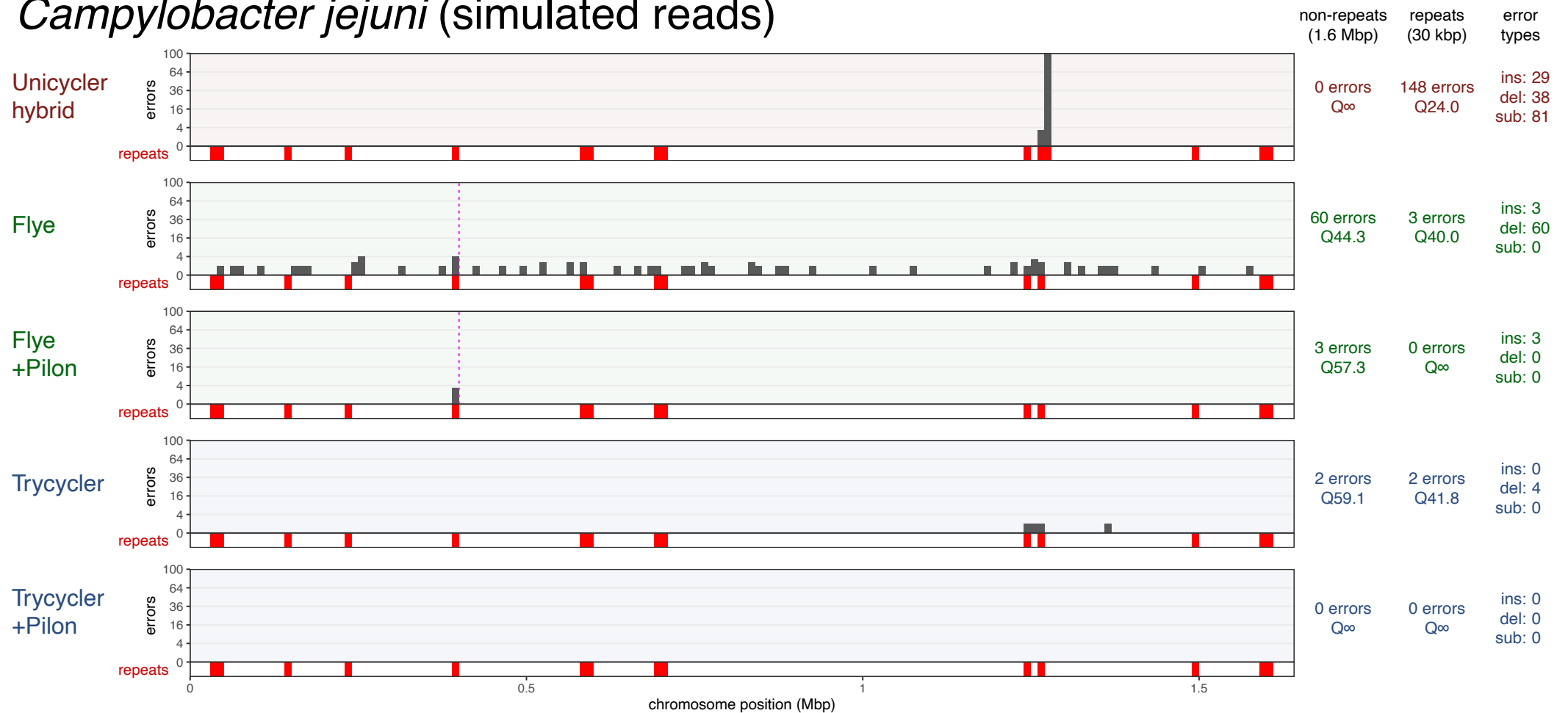**Fig. S1**: steps in the Trycycler assembly pipeline.

**Fig. S2**: results for the simulated read tests for long-read-only assemblers. This figure contains the same analyses as are shown in Figure 2, but it includes assemblies from all long-read-only assemblers (both before and after Pilon polishing).
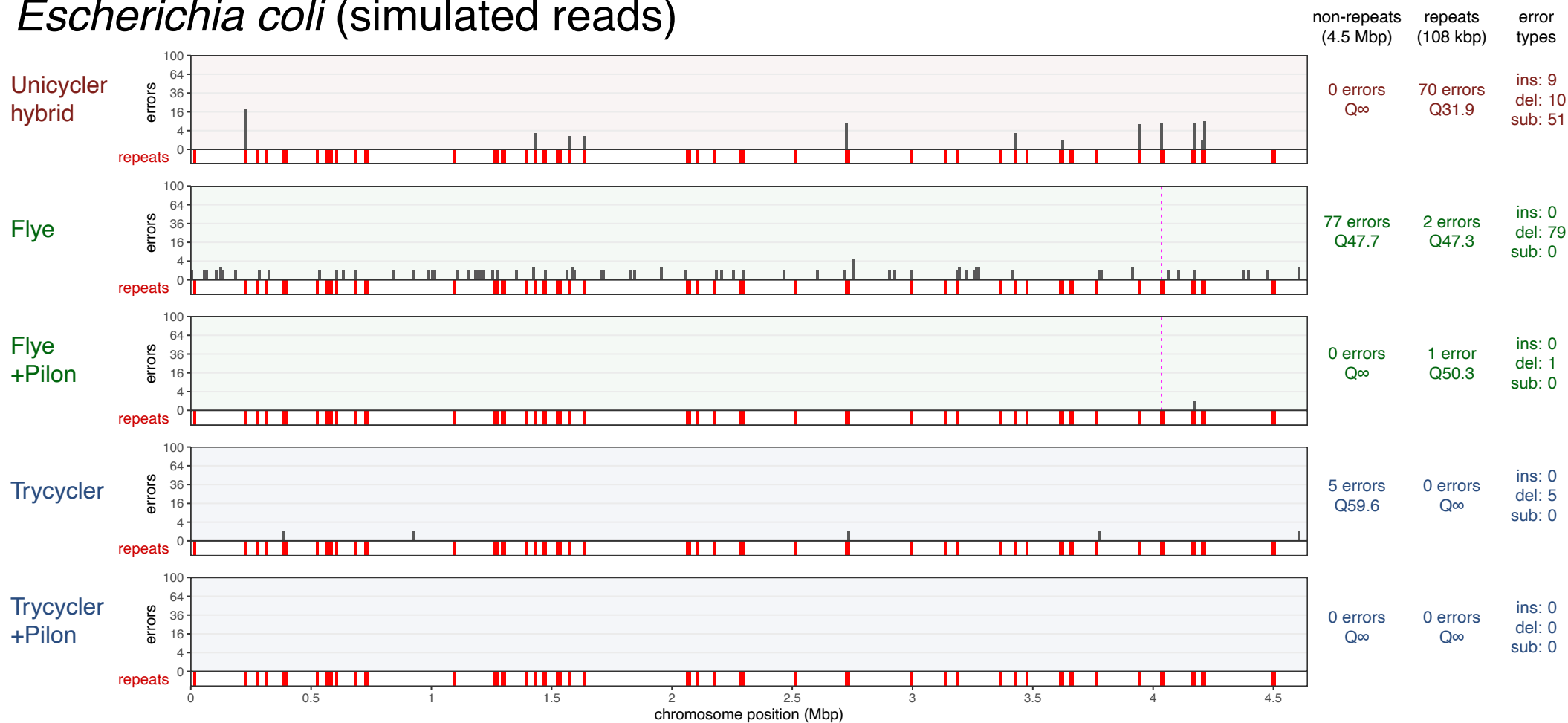
**Fig. S3**: results for the real read tests for long-read-only assemblers. This figure contains the same analyses as are shown in Figure 3, but it includes assemblies from all long-read-only assemblers (unpolished, Medaka-polished and Medaka+Pilon-polished for each). In two of the six genomes, Miniasm failed to produce a completed chromosome for both read sets, resulting in only four data points. In one of the six genomes, Raven failed to produce a completed chromosome for both read sets, resulting in only five data points.
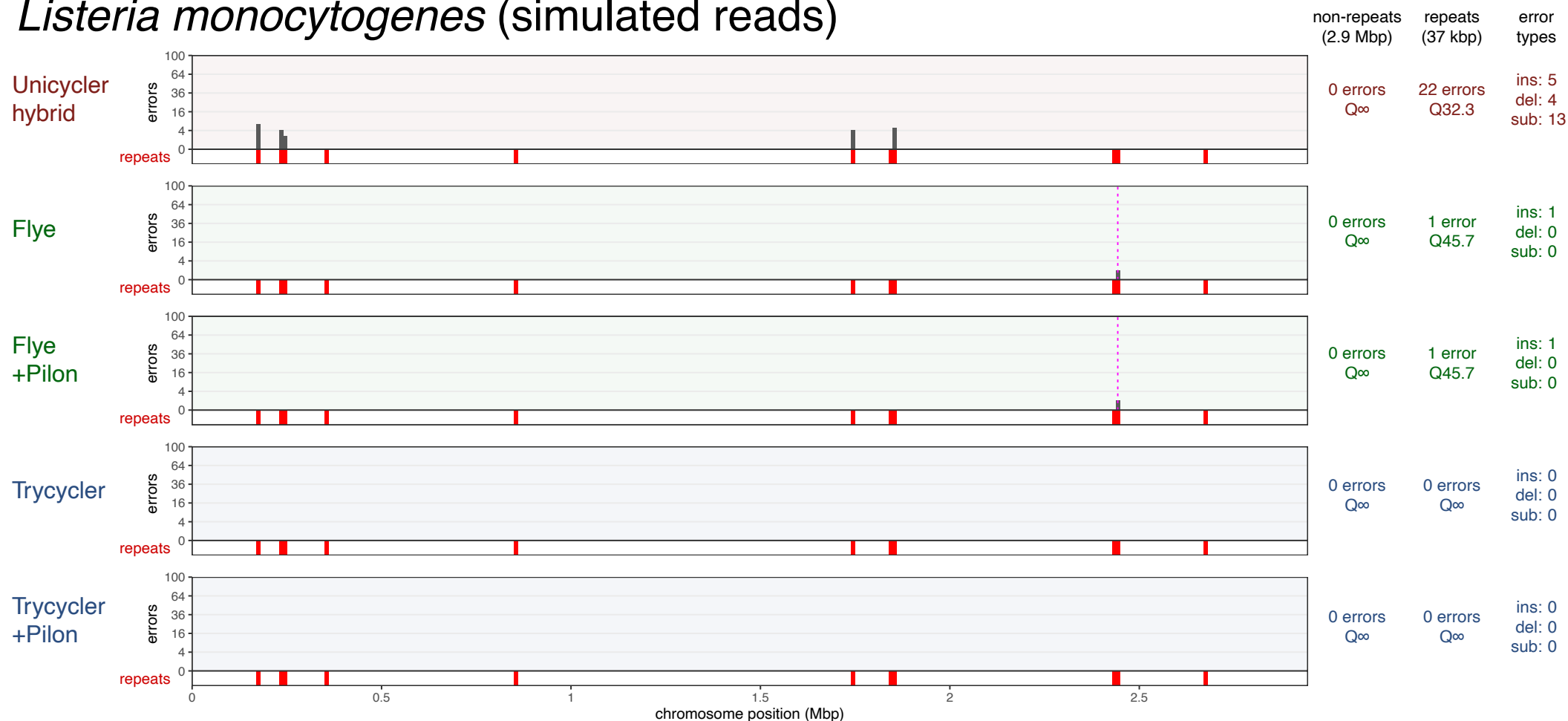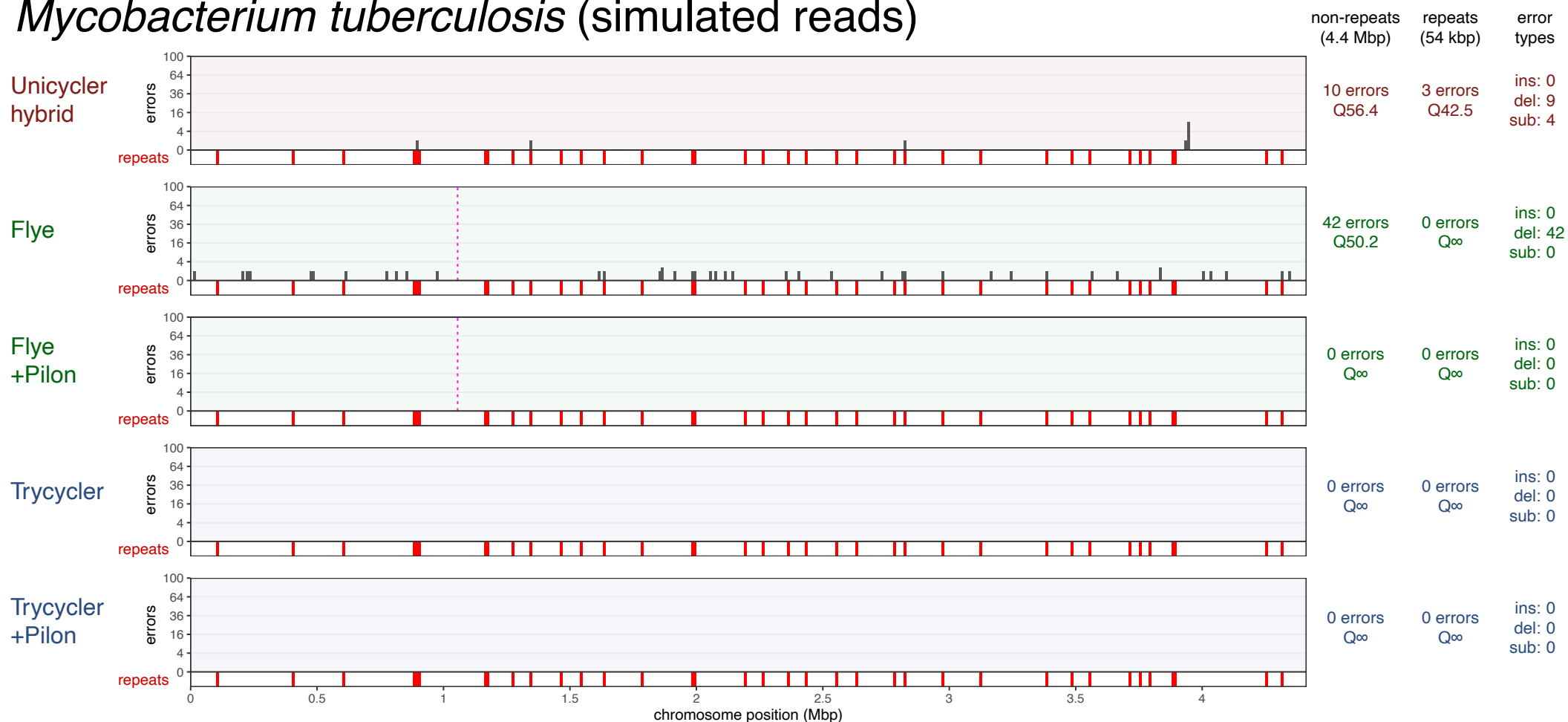
**Fig. S4-a**: error distributions (as defined by alignment to the reference sequence) for the simulated-read assemblies of the *Campylobacter jejuni* chromosome. Repetitive regions of the genome are shown as red blocks at the bottom of each plot. For Flye assemblies, the contig start/end position (where the sequence began before it was rotated to be consistent with the reference) is shown by the dashed line.

**Fig. S4-b**: error distributions (as defined by alignment to the reference sequence) for the simulated-read assemblies of the *Escherichia coli* chromosome. Repetitive regions of the genome are shown as red blocks at the bottom of each plot. For Flye assemblies, the contig start/end position (where the sequence began before it was rotated to be consistent with the reference) is shown by the dashed line.

**Fig. S4-c**: error distributions (as defined by alignment to the reference sequence) for the simulated-read assemblies of the *Klebsiella pneumoniae* chromosome. Repetitive regions of the genome are shown as red blocks at the bottom of each plot. For Flye assemblies, the contig start/end position (where the sequence began before it was rotated to be consistent with the reference) is shown by the dashed line.
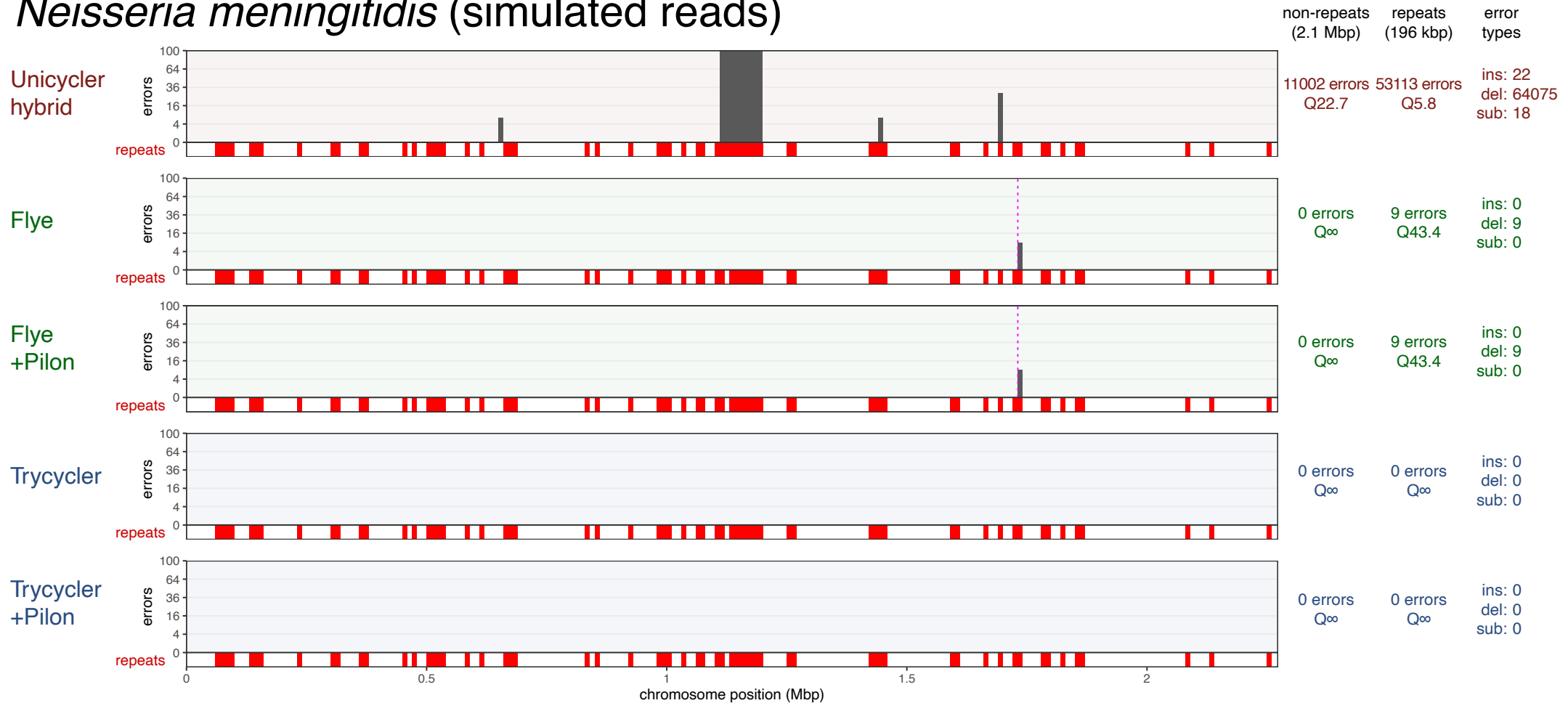
# *Listeria monocytogenes* (simulated reads)



**Fig. S4-d**: error distributions (as defined by alignment to the reference sequence) for the simulated-read assemblies of the *Listeria monocytogenes* chromosome. Repetitive regions of the genome are shown as red blocks at the bottom of each plot. For Flye assemblies, the contig start/end position (where the sequence began before it was rotated to be consistent with the reference) is shown by the dashed line.
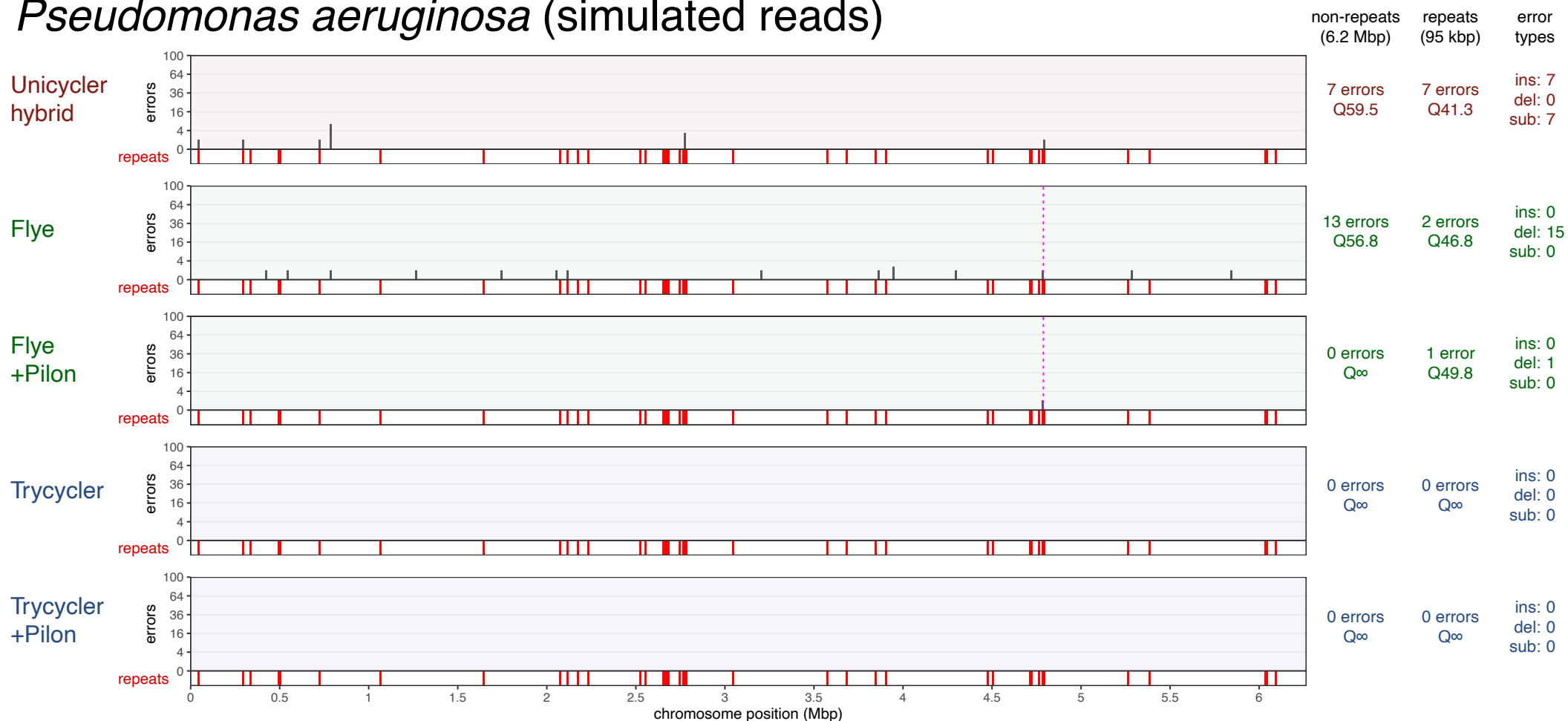
**Fig. S4-e**: error distributions (as defined by alignment to the reference sequence) for the simulated-read assemblies of the *Mycobacterium tuberculosis* chromosome. Repetitive regions of the genome are shown as red blocks at the bottom of each plot. For Flye assemblies, the contig start/end position (where the sequence began before it was rotated to be consistent with the reference) is shown by the dashed line.
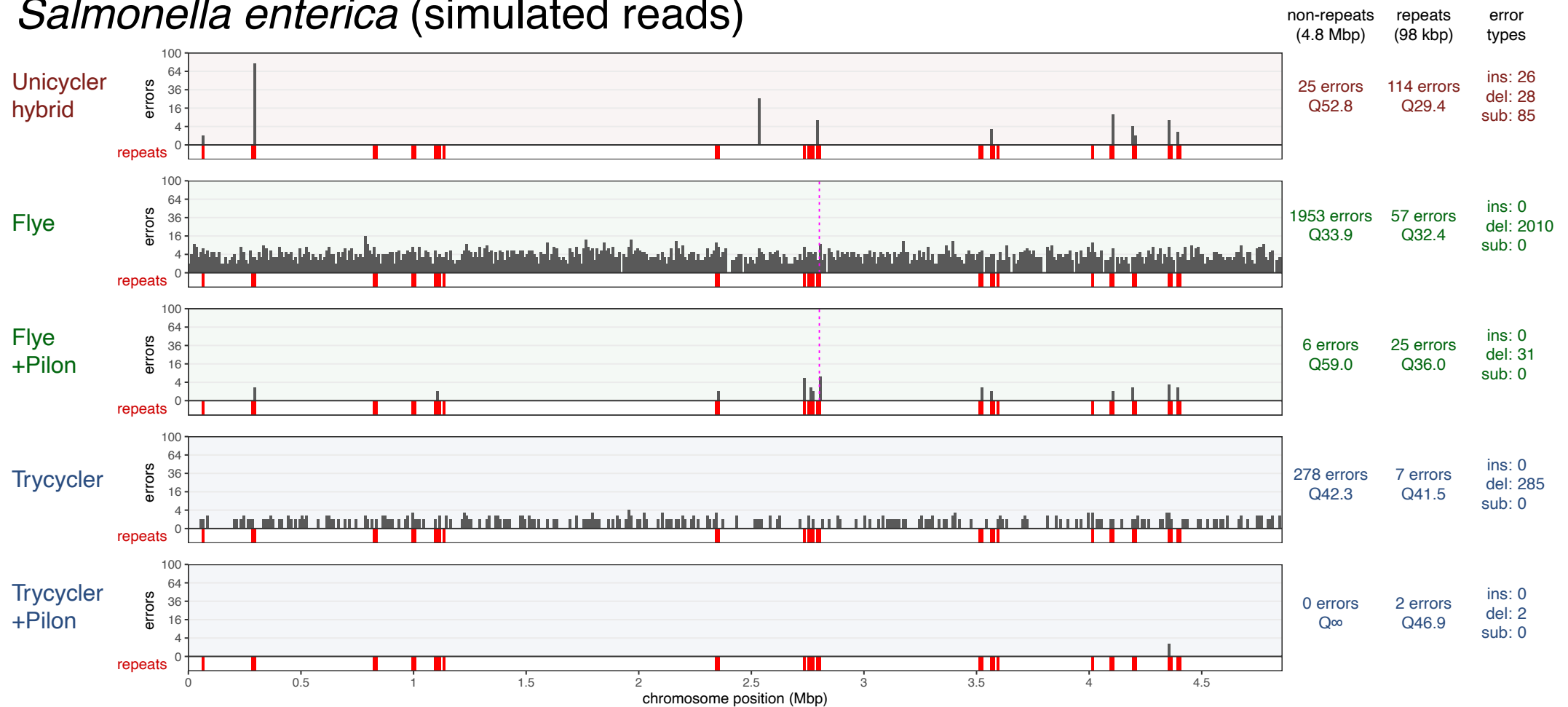
# *Neisseria meningitidis* (simulated reads)
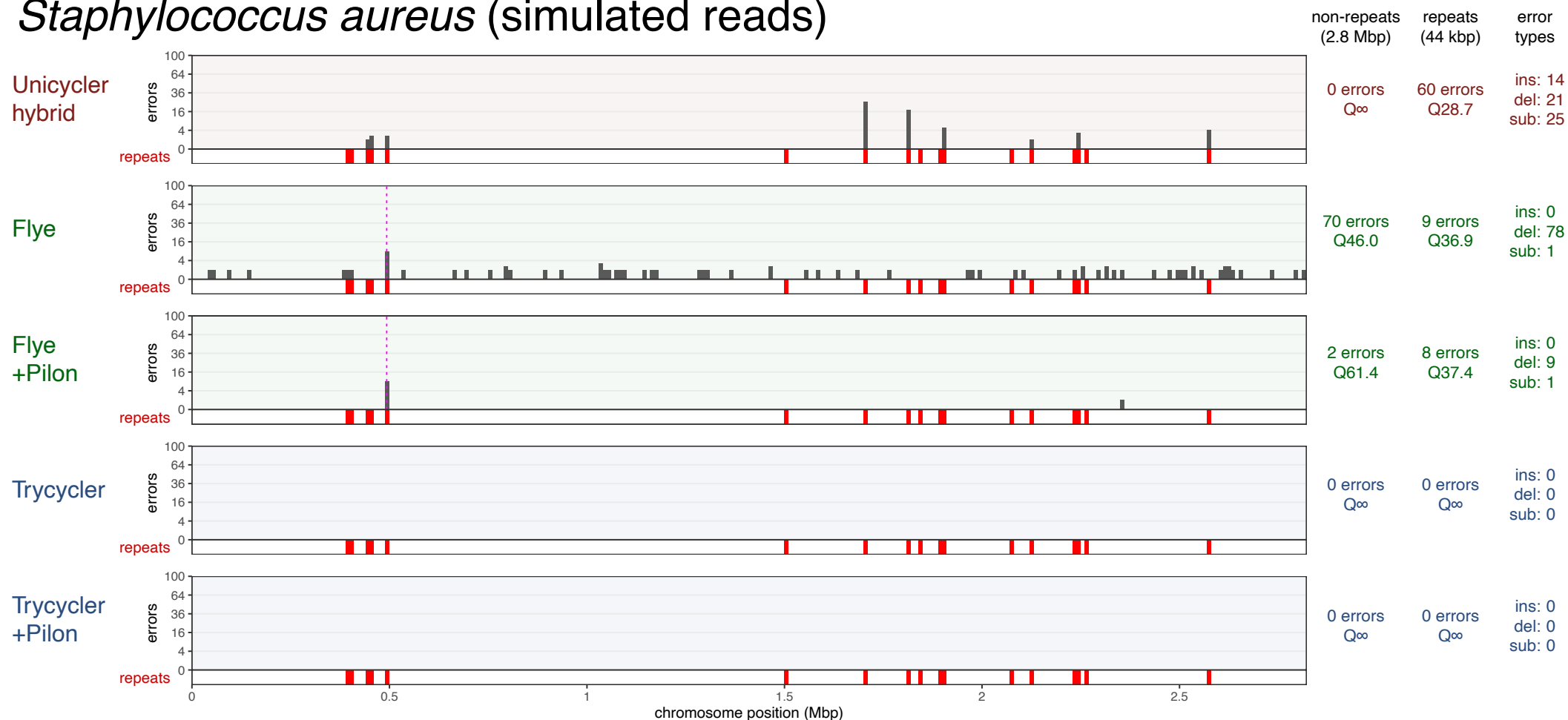


**Fig. S4-f**: error distributions (as defined by alignment to the reference sequence) for the simulated-read assemblies of the *Neisseria meningitidis* chromosome. Repetitive regions of the genome are shown as red blocks at the bottom of each plot. For Flye assemblies, the contig start/end position (where the sequence began before it was rotated to be consistent with the reference) is shown by the dashed line.
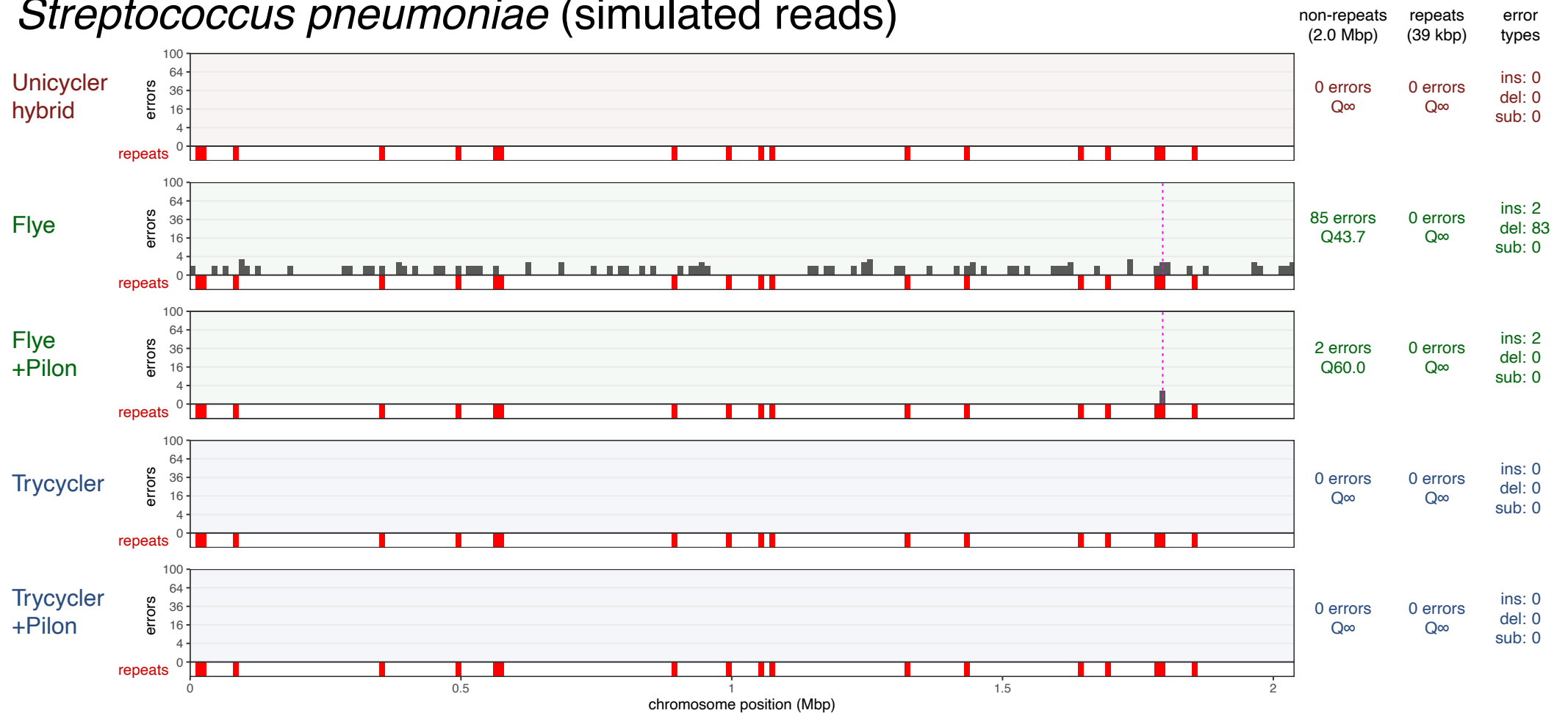
**Fig. S4-g**: error distributions (as defined by alignment to the reference sequence) for the simulated-read assemblies of the *Pseudomonas aeruginosa* chromosome. Repetitive regions of the genome are shown as red blocks at the bottom of each plot. For Flye assemblies, the contig start/end position (where the sequence began before it was rotated to be consistent with the reference) is shown by the dashed line.

**Fig. S4-h**: error distributions (as defined by alignment to the reference sequence) for the simulated-read assemblies of the *Salmonella enterica* chromosome. Repetitive regions of the genome are shown as red blocks at the bottom of each plot. For Flye assemblies, the contig start/end position (where the sequence began before it was rotated to be consistent with the reference) is shown by the dashed line.
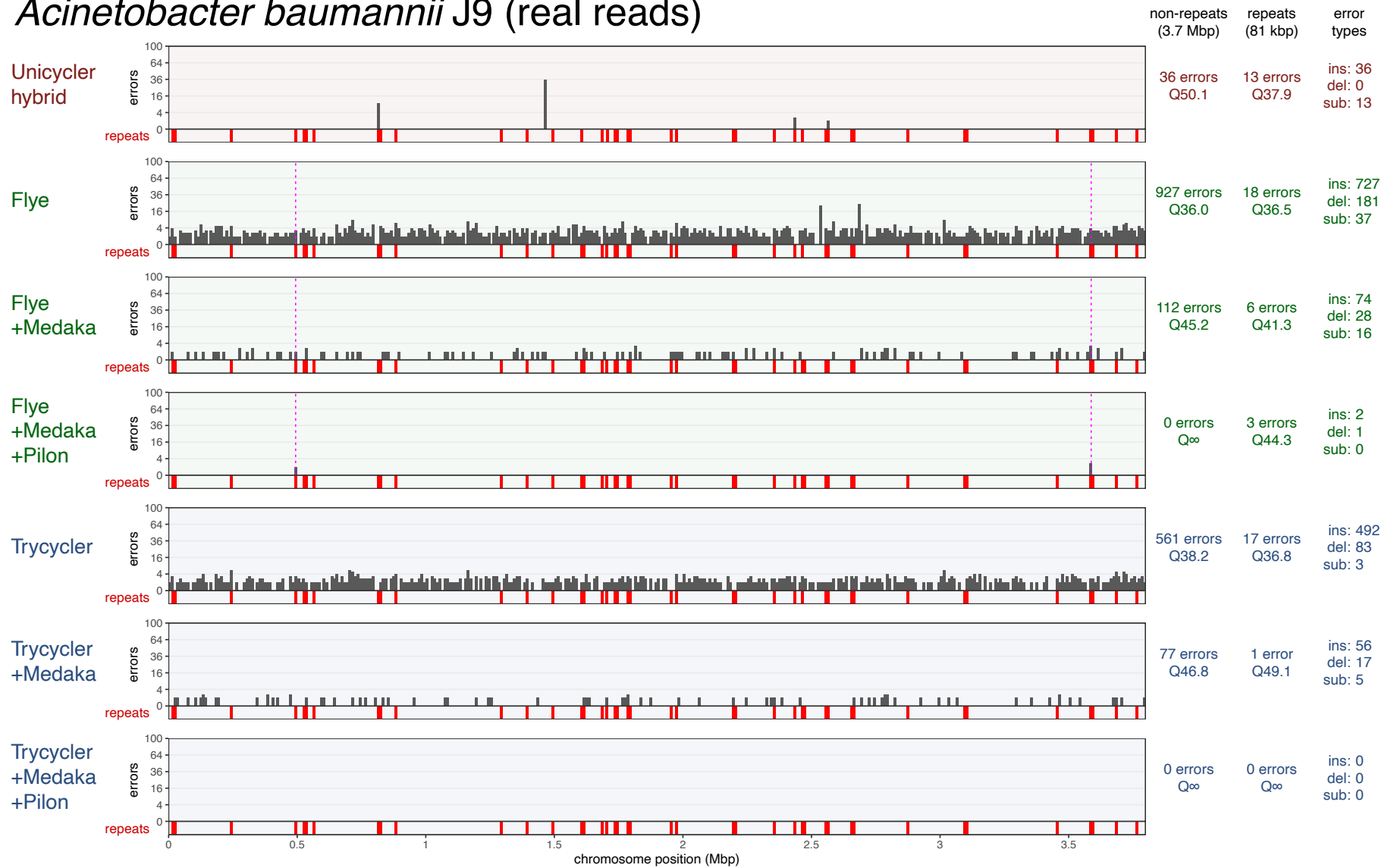
**Fig. S4-i**: error distributions (as defined by alignment to the reference sequence) for the simulated-read assemblies of the *Staphylococcus aureus* chromosome. Repetitive regions of the genome are shown as red blocks at the bottom of each plot. For Flye assemblies, the contig start/end position (where the sequence began before it was rotated to be consistent with the reference) is shown by the dashed line.
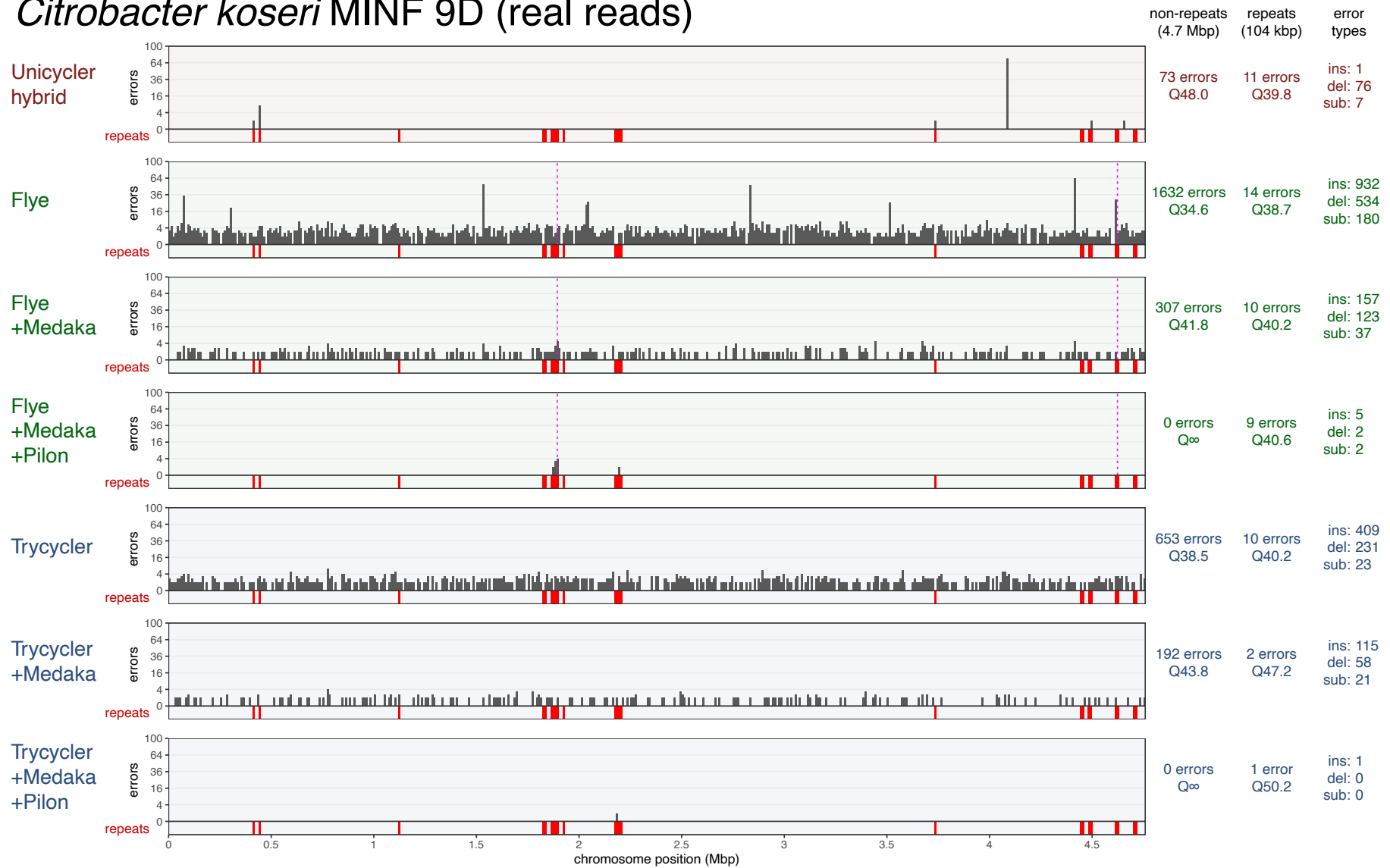
**Fig. S4-j**: error distributions (as defined by alignment to the reference sequence) for the simulated-read assemblies of the *Streptococcus pneumoniae* chromosome. Repetitive regions of the genome are shown as red blocks at the bottom of each plot. For Flye assemblies, the contig start/end position (where the sequence began before it was rotated to be consistent with the reference) is shown by the dashed line.
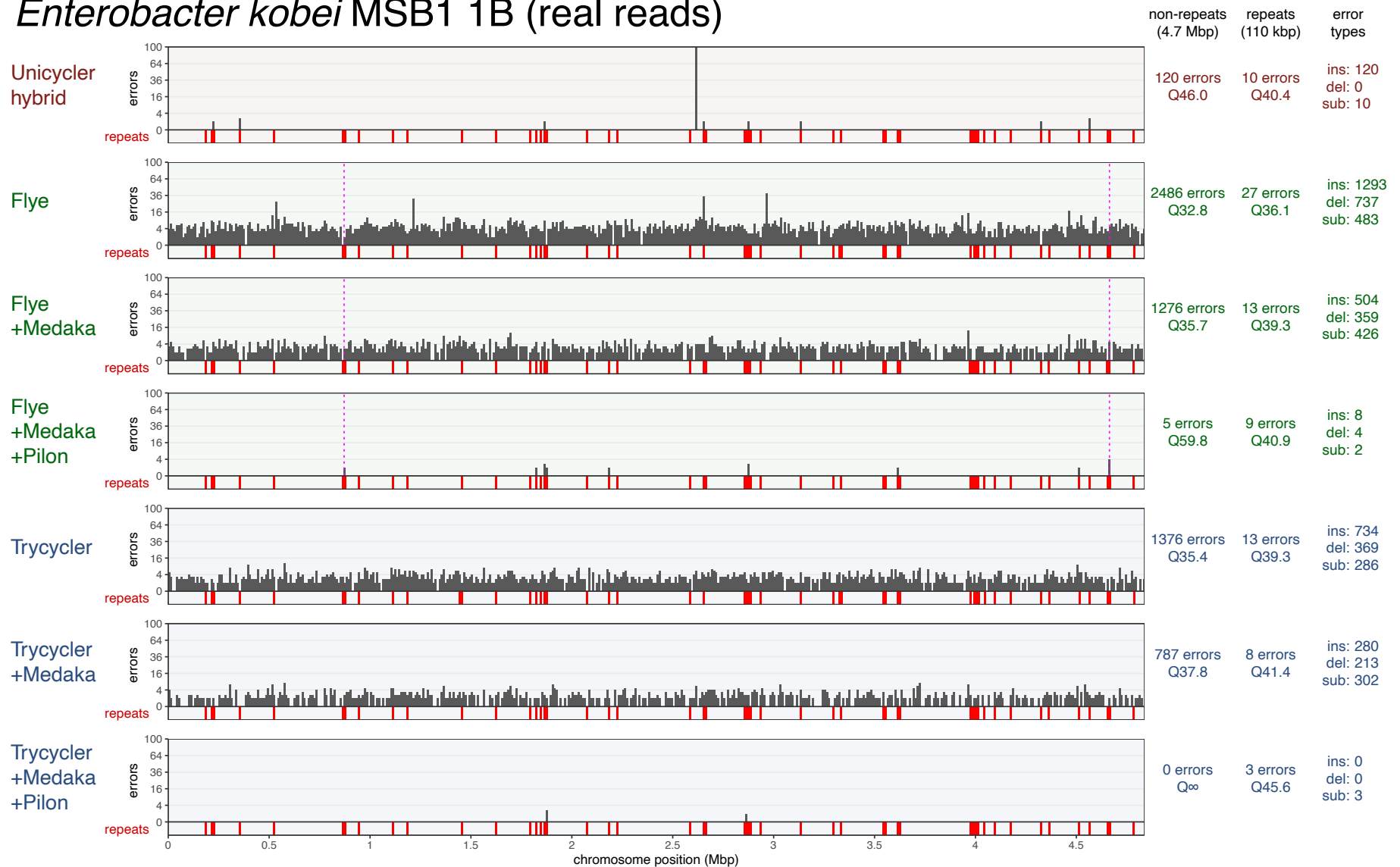
**Fig. S4-k**: error distributions (as defined by alignment between two independent assemblies) for the real-read assemblies of the *Acinetobacter baumannii* chromosome. Repetitive regions of the genome are shown as red blocks at the bottom of each plot. For Flye assemblies, the contig start/end positions (where the sequences began before they were rotated to be consistent with the reference) are shown by the dashed lines.

**Fig. S4-l**: error distributions (as defined by alignment between two independent assemblies) for the real-read assemblies of the *Citrobacter koseri* chromosome. Repetitive regions of the genome are shown as red blocks at the bottom of each plot. For Flye assemblies, the contig start/end positions (where the sequences began before they were rotated to be consistent with the reference) are shown by the dashed lines.
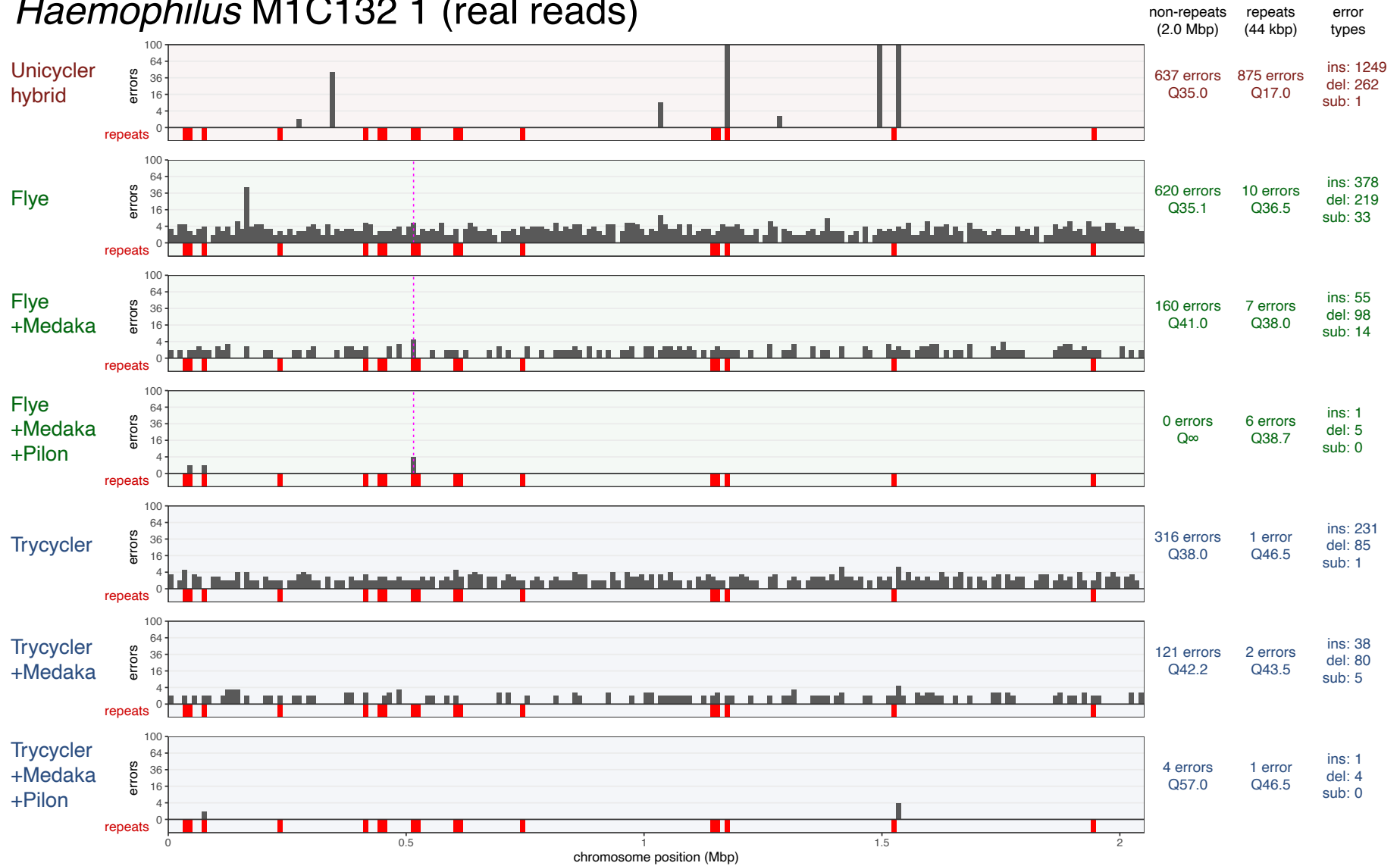
**Fig. S4-m**: error distributions (as defined by alignment between two independent assemblies) for the real-read assemblies of the *Enterobacter kobei* chromosome. Repetitive regions of the genome are shown as red blocks at the bottom of each plot. For Flye assemblies, the contig start/end positions (where the sequences began before they were rotated to be consistent with the reference) are shown by the dashed lines.
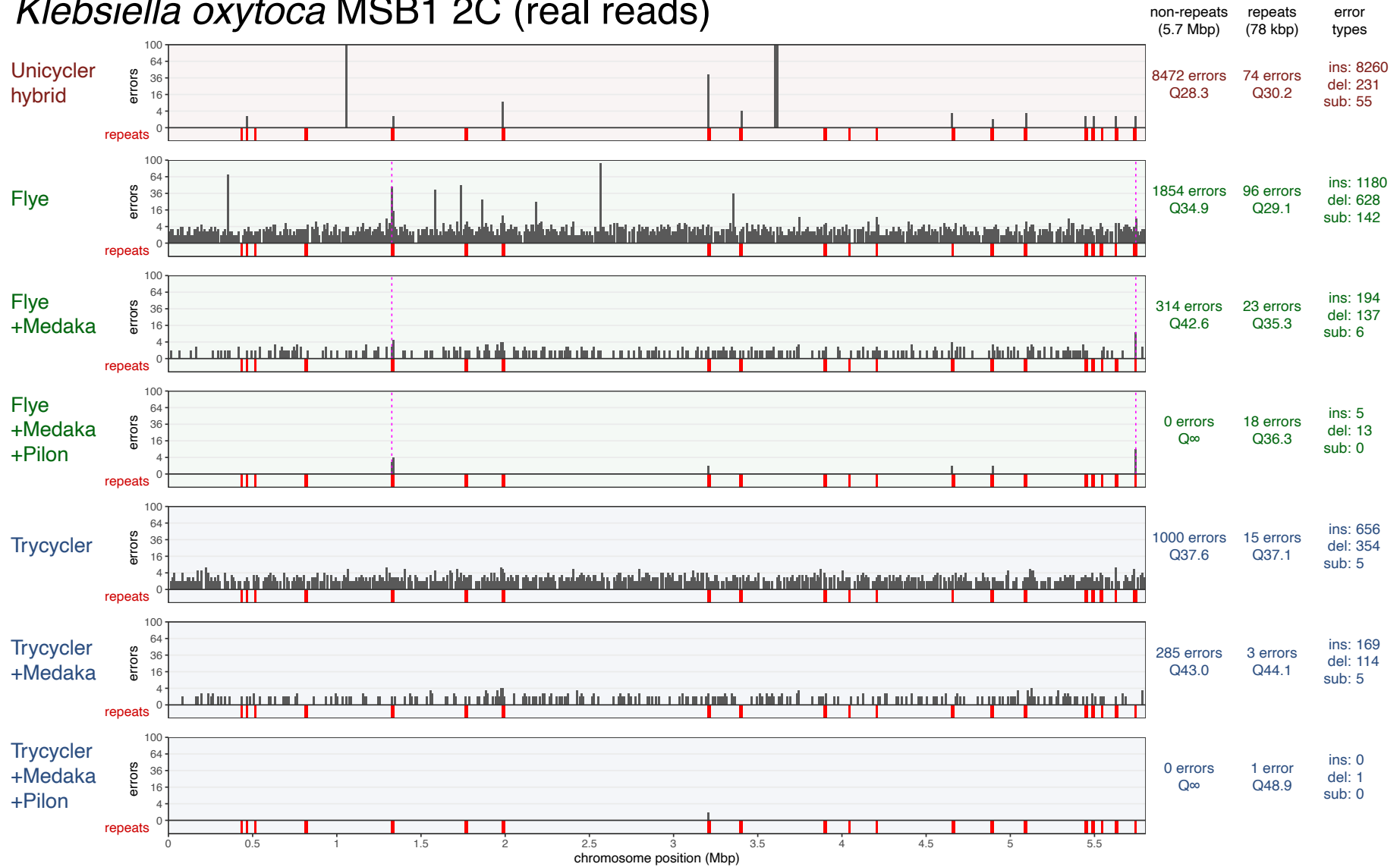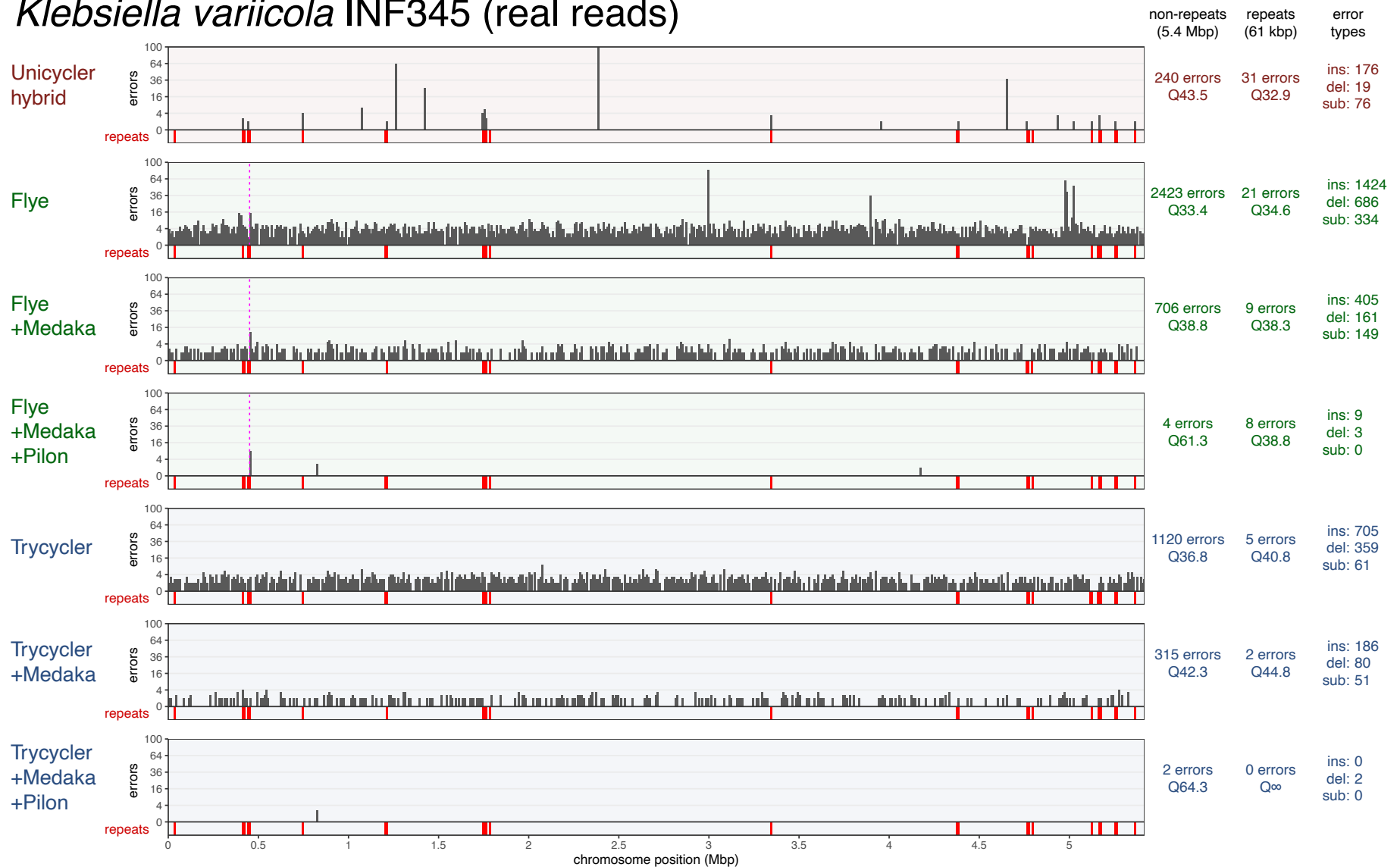
**Fig. S4-n**: error distributions (as defined by alignment between two independent assemblies) for the real-read assemblies of the *Haemophilus* chromosome. Repetitive regions of the genome are shown as red blocks at the bottom of each plot. For Flye assemblies, the contig start/end positions (where the sequences began before they were rotated to be consistent with the reference) are shown by the dashed lines.
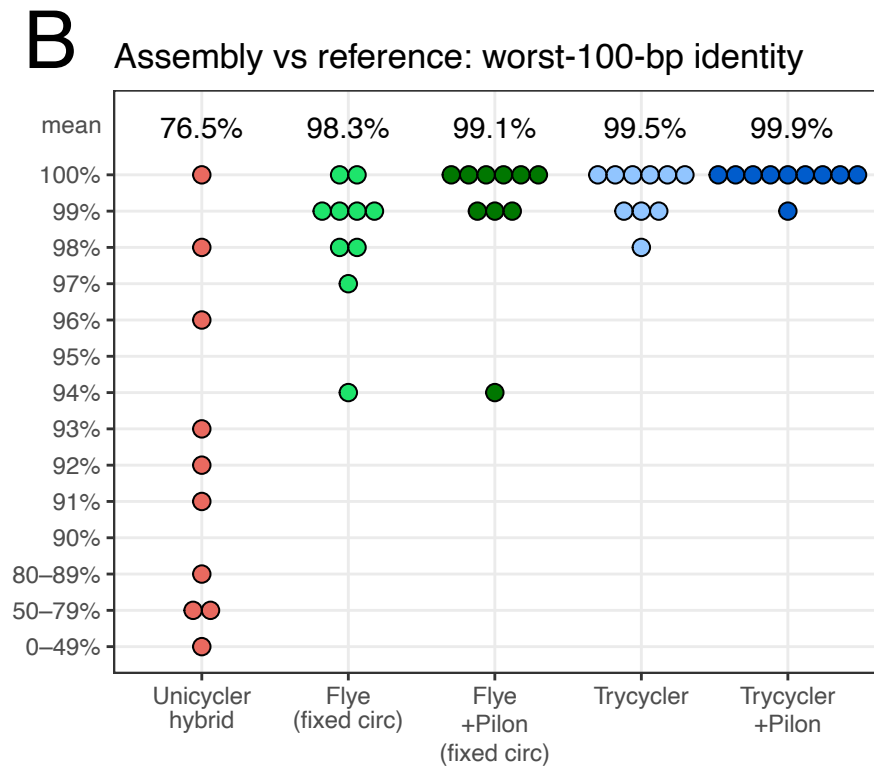
**Fig. S4-o**: error distributions (as defined by alignment between two independent assemblies) for the real-read assemblies of the *Klebsiella oxytoca* chromosome. Repetitive regions of the genome are shown as red blocks at the bottom of each plot. For Flye assemblies, the contig start/end positions (where the sequences began before they were rotated to be consistent with the reference) are shown by the dashed lines.
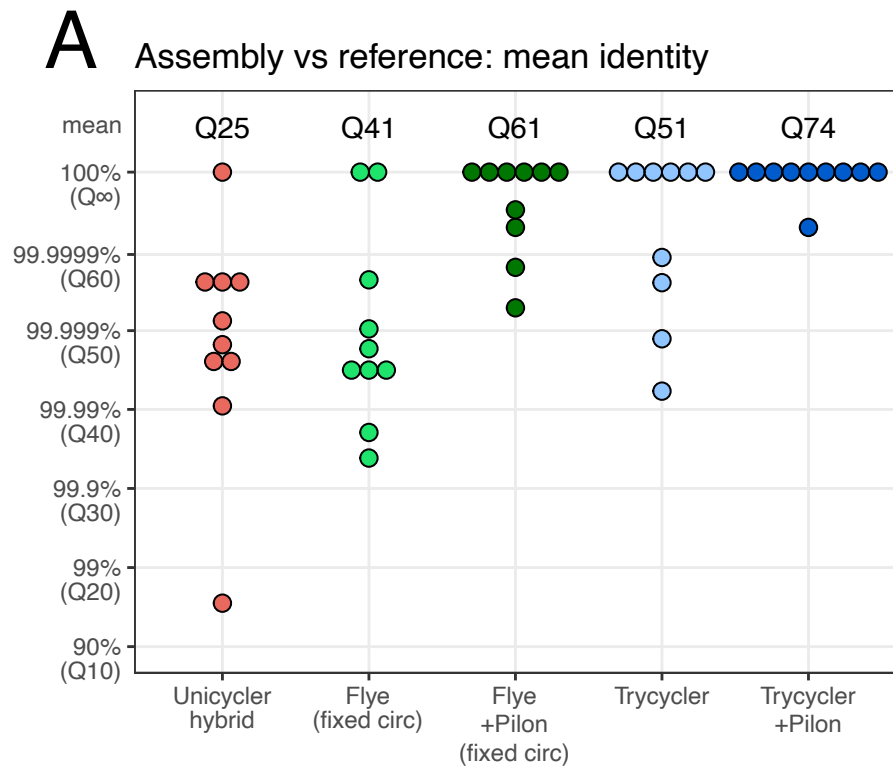
**Fig. S4-p**: error distributions (as defined by alignment between two independent assemblies) for the real-read assemblies of the *Klebsiella variicola* chromosome. Repetitive regions of the genome are shown as red blocks at the bottom of each plot. For Flye assemblies, the contig start/end positions (where the sequences began before they were rotated to be consistent with the reference) are shown by the dashed lines.
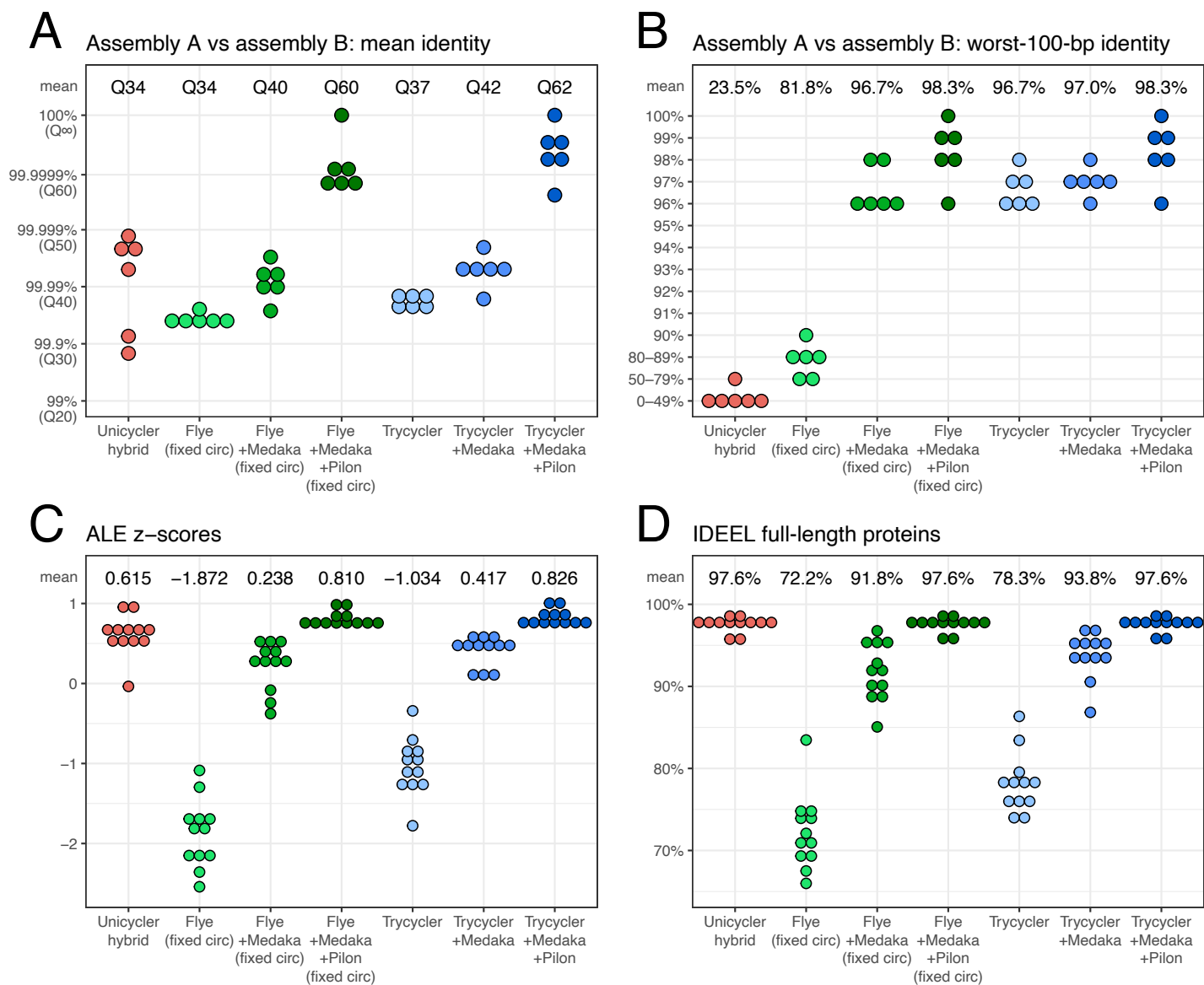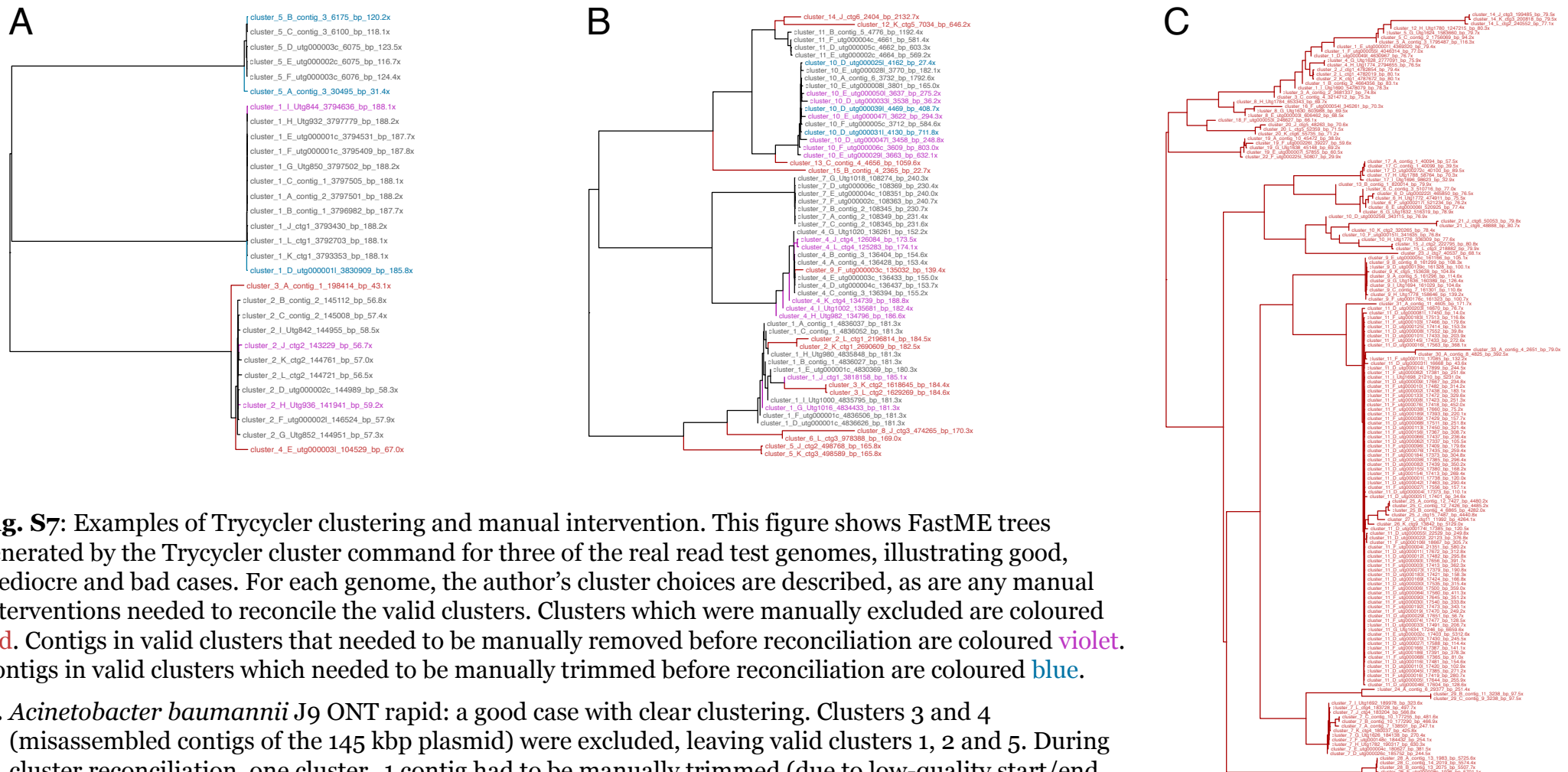
**Fig. S5**: results for the simulated read tests. This figure contains the same analyses as are shown in Figure 2, but the Flye assemblies have had their circularisation manually repaired.

**Fig. S6**: results for the real read tests. This figure contains the same analyses as are shown in Figure 3, but the Flye assemblies have had their circularisation manually repaired.

**Fig. S7**: Examples of Trycycler clustering and manual intervention. This figure shows FastME trees generated by the Trycycler cluster command for three of the real read test genomes, illustrating good, mediocre and bad cases. For each genome, the author's cluster choices are described, as are any manual interventions needed to reconcile the valid clusters. Clusters which were manually excluded are coloured red. Contigs in valid clusters that needed to be manually removed before reconciliation are coloured violet. Contigs in valid clusters which needed to be manually trimmed before reconciliation are coloured blue.

**A.** *Acinetobacter baumannii* J9 ONT rapid: a good case with clear clustering. Clusters 3 and 4 (misassembled contigs of the 145 kbp plasmid) were excluded, leaving valid clusters 1, 2 and 5. During cluster reconciliation, one cluster_1 contig had to be manually trimmed (due to low-quality start/end sequence), one cluster_1 contig had to be manually removed (due to poor pairwise alignment), two cluster_2 contigs had to be manually removed (due to being incomplete) and two cluster_5 contigs had to be manually trimmed (due to excessive length).

**B.** *Enterobacter kobei* MSB1_1B ONT rapid: a mediocre case with more complex clusters. Clusters 2, 3, 5, 6, 8 (misassembled contigs of the chromosome), 9 (misassembled contig of the 136 kbp plasmid), 12, 13, 14, and 15 (misassembled contigs of the small plasmids) were excluded, leaving valid clusters 1, 4, 7, 10 and 11. During cluster reconciliation, two cluster_1 contigs had to be manually removed (due to being incomplete or poor pairwise alignment), five cluster_4 contigs had to be manually removed (due to being incomplete), three cluster_10 contigs had to be manually trimmed (due to excessive length) and six cluster_10 contigs had to be manually removed (due to being incomplete or unable to circularise).

**C.** *Serratia marcescens* 17-147-1671 ONT rapid: a bad case where valid clusters were unclear. Insufficient read length and genome heterogeneity both contributed to the poor results. Without good clusters, it was not possible to proceed with Trycycler assembly.